

# Two Envelopes and Binding

Casper Storm Hansen<sup>a</sup>

<sup>a</sup>The Polonsky Academy for Advanced Study in the Humanities and Social Sciences, The Van Leer Jerusalem Institute

## ABSTRACT

This paper describes a way of defending a modification of Eckhardt's [2013] solution to the Two Envelopes Paradox. The defence is based on ideas from Arntzenius, Elga, and Hawthorne [2004].

**KEYWORDS** two envelopes paradox; binding; decision theory

## 1. Introduction

In their paper 'Bayesianism, Infinite Decisions, and Binding' Arntzenius, Elga, and Hawthorne [2004] provided answers to a long list of decision-theoretic puzzles, with the aim of giving a coherent and unified account of this family of puzzles based on a small number of principles. I will argue that, in the case of the Two Envelopes Paradox of Broome [1995], they failed to apply their principles in the manner required to reach the most interesting conclusion that those principles imply. I will further suggest that, if they are applied the right way, the result is a modification of the answer to the Two Envelopes Paradox that was formulated by Eckhardt [2013]. Eckhardt, in turn, did not, in my opinion, argue persuasively for his answer. This paper aims to do better in that regard by creating a synthesis of the ideas of Eckhardt on the one hand and Arntzenius, Elga, and Hawthorne on the other.

Here is one way of formulating the Two Envelopes scenario and the problem it gives rise to. Two envelopes, marked 'A' and 'B' respectively, contain one voucher apiece for a number of utils, redeemable by a divine agency. One envelope contains twice as many utils as the other, and the amounts have been decided by chance. There is a probability of  $\frac{1}{4}$  that one of the amounts is 1 and the other 2, and that probability is divided evenly between A and B containing the larger. Of the remaining  $\frac{3}{4}$  probability,  $\frac{1}{4}$  is assigned to the possibility that one of the amounts is 2 and the other 4, and that probability is also divided evenly between A and B containing the larger amount. Of the remaining  $(\frac{3}{4})^2$  probability,  $\frac{1}{4}$  is assigned to the possibility that one of the amounts is 4 and the other 8, and again that probability is divided evenly between A and B containing the larger; and so on. A player is given envelope A, allowed to look inside it to learn the value,  $a$ , of its voucher, and then asked whether he would like to exchange it for envelope B and its  $b$  utils. If he learns that  $a = 1$ , then he will know that  $b = 2$ , so he will want to swap. And he will also want to swap if  $a$  is any other amount, because there is an expected gain from doing so: although there is a slightly larger chance that he will lose utils instead of gaining more, namely  $\frac{4}{7}$  versus  $\frac{3}{7}$ , his potential gain of  $a$  extra utils is twice as big as his potential loss of  $\frac{a}{2}$  utils, making for a positive expected gain of  $\frac{3}{7} \cdot a - \frac{4}{7} \cdot \frac{a}{2} = \frac{a}{7}$ . Hence, it seems to be desirable for him to swap, irrespective of the content of envelope A. However, that conclusion conflicts with a very simple piece of symmetry reasoning: that, because the content of B has the same probability distribution as the content of A, it cannot be desirable to swap for all values of  $a$ .

In section 2, I will present Eckhardt's answer to this paradox, a summary of his arguments, and my critique thereof. Section 3 rehearses the relevant parts of Arntzenius, Elga, and Hawthorne's paper. And in section 4, I present my argument in favour of a variation on Eckhardt's answer, based on Arntzenius et al.'s principles.

While I will improve upon Eckhardt's thesis and its justification, I will not wholeheartedly endorse the improved thesis. Although I will argue that we seem to be forced to accept the thesis, there is something unsatisfying and puzzling about it. I will explain my worries in sections 5 and 6.

## 2. Eckhardt

Some contributors to the discussion of Two Envelopes have claimed that the player ought to swap after learning the value of  $a$ , no matter what it is, in spite of the symmetry argument [Scott and Scott 1997; Arntzenius and McCarthy 1997]. Others have held that he does not have a reason to swap (except when  $a = 1$ ), in spite of the positive expected gain [Clark and Shackel 2000]. Eckhardt [2013: ch. 8 and sec. 9.6] took a third position: that swapping for all values of  $a$  and keeping for all values of  $a$  are equally bad strategies, and both inferior to any strategy that involves swapping in the case of 'low' values of  $a$  and keeping in the case of 'high' values.

Denoting by ' $M(2^n)$ ' the strategy of swapping A for B if  $a \leq 2^n$  and keeping A if  $a > 2^n$ , Eckhardt calculated the expected gain from following  $M(2^n)$  compared to always keeping A to be  $(\frac{3}{2})^n/8$ , a result arrived at simply by taking the weighted average over a finite number of possibilities.<sup>1</sup> This means that  $M(2^n)$  is a better strategy than  $M(2^m)$  when  $n > m$ . However, he also maintained that the limit of the sequence  $M(2^0), M(2^1), M(2^2), \dots$  of better and better strategies, that is, the strategy of swapping in all cases, is just as bad as never swapping. In other words, there is no best strategy, so what the player ought to do is pick some large  $n$ , more or less arbitrarily, and then decide to play according to the  $M(2^n)$  strategy.

The obvious challenge for this view is the difficulty of reconciling the implication that if the player decides on  $M(2^n)$  and subsequently learns that  $a > 2^n$ , he should not swap, with the fact that there is a well-defined and positive expected gain from doing so when the value of  $a$  is known. In my opinion, Eckhardt fails to explain this in a convincing manner, and the remainder of this section is devoted to that topic.

Part of Eckhardt's justification consists of alleged experimental corroboration. He had a computer play the Two Envelopes game a thousand times, using the  $M(2^5)$  strategy, the always-swap strategy, and the never-swap strategy. The theory's prediction for  $M(2^5)$  is a gain of 949, and the experimental result was 881. That, I agree, is close enough to count as confirmation. But Eckhardt further claimed that the experiments confirmed his thesis that the two other strategies are equally good (or bad); and that, I do not believe. He did not elaborate on this latter claim, neither presenting any concrete data nor explaining in what sense he took the data to constitute confirmation. That leaves us with no option but to try to find such confirmation in simulations we run ourselves. And when we do, we instead find that the results are – in a very peculiar way – systematically inconclusive.

---

<sup>1</sup> Insert  $r = \frac{3}{4}$  in the result on page 67 of Eckhardt [2013].

Table 1 shows the results of one series of 10,000 experiments: the average gain from swapping A for B, both separately for each value of  $a$  and the overall average. In total, I completed 10 series of 1,000 experiments, 10 series of 10,000, 10 series of 100,000, and 10 series of 1,000,000, and this is what can be observed:

- I. The average gain for individual values of  $a$  is reliably positive for small  $a$ .
- II. When the number of experiments in a series increases, larger values count as ‘small’ for the purpose of observation I.
- III. The overall average gain in a series is as likely to be positive as negative.<sup>2</sup>
- IV. The overall average gain in a series tends to be heavily influenced by one or a handful of extreme outcomes.
- V. The absolute value of the overall average gain tends to increase with the sample size.
- VI. (A corollary of item V) Taking the average of the average gains of a number of series with the same sample size tends to yield a result that is numerically larger than the majority of the averages for the individual series.

Observations I and II indicate that, for any value of  $a$ , there is a particular number of experiments such that, if we had run a series with that many experiments, we would with near certainty have obtained empirical-statistical evidence for the desirability of swapping upon observing the content of A to be  $a$ . Similarly, we could, for any value of  $b$ , have obtained such evidence for the desirability of not swapping upon observing the content of B to be  $b$ .

The remaining items in the list are relevant to my objection to Eckhardt. Observation III is the only feature of the dataset that could *prima facie* be taken as an indication that always swapping is no better than always keeping. However, it is easy to see that III does not, by itself, justify that conclusion. Rather, III is consistent with there being a probability of 0.5 that the average gain from always swapping is 1,000,000, and a probability of 0.5 that it is -1, in which case it would be preferable to always swap. Only averages could, potentially, be used to justify Eckhardt’s conclusion using empirical evidence. And observations IV, V, and VI indicate that such justification cannot be expected from statistical investigation, however extensive and thorough it may be: we would not get the convergence to an average gain of 0 that would confirm the thesis.

The problem is that what normally justifies the use of statistical-empirical methods is absent here. Normally, when statistics constitute evidence for an unobserved event, it is by virtue of the Law of Large Numbers, which states that the average of the results from a number of trials will tend to approximate the theoretical mean better and better as the sample size increases. But the Law of Large Numbers does not apply to the Two Envelopes game; see Norton [1998: 50].<sup>3</sup>

So, justification of a more theoretical nature is needed. Eckhardt also attempted to supply that, with arguments revolving around these three propositions:

---

<sup>2</sup> I observed 21 positive and 19 negative.

<sup>3</sup> Using an empirical-statistical method is justified for the purpose of items I and II, because if, for each  $n \in \mathbb{N}_0$ , we define the stochastic variable  $b - a_n$  to be equal to  $b - a$  if  $a = 2^n$  and equal to 0 otherwise, and similarly for  $a - b_n$ , the Law of Large Numbers holds for each  $b - a_n$  and  $a - b_n$ . This follows from the fact that these stochastic variables have well-defined, finite expected values; see Feller [1968: 260].

- EACH CASE (EC): For each value of  $a$  considered individually, it is favourable to swap.
- SYMMETRY (SYM): The strategy of swapping for all  $a$  is exactly as (un)favourable as the strategy of keeping for all  $a$ .
- ALWAYS (AL): The strategy of swapping for all  $a$  is more favourable than the strategy of keeping for all  $a$ .

He accepts EC and SYM but denies AL, and he diagnoses the paradox as arising from uncritical acceptance of  $EC \rightarrow AL$ . His main argument against this conditional is as follows:

EC and SYM are well established through independent arguments (expected value calculations and the symmetric ignorance principle,<sup>4</sup> respectively) while  $EC \rightarrow AL$  places a wedge of contradiction between them. Moreover the only reason to believe AL is the combination of EC and  $EC \rightarrow AL$ . Therefore, EC and SYM are true; AL and  $EC \rightarrow AL$  are false. [Eckhardt 2013: 51]

What we have here are three propositions, namely EC, SYM, and  $EC \rightarrow AL$ , each of which seems obviously true, but which together form an inconsistent set. Starting from any two of them, one can ‘prove’ that the third is false. But doing so would only be useful if one had good reasons to think that those two were the true propositions and the third, the false one; and if one did have such reasons, there would not be a problem in the first place. Relying on such a ‘proof’ is to underestimate the challenge of the special dialectical situation that one is in when trying to solve a paradox.

Of course, Eckhardt does present arguments in favour of EC and SYM, namely those referred to in the parentheses in the quoted passage above. But equally plausible arguments can be given for  $EC \rightarrow AL$  (I will come to one in a moment). What can bring the discussion forward in this kind of special dialectical situation is not an argument *for* some of the obvious-seeming propositions, but an argument *against* the proposition that one wants to reject. And that argument cannot just be the *reductio* based on the other contentious propositions, for this is just to wield the big stick instead of offering an explanation, as Dummett [1991: 316] so beautifully put it.

So what *reductio*-independent arguments does Eckhardt give? I can only discern one, and it is highly underdeveloped. He writes [2013: 55] that ‘ $EC \rightarrow AL$  is based on a fallacy of composition’.

*Prima facie*, it does not seem right to claim that to infer from ‘for each value of  $a$  it is favorable to swap’ to ‘for all values of  $a$  it is favorable to swap’ is a fallacy. Inferring from ‘for each  $x \in X, P(x)$ ’ to ‘ $P(X)$ ’ is not a fallacy if  $P$  is a distributive predicate. For instance, it is a valid inference to proceed from ‘each of the two persons is a philosopher’ to ‘the two persons are philosophers’. That some persons are philosophers does not mean anything over and above that they are each a philosopher. And the favourability of swapping likewise seems to be a distributive property, as it is not clear what the truth criterion is for ‘for all values of  $a$  it is favourable to swap’,

---

<sup>4</sup> The symmetric ignorance principle is defined in Eckhardt [2013: 48] as follows: ‘an agent is symmetrically ignorant with respect to the choice of two options, if everything the agent knows about either option applies equally to both of them [...] The symmetrical ignorance principle states that symmetrical ignorance forestalls rational preference.’

over and above that it is favourable to swap for each value of  $a$ . They amount to the same thing when the player has opened envelope A: he should swap!

Part of the problem is that the exact meaning of EC is very unclear. How should ‘considered individually’ be understood? After having opened the envelope, the player is looking at *one* value of  $a$ , so whether he considers the (previously) possible values of  $a$  ‘individually’ or ‘collectively’ seems to be a distinction without a difference.

Let us say that, having read Eckhardt’s book and found nothing in it to disagree with, the player has decided to follow one of the discriminating  $M(2^n)$  strategies. If he opens envelope A and finds that  $a > 2^n$ , he nevertheless knows of an argument with the conclusion that he should swap, and premises that he cannot bring himself to deny—even though the argument shows them to contradict Eckhardt’s position, with which he sympathizes. The first premise is that the expected gain from swapping is positive, while the expected gain from keeping A is 0, and that these are his only two options. When the value of  $a$  is known, all of this is undeniable. The second premise is that whenever one is faced with a choice between finitely many options, all of which have finite expected gains, one ought to choose the option with the highest expected gain (or one of the options with maximum expected gain, if there are several). This is the core principle of ordinary, finite decision theory, and Eckhardt seems committed to denying it, despite his statement that the Two Envelopes problem ‘requires not innovation but application of long accepted principles’ [2013: 55]. I would suggest that it does require quite a bit of innovation. This is where Arntzenius, Elga, and Hawthorne come in.

### 3. Arntzenius, Elga, and Hawthorne

Arntzenius et al.’s conclusions concerning the Two Envelopes case are products of an argumentative strategy that can be characterized as ‘first up, then down’. That is, they first argue for a number of abstract principles on the basis of some simple example cases, and then apply those abstract principles to more complicated cases, Two Envelopes among them. These abstract principles are innovative, not long accepted, and important to Two Envelopes.

One of the simple cases features Eve in the Garden of Eden. Satan has cut an apple into infinitely many pieces, one for each of the natural numbers. In a supertask, Eve is first asked if she would like to have piece no. 1, then if she would like to have piece no. 2, and so on. Any combination of positive and negative answers will be accommodated by Satan. Eve would like to have as many pieces as possible, but iff she takes an infinite number she will be ejected from the Garden; and staying there is more important to her than any amount of apple. The problem is that for each question, it is favourable for Eve to answer affirmatively: the answer to an individual question does not have an impact on whether she stays in the Garden, and whether she does or not, having that extra piece of apple is preferable to not having it. Yet, if she answers all the questions in the affirmative, she ends up with a very unfavourable result.

That is the diachronic version of ‘Satan’s Apple’. In the synchronic variant, Eve makes her decisions regarding all the pieces at the same time. The synchronic version is so simple that there is little room for disagreement about what Eve ought to do, in spite of the fact that this recommendation is somewhat surprising. Eve should accept

some finite (but large) number of pieces. This is obvious because the alternative—accepting infinitely many—is a worse choice. But it is a little surprising that the rational thing to do is to accept some option when there is another option (taking one more piece) that is better.

A lesson that can be learned from the synchronic version is that the principles of decision theory have to be applied to the choice between alternative complete strategies, and not to elements thereof in a piecewise fashion. If Eve considers what she should do with regard to any single piece of apple, the verdict of standard decision theory applied in isolation is that she should take it. But looking at the big picture, it is rational to say ‘no’ to some (indeed most) of them.

Arntzenius et al. draw an additional moral from the diachronic version of Satan’s Apple. The shift to diachronicity makes no difference to how good or bad the different strategies are, so again Eve ought to take a finite number of pieces and reject the rest. However, at any given instant in which she is to take a decision, if she is sure that her decision at that time will have no causal influence on the other individual choices she has to make, the rational choice is—for the reason given in the second paragraph of this section—to accept the piece in question. To achieve a good result, Eve needs to decide on a complete strategy from the beginning and then stick to it. She must *bind* herself to a plan and not let rationality override it afterwards, because acting fully rationally will, in the absence of this binding of her own future self, lead to her fall.

So Arntzenius et al.’s second lesson is that an agent who can bind herself to a plan can achieve a better result than someone who at every instant acts according to the verdicts of rationality at that instant.<sup>5</sup>

Let us return to the Two Envelopes Paradox. Arntzenius et al. consider this variation of the scenario that was presented in section 1: Swapping from envelope A to envelope B will cost the player a fee of 0.01 utils. If he decides to swap, his memory of the value of  $a$  will be deleted; he will be allowed to learn the value of  $b$ ; and he will be given the option of swapping back to envelope A in exchange for another 0.01 utils.

No matter what the player finds in envelope A, there will be an expected gain from swapping despite the price tag, so ‘in-the-moment rationality’ prescribes that he swap. And the same holds for swapping back. But in combination, this results in a guaranteed loss of 0.02 utils, so not swapping at all is the better overall strategy. Applying their principles, Arntzenius et al. conclude, first, that this conflict should be resolved in favour of the big-picture view; and second, that to deal with the situation in the best possible way, the player ought (if he can) to bind himself from the outset to a course of action which is in line with that recommendation: that is, which prevents the verdicts of his rationality from influencing his actions at the time when he knows the value of  $a$  and at the time when he knows the value of  $b$ .

---

<sup>5</sup> In the interest of brevity, I based this summary on only one of Arntzenius et al.’s examples. Another that is particularly relevant for comparison with Two Envelopes is ‘Trouble in St. Petersburg’, but I leave that to the reader.

## 4. Synthesis

Arntzenius et al.'s principles can support a much more interesting conclusion: Eckhardt's. Let us first use Arntzenius et al.'s insights to clarify the difference between EC and AL. We can identify a strategy in Satan's Apple with a subset of the set of natural numbers, namely the set of the numbers of the pieces of apple to be taken. This leads us to the following two propositions, which both enjoy the virtues of being precise and having uncontroversial truth values—true in the first case, false in the second:

- EC': Given any strategy  $N \subset \mathbb{N}$  and an  $n \notin N$ , the strategy  $N \cup \{n\}$  is better than  $N$ .
- AL':  $\mathbb{N}$  is the best strategy.

The strategies of Two Envelopes can also be identified with subsets of the set of natural numbers (if we now take that to include 0), namely the set of those numbers  $n$  such that the player swaps in cases where  $a$  equals  $2^n$ . If we interpret EC' and AL' as being about Two Envelopes, then we have made progress in comparison to Eckhardt's unclear formulations. In addition to rendering the difference between the two propositions clear, it can also be seen that AL' does not follow from EC' by logic and set theory alone, which is an important first step towards being able to maintain that  $EC \rightarrow AL$  is false.

Given these explications of EC and AL, the inference from the former to the latter is not a fallacy of composition. Both concern full strategies, and the 'sum' of two strategies that give opposite instructions concerning at least one  $n$  is an inconsistent strategy. So, what is really in play is a fallacy of limits: inferring from the premise that each element of some infinite sequence has a given property (and no further premises) to the conclusion that the limit also has that property.

The first lesson we took from Arntzenius et al. was that the principles of decision theory should be applied to complete strategies. And complete strategies are what EC' and AL' are about. Having absorbed this lesson, we can consistently believe that AL' is false, and yet that for each value of  $a$  taken in isolation, it is preferable for the player to swap, for the latter is not an assertion about complete strategies. We are one step closer to accepting the falsity of  $EC \rightarrow AL$ . For while Eckhardt left us with the impression that accepting his claim committed us to rejecting ordinary, finite decision theory (cf. the final paragraph of section 2 above), Arntzenius et al. lessen the cost. What we have to revise is merely our beliefs about how those principles should be applied: an agent who is able to bind herself should apply the principles to complete strategies, not to sub-strategies.

Binding was the subject of the second lesson from Arntzenius et al. When the player looks into envelope A it will become rational for him to swap no matter what he sees there, because there is, in that moment, an expected gain from doing so. But if, before opening the envelope, he does not prevent his future self from trying to seize *any* such opportunity for a gain, his present self cannot expect a gain. So the truth of EC', which leads a rational but non-bound agent who knows the content of envelope A to swap, does not imply AL', the truth of which would imply that an agent who does not know the content of A ought to swap for an envelope with an identical probability distribution.

EC and SYM have theoretical support: an argument about finite and well-defined expected gain in the case of EC, and an argument about identity of probability distributions in the case of SYM. However, that theoretical support is impotent when contrasted against what seems to be equally strong if not stronger theoretical support for  $EC \rightarrow AL$ , namely, the argument that EC and AL amount to the same thing. That support has now been undermined, using Arntzenius et al.'s principles.

My formulation of EC' above should not be interpreted as an exegetical thesis about what Eckhardt *really* meant with EC. As far as I can discern, Eckhardt was sliding back and forth between several different propositions, all squeezed together under a single label. Below, I list four different ways that EC *could* be interpreted (though again, no exegetical thesis is being advanced):

- For each possible value of  $a$ , the expected gain from swapping, conditional on that value of  $a$ , is positive.
- EC'.
- For each possible value of  $a$ , if and when the content of envelope A is revealed to be that value, it is rational to swap.
- For each possible value of  $a$ , it is irrational to commit to any strategy that implies keeping envelope A, if and when it is revealed to contain that value.

The last option does indeed imply AL; so, assuming SYM and the negation of AL, it has to be separated from whichever of the three first options that one wants to assert. But Eckhardt did not develop the theory needed to do so in a clear and intelligible way, and that is why Eckhardt's answer to the paradox has to be modified in the course of justifying it on the basis of Arntzenius et al.'s principles. The modified answer involves (1) accepting the first three propositions listed above, (2) rejecting the last one, and (3) distinguishing carefully between (i) whether an action  $\alpha$  is rational at an instant of time  $t_2$  and (ii) whether it is rational for the agent in question to bind himself at an earlier instant of time  $t_1$  against performing  $\alpha$  at  $t_2$ , a distinction which is absent from Eckhardt's answer.

Having cleaned up the conceptual confusion, we can formulate a solution. Arntzenius et al. managed to provide a clear analysis of Satan's Apple by first addressing the simple synchronic version and then the more complicated diachronic version. Let us follow their lead. The analogue of the synchronic version, in the case of Two Envelopes, is a version of the scenario in which the player decides on a strategy and receives envelope A or B in accordance with the strategy, instead of being asked to make his decision after having seen the content of A. This simplifies the puzzle, because there is nothing to consider except the relative merits of strategies. And it seems clear that the always-swap and never-swap strategies are equally bad (because of their identical probability distributions); that  $M(2^n)$  is better than  $M(2^m)$  for all  $n > m$  (because of the well-defined expected gain); and that it is rationally permissible to choose a strategy for which a better alternative exists (because that is true of every strategy and the player must choose one).

The analogue of the diachronic version of Satan's Apple, then, is the original version of our scenario. The difference between the synchronic and diachronic versions has no effect on the merits of the various strategies in the case of the Two Envelopes, any more than it does in the case of Satan's Apple. So the conclusion from the previous paragraph still holds: it is rational to follow, for example, the  $M(2^5)$  strategy, and better to do so than to follow the always-swap strategy. When the player is informed



of the value of  $a$ , it becomes rational for him to swap (because of the positive expected gain), but if he allows himself to do that irrespective of whether doing so is in conformity with  $M(2^5)$ , he is really following the always-swap strategy. Therefore, it is better for the player to bind his future self to the  $M(2^5)$  strategy before he learns the value of  $a$ , than to allow himself to act in conformity with his rationality later. That is the modification of Eckhardt's position that follows from Arntzenius et al.'s principles.

In section 2, I noted that a problem for Eckhardt's view is that it is difficult to reconcile the implication that when the player has settled on an  $M(2^n)$  strategy and subsequently learns that  $a > 2^n$ , he should not swap, with the fact that there is a well-defined and positive expected gain from doing so whenever the value of  $a$  is known. In my modification, this contradiction has been dissolved via disambiguation with respect to the time index: after looking in the envelope, the player ought to swap according to 'in-the-moment rationality'; but earlier, he should have ensured that he *cannot* swap, because it is better for him to act according to the bird's-eye view available to him before he learns that  $a > 2^n$ .

Let me make it explicit what Eckhardt brings to the table that Arntzenius et al. miss. The latter's insights about binding can only be applied in scenarios where multiple decisions are to be made. In the case of Two Envelopes, these insights can be applied at two levels. The first is a case in which multiple decisions about whether to swap must be made. That is what Arntzenius et al. consider: should the player exchange A for B, and after doing so, exchange back? A player with self-binding capability who has read their paper can utilize the knowledge thereby gained to avoid paying a fee to maintain the status quo, by preventing himself from swapping back after having accepted the first swap. But why not accept that first swap? Arntzenius et al. [2004: 273] claim that he should not, but that is just the intuitive verdict, and not something that follows from their analysis. To reach that conclusion, the principles concerning binding must be applied at the second level: a single decision about whether to swap *analyzed* as an infinity of decisions about whether to swap conditional on each of the different possible values of  $a$ .

## 5. Incomparable Strategies

As mentioned above, each subset of  $\mathbb{N}$  corresponds to a strategy that the player can follow. Of these, Eckhardt only considers a few, namely the ' $M$ -strategies'; the strategies of only swapping for a single value of  $a$ ; and the strategy of swapping for all values of  $a$ . By restricting his treatment to those strategies, he hides an important issue.

The expected gain from following one strategy compared to another strategy is well-defined iff the two corresponding sets of natural numbers have a finite symmetric difference; that is, iff the set of elements belonging to one but not both of the two given sets is finite. Therefore, the principle that strategy  $S_1$  is better than/equally as good as/worse than  $S_2$  if  $S_1$  has positive/zero/negative expected gain compared to  $S_2$  ensures that all the strategies discussed by Eckhardt can be so compared. That is, except for the always-swap strategy, but for that one, he in effect adds another principle: if two strategies have the same outcome probability distribution, then they are equally good. That allows for direct comparison of the  $M(2^0)$  and always-swap strategies, and indirect comparison of the always-swap strategy with all of Eckhardt's

other strategies. However, in the sense of ‘direct comparison’ just used, this extra principle only applies directly to that one pair of strategies and to nothing else.<sup>6</sup>

That means that many pairs of strategies are incomparable. To be precise, the set of strategies identified with the power set of  $\mathbb{N}$  is divided into equivalence classes of comparable strategies, and two strategies are in the same equivalence class iff they have a finite symmetric difference or one of them has a finite symmetric difference with  $\emptyset$  (that is,  $M(2^0)$ ) and the other has a finite symmetric difference with  $\mathbb{N}$  (that is, the always-swap strategy).<sup>7</sup>

It might be possible to add principles, on top of the two already mentioned, that would place more pairs of strategies into the same equivalence class. However, it is difficult to see how one could come up with enough plausible principles to collapse it all down to just one equivalence class. And short of that, there is a significant bullet to bite if one wants to give the Eckhardian answer to the Two Envelopes Paradox explicated in this paper: some pairs of strategies are rationally incomparable. Should I prefer to always keep envelope A, or swap exactly when the contents of envelope A is  $2^n$  with even  $n$ ? No answer!

## 6. The Pull of Paradoxes

A fully satisfying solution to a paradox not only removes a contradiction in favour of a well-argued answer to the issue that the paradox concerns, but also provides the philosophical therapy needed to dissolve the psychological *pull* of the paradox. By that standard, I do not think that a fully satisfying solution to the Two Envelopes Paradox has been achieved here. The well-argued removal of the contradiction can be summarized as follows: there is only one reasonable stance one can take regarding Satan’s Apple; that stance points to some necessary revisions of the principles of decision theory; and those revisions remove the contradiction in the case of Two Envelopes. However, the stance regarding Satan’s Apple is an anti-rationality stance (involving a recommendation to bind oneself against the verdicts of one’s future self’s rationality, and implying the existence of rationally incomparable strategies). And even for someone who accepts the unavailability of the stance, that can be hard to swallow.

This, I would suggest, is because Arntzenius et al. fail to provide therapy: they convince us *that* following the verdict of rationality at a given point in time can be a bad thing, but they do not provide us with the insight needed to fully understand *how*

---

<sup>6</sup> Proof: For any  $n > 0$ , the probability that the player ends up with  $2^n$  utils after having followed a strategy depends only on what that strategy prescribes regarding the situations where  $a$  is  $2^{n-1}$ ,  $2^n$ , or  $2^{n+1}$ . For each of those three situations, the player can either swap or keep, which gives eight combinations. Calculating the probability of ending up with  $2^n$  utils for each of those combinations will show that they are all different, except for the combination of swapping for all three and the combination of keeping for all three, which both yield a probability of  $2^{-1}(3^n \cdot 4^{-n-1} + 3^{n-1} \cdot 4^{-n})$ . Something similar holds for  $n = 0$ , even though that case only gives rise to four combinations: only swapping for both  $a = 2^0$  and  $a = 2^1$  and keeping for both give the same probability of ending up with  $2^0$  utils. From this, it can be seen that the only two different strategies that have the same probability for *all* possible outcomes are the always-swap and the always-keep strategies; all other strategy-pairs will differ with respect to at least one possible outcome’s probability.

<sup>7</sup> This implies that for every strategy  $N \subset \mathbb{N}$  there is a better strategy and there is a worse strategy: if  $N \neq \mathbb{N}$  and  $n \in \mathbb{N} \setminus N$ , then  $N \cup \{n\}$  is better; if  $N = \mathbb{N}$  and  $n \in \mathbb{N}$ , then  $\{n\}$  is better; if  $N \neq \emptyset$  and  $n \in \mathbb{N}$ , then  $N \setminus \{n\}$  is worse; if  $N = \emptyset$  and  $n \in \mathbb{N}$ , then  $\mathbb{N} \setminus \{n\}$  is worse.

that can be: *how can it be correct that there is something wrong with ideal rationality?!* The lack of an answer to this question means that there is still a strong pull towards thinking that, when the player learns the value of  $a$ , any binding that might have been imposed to prevent him from swapping was a mistake. To remove that pull, it would seem that a much deeper analysis of the nature of rationality will be needed.

## Acknowledgements

I am grateful to Federico Luzzi and William Eckhardt for helpful discussion and to the editor and reviewers for useful suggestions.

## REFERENCES

Arntzenius, F., A. Elga, and J. Hawthorne 2004. Bayesianism, Infinite Decisions, and Binding, *Mind* 113/450: 251–83.

Arntzenius, F. and D. McCarthy 1997. The Two Envelope Paradox and Infinite Expectations, *Analysis* 57/1: 42–50.

Broome, J. 1995. The Two-Envelope Paradox, *Analysis* 55/1: 6–11.

Clark, M. and N. Shackel 2000. The Two-Envelope Paradox, *Mind* 109/435: 415–42.

Dummett, M. 1991. *Frege: Philosophy of Mathematics*, London: Duckworth.

Eckhardt, W. 2013. *Paradoxes in Probability Theory*, Dordrecht: Springer.

Feller, W. 1968. *An Introduction to Probability Theory and Its Applications (3rd ed.)*, Volume 1, New York: John Wiley and Sons.

Norton, J.D. 1998. When the Sum of Our Expectations Fails Us: The Exchange Paradox, *Pacific Philosophical Quarterly* 79/1: 34–58.

Scott, A.D. and M. Scott 1997. What's in the Two Envelope Paradox? *Analysis* 57/1: 34–41.

| $a$      | freq. | avg. $b - a$ | $a$      | freq. | avg. $b - a$  |
|----------|-------|--------------|----------|-------|---------------|
| $2^0$    | 1234  | 1.00         | $2^{15}$ | 37    | 2214.05       |
| $2^1$    | 2286  | 0.31         | $2^{16}$ | 29    | 24858.48      |
| $2^2$    | 1616  | 0.56         | $2^{17}$ | 18    | 10922.67      |
| $2^3$    | 1228  | 1.05         | $2^{18}$ | 9     | 43690.67      |
| $2^4$    | 885   | 2.52         | $2^{19}$ | 9     | 262144.00     |
| $2^5$    | 691   | 4.84         | $2^{20}$ | 10    | 262144.00     |
| $2^6$    | 521   | 8.72         | $2^{21}$ | 9     | 0.00          |
| $2^7$    | 355   | 27.40        | $2^{22}$ | 5     | 2936012.80    |
| $2^8$    | 301   | 42.95        | $2^{23}$ | 2     | -4194304.00   |
| $2^9$    | 240   | 108.80       | $2^{24}$ | 2     | 4194304.00    |
| $2^{10}$ | 180   | 110.93       | $2^{25}$ | 2     | -16777216.00  |
| $2^{11}$ | 117   | 341.33       | $2^{26}$ | 1     | 67108864.00   |
| $2^{12}$ | 87    | 776.83       | $2^{27}$ | 0     | —             |
| $2^{13}$ | 70    | 1345.83      | $2^{28}$ | 1     | -134217728.00 |
| $2^{14}$ | 54    | 455.11       | $2^{29}$ | 1     | -268435456.00 |
|          |       |              | tot.     | 10000 | -34773.60     |

Table 1: Frequency of different values of  $a$  in a series of 10,000 experiments and the average gain achieved by swapping from A to B