C. L. SHENG

A NOTE ON THE PRISONER'S DILEMMA

ABSTRACT. This paper clarifies some basic concepts or assumptions of the prisoner's dilemma, asserts the independence between the two agents A and B, and advocates the application of the dominance principle of decision theory to the prisoner's dilemma. It discusses several versions of the prisoner's dilemma, including the one-shot and repeated cases of a noncooperative game from a purely egoistic point of view. The main part of this paper, however, is a study of the problem from a moral point of view through a special decision-theoretic approach. Morality is taken into account by incorporating the utility of the feeling of moral satisfaction for the agent, as a part of the total utility for the agent, into the decision-theoretic model. In this way the problem will appear as a purely technical decision problem, and the conflicts between various assumptions, or the dilemma caused by the problem, will no longer exist. It is also pointed out that in a more general case, for some values of the coefficient of morality k, dominance will not exist so that the dominance principle will not be applicable.

Keywords: Prisoner's dilemma, dominance principle, maximization of utility, social utility, feeling of moral satisfaction, coefficient of morality.

1. INTRODUCTION

In this paper I propose a new approach to the well-known prisoner's dilemma – an explanation of the problem in terms of the feeling of moral satisfaction from the moral point of view. Recently Randall K. Campbell (1989) discussed this dilemma, compared Lawrence Davis's (1977) 'symmetry argument for cooperation' (called 'basic argument for cooperation' by Davis himself) with the familiar dominance argument, and concluded that the symmetry argument for cooperation seems to fail. As pointed out by Richmond Campbell (1985), this is in fact a main difficulty with the study of the prisoner's dilemma, and opinions are very controversial.¹

The prisoner's dilemma symbolizes the basic justification for an ethical theory, especially utilitarianism, because it ties up the maximization of utility with morality or concerns the gap between 'is' and 'ought'. Thus the problem of the prisoner's dilemma has two dimensions for consideration: one is the maximization of self-interest or personal utility. This is a nonmoral problem of pure rationality and decision. A second dimension is the maximization of aggregate or social utility, or the common interest of the two persons A and B involved in the problem of prisoner's dilemma together as a society. This is related to the principle of utility and symbolizes the essence of morality from the societal point of view.

I study the prisoner's dilemma mainly from the second dimension, namely a moral point of view, by considering both the interest of the prisoner who is to make a decision her/himself and also the interest of the other prisoner.

Before this main part, in Section 2 I first discuss the various versions of the prisoner's dilemma as a game. As a pure game, morality of the interest of the other prisoner need not be taken into account. That is, the problem is studied from a purely egoistic point of view by considering the interest of the agent her/himself alone. Then there are again two different versions, namely the one-shot case and the dynamic or repeated case. For the one-shot case I hold that the dominance principle is applicable, even though a decision according to the dominance principle does not necessarily lead to maximal utility. For the repeated case I proceed according to game theory, namely to find the equilibrium point or probability distribution for maximal utility.

In Section 3 I study the prisoner's dilemma from a moral point of view, using a special approach of my own. I take morality into account by incorporating the utility of the feeling of moral satisfaction for the agent, into the decision-theoretic model. In this way the problem will appear as a purely technical decision problem, and the conflicts between various assumptions, or the dilemma caused by the problem, will no longer exist.

Finally, in Section 4 I study a more general case where the entries in cell (I, I) are smaller than (9, 9). Then for some values of k dominance does not exist so that the dominance principle will not be applicable.

2. THE VARIOUS VERSIONS OF THE PRISONER'S DILEMMA

Theoretically, the prisoner's dilemma has many versions or variations. It is not necessary, nor practical, to study all the variations. My study of it in this paper is restricted to a specific version of it. However, it seems to be in order to have a quick look at these versions.

First, the prisoner's dilemma is usually studied as an application of game theory, and games may be classified into *cooperative* ones and *noncooperative* ones. Since the two prisoners are supposed to be jailed separately and no communication between them is permitted, the prisoner's dilemma is obviously a noncooperative game.

Second, although the prisoner's dilemma is a noncooperative game, each prisoner can solve the problem either in terms of her/his own interest only, i.e. from a purely egoistic point of view, or in terms of her/his own interest and the interest of the other prisoner as well, i.e., from a moral point of view. Normally game theory is restricted to the egoistic point of view only. However, since the prisoner's dilemma is studied as a problem of moral philosophy, each prisoner is supposed to consider the interest of other people as well. (In a prisoner's dilemma the other people are restricted to the other prisoner.) This is exactly my emphasis in this paper. Before the presentation of my approach to the problem, it is in order to discuss briefly my view of the prisoner's dilemma as a pure game, i.e. from a purely egoistic point of view by considering the interest of the prisoner who is to make the decision her/himself alone. Then there are again two subcases: one being a one-shot case and the other a dynamic or repeated case.

Normally the prisoner's dilemma should be regarded as a one-shot case, because it is not likely or practical to have two prisoners jailed and released for a large number of times.

For the one-shot case, I hold that the prisoners A and B are both independent of each other in making decisions and that the dominance principle is still applicable, even if it will not necessarily result in maximal utility. My view in this respect is different from the usual assumption that A and B are interdependent to the extent that if Achooses alternative I then B will choose alternative I too and if Achooses alternative II then B will choose alternative II too. A and B, as two automonous persons having free will, should be regarded as independent of each other so far as decision-making is concerned. The identity or closeness of their choices stems purely from the fact that they are symmetrically situated and that they are both rational. For A, B's two alternatives of choice, I and II, represent the states of the

world, which should be completely independent of A and over which A has no control at all. Similarly, for B, A's two alternatives of choice, I and II, also represent the state of the world, which should be completely independent of B and over which B has no control at all. The symmetry of situations of A and B implies that A and B will have similar reasoning, but does not imply that whenever A chooses I, Bwill necessarily choose I too, nor that whenever A chooses II, B will necessarily choose II too. Therefore the analysis is extremely simple, i.e., either A's or B's decision is to make a choice of II, according to the dominance principle. That the situation of both A and B choosing II will end up with a consequence of nonmaximal personal utility is something that is out of control and cannot be helped. The dominance principle is still a rational guide to the decision-making of the agent for the obtaining of maximal utility, but it is not the responsibility of the dominance principle to guarantee maximal utility, because the obtaining of maximal utility depends not only on the rationality of choice alone, but also on the state of the world, which is beyond the control of the agent.

I am not alone in holding this view. James W. Friedman's opinion is exactly the same as mine. He says (Friedman, 1986):

The prisoner's dilemma is a famous example of a game with a dominant strategy equilibrium. As Table 4.1 (p. 109) reveals, *confess* dominates *not confess* for both players, and the equilibrium payoff is not Pareto optimal; both players have higher payoffs if both choose *not confess*. Thus, a dominant strategy equilibrium is a compelling outcome at which inefficient payoffs can be received by the players. A final point to note on this topic is that, if a player has two or more dominant strategies in a game, then these dominant strategies must be equivalent.

Therefore we are not entitled to say that a rational person will be able to make choices that will always result in maximal utility, but can only say that a rational person will make choices always with an *intention* to obtain maximal utility. Actually s/he may make mistakes because of insufficient knowledge, impossibility of having complete knowledge, and/or an adverse state of the world.

As to the latter case, namely the dynamic or repeated case, although I think it unreasonable to consider it a version of the prisoner's dilemma, the case can still be studied as a separate problem. In fact it has been studied by many game theorists, including J. W. Friedman (1986), R. Axelrod (1986), and several others in the anthology edited by A. Diekmann and P. Mitter (1986).

I give a brief discussion of it here. Consider the payoff matrix of Figure 1 representing the prisoner's dilemma, which is reproduced from Randall K. Campbell (1989) with the addition of probabilities of the states of the world p_1 and $p_2 = 1 - p_1$.

Since the matrix is a symmetrical one, prisoners A and B are symmetrically situated and, therefore, the probability distribution of prisoner A's choice, which is also the probability distribution of prisoner B's states of the world, is assumed to be equal to the probability distribution of prisoner B's choices, which is also the probability distribution of prisoner A's states of the world.

The expected utilities of A and B are equal and may be expressed as

$$U = 9p_1^2 + 10p_1(1 - p_1) + (1 - p_1)^2$$

= $9p_1^2 + 10p_1 - 10p_1^2 + 1 - 2p_1 + p_1^2$
= $8p_1 + 1$.

It can readily be seen that as p_1 increases from 0 to 1, U monotonically increases, too. Therefore maximal utility occurs at $p_1 = 1$, where the probability distribution is (1, 0). Or both A and B should choose action I in all the repeated games.

 P_2 will be nonzero when the nonzero entries in cells (II, I) and (I, II) increase to more than 18, or the entries in cell (I, I) decrease to

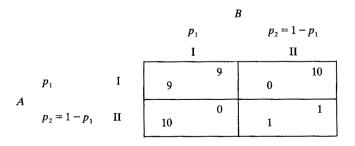


Fig. 1. A typical payoff matrix representing the prisoner's dilemma.

C. L. SHENG

less than 5. Let us consider the payoff matrix in Figure 2, where the entry 10 in cells (II, I) and (I, II) is increased to 30.

For the matrix in Figure 2 the expected utility for either prisoner A or prisoner B is

$$\begin{split} U &= 9p_1^2 + 30p_1(1-p_1) + (1-p_1)^2 \\ &= 9p_1^2 + 30p_1 - 30p_1^2 + 1 - 2p_1 + p_1^2 \\ &= -20p_1^2 + 28p_1 + 1 \\ \frac{\mathrm{d}U}{\mathrm{d}p_1} &= -40p_1 + 28 \;. \end{split}$$

Setting $dU/dp_1 = 0$, we have

$$p_1 = 28/40 = 0.7$$
.

Therefore this matrix has an equilibrium point for maximal utility and the probability distribution for maximal utility is (0.7, 0.3). The maximal utility is

$$U(\max) = -20 \times 0.7^2 + 28 \times 0.7 + 1 = 10.8$$
.

Although the problem of the dynamic or repeated case is solvable, I still do not regard it as a version of the prisoner's dilemma, because the solution of the problem hinges on the prisoners' thinking in terms of game theory, which is far beyond the capacity of an ordinary person and, hence, also far beyond the general assumption of rationality. If one does not know game theory at all, one would not be able to

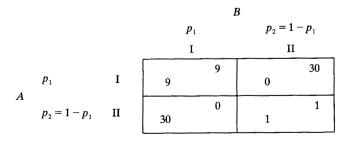


Fig. 2. Another payoff matrix with different entries in cells (II, I) and (I, II).

calculate the optimum probability distribution and to make choices between actions I and II according to this probability distribution in a random way.

3. AN ANALYSIS FROM A MORAL POINT OF VIEW

Since the prisoner's dilemma is also a moral problem, I now discuss it from a moral point of view. The prisoner's dilemma is regarded by many philosophers as a group action to symbolize the essence of morality. In fact it is the thesis of many books and articles. For instance, Donald H. Regan (1980) emphasizes cooperation and David Gauthier (1986) emphasizes agreement among members of society. A contractarian theory, such as John Rawls' (1971) theory of justice, emphasizes a contract to be established among individuals or between an individual and society in terms of law and morality, so as to ensure the maximization of aggregate or social utility. I (Sheng, 1987) have discussed elsewhere the constraints on group actions and have proposed a utilitarian interpretation of fairness in such a situation.

The prisoner's dilemma, although a moral problem, can still be treated quantitatively. However, if the decision-theoretic model is still used, then there arise many difficulties. I think the main difficulty lies in the obscurity of social utility or how to count social utility, because there can be several different interpretations of it.

First, social utility may be taken to be the sum of the utilities for all members of society. Then, in the case of the prisoner's dilemma, in each cell of the matrix in Figure 1, the social utility is the sum of the two entries. For instance, in cell (I, I) of Figure 1, the social utility is 9 + 9 = 18.

To take the sum of personal utilities to be the social utility is certainly right for public actions, particularly in the study of distributive justice. But for a personal action it is too high a requirement to be met, because it requires one to count the utility for any other person as important as the utility for oneself. Very few persons, even with high morality, are able to attain this high level. It may be regarded as a moral ideal, but certainly should not be set as a moral requirement.

On the other hand, however, some charitable and supererogatory

actions are far beyond that. An agent sometimes takes an action to produce a utility for some recipient(s) at the expense of the utility for her/himself, which may be much greater than the utility produced for the recipient(s). In that case the social utility, if taken to be the sum of utilities for both the agent and the recipient(s) (assuming that the utilities for the remaining members of society are unaffected by the action), is decreased rather than maximized. Yet such a supererogatory action is still considered not only right, but also very good and very virtuous. To conform to the principle of utility, it seems that now social utility should be taken to be the sum of the utilities for all members of society excluding the agent her/himself. This seems to be a still higher requirement than that dictated by the social utility as the sum of utilities for all members including the agent her/himself.

The two views discussed above are in opposite directions. My explanation of this discrepancy is that it is due to what I call the 'flexible nature of morality' (1986). I hold that, for certain charitable and supererogatory actions, it is very difficult, and sometimes impossible, to set a clear-cut moral requirement. In fact, in the descending order of moral requirement, there are moral duties, moral actions that are not moral duties but actions that one ought to take, charitable actions that one does not ought to take, and supererogatory actions.

Because of this difficulty, instead of setting a moral requirement, I emphasize the utility of the feeling of moral satisfaction for the agent and regard it simply as a fact or phenomenon. I then study moral actions still using the decision-theoretic model from the *self-interest* point of view, but with the utility produced by the feeling of moral satisfaction for the agent counted as an extra part of the total utility for the agent. I believe that this approach, naively simple and straightforward as it may appear to be, still serves the purpose of solving the prisoner's dilemma and, hence, explaining morality.

When one takes a morally good action, one has a feeling of moral satisfaction, and when one takes a morally bad action, one has a feeling of moral dissatisfaction. The utility/disutility of the feeling of moral satisfaction/dissatisfaction for the agent is a mathematical function of the utility/disutility of the consequences produced by the action for the recipient(s) of the action. Let U_c be the utility of the consequences of the action for recipient(s), and U_m be the utility of the

feeling of moral satisfaction for the agent. Then we have

(1)
$$U_{\rm m} = f(U_{\rm c})$$
.

It is quite natural that when the utility of the consequences of an action is small, the feeling of moral satisfaction is also small, and that when the utility is large, the feeling of moral satisfaction is also large. Thus it is reasonable to assume that $U_{\rm m}$ is directly proportional to $U_{\rm c}$. Then we have

(2)
$$U_{\rm m} = k U_{\rm c}$$

where k is a coefficient depending on the morality of the agent and is called *coefficient of morality*. It is seen that k can serve as a general index for the morality of the agent. For normal people, it seems reasonable to set the range of k to be from 0 to 1, i.e.,

$$(3) \qquad 0 \leq k \leq 1.$$

The lower bound (k = 0) represents complete indifference or no sympathy at all, which means that the agent feels nothing when the recipient receives a utility U_c . It represents the lower limit of morality.

The upper bound (k=1) represents complete sympathy, which means that, when the recipient receives a utility U_c , the agent feels as if s/he her/himself received it. It represents the upper limit of morality.

The upper and lower limits of morality roughly correspond to R. M. Hare's two levels of moral thinking. Moreover, that the actual levels of morality form a continuous spectrum ranging from k = 1 to k = 0 is also compatible with the following statement of Hare (1981, p. 45):

Although the archangel and the prole are exaggerated versions of the top and bottom classes in Plato's Republic, it is far from my intention to divide up the human race into archangels and proles; we all share the characteristics of both to limited and varying degrees and at different times.

Of course k can be negative. In that case the agent enjoys doing or having done a bad thing and suffers when doing or having done a good thing. This is extremely abnormal and also very exceptional. So the cases of negative k will not be considered.

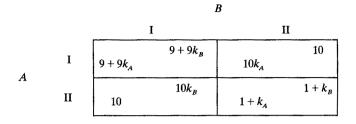


Fig. 3. Payoff matrix of example of Figure 1, with total utilities for A and B including the utilities derived from the feeling of moral satisfaction.

I shall now analyze the prisoner's dilemma by taking total utility to be the sum of the utility for the agent her/himself and k times the utility for the other person. Consider the example of Figure 1. Let the entries in each cell of the pay-off matrix be replaced by the two total utilities for A and B, respectively, each including the utility derived from the factor of the feeling of moral satisfaction. The matrix is shown in Figure 3.

Now the choice of A (similarly of B) will depend on the value of k_A (similarly k_B). When $k_A = 0$, which means that A is completely indifferent to the utility for B, the choice according to the dominance argument will be II. When $k_A = 1$, which means that A has full sympathy with B or that A feels as if the utility for B were for A, the choice will be I. Even if k_A is as low as 0.2, the choice is still I. Incidentally, for the matrix of Figure 3, the dominance argument is always applicable. There is a threshold at $k_A = 0.111 \dots$ When $k_A > 0.111 \dots$, the dominance is for choice I, but when $k < 0.111 \dots$, the dominance II.

4. A MORE GENERAL CASE

It is seen from the above analysis that, in the example of Figure 3, there is a threshold of 0.111..., above which dominance exists and the choice is I, and below which dominance also exists but the choice is II. In other words, the threshold 0.111... divides the decisions into two regions where dominance exists, but to different extents.

It is of interest to note that, if the entries in cell (I, I) have lower

values, say (8, 8) instead of (9, 9), then three regions will exist instead of two. In one region there is dominance for choice I, in a second region there is no dominance, and in a third region there is dominance for choice II. The payoff matrix with entries (8, 8) in cell (I, I) is shown in Figure 4 and the corresponding payoff matrix employing the coefficient of morality k is shown in Figure 5.

For the sake of simplicity, assume that $k_A = k_B$. When $k_A = k_B = 0.3$, the payoff matrix is shown in Figure 6, and it can be seen that the dominance is for choice I.

When $k_A = k_B = 0.2$, the payoff matrix is shown in Figure 7, and it can be seen that there is no dominance.

When $k_A = k_B = 0.1$, the payoff matrix is shown in Figure 8, and it can be seen that there is dominance again, but it is for choice II.

It can readily be shown that the thresholds are 0.25 and 0.111..., respectively. When k > 0.25, there is dominance and it is for choice I.

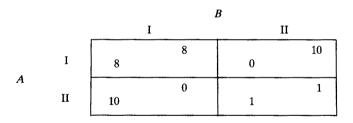


Fig. 4. A payoff matrix representing a more general case of the prisoner's dilemma.

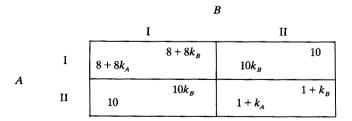


Fig. 5. Payoff matrix of example of Figure 3, with total utilities for A and B including the utilities derived from the feeling of moral satisfaction.

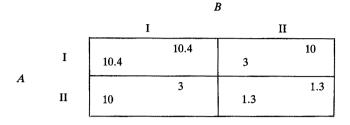


Fig. 6. A particular case of Figure 4, where $k_A = k_B = 0.3$.

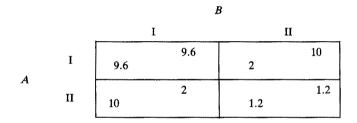


Fig. 7. A particular case of Figure 4, where $k_A = k_B = 0.2$.

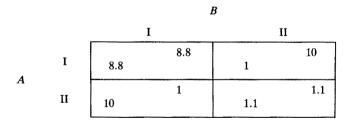


Fig. 8. A particular case of Figure 4, where $k_A = k_B = 0.1$.

When 0.111... < k < 0.25, there is no dominance. When k < 0.111..., there is dominance but it is for choice II.

It is seen that the greater the value of k, the larger the region of choice I will be. This is in conformity with my unified utilitarian theory (Sheng, 1991).

It is of interest to note that, as the entries in cell (I, I) decrease

further, the first threshold will rise. When the entries in cell (I, I) are (5, 5) or smaller, the region of dominance for choice I will disappear, and there will be only two regions again, one region of no dominance when k > 0.111..., and a second region of dominance for choice II when k < 0.111...

When dominance does not exist, then the choice has to be made either according to expected utility, if the probabilities of the states of the world are known, or according to some ordinal criterion, such as the maximax principle, the maximin principle, etc., if the probabilities are not known. Nevertheless, with the utility derived from the feeling of moral satisfaction taken into account as a part of the utility for the agent her/himself, the decision-theoretic model is still valid.

From the above analysis it can be seen that the decision to be made depends not only on the entries in the payoff matrix, but also on the coefficient of morality k. It is in this way that my analysis serves to explain the prisoner's dilemma. My explanation, although but one of many possible explanations, seems to be a feasible and reasonable one.

NOTE

¹ Previously, the prisoner's dilemma has been studied using probabilities in a quite different manner – from the point of view of conditional probability and statistical dependence between the choices of action by A and B. The identity or closeness of the choices by A and B are indicated in the conditional probabilities. Since I argue for the independence between A and B, I contend that the use of conditional probabilities is unjustified. This is perhaps the reason why Richmond Campbell distinguishes between causal dependence and probabilistic dependence and admits of causally independent but probabilistically dependent actions, which, incidentally, cause conflicts. As far as I see, probabilistic dependence is irrelevant to the problem.

REFERENCES

Axelrod, R.: 1986, The Evolution of Cooperation, Basic Books, New York.

- Campbell, Randall K.: 1989, 'The Prisoner's Dilemma and the Symmetry Argument for Cooperation', Analysis 49, 60-65.
- Campbell, Richmond: 1985, 'Introduction: Background for the Uninitiated', in *Paradoxes of Rationality and Cooperation*, Richmond Campbell and Lanning Sowden (Eds.), The University of British Columbia Press, Vancouver, Canada, pp. 3-41.

Davis, Lawrence: 1977, 'Prisoner's Paradox and Rationality', American Philosophical Quarterly 14, 319-327.

C. L. SHENG

- Dickmann, A. and Mitter, P. (Eds.): 1986, Paradoxical Effects of Social Behavior: Essays in Honor of Anatol Rapoport.
- Friedman, J. W.: 1986, *Game Theory with Applications to Economy*, Oxford University Press, New York.
- Gauthier, David: 1986, Morals by Agreement, Oxford University Press, Oxford, England.
- Hare, R. M.: 1981, Moral Thinking, Oxford University Press, Oxford, England.
- Rawls, John: 1971, A Theory of Justice, Harvard University Press, Cambridge, Massachusetts, U.S.A.
- Regan, Donald H.: 1980, Utilitarianism and Cooperation, Oxford University Press, Oxford, England.
- Sheng, C. L.: 1986, 'On the Flexible Nature of Morality', *Philosophy Research Archives* 12, 125–142.
- Sheng, C. L.: 1987, 'Constraints on Utilitarian Prescriptions for Group Actions', *Theory* and Decision 23, 301-316.
- Sheng, C. L.: 1991, A New Approach to Utilitarianism: A Unified Utilitarian Theory and its Application to Distributive Justice, Dordrecht, The Netherlands: Kluwer Academic Publishers.

Graduate Institute of Management Sciences, Tamkang University, Taipei, Taiwan, Republic of China.