

## UNPREDICTABILITY, PROBABILITY UPDATING AND THE THREE PRISONERS PARADOX

Rosangela H. Loschi

Universidade Federal de Minas Gerais, Brazil

Pilar L. Iglesias

Pontificia Universidad Católica, Chile

Sergio Wechsler

Universidade de São Paulo, Brazil

loschi@est.ufmg.br

*This paper discusses the Three Prisoners paradox in the light of three different procedures for the updating of probabilities - Bayesian conditioning, superconditioning and Jeffrey's rule - as well as assuming the unpredictability of receipt of information by prisoner A. The formulation of the paradox in this temporal setting brings new insight to the problem and, on the other hand, the paradox is a good way to explain the different updating probability procedures and the difference between conditional probabilities and posterior distributions.*

### INTRODUCTION

*The three prisoners. Two are to be shot and the other freed; none is to know his fate until the morning. Prisoner A asks the warden to confide the name of one other than himself who will be shot, explaining that as there must be at least one, the warden won't be giving away anything relevant to A's own case. The warden agrees, and tells him that B will be shot. This cheers A up a little, by making his judgement probability for being freed rise from 1/3 to 1/2. But that's silly: A knew already that one of the others would be shot, and (as he told the warden) he's no wiser about his own fate for knowing the name of some other victim." (Jeffrey, 1992, p. 122)*

The Three Prisoners paradox, just presented, is an old problem from Probability Calculus which has been discussed from many different points-of-view. This "paradox" is also an excellent tool in teaching some different procedures for the updating of probabilities and the difference between posterior distributions and conditional probabilities. On the other hand, these updating procedures provide new insight in the paradox.

Apparently the answer provided by prisoner A is a contradiction since, as prisoner A tells the warden, the information given about the other two prisoners does not apprise prisoner A of his own condition. Thus, the prisoner A's opinion about the event "A will live" ought to be the same after the receipt of the information provided by the warden, i.e., the posterior probability of this event should also be 1/3. On the other hand, we should notice that the solution presented by prisoner A does not change the prior opinion of indifference among the prisoners revealed by the prisoner A's prior distribution. What is the right answer?

According to Jeffrey (1992), prisoner A follows erroneously the council of parochialism - that is, prisoner A constructs his posterior distribution using Bayesian conditioning. Prisoner A considers the sample space  $\Omega_1 = \{A, B, C\}$ , where  $X$  represents "X will live,"  $X = A, B, C$ . On this space the information provided by the warden does not generate a sufficient partition for  $\{P, P^*\}$ , which makes the use of Bayesian conditioning inappropriate. (See de Finetti (1972), for the difference between Bayes' formula and Bayesian Conditioning.)

Some alternative procedures for the updating of probabilities are proposed in the literature. See, for example, Diaconis and Zabell (1982), Jeffrey (1992), Howson and Urbach (1993) and others. However, there is no guidance for coherent temporal behavior which produces an inevitable probability updating procedure (Goldstein, 1985; Dawid, 1985). Consequently, prisoner A may update his/her prior opinion by means of a complete reassessment of his/her opinion.

This paper extends previous works by presenting alternative explanations to the Three Prisoners paradox using Bayesian conditioning, superconditioning and Jeffrey's rule as well as by

considering (when it is possible) two different settings: i) prisoner A plans to ask the warden for information and ii) prisoner A receives unexpectedly the information. It will be argued that the solution presented by A may be correct if we do not arbitrate, as it is usually done, that: the receipt of information about B and C is planned from the beginning; the received information always takes prisoner A to the certainty that B will die; and prisoner A judges that the warden has the same chance to reveal the names of prisoners B and C, in case A is the one to be freed. Therefore, we will point out some limitations to Bayesian conditioning as well as highlight the influence of the way by which the information is received on the construction of the posterior distribution. In this paper, it will be assumed that the prior probabilities declared by prisoner A for the events “X will live,”  $X = A, B, C$ , are the same. We denote by  $P$  and  $P^*$  two probability measures defined on the measurable space  $(\Omega, \mathcal{A})$ , where  $\Omega$  is a countable set, and interpret  $P$  and  $P^*$  as the prior and posterior opinions of prisoner A about events in  $\mathcal{A}$ , respectively.

Next, we will briefly present Bayesian conditioning and responses to that paradox using this procedure for the updating of probabilities. The difference between posterior distribution and conditional probabilities is briefly discussed.

### BAYESIAN CONDITIONING AND THE CALCULUS OF PRISONER A

Bayesian conditioning is the procedure for the updating of probabilities which links prior and posterior distributions using Bayes’ formula, that is, for all event  $E$  in  $\mathcal{A}$  such that  $P(E) > 0$ , the posterior distribution  $P^*$  is obtained from the prior distribution  $P$  using the expression:

$$P^*(\cdot) = P(\cdot | E). \quad (1)$$

Using some properties of probability measure, Jeffrey (1992) states some conditions under which Bayesian conditioning can be performed. These conditions are presented in the following theorem.

*Theorem 1: Let  $E \in \mathcal{A}$  be an event such that  $P(E) > 0$ . Then, for every  $A \in \mathcal{A}$ ,  $P^*(A) = P(A|E)$  if, and only if*

1.  $P^*(E) = 1$
2.  $P^*(A|E) = P(A|E)$ .

Theorem 1 is named by Howson and Urbach (1993) *The Principle of Bayesian Conditionalisation* and conditions (1) and (2) are respectively known as *certainty* and *sufficiency*. Using this terminology it can be stated that Bayesian conditioning is an acceptable procedure for the updating of probabilities if the received information makes You move from an initial state of uncertainty about the conditioning event  $E$  to the posterior certainty of its occurrence and if, beyond this, the partition  $\{E, \bar{E}\}$ , generated by the received information, contains all relevant information to the construction of your posterior distribution - i.e.,  $\{E, \bar{E}\}$  is a sufficient partition to the family  $\{P, P^*\}$ . Notice, moreover, that certainty and sufficiency (conditions (i) and (ii), respectively) are subjective evaluations and depends on your judgement.

To analyze the Three Prisoners paradox using Bayesian conditioning, firstly, consider that prisoner A plans to ask the warden about the situation of the other two prisoners from the beginning. The sample space which appropriately describes the experiment performed by prisoner A is the space  $\Omega_2 = \{Ab, Ac, Bc, Cb\}$  where  $Xy$  denotes the event “X will live and the warden names y.” Thus, initially we have that  $P(Ab) = p = 1/3 - P(Ac)$  and  $P(Bc) = P(Cb) = 1/3$ , where  $p \in [0, 1/3)$ . Notice that, in this case, the prior conditional probability that the warden informs that B will die, supposing that A will be freed is  $3p$ .

At the very moment A declares his prior distribution, he also reveals his probabilities, supposing that the event “the warden informs that B will die” occurs. Denote this event by  $E$  and notice that  $E = Ab \cup Cb$ . Supposing that  $E$  occurs and being coherent from the static point-of-view de Finetti (1937), the prior conditional probability of the event “A will live” is obtained from Bayes’ formula as follows  $P(Ab \cup Cb | E) = 3p/(3p + 1)$ .

If prisoner A judges that the conditional probabilities in  $E$  declared before are kept after consulting the warden and if his new opinion about the event  $E$  is  $P^*(E)=1$ , Bayesian conditioning can be performed and the posterior probability for the event “A will live” is the probability obtained in the expression (1), i.e.,  $P^*(A)=P(Ab \cup Cb | E) = 3p/(3p + 1)$ .

Notice that the posterior probability of  $A$  is  $1/3$  only if  $p = 1/6$ . This choice of  $p$  shows that, for prisoner A, the warden has equal chances to tell the names of either B or C, in case A is the prisoner who will live. On the other hand, if A thinks the warden would never tell the name of C, in case A were the survivor, that is, if he declares  $p = 1/3$ , the posterior distribution provided in (1) would agree with the posterior distribution divulged by prisoner A.

Conversely, suppose now that prisoner A has not planned to ask the warden about the conditions of prisoners B and C. In this case, the suitable sample space to describe the experiment performed by A is the space  $\Omega_1$  defined in above.

As it has already been assumed, consider that each prisoner has a  $1/3$  chance of being the survivor. Consequently, the prior conditional probability for the event “A will live,” supposing that B will die, is given by:

$$P(A|E) = P(A) / [P(A \cup C)] = 1/2, \quad (2)$$

in which  $E = A \cup C$ .

After A states his opinion about his being the one who will survive, the warden tells him that B will die. If in possession of the information given by the warden, A judges that  $E$  is a certain event and that all the conditional probabilities in  $E$  are not modified, the posterior probability for the event “A will live” is given by the expression (2), confirming prisoner A’s initial statement.

In both situations we assume that the received information is the same and Bayesian conditioning is the procedure adopted to update probabilities. Yet only in some situations the posterior distribution for the event “A will live” coincides with what seemed intuitive and logical at first. Besides, if the information is received unexpectedly, prisoner A will always be right.

Notice that the way by which the information is received influences the posterior probability calculus for the event “A will live” as it interferes directly with the construction of the sample space appropriate for the problem. In case the receipt of information is not previously anticipated, the change in the value of the posterior probability declared by prisoner A is plainly justifiable (in this case there is a change in the expectation), and the assessment made by prisoner A is not contradictory.

In the next section we present a different experimental design which could be performed by prisoner A. For this “new” design Bayesian conditioning can not be used as updating probabilities procedure.

#### EXPLAINING THE PARADOX VIA SUPERCONDITIONING

Suppose that prisoner A has not planned to ask the warden for information - i.e., A considers the space  $\Omega_1 = \{A, B, C\}$ . Admit that prisoner A’s prior probabilities for the survival of each prisoner is  $1/3$ .

After eliciting his prior distribution for the events of  $\Omega_1$  (and before declaring his posterior distribution), prisoner A decides to ask the warden which of the other two prisoners will be sentenced, alleging that this information does not tell anything about his own condition. How can the posterior distribution on  $\Omega_1$  be determined? By doing so, prisoner A performs an experiment whose possible results are in the sample space  $\Omega_2 = \{Ab, Ac, Bc, Cb\}$ . Notice that this situation is slightly different from that one described previously, where prisoner A has planned to ask for information before declaring his prior distribution. Here Bayesian conditioning is not applicable.

The superconditioning introduced by Diaconis and Zabell (1982) offers us an alternative way to explain the Three Prisoners paradox, in this situation.

Definition 1: (*Superconditioning*). The posterior distribution  $P^*$  can be obtained from the prior distribution  $P$  by Superconditioning if there exist a probability space  $(\bar{\Omega}, \bar{A}, Q)$  and a class of events  $D = \{E_w \in \bar{A}\}$ ,  $w \in \Omega$ , such that:

- (i)  $E_w$  occurs if, and only if  $w$  occurs, for all  $w \in \Omega$ ;
- (ii)  $Q(E_w) = P(\{w\})$ , for all  $w \in \Omega$  and
- (iii)  $P^*(\{w\}) = Q(E_w | E)$  for every  $w \in \Omega$  and for an event  $E \in \bar{A}$  such that  $Q(E_w) > 0$ .

Notice that, as for Bayesian conditioning, the updating of probabilities via superconditioning is obtained multiplying the prior distribution by an appropriate likelihood function.

To analyze the Three Prisoners paradox considering the design described before in this section, suppose that prisoner A specifies the following probability measure  $Q$  on  $\Omega_2$ :  $Q(Ab) = q_1$ ,  $Q(Ac) = q_2$ ,  $Q(Bc) = q_3$  and  $Q(Cb) = q_4$ , where  $q_i \in [0, 1]$ , for all  $i$  and  $\sum_i q_i = 1$ , before asking the warden. Define  $E_A = Ab \cup Ac$ ,  $E_B = Bc$ ,  $E_C = Cb$  and  $E = Ab \cup Cb$ . Then we have that  $Q(E_A) = q_1 + q_2 = 1/3$ ,  $Q(E_B) = q_3 = 1/3 = q_4 = Q(E_C)$  and  $Q(E) = q_1 + 1/3$ . Thus, using superconditioning, the posterior probability for the event "A will live" is  $P^*(A) = Q(E_A | E) = Q(Ab) / Q(E) = q_1 / [q_1 + 1/3]$ .

At first, suppose that prisoner A thinks the warden has an equal preference for both prisoners B and C, in such a way that  $Q(Ab) = Q(Ac) = 1/6$  and  $Q(Bc) = Q(Cb) = 1/3$ . From (iii) we have that  $P^*(A) = (1/6) / (1/2) = 1/3$ , which coincides with prisoner A's prior opinion about this event.

On the other hand, if A suspects that the warden will never tell that prisoner C will die, in case he is the survivor - which makes A declare  $Q(Ac) = 0$  and  $Q(Ab) = Q(Bc) = Q(Cb) = 1/3$  - the posterior probability of A being the survivor is  $P^*(A) = 1/2$ , what confirms the reason for prisoner A's excitement after talking with the warden.

If prisoner A adopts the superconditioning and has not uniform prior distribution on singleton events of  $\Omega_1$ , A will only change his initial opinion to 1/2 if there is any reason to judge that  $Q(Ab) = Q(Cb)$ . In case prisoner A considers  $Q(Cb) = 2Q(Ab)$  the posterior probability of the event "A will live" will be  $P^*(A) = P(A) = 1/3$ .

However, we must notice that the posterior distribution  $P^*$  is not always obtained from the prior distribution  $P$  by superconditioning. In the Three Prisoners paradox (if the experimental design performed by prisoner A is that one described in this section) we can always use this updating procedure because the sample space is countable, as we can see in the next theorem from Diaconis and Zabell (1982).

Theorem 2: Let  $P$  and  $P^*$  be probability measures defined on  $(\Omega, \bar{A})$ , where  $\Omega$  is a countable set.  $P^*$  is obtained from  $P$  by superconditioning if, and only if, there is a constant  $B \neq 1$  such that  $P^*(w) \leq BP(w)$ , for all  $w \in \Omega$ .

More about superconditioning can be found in Jeffrey (1992). In the next section, another way to explain the calculus done by prisoner A will be shown by considering Jeffrey's rule.

#### JEFFREY'S RULE AND THE THREE PRISONER PARADOX

Another possible experimental design which could be performed by prisoner A is described in the following. Suppose it is not prisoner A's intention to ask the warden for information about prisoners B and C, i.e., admit the sample space  $\Omega_1 = \{A, B, C\}$ . As before, admit that the prior probability distribution stated by A for the singleton events of  $\Omega_1$  are all 1/3.

After having declared his prior distribution on  $\Omega_1$ , assume that prisoner A is unexpectedly informed by the warden that B will die, and that this information makes A change his opinion about the event  $E = A \cup C$  arbitrarily, establishing that  $P^*(E) = p^* < 1$ . Notice that, as for Bayesian conditioning, the events  $E$  and  $\bar{E}$  define a partition of the sample space, but here the information provided by the warden does not make prisoner A certain about the truth of  $E$ , that is, prisoner A believes that the warden can be lying.

Jeffrey's rule, presented in Jeffrey (1992), permits the construction of the posterior distribution in situations similar to the one we have just described.

*Definition 2: (Jeffrey's Rule). Let  $\Omega$  be a countable sample space. The posterior distribution  $P^*$  is obtained from the prior distribution  $P$  by Jeffrey's Rule if*

$$P^*(.) = \sum_i P(. | E_i) P^*(E_i), \quad E_i \in E, \quad (3)$$

where  $E$  is a partition of  $\Omega$ ,  $P^*(E_i) \geq 0$  for every  $i$  and  $\sum_i P^*(. | E_i) = 1$ .

Jeffrey's rule also permits the arbitrary updating of the prior probabilities attributed to the elements of the partition. Its difference from Bayesian conditioning is that those arbitrary reassessments  $P^*(E_i)$  may assume values smaller than 1 for every  $i$ . Notice that, since there is not a procedure to obtain the posterior probabilities  $P^*(E_i)$  to elicit them can be such a psychologically complex task as stating the prior distribution.

Notice that, as defined in (3),  $P^*$  does not follow the laws of Probability Calculus, that is,  $P^*$  is not necessarily coherent in the sense defined by de Finetti (1937). A posterior probability measure will be obtained by Jeffrey's rule, if the partition of  $\Omega$  generated by the received information is sufficient to  $\{P, P^*\}$  as we can see in Theorem 3 in the following.

*Theorem 2: Let  $(\Omega, A)$  be a probability space where  $\Omega$  is a countable set. Consider an  $A$ -measurable event  $A$  and suppose that for every  $E_i \in E$ ,  $P^*(E_i) > 0$ . Then  $P^*(A) = \sum P(A | E_i) P^*(E_i)$  if, and only if, for every  $A \in A$  and for every  $E_i \in E$ ,*

$$P(A | E_i) = P^*(A | E_i). \quad (4)$$

The proof of this theorem and a valuable discussion on Jeffrey's rule can be found in Jeffrey (1992). See more in Diaconis and Zabell (1982) and Loschi, Iglesias and Arellano-Valle (2002).

Let us admit, however, that the condition (4) (the *J-condition* or sufficiency condition) above is verified for prisoner A's prior and posterior opinions - i.e., for A the partition  $E = \{E, \bar{E}\}$ , generated from the information provided by the warden, is sufficient to the family of probability measures  $\{P, P^*\}$ . Since  $P^*(E) = p^* < 1$ , the posterior probability obtained by Jeffrey's rule for the event "A will live" is:  $P^*(A) = P(A | E) p^* + P(A | \bar{E}) (1 - p^*) = p^*/2$ .

Notice that the prior probability of event "A will live" remains unchanged afterwards if the information provided by the warden makes prisoner A less uncertain about the event  $E$ , but not totally convinced of its occurrence. In fact, A would have to declare  $p^* = 2/3$  to have his prior distribution unchanged. In any other circumstances the posterior probability for the event "A will live" will be different from the prior probability and also different from the posterior probability established by prisoner A. Notice also that the only reason for prisoner A's excitement about having his chance of survival increased to  $1/2$  is when  $p^* = 1$ , which would be the same as using Bayesian conditioning.

On the other hand, if prisoner A intends to inquire the warden about his companions before stating his prior distribution, the information that B will die causes the occurrence of the event  $Ab \cup Cb$  of the sample space  $\Omega_2$ . In this case, updating probabilities via Jeffrey's rule is equivalent to an updating by Bayesian conditioning; thus the posterior probability of A being the

survivor is given in (1). Notice that Jeffrey's rule also does not apply to the problem considered by Morgan *et al.* (1991). In that problem, to use Jeffrey's rule always corresponds considering Bayesian conditioning.

#### FINAL COMMENTS

In this paper we present several scenarios to explain the Three Prisoners paradox and different approaches are considered to explain it. However, the problem is more general and often occurs in practical situations. Suppose, for example, we are predicting a time series with a truly dynamic model. Then, before, seeing more data, news broke that a new terrorism attack has just taken place around the world. How are we going to make use of this relevant piece of information? In order to improve our predictions, this unpredictable information must be incorporated in the model. Notice that this real problem is very similar to the problem lived by prisoner A.

The Three Prisoners paradox, as well as any other problem which involves the updating of probabilities, does not have a unique solution (for the Three Prisoners paradox we have seen that the number of solutions can be very large, and not only  $1/2$  and  $1/3$  as it is usually considered). One of the reasons is the lack of a normative rule for the choice of the procedure to be used for the updating of probabilities. In the absence of such a rule we may choose the procedure we judge the most adequate to the construction of our posterior distribution. We may even construct the posterior distribution by means of a completely arbitrary reassessment, that is, without using any mathematical formula.

Besides, the way by which the information is received influences our judgement. For all the situations discussed, the information was the same. However, how it was obtained as well as its influence on prisoner A's way of thinking generated distinct posterior distributions. We do not know how prisoner A designed his experiment, nor do we know which procedure he used to construct his posterior distribution. Therefore, we cannot affirm that A made a mistake declaring  $1/2$  as the posterior probability of his being the survivor.

#### ACKNOWLEDGEMENTS

R. H. Loschi received partial financial support from CNPq (grants 300325/2003-7 and 472066/2004-8), Brazil, and P.L. Iglesias by FONDECYT, Chile (grant 8000004).

#### REFERENCES

- Dawid, A. P.(1985). The impossibility of inductive inference. *Journal of the American Statistical Association*, 80(390), 340-341.
- De Finetti, B. (1972). *Probability, Induction and Statistics*. New York: John Wiley and Sons Ltd.
- De Finetti, B.(1937). La prévision: Ses lois logiques, ses sources subjectives. *Annals of the Institute Henri Poincaré* 7, Paris.
- Diaconis, P. and Zabell, S. L.(1982). Updating subjective probability, *Journal of the American Statistical Association*, 77(380), 822-830.
- Goldstein, M. (1985). Temporal coherence. In J. M. Bernardo, M. H. de Groot, D. V. Lindley and A. F. M. Smith (Eds.), *Bayesian Statistics 2*, (pp. 231-148). North Holland: Elsevier Science.
- Howson, C. and Urbach, P.(1993). *Scientific Reasoning: The Bayesian Approach* (2nd edition). Chicago: Open Court.
- Jeffrey, R. (1992). *Probability and the Art of Judgement*. Cambridge: Cambridge University Press.
- Loschi, R. H., Iglesias, P. L. and Arellano-Valle, R. B. (2002). Conditioning on uncertain events: Extensions to Bayesian inference. *Test*, 11(2), 365-382.
- Morgan, J. P., Chaganty, N. R., Dahiya, R. C. and Doviak, M. J. (1991). Let's make a deal: The player's dilemma. *The American Statistician*, 45(4), 284-289.