

Yale University

## EliScholar – A Digital Platform for Scholarly Publishing at Yale

---

Cowles Foundation Discussion Papers

Cowles Foundation

---

12-1-1957

### Bargaining in Ignorance of the Opponents' Utility Function

John C. Harsanyi

Follow this and additional works at: <https://elischolar.library.yale.edu/cowles-discussion-paper-series>



Part of the [Economics Commons](#)

---

#### Recommended Citation

Harsanyi, John C., "Bargaining in Ignorance of the Opponents' Utility Function" (1957). *Cowles Foundation Discussion Papers*. 263.

<https://elischolar.library.yale.edu/cowles-discussion-paper-series/263>

This Discussion Paper is brought to you for free and open access by the Cowles Foundation at EliScholar – A Digital Platform for Scholarly Publishing at Yale. It has been accepted for inclusion in Cowles Foundation Discussion Papers by an authorized administrator of EliScholar – A Digital Platform for Scholarly Publishing at Yale. For more information, please contact [elischolar@yale.edu](mailto:elischolar@yale.edu).

COWLES FOUNDATION DISCUSSION PAPER NO. 46

Note: Cowles Foundation Discussion Papers are preliminary materials circulated privately to stimulate private discussion and critical comment. References in publications to Discussion Papers (other than mere acknowledgment by a writer that he has access to such unpublished material) should be cleared with the author to protect the tentative character of these papers.

Bargaining

in Ignorance of the Opponent's Utility Function.\*

John C. Harsanyi

December 11, 1957

---

\* Research undertaken by the Cowles Commission for Research in Economics under Contract Nonr-358(01), NR 047-006 with the Office of Naval Research.

---

1.

The Zeuthen-Nash theory of bargaining\* is based on the assumption

---

\* See John Nash, "The Bargaining Problem," Econometrica, vol. 18 (1950), pp. 115-162 -- and "Two-Person Cooperative Games," Ibid., vol. 21 (1953), pp. 128-140; F. Zeuthen, Problems of Monopoly and Economic Warfare, London, 1930, chap. IV.

On the mathematical equivalence of Nash's and Zeuthen's theories, see John C. Harsanyi, "Approaches to the Bargaining Problem Before and After the Theory of Games," Econometrica, vol. 24 (1956), pp. 144-157.

---

that the two bargaining parties fully know each other's cardinal utility functions, which essentially means that they know each other's preferences as well as each other's attitudes towards risk. On the basis of this assumption and some others, the theory defines optimal strategies for both parties and also undertakes to uniquely predict the terms of agreement between two rational bargainers. I now propose to discuss the generalization of the Zeuthen-Nash theory for the more realistic case where the two bargainers do not know each other's utility functions. I should also like to consider how much light is thrown by this more general model upon bargaining behaviour as observed in the real world.

2.

In bargaining each party's behaviour will primarily depend on what he expects the other party will do. Each party will try to estimate the probability that his opponent will eventually agree to various possible conditions; and then will insist upon obtaining that set of conditions which maximize his expected utility. If a given party has

a definite hypothesis (a single-valued expectation) as to what are the best terms he can obtain from his opponent, his bargaining policy will be simply to insist upon these particular terms.

What the Zeuthen-Nash theory of bargaining essentially proposes to do is to specify what are the expectations that two rational bargainers can consistently entertain as to each other's bargaining strategies if they know each other's utility functions. The fundamental postulate of the theory is a symmetry axiom, which states that the functions defining the two parties' optimal strategies in terms of the data (or, equivalently, the functions defining the two parties' final payoffs) have the same mathematical form, except that, of course, the variables associated with the two parties have to be interchanged. Intuitively the assumption underlying this axiom is that a rational bargainer will not expect a rational opponent to grant him larger concessions than he would make himself under similar conditions.

The prediction made by the Zeuthen-Nash theory is that both bargainers will agree, and will expect each other to agree, to a pair of payoffs  $u_1$  and  $u_2$  such that maximize the product

$$\pi = \pi(U_1, U_2) = \{ U_1(t) - U_1(c) \} \cdot \{ U_2(t) - U_2(c) \}$$

where  $U_1(t)$  and  $U_2(t)$  are the utilities that the two parties assign to the terms  $t$  while  $U_1(c)$  and  $U_2(c)$  are the utilities they assign to the conflict situation. The maximization of this product is subject to the conditions:

$$U_1(t) \geq U_1(c)$$

$$U_2(t) \geq U_2(c)$$

$$t \in T$$

where  $T$  is the set of all possible terms available to the two parties. The terms  $t$  that actually maximize the product  $\pi$  will be denoted by  $t_0$ .

3.

Extension of this result to the case where the two bargainers do not know each other's utility functions but do form determinate estimates on each other's utility functions is fairly straightforward.

Let  $U_1$  and  $U_2$  be the two parties' "true" utility functions as before. We define the operators  $e_1$  and  $e_2$  with the meanings: "the first party's estimate of ..." and "the second party's estimate of ..."

Then, it seems to be rational, at least as a matter of first approximation, for the first party to assume that the best terms he can obtain from the second party are the terms  $t_0$ , which represent the solution of the bargaining problem according to the Zeuthen-Nash theory. But he will not be able to calculate  $t_0$  because this would require maximizing the product  $\pi$ , defined in terms of his opponent's utility function  $U_2$ , which he does not know. The best he can do is to estimate  $t_0$  by maximizing the product

$$\pi_1 = \pi(U_1, e_1 U_1) = \{U_1(t) - U_1(c)\} \cdot \{e_1 U_2(t) - e_1 U_2(c)\}$$

which is obtained from  $\pi$  by substituting the estimated utility function  $e_1 U_2$  for the true utility function  $U_2$ . The maximization of  $\pi_1$  will be subject to the conditions:

$$\begin{aligned} U_1(t) &\geq U_1(c) \\ e_1 U_2(t) &\geq e_1 U_2(c) \\ t &\in T. \end{aligned}$$

The terms  $t$  that maximize the product  $\pi_1$  can be denoted by  $e_1 t_0$  as they represent the first party's estimate of  $t_0$ . The rational policy for the first party will be to demand these terms  $e_1 t_0$  from his opponent. Similarly, for the second party it will be rational in first approximation to insist upon those terms  $t$  that maximize the product

$$\pi_2 = \pi(e_2 U_1, U_2) = \{e_2 U_1(t) - e_2 U_1(c)\} \cdot \{U_2(t) - U_2(c)\}$$

subject to analogous conditions. The terms  $t$  that maximize  $\pi_2$  can be regarded as the second party's estimate of  $t_0$  and can be denoted by  $e_2 t_0$ .

But more sophisticated bargainers can do better than that. If the first party knows that the second party's policy is to insist upon the terms  $e_2 t_0$ , rather than upon the terms  $t_0$ , his best reply is to estimate  $e_2 t_0$  and demand the terms  $e_1 e_2 t_0$  corresponding to his estimate of  $e_2 t_0$ . He can calculate  $e_1 e_2 t_0$  by maximizing the product

$$\pi_{21} = \pi(e_1 e_2 U_1, e_1 U_2) = \{e_1 e_2 U_1(t) - e_1 e_2 U_1(c)\} \cdot \{e_1 U_2(t) - e_1 U_2(c)\}$$

Similarly, for the second bargainer it will be rational to insist upon the terms  $e_2 e_1 t_0$  that represent his estimate of  $e_1 t_0$  and are calculated by maximizing the product

$$\pi_{12} = \pi(e_2 U_1, e_1 e_2 U_2) = \{e_2 U_1(t) - e_2 U_1(c)\} \cdot \{e_2 e_1 U_2(t) - e_2 e_1 U_2(c)\}$$

But once the first bargainer realizes that this is the strategy of the second bargainer, he may do better if he demands the terms  $e_1 e_2 e_1 t_0$ , which represent his estimate of  $e_2 e_1 t_0$  and are calculated

by maximizing the product

$$\pi_{121} = \pi(e_1 e_2 U_1, e_1 e_2 e_1 U_2) = \\ \{e_1 e_2 U_1(t) - e_1 e_2 U_1(c)\} \cdot \{e_1 e_2 e_1 U_2(t) - e_1 e_2 e_1 U_2(c)\}$$

and so on.

Thus the two bargainers will have well-defined optimal strategies only if the two series

$$e_1 t_0, e_1 e_2 t_0, e_1 e_2 e_1 t_0, \dots, e_1 (e_2 e_1)^k t_0, (e_1 e_2)^{k+1} t_0, \dots \\ e_2 t_0, e_2 e_1 t_0, e_2 e_1 e_2 t_0, \dots, e_2 (e_1 e_2)^k t_0, (e_2 e_1)^{k+1} t_0, \dots$$

converge. Actually, two bargaining parties will seldom have enough information (and enough computing ability) to make it worth their while to carry on this estimating process for more than a very small number of steps, which can be formally interpreted as meaning that each bargainer assigns the same constant value to his own series after (say) the  $k$ th member.\* (Of course,  $k$  need not be the same for

\* It should be noted that the preceding argument is at each step subject to the overriding restriction that

$$U_i(t_i^{(k)}) \geq U_i(c) \quad i = 1, 2$$

where  $t_i^{(k)}$  is the  $i$ th party's demand at the  $k$ th step. In other words, irrespective of what he thinks that the other party's bargaining policy will be, neither party will agree to terms less favourable to him than would be the conflict situation itself.

both parties.)

This model has the following interesting implication. In most other parts of economic theory it is a general rule that being in error can never improve one's position. But in the theory of bargaining this proposition is no longer true. Suppose that a person wants to sell a house worth \$10,000 to him, to another person to whom the house is worth \$30,000. Suppose also that both persons' utility functions are linear in money. Then, if both know each other's limit prices, the Zeuthen-Nash theory predicts that they will split the difference and the house will be sold for \$20,000. But now suppose that the seller mistakenly believes that the house is worth as much as \$48,000 to the buyer and, splitting the imagined difference, insists upon a price of \$29,000. Then it will be rational for the buyer to accept this price (if he cannot convince the seller that his true demand price for the house is less than the latter thinks) rather than go without the house, as \$29,000 is still less than the \$30,000, which represent the value of the house to him.

In general, if a given bargaining party overestimates the strength of his own bargaining position, and if this is known to the other party, his error will tend to benefit him. Of course, if the second party does not know about the first party's error, or does not accept it as a genuine error but regards it as a mere bluff, then both parties will lose because agreement will fail. On the other hand, if a given party underestimates the strength of his own position, he will lose and the other party will gain as the outcome will shift in favour



of the second party.\*

---

\* The fact that being in error may improve one's bargaining position is a special case of the more general fact, recently stressed by Thomas C. Schelling, that in bargaining any "weakness" may prove to be a source of strength. See his "An Essay on Bargaining," American Economic Review, vol. XLVI, June 1956, pp. 281-306.

---

4.

In the more general case the two parties will not form determinate estimates on each other's utility functions and on each other's estimates concerning these utility functions etc. Instead, they will act on the basis of a priori probability distributions they assign to these variables.

Consider the decision problem confronting the first party. As a matter of first approximation he can again assume that the best terms he can get from the other party are the terms  $t_0$  maximizing the product  $\pi$ , i.e. the terms corresponding to the Zuthen-Nash solution in terms of the two parties' "true" utility functions  $U_1$  and  $U_2$ . Of course, the first party does not know  $U_2$  and, on our present assumption, does not form even a unique estimate of  $U_2$ . Instead, he has in mind a number of alternative hypotheses on  $U_2$ , and has an a priori probability distribution on all these alternative possibilities. Let  $U_2^H$  be the second player's utility function on hypothesis  $H$ . Then the terms  $t_0^H$  which maximize the product  $\pi^H = \pi(U_1, U_2^H)$  will represent the best terms that the first bargainer

can expect to obtain from the second on hypothesis  $H$ , and  $\bar{u}_1 = U_1(t_0^H)$  will be the utility level that he can expect to achieve on hypothesis  $H$ . To each hypothesis  $H$  there belongs a unique utility  $\bar{u}_1$ . Therefore the first bargainer's a priori probability distribution over all alternative hypotheses  $H$  will generate a probability distribution  $F_1(\bar{u})$  for the variable  $\bar{u}_1 = U_1(t_0)$ , which represents the highest utility level that the first bargainer can demand without risking a conflict, on the assumption that the second bargainer's final terms are  $t_0$ . Thus we can define

$$F_1(u_1) = \text{Prob}\{\bar{u}_1 \leq u_1\}$$

Accordingly, as a first approximation the best policy for the first bargainer will be to insist on such terms  $t = t_1$  which maximize his expected utility. If the highest utility level that his opponent is prepared to grant him is  $\bar{u}_1$ , then if he insists upon such terms  $t$  for which  $U_1(t) > \bar{u}_1$  a conflict will result and he will achieve only the utility level  $U_1(c)$ . If he chooses terms  $t$  such that  $U_1(t) \leq \bar{u}_1$  he will achieve the utility level  $U_1(t) = u_1$ , his expected utility will be

$$u_1 \cdot F_1(u) + U_1(c) \cdot [1 - F_1(u)]$$

and his best policy under this first-order approximation will be to demand such terms  $t = t_1$  which maximize this expression.

A similar definition can be given to the terms  $t_2$ , which represent the second bargainer's demand if he follows his first-order approximation optimal policy.  $t_2$  will be defined in terms of the

a priori probability distribution  $F_2$  which the second bargainer attaches to alternative possible values of the variable  $\bar{u}_2 = U_2(t_0)$ , i.e. of the utility level that he can demand without bringing about a conflict.

But it will be true once more that either bargainer can improve his strategy by going over to a second-order approximation. For instance, the first bargainer will do better if he maximizes his expected utility in terms of his expectations concerning his opponent's choice of  $t_2$  rather than in terms of the "true" Zeuthen-Nash solution  $t_0$ . He can define an a priori probability distribution  $F_1^{(2)}$  for the variable  $\bar{u}_1^{(2)} = U_1(t_2)$  on the basis of the a priori probabilities that he attaches to various alternative hypotheses concerning his opponent's utility function  $U_2$  as well as the a priori probability distribution  $F_2$  used by his opponent. The best policy for the first bargainer in second approximation will be accordingly to insist upon such terms  $t_1^{(2)}$  which maximize his expected utility in terms of the probability distribution  $F_1^{(2)}$ , and so on.

Again, we obtain an infinite series of strategies for each bargainer, and his optimal strategy will be the limit of this series if such exists. But reasons similar to those mentioned in Section 3

make it likely that this series will usually converge.\*

---

\* The two models of Sections 3 and 4 can be easily extended, without the use of any new principle, to the case where the bargainers do not precisely know the utilities that they would attach themselves to various outcomes, or at least do not know what utilities to attach to the conflict situation e.g., because they do not know each other's physical damaging power.

---

5.

In the case where two bargainers know each other's utility functions (and know the strategies available to each other), the Zeuthen-Nash theory predicts that they will always reach an agreement -- at least if both of them follow rational policies. The only exception is the case where the two parties have no real interest in cooperation because the conflict situation is preferable, at least for one of them, to any possible agreement.

But in the case where the bargainers do not know each other's utility functions, there is always the possibility that agreement will fail even if both parties use their optimal strategies. This is so because each party has to select his optimal policy on the basis of the information available to him, and if on the basis of this information he overestimates the concessions he can obtain from his opponent, then the two parties will insist upon mutually incompatible demands, and agreement must fail.

Actually, most instances of bargaining seem to end with agreement, which seems to show that the two parties are free at least from

gross errors in judging each other's utility functions\*, and probably

---

\* More exactly, in judging their opponents' utility functions, they do not err in the direction of overestimating the concessions they can obtain from their opponents.

---

also shows that they tend to follow fairly cautious bargaining policies, i.e., tend to have conservative attitudes towards risk-taking, corresponding to cardinal utility functions convex to above.

But it is hard to imagine that bargainers should always be able to form exact estimates of the limit of their opponent's willingness to yield. Very likely, in most cases where agreement is reached the two parties would have been ready to make larger -- maybe much larger -- concessions than they have actually made. Or, in other words, the two parties must have made better bargains than they really expected.

It would be very interesting to find out empirically how the terms actually agreed upon tend to compare with the terms anticipated by each party before the negotiations.

6.

How do bargaining parties form their estimates (or a priori probability distributions) concerning each other's utility functions and other psychological variables?

Obviously they often possess a considerable amount of information -- and/or misinformation -- on each other's interests and attitudes when they arrive at the negotiation table. But how much additional

information can they collect by observing their opponent's behaviour during the negotiations?

We shall argue that, if both bargainers act rationally, neither of them will be able to obtain any information at all on the other party's true attitudes, during the negotiations -- though they may obtain information on facts subject to objective evidence (they may be shown ledgers, balance sheets, expert opinions etc.) This is so because neither party, if he acts rationally, will disclose any information that would weaken his own bargaining position. On the other hand, if a bargaining party tries to impart to the other party any information that would strengthen his own bargaining position, then this other party will have no reason to give any credit to this information received (unless some objective evidence is produced in its support) and will have every reason to regard it as a mere bluff.

This conclusion applies to the two parties' actual bargaining moves (their offers and counter-offers) no less than it applies to their other behaviour (their verbal statements, gestures etc.) during the negotiations. If a given bargainer refuses to make any worthwhile concessions during the negotiations (up to the very last move), this gives no information on his real attitudes as it may always be a mere bluff. On the other hand, if a bargainer makes more generous concessions than his opponent anticipated, he does give away information to his opponent, but this is obviously a tactical mistake on his part as he achieves nothing else than a worsening of his own bargaining position.

## 7.

We may distinguish two conceivable types of bargaining. One consists in a series of offers and counter-offers by the two parties. This may be called "extensive" bargaining, and this is what we normally think of when we speak of bargaining. The other type is usually not called "bargaining" at all, but for convenience we shall label it "single-move" bargaining. It could conceivably take two forms. Either, one party makes a single offer, indicating that this is his last offer, and the other party can do no more than either take it or leave it. Or, both parties simultaneously and independently of each other make a single bid (say, in closed envelopes), and then the two bids are compared: if they are compatible, then an agreement has been reached, while if they are incompatible, then agreement has failed. In the former case we may speak of "asymmetric" single-move bargaining while in the latter case of "symmetric" single-move bargaining.

As a matter of common knowledge, in the Western culture there has been a strong trend away from extensive bargaining and towards single-move bargaining when the system of fixed prices has gradually replaced the custom of bargaining in the usual sense. One of the reasons has been that extensive bargaining is a very time-consuming process. But a shift in the society's moral views concerning what is proper behaviour in business affairs has also played a role.

In any case, extensive bargaining still persists in a number of fields, one of the most important ones being collective bargaining on the labour market. There is a lot of extensive bargaining also

in politics. All these fields seem to have in common that asymmetric single-move bargaining is ruled out by the fact that neither party is strong enough to "dictate" terms to the other. But why do not they resort to some form of symmetric single-move bargaining?

The most obvious explanation seems to be that, before finally committing themselves, they wish to "test out" the opponent's attitudes, i.e. they want to gather information. Very likely, in addition, extensive bargaining also has a sort of "ceremonial" function. Public opinion tends to look with disfavour upon a bargainer who sticks to his first offer and makes no further concession -- it does not help if his first offer is already so moderate as to include all concessions he would be prepared to make. Conversely, if a bargainer negotiates not for himself alone but for a whole interest group (say, for an employer or an employee organization), his constituents will insist that he should "show a fight" and should not accept his opponents' first offer immediately.

But, if we imagine fully rational bargainers (as well as a fully rational public opinion), there would never be a case for extensive bargaining, and bargainers would always make their "last" offer already at the first step -- as both parties would know that they cannot possibly find out about the other party's attitudes anything that they did not already know, and as the "ceremonial" considerations would be absent.



Of course, if a rational bargainer is confronted with an opponent about whose rationality he is in doubt, i.e., whom he considers to be likely to make mistakes, it is fully rational for him to undertake extensive bargaining with him in an attempt to obtain information about his true attitudes in case he makes the mistake of disclosing some information of this sort.

8.

To sum up, in the case of bargaining with an opponent whose true utility function one does not know, the problem of decision making in principle requires an infinite series of steps, involving estimation of what the opponent's true utility function is, estimation of what he thinks of one's own utility function, estimation of what he thinks one thinks oneself on his utility function etc. This is true both in the case where the two parties think in terms of single-valued estimates and in the case where they think in terms of a priori probability distributions. The optimal strategy of each party is the limit of an infinite series of approximations if such limit exists.

Errors in judging each other's utility functions are the main explanation for the fact that even very intelligent bargainers may fail to reach an agreement even if this would be very important to both of them.

In bargaining between two fully rational bargainers neither party will disclose any information weakening his own bargaining position -- and cannot successfully convey to his opponent any information that

would strengthen his own bargaining position. (Only facts subject to objective evidence form an exception.)

Therefore bargaining between fully rational bargainers would always consist in one single move and counter-move; any offer made would amount to an ultimatum.

In the real world, there is a good deal of many-move bargaining. This may be due to the fact that the two parties expect each other to commit tactical mistakes and disclose information to their own disadvantage -- or it may have a mere ceremonial function.

It is an interesting problem for empirical research to find out

1. how definite and how realistic two bargainers' ideas usually are before the negotiations on the terms they can eventually achieve;
2. whether bargainers actually expect to obtain information on each other's attitudes during the negotiations;
3. whether in actual fact they do receive worthwhile information during the negotiations, i.e., whether they often change their views in bargaining on what terms they can expect to achieve, and whether these changes in their views tend to make their views more realistic or the other way round;
4. and what differences<sup>there</sup> are in all these respects between experienced and inexperienced negotiators.