

Strategies in Social Software

Jan van Eijck

CWI and ILLC

Abstract. Viewing the way society has defined its rules and mechanisms as “social software”, we want to understand how people behave given their understanding of the societal rules and given their wish to further their interest *as they conceive it*, and how social mechanisms should be designed to suit people furthering their interest *as they conceive it*. This chapter is written from the perspective of strategic game theory, and uses strategic game scenarios and game transformations to analyze societal mechanisms.

Know the enemy and know
yourself; in a hundred battles you
will never be in peril.

When you are ignorant of the
enemy but know yourself, your
chances of winning and losing are
equal.

If ignorant both of your enemy
and of yourself, you are certain in
every battle to be in peril.

Sun Tzu [450BC]

1 What is Social Software?

Social software is a term coined by Parikh [2002] for social procedures designed to regulate societal behaviour. Many of these mechanisms are connected to strategic reasoning. Parikh’s paper is a plea to view social procedures as algorithms, and to study them with the methods of logic and theoretical computer science. See Chwe [2001] for illuminating use of this methodology to explain rituals in societies. The discourses in Van Eijck and Verbrugge [2009] give further informal introduction.

In fact, design and analysis of social software is at the intersection of various academic disciplines. It is related to what is called mechanism design in game theory and economics [Hurwicz and Reiter, 2006], to behavioral architecture in political theory [Thaler and Sunstein, 2008], to

rational decision making in decision theory [Gilboa, 2010, Körner, 2008], to multi-agent theory in artificial intelligence [Shoham and Leyton-Brown, 2008], and to auction theory in economics [Milgrom, 2004, Krishna, 2009], to name but a few. If it is different from any of these, then the difference lies in the emphasis on the use of tools from logic and theoretical computer science, while bearing in mind that the objects of study are humans rather than microprocessors.

Indeed, the human participants in a social procedure are quite different from microprocessors interacting in a calculation. Unlike microprocessors, humans are, to some extent at least, aware of the social mechanisms they are involved in. This awareness may inspire them to act strategically: to use their knowledge of the mechanism to improve their welfare. Conversely, social mechanisms may be designed with this behaviour in mind. Ideally, the design of the mechanism should ensure that it cannot be exploited, or at least that the mechanism is resistant to exploitation attempts.

A central topic in economics, and in a branch of game theory called evolutionary game theory, is to explain how selfish behaviour can lead to beneficial outcomes on the societal level [Sigmund, 2010]. On the other hand, some economists have argued convincingly that modelling human beings as selfish misses a point: the scope of economics should be broadened to the study of interacting agents maximizing welfare *as they conceive it* [Becker, 1993].

Undoubtedly the most famous social mechanism that employs strategic behaviour is the mechanism of the free market, supposedly guided by Adam Smith's invisible hand, the hidden mechanism that fuses actions motivated by individual interests into a self-regulating social mechanism beneficent to all. In Smith's famous words:

It is not from the benevolence of the butcher, the brewer or the baker that we expect our dinner, but from their regard to their own interest. We address ourselves not to their humanity but to their self-love, and never talk to them of our necessities but of their advantages. Nobody but a beggar chooses to depend chiefly upon the benevolence of their fellow-citizens. [Smith, (1776, Book 1, Chapter II)]

The mechanism of the free market is designed, so to speak, to put individual self-interest at the service of society. But in other cases, social mechanism design has as one of its goals to discourage 'strategic behaviour', which now is taken to mean *misuse* of the mechanism to guarantee a

better individual outcome. An example of this is auction design. Vickrey auctions, where bids are made simultaneously, bidders do not know the values of the other bids, and the highest bidder wins but pays the second-highest bid [Vickrey, 1961], are an example. The design of the procedure discourages the bidders from making bids that do not reflect their true valuation. Google and Yahoo use variations on this when auctioning advertisement space.

Many societal mechanisms are set up so as to ensure a desirable outcome for society. What is the social procedure that ensures that soldiers that are sent into battle actually fight? One force is the fact that the other soldiers are fighting — a chain-effect. An additional factor may be the public announcement to the effect that deserters will be shot. This reduces the choice that a soldier has to that between facing the risk of death while engaging in battle versus certain death when avoiding to fight. Surely, other factors are involved, having to do with how soldiers perceive their own behaviour and their relationship to each other and to the community they fight for: comradeship, the desire to return home as a hero rather than a coward. The point is that society has mechanisms in play to ensure that individuals behave in ways that, at first sight, are squarely against their self-interest.

Some societies stage public executions of criminals. Such stagings serve a strategic social goal: to inspire terror in aspiring criminals. Or maybe, to inspire terror in the population at large, in case the victims are convicted for political crimes. Some societies keep their citizens in check with the threat of corporal punishment — in Singapore you risk a blow with the stick for a relatively minor offense — while other societies consider such methods off-limits and barbarian. Someone interested in design and analysis of social software will want to understand such radical differences in attitude.

The mechanisms of performance-related pay and of investment bankers' bonuses are other examples of social procedures designed with the goal of influencing the strategies of workers, in order to increase productivity or profitability. While the current financial crisis is evidence of the dangers of this system, experimental economics also points out that very high rewards are in fact detrimental to performance [Ariely et al., 2009].

On the other hand, in another setting such a mechanism may have considerable advantages. The bonus-malus system (BMS) used in the insurance industry adjusts the premium that a customer has to pay to insure a certain risk according to the individual claim history of the customer. This inspires caution in claiming damage on an insured vehicle, for in-

stance, for a claim means that the customer loses her no-claim discount. Thus, the system is designed to induce strategic behaviour in the customers, for it now makes sense to not report minor damages to your car to the insurance company, as the loss of the no-claim discount outweighs the cost of the damage. This is in the interest of the insurance company, and indirectly in the interest of the public in need of car insurance, for it helps to keep premiums low. So why does the bonus system work well in this situation while creating disaster in other settings?

A general problem with social mechanism design is that technological fixes to societal problems have a tendency to misfire and create new problems, because the technological solution leads to paradoxical and unintended consequences. New roads lead to bigger traffic jams and uncontrolled growth of suburbia. Rent regulation, intended to protect tenants, may lead to poorer housing conditions for the less affluent. Worker protection laws may be a factor in causing unemployment because they make employers reluctant to hire new staff. Performance related pay may induce bankers to take unreasonable risks. Tenner [1996] gives many other examples, with insightful comments.

Still, insights from the design and analysis of algorithms can and should be applied to analysis and design in social interaction, all the time bearing in mind that agents involved in social interaction are *aware* of the societal mechanisms. This awareness often takes the form of strategic reasoning by people about how to further their best interest, *as they conceive it*, given the way society has defined its rules and mechanisms. The analysis and design of social software should take this awareness of participants into account.

The focus of this chapter is on strategic games rather than dynamic games. The structure of the chapter is as follows. In Section 2 we distinguish three levels at which strategizing might occur. In Section 3, we use the situation of the well-known prisoner's dilemma as a starting point for a discussion of what goes on in social software design and analysis. Section 4 discusses, in a game-theoretic setting, how strategies are affected by punishment, and Section 5 focusses on the influence of rewards on strategies. In Section 6, these same topics return in the tragedy of the commons scenario, closely related to the prisoner's dilemma. The theme of individual versus collective is brought out even more poignantly in renunciation games, presented in Section 7. Section 8 discusses the use of game scenarios in experiments about knowledge and trust in social protocols, and Section 9 concludes with a mention of logics for strategic games, together with some remarkable arguments for democracy.

2 Strategizing at Various Levels

The social problem of collective decision making involves strategizing at various levels. Consider as an example a scientific advisory board that has to rank a number of research proposals in order of quality. Participants in such a meeting strategize at various levels, and strategizing also takes place at the level of the scientific community at large.

Strategizing at the micro-level How much should I, as a participant in such a meeting, reveal about my own true preferences (or: of my own knowledge and ignorance), in order to make me maximally effective in influencing the other participants?

Strategizing at intermediate level How should the chair structure the decision making process, so as to ensure that consensus is reached and that the meeting terminates within a reasonable period of time? The chair could propose rules like “For any two proposals X and Y, once we have reached a decision on their relative merit, this order will remain fixed.” Or: “A meeting participant who has close working relationships with the writer of a research proposal should leave the room when the merit of that proposal is discussed.” Slightly more general, but still at the intermediate level: How should the general rules for ranking research proposals be designed? E.g., collect at least three independent reviews per proposal, and use the reviews to get at a preliminary ranking. Ask participants to declare conflicts of interest before the meeting starts. And so on.

Strategizing at the macro-level How does the scientific community at large determine quality of research? How do the peer review system and the impact factor rankings of scientific journals influence the behaviour of researchers or research groups?

In other cases we can make similar distinctions. Take the case of voting as a decision making mechanism.

Micro-level At the micro level, individual voters decide what to do, given a particular voting rule, given their true preferences, and given what they know about the preferences of the other voters. The key question here is: “Should I vote according to my true preferences or not?” This is the question of *strategic voting*: deviating from one’s true preference in the hope of a better outcome. In a school class, during an election for a football captain, excellent sportsmen may cast a strategic vote on a mediocre player to further their own interests.

Intermediate level At the intermediate level, organizers of meetings decide which voting procedures to adopt in particular situations. How to fix the available set of alternatives? Does the situation call for secret ballots or not? How to settle the order for reaching decisions about sub-issues?

Macro-level At the macro-level, there is the issue of the design and analysis of voting procedures, or the improvement of voting procedures that fail to serve their political goal of rational collective decision making in a changing society. Think of the discussion of the merits of “first past the post” election systems in single-member districts, which favour the development of a small number of large parties, versus proportional representation systems which make it possible for smaller parties to survive in the legislature, but also engender the need for coalitions of several parties to aggregate in a working majority.

Writers about strategizing in warfare make similar level distinctions. Carl von Clausewitz, who defines war as “an act of violence or force intended to compel our enemy to do our will,” makes a famous distinction between tactics, the doctrine of the use of troops in battle, and strategy, the doctrine of the use of armed engagements to further the aims of the war [von Clausewitz, 1832–1834]. In our terminology, these issues are at the micro- and at the intermediate level, while the political choices between war and peace are being made at the macro-level.

3 The Prisoner’s Dilemma as an Exemplar

The game known as the “prisoner’s dilemma” is an evergreen of game theory because it is a top-level description of the plight of two people, or countries, who can either act trustfully or not, with the worst outcome that of being a sucker whose trust gets exploited by the other player.

One particular choice for a player in such a situation is called a strategy. Further on, we will discuss how the choice of strategies in this sense is affected by redesign of the scenario, thus shifting attention to possible strategies for the social mechanism designer, so to speak.

But first, here is a brief recap. Agents I and II are imprisoned, in different cells, and are faced with a choice between cooperating with each other or defecting. If the two cooperate they both benefit, but, unfortunately, it pays to defect if the other cooperates. If they both defect they both lose. This situation can be described in the following payoff matrix.

	II cooperates	II defects
I cooperates	3, 3	0, 4
I defects	4, 0	1, 1

This table displays a two person non-zero-sum game. The first member of each payoff pair gives I's payoff, the second member II's payoff. The table outcome 4, 0 indicates that if I defects while II cooperates, the payoff for I is 4 and the payoff for II is 0. This indicates that for I it is profitable to cheat if II stays honest. In fact it is more profitable for I to cheat than to stay honest in this case, for honesty gives him a payoff of only 3.

For the prison setting, read the two options as 'keep silent' or 'betray (by talking to the prison authorities)'. For an armament race version, read the two options as 'disarm' or 'arm'.

What matters for the payoffs is the preference order that is implied by the numbers in the payoff matrix. Abbreviate the strategy pair where I cooperates and II defects as (c, d) , and so on. Then the preference order for I can be given as $(d, c), (c, c), (d, d), (c, d)$, while the preference order for II swaps the elements of the pairs: $(c, d), (c, c), (d, d), (d, c)$. Replacing the payoffs by different numbers reflecting the same preference order does not change the nature of the game.

Suppose player I decides to follow a particular strategy. If player II has also made up her mind about what to do, this determines the outcome. Does player I get a better payoff if he changes his strategy, given that II sticks to hers? Player II can ask the same question. A situation where neither player can improve his outcome by deviating from his strategy while it is given that the other player sticks to hers is called a *Nash equilibrium*, after John Nash [Kuhn and Nasar, 2002].

Observe that the strategy pair (d, d) is a Nash equilibrium, and no other strategy pair is. This is what makes the situation of the game a dilemma, for the outcome (c, c) would have been better for both.

Not only is (d, d) a Nash equilibrium of the game, but it holds that (d, d) is the *only* Nash equilibrium. What follows is that for each player, d is the optimal action, *no matter what the other player does*. Such a strategy is called a *dominant* strategy.

Using u_I for I's utility function, and u_{II} for II's utility function, we can say that what makes the game into a dilemma is the fact that $u_I(d, d) > u_I(c, d)$ and $u_I(d, c) > u_I(c, c)$, and similarly for II: $u_{II}(d, d) > u_{II}(d, c)$ and $u_{II}(c, d) > u_{II}(c, c)$.

The two-player prisoner's dilemma can be generalized to an n player prisoner's dilemma (NPD), which can be used to model situations where the invisible hand does *not* work to the benefit of all. See Section 6 below.

Here we remind the reader of some formal terminology for strategic n -person games. A strategic game G is a tuple

$$(\{1, \dots, n\}, \{S_i\}_{i \in \{1, \dots, n\}}, \{u_i\}_{i \in \{1, \dots, n\}}),$$

where $\{1, \dots, n\}$ with $n > 1$ is the set of players, each S_i is a set of strategies, and each u_i is a function from $S_1 \times \dots \times S_n$ to \mathbb{R} (the utility function for player i). I use N for $\{1, \dots, n\}$, S for $S_1 \times \dots \times S_n$ and u for $\{u_i\}_{i \in \{1, \dots, n\}}$, so that I can use (N, S, u) to denote a game.

A member of $S_1 \times \dots \times S_n$ is called a strategy profile: each player i picks a strategy $s_i \in S_i$. I use s to range over strategy profiles, and s_{-i} for the strategy profile that results by deleting strategy choice s_i of player i from s . Let (s'_i, s_{-i}) be the strategy profile that is like s for all players except i , but has s_i replaced by s'_i . Let S_{-i} be the set of all strategy profiles minus the strategy for player i (the product of all strategy sets minus S_i). Note that $s_{-i} \in S_{-i}$. A strategy s_i is a *best response* in s if

$$\forall s'_i \in S_i \quad u_i(s) \geq u_i(s'_i, s_{-i}).$$

A strategy profile s is a (pure) Nash equilibrium if each s_i is a best response in s :

$$\forall i \in N \quad \forall s'_i \in S_i \quad u_i(s) \geq u_i(s'_i, s_{-i}).$$

Let $\text{nash}(G) = \{s \in S \mid s \text{ is a Nash equilibrium of } G\}$.

A game G is *Nash* if G has a (pure) Nash equilibrium.

A strategy $s^* \in S_i$ weakly dominates another strategy $s' \in S_i$ if

$$\forall s_{-i} \in S_{-i} \quad u_i(s^*, s_{-i}) \geq u_i(s', s_{-i}).$$

A strategy $s^* \in S_i$ strictly dominates another strategy $s' \in S_i$ if

$$\forall s_{-i} \in S_{-i} \quad u_i(s^*, s_{-i}) > u_i(s', s_{-i}).$$

If a two-player game has a strictly dominant strategy for each player, both players will play that strategy no matter what the other player does, and the dominant strategy pair will form the only Nash equilibrium of the game. This is what happens in the prisoner's dilemma game.

Define a *social welfare function* $W : S_1 \times \dots \times S_n \rightarrow \mathbb{R}$ by setting

$$W(s) = \sum_{i=1}^n u_i(s).$$

A strategy profile s of a game $G = (N, S, u)$ is a *social optimum* if

$$W(s) = \sup\{W(t) \mid t \in S\}.$$

For a finite game, s is a social optimum if $W(s)$ is the maximum of the welfare function for that game.

In the case of the prisoner's dilemma game, the social optimum is reached at (c, c) , with outcome $W(c, c) = 3 + 3 = 6$.

We can turn the prisoner's dilemma setting into a playground for social software engineering, in several ways. In Section 4 we explore punishment mechanisms, while in Section 5 we look at welfare redistribution.

4 Appropriate Punishment

Suppose the social software engineer is confronted with a PD situation, and has to design a policy that makes defection less profitable. One way of doing that is to put a penalty P on defection. Notice that we now talk about engineering a strategy at a level different from the level where I and II choose their strategies in the game.

A penalty P on cheating does not have an immediate effect, for it can only be imposed if the one who cheats gets caught. Suppose the probability of getting caught is p . In case the cheater gets caught, she gets the penalty, otherwise she gets what she would have got in the original game.

Then adopting the policy amounts to a *change in the utility functions*. In other words, the policy change can be viewed as a *game transformation* that maps strategic game G to strategic game G^{pP} , where G^{pP} is like G except for the fact that the utility function is replaced by:

$$\begin{aligned} u_I^{pP}(c, c) &= u_I(c, c), \\ u_I^{pP}(d, c) &= pP + (1 - p)u_I(d, c), \\ u_I^{pP}(c, d) &= u_I(c, d), \\ u_I^{pP}(d, d) &= pP + (1 - p)u_I(d, d), \end{aligned}$$

and similarly for u_{II}^{pP} . The utility for behaving honestly if the other player is also honest does not change. The new utility of cheating if the other is honest amounts to P in case you get caught, and to the old utility of cheating in case you can get away with it. The probability of getting caught is p , that of getting away unpunished is $1 - p$. Hence $u_I^{pP} = pP + (1 - p)u_I(d, c)$.

This allows us to compute the 'right' amount of punishment as a function of the utilities of being honest and of cheating without being caught, and the probability of being caught. Recall the assumption that the utility of staying honest while the other player cheats has not changed.

Call this H . Let C be the reward for cheating without being caught. Let p is the probability of getting caught. Then a punishment of $\frac{H+pC-C}{p}$ is “just right” for making cheating lose its appeal. Technically, this is the least amount of punishment that turns the social optimum of the game into a Nash equilibrium.

For example, suppose the probability of getting caught cheating is $\frac{1}{9}$. Then the punishment that ensures that cheating loses its appeal in case the other player is honest, for the utilities shown above, equals $\frac{3+(1/9)4-4}{1/9} = -5$. This amounts to the following transformation of the prisoner’s dilemma game:

$$\begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 3, 3 & 0, 4 \\ \hline d & 4, 0 & 1, 1 \\ \hline \end{array} \Rightarrow \left(-5, \frac{1}{9}\right) \Rightarrow \begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 3, 3 & 0, 3 \\ \hline d & 3, 0 & \frac{1}{3}, \frac{1}{3} \\ \hline \end{array}$$

The new game has two Nash equilibria: (c, c) with payoff $(3, 3)$, and (d, d) , with payoff $(\frac{1}{3}, \frac{1}{3})$. If the other player is honest, cheating loses its appeal, but if the other player cheats, cheating still pays off.

The punishment that ensures that cheating loses its appeal in case the other player is cheating (assuming the probability of getting caught is still the same) is higher. It equals $\frac{(1/9)-1}{1/9} = -8$. This corresponds to the following game transformation:

$$\begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 3, 3 & 0, 4 \\ \hline d & 4, 0 & 1, 1 \\ \hline \end{array} \Rightarrow \left(-8, \frac{1}{9}\right) \Rightarrow \begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 3, 3 & 0, 2\frac{2}{3} \\ \hline d & 2\frac{2}{3}, 0 & 0, 0 \\ \hline \end{array}$$

In the result of this new transformation, the social optimum (c, c) is the only Nash equilibrium.

There are many possible variations on this. One reviewer suggested that in the case where cheating gets detected, there should also be an implication for the honest player. The cheating player should get the penalty indeed, but maybe the honest player should get what she would get in case both players are honest.

Another perspective on this is that punishment presupposes an agent who administers it, and that doling out punishment has a certain cost. Think of real-life examples such as confronting a queue jumper in a supermarket line. The act of confrontation takes courage, but if it succeeds all people in the queue benefit [Fehr and Gächter, 2002].

Game theory does not directly model what goes on in society, but game-theoretical scenarios can be used to *illuminate* what goes on in society. The transformation mechanism for the prisoner’s dilemma scenario

illustrates, for example, why societies with widespread crime need more severe criminal laws than societies with less crime. Also, the calculations suggest that if a society wants to avoid severe punishments, it has to invest in measures that ensure a higher probability of getting caught.

A game-theoretical perspective on crime and punishment is in line with rational thinking about what constitutes ‘just punishment’, which goes back (at least) to Beccaria [1764]. What the analysis is still missing is the important principle that the punishment should somehow be in proportion to the severity of the crime. Such proportionality is important:

If an equal punishment be ordained for two crimes that injure society in different degrees, there is nothing to deter men from committing the greater as often as it is attended with greater advantage. [Beccaria, 1764, Ch 6]

Let us define a measure for social harm caused by the strategy of an individual player. Let a game $G = (N, S, u)$ be given. For any $i \in N$, define the individual harm function $H_i : S \rightarrow \mathbb{R}$, as follows:

$$H_i(s) = \sup_{s'_i \in S_i} W(s'_i, s_{-i}) - W(s).$$

This gives the difference between the best outcome for society as i unilaterally deviates from her current strategy and the present outcome for society. Assuming that the set S_i is finite, we can replace this by:

$$H_i(s) = \max_{s'_i \in S_i} W(s'_i, s_{-i}) - W(s).$$

That is, $H_i(s)$ gives a measure for how much player i harms society by playing s_i rather than the alternative s'_i that ensures the maximum social welfare. Clearly, in case s is a social optimum, $H_i(s) = 0$ for any i .

In the case of the prisoner’s dilemma, if player II is honest, the cheating behaviour of player I causes 2 units of societal harm:

$$H_I(c, c) = 0, H_I(d, c) = W(c, c) - W(d, c) = 6 - 4 = 2.$$

Also in case II cheats, the cheating behaviour of I causes 2 units of societal harm:

$$H_I(d, d) = H(c, d) - H(d, d) = 4 - 2 = 2.$$

Finally, $H_I(c, d) = 0$, for playing honest if the other player is cheating is better for society than cheating when the other player is cheating.

Punishment can now be made proportional to social harm, as follows. If $G = (N, S, u)$ is a strategic game, $p \in [0, 1]$, $\beta \in \mathbb{R}_{\geq 0}$, then $G^{p\beta}$ is the game $(N, S, u^{p\beta})$, where $u^{p\beta}$ is given by:

$$u_i^{p\beta}(s) := u_i(s) - p\beta H_i(s).$$

To see what this means, first consider cases of s and i with $H_i(s) = 0$. In these cases we get that $u_i^{p\beta}(s) = u_i(s)$. In cases where $H_i(s) > 0$ we get that the penalty for harming society is proportional to the harm. Observe that

$$u_i(s) - p\beta H_i(s) = (1 - p)u_i(s) + p(u_i(s) - \beta H_i(s)).$$

So, with probability $1 - p$ the crime gets undetected, and the player gets $u_i(s)$. With probability p , the crime gets detected, and the player gets $u_i(s) - \beta H_i(s)$, the original reward minus the penalty.

Now we have to find the least β that deters players from harming society. Games where no player has an incentive for harming society are the games that have a social optimum that is also a Nash equilibrium. For any game G it holds that $G^{0\beta} = G$, for all β , for if there is no possibility of detection, it does not matter what the penalty is. If the probability p of detection is non-zero, we can investigate the class of games $\{G^{p\beta} \mid \beta \in \mathbb{R}_{\geq 0}\}$.

Note that G and $G^{p\beta}$ have the same social optima, for in a social optimum s it holds for any player i that $H_i(s) = 0$. Moreover, if s is a social optimum of G , then $W(s) = W^{p\beta}(s)$.

As an example, consider the prisoner's dilemma again. We get:

$$u_I^{p\beta}(d, c) = u_I(d, c) - p\beta H_I(d, c) = 4 - 2p\beta.$$

To make (c, c) Nash, we need $4 - 2p\beta \leq 3$, whence $\beta \geq \frac{1}{2p}$ (recall that $p > 0$).

Nash equilibrium can be viewed as the outcome of the agents' strategic reasoning. It is the most commonly used notion of equilibrium in game theory, but that does not mean that this is the obviously right choice in any application. Here is one example of a modification. Call a strategy s_i a *social best response* in s if

$$\forall s'_i \in S_i (W(s) \leq W(s'_i, s_{-i}) \rightarrow u_i(s) \geq u_i(s'_i, s_{-i})).$$

What this means is that no other response for i from among the responses that do not harm the social welfare payoff is strictly better than the current response for i .

Call a strategy profile a *social equilibrium* if each s_i is a social best response in s :

$$\forall i \in \{1, \dots, n\} \forall s'_i \in S_i (W(s) \leq W(s'_i, s_{-i}) \rightarrow u_i(s) \geq u_i(s'_i, s_{-i})).$$

The PD game has two social equilibria: (c, c) and (d, d) . The strategy pair (c, c) is a social equilibrium because for each of the players, deviating from this profile harms the collective. The strategy pair (d, d) is a social equilibrium because it holds for each player that deviating from it harms that player.

A strategy profile is called Pareto optimal if it is impossible in that state to make a player better off without making at least one other player worse off. The profiles (c, c) , (c, d) and (d, c) in the PD game are Pareto optimal, while (d, d) is Pareto-dominated by (c, c) . This shows that Pareto optimality is different from being a social equilibrium.

If one would be allowed to assume that players are ‘social’ in the sense that they would always refrain from actions that harm society as a whole, then meeting out punishment proportional to the social harm that is caused would make no sense anymore, for players would not cause any social harm in the first place. In a more realistic setting, one would assume a certain mix of socially responsible and socially irresponsible players, and study what happens in repeated game playing for populations of such player types [Sigmund, 2010]. If the distinction between socially responsible players and selfish players makes sense, the distinction between Nash equilibria and social equilibria may be useful for an analysis of social responsibility. I must leave this for future work.

5 Welfare Redistribution

For another variation on the prisoner’s dilemma game, we can think of reward rather than punishment. The idea of using *welfare redistribution* to make a game more altruistic can be found in many places, and has made it to the textbooks. Consider the following exercise in Osborne [2004], where the student is invited to analyze a variation on the prisoner’s dilemma:

The players are not “selfish”; rather the preferences of each player i are represented by the payoff function $m_i(a) + \alpha m_j(a)$, where $m_i(a)$ is the amount of money received by player i when the action profile is a , j is the other player, and α is a given non-negative number.

[Osborne, 2004, exercise 27.1 on page 27]

This idea is worked out in Apt and Schaefer [2012] for the general case of n player strategic games, where the *selfishness level* of a game is computed by transforming a game G to a different game $G(\alpha)$, with α a positive real number, and $G(\alpha)$ the result of modifying the payoff function of G by adding $\alpha W(s)$ to each utility ($W(s)$ is the social welfare outcome for the strategy profile s).

As an example, using $\alpha = 1$, we can transform the prisoner's dilemma game PD into $PD(1)$, as follows:

$$\begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 3, 3 & 0, 4 \\ \hline d & 4, 0 & 1, 1 \\ \hline \end{array} \Rightarrow (\alpha = 1) \Rightarrow \begin{array}{|c|c|c|} \hline & c & d \\ \hline c & 9, 9 & 4, 8 \\ \hline d & 8, 4 & 3, 3 \\ \hline \end{array}$$

This gives a new game, and in this modified game, the only Nash equilibrium is (c, c) . This means that the social optimum now is a Nash equilibrium. The selfishness level of a game G is defined as the least α for which the move from G to $G(\alpha)$ yields a game for which a social optimum is a Nash equilibrium. For the prisoner's dilemma game with the payoffs as given in the table on page 6, the selfishness level α can be computed by equating the payoff in the social optimum with the best payoff in the Nash equilibrium: $3 + 6\alpha = 4 + 4\alpha$, which gives $\alpha = \frac{1}{2}$.

There are also games G that have at least one social optimum, but that cannot be turned into a game $G(\alpha)$ with a socially optimal Nash equilibrium for any α . Apt and Schaefer [2012] stipulate that such games have a selfishness level equal to ∞ .

Instead of the selfishness level, I will use a reformulation of this idea which consists in computing what is the least amount of *welfare redistribution* that is necessary to convert a social optimum into a Nash equilibrium. In other words: *how far* do you have to move on the scale from pure capitalism to pure communism to ensure that a social optimum is a Nash equilibrium? (But whether this is more perspicuous remains a matter of taste, for I have tried in vain to convince the authors to adjust their definition.)

The map for welfare redistribution is $G \mapsto G[\gamma]$, where $\gamma \in [0, 1]$ (our γ is a *proportion*), and the payoff u_i^γ in the new game $G[\gamma]$ is computed from the payoff u_i in G (assuming there are n players) by means of:

$$u_i^\gamma(s) = (1 - \gamma)u_i(s) + \gamma \frac{W(s)}{n}.$$

Here $W(s)$ gives the result of the welfare function on s in G .

Thus, player i is allowed to keep $1 - \gamma$ of her old revenue $u_i(s)$, and gets an equal share $\frac{1}{n}$ of $\gamma W(s)$, which is the part of the welfare that gets redistributed. This definition is mentioned (but not used) in Chen and Kempe [2008].

Notice the similarity to the probability of punishment computation on page 9. Also notice that if $\gamma = 0$, no redistribution of wealth takes place (pure capitalism), and if $\gamma = 1$, all wealth gets distributed equally (pure communism).

The *civilization cost* of a game G is the least γ for which the move from G to $G[\gamma]$ turns a social optimum into a Nash equilibrium. In case G has no social optimum, the civilization cost is undefined.

Note the difference with the notion of the selfishness level of a game, computed by means of $u_i^\alpha(s) = u_i(s) + \alpha W(s)$. Summing over all the players this gives a new social welfare $W' = (1 + n\alpha)W$. If we rescale by dividing all new payoffs by $1 + n\alpha$, we see that this uses a different recipe: $q_i(s) = \frac{u_i(s) + \alpha W}{1 + n\alpha}$. Thus, the definitions of selfishness level and civilisation cost are *not* related by rescaling (linear transformation). Rather, they are related, for the case where $\gamma \in [0, 1)$, by the nonlinear transformation $\alpha = \frac{\gamma}{n(1-\gamma)}$. This transformation is undefined for $\gamma = 1$. Note that the map $G, \gamma \mapsto G[\gamma]$ is more general than the map $G, \alpha \mapsto G(\alpha)$, for the game $G[1]$ where all welfare gets distributed equally has no counterpart $G(\cdot)$. Setting α equal to ∞ would result in a ‘game’ with infinite payoffs.

An example of a game for which the selfishness level and the civilization cost are 0 is the stag hunting game (first mentioned in the context of the establishment of social convention, in Lewis [1969], but the example goes back to Rousseau [1755]), with s for hunting stag and h for hunting hare.

	s	h
s	2, 2	0, 1
h	1, 0	1, 1

These payoffs are meant to reflect the fact that stag hunting is more rewarding than hunting hare, but one cannot hunt stag on one’s own.

Note the difference in payoffs with the prisoner’s dilemma game: if your strategy is to hunt hare on your own, it makes no difference for your payoff whether the others also hunt hare or not. This game has two Nash equilibria, one of which is also a social optimum. This is the strategy tuple where everyone joins the stag hunt. So the selfishness level and the civilisation cost of this game are 0.

Here is how the result of redistribution of proportion γ of the social welfare is computed for the PD game, for the case of I (the computation for II is similar):

$$\begin{aligned} u_I^\gamma(c, c) &= u_I(c, c), \\ u_I^\gamma(d, c) &= (1 - \gamma)u_I(d, c) + \gamma \frac{W(d, c)}{2} \\ u_I^\gamma(c, d) &= (1 - \gamma)u_I(c, d) + \gamma \frac{W(c, d)}{2} \\ u_I^\gamma(d, d) &= u_I(d, d). \end{aligned}$$

In the cases $u_I^\gamma(c, c)$ and $u_I^\gamma(d, d)$ nothing changes, for in these cases the payoffs for I and II were already equal.

The civilisation cost of the prisoner's dilemma, with the payoffs of the example, is computed by means of $3 = 4(1 - \gamma) + \frac{\gamma}{2}4$, which yields $\gamma = \frac{1}{2}$. This is the same value as that of the selfishness level, because substitution of $\frac{1}{2}$ for γ in the equation $\alpha = \frac{\gamma}{2-2\gamma}$ yields $\alpha = \frac{1}{2}$.

If we change the payoffs by setting $u_I(d, c) = u_{II}(c, d) = 5$, while leaving everything else unchanged, the cost of civilization is given by $3 = 5(1 - \alpha) + \frac{\alpha}{2}5$, which yields $\alpha = \frac{4}{5}$. The selfishness level in this case is given by $3 + 6\alpha = 5 + 5\alpha$, which yields $\alpha = 2$.

These were mere illustrations of how the effects of welfare redistribution can be studied in a game-theoretic setting. This analysis can help to understand social interaction in real life, provided of course that the game-theoretic model fits the situation. In the next section we will look at a more sophisticated model for the conflict between individual and societal interest than the prisoner's dilemma game model.

6 Tragedy-of-the-Commons Scenarios and Social Engineering

The tragedy of the commons game scenario that applies to games of competition for shares in a commonly owned resource was first analyzed in Gordon [1954] and was made famous in an essay by Garrett Hardin:

The tragedy of the commons develops in this way. Picture a pasture open to all. It is to be expected that each herdsman will try to keep as many cattle as possible on the commons. Such an arrangement may work reasonably satisfactorily for centuries because tribal wars, poaching, and disease keep the numbers of both

man and beast well below the carrying capacity of the land. Finally, however, comes the day of reckoning, that is, the day when the long-desired goal of social stability becomes a reality. At this point, the inherent logic of the commons remorselessly generates tragedy. [Hardin, 1968]

Bringing more and more goats to the pasture will in the end destroy the commodity for all. Still, from the perspective of an individual herdsman it is profitable until almost the very end to bring an extra goat.

The tragedy of the commons can be analyzed as a multi-agent version of the prisoner's dilemma. The players' optimal selfish strategies depend on what the other players do, and the outcome if all players pursue their individual interest is detrimental to the collective. One can also view this as a game of an individual herdsman I against the collective II. Then the matrix is:

	m	g
m	2, 2	0, 3
g	3, 0	-1, -1

Each player has a choice between g (adding goats) and m (being moderate). Assuming that the collective is well-behaved, it pays off to be a free rider. But if everyone acts like this, system breakdown will result.

In a more sophisticated multi-player version, assume there are n players. I use the modelling of Chapter 1 of Vazirani et al. [2007]. The players each want to have part of a shared resource. Setting the value of the resource to 1, each player i has to decide on the part of the resource x_i to use, so we can assume that $x_i \in [0, 1]$. Note that in this model, each player can choose from an infinite number of possible strategies.

Let us stipulate the following payoff function. Let N be the set of agents. If $\sum_{j \in N} x_j < 1$ then the value for player i is $u_i = x_i(1 - \sum_{j \in N} x_j)$: the benefit for i decreases as the resource gets exhausted. If $\sum_{j \in N} x_j \geq 1$ (the demands on the resource exceed the supply), the payoff for the players becomes 0.

So what are equilibrium strategies? Take the perspective of player i . Let D be the total demand of the other players, i.e., $D = \sum_{j \in N, j \neq i} x_j < 1$. Then strategy x_i gives payoff $u_i = x_i(1 - (D + x_i))$, so the optimal solution for i is $x_i = (1 - D)/2$. Since the optimal solution for each player is the same, this gives $x = \frac{1 - (n-1)x}{2}$, and thus $x = \frac{1}{n+1}$ as the optimal strategy for each player. This gives $D + x = \frac{n}{n+1}$, and payoff for x of $u = \frac{1}{n+1}(1 - \frac{n}{n+1}) = \frac{1}{(n+1)^2}$, and a total payoff of $\frac{n}{(n+1)^2}$, which is roughly

$\frac{1}{n}$. This means that the social welfare in the Nash equilibrium for this game depends inversely on the number of players.

If the players had agreed to leave the resource to a single player, the total payoff would have been $u = x(1 - x)$, which is optimal for $x = \frac{1}{2}$, yielding payoff $u = \frac{1}{4}$. If the players had agreed to use only $\frac{1}{2}$ of the resource, they would have had a payoff of $\frac{1}{4n}$ each, which is much more than $\frac{1}{(n+1)^2}$ for large n . Tragedy indeed.

Can we remedy this by changing the payoff function, transforming the ToC into ToC $[\gamma]$ with a Nash equilibrium which also is a social optimum? It turns out we can, but only at the cost of complete redistribution of welfare. The civilization cost of the ToC is 1. Here is why. If all players decide to leave the resource to a single player i , the payoff for i is given by $u_i = x_i(1 - x_i)$. This is optimal for $x_i = \frac{1}{2}$, and the payoff for this strategy, in the profile where all other players play 0, is $\frac{1}{4}$. This is the social optimum.

Suppose we are in a social optimum s . Then $W(s) = \frac{1}{4}$. Player i deviates by moving from x_i to $x_i + y$. The new payoff is $(x_i + y)(\frac{1}{2} - y) = \frac{1}{2}(x_i + y) - y(x_i + y)$. The deviation is tempting if $(x_i + y)(\frac{1}{2} - y) > \frac{1}{2}x_i$. Solving for y gives: $y < \frac{1}{2}$.

Let s' be the profile where i plays $x_i + y$. Then $W(s') = (\frac{1}{2} + y)(\frac{1}{2} - y) = \frac{1}{4} - y^2$, so $W(s) - W(s') = y^2$.

$$u_i(s') - u_i(s) = \frac{1}{2}(x_i + y) - y(x_i + y) - \frac{1}{4} = y(\frac{1}{2} - x_i - y).$$

We can now calculate just how much welfare we have to distribute for a given alternative to social optimum s to lose its appeal for i . A tempting alternative s' for i in s loses its appeal for i in s when the following holds:

$$u_i^\gamma(s') \leq u_i^\gamma(s).$$

Write out the definition of u_i^γ :

$$(1 - \gamma)u_i(s') + \gamma \frac{W(s')}{n} \leq (1 - \gamma)u_i(s) + \gamma \frac{W(s)}{n}.$$

Solve for γ :

$$\frac{n(u_i(s') - u_i(s))}{n(u_i(s') - u_i(s)) + W(s) - W(s')} \leq \gamma.$$

In our particular case, this gives:

$$\frac{ny(\frac{1}{2} - x_i - y)}{ny(\frac{1}{2} - x_i - y) + y^2} = \frac{nx_i + ny - n}{nx_i + ny - n - y^2}.$$

We have that $0 \leq x_i \leq \frac{1}{2}$, $0 \leq y < \frac{1}{2}$, Plugging these values in, we get:

$$\sup_{0 \leq x_i \leq \frac{1}{2}, 0 \leq y < \frac{1}{2}} \frac{nx_i + ny - n}{nx_i + ny - n - y^2} = 1.$$

Since the social optimum s was arbitrary, it follows that the cost of civilization for the tragedy of the commons game is 1. (This corresponds to selfishness level ∞ .)

Now for the key question: what does this all mean for policy making in ToC situations? One can ask what a responsible individual should do in a ToC situation to optimize social welfare. Let $D = \sum_{i \in N} x_i$, i.e., D is the total demand on the resource. Suppose j is a new player who wants to act responsibly. What should j do? If $D < \frac{1}{2}$, j should demand

$$x_j = \frac{1}{2} - D.$$

This will make the new demand equal to $\frac{1}{2}$, and the welfare equal to $D - D^2 = \frac{1}{4}$, which is the social optimum.

If $D = \frac{1}{2}$, any positive demand of j would harm the social welfare, so in this case j should put $x_j = 0$. An alternative would be to persuade the n other players to each drop their individual demands from $\frac{1}{2n}$ (on average) to $\frac{1}{2n+2}$. If this plea succeeds, j can also demand $\frac{1}{2n+2}$, and the new total demand becomes $\frac{n+1}{2n+2} = \frac{1}{2}$, so that again the social optimum of $\frac{1}{4}$ is reached.

If $D > \frac{1}{2}$, any positive demand of j would harm the social welfare, so again j should put $x_j = 0$. In this case, the prospect of persuading the other players to lower their demands may be brighter, provided the players agree that they all have equal rights. Once this is settled, it is clear what the individual demands should be for optimum welfare. The optimum individual demand is $\frac{1}{2n}$ if there are n players, and $\frac{1}{2n+2}$ if there are $n+1$ players. Allowing in one extra player would cost each player $\frac{1}{4n} - \frac{1}{4n+4}$.

To change to the punishment perspective, suppose D is the demand in the old situation s . A new player comes in and demands x . Call the new situation s' . Let D be the total demand in s . Then $W(s) = D - D^2$. If $D + x > 1$ then $W(s') = 0$. So in this case, the social damage equals the original welfare, and the appropriate punishment is $-W(s)$.

In the case where $x + D \leq 1$, the excess demand is anything in excess of $\frac{1}{2}$, so the appropriate punishment is the welfare deterioration caused by the excess demand y . Thus, the appropriate punishment is given by:

$$\frac{1}{4} - W(s') = \frac{1}{4} - \left(\frac{1}{2} + y\right)\left(\frac{1}{2} - y\right) = y^2.$$

If this is combined with the probability p of catching offenders, the penalty for excess demand y should be $\frac{y^2}{p}$.

Take an example case. Two players each demand $\frac{1}{5}$, so each gets $\frac{1}{5}(1 - \frac{2}{5}) = \frac{3}{25}$. We have $D = \frac{2}{5}$, and $W = D - D^2 = \frac{6}{25}$. A third player comes along and demands $\frac{1}{3}$. Then the new demand D' becomes $\frac{11}{15}$, which results in new welfare $W' = D' - D'^2 = \frac{44}{225}$. The welfare in the social optimum is $\frac{1}{4}$. The excess demand is $(\frac{2}{5} + \frac{1}{3}) - \frac{1}{2} = \frac{7}{30}$. The deterioration in welfare is $\frac{1}{4} - W' = \frac{1}{4} - \frac{44}{225} = \frac{49}{900}$. This is exactly equal to the square of the excess demand $\frac{7}{30}$.

A modern and pressing case of the tragedy of the commons is presented in the Fourth IPCC Assessment report:

The climate system tends to be overused (excessive GHG concentrations) because of its natural availability as a resource whose access is open to all free of charge. In contrast, climate protection tends to be underprovided. In general, the benefits of avoided climate change are spatially indivisible, freely available to all (non-excludability), irrespective of whether one is contributing to the regime costs or not. As regime benefits by one individual (nation) do not diminish their availability to others (non-rivalry), it is difficult to enforce binding commitments on the use of the climate system [Kaul et al., 1999, 2003]. This may result in “free riding”, a situation in which mitigation costs are borne by some individuals (nations) while others (the “free riders”) succeed in evading them but still enjoy the benefits of the mitigation commitments of the former. [Rogner et al., 2007, page 102]

The problem of collective rationality has been a key issue in practical philosophy for more than two millennia. Aristotle discusses it at length, in the *Politics*:

For that which is common to the greatest number has the least care bestowed upon it. Every one thinks chiefly of his own, hardly at all of the common interest; and only when he is himself concerned as an individual. For besides other considerations, everybody is more inclined to neglect the duty which he expects another to fulfill; as in families many attendants are often less useful than a few. [Aristotle, (330 BC, paragraph 403, Book II)]

What is important about the game-theoretical analysis is the insight that there are situations where lots of individual actions of enlightened

self-interest may endanger the common good. There is not always an invisible hand to ensure a happy outcome.

The phenomenon that Aristotle alludes to is called the ‘bystander effect’ in Darley and Letane [1968]: solitary people usually intervene in case of an emergency, whereas a large group of bystanders may fail to intervene — everyone thinks that someone else is bound to have called the emergency hotline already (pluralistic ignorance), or that someone else is bound to be more qualified to give medical help (diffused responsibility). See Osborne [2004] for an account of this social phenomenon in terms of game theory, Pacuit et al. [2006] for a logical analysis, and Manning et al. [2007] for historical nuance about the often quoted and much discussed case of Kitty Genovese (who, according to the story, was stabbed to death in 1964 while 38 neighbours watched from their windows but did nothing).

Garrett Hardin, in his famous essay, also discusses how tragedy of the commons situations can be resolved. He makes a plea for the collective (or perhaps: enlightened individuals within the collective) to impose “mutual constraints, mutually agreed upon,” and he quotes Sigmund Freud’s *Civilisation and its Discontents* [Freud, 1930] to put the unavoidable tension between civilisation and the desires or inclinations of individuals in perspective.

On the other hand, Ostrom [1990] warns against the temptation to get carried away by the game-theoretical analysis of ToC situations, and shows by careful study of real-world cases of institutions (fisheries, irrigation water allocation schemes) how — given appropriate circumstances — effective collective action can be organized for governing common pool resources without resorting to a central authority. See Baden and Noonan [1998] for further discussion.

7 Renunciation Games

Let me depart now from the standard game theory textbook fare, and introduce three new games where an individual is pitted against a collective. The setup of the games is such that the social optimum of the game can only be reached at the expense of *one single* individual. I call such games renunciation games. When will an individual sacrifice his or her own interest to save society? It turns out that the nature of the renunciation game changes crucially depending on the temptation offered to the renouncer.

Pure Renunciation Game The *pure renunciation game* has n players, who each choose a strategy in $[0, 1]$, which represents their demand. If

at least one player renounces (demands 0), then all other players get as payoff what they demand. Otherwise, nobody gets anything. The payoff function for i is given by:

$$u_i(s) = \begin{cases} s_i & \text{if } \exists j \neq i : s_j = 0 \\ 0 & \text{otherwise.} \end{cases}$$

This game has n social optima $(0, 1, \dots, 1), (1, 0, 1, \dots, 1), \dots, (1, \dots, 1, 0)$, where the social welfare W equals $n - 1$. The social optima are also Nash equilibria. No need for welfare redistribution, no need for punishment. The situation changes if there is a temptation for the renouncer in the game.

Renunciation Game With Mild Temptation This renunciation game has n players, who each choose a strategy in $[0, 1]$, which represents their demand. If at least one player renounces (demands 0), then all other players get as payoff what they demand. Otherwise, if there is one player i who demands less than any other player, i gets what she demands, and the others get nothing. In all other cases nobody gets anything. The payoff function for i is given by:

$$u_i(s) = \begin{cases} s_i & \text{if } \exists j \neq i : s_j = 0 \\ & \text{or } \forall j \neq i : 0 < s_i < s_j \\ 0 & \text{otherwise.} \end{cases}$$

This game has n social optima. There are no Nash equilibria. The cost of civilization for the Renunciation Game is $\gamma = \frac{1}{2^{n-2}}$. Indeed, this game has n social optima $(0, 1, \dots, 1), (1, 0, 1, \dots, 1), \dots, (1, \dots, 1, 0)$, where the social welfare W equals $n - 1$. In particular, the social optima are not Nash equilibria. For in a social optimum, the player who renounces (and receives nothing) can get any q with $0 < q < 1$ by playing q . That's the temptation.

Now focus on player 1 and compute the least γ for which the social optimum $(0, 1, \dots, 1)$ turns into a Nash equilibrium in $G[\gamma]$. The payoff function for player 1 in $G[\gamma]$ satisfies:

$$u_1^\gamma(0, 1, \dots, 1) = \gamma \frac{n-1}{n}.$$

For the social optimum to be Nash, this value has to majorize

$$u_1^\gamma(q, 1, \dots, 1) = (1 - \gamma)q + \frac{\gamma}{n}q.$$

Since q can be arbitrarily close to 1, we get $u_1^\gamma(q, 1, \dots, 1) < (1 - \gamma) + \frac{\gamma}{n}$. Therefore $(0, 1, \dots, 1)$ is a social optimum in $G[\gamma]$ iff $\gamma \frac{n-1}{n} \geq (1 - \gamma) + \frac{\gamma}{n}$. Solving this for γ gives $\gamma \geq \frac{1}{2n-2}$.

The situation changes drastically if there is heavy temptation.

Renunciation Game With Heavy Temptation This renunciation game has n players, who each choose a strategy q in $[0, 1]$, which represents their demand. If at least one player renounces (demands 0), then all other players get as payoff what they demand. Otherwise, if there is one player i who demands less than any other player, i gets $n - 1$ times what she demands, and the others get nothing. In all other cases nobody gets anything. The payoff function for i is given by:

$$u_i(s) = \begin{cases} s_i & \text{if } \exists j \neq i : s_j = 0 \\ (n-1)s_i & \text{if } \forall j \neq i : 0 < s_i < s_j \\ 0 & \text{otherwise.} \end{cases}$$

The civilization cost for Renunciation With Heavy Temptation is 1. Social optima are the same as before. We have to compute the least γ that turns social optimum $(0, 1, \dots, 1)$ into a Nash equilibrium in $G[\gamma]$. The constraint on the payoff function for player 1 is:

$$u_1^\gamma(q, 1, \dots, 1) = (1 - \gamma)(n - 1)q + \frac{\gamma}{n}(n - 1)q.$$

Since q can be arbitrarily close to 1, this gives

$$u_1^\gamma(q, 1, \dots, 1) < (1 - \gamma)(n - 1) + \frac{\gamma}{n}(n - 1).$$

This puts the following constraint on γ :

$$\gamma \frac{n-1}{n} \geq (1 - \gamma)(n - 1) + \frac{\gamma}{n}(n - 1).$$

Solving for γ gives $n\gamma \geq n$, and it follows that $\gamma = 1$.

These games are offered here as examples of new metaphors for social interaction, showing that the store-room of game-theoretic metaphors is far from exhausted. I hope to analyse renunciation games in future work.

8 Experiments with Knowledge and Trust

In many social protocols (scenarios for social interaction) the knowledge that the participants have about each other and about the protocol itself play a crucial role.

The prisoner's dilemma scenario, e.g., assumes that there is common knowledge among the players about the utilities. Also, it is assumed that there is common knowledge that the players cannot find out what the other player is going to do. If we change the scenario, by letting the players move one by one, or by communicating the move of the first player to the second player, this changes the nature of the game completely.

Suppose two players meet up with a host, who hands over a bill of ten euros to each of them, and then explains that they will each be asked whether they are willing to donate some or all of the money to the other player. The host adds the information that donated amounts of money will be doubled.

What will happen now depends on the set-up. If each player communicates in private to the host, we are back with the prisoner's dilemma situation. If the players are allowed to coordinate their strategies, and if they act under mutual trust, they will each donate all of their money to the other player, so that they each end up with 20 euros. If the first player is asked in public what she will do, it depends on what she believes the other player will do if it is his turn, and so on.

Experiments based on this kind of scenario have been staged by game theorists, to explain the emergence of trust in social situations. A relevant game is the so-called ultimatum game, first used in Güth et al. [1982].

Player I is shown a substantial amount of money, say 100 euros. He is asked to propose a split of the money between himself and player II. If player II accepts the deal, they both keep their shares, otherwise they both receive nothing. If this game is played once, a split (99, 1) should be acceptable for II. After all, receiving 1 euro is better than receiving nothing. But this is not what we observe when this game is played. What we see is that II rejects the deal, often with great indignation [Camerer, 2003].

Evidence from experiments with playing the ultimatum game and repeated prisoner's dilemma games suggests that people are willing to punish those who misbehave, even if this involves personal cost.

Another game that was used in actual experiments, the investment game, suggests that people are also willing to reward appropriate behaviour, even if there is no personal benefit in giving the reward.

The investment game is played between a group of people in a room A and another group of people in a room B, and can be summarized as follows. Each person in room A and each person in room B has been given 10 euros as show up money. The people in room A will have the opportunity to send some or all of their money to an anonymous receiver

in group B. The amount of money sent will be tripled, and this is common knowledge. E.g., an envelope sent with 9 euros will contain 27 euros when it reaches its recipient in room B. The recipient in group B, who knows that someone in group A parted with one third of the amount of money she just received, will then decide how much of the money to keep and how much to send back to the giver in room A. Consult Berg et al. [1995] for the results of the experiment.

Reputation systems such as those used in Ebay are examples of engineered social software. The design aim of these public ratings of past behaviour is to make sure that trust can emerge between players that exchange goods or services. Reputation can be computed: Kleinberg [1999] gives a now-famous algorithm which ranks pages on the internet for authoritativeness in answering informative questions. One of the ways to strategically misuse reputation systems is by creating so-called “sybils”: fake identities which falsely raise the reputation of an item by means of fake links. So a design aim can be to create reputation mechanisms that are *sybil proof*; see Cheng and Friedman [2005]. For further general information on reputation systems, consult Resnick et al. [2000].

These systems can also be studied empirically: how does the designed reputation system influence social behaviour? The same holds for the renunciation game scenarios from the previous section. Empirical studies using these scenarios might yield some revealing answers to the question “What do people actually do when being asked to renounce for the benefit of society?”

9 Conclusion

Several chapters in this book present relevant logics for strategic reasoning. Van Benthem [2012] makes a plea for applying the general perspective of action logic to reasoning about strategies in games. In Van Eijck [2013] it is demonstrated how propositional dynamic logic or PDL [Pratt, 1976, Kozen and Parikh, 1981] can be turned into a logic for reasoning about finite strategic games. Such logics can be used to study, e.g., voting rules or auction protocols from a logical point of view. In voting, think of casting an individual vote as a strategy. Now fix a voting rule and determine a payoff function, and you have an n player voting game. Next, represent and analyze this in PDL, or in any of the logic formalisms taken from this book.

Voting is a form of collective decision making. A key distinction in decision making is between cases where there is a correct outcome, and

the challenge for the collective is to find that outcome, and cases where the notion of correctness does not apply, and the challenge for the collective is to arrive at a choice that everyone can live with.

A famous result from the early days of voting theory is Condorcet's jury theorem [Condorcet, 1785]. The case of a jury that has to reach a collective decision 'guilty or not', say in a murder trial, has a correct answer. For either the accused has committed the murder, or he has not. The trouble is that no member of the jury knows for sure what the answer is. Condorcet's jury theorem states the following:

Suppose each voter has an independent probability p of arriving at the correct answer. If p is greater than $\frac{1}{2}$ then adding more voters increases the probability of a correct majority decision. If p is smaller than $\frac{1}{2}$ then it is the other way around, and an optimal jury consists of a single voter.

To see why this is true, assume there are n voters. For simplicity, we assume n is odd. Assume that m voters have made the correct decision. Consider what happens when we add two new voters. Then the majority vote outcome changes in only two cases.

1. m was one vote short to get a majority of the n votes, and both new voters voted correctly. In this case the vote outcome changes from incorrect to correct.
2. m was just equal to a majority of the n votes, but both new voters voted incorrectly. In this case the vote outcome changes from correct to incorrect.

In both of these cases we can assume that it is the last of the n voters who casts the deciding vote. In the first case, voter n voted correctly, in the second case voter n voted incorrectly. But we know that voter n has probability p of arriving at a correct decision, so we know that in case there is just a difference of a single vote between the correct and the incorrect decision among n voters, the probability of the n voters arriving at a correct decision is p . Now add the two new voters. The probability of case (1), from incorrect to correct, is $(1 - p)p^2$, and the probability of case (2), from correct to incorrect, is $p(1 - p)^2$. Observe that $(1 - p)p^2 > p(1 - p)^2$ iff $p > \frac{1}{2}$. The case where there is an even number of voters is similar, but in this case we have to assume that ties are broken by a fair coin flip, with probability equal to $\frac{1}{2}$ of arriving at the correct decision.

Condorcet's jury theorem is taken by some as an argument for democracy; whether the argument cuts wood depends of course on whether one

believes in the notion of ‘correct societal decisions’. See List and Goodin [2001] for further discussion.

Let me finish, light-heartedly, with another famous argument for democracy, by Sir Francis Galton, in an amusing short paper ‘Vox Populi’ in *Nature*. Galton’s narrative is one of the key story lines in Surowiecki [2004]. Galton [1907] starts as follows:

In these democratic days, any investigation into the trustworthiness and peculiarities of popular judgments is of interest. The material about to be discussed refers to a small matter, but is much to the point.

A weight-judging competition was carried on at the annual show of the West of England Fat Stock and Poultry Exhibition recently held at Plymouth (England). A fat ox having been selected, competitors bought stamped and numbered cards, for 6d. each, on which to inscribe their respective names, addresses, and estimates of what the ox would weigh after it had been slaughtered and “dressed.” Those who guessed most successfully received prizes. About 800 tickets were issued, which were kindly lent me for examination after they had fulfilled their immediate purpose.

Galton then goes on to tell what he found. As it turned out, 13 tickets were defective or illegible, but the median of the 787 remaining ones contained the remarkably accurate guess of 1207 pounds, which was only 9 pounds above the actual weight of the slaughtered ox: 1198 pounds. The majority plus one rule gave the approximately correct answer.

What does this have to do with strategies and strategic reasoning, the reader might ask. The strategic reasoning is lifted to the meta-level now: Are we in a decision-making situation that is like weight-judging, or are we not? Is this a social situation where many know more than one, or isn’t it? Does the optimal jury for this consist of a single person, or does it not? Which brings us to the key strategic question we all face when about to make the decisions in life that really matter: “Should I take this decision on my own, or is it better to consult others before making my move?”

Acknowledgement Johan van Benthem, Robin Clark and Rainer Kessler sent their written comments on an early draft, in which Floor Sietsma was also involved. Later, I received helpful comments from Barteld Kooi.

The final version has benefitted from extensive reports by two anonymous reviewers. Since these reports were excellent, I have tried to imple-

ment almost all referee suggestions. Inspiring conversations with Mamoru Kaneko have led to further improvements.

Thanks are also due to book editors Rineke Verbrugge and Sujata Ghosh for detailed comments on the final version and for help with proof-reading, and to the other participants in the Lorentz workshop on Modeling Strategic Reasoning for inspiring feedback. Finally, I acknowledge communication with Krzysztof Apt and Guido Schaefer.

Bibliography

- K. R. Apt and G. Schaefer. Selfishness level of strategic games. In M. Serna, editor, *Algorithmic Game Theory*, volume 2012 of *Lecture Notes in Computer Science*, pages 13–24. Springer, 2012.
- D. Ariely, U. Gneezy, G. Loewenstein, and N. Mazar. Large stakes and big mistakes. *Review of Economic Studies*, pages 451–469, 2009.
- Aristotle. *The Politics of Aristotle: Translated into English with Introduction, Marginal Analysis, Essays, Notes and Indices*, volume 1. Clarendon Press, Oxford, (330 BC). Translated and annotated by B. Jowett (1885).
- R. Aumann. Agreeing to disagree. *Annals of Statistics*, 4(6):1236–1239, 1976.
- J. A. Baden and D. S. Noonan, editors. *Managing the Commons – 2nd edition*. Indiana University Press, 1998.
- C. Beccaria. *On Crimes and Punishment*. Marco Coltellini, 1764. In Italian: Dei delitti e delle pene.
- G. S. Becker. Nobel lecture: The economic way of looking at behavior. *Journal of Political Economy*, 101(3):385–409, 1993.
- J. v. Benthem. In praise of strategies. In J. v. Eijck and R. Verbrugge, editors, *Games, Actions, and Social Software*, volume 7010 of *Texts in Logic and Games, LNAI*, pages 105–125. Springer Verlag, Berlin, 2012.
- J. Berg, J. Dickhaut, and K. McCabe. Trust, reciprocity, and social history. *Games and Economic Behavior*, 10(1):122–142, 1995.
- S. Brams. *Mathematics and Democracy: Designing Better Voting and Fair Division Procedures*. Princeton University Press, 2008.
- C. F. Camerer. *Behavioral Game Theory: Experiments in Strategic Interaction*. Princeton University Press, 2003.
- P.-A. Chen and D. Kempe. Altruism, selfishness, and spite in traffic routing. In *Proceedings 10th ACM Conference on Electronic Commerce*, pages 140–149, 2008.
- A. Cheng and E. Friedman. Sybilproof reputation mechanisms. In *Proceedings of the 2005 ACM SIGCOMM workshop on Economics of peer-to-peer systems*, P2PECON '05, pages 128–132, New York, NY, USA, 2005. ACM.
- M. S.-Y. Chwe. *Rational Ritual*. Princeton University Press, Princeton and Oxford, 2001.

- M. Condorcet. *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*. Imprimerie Royale, Paris, 1785.
- J. Darley and B. Latane. Bystander intervention in emergencies: Diffusion of responsibility. *Journal of Personality and Social Psychology*, 8:377–383, 1968.
- J. v. Eijck. PDL as a multi-agent strategy logic. In B. C. Schipper, editor, *TARK 2013 – Theoretical Aspects of Reasoning About Knowledge, Proceedings of the 14th Conference – Chennai, India*, pages 206–215, 2013.
- J. v. Eijck and R. Verbrugge, editors. *Discourses on Social Software*, volume 5 of *Texts in Logic and Games*. Amsterdam University Press, Amsterdam, 2009.
- E. Fehr and S. Gächter. Altruistic punishment in humans. *Nature*, 415:137–140, 2002.
- S. Freud. *Das Unbehagen in der Kultur (Civilization and Its Discontents)*. Internationaler Psychoanalytischer Verlag, Wien (Vienna), 1930.
- F. Galton. Vox populi. *Nature*, pages 450–451, March 1907.
- I. Gilboa. *Rational Choice*. MIT Press, Cambridge, Massachusetts, 2010.
- H. S. Gordon. The economic theory of a common-property resource: The fishery. *Journal of Political Economy*, 62:124–142, 1954.
- W. Güth, R. Schmittberger, and B. Schwarze. An experimental analysis of ultimatum bargaining. *Journal of Economic Behavior and Organization*, 3(4):367–388, 1982.
- G. Hardin. The tragedy of the commons. *Science*, 162:1243–48, 1968.
- L. Hurwicz and S. Reiter. *Designing Economic Mechanisms*. Cambridge University Press, 2006. ISBN 9780521836418.
- I. Kaul, I. Grunberg, and M. Stern. *Global Public Goods*. Oxford University Press, 1999.
- I. Kaul, P. Conceicao, K. L. Gouven, and R. Mendoz. *Providing Global Public Goods*. Oxford University Press, 2003.
- J. M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, Sept. 1999. ISSN 0004-5411.
- T. Körner. *Naive Decision Making: Mathematics Applied to the Social World*. Cambridge University Press, 2008.
- D. Kozen and R. Parikh. An elementary proof of the completeness of PDL. *Theoretical Computer Science*, 14:113–118, 1981.
- V. Krishna. *Auction Theory*. Elsevier Science, 2009.
- H. W. Kuhn and S. Nasar, editors. *The Essential John Nash*. Princeton University Press, Princeton and Oxford, 2002.

- D. Lewis. *Convention: A Philosophical Study*. Harvard University Press, 1969.
- C. List and R. E. Goodin. Epistemic democracy: Generalizing the Condorcet jury theorem. *Journal of Political Philosophy*, 9(3):277–306, 2001.
- R. Manning, M. Levine, and A. Collins. The Kitty Genovese murder and the social psychology of helping: The parable of the 38 witnesses. *American Psychologist*, 62:555–562, 2007.
- Sun Tzu. *The Art of War, translated and with an introduction by Samuel B. Griffith*. Oxford University Press, 450BC. Translation from 1963.
- P. Milgrom. *Putting Auction Theory to Work*. Churchill Lectures in Economics. Cambridge University Press, 2004.
- M. J. Osborne. *An Introduction to Game Theory*. Oxford University Press, New York, Oxford, 2004.
- E. Ostrom. *Governing the Commons: The Evolution of Institutions for Collective Action*. Political Economy of Institutions and Decisions. Cambridge University Press, 1990.
- E. Pacuit, R. Parikh, and E. Cogan. The logic of knowledge based obligation. *Synthese*, 31:311–341, 2006.
- R. Parikh. Social software. *Synthese*, 132:187–211, 2002.
- V. Pratt. Semantical considerations on Floyd–Hoare logic. *Proceedings 17th IEEE Symposium on Foundations of Computer Science*, pages 109–121, 1976.
- P. Resnick, K. Kuwabara, R. Zeckhauser, and E. Friedman. Reputation systems. *Commun. ACM*, 43(12):45–48, Dec. 2000.
- H.-H. Rogner, R. Zhou, R. Bradley, P. Crabbé, O. Edenhofer, B. Hare, L. Kuijpers, and M. Yamaguchi. Introduction. In B. M. et al., editor, *Climate Change 2007: Mitigation. Contribution of Working Group III to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change*. Cambridge University Press, Cambridge, 2007.
- J. J. Rousseau. *Discours sur l'origine et les fondements de l'inégalité parmi les hommes*. Marc Michel Rey, Amsterdam, 1755.
- Y. Shoham and K. Leyton-Brown. *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press, 2008.
- K. Sigmund. *The Calculus of Selfishness*. Princeton Series in Theoretical and Computational Biology. Princeton University Press, Princeton and Oxford, 2010.
- A. Smith. *An Inquiry into the Nature and Causes of the Wealth of Nations*. Liberty Fund, Indianapolis, (1776). This edition: 1982.

- J. Surowiecki. *The Wisdom of Crowds: Why the Many Are Smarter Than the Few and How Collective Wisdom Shapes Business, Economies, Societies and Nation*. Random House, 2004.
- A. D. Taylor. *Social Choice and the Mathematics of Manipulation*. Cambridge University Press, 2005.
- E. Tenner. *Why Things Bite Back — Technology and the Revenge Effect*. Fourth Estate, 1996.
- R. Thaler and C. Sunstein. *Nudge: Improving Decisions About Health, Wealth, and Happiness*. A Caravan book. Yale University Press, 2008.
- V. V. Vazirani, N. Nisan, T. Roughgarden, and E. Tardos. *Algorithmic Game Theory*. Cambridge University Press, 2007.
- W. Vickrey. Counterspeculation, auctions, and competitive sealed tenders. *The Journal of Finance*, 16(1):8–37, 1961.
- M. von Clausewitz, editor. *Vom Kriege, Hinterlassenes Werk des Generals Carl von Clausewitz*. Ferdinand Dümmler, Berlin, 1832–1834.