# LETTER TO DAVID LEWIS May 21, 1972

It does seem to me worth noting that if P is a probability distribution, and if for any A and B,  $P_B(A) = P(B > A)$ , then  $P_B$  is a probability distribution too (excepting the absurd case). What it is good for, I would like to suggest, is deliberation – the calculation of expected utilities.

Let  $S_1, \ldots, S_n$  be an exhaustive set of mutually exclusive propositions characterizing the alternative possible outcomes of some contemplated action. Let A be the proposition that I perform the action. My suggestion is that expected utility should be defined as follows:

$$\mathbf{u}(A) = \mathbf{P}(A > S_1) \times \mathbf{u}(S_1) + \cdots + \mathbf{P}(A > S_n) \times \mathbf{u}(S_n).$$

Why  $P(A > S_i)$  rather than  $P(S_i/A)$ ? Because what is relevant to deliberation is a comparison of what will happen if I perform some action with what would happen if I instead did something else. A difference between P(S/A)and P(S) represents a belief that A is evidentially relevant to the truth of S, but not necessarily a belief that the action has any causal influence on the outcome. That a person performs a certain kind of action can be evidence that makes some state subjectively more probable, even when the action in no way contributes to the state. Suppose that this is true for some action A and desirable state S. Then P(S/A) > P(S), but only an ostrich would count this as any sort of reason inclining one to bring it about that A. To do so would be to act so as to change the evidence, knowing full well that one is in no way changing the facts for which the evidence is evidence.

I am thinking of Nozick's puzzle ("Newcomb's problem", in the Hempel festschrift), which I just discovered, but which I assume you know. My intuitive reaction to this puzzle was the following: there is only one rational choice (assuming there is no backwards causation in the case), and that is to choose the dominating action. But this seems to conflict with the principle of maximizing expected utility. But from my suggested version of the principle, the rational choice follows. The principle of expected utility may be held to be universally applicable.

Since quotient conditionalization is the way to revise your beliefs, it is also rational in the Newcomb problem to bet, after having made the rational choice, that you will fail to get the million dollars. Had you made the other

W. L. Harper, R Stalnaker, and G. Pearce (eds.), Ifs, 151-152. Copyright © 1980 by D. Reidel Publishing Company.

151

choice, it would have been rational to bet that you would succeed in getting the million dollars. But *this* is no reason to wish that you had chosen differently, since you could have changed only the fair betting odds, not the facts, by acting differently.

The suggested version of the expected utility principle makes it possible for a single principle to account for various mixed cases: the probabilistic dependence may have two components, one causal and one non-causal. The components may reinforce each other, or counteract each other. They might cancel out, leaving the evidence irrelevant, even though there is a believed causal dependence. Also, it may be unknown whether the probabilistic dependence is causal or not. Imagine a man deliberating about whether or not to smoke. There are two, equally likely hypotheses (according to his beliefs) for explaining the statistical correlation between smoking and cancer: (1) a genetic disposition to cancer is correlated with a genetic tendency to the sort of nervous disposition which often inclines one to smoke. (2) Smoking, more or less, causes cancer in some cases. If hypothesis (1) is true, he has no independent way to find out whether or not he has the right sort of nervous disposition. In such a case, it seems clear that the probability of the conditional (if I were to smoke, I would get cancer), and not the conditional probability is what is relevant ....

# COUNTERFACTUALS AND TWO KINDS OF EXPECTED UTILITY\*

## **1. INTRODUCTION**

We begin with a rough theory of rational decision-making. In the first place, rational decision-making involves conditional propositions: when a person weighs a major decision, it is rational for him to ask, for each act he considers, what would happen if he performed that act. It is rational, then, for him to consider propositions of the form 'If I were to do *a*, then *c* would happen'. Such a proposition we shall call a *counterfactual*, and we shall form counterfactuals with a connective ' $\Box$ -' on this pattern: 'If I were to do *a*, then *c* would happen' is to be written 'I do  $a \Box \rightarrow c$  happens'.

Now ordinarily, of course, a person does not know everything that would happen if he performed a given act. He must resort to probabilities: he must ascribe a probability to each pertinent counterfactual 'I do  $a \Box \rightarrow c$  happens'. He can then use these probabilities, along with the desirabilities he ascribes to the various things that might happen if he did a given act, to reckon the expected utility of a. If a has possible outcomes  $o_1, \ldots, o_n$ , the expected utility of a is the weighted sum

 $\Sigma_i \operatorname{prob} (\operatorname{I} \operatorname{do} a \Box \rightarrow o_i \operatorname{obtains}) \mathcal{D} o_i,$ 

where  $\mathcal{D}o_i$  is the desirability of  $o_i$ . On the view we are sketching, then, the probabilities to be used in calculating expected utility are the probabilities of certain counterfactuals.

That is not the story told in familiar Bayesian accounts of rational decision; those accounts make no overt mention of counterfactuals. We shall discuss later how Savage's account (1972) does without counterfactuals; consider first an account given by Jeffrey (1965, pp. 5–6).

A formal Bayesian decision problem is specified by two rectangular arrays (matrices) of numbers which represent probability and desirability assignments to the act-condition pairs. The columns represent a set of incompatible conditions, an unknown one of which actually obtains. Each row of the desirability matrix,

 $d_1 d_2 \ldots d_n$ 

represents the desirabilities that the agent attributes to the n conditions described by the

153

W. L. Harper, R. Stalnaker, and G. Pearce (eds.), Ifs, 153–190. Copyright © 1978 by D. Reidel Publishing Company. column headings, on the assumption that he is about to perform the act described by the row heading; and the corresponding row of the probability matrix,

$$p_1 p_1 \ldots p_n$$

represents the probabilities that the agent attributes to the same n conditions, still on the assumption that he is about to perform the act described by the row heading. To compute the expected desirability of the act, multiply the corresponding probabilities and desirabilities, and add:

$$p_1d_1+p_2d_2+\ldots+p_nd_n.$$

On the Bayesian model as presented by Jeffrey, then, the probabilities to be used in calculating 'expected desirability' are 'probabilities that the agent attributes' to certain conditions 'on the assumption that he is about to perform' a given act. These, then, are conditional probabilities; they take the form *prob* (S/A), where A is the proposition that the agent is about to perform a given act and S is the proposition that a given condition holds.

On the account Jeffrey gives, then, the probabilities to be used in decision problems are not the unconditional probabilities of certain counterfactuals, but are instead certain conditional probabilities. They take the form *prob* (S/A), whereas on the view we sketched at the outset, they should take the form *prob*  $(A \square \rightarrow S)$ . Now perhaps, for all we have said so far, the difference between these accounts is merely one of presentation. Perhaps for every appropriate A and S, we have

(1) 
$$prob(A \Box \rightarrow S) = prob(S/A);$$

the probability of a counterfactual  $A \square \rightarrow S$  always equals the corresponding conditional probability. That would be so if (1) is a logical truth. David Lewis, however, has shown (1976) that on certain very weak and plausible assumptions, (1) is not a logical truth: it does not hold in general for arbitrary propositions A and S.<sup>1</sup> That leaves the possibility that (1) holds at least in all decision contexts: that it holds whenever A is an act an agent can perform and *prob* gives that agent's probability ascriptions at the time.

In Section 3, we shall state a condition that guarantees the truth of (1) in decision contexts. We shall argue, however, that there are decision contexts in which this condition is violated. The context we shall use as an example is patterned after one given by Stalnaker. We shall follow Stalnaker in arguing that in such contexts, (1) indeed fails, and it is probabilities of counterfactuals rather than conditional probabilities that should be used in calculations of expected utility. The rest of the paper takes up the ramifications for decision theory of the two ways of calculating expected utility. In particular, the two opposing answers to Newcomb's problem (Nozick, 1969) are supported

respectively by the two kinds of expected utility maximization we are discussing.

We are working in this paper within the Bayesian tradition in decision theory, in that the probabilities we are using are subjective probabilities, and we suppose an agent to ascribe values to all probabilities needed in calculations of expected utilities. It is not our purpose here to defend this general tradition, but rather to work within it, and to consider two divergent ways of developing it.

## 2. COUNTERFACTUALS

What we shall be saying requires little in the way of an elaborate theory of counterfactuals. We do suppose that counterfactuals are genuine propositions. For a proposition to be a counterfactual, we do not require that its antecedent be false: on the view we are considering, a rational agent entertains counterfactuals of the form 'I do  $a \Box \rightarrow S$ ' both for the act he will turn out to perform and for acts he will turn out not to perform. To say  $A \Box \rightarrow S$  is not to say that A's holding would bring about S's holding:  $A \Box \rightarrow S$  is true also if S would hold regardless of whether A held.

These comments by no means constitute a full theory of counterfactuals. In what follows, we shall appeal not to a theory of counterfactuals, but to the reader's intuitions about them — asking the reader to bear clearly in mind that 'I do  $a \Box \rightarrow S$ ' is to be read 'If I were to do a, then S would hold'.

It may nevertheless be useful to sketch a theory that would support what we shall be saying; the theory we sketch here is somewhat like that of Stalnaker and Thomason (Stalnaker, 1968; Stalnaker and Thomason, 1970). Let *a* be an act which I might decide at time *t* to perform. An *a*-world will be a possible world which is like the actual world before *t*, in which I decide to do *a* at *t* and do it, and which obeys physical laws from time *t* on. Let  $W_a$  be the *a*world which, at *t*, is most like the actual world at *t*. Thus  $W_a$  is a possible world which unfolds after *t* in accordance with physical law, and whose initial conditions at time *t* are minimally different from conditions in the actual world at *t* in such a way that 'I do *a*' is true in  $W_a$ . The differences in initial conditions should be entirely within the agent's decision-making apparatus. Then 'I do  $a \Box \rightarrow S$ ' is true iff *S* is true in  $W_a$ .<sup>2</sup>

Two axioms that hold on this theory will be useful in later arguments. Our first axiom is just a principle of modus ponens for the counterfactual.

AXIOM 1.  $(A \& (A \Box \rightarrow S)) \supset S$ .

Our second axiom is a Stalnaker-like principle.

AXIOM 2. 
$$(A \Box \rightarrow \overline{S}) \equiv (\overline{A \Box \rightarrow S})$$
.

The rationale for this is that 'I do  $a \Box \rightarrow S$ ' is true iff S holds in  $W_a$  and 'I do  $a \Box \rightarrow \overline{S}$ ' is true iff  $\overline{S}$  holds in  $W_a$ . We shall also appeal to a consequence of these axioms.

CONSEQUENCE 1.  $A \supset [(A \Box \rightarrow S) \equiv S]$ .

We do not regard Axiom 2 and Consequence 1 as self-evident. Our reason for casting the rough theory in a form which gives these principles is that circumstances where these can fail involve complications which it would be best to ignore in preliminary work.<sup>3</sup> Our appeals to these Axioms will be rare and explicit. For the most part in treating counterfactuals we shall simply depend on a normal understanding of the way counterfactuals apply to the situations we discuss.

## 3. Two kinds of expected utility

We have spoken on the one hand of expected utility calculated from the probabilities of counterfactuals, and on the other hand of expected utility calculated from conditional probabilities. In what follows, we shall not distinguish between an act an agent can perform and the proposition that says that he is about to perform it; acts will be expressed by capital letters early in the alphabet. An act will ordinarily have a number of alternative outcomes, where an *outcome* of an act is a single proposition which, for all the agent knows, expresses *all* the consequences of that act which he cares about. An outcome, then, is a specification of what might eventuate which is complete in the sense that any further specification of detail is irrelevant to the agent's concerns, and it specifies something that, for all the agent knows, might really happen if he performed the act. The agent, we shall assume, ascribes a magnitude  $\mathcal{DO}$  to each outcome O. He knows that if he performed the act, one and only one of its outcomes would obtain, although he does not ordinarily know which of its outcomes that would be.

Let  $O_1, \ldots, O_m$  be the outcomes of act A. The expected utility of A calculated from probabilities of counterfactuals we shall call  $\mathcal{U}(A)$ ; it is given by the formula

 $\mathscr{U}(A) = \Sigma_i \operatorname{prob} (A \Box \to O_i) \mathscr{D}O_i.$ 

The expected utility of A calculated from conditional probabilities we shall call  $\mathscr{V}(A)$ ; it is given by the formula

$$\mathscr{V}(A) = \Sigma_i \operatorname{prob} \left( O_i | A \right) \mathscr{D} O_i.$$

Perhaps the best mnemonic for distinguishing  $\mathcal{U}$  from  $\mathscr{V}$  is this: we shall be advocating the use of counterfactuals in calculating expected utility, and we shall claim that  $\mathcal{U}(A)$  is the genuine expected utility of A.  $\mathscr{V}(A)$ , we shall claim, measures instead the welcomeness of the news that one is about to perform A. Remember  $\mathscr{V}(A)$ , then, as the value of A as news, and remember  $\mathscr{U}(A)$  as what the authors regard as the genuine expected utility of A.

Now clearly  $\mathcal{U}(A)$  and  $\mathscr{V}(A)$  will be the same if

(2) 
$$prob (A \Box \rightarrow O_i) = prob (O_i | A)$$

for each outcome  $O_j$ . Unless (2) holds for every  $O_j$  such that  $\mathcal{D}O_j \neq 0$ ,  $\mathcal{U}(A)$  and  $\mathscr{V}(A)$  will be the same only by coincidence. We know from Lewis's work (1976) that (2) does not hold for all propositions A and  $O_j$ ; can we expect that (2) will hold for the appropriate propositions?

One assumption, together with the logical truth of Consequence 1, will guarantee that (2) holds for an act and its outcomes. Here and throughout, we suppose that the function *prob* gives the probability ascriptions of an agent who can immediately perform the act in question, and that *prob*  $\phi = 1$  for any logical truth  $\phi$ .

CONDITION 1 on act A and outcome  $O_i$ . The counterfactual  $A \square \rightarrow O_i$  is stochastically independent of the act A. That is to say,

 $prob(A \Box \rightarrow O_i/A) = prob (A \Box \rightarrow O_i).$ 

(Read prob  $(A \Box \rightarrow O_i/A)$  as the conditional probability of  $A \Box \rightarrow O_i$  on A.)

ASSERTION 1. Suppose Consequence 1 is a logical truth. If A and  $O_i$  satisfy Condition 1, and *prob* (A) > O, then

prob  $(A \Box \rightarrow O_i) = prob (O_i/A).^4$ 

*Proof.* Since Consequence 1 is a logical truth, for any propositions P and Q,

prob  $(P \supset [(P \Box \rightarrow Q) \equiv Q]) = 1.$ 

Hence if *prob* P > O, then

prob  $([(P \Box \rightarrow Q) \equiv Q]/P) = 1;$ 

 $\therefore$  prob  $(P \Box \rightarrow Q/P) = prob (Q/P).$ 

From this general truth we have

 $prob (A \Box \rightarrow O_i / A) = prob (O_i / A),$ 

and from this and Condition 1, it follows that

prob  $(A \Box \rightarrow O_i) = prob (O_i/A).$ 

That proves the Assertion.

Condition 1 is that the counterfactuals relevant to decision be stochastically independent of the acts contemplated. Stochastic independence is the same as epistemic independence. For prob  $(A \Box \rightarrow O_i/A)$  is the probability it would be rational for the agent to ascribe to the counterfactual  $A \Box \rightarrow O_i$  on learning A and nothing else – on learning that he was about to perform that act. Thus to say that prob  $(A \Box \rightarrow O_i/A) = prob (A \Box \rightarrow O_i)$  is to say that learning that one was about to perform the act would not change the probability one ascribes to the proposition that if one were to perform the act, outcome  $O_i$  would obtain. We shall use the terms 'stochastic independence' and 'epistemic independence' interchangeably.

The two kinds of expected utility  $\mathcal{U}$  and  $\mathcal{V}$  can also be characterized in a way suggested by Jeffrey's account of the Bayesian model. Let acts  $A_1, \ldots, A_m$  be open to the agent. Let states  $S_1, \ldots, S_n$  partition the possibilities in the following sense. For any propositions  $S_1, \ldots, S_n$ , the truth-function  $aut(S_1, \ldots, S_n)$  will be their exclusive disjunction:  $aut(S_1, \ldots, S_n)$  holds in and only in circumstances where exactly one of  $S_1, \ldots, S_n$  is true. Let the agent know  $aut (S_1, \ldots, S_n)$ . For each act  $A_i$  and state  $S_j$ , let him know that if he did  $A_i$  and  $S_j$  obtained, the outcome would be  $O_{ij}$ . Let him ascribe each outcome  $O_{ij}$  a desirability  $\mathcal{D}O_{ij}$ . This will be a matrix formulation of a decision problem; its defining features are that the agent knows that  $S_1, \ldots, S_n$ , each act open to the agent has a unique outcome. A set  $\{S_1, \ldots, S_n\}$  of states which satisfy these conditions will be called the states of a matrix formulation of the decision problem in question.

Both  $\mathcal U$  and  $\mathcal V$  can be characterized in terms of a matrix formulation:

$$\mathcal{U}(A_i) = \sum_j \operatorname{prob} (A_i \Box \rightarrow S_j) \mathcal{D}O_{ij};$$
  
$$\mathcal{V}(A_i) = \sum_j \operatorname{prob} (S_j | A_i) \mathcal{D}O_{ij}.$$

If  $\mathcal{D}O_{ij}$  can be regarded as the desirability the agent attributes to  $S_j$  on the assumption that' he will do  $A_i$ , then  $\mathscr{V}(A_i)$  is the desirability of  $A_i$  as characterized in the account we quoted from Jeffrey.

On the basis of these matrix characterizations of  $\mathscr{U}$  and  $\mathscr{V}$ , we can state another sufficient condition for the  $\mathscr{U}$ -utility and  $\mathscr{V}$ -utility of an act to be the same.

CONDITION 2 on act  $A_i$ , states  $S_1, \ldots, S_n$ , and the function prob. For each  $A_i$  and  $S_i$ ,

prob 
$$(A_i \Box \rightarrow S_i | A_i) = prob (A_i \Box \rightarrow S_i).$$

ASSERTION 2. Suppose Consequence 1 is a logical truth. If a decision problem satisfies Condition 2 for act  $A_i$ , then  $\mathcal{U}(A_i) = \mathscr{V}(A_i)$ . The proof is like that of Assertion 1.

## 4. ACT-DEPENDENT STATES IN THE SAVAGE FRAMEWORK

Savage's representation of decision problems (1954) is roughly the matrix formulation just discussed. Ignorance is represented as ignorance about which of a number of states of the world obtains. These states are mutually exclusive, and as specific as the problem requires (p. 15). The agent ascribes desirability to 'consequences', or what we are calling *outcomes*. For each act open to the agent, he knows what outcome obtains for each state of the world; if he does not, the problem must be reformulated so that he does. Savage indeed defines an act as a function from states to outcomes (Savage, 1954, p. 14).

It is a consequence of the axioms Savage gives that a rational agent is disposed to choose as if he ascribed a numerical desirability to each outcome and a numerical probability to each state, and then acted to maximize expected utility, where the expected utility of an act A is

(3)  $\Sigma_S \operatorname{prob}(S) \mathscr{D}O(A, S).$ 

(Here O(A, S) is the outcome of act A in state S.) Another consequence of Savage's axioms is the principle of dominance: If for every state S, the outcome of act A in S is more desirable than the outcome of B in S, then A is preferable to B.

Consider this misuse of the Savage apparatus; it is of a kind discussed by Jeffrey (1965, pp. 8-10).

CASE 1. David wants Bathsheba, but since she is the wife of Uriah, he fears that summoning her to him would provoke a revolt. He reasons to himself as follows: 'There are two possibilities: R, that there will be a revolt, and  $\overline{R}$ , that there won't be. The outcomes and their desirabilities are given in Matrix 1,

where B is that I take Bathsheba and A is that I abstain from her. Whether or not there is a revolt, I prefer having Bathsheba to not having her, and so taking Bathsheba dominates over abstaining from her.

$$\begin{array}{c|cccc}
R & \overline{R} \\
\hline
A & R\overline{B}(0) & \overline{R}\overline{B}(9) \\
B & RB(1) & \overline{R}B(10)
\end{array}$$

Matrix 1

This argument is of course fallacious: dominance requires that the states in question be independent of the acts contemplated, whereas taking Bathsheba may provoke revolt. To apply the Savage framework to a decision problem, one must find states of the world which are in some sense actindependent.

We now pursue a suggestion by Jeffrey on how to deal with states that are act-dependent. Construct four new conditionalized<sup>5</sup> states:

- $S_{00}$ : There would be no revolt whatever I did.
- $S_{01}$ : A would not elicit revolt, whereas B would.
- $S_{10}$ : A would elicit revolt, whereas B would not.
- $S_{11}$ : There would be a revolt whatever I did.

If these states hold independently of A and B, we can now work from Matrix 2 without fallacy. Since in Matrix 2 neither row dominates, the decision must be made on the basis of probabilities ascribed to the states  $S_{00}, \ldots, S_{11}$ .

|   | S <sub>00</sub> | S <sub>01</sub> | <i>S</i> <sub>10</sub> | <i>S</i> <sub>11</sub> |  |
|---|-----------------|-----------------|------------------------|------------------------|--|
| A | <i>RB</i> (9)   | <i>RB</i> (9)   | $R\overline{B}(0)$     | $R\overline{B}(0)$     |  |
| B | $\bar{R}B(10)$  | <i>RB</i> (1)   | <i>RB</i> (10)         | <i>RB</i> (1)          |  |
|   |                 | Matrix          | 2                      |                        |  |

What should the probabilities of these states be? One possible answer would be this: Each of the four states  $S_{00}, \ldots, S_{11}$  can be expressed as a conjunction of counterfactuals.  $S_{01}$ , for instance, is the proposition  $(A \square \rightarrow \overline{R})$ &  $(B \square \rightarrow R)$ . The probability of  $S_{01}$ , then, is simply the probability of this proposition, prob ( $[A \square \rightarrow \overline{R}]$  &  $[B \square \rightarrow R]$ ).

The expected utility of an act can now be calculated in the standard way given by (3). The expected utility of A, for instance, will be

(4)  $\Sigma_S \operatorname{prob}(S) \mathcal{D}O(A,S),$ 

where the summation is over the new states S, and O(A, S) is the outcome of A in state S.

Does this procedure give the correct expected utility for the act? What it gives as the expected utility of A, we can show, is  $\mathcal{U}(A)$  – at least that is what it gives if Axiom 2 is part of the logic of counterfactuals. For (4) expands to

$$prob (S_{00}) \mathscr{D}O(A, S_{00}) + prob (S_{01}) \mathscr{D}O(A, S_{01}) + prob (S_{10}) \mathscr{D}O(A, S_{10}) + prob (S_{11}) \mathscr{D}O(A, S_{11}).$$

We have

$$O(A, S_{00}) = O(A, S_{01}) = \overline{R}\overline{B};$$
  

$$O(A, S_{10}) = O(A, S_{11}) = R\overline{B}.$$

Thus since  $S_{00}$  and  $S_{01}$  are mutually exclusive, (4) becomes

 $prob (S_{00} \lor S_{01}) \mathscr{D}\overline{R}\overline{B} + prob (S_{10} \lor S_{11}) \mathscr{D}R\overline{B}.$ 

Now  $S_{00} \vee S_{01}$  is  $([A \Box \rightarrow \overline{R}] \& [B \Box \rightarrow \overline{R}]) \vee [A \Box \rightarrow \overline{R}] \& [B \Box \rightarrow R])$ , and in virtue of the logical truth of Axiom 2, this is  $A \Box \rightarrow \overline{R}$ . Similarly,  $S_{10} \vee S_{11}$ is  $A \Box \rightarrow R$ . Thus (4) becomes

prob 
$$(A \Box \rightarrow \overline{R}) \mathscr{D} \overline{R} \overline{B} + prob (A \Box \rightarrow R) \mathscr{D} R \overline{B},$$

which is  $\mathcal{U}(A)$ . This proof can of course be generalized.

We have considered one way to construct conditionalized states from actdependent states; it is a way that makes use of counterfactuals. Suppose, though, we want to avoid the use of counterfactuals and rely instead on conditional probabilities. Jeffrey, as we understand him, suggests the following: ascribe to each new, conditionalized state the product of the pertinent conditional probabilities. We shall call this probability *prob*\*; thus, for instance,

$$prob^*(S_{01}) = prob(\overline{R}/A) prob(R/B),$$

and corresponding formulas hold for the other new states  $S_{00}$ ,  $S_{10}$ , and  $S_{11}$ .

Using  $prob^*$ , we can again calculate expected utility in the standard way given by (3). The expected utility of A, for instance, will be

(5)  $\Sigma_S \operatorname{prob}^*(S) \mathcal{D}O(A, S),$ 

where again the summation is over the new states  $S_{00}$ ,  $S_{01}$ ,  $S_{10}$ , and  $S_{11}$ . Now (5), it can be shown, has the value  $\mathscr{V}(A)$ . For (5) is the sum of terms

$$prob^* (S_{00}) \mathcal{D}O(A, S_{00}) = prob (\bar{R}/A) prob (\bar{R}/B) \mathcal{D}\bar{R}\bar{B},$$
  
$$prob^* (S_{01}) \mathcal{D}O(A, S_{01}) = prob (\bar{R}/A) prob (R/B) \mathcal{D}\bar{R}\bar{B},$$

 $prob^* (S_{10}) \mathcal{D}O(A, S_{10}) = prob (R/A) prob (\overline{R}/B) \mathcal{D}R\overline{B},$  $prob^* (S_{11}) \mathcal{D}O(A, S_{11}) = prob (R/A) prob (R/B) \mathcal{D}R\overline{B}.$ 

Thus (5) equals

$$[prob (\overline{R}|B) + prob (R|B)] prob (\overline{R}|A) \mathscr{D}\overline{R}\overline{B} + [prob (\overline{R}|B) + prob (R|B)] prob (R|A) \mathscr{D}R\overline{B} = prob (\overline{R}|A) \mathscr{D}\overline{R}\overline{B} + prob (R|A) \mathscr{D}R\overline{B},$$

and this is  $\mathscr{V}(A)$ .

Where, then, a decision problem is misformulated in the Savage framework with act-dependent states, we now have two ways of reformulating the problem with conditionalized states. The first way is to express each conditionalized state as a conjunction of counterfactuals. If the expected utility of an act A is then calculated in the standard manner and Axiom 2 holds, the result is  $\mathcal{U}(A)$ . The second way to reformulate the problem is to ascribe to each new conditionalized state the product of the pertinent conditional probabilities. If the expected utility of an act A is then calculated in the standard manner, the result is  $\mathcal{U}(A)$ . If Axiom 2 holds, then, the two reformulations yield respectively the two kinds of expected utility we have been discussing.

So far we have given the two reformulations only for an example. Here is the way the two methods of reformulation work in general. Let acts  $A_1$ ,  $\ldots$ ,  $A_m$  be open to the agent, let states  $S_1, \ldots, S_n$  not all be act-independent, and for each  $A_i$  and  $S_j$ , let the outcome of act  $A_i$  in  $S_j$  be  $O_{ij}$ . For each possible sequence  $T_1, \ldots, T_m$  consisting of states in  $\{S_1, \ldots, S_n\}$ , there will be a new, conditionalized state  $S(T_1, \ldots, T_m)$ . The outcome of an act  $A_i$  in the new state  $S(T_1, \ldots, T_m)$  will simply be the outcome of  $A_i$  in the old state  $T_i$ . What has been said so far applies to both methods. Now, according to the first method of reformulation, this new state  $S(T_1, \ldots, T_m)$  will be

$$(A_1 \Box \to T_1) \& \ldots \& (A_m \Box \to T_m),$$

and hence, of course, its probability will be the probability of this proposition. According to the second method of reformulation, the probability of new state  $S(T_1, \ldots, T_m)$  will be

$$prob(T_1|A_1) \times \ldots \times prob(T_m|A_m)$$

Once the problem is reformulated, expected utility is to be calculated in the standard way by formula (3).

Are these two ways of reformulating a decision problem equivalent or distinct? They are, of course, equivalent if Axiom 2 hold and  $\mathcal{U}(A_i) = \mathscr{V}(A_i)$ 

for each act  $A_i$ , since the first method yields  $\mathcal{U}(A_i)$  if Axiom 2 holds and the second method yields  $\mathcal{V}(A_i)$ . We already know that Condition 2 and the logical truth of Consequence 1 guarantee that  $\mathcal{U}(A_i) = \mathcal{V}(A_i)$ . Therefore, we may conclude that if Condition 2 holds and Axioms 1 and 2 are logical truths, the two reformulations are equivalent. Condition 2, recall, is that the counterfactuals  $A_i \Box \rightarrow S_j$  are epistemically act-independent: that for each of the old, act-dependent states in terms of which the problem is formulated, learning that one is about to perform a given act will not change the probability one ascribes to the proposition that if one were to perform that act, that state would obtain.

The upshot of the discussion is this. For the Savage apparatus to apply to a decision problem, the states of the decision matrix must be independent of the acts. We have considered two ways of dealing with a problem stated in terms of act-dependent states; both ways involve reformulating the problem in terms of new states which are act-independent. Given the logical truth of Axioms 1 and 2, a sufficient condition for the equivalence of the two reformulations is that the counterfactuals  $A_i \Box \rightarrow S_j$  be epistemically actindependent.

## 5. ACT-DEPENDENT COUNTERFACTUALS

Should we expect Condition 2 to hold? In the case of David, it seems that we should. Suppose David somehow learned that he was about to send for Bathsheba; that would give him no reason to change the probability he ascribes to the proposition 'If I were to send for Bathsheba, there would be a revolt'. Similarly, if David learned that he was about to abstain from Bathsheba, that would give him no reason to change the probability he ascribes to the proposition 'If I were to abstain from Bathsheba, there would be a revolt'. In the case of David, it seems, the pertinent counterfactuals are epistemically act-independent, and hence for each act he can perform, the  $\mathscr{U}$ -utility and the  $\mathscr{V}$ -utility are the same.

When, however, a common factor is believed to affect both behaviour and outcome, Condition 2 may fail, and  $\mathcal{U}$ -utility may diverge from  $\mathscr{V}$ -utility. The following case is patterned after an example used by Stalnaker to make the same point.<sup>6</sup>

CASE 2. Solomon faces a situation like David's, but he, unlike David, has studied works on psychology and political science which teach him the following: Kings have two basic personality types, charismatic and uncharismatic.

A king's degree of charisma depends on his genetic make-up and early childhood experiences, and cannot be changed in adulthood. Now charistmatic kings tend to act justly and uncharismatic kings unjustly. Successful revolts against charismatic kings are rare, whereas successful revolts against uncharismatic kings are frequent. Unjust acts themselves, though, do not cause successful revolts; the reason that uncharismatic kings are prone to successful revolts is that they have a sneaky, ignoble bearing. Solomon does not know whether or not he is charismatic; he does know that it is unjust to send for another man's wife.

Now in this case, Condition 2 fails for states R and  $\overline{R}$ . The counterfactual  $B \Box \rightarrow R$  is not epistemically independent of B: we have

prob 
$$(B \Box \rightarrow R/B) > prob (B \Box \rightarrow R)$$
.

For the conditional probability of anything on *B* is the probability Solomon would rationally ascribe to it if he learned that *B*. Since he knows that *B*'s holding would in no way tend to bring about *R*'s holding, he always ascribes the same probability to  $B \Box \rightarrow R$  as to *R*. Hence both *prob*  $(B \Box \rightarrow R) =$ *prob* (*R*) and *prob*  $(B \Box \rightarrow R/B) = prob$  (*R/B*). Now if Solomon learned that *B*, he would have reason to think that he was uncharismatic, and thus revoltprone. Hence *prob* (*R/B*) > *prob* (*R*), and therefore

(6)  $prob (B \Box \rightarrow R/B) = prob (R/B) > prob (R) = prob (B \Box \rightarrow R).$ 

Here, then, the counterfactual is not epistemically act-independent.

(6) states also that prob  $(B \Box \rightarrow R) < prob (R/B)$ , so that in this case, the probability of the counterfactual does not equal the corresponding conditional probability. By similar argument we could show that  $prob (A \Box \rightarrow R) > prob (R/A)$ . Indeed in this case a  $\mathscr{U}$ -maximizer will choose to send for his neighbour's wife whereas a  $\mathscr{U}$ -maximizer will choose to abstain from her – although we shall need to stipulate the case in more detail to prove the latter.

Consider first *Q*-maximization. We have that

$$\mathcal{U}(B) = prob \ (B \square \rightarrow \overline{R}) \mathcal{D}\overline{R}B + prob \ (B \square \rightarrow R) \mathcal{D}RB;$$
  
$$\mathcal{U}(A) = prob \ (A \square \rightarrow \overline{R}) \mathcal{D}\overline{R}\overline{B} + prob \ (A \square \rightarrow R) \mathcal{D}R\overline{B}.$$

We have argued that  $prob (B \Box \rightarrow R) = prob (R)$ . Similarly,  $prob (A \Box \rightarrow R) = prob (R)$ , and so  $prob (A \Box \rightarrow R) = prob (B \Box \rightarrow R)$ . Likewise  $prob (A \Box \rightarrow \overline{R}) = prob (B \Box \rightarrow \overline{R})$ . We know that  $\mathcal{D}\overline{R}B > \mathcal{D}\overline{R}\overline{B}$  and  $\mathcal{D}RB > \mathcal{D}R\overline{B}$ . Therefore  $\mathcal{U}(B) > \mathcal{U}(A)$ . This is in effect an argument from dominance, as we shall discuss in Section 8.

Now consider  $\mathscr{V}$ -maximization. Learning that A would give Solomon reason to think he was charismatic and thus not-revolt prone, whereas learning that B would give him reason to think that he was uncharismatic and revolt-prone. Thus prob (R/B) > prob (R/A). Suppose the difference between these probabilities is greater than 1/9, so that where  $prob (R/A) = \alpha$  and prob  $(R/B) = \alpha + \epsilon$ , we have  $\epsilon > 1/9$ . From Matrix 1, we have

$$\begin{aligned} \mathscr{V}(A) &= prob \ (\bar{R}/A) \mathscr{D}\bar{R}\bar{B} + prob \ (R/A) \mathscr{D}R\bar{B} = 9(1-\alpha) + 0. \\ \mathscr{V}(B) &= prob \ (\bar{R}/B) \mathscr{D}\bar{R}B + prob \ (R/B) \mathscr{D}RB \\ &= 10(1-\alpha-\epsilon) + 1(\alpha+\epsilon). \end{aligned}$$

Therefore  $\mathscr{V}(A) - \mathscr{V}(B) = 9\epsilon - 1$ , and since  $\epsilon > 1/9$ , this is positive. We have shown that if  $\epsilon > 1/9$ , then although  $\mathscr{U}(B) > \mathscr{U}(A)$ , we have  $\mathscr{V}(A) > \mathscr{V}(B)$ . Thus  $\mathscr{U}$ -maximization and  $\mathscr{V}$ -maximization in this case yield conflicting prescriptions.

Which of these prescriptions is the rational one? It seems clear that in this case it is rational to perform the  $\mathcal{U}$ -maximizing act: unjustly to send for the wife of his neighbor. For Solomon cares only about getting the woman and avoiding revolt. He knows that sending for the woman would not cause a revolt. To be sure, sending for her would be an indication that Solomon lacked charisma, and hence an indication that he will face a revolt. To abstain from the woman for this reason, though, would be knowingly to bring about an indication of a desired outcome without in any way bringing about the desired outcome itself. That seems clearly irrational.

For those who find Solomon too distant in time and place or who mistrust charisma, we offer the case of Robert Jones, rising young executive of International Energy Conglomerate Incorporated. Jones and several other young executives have been competing for a very lucrative promotion. The company brass found the candidates so evenly matched that they employed a psychologist to break the tie by testing for personality qualities that lead to long run successful performance in the corporate world. The test was administered to the candidates on Thursday. The promotion decision is made on the basis of the test and will be announced on Monday. It is now Friday. Jones learns, through a reliable company grapevine, that all the candidates have scored equally well on all factors except ruthlessness and that the promotion will go to whichever of them has scored highest on this factor, but he cannot find out which of them this is.

On Friday afternoon Jones is faced with a new problem. He must decide whether or not to fire poor old John Smith, who failed to meet his sales quota this month because of the death of his wife. Jones believes that Smith will come up to snuff after he gets over his loss provided that he is treated leniently, and that he can convince the brass that leniency to Smith will benefit the company. Moreover, he believes that this would favorably impress the brass with his astuteness. Unfortunately, Jones has no way to get in touch with them until after they announce the promotion on Monday.

Jones knows that the ruthlessness factor of the personality test he has taken accurately predicts his behaviour in just the sort of decision he now faces. Firing Smith is good evidence that he has passed the test and will get the promotion, while leniency is good evidence that he has failed the test and will not get the promotion. We suppose that the utilities and probabilities correspond to those facing Solomon.  $\mathscr{V}$ -maximizing recommends firing Smith, while  $\mathscr{U}$ -maximizing recommends leniency. Firing Smith would produce evidence that Jones will get his desired promotion. It seems clear, however, that to fire Smith for this reason despite the fact that to do so would in no way help to bring about the promotion and would itself be harmful, is irrational.

# 6. The significance of $\mathscr{U}$ and $\mathscr{V}$

From the Solomon example, it should be apparent that the  $\mathscr{V}$ -utility of an act is a measure of the welcomeness of the news that one is about to perform that act. Such news may tend to be welcome because the act is likely to have desirable consequences, or tend to be unwelcome because the act is likely to have disagreeable consequences. Those, however, are not the only reasons an act may be welcome or unwelcome: an act may be welcome because its being performed is an indication that the world is in a desired state. Solomon, for instance, would welcome the news that he was about to abstain from his neighbor's wife, but he would welcome it not because he thought just acts any more likely to have desirable consequences than unjust acts, but because he takes just acts to be a sign of charisma, and he thinks that charisma may bring about a desired outcome.

 $\mathscr{U}$ -utility, in contrast, is a measure of the expected efficacy of an act in bringing about states of affairs the agent desires; it measures the expected value of the consequences of an act. That can be seen in the case of Solomon. The  $\mathscr{U}$ -utility of sending for his neighbor's wife is greater than that of abstaining, and that is because he knows that sending for her will bring about a consequence he desires — having the woman — and he knows that it will not bring about any consequences he wishes to avoid: in particular, he knows that it will not bring about a revolt.

166

What is it for an act to bring about a consequence? Here are two possible answers, both formulated in terms of counterfactuals.

In the first place, roughly following Sobel (1970, p. 400) we may say that act A brings about state S if  $A \square \rightarrow S$  holds, and for some alternative  $A^*$  to  $A, A^* \square \rightarrow S$  does not hold.<sup>7</sup> (An *alternative* to A is another act open to the agent on the same occasion). Now on this analysis, the **U**-utility of an act as we have defined it is the sum of the expected value of its consequences plus a term which is the same for all acts open to the agent on the occasion in question; this latter term is the expected value of unavoidable outcomes. A state S is *unavoidable* iff for every act  $A^*$  open to the agent,  $A^* \square \rightarrow S$  holds. Thus  $A \square \rightarrow S$  holds iff S is a consequence of A or S is unavoidable. Hence in particular, for any outcome O,

$$prob (A \Box \rightarrow O) = prob (O \text{ is a consequence of } A) + prob (O \text{ is unavoidable}),$$

and so we have

$$\mathcal{U}(A) = \sum_{O} prob \ (A \square \rightarrow O) \mathcal{D}O$$
  
=  $\sum_{O} prob \ (O \text{ is a consequence of } A) \mathcal{D}O$   
+  $\sum_{O} prob \ (O \text{ is unavoidable}) \mathcal{D}O.$ 

The first term is the expected value of the consequences of A, and the second term is the same for all acts open to the agent. Therefore on this analysis of the term 'consequence',  $\mathcal{U}$ -utility is maximal for the act or acts whose consequences have maximal expected value.

Here is a second possible way of analyzing what it is to be a consequence. When an agent chooses between two acts A and B, what he really needs to know is not what the consequences of A are and what the consequences of Bare, but rather what the consequences are of A as opposed to B and vice versa. Thus for purposes of decision-making, we can do without an analysis of the clause 'S is a consequence of A', and analyze instead the clause 'S is a consequence of A as opposed to B'. This we can analyze as

$$(A \Box \to S) \& \sim (B \Box \to S).$$

Now on this analysis,  $\mathcal{U}(A) > \mathcal{U}(B)$  iff the expected value of the consequences of A as opposed to B exceeds the expected value of the consequences of B as opposed to A. For any state S,  $A \square \rightarrow S$  holds iff either S is a consequence of A as opposed to B or  $(A \square \rightarrow S) \& (B \square \rightarrow S)$  holds.

Thus

$$\begin{aligned} &\mathcal{U}(A) = \sum_{O} prob \ (A \square \to O) \mathcal{D}O \\ &= \sum_{O} prob \ (O \text{ is a consequence of } A \text{ as opposed to } B) \mathcal{D}O \\ &+ \sum_{O} prob \ ([A \square \to O]] \& [B \square \to O]) \mathcal{D}O \\ &\mathcal{U}(B) = \sum_{O} prob \ (O \text{ is a consequence of } B \text{ as opposed to } A) \mathcal{D}O \\ &+ \sum_{O} prob \ ([A \square \to O]] \& [B \square \to O]) \mathcal{D}O. \end{aligned}$$

The second term is the same in both cases, and so  $\mathcal{U}(A) > \mathcal{U}(B)$  iff

 $\Sigma_O \operatorname{prob} (O \text{ is a consequence of } A \text{ as opposed to } B) \mathcal{D}O > \Sigma_O \operatorname{prob} (O \text{ is a consequence of } B \text{ as opposed to } A) \mathcal{D}O.$ 

The left side is the expected value of the consequences of A as opposed to B; the right side is the expected value of the consequences of B as opposed to A. Thus for any pair of alternatives, to prefer the one with the higher  $\mathcal{U}$ -utility is to prefer the one the consequences of which as opposed to the other have the greater expected value.

We can now ask whether  $\mathscr{U}$  or  $\mathscr{V}$  is more properly called the 'utility' of an act. The answer seems clearly to be  $\mathscr{U}$ . The 'utility' of an act should be its expected genuine efficacy in bringing about states of affairs the agent wants, not the degree to which news of the act ought to cheer the agent. Since  $\mathscr{U}$ -utility is a matter of what the act can be expected to bring about whereas  $\mathscr{V}$ -utility is a matter of the welcomeness of news,  $\mathscr{U}$ -utility seems best to capture the notion of utility.

Jeffrey (1965, pp. 73-4) writes, 'If the agent is deliberating about performing act A or act B, and if AB is impossible, there is no effective difference between asking whether he prefers A to B as a news item or as an act, for he makes the news'. It should now be clear why it may sometimes be rational for an agent to choose an act B instead of an act A, even though he would welcome the news of A more than that of B. The news of an act may furnish evidence of a state of the world which the act itself is known not to produce. In that case, though the agent indeed makes the news of his act, he does not make all the news his act bespeaks.

## 7. Two sure thing principles

CASE 3. Upon his accession to the throne, Reoboam wonders whether to announce that he will reign severely or to announce that he will reign leniently. He will be bound by what he announces. He slightly prefers a short severe reign to a short lenient reign, and he slightly prefers a long severe reign to a long lenient reign. He strongly prefers a long reign of any kind to a

168

short reign of any kind. Where L is that he is lenient and D, that he is deposed early, his utilities are as in the Matrix 3.

$$\begin{array}{c|c}
D & \overline{D} \\
\hline L & 0 & 80 \\
\overline{L} & 10 & 100 \\
\hline
Matrix 3
\end{array}$$

The wise men of the kingdom give him these findings of behavioural science: There is no correlation between a king's severity and the length of his reign. Severity, nevertheless, often causes early deposition. The reason for the lack of correlation between severity and early deposition is that on the one hand, charismatic kings tend to be severe, and on the other hand, lack of charisma tends to elicit revolts. A king's degree of charisma cannot be changed in adulthood. There is at present no indication of whether Reoboam is charismatic or not.

These findings were based on a sample of 100 kings, 48 of whom had their reigns cut short by revolt. On post mortem examination of the pineal gland, 50 were found to have been charismatic and 50 uncharismatic. 80% of the charismatic kings had been severe and 80% of the uncharismatic kings had been lenient. Of the charismatic kings, 40% of those who were severe were deposed whereas only 20% of those who were lenient were deposed. Of the uncharismatic kings, 80% of those who were severe were deposed whereas only 55% of those who were lenient were deposed. The totals were as in Table 1. This is Reoboam's total evidence on the subject.<sup>8</sup>

|         | Charismatic      | Uncharismatic    | Total            |
|---------|------------------|------------------|------------------|
| Severe  | 16 deposed (40%) | 8 deposed (80%)  | 24 deposed (48%) |
|         | 24 long-reigned  | 2 long-reigned   | 26 long reigned  |
| Lenient | 2 deposed (20%)  | 22 deposed (55%) | 24 deposed (48%) |
|         | 8 long-reigned   | 18 long-reigned  | 26 long reigned  |

TABLE 1

Reoboam's older advisors argue from a sure thing principle. There are two possibilities, they say: that Reoboam is charismatic and that he is uncharismatic; what he does now will not affect his degree of charisma. On the assumption that he is charismatic, it is rational to prefer lenience. For since 40% of severe charismatic kings are deposed, the expected utility of severity in that case would be

$$0.4\mathscr{D}SD + 0.6\mathscr{D}S\overline{D} = 0.4 \times 10 + 0.6 \times 80 = 52$$

whereas since only 20% of lenient charismatic kings are deposed, the expected utility of lenience in that case would be

$$0.2\mathscr{D}LD + 0.8\mathscr{D}L\overline{D} = 0.2 \times 0 + 0.8 \times 10 = 64.$$

On the assumption that he is uncharismatic, it is again rational to prefer lenience. For since 80% of severe uncharismatic kings are deposed, the expected utility of severity in this case would be

$$0.8\mathscr{D}SD + 0.2\mathscr{D}S\overline{D} = 0.8 \times 10 + 0.2 \times 100 = 28$$

whereas since only 55% of lenient uncharismatic kings are deposed, the expected utility of lenience in this case would be

$$0.55\mathcal{D}LD + 0.45\mathcal{D}L\overline{D} = 0.55 \times 0 + 0.45 \times 80 = 36.$$

Thus in either case, lenience is to be preferred, and so by a sure thing principle, it is rational to prefer lenience in the actual case.

Reoboam's youthful friends argue that on the contrary, sure thing considerations prescribe severity. Severity is indeed the dominant strategy. There are two possibilities: D, that Reoboam will be deposed, and  $\overline{D}$ , that he will not be. These two states are stochastically independent of the acts contemplated: both *prob* (D/S) and *prob* (D/L) are 0.48. Therefore, his youthful friends urge, one can without fallacy use the states D and  $\overline{D}$  in an argument from dominance. On the assumption that he will be deposed, he prefers to be severe, and likewise on the assumption that he will not be deposed, he prefers to be severe. Thus by dominance, it is rational for him to prefer severity.

Here, then, are two sure thing arguments which lead to contrary prescriptions. One argument appeals to the finding that charisma is causally independent of the acts contemplated; the other appeals to the finding that being deposed is stochastically independent of the acts. The old advisors and youthful companions are in effect appealing to different versions of a sure thing principle, one of which requires causal independence and the other of which requires stochastic independence. The two versions lead to incompatible conclusions.

The sure thing principle is this: if a rational agent knows *aut*  $(S_1, \ldots, S_n)$  and prefers A to B in each case, then he prefers A to B. If the propositions

 $S_1, \ldots, S_n$  are required to be states in a matrix formulation of the decision problem, so that each pair of state and act determine a unique outcome, the sure thing principle becomes the principle of dominance to be discussed in Section 8; the principle of dominance is thus a special case of the sure thing principle. Now the principle of dominance, we have said, requires a proviso that the states in question be act-independent. The sure thing principle should presumably include the same proviso. The sure thing principle, then, should be this: If a rational agent knows that precisely one of the propositions  $S_1, \ldots, S_n$  holds and prefers act A to act B in each case, and if in addition the propositions  $S_1, \ldots, S_n$  are independent of the acts A and B, then he prefers A to B.

The problem in the case of Reoboam is that his two groups of advisors appeal to different kinds of independence to reach opposing conclusions. The older advisors appeal to causal independence; they cite the finding that a king's degree of charisma is unaffected by his adult actions. His youthful companions appeal to stochastic independence; they cite the finding that there is no correlation between severity in kings and revolt. The two appeals yield opposite conclusions.

It seems, then, that the sure thing principle comes in two different versions, one of which requires that the propositions in question be causally independent of the acts, and the other of which requires the propositions to be stochastically independent of the acts.

The principle to which the youthful companions appeal can be put as follows.

DEFINITION. Act A is sure against act B with stochastic independence of  $S_1, \ldots, S_n$  iff the following hold. The agent knows that independently of the choice between A and B, propositions  $S_1, \ldots, S_n$  partition the possibilities; that is to say, prob (aut  $(S_1, \ldots, S_n)/A) = 1$  and prob (aut  $(S_1, \ldots, S_n)/B) = 1$ . The propositions  $S_1, \ldots, S_n$  are epistemically independent of the choice between A and B, in the sense that for each, prob  $(S_i/A) = prob (S_i/B)$ . Finally, for each of these propositions  $S_i$  it would be rational to prefer A to B if it were known that  $S_i$  held.

Sure-thing with Stochastic Independence. If act A is sure against act B with stochastic independence, then it is rational to prefer A to B.

The principle to which the older advisors appeal will take longer to formulate. The proviso for this version will be that the propositions  $S_1$ , ...,  $S_n$  be causally independent of the choice between A and B; this can be

formulated in terms of counterfactuals. To say that a state  $S_i$  is causally independent of the choice between A and B is to say that  $S_i$  would hold if A were performed iff  $S_i$  would hold if B were performed:  $(A \Box \rightarrow S_i) \equiv (B \Box \rightarrow S_i)$ . We now want to suppose that for each state  $S_i$ , A would be preferred to B given, in some sense, knowledge of  $S_i$ . This knowledge of  $S_i$  should not simply be knowledge that  $S_i$  holds, but knowledge that  $S_i$  holds independently of the choice between A and B: that  $(A \Box \rightarrow S_i) \& (B \Box \rightarrow S_i)$ . We can now state the principle.

For each  $S_i$  let  $S_i^*$  be  $(A \Box \rightarrow S_i) \& (B \Box \rightarrow S_i)$ .

DEFINITION. A is sure against B with causal independence of  $S_1, \ldots, S_n$  iff the following hold. The agent knows  $aut(S_1^*, \ldots, S_n^*)$ , and for each  $S_i$  it would be rational to prefer A to B if  $S_i^*$  were known to hold.<sup>9</sup>

(Note that since for each  $S_i$ ,  $(A \Box \rightarrow S_i) \equiv (B \Box \rightarrow S_i)$  follows from *aut*  $(S_i^*, \ldots, S_n^*)$ , this guarantees that our agent knows that each  $S_i$  is causally independent of the choice between A and B.) We can now state the principle to which the older advisors appeal.

Sure-thing with Causal Independence. If A is sure against B with causal independence, then it is rational to prefer A to B.

In the case of Reoboam, we have seen, Sure-thing with Stochastic Independence prescribes severity and Sure-thing with Causal Independence prescribes lenience. Now to us it seems clear that the only rational action in this case is that prescribed by Sure-thing with Causal Independence. It is rational for Reoboam to prefer lenience because severity tends to bring about deposition and he wants not to be deposed much more strongly than he wants to be severe. To be guided by Sure-thing with Stochastic Independence in this case is to ignore the finding that severity tends to bring about revolt — to ignore that finding simply because severity is not on balance a *sign* that revolt will occur. To choose to be severe is to act in a way that tends to bring about a dreaded consequence, simply because the act is not a sign of the consequence. That seems to us to be irrational.

The two versions of the sure thing principle we have discussed correspond to the two kinds of utility discussed earlier. Sure Thing with Stochastic Independence follows from the principle that an act is rationally preferred to another iff it has greater  $\mathscr{V}$ -utility, whereas Sure Thing with Causal Independence follows from the principle that an act is rationally preferred to another iff it maximizes  $\mathscr{U}$ -utility. ASSERTION. Suppose that in any possible situation, it is rational to prefer an act A to an act B iff the  $\mathcal{U}$ -utility of A is greater than that of B. Then Sure Thing with Causal Independence holds.

**Proof.** Suppose A is sure against B with causal independence of  $S_1, \ldots, S_n$ , and that in any possible circumstance, it would be rational to prefer A to B iff A's  $\mathcal{U}$ -utility were greater than B's. The Assertion will be proved if we show from these assumptions that  $\mathcal{U}(A) > \mathcal{U}(B)$ .

Since A is sure against B with causal independence of  $S_1, \ldots, S_n$ , for each  $S_i$  it would be rational to prefer A to B if  $S_i^*$  were known to hold. Therefore if  $S_i^*$  were known to hold, the  $\mathcal{U}$ -utility of A would be greater than that of B. Now the  $\mathcal{U}$ -utility that A would have if  $S_i^*$  were known is

 $\Sigma_O \operatorname{prob} (A \Box \rightarrow O/S_i^*) \mathcal{D}O.$ 

Call this  $\mathcal{U}_i^*(A)$ , and define  $\mathcal{U}_i^*(B)$  in a like manner. We have supposed that for each  $S_i, \mathcal{U}_i^*(A) > \mathcal{U}_i^*(B)$ .

Now by definition of the function  $\mathcal{U}$ ,

$$\mathcal{U}(A) = \Sigma_O \operatorname{prob} (A \Box \to O) \mathcal{D}O$$

Since A is sure against B with causal independence of  $S_1, \ldots, S_n$ , it is known that *aut*  $(S_1^*, \ldots, S_i^*)$  holds. By the probability calculus, then, for each outcome O

$$prob \ (A \square \to O) = \Sigma_i \ prob \ (A \square \to O/S_i^*) \ prob/S_i^*.$$

Therefore

$$\begin{aligned} \boldsymbol{\mathscr{U}}(A) &= \sum_{O} [\sum_{i} \operatorname{prob} (A \square \to O/S_{i}^{*}) \operatorname{prob} S_{i}^{*}] \mathcal{D}O \\ &= \sum_{i} \operatorname{prob} S_{i}^{*} [\sum_{O} \operatorname{prob} (A \square \to O/S_{i}^{*}) \mathcal{D}O] , \\ &= \sum_{i} \boldsymbol{\mathscr{U}}_{i}^{*}(A) \operatorname{prob} S_{i}^{*}. \end{aligned}$$

By a like argument,

$$\mathscr{U}(B) = \Sigma_i \mathscr{U}_i^*(B) \operatorname{prob} S_i^*.$$

Since for each  $S_i, \mathcal{U}_i^*(A) > \mathcal{U}_i^*(B)$ , it follows that  $\mathcal{U}(A) > \mathcal{U}(B)$ , and the Assertion is proved.

ASSERTION. Suppose that in any possible circumstance, it is rational to prefer an act A to an act B iff the  $\mathscr{V}$ -utility of A is greater than that of B. Then Sure Thing with Stochastic Independence holds.

**Proof.** Suppose A is sure against B with causal independence of  $S_1, \ldots, S_n$ , and that in any possible circumstances, it would be rational to prefer A

to B iff A's  $\mathscr{V}$ -utility is greater than B's. The Assertion will be proved if we show from these assumptions that  $\mathscr{V}(A) > \mathscr{V}(B)$ .

Now since A is sure against B with stochastic independence of  $S_1, \ldots, S_n$ , for each  $S_i$  it would be rational to prefer A to B if  $S_i$  were known to hold. Therefore, if  $S_i$  were known to hold, then the  $\mathscr{V}$ -utility of A would be greater than that of B. Now the  $\mathscr{V}$ -utility that A would have if  $S_i$  were known to hold is

 $\Sigma_O prob (O|AS_i) \mathcal{D}O.$ 

Call this  $\bar{\mathscr{V}}_i^*(A)$ , and define  $\mathscr{V}_i^*(B)$  correspondingly. We have that for each  $S_i$ ,  $\mathscr{V}_i^*(A) > \mathscr{V}_i^*(B)$ . Now by definition of the function  $\mathscr{V}$ ,

 $\mathscr{V}(A) = \Sigma_O prob (O|A) \mathscr{D}O.$ 

Since A is sure against B with stochastic independence of  $S_1, \ldots, S_n$ , we have prob (aut  $(S_1, \ldots, S_n)/A$ ) = 1, and so by the probability calculus, for each O,

$$\mathscr{V}$$
 prob  $(O|A) = \Sigma_i$  prob  $(O|AS_i)$  prob  $(S_i|A)$ .

Hence

$$\begin{aligned} \Psi(A) &= \sum_{O} [\sum_{i} \operatorname{prob} (O|AS_{i}) \operatorname{prob} (S_{i}|A)] \mathcal{D}O \\ &= \sum_{i} \operatorname{prob} (S_{i}|A) [\sum_{O} \operatorname{prob} (O|AS_{i}) \mathcal{D}O] \\ &= \sum_{i} \operatorname{prob} (S_{i}|A) \mathcal{\Psi}_{i}^{*}(A). \end{aligned}$$

By a like argument,

$$\mathscr{V}(B) = \Sigma_i \operatorname{prob} (S_i/B) \mathscr{V}_i^*(B).$$

Since for each  $S_i$ , prob  $(S_i|A) = prob(S_i|B)$  and  $\mathscr{V}_i^*(A) > \mathscr{V}_i^*(B)$  it follows that  $\mathscr{V}(A) > \mathscr{V}(B)$ , and the Assertion is proved.

## 8. Two kinds of dominance

We have said that the principle of dominance is the sure thing principle restricted to a special case, and that the sure thing principle has two versions, one of which holds for  $\mathcal{U}$ -maximization and the other for  $\mathscr{V}$ -maximization. There should, then, be two versions of the principle of dominance, one for each kind of utility maximization. The principles can be formulated as follows.

DEFINITION. Let  $S_1, \ldots, S_n$  be the states of a standard decision matrix, and let A and B be acts. Then A strongly dominates B with respect to

 $S_1, \ldots, S_n$  if for each  $S_i$ , the outcome of A in  $S_i$  is more desirable than the outcome of B in  $S_i$ .

Principle of Dominance with Causal Independence. Suppose act A strongly dominates act B with respect to states  $S_1, \ldots, S_n$ . If for each state  $S_i$ , the agent knows that  $(A \Box \rightarrow S_i) \equiv S_i$  and  $(B \Box \rightarrow S_i) \equiv S_i$ , then it is rational for him to prefer A to B.

Principle of Dominance with Stochastic Independence. Suppose act A strongly dominates act B with respect to states  $S_1, \ldots, S_n$ . If for each state  $S_i$ , prob  $(S_i|A) = prob (S_i) = prob (S_i|B)$ , then it is rational for him to prefer A to B.

The Principle of Dominance with Causal Independence holds if rationality requires maximization of  $\mathcal{U}$ , and the Principle of Dominance with Stochastic Independence holds if rationality requires maximization of  $\mathscr{V}$ .<sup>10</sup>

Although these two principles are respective consequences of two principles of expected utility maximization which may conflict, they cannot themselves conflict. For suppose A strongly dominates B with respect to some set of states  $S_1, \ldots, S_n$ . Then the worst outcome of A is more desirable than some outcome of B. For the worst outcome of A is the outcome of A in some state  $S_i$ , and since A strongly dominates B with respect to  $S_1, \ldots, S_n$ , the outcome of A is more desirable than the outcome of A is more desirable than the vorst outcome of B. It cannot be the case, then, that B strongly dominates A with respect to some other set of states  $T_1, \ldots, T_n$ . For if that indeed were the case, then, we have seen, the worst outcome of B would be more desirable than the worst outcome of A. We have seen that if A strongly dominates B with respect to a set of states, then there is no set of states with respect to which B strongly dominates A. For that reason, the two principles of dominance we have stated will never yield conflicting prescriptions for a simple decision problem.

In a weaker form, however, dominance indeed can be exploited to yield conflicting prescriptions.

DEFINITION. Let  $S_1, \ldots, S_n$  be the states of a standard decision matrix, and let A and B be acts. A weakly dominates B with respect to  $S_1, \ldots, S_n$ iff for each state  $S_i$ , the outcome of A in  $S_i$  is at least as desirable as the outcome of B in  $S_i$ , and for some state  $S_i$  with prob  $(S_i) > O$ , the outcome of A in  $S_i$  is more desirable than the outcome of B in  $S_i$ . We now get two Principles of Weak Dominance by substituting 'weakly dominates' for 'strongly dominates' in the two Principles of Dominance stated above.

CASE 4. A subject is presented with two boxes, one to the left and one to the right. He must choose between two acts:

- $A_L$  Take the box on the left.
- $A_R$  Take the box on the right.

The experimenter has already done one of the following.

- $M_{11}$  Place a million dollars in each box.
- $M_{01}$  Place a million dollars in the box on the right and nothing in the box on the left.
- $M_{00}$  Place nothing in either box.

He has definitely not placed money in the left box without placing money in the right box. Now the experimenter has predicted the behavior of the subject, and before making his prediction, he has used a random device to select one of the following three strategies.

- (i) Reward choice of left box:  $M_{11}$  if  $A_L$  is predicted;  $M_{00}$  if  $A_R$  is predicted.
- (ii) Ensure payment:  $M_{11}$  if  $A_L$  is predicted;  $M_{01}$  if  $A_R$  is predicted.
- (iii) Ensure non-payment:  $M_{01}$  if  $A_L$  is predicted;  $M_{00}$  if  $A_R$  is predicted.

The subject knows all this, and believes in the accuracy of the experimenter's predictions with complete certainty.

The Principle of Weak Dominance with Causal Independence prescribes taking the box on the right. The three states  $M_{11}$ ,  $M_{01}$ , and  $M_{00}$  are causally independent of the act the subject performs. The possible outcomes are shown in the table, where 1 is getting the million dollars and 0 is not getting it.

|       | <i>M</i> <sub>11</sub> | <i>M</i> <sub>01</sub> | $M_{00}$ |  |
|-------|------------------------|------------------------|----------|--|
| $A_L$ | 1                      | 0                      | 0        |  |
| $A_R$ | 1                      | 1                      | 0        |  |

 $M_{01}$  has non-zero probability, since if  $A_L$  was predicted it would result from the experimenter's using strategy (iii) and if  $A_R$  was predicted, it would result from the experimenter's using strategy (ii). Thus  $A_R$  weakly dominates

 $A_L$  with respect to  $M_{11}$ ,  $M_{01}$ ,  $M_{00}$ , and the Principle of Weak Dominance with Causal Independence prescribes taking the box on the right.

The Principle of Weak Dominance with Stochastic Independence, in contrast, prescribes taking the box on the left.

The possibilities can be partitioned as follows:

 $S_1$  the experimenter predicts correctly and follows strategy (i).

 $S_2$   $S_1$  does not hold and the subject wins a million dollars.

 $S_3$   $S_1$  does not hold and the subject wins nothing.

The payoffs are given in the table.

|       | $S_1$ | $S_2$ | $S_3$ |  |
|-------|-------|-------|-------|--|
| $A_L$ | 1     | 1     | 0     |  |
| $A_R$ | 0     | 1     | 0     |  |

Now prob  $(S_1) \neq 0$ , and hence  $A_L$  weakly dominates  $A_R$  with respect to  $S_1, S_2, S_3$ . Moreover, the states  $S_1, S_2$ , and  $S_3$  are stochastically independent of  $A_L$  and  $A_R$ . For the subject knows that the experimenter has selected his strategy independently of his prediction, by means of a random device; hence learning that he was about to perform  $A_L$ , say, would not affect the probability he ascribes to the experimenter's having had any given strategy. By the subject's probability function, then, which strategy the experimenter has used is stochastically independent of the subject's act. Now the subject believes that the experimenter has predicted correctly and used strategy (i), (ii), or (iii). Hence he thinks that  $S_1$  holds iff the experimenter has used strategy (ii), and that  $S_3$  holds if the experimenter has used strategy (iii). Hence under his probability function, states  $S_1, S_2$ , and  $S_3$  are stochastically independent of  $A_L$  and  $A_R$ . Thus the Principle of Weak Dominance with Stochastic Independence applies, and it prescribes taking the box on the left.

Some readers may object in Case 4 to the subject's complete certainty that the experimenter has predicted correctly. It is possible to construct a conflict between the two principles of weak dominance without requiring such certainty, but the example becomes more complicated.

CASE 5. Same as Case 4, except for the following.

The subject ascribes a probability of 0.8 to the experimenter's having predicted correctly, and this probability is independent of the subject's choice of  $A_L$  or  $A_R$ . Thus where C is 'the experimenter has predicted correctly',

 $prob(C/A_L) = 0.8$  and  $prob(C/A_R) = 0.8$ .

The experimenter has chosen among the following three strategies by means of a random device.

- (i)  $M_{11}$  if  $A_L$  is predicted;  $M_{00}$  if  $A_R$  is predicted.
- (ii\*)  $M_{11}$  if  $A_L$  is predicted;  $M_{01}$  or  $M_{00}$ , with equal probability, if  $A_R$  is predicted.
- (iii\*)  $M_{11}$  or  $M_{01}$ , with equal probability, if  $A_L$  is predicted;  $M_{00}$  if  $A_R$  is predicted.

He has followed (i) with a probability 0.5, (ii\*) with a probability 0.25, and (iii\*) with a probability 0.25.

In Case 5, as in Case 4, the states  $M_{11}$ ,  $M_{01}$ , and  $M_{00}$ , are causally independent of the acts  $A_R$  and  $A_L$ , and from the Principle of Weak Dominance with Causal Independence and the facts of the case, it follows that it is rational to prefer  $A_R$  to  $A_L$ .

Now let states  $S_1$ ,  $S_2$ , and  $S_3$  be as before:  $S_1$  is that the experimenter predicts correctly and follows strategy (i);  $S_2$  is that  $S_1$  does not hold and the subject receives a million dollars;  $S_3$  is that  $S_1$  does not hold nd the subject receives nothing. As in Case 4, if  $S_1$ ,  $S_2$ , and  $S_3$  are stochastically independent of  $A_L$  and  $A_R$ , then from the Principle of Weak Dominance with Stochastic Independence and the facts of the case, it follows that it is rational to prefer  $A_L$  to  $A_R$ . It is clear that  $S_1$  is stochastically independent of the acts  $A_L$  and  $A_R$ ; we now show that  $S_2$  and  $S_3$  are as well: that prob  $(S_2/A_L) = prob$  $(S_2/A_R)$  and prob  $(S_3/A_L) = prob (S_3/A_R)$ .

There are two possible acts, two possible experimenter's predictions, and three possible experimenter's strategies, some of which may involve the flip of a coin. Call a combination of act, prediction, experimenter's strategy, and result of coin flip if it matters, a *case*. For each case, the Table 2 shows.

- (1) The state  $M_{11}$ ,  $M_{01}$ , or  $M_{00}$  which would hold in that case.
- (2) The conditional probability of the case given the act.
- (3) The outcome in that case: 1 for getting the million dollars, 0 for not.
- (4) The state  $S_1$ ,  $S_2$ , or  $S_3$  which holds in that case.

The conditional probability prob  $(S_2/A_L)$  is then obtained by adding up the conditional probabilities given  $A_L$  of cases in which  $S_2$  holds; a like procedure gives prob  $(S_3/A_L)$ , prob  $(S_2/A_R)$ , and prob  $(S_3/A_R)$ .

The conclusion of Table 2 is that the states  $S_1$ ,  $S_2$ , and  $S_3$  are indeed

epistemically independent of the acts  $A_L$  and  $A_R$ . Since  $A_L$  weakly dominates  $A_R$  with respect to states  $S_1$ ,  $S_2$ , and  $S_3$ , it follows that  $A_L$  weakly dominates  $A_R$  with respect to stochastically independent states. We already know that  $A_R$  weakly dominates  $A_L$  with respect to causally independent states  $M_{11}$ ,  $M_{01}$ , and  $M_{00}$ . In Case 5, then, the two principles of weak dominance are in conflict.

|                            | $A_L$ Performed                                |                         |   | $\bar{A_R}$ Performed                              |   |  |                                       |                              |
|----------------------------|--|-------------------------|---|--|---|--|---------------------------------------|------------------------------|
|                            | $\begin{array}{c} - \\ A_L \\ 0.8 \end{array}$ | redicted                | A <sub>R</sub> P<br>0.2                   | redicted   | $\begin{array}{c} A_L \\ 0.2 \end{array}$ | redicted   | <i>A<sub>R</sub></i> P 0.8            | redicted                     |
| Strategy<br>(i)<br>0.5     | <i>M</i> <sub>11</sub><br>1                    | 0.4<br>S <sub>1</sub>   | М <sub>00</sub><br>О                      | 0.1<br>S <sub>3</sub>                              | <i>M</i> <sub>11</sub><br>1               | 0.1<br>S <sub>2</sub>  | <i>М</i> <sub>00</sub><br>0           | 0.4<br>S <sub>1</sub>        |
| Strategy<br>(ii*)<br>0.25  | <i>M</i> <sub>11</sub><br>1                    | 0.2<br>S <sub>2</sub>   | $ \frac{M_{01}}{0} $ $ \frac{M_{00}}{0} $ | 0.025<br>S <sub>3</sub><br>0.025<br>S <sub>3</sub> | <i>M</i> <sub>11</sub><br>1               | 0.05<br>S <sub>2</sub>   | $\frac{M_{01}}{1}$ $\frac{M_{00}}{0}$ | 0.1<br>$S_2$<br>0.1<br>$S_3$ |
| Strategy<br>(iii*)<br>0.25 | $\frac{M_{11}}{1}$ $\frac{M_{01}}{0}$          | $0.1$ $S_2$ $0.1$ $S_3$ | <i>M</i> <sub>00</sub><br>0               | 0.05<br>S <sub>3</sub>                             |   | $ \begin{array}{c} 0.025 \\ S_2 \\ 0.025 \\ S_2 \\ \end{array} $ | <i>M</i> <sub>00</sub><br>0           | 0.2<br>S <sub>3</sub>        |
| Totals                     | $prob (S_2/A_L) = 0.3prob (S_3/A_L) = 0.3$     |                         |   | $prob (S_2/A_R) = 0.3$<br>$prob (S_3/A_R) = 0.3$   |   |  |                                       |                              |

TABLE 2

## 9. ACT-INDEPENDENCE IN THE SAVAGE FORMULATION

In Section 4, we said that to apply the Savage framework to a decision problem, one must find states of the world that are in some sense act-independent. In the last section, we distinguished two kinds of independence, causal and epistemic. Which kind is needed in the Savage formulation of decision problems? ALLAN GIBBARD AND WILLIAM L. HARPER

The answer is that the Savage formulation has both a  $\mathcal{U}$ -maximizing interpretation and a  $\mathscr{V}$ -maximizing interpretation. On the  $\mathcal{U}$ -maximizing interpretation, the states must be causally independent of the acts, whereas on the  $\mathscr{V}$ -maximizing interpretation, the states must be epistemically independent of the acts. That is to say, if the states are causally act-independent, then utility as calculated by the Savage method is  $\mathscr{U}$ -utility, whereas if the states are epistemically act-independent, then utility as calculated by the Savage method is  $\mathscr{V}$ -utility. If the states are both causally and epistemically act-independent, then the  $\mathscr{U}$ -utility of each act equals its  $\mathscr{V}$ -utility. Thus the Savage formulation itself is not committed to either kind of utility: the kind of utility it yields depends on the way it is applied to decision problems.

The expected utility of an act A in the Savage theory is

(3)  $\Sigma_S \operatorname{prob}(S) \mathcal{D}O(A, S).$ 

If the states S are all known to be causally independent of A, so that for each state S, the agent knows that  $(A \square \rightarrow S) \equiv S$ , then for each S, we have prob  $(S) = prob (A \square \rightarrow S)$ . (3) thus becomes

 $\Sigma_S \text{ prob } (A \Box \rightarrow S) \mathcal{D}O(A, S),$ 

and this, we said in Section 3, is  $\mathcal{U}(A)$ . If, on the other hand, the states S are stochastically independent of A, so that for each S, prob (S) = prob (S/A), then (3) becomes

 $\Sigma_S \operatorname{prob}(S|A) \mathcal{D}(A, S),$ 

which is  $\mathscr{V}(A)$ .

## 10. Newcomb's problem

The Newcomb paradox discussed by Nozick (1969) has the same structure as the case of Solomon discussed in Section 3. Nozick treats it as a conflict between the principle of expected utility maximization and the principle of dominance. On the views we have propounded in this paper, the problem is rather a conflict between two kinds of expected utility maximization. The problem is this. There are two boxes, transparent and opaque; the transparent box contains a thousand dollars. The agent can perform  $A_1$ , taking just the contents of the opaque box, or  $A_2$ , taking the contents of both boxes. A predictor has already placed a million dollars in the opaque box if he predicted  $A_1$  and nothing if he predicted  $A_2$ . The agent knows all this, and he knows the predictor to be highly reliable in that both *prob* (he has

180

predicted  $A_1/A_1$  and prob (he has predicted  $A_2/A_2$ ) are close to one.

To show how the expected utility calculations work, we must add detail to the specification of the situation. Suppose, somewhat unrealistically, that getting no money has a utility of zero, getting \$1000 a utility of 10, that getting \$1,000,000 has a utility of 100, and that getting \$1,001,000 has a utility of 101. Let M be 'there are a million dollars in the opaque box', and suppose prob  $(M/A_1) = 0.9$  and prob  $(M/A_2) = 0.1$ . The calculation of  $\mathscr{V}(A_1)$  and  $\mathscr{V}(A_2)$  is familiar.

$$\begin{aligned} \mathscr{V}(A_1) &= prob \ (M/A_1) \mathscr{D}\$1,000,000 + prob \ (\bar{M}/A_1) \mathscr{D}\$0 \\ &= 0.9 \ (100) + 0.1(0) = 90. \\ \mathscr{V}(A_2) &= prob \ (M/A_2) \mathscr{D}\$1,001,000 + prob \ (\bar{M}/A_2) \mathscr{D}\$1000 \\ &= 0.1(101) + 0.9(10) = 19.1 \end{aligned}$$

Maximization of  $\mathscr{V}$ , as is well known, prescribes taking only the contents of the opaque box.<sup>11</sup>

 $\mathcal{U}(A_1)$  and  $\mathcal{U}(A_2)$  depend on the probability of M, which in turn depends on the probabilities of  $A_1$  and  $A_2$ . For any probability of M, though, we have  $\mathcal{U}(A_2) > \mathcal{U}(A_1)$ . For let the probability of M be  $\mu$ ; then since M is causally act-independent, prob  $(A_1 \Box \rightarrow M) = \mu$  and prob  $(A_2 \Box \rightarrow M) = \mu$ . Therefore

$$\begin{aligned} \mathscr{U}(A_1) &= prob \; (A_1 \square \to M) \mathscr{D}\$1,000,000 + prob \; (A_1 \square \to \overline{M}) \mathscr{D}\$0 \\ &= 100\mu + 0(1-\mu) = 100\mu. \\ \mathscr{U}(A_2) &= prob \; (A_2 \square \to M) \mathscr{D}\$1,001,000 + prob \; (A_2 \square \to \overline{M}) \mathscr{D}\$1000 \\ &= 101\mu + 10(1-\mu) = 91\mu + 10. \end{aligned}$$

Thus  $\mathcal{U}(A_2) - \mathcal{U}(A_1) = 10 - 9\mu$ , and since  $\mu \leq 1$ , this is always positive. Therefore whatever probability M may have,  $\mathcal{U}(A_2) > \mathcal{U}(A_1)$ , and  $\mathcal{U}$ -maximization prescribes taking both boxes.

To some people, this prescription seems irrational.<sup>12</sup> One possible argument against it takes roughly the form 'If you're so smart, why ain't you rich?'  $\mathscr{V}$ -maximizers tend to leave the experiment millionaires whereas  $\mathscr{U}$ -maximizers do not. Both very much want to be millionaires, and the  $\mathscr{V}$ -maximizers usually succeed; hence it must be the  $\mathscr{V}$ -maximizers who are making the rational choice. We take the moral of the paradox to be something else: If someone is very good at predicting behavior and rewards predicted irrationality richly; then irrationality will be richly rewarded.

To see this, consider a variation on Newcomb's story: the subject of the experiment is to take the contents of the opaque box first and learn what it is; he then may choose either to take the thousand dollars in the second box or not to take it. The predictor has an excellent record, and a thoroughly

accepted theory to back it up. Most people find nothing in the first box and then take the contents of the second box. Of the million subjects tested, 1% have found a million dollars in the first box, and strangely enough only 1% of these -100 in 10,000 – have gone on to take the thousand dollars they could each see in the second box. When those who leave the thousand dollars are later asked why they do so, they say things like 'If I were the sort of person who would take the thousand dollars in that situation, I wouldn't be a millionaire'.

On both grounds of  $\mathcal{U}$ -maximization and of  $\mathscr{V}$ -maximization, these new millionaires have acted irrationally in failing to take the extra thousand dollars. They know for certain that they have the million dollars; therefore the  $\mathscr{V}$ -utility of taking the thousand as well is 101, whereas the  $\mathscr{V}$ -utility of not taking it is 100. Even on the view of  $\mathscr{V}$ -maximizers, then, this experiment will almost always make irrational people and only irrational people millionaires. Everyone knows so at the outset.

Return now to the unmodified Newcomb situation, where the subject must take or pass up the thousand dollars before he sees whether the opaque box is full or empty. What happens if the subject knows not merely that the predictor is highly reliable, but that he is infallible? The argument that the  $\mathcal{U}$ -utility of taking both boxes exceeds that of taking only one box goes through unchanged. To some people, however, it seems especially apparent in this case that it is rational to take only the opaque box and irrational to take both. For in this case the subject is certain that he will be a millionaire if and only if he takes only the opaque box. If in the case where the predictor is known to be infallible it is irrational to take both boxes, then,  $\mathcal{U}$ -maximization is not always the rational policy.

We maintain that  $\mathfrak{A}$ -maximization is rational even in the case where the predictor is known to be infallible. True, where R is 'I become a millionaire', the agent knows in this case that R holds if  $A_1$  holds: he knows the truth-functional proposition  $R \equiv A_1$ . From this proposition, however, it does not follow that he would be a millionaire if he did  $A_1$ , or that he would be a non-millionaire if he did  $A_2$ .

If the subject knows for sure that he will take just the opaque box, then he knows for sure that the million dollars is in the opaque box, and so he knows for sure that he will be a millionaire. But since he knows for sure that the million dollars is already in the opaque box, he knows for sure that even if he were to take both boxes, he would be a millionaire. If, on the other hand, the subject knows for sure that he will take both boxes, then he knows for sure that the opaque box is empty, and so he knows for sure that he will be a non-millionaire. But since in this case he knows for sure that the opaque box is empty, he knows for sure that even if he were to take just the opaque box, he would be a non-millionaire.

If the subject does not know what he will do, then what he knows is this: either he will take just the opaque box and be a millionaire, or he will take both boxes and be a non-millionaire. From this, however, it follows neither that (i) if he took just the opaque box, he would be a millionaire, nor that (ii) if he took both boxes he would be a non-millionaire. For (i), the subject knows, is true iff the opaque box is filled with a million dollars, and (ii), the subject knows, is true iff the opaque box is empty. Thus, if (i) followed from what the agent knows, he could conclude for certain that the opaque box contains a million dollars, and if (ii) followed from what the agent knows, he could conclude that the opaque box is empty. Since the subject, we have supposed, does not know what he will do, he can conclude neither that the opaque box contains a million dollars nor that it is empty. Therefore neither (i) nor (ii) follows from what the subject knows.

Rational choice in Newcomb's situation, we maintain, depends on a comparison of what would happen if one took both boxes with what would happen if one took only the opaque box. What the agent knows for sure is this: if he took both boxes, he would get a thousand dollars more than he would if he took only the opaque box. That, on our view, makes it rational for someone who wants as much much as he can get to take both boxes, and irrational to take only one box.

Why, then, does it seem obvious to many people that if the predictor is known to be infallible, it is rational to take only the opaque box and irrational to take both boxes? We have three possible explanations. The first is that a person may have a tendency to want to bring about an indication of a desired state of the world, even if it is known that the act that brings about the indication in no way brings about the desired state itself. Taking just the opaque box would be a sure indication that it contained a million dollars, even though taking just the opaque box in no way brings it about that the box contains a million dollars.

The second possible explanation lies in the force of the argument 'If you're so smart, why ain't you rich?' That argument, though, if it holds good, should apply equally well to the modified Newcomb situation, with a predictor who is known to be highly accurate but fallible. There the conclusion of the argument seems absurd: according to the argument, having already received the million dollars, one should pass up the additional thousand dollars one is free to take, on the grounds that those who are disposed to pass it up tend to become millionaires. Since the argument leads to an absurd conclusion in one case, there must be something wrong with it.

The third possible explanation is the fallacious inference we have just discussed, from

Either I shall take one box and be a millionaire, or I shall take both boxes and be a non-millionaire

to the conclusion

If I were to take one box, I would be a millionaire, and if I were to take both boxes, I would be a non-millionaire.

If, to someone who is free of fallacies, it is still intuitively apparent that the subject should take only the opaque box, we have no further arguments to give him. If in addition he thinks the subject should take only the opaque box even in the case where the predictor is known to be somewhat fallible, if he also thinks that in the modified Newcomb situation the subject, on receiving the extra million dollars, should take the extra thousand, if he also thinks that it is rational for Reoboam to be severe, and if he also thinks it is rational for Solomon to abstain from his neighbor's wife, then he may genuinely have the intuitions of a  $\mathscr{V}$ -maximizer:  $\mathscr{V}$ -maximization then provides a systematic account of his intuitions. If he thinks some of these things but not all of them, then we leave it to him to provide a systematic account of his views. Our own views are systematically accounted for by  $\mathscr{U}$ -maximization.

## 11. STABILITY OF DECISION

When a person decides what to do, he has in effect learned what he will do, and so he has new information. He will adjust his probability ascriptions accordingly. These adjustments may affect the  $\mathcal{U}$ -utility of the various acts open to him.

Indeed, once the person decides to perform an act A, the  $\mathscr{U}$ -utility of A will be equal to its  $\mathscr{V}$ -utility.<sup>13</sup> Or at least this holds if Consequence 1 in Section 2, that  $A \supset [(A \square \rightarrow C) \equiv C]$ , is a logical truth. For we saw in the proof of Assertion 1 that if Consequence 1 is a logical truth, then for any pair of propositions P and Q, prob  $(P \square \rightarrow Q/P) = prob(Q/P)$ . Now let  $\mathscr{U}_A$  (A) be the  $\mathscr{U}$ -utility of act A as reckoned by the agent after he has decided

for sure to do A, let *prob* give the agent's probability ascriptions before he has decided what to do. Let  $prob_A$  give the agent's probability ascriptions after he has decided for sure to do A. Then for any proposition P,  $prob_A$  (P) = prob (P/A). Thus

$$\begin{aligned} \boldsymbol{\mathscr{U}}_{A}(A) &= \boldsymbol{\Sigma}_{O} prob_{A} \ (A \Box \rightarrow O) \boldsymbol{\mathscr{D}} O \\ &= \boldsymbol{\Sigma}_{O} prob \ (A \Box \rightarrow O/A) \boldsymbol{\mathscr{D}} O \\ &= \boldsymbol{\Sigma}_{O} prob \ (O/A) \boldsymbol{\mathscr{D}} O \\ &= \boldsymbol{\mathscr{V}}(A). \end{aligned}$$

The  $\mathscr{V}$ -utility of an act, then, is what its  $\mathscr{Q}$ -utility would be if the agent knew he were going to perform it.

It does not follow that once a person knows what he will do,  $\mathscr{V}$ -maximization and  $\mathscr{U}$ -maximization give the same prescriptions. For although for any act  $A, \mathscr{U}_A(A) = \mathscr{V}(A)$ , it is not in general true that for alternatives B to A,  $\mathscr{U}_A(B) = \mathscr{V}(B)$ . Thus in cases where  $\mathscr{U}(A) < \mathscr{U}(B)$  but  $\mathscr{V}(A) > \mathscr{V}(B)$ , it is consistent with what we have said to suppose that  $\mathscr{U}_A(A) < \mathscr{U}_A(B)$ . In such a case,  $\mathscr{V}$ -maximization prescribes A regardless of what the agent believes he will do, but even if he believes he will do  $A, \mathscr{U}$ -maximization prescribes B. The situation is this:

$$\mathscr{U}_{A}(B) > \mathscr{U}_{A}(A) = \mathscr{V}(A) > \mathscr{V}(B).$$

Even though, once an agent knows what he will do, the distinction between the  $\mathcal{U}$ -utility of that act and its  $\mathscr{V}$ -utility disappears, the distinction between  $\mathcal{U}$ -maximization and  $\mathscr{V}$ -maximization remains.

That deciding what to do can affect the  $\mathcal{U}$ -utilities of the acts open to an agent raises a problem of stability of decision for  $\mathcal{U}$ -maximizers. Consider the story of the man who met death in Damascus.<sup>14</sup> Death looked surprised, but then recovered his ghastly composure and said, 'I am coming for you to-morrow'. The terrified man that night bought a camel and rode to Aleppo. The next day, death knocked on the door of the room where he was hiding and said 'I have come for you'.

'But I thought you would be looking for me in Damascus', said the man.

'Not at all', said death 'that is why I was surprised to see you yesterday. I knew that today I was to find you in Aleppo'.

Now suppose the man knows the following. Death works from an appointment book which states time and place; a person dies if and only if the book correctly states in what city he will be at the stated time. The book is made up weeks in advance on the basis of highly reliable predictions. An appointment on the next day has been inscribed for him. Suppose, on this basis, the man would take his being in Damascus the next day as strong evidence that his appointment with death is in Damascus, and would take his being in Aleppo the next day as strong evidence that his appointment is in Aleppo.

Two acts are open to him: A, go to Aleppo, and D, stay in Damascus. There are two possibilities:  $S_A$ , death will seek him in Aleppo, and  $S_D$ , death will seek him in Damascus. He knows that death will find him if and only if death looks for him in the right city, so that, where L is that he lives, he knows  $(D \Box \rightarrow L) \equiv S_A$  and  $(A \Box \rightarrow L) \equiv S_D$ . He ascribes conditional probabilities  $prob(S_A/A) \approx 1$  and  $prob(S_D/D) \approx 1$ ; suppose these are both 0.99 and that  $prob(S_D/A) = 0.01$  and  $prob(S_A/D) = 0.01$ . Suppose  $\mathcal{D}(L) = -100$  and  $\mathcal{D}(L) = 0$ . Then where  $\alpha$  is prob(A), his probability of going to Aleppo, and  $1 - \alpha$  is his probability of going to Damascus,

$$prob (A \Box \rightarrow L) = prob (S_D) = \alpha prob (S_D/A) + (1-\alpha) prob (S_D/D) = 0.01\alpha + 0.99(1-\alpha) = 0.99 - 0.98\alpha prob (A \Box \rightarrow \overline{L}) = prob (S_A) = 1 - prob (S_D) = 0.01 + 0.98\alpha.$$

Thus

$$\mathcal{U}(A) = prob \ (A \Box \rightarrow L)\mathcal{D}(L) + prob \ (A \Box \rightarrow \overline{L})\mathcal{D}(L)$$
$$= (0.01 + 0.98\alpha) \ (-100) = -1 \ -98\alpha.$$

By a like calculation,  $\mathcal{U}(D) = -99 + 98\alpha$ . Thus if  $\alpha = 1$ , then  $\mathcal{U}(D) = -1$  and  $\mathcal{U}(A) = -99$ , and thus  $\mathcal{U}(D) > \mathcal{U}(A)$ . If  $\alpha = 0$ , then  $\mathcal{U}(D) = -99$  and  $\mathcal{U}(A) = -1$ , so that  $\mathcal{U}(A) > \mathcal{U}(D)$ . Indeed we have  $\mathcal{U}(D) > \mathcal{U}(A)$  whenever prob (A) > 1/2, and  $\mathcal{U}(A) > \mathcal{U}(D)$  whenever prob (D) > 1/2.

What are we to make of this? If the man ascribes himself equal probabilities of going to Aleppo and staying in Damascus, he has equal grounds for thinking that death intends to seek him in Damascus and that death intends to seek him in Aleppo. If, however, he decides to go to Aleppo, he then has strong grounds for expecting that Aleppo is where death already expects him to be, and hence it is rational for him to prefer staying in Damascus. Similarly, deciding to stay in Damascus would give him strong grounds for thinking that he ought to go to Aleppo: once he knows he will stay in Damascus, he can be almost sure that death already expects him in Damascus, and hence that if he had gone to Aleppo, death would have sought him in vain.

 $\mathscr{V}$ -maximization does not lead to such instability. What happens to  $\mathscr{V}$ utility when an agent knows for sure what he will do is somewhat unclear. Standard probability theory offers no interpretation of  $prob_A(O/B)$  where prob(B/A) = O, and so on the standard theory, once an agent knows for sure
what he will do, the  $\mathscr{V}$ -utility of the alternatives ceases to be well-defined.

186

What we can say about  $\mathscr{V}$ -utility is this: as long as an act's being performed has non-zero probability, its  $\mathscr{V}$ -utility is independent of its probability and the probabilities of alternatives to it. For the  $\mathscr{V}$ -utility of an act A depends on conditional probabilities of the form *prob* (O/A). This is just the probability the agent would ascribe to O on learning A for sure, and that is independent of how likely he now regards A. Whereas, then, the  $\mathscr{U}$ -utility of an act may vary with its probability of being performed, its  $\mathscr{V}$ -utility does not.  $\mathscr{U}$ -maximization, then, may give rise to a kind of instability which  $\mathscr{V}$ -maximization precludes: in certain cases, an act will be  $\mathscr{U}$ -maximal if and only if the probability of its performance is low.

Is this a reason for preferring  $\mathscr{V}$ -maximization? We think not. In the case of death in Damascus, rational decision does seem to be unstable. Any reason the doomed man has for thinking he will go to Aleppo is a reason for thinking he would live longer if he stayed in Damascus, and any reason he has for thinking he will stay in Damascus is reason for thinking he would live longer if he went to Aleppo. Thinking he will do one is reason for doing the other. That there can be cases of unstable  $\mathscr{U}$ -maximization seems strange, but the strangeness lies in the cases, not in  $\mathscr{U}$ -maximization: instability of rational decision seems to be a genuine feature of such cases.

#### 12. APPLICATIONS TO GAME THEORY

Game theory provides many cases where  $\mathcal{U}$ -maximizing and  $\checkmark$ -maximizing diverge; perhaps the most striking of these is the prisoners' dilemma, for which a desirability matrix is shown.



Here  $A_0$  and  $B_0$  are respectively A's and B's options of confessing, while  $A_1$  and  $B_1$  are the options of not confessing. The desirabilities reflect these facts: (1) if both confess, they both get long prison terms; (2) if one confesses and the other doesn't, then the confessor gets off while the other gets an even longer prison term; (3) if neither confesses, both get off with very light sentences.

Suppose each prisoner knows that the other thinks in much the same way he does. Then his own choice gives him evidence for what the other will do. Thus, the conditional probability of a long prison term on his confessing is greater than the conditional probability of a long prison term on his not confessing. If the difference between these two conditional probabilities is sufficiently great, then  $\mathscr{V}$ -maximizing will prescribe not confessing.

The  $\mathscr{V}$ -utilities of the acts open to B will be as follows.

$$\Psi(B_0) = \operatorname{prob} (A_0/B_0) \times 1 + \operatorname{prob} (A_1/B_0) \times 10$$
  
$$\Psi(B_1) = \operatorname{prob} (A_0/B_1) \times 0 + \operatorname{prob} (A_1/B_1) \times 9.$$

If prob  $(A_1/B_1) - prob (A_1/B_0)$  is sufficiently great (in this case 1/9 or more), then  $\mathscr{V}$ -maximizing recommends that B take option  $B_1$  and not confess. If the probabilities for A are similar, then  $\mathscr{V}$ -maximizing also recommends not confessing for A. The outcome if both  $\mathscr{V}$ -maximize is  $A_1B_1$ , the optimal one of mutual co-operation.<sup>15</sup>

For a  $\mathcal{U}$ -maximizer, dominance applies because his companion's choice is causally independent of his own. Therefore,  $\mathcal{U}$ -maximizing yields the classical outcome of the prisoners' dilemma. This suggests that  $\mathcal{U}$ -maximizing and not  $\mathscr{V}$ -maximizing corresponds to the kind of utility maximizing commonly assumed in game theory.

University of Michigan and University of Western Ontario.

#### NOTES

\* An earlier draft of this paper was circulated in January 1976. A much shorter version was presented to the 5th International Congress of Logic, Methodology, and Philosophy of Science, London, Ontario, August 1975. There, and at the earlier University of Western Ontario research colloquium on Foundations and Applications of Decision Theory we benefited from discussions with many people; in particular we should mention Richard Jeffrey, Isaac Levi, Barry O'Neill and Howard Sobel.

 $^{1}$  Lewis first presented this result at the June 1972 meeting of the Canadian Philosophical Association.

<sup>2</sup> Although the rough treatment of counterfactuals we propose is similar is many respects to the theories developed by Stalnaker and Lewis, it differs from them in some important respects. Stalnaker and Lewis each base their accounts on comparisons of overall similarity of worlds. On our account, what matters is comparative similarity of worlds at the instant of decision. Whether a given *a*-world is selected as  $W_a$  depends not at all on how similar the future in that world is to the actual future; whatever similarities the future in  $W_a$  may have to the actual future will be a semantical consequence of laws of nature, conditions in  $W_a$  at the instant of decision, and actual conditions at that

188

instant. (Roughly, then, they will be consequences of laws of nature and the similarity of  $W_a$  to the actual world at the instant of decision.) We consider only worlds in which the past is exactly like the actual past, for since the agent cannot now alter the past, those are the only worlds relevant to his decision. Lewis (1973, p. 566 and in conversation) suggests that a proper treatment of overall similarity will yield as a deep consequence of general facts about the world the conditions we are imposing by fiat.

<sup>3</sup> In characterizing our conditional we have imposed the Stalnaker-like constraint that there is a unique world  $W_a$  which would eventuate from performing *a* at *t*. Our rationale for Axiom 2 depends on this assumption and on the assumption that if *a* is actually performed then  $W_a$  is the actual world itself. Consequence 1 is weaker than Axiom 2, and only depends on the second part of this assumption. In circumstances where these assumptions break down, it would seem to us that using conditionals to compute expected utility is inappropriate. A more general approach is needed to handle such cases.

<sup>4</sup> This is stated by Lewis (1975, note 10).

<sup>5</sup> This is our understanding of a proposal made by Jeffrey at the colloquium on Foundations and Applications of Decision Theory, University of Western Ontario, 1975. J. H. Sobel shows (in an unpublished manuscript) that, for all we have said, these new, conditionalized states may not themselves be act-independent. This section is slightly changed in light of Sobel's result.

<sup>6</sup> Meeting of the Canadian Philosophical Association, 1972. Nozick gives a similar example (1969, p. 125).

<sup>7</sup> Sobel actually uses  $A^* \Box \to \overline{S}$  does hold' where we use  $A^* \Box \to S$  does not hold'. With Axiom 2, these are equivalent.

<sup>8</sup> We realize that a Bayesian king presented with these data would not ordinarily take on degrees of belief that exactly match the frequencies given in the table; nevertheless, with appropriate prior beliefs and evidence, he would come to have those degrees of belief. Assume that he does.

<sup>9</sup> Under these conditions, if A and B are the only alternatives, then  $S_i^*$  holds if and only if  $S_i$  holds. If there are other alternatives, it may be that neither A nor B is performed and  $S_i$  holds without either  $A \square \rightarrow S_i$  or  $B \square \rightarrow S_i$ . In that case, what matters is not whether it would be rational to prefer A to B knowing that  $S_i$  holds, but whether it would be rational to prefer A to B knowing  $(A \square \rightarrow S_i) \& (B \square \rightarrow S_i)$ .

<sup>10</sup> Nozick (1969) in effect endorses the Principle of Dominance with Stochastic Independence (p. 127), but not  $\mathscr{V}$ -maximization: in cases of the kind we have been considering, he considers the recommendations of  $\mathscr{V}$ -maximization 'perfectly wild' (p. 126). Nozick also states and endorses the principle of dominance with causal independence (p. 132).

<sup>11</sup> For *I*-maximizing treatments of Newcomb's problem, see Bar Hillel and Margalit (1972) and Levi (1975).

<sup>12</sup> Levi (1975) reconstructs Nozick's argument for taking both boxes in a way which uses prob(M) rather than  $prob(M/A_1)$  and  $prob(M/A_2)$  as the appropriate probabilities for computing expected utility in Newcomb's problem. This agrees with  $\mathcal{U}$ -maximizing in that the same probabilities are used for computing expected utility for  $A_1$  as for  $A_2$ , and results in the same recommendation to take both boxes. Levi is one of the people to whom this recommendation seems irrational. <sup>13</sup> We owe this point to Barry O'Neill.

<sup>14</sup> A version of this story quoted from Somerset Maugham's play Sheppey (New York, Doubleday 1934) appears on the facing page of John O'Hara's novel Appointment in Samarra. (New York, Random House 1934). The story is undoubtedly much older.

<sup>15</sup> Nozick (1969), Brams (1975), Grofman (1975) and Rapoport (1975), have all suggested a link between Newcomb's problem and the Prisoners Dilemma. Brams, Grofman and Rapoport all endorse co-operative solutions, Rapoport (1975, p. 619) appears to endorse *#*-maximizing.

#### BIBLIOGRAPHY

- Bar Hillel, M. and Margalit, A., 'Newcomb's Paradox Revisited', British Journal for the Philosophy of Science 23 (1972), 295-304.
- Brams, S.J., 'Newcomb's Problem and Prisoners' Dilemma', *Journal of Conflict Resolution* 19, 4, December 1975.
- Grofman, B., 'A Comment on "Newcomb's Problem and the Prisoners' Dilemma" ', Manuscript, September 1975.
- Jeffrey, R. C., The Logic of Decision, McGraw Hill, New York, 1965.
- Levi, I., 'Newcomb's Many Problems', Theory and Decision 6 (1975), 161-75.
- Lewis, D. K., 'Probabilities of Conditionals and Conditional Probabilities', *Philosophical Review* 85 (1976), 297-315.
- Lewis, D. K., Counterfactuals, Harvard University Press, Cambridge, Massachusetts, 1973a.
- Lewis, D. K., 'The Counterfactual Analysis of Causation', *Journal of Philosophy* Volume LXX, Number 17, (1973b), 556-567.
- Nozick, R., 'Newcomb's Problem and Two Principles of Choice', in Nicholas Rescher (ed.), Essays in Honor of Carl G. Hempel, Reidel, Dordrecht-Holland, 1969.
- Rapoport A., 'Comment on Brams's Discussion of Newcomb's Paradox', Journal of Conflict Resolution 19, 4, December 1975.
- Savage, L. J., *The Foundations of Statistics*, Dover, New York, 1972 (original edition 1954).
- Sobel, J. H., 'Utilitarianisms: Simple and General', Inquiry 13 (1970), 394-449.
- Stalnaker, R., 'A Theory of Conditionals', in *Studies in Logical Theory*, American Philosophical Quarterly Monograph Series, No. 2, 1968.
- Stalnaker, R. and Thomason, R., 'A Semantic Analysis of Conditional Logic', *Theoria* 36 (1970), 23-42.

190