# A note on the "index of cooperation" for Prisoner's Dilemma[1]

ANATOL RAPOPORT
*Mental Health Research Institute, The University of Michigan*

Numerous experiments have shown that the frequency of cooperative (C) responses observed in iterated plays of Prisoner's Dilemma depends on the particular payoff structure of the game. The generalized payoff matrix of Prisoner's Dilemma is represented by:

|   | C | D |
|---|---|---|
| C | R, R | S, T |
| D | T, S | P, P |

where, in accordance with the definition of Prisoner's Dilemma, the following inequality must be satisfied:

$$T > R > P > S . \qquad (1)$$

Assuming a simple reinforcement effect in repeated plays, we would expect that the frequency of C choices would increase as R and S increase (since these payoffs are associated with the C choice), and correspondingly, the frequency of C should decrease as T and P increase, because these payoffs are associated with the D choice. On the whole, experimental evidence tends to confirm this expectation. Consequently, if we imagine the "index of cooperation" for Pris-

oner's Dilemma to be a function of the four payoffs, namely,

$$K = K (T, T, P, S) , \qquad (2)$$

the first requirement we should impose on such a function is

$$\frac{\partial K}{\partial R} \geqslant 0; \; \frac{\partial K}{\partial S} \geqslant 0; \; \frac{\partial K}{\partial T} \leqslant 0; \; \frac{\partial K}{\partial P} \leqslant 0. \quad (3)$$

Clearly, this restriction is still too mild to allow us to choose among the innumerable functions which satisfy inequalities (3). Can we impose any additional restrictions? From the game theoretical point of view, it is tempting to demand

$$K (mT + n, mR + n, mS + n)$$
$$= K (T, R, P, S), \text{ for } m \geqslant 0, \quad (4)$$

in accordance with the assumption usually made in game theory that the payoffs are utilities, given only up to a positive linear transformation.

Evidence for or against (4) is scanty, but what there is seems to be in favor.

If we impose condition (4), the class of functions to choose from becomes much smaller. Of these there is a simplest subclass, namely quotients of linear polynomials in the differences $(T - R)$, $(T - P)$, $(R - P)$, etc.

It can be easily shown that all such homogeneous bilinear functions can be represented as

$$K (T, R, P, S)$$
$$= \frac{a'T + b'R + c'P - (a' + b' + c')S}{\alpha'T + \beta'R + \gamma'P - (\alpha' + \beta' + \gamma')S} \quad (5)$$
$$= \frac{a'(T - R) + b'(R - S) + c'(P - S)}{\alpha'(T - R) + \beta'(R - S) + \gamma'(P - S)}.$$

There are, of course, many other possible representations of formula (5). For example, we could write (5) as

$$K = \frac{a(R - P) + b(T - R) + c(P - S)}{\alpha(R - P) + \beta(T - R) + \gamma(P - S)}, \quad (6)$$

where $b = a'$, $a - b = b'$, $c - a = c'$, $\beta = \alpha'$, $\alpha - \beta = \beta'$, $\gamma - \alpha = \alpha'$.

The differences in (6) are especially suggestive in the context of Prisoner's Dilemma. Thus, assuming that the other player will choose C, the (positive) difference $(T - R)$ is a reasonable measure of the pressure to play D. Assuming that the other player will choose D, $(P - S)$ is the corresponding measure of the pressure to play D. There is also a counterpressure to play C, namely, the degree to which R (the reward for double C) is larger than P (the punishment for double D); in other words, the difference $(R - P)$. Since a linear combination of the six differences can be expressed as a linear combination of any three of them, the three differences in (6) seem to be the most appropriate for the reasons stated. Moreover, without loss of generality, we can set $\alpha = 1$ by properly choosing the payoff units.

Let us now see whether additional restrictions can be imposed, so as to simplify our index still further. We shall suppose K to be the actual probability of choosing C, as manifested by the relative frequency of this choice.

Consider a "degenerate" Prisoner's Dilemma game, in which $T = R$, $P = S$. In this game, there ought to be no pressure to play D, because no gain accrues to the "defector," whether the other plays C or D. On the other hand, if $R > P$, both gain by playing C. Therefore, logically we may expect

$$K (R, R, P, P) = 1. \quad (7)$$

Next, consider another degenerate game where $R = P$. Here there is nothing to gain from playing C. Hence, logically, we ought to have

$$K (T, R, R, S) = 0, \quad (8)$$

which implies $b = c = 0$.

Combining conditions (7) and (8), we should have

$$K = \frac{R - P}{(R - P) + \beta(T - R) + \gamma(P - S)}, \quad (9)$$

which is the simplest index of cooperation compatible with the four postulates:

$A_1$: The index is invariant with respect to positive linear transformations of the payoffs.

$A_2$: There is no defection (D) if the corresponding pressure is zero while the pressure for cooperation remains positive.

$A_3$: There is no cooperation (C) if the corresponding pressure is zero while pressure for defection remains positive.

$A_4$: Inequalities (3) hold.

Formula (9) has just two free parameters, $\beta$ and $\gamma$. The only restriction on these is that both be nonnegative, for this is the necessary and sufficient condition to satisfy $A_4$. Note that if $\beta = \gamma = 1$, (9) reduces to

$$K = \frac{R - P}{T - S}, \quad (10)$$

which is one of the two indices suggested elsewhere (Rapoport and Chammah, 1965).

If we take a linear combination of the two indices suggested in the above mentioned monograph, we obtain

$$K = \frac{m(R - P) + n(R - S)}{T - S}, \quad (11)$$

which is a special case of (6) with $m + n = a - b$, $a = c = n$, $\beta = \gamma = 1$.

There is evidence that even if $T = R$, $P = S$, defections occur. These may be due to the fact that subjects pay attention to the payoff *difference*, rather than to their own payoffs (the so-called "zero-sum bias"), as suggested by various authors (e.g., Scodel *et al.*, 1959); or they may be due to boredom (in long sequences of repeated plays), to error, or to malice.

Similarly, there is evidence that even if $R = P$, there will be some C choices. These may be due to "pure altruism," as well as to boredom or to error. To account for these cases, postulates $A_2$ and $A_3$ must be dropped. We have, then, our original formula (6) with $\alpha = 1$, on which, however, the restrictions embodied in $A_4$ must be imposed. It can be easily verified that these restrictions are as follows:

$0 < a < 1; \quad b\beta \geqslant 0; \quad b\gamma \geqslant 0; \quad c\gamma \geqslant 0;$
$b\gamma \geqslant c\beta; \quad a\gamma \geqslant c; \quad \beta \geqslant b; \quad a\beta > b;$
$a\gamma + \beta c \geqslant b\gamma + c; \quad a\beta + \beta c \geqslant b + b\gamma;$
$a\beta + \gamma b \geqslant c\beta + b .$                    (12)

Inequalities (12) are not inconsistent, as is shown by the following numerical example, which satisfies $A_1$ through $A_4$:

$$K = \frac{.8\,(R-P) + (T-R) + (P-S)}{R-P+2(T-R)+3(P-S)}. \quad (13)$$

Robert Axelrod (in this issue) has proposed an "index of conflict," namely

$$K' = \frac{(T-R)(T-S)}{(T-P)^2}. \quad (14)$$

Being an index of "conflict" rather than cooperation, Axelrod's index should be related inversely to our proposed index K. Indeed, it is seen by inspection that

$$\frac{\partial K'}{\partial R} \leqslant 0; \quad \frac{\partial K'}{\partial S} \leqslant 0; \quad \frac{\partial K'}{\partial P} \geqslant 0, \quad (15)$$

identically. However, the inequality $\partial K'/\partial T \geqslant 0$ may be violated for values of R close to P.

Clearly, if we are willing to consider more complex indices than bilinear functions of the payoffs, the choice of function increases enormously.

It is not easy to decide the relative merits of the various indices proposed. The formula (9) above has the advantage of simplicity and of a straightforward interpretation of the parameters $\beta$ and $\gamma$ (weights assigned to the "greed pressure" $[T-R]$ and the "fear pressure" $[P-S]$). On the other hand, the evidence obtained from "degenerate" games suggests that this index is not adequate.

Formula (6), involving five "semi-free" parameters subject only to inequality restrictions, is extremely flexible and so can be "fitted" to the varied sets of data. The flexibility, however, is obtained at a cost of losing the straightforward interpretation of the parameters.

Axelrod's formula, derived from the theory of negotiable games, has a suggestive geometric interpretation.

Criteria for choosing among the various formulae can be established only after very large masses of data under a variety of conditions have been accumulated.

It is important to keep in mind the *ad hoc* character of all such indices. In particular, if the C frequencies are estimated from the protocols of a whole population of players, it is not legitimate to conclude that the index for the whole population is the same function of the payoffs as the index for an individual player.

Suppose, for example, the index for each individual player is given by (6), and suppose the parameters $\beta$ and $\gamma$ are distributed uniformly in the population of players in the range $0 < \beta < 1; 0 < \gamma < 1$. Then the "average" index observed in the population will be

$$\bar{K} =$$
$$\int_0^{\prime} \int_0^{\prime} \frac{R-P}{(R-P)+(T-R)+(P-S)} \, d\beta \, d\gamma$$
$$= \frac{R-P}{(T-R)(P-S)}[(T-S)\log(T-S) -$$
$$- (T-P)\log(T-P) - (R-S)\cdot$$
$$\log(R-S) + (R-P)\log(R-P)]. \quad (16)$$

We can verify that $\bar{K}$ satisfies postulates $A_1$ through $A_4$ and so is a "legitimate" index, derived by postulating (6) for each individual in the population and uniform independent distributions of the parameters. Clearly, different distributions of the parameters (e.g., Gaussian) will yield still different forms of $\bar{K}$. Since an estimate of the parameters of individuals, and especially of their distributions in the population, is usually extremely difficult, the limitations of such a "theoretical" approach are apparent. For this reason one is forced for the time being to choose indices more or less arbitrarily and to justify them on "untheoretical," purely pragmatic grounds.

## REFERENCES

AXELROD, ROBERT. "Conflict of Interest: An Axiomatic Approach," *Journal of Conflict Resolution* (this issue).

RAPOPORT, A., and A. M. CHAMMAH. *Prisoner's Dilemma: A Study in Conflict and Cooperation.* Ann Arbor: University of Michigan Press, 1965.

SCODEL, A., J. S. MINAS, P. RATOOSH, and M. LIPETZ. "Some Descriptive Aspects of Two-Person Non-Zero-Sum Games," *Journal of Conflict Resolution*, 3, 2 (June 1959), 114–19.