

# Trust the Evidence: Two Deference Principles for Imprecise Probabilities

Giacomo Molinari

GIACOMO.MOLINARI@BRISTOL.AC.UK

Department of Philosophy, University of Bristol, UK

## Abstract

Our intuition that rational agents should value the evidence can be captured by a well-known theorem due to I. J. Good. However, Good’s theorem fails when agents have imprecise credences, raising the worry that agents with imprecise credences don’t value the evidence. This essay shows a different way to capture our starting intuition, as the claim that rational agents defer to their informed selves. I introduce and motivate two deference principles for imprecise probabilities, and show that rational imprecise agents defer to their informed selves according to these principles. This shows a sense in which imprecise agents value the evidence. I end by comparing the deference principles introduced here with an alternative from the literature.

**Keywords:** imprecise probability, deference principles, reflection, Good’s theorem

## 1. Introduction

We like to think that rational agents *value the evidence*. Precise Bayesians can capture this intuition in terms of sequential decision-making by appealing to a theorem due to I. J. Good [8]. The theorem shows that, if a rational agent is offered the option to learn some new evidence for free before facing a decision problem, they are not willing to pay to turn down this offer. In other words: rational agents never pay to avoid free evidence.

Good’s theorem fails if we allow agents to have imprecise credences. Rational imprecise agents are sometimes allowed to pay in order to avoid free evidence, and depending on the imprecise decision theory one picks, they may even be required to do so [2, 10]. This raises the worry that the coherence and decision-making norms of imprecise probability are somewhat faulty: agents who follow them don’t value the evidence, and are therefore not rational.

In this essay I will respond to this worry by showing a different way to capture our starting intuition that rational agents value the evidence, which appeals to a *deference principle*. Deference principles tell us, given an agent’s credal state, which credal states they consider as experts. The intuition that rational agents value evidence can be captured as the claim that rational agents defer to their informed selves, treating their informed selves like an epistemic authority.

I will introduce and defend two deference principles for imprecise credences, called Strong and Weak Total Trust. I will then show that imprecise agents defer to their informed selves according to these principles, thus showing a way in which imprecise agents can be said to value the evidence, independently of any theory of sequential decision-making. I end the essay by comparing Strong and Weak Total Trust to an alternative deference principle for imprecise credences that has been proposed in the literature, called Identity Reflection. Identity Reflection is shown to be strictly stronger than Strong Total Trust. In particular, Identity Reflection does not allow an agent to defer to a modest expert, whereas this is possible under Strong Total Trust.

## 2. Some Notation

A *finite probability space* is a pair  $(\Omega, p)$  where  $\Omega$  is a finite possibility space, and  $p : 2^\Omega \rightarrow \mathbb{R}$  is a probability function. I will assume throughout that the possibility space  $\Omega$  is finite, and will speak of a probability function instead of a probability space whenever there is no risk of confusion regarding the function’s domain. I use  $\mathcal{P}_\Omega$  to describe the set of all probability functions on a given  $\Omega$ .

We model an agent’s individual probabilistic judgements as sets of probability functions. For example, let  $H$  be the event that a coin lands heads, and  $T$  the event that it lands tails. Then the judgement that the coin is fair can be captured by the set  $\{p \in \mathcal{P}_\Omega : p(H) = 1/2\}$  of probability functions which assign probability  $1/2$  to  $H$ . The judgement that the coin is biased towards heads can be captured by the set  $\{p \in \mathcal{P}_\Omega : p(H) > 1/2\}$  of all functions which assign greater probability to  $H$  than to  $T$ .

We model an agent’s entire doxastic state by a single nonempty set  $P \subseteq \mathcal{P}_\Omega$ , known as the agent’s *credal set*. The idea is that an agent makes a probabilistic judgement, such as the judgement that a coin is biased towards heads, iff every probability in  $P$  makes that judgment, meaning that  $P \subseteq \{p \in \mathcal{P}_\Omega : p(H) > 1/2\}$ . More generally, an agent makes a probabilistic judgement iff their credal set  $P$  is contained in the set of probability functions corresponding to that judgement.

In this essay I will restrict myself to *regular* credal sets, that is, credal sets whose members assign some positive

probability to each possibility in their domain  $\Omega$ . As pointed out later on, this assumption considerably simplifies the relationship between an agent's credal set and the set of gambles they find desirable.

When  $P \subseteq \mathcal{P}_\Omega$  is a credal set and  $A \subseteq \Omega$  an event, I write  $P(A)$  to denote the value set  $\{r \in \mathbb{R} : (\exists p \in P)p(A) = r\}$ , and I denote by  $P(\cdot|A)$  the following conditional credal set:

$$P(\cdot|A) = \{p(\cdot|A) : p \in P\}$$

which is defined whenever  $A \neq \emptyset$ , under the assumption that  $P$  is regular.

### 3. The Value of Evidence: Sequential Choice Characterisation

We like to think that rational agents value the evidence. In this section, I look at one way we may capture this intuition: by requiring that rational agents never pay to avoid learning free evidence. A famous theorem due to I.J. Good shows that agents with precise probabilistic credences obey this requirement. However, for agents with imprecise credences, it is sometimes admissible to pay to avoid learning free evidence.

#### 3.1. Good's Theorem

Let  $\Omega = \{\omega_1, \dots, \omega_n\}$  be a finite possibility space. Consider an agent facing a decision problem  $\mathcal{A} = \{a_1, \dots, a_m\}$ . Let  $U : \mathcal{A} \times \Omega \rightarrow \mathbb{R}$  be the agent's utility function, so the utility for option  $a_j$  when  $\omega_i$  is the case is given by  $U(a_j, \omega_i)$ . Let  $\mathcal{E} = \{E_1, \dots, E_k\}$  be an arbitrary partition of events, such that learning which  $E_i$  is true does not affect the agent's utility assignment over  $\mathcal{A}$ .<sup>1</sup> Imagine that, at  $t = 0$ , the agent is offered the following choice: she can either pick some option from  $\mathcal{A}$  now (at  $t = 0$ ), or learn which  $E_s \in \mathcal{E}$  is true, and then pick some option from  $\mathcal{A}$  (at  $t = 1$ ). I call a scenario of this kind a *sequential learning problem*  $D(\mathcal{E}, \mathcal{A})$ .

Intuitively, if an agent values the evidence, she won't be willing to pay to avoid free evidence in a sequential learning problem of the kind described above. This suggests the following characterisation of what it means for an agent to value the evidence:<sup>2</sup>

- **Value of Evidence - Sequential Choice (VE-SC):** An agent values the evidence when for any sequential learning problem  $D(\mathcal{E}, \mathcal{A})$ , she is not willing to pay to avoid learning which  $E_i$  is true before choosing from  $\mathcal{A}$ .

<sup>1</sup>If learning the evidence alters the agent's utility function, then the evidence is not "free". For an in-depth discussion of this assumption, and of Good's theorem more generally, see [10].

<sup>2</sup>A similar characterisation is given in Dorst [4], where it is generalised to allow for cases where  $\mathcal{E}$  is not a partition.

Good's theorem [8] shows that for any sequential learning problem  $D(\mathcal{E}, \mathcal{A})$ , if the agent facing the problem has precise probabilistic credences, her expected utility for choosing from  $\mathcal{A}$  after learning which  $E_i$  is true is at least as great as her expected utility for choosing from  $\mathcal{A}$  without learning. Therefore, if she makes choices by maximising her expected utility, such an agent will not be disposed to pay to avoid learning. Hence, rational agents with precise probabilistic credences value the evidence according to (VE-SC).

#### 3.2. Imprecision and Sequential Choice

Trying to show that agents with imprecise credences value the evidence as described by (VE-SC), one is faced with a number of difficulties. First of all, while it's commonly assumed that precise Bayesian agents make choices by maximising expected utility, a number of different decision rules exist for imprecise agents [17]. Furthermore, while it's straightforward to extend expected utility maximisation to sequential problems, not all IP decision rules are so easily extended. Consider the following example:<sup>3</sup>

**Example 1 (Coin Toss Puzzle)** *Jack has a coin which you know is fair. You know that Jack knows whether  $A$  is true. You know nothing about  $A$ , but judge that whether  $A$  is true is independent of the result of the coin toss. Jack paints the two sides of the coin so you can't tell which one is heads. If  $A$  is true, he writes " $A$ " on the heads side, and " $\neg A$ " on the tails side. If  $A$  is false, he writes " $\neg A$ " on the heads side, and " $A$ " on the tails side.*

Let  $H$  be the event that the painted coin lands with the heads face up. Since you know the coin is fair, your starting credence in  $H$  should be  $1/2$ . That is, your credal set  $P$  should be such that for every  $p \in P$ ,  $p(H) = 1/2$ , i.e.  $P(H) = \{1/2\}$ . Since you know nothing about  $A$ , you can (and perhaps should) have maximally imprecise credence in  $A$ . That is,  $P(A) = (0, 1)$ . Furthermore, you judge the coin toss to be independent of  $A$ . That is, if we let  $E_A$  be the event that the coin lands with the face on which " $A$ " is painted facing up, then for every  $p \in P$  it should be  $p(E_A|A) = p(E_A) = 1/2$ .<sup>4</sup>

Consider what happens after observing the painted coin toss. If the coin were to land with the " $A$ " side up, then

<sup>3</sup>As we shall see in Section 5, this example was originally given by White [21] to illustrate a conflict between deference and credal dilation. A similar example is given by Walley [20, pp. 298-299].

<sup>4</sup>Here I'm using  $P(A) = (0, 1)$  instead of  $P(A) = [0, 1]$  to ensure the resulting credal set is regular. This also allows me to express the judgement that the coin toss is independent of  $A$  as the fact that every  $p \in P$  has  $p(E_A|A) = p(E_A)$ . If we had  $P(A) = [0, 1]$  there would be some  $p \in P$  that assigns probability 0 to  $A$ , for which the conditional probability  $p(E_A|A)$  is not defined. The example can be adapted to work for any starting credal set with  $P(A) = [x_1, x_2]$ , where  $x_1, x_2 \in (0, 1)$  and  $x_1 < x_2$ .

each  $p \in P$  would take this as either evidence in favour of, or as evidence against, the fact that the coin landed heads, depending on  $p(A)$ . For each  $p \in P$  we have that:

$$p(H|E_A) = \frac{p(H \cap E_A)}{p(E_A)} \quad (1)$$

$$= \frac{p(H \cap E_A|A)p(A) + p(H \cap E_A|\neg A)p(\neg A)}{p(E_A)} \quad (2)$$

$$= \frac{p(E_A|A)p(A)}{p(E_A)} = p(A) \quad (3)$$

since conditional on  $A$ , the events  $E_A$  and  $H$  are equivalent. Thus after observing  $E_A$ , your updated credal set would be maximally imprecise about  $H$ , in the sense that for every  $r \in (0, 1)$ , there is some  $p \in P(\cdot|E_A)$  with  $p(H) = r$ . Thus we say that your credence in  $H$  *dilates* after observing  $E_A$ .

The key feature of this example is that the two possible outcomes of the coin toss, landing with the "A" side up or with the " $\neg A$ " side up, are symmetrical. Your starting credal set is also maximally imprecise about  $\neg A$ , meaning that for each  $r \in (0, 1)$  there is some  $p \in P$  with  $p(\neg A) = r$ . And if the coin lands with the " $\neg A$ " side up, we would have:

$$p(H|\neg E_A) = \frac{p(\neg E_A|\neg A)p(\neg A)}{p(\neg E_A)} = p(\neg A). \quad (4)$$

So in this case too, your credence in  $H$  will dilate.

We can use this fact to construct a sequential decision problem where the intuition that evidence is valuable seems to fail.

### Example 1 (continued, sequential decision problem)

Let  $a_H$  be the option of making a bet on the next coin toss which gains 100\$ if the painted coin comes up heads, and loses 90\$ otherwise. Let  $a_0$  be the option of making no bet. Consider the following sequential decision problem: you can either learn nothing, and choose from  $\mathcal{A} = \{a_0, a_H\}$  now (at  $t = 0$ ); or you can observe the painted coin toss, see whether it lands with the "A" or " $\neg A$ " face up, and then choose from  $\mathcal{A} = \{a_0, a_H\}$  afterwards (at  $t = 1$ ). The situation is represented in Figure 1.

A popular decision rule for agents with imprecise credences is Maximality.

**Definition 1 (Maximality)** Let  $\mathcal{A}$  be a decision problem and  $P$  a credal set. Then  $a_j$  is admissible for  $P$  from  $\mathcal{A}$  iff there is no  $a_i \in \mathcal{A}$  such that:

$$(\forall p \in P)EU_p(a_i) > EU_p(a_j).$$

It's not obvious how we should apply Maximality to sequential decision problems. Consider the example in Figure 1. The agent knows that at decision node 1, she would choose  $a_H$ , so at node 0 she can identify the option  $\sim Learn$  with  $a_H$ . But at decision nodes 2 and 3, both  $a_H$  and  $a_0$  are

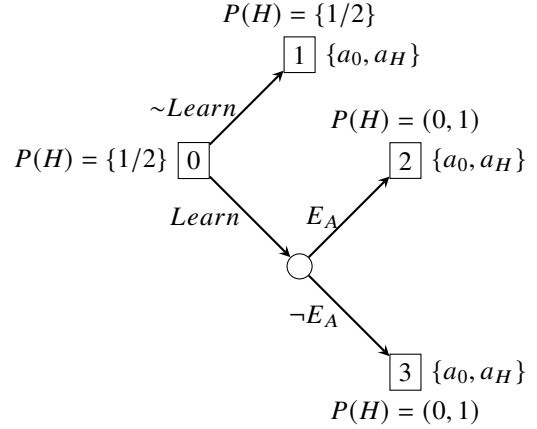


Figure 1: Representation of the Coin Puzzle as a sequential decision problem.

admissible to her. So at node 0, option  $Learn$  is not straightforwardly equivalent to one of the terminal options, and hence it's not obvious how it should be compared against  $\sim Learn$ . One way to tackle this would be to identify  $Learn$  with the option set  $\{a_0, a_H\}$ , and adapt the definition of Maximality to accommodate choices among option sets:

**Definition 2 (Set-Maximality)** Let  $\Theta = \{\mathcal{A}_1, \dots, \mathcal{A}_n\}$  be a set of option sets and  $P$  a credal set. Then  $\mathcal{A}_j$  is admissible for  $P$  from  $\Theta$  iff there is no  $\mathcal{A}_i \in \Theta$  such that the following two conditions hold:

1.  $(\forall x \in \mathcal{A}_j)(\exists y \in \mathcal{A}_i)(\forall p \in P)EU_p(y) > EU_p(x)$ ,
2.  $(\forall y \in \mathcal{A}_i)(\exists x \in \mathcal{A}_j)(\forall p \in P)EU_p(y) > EU_p(x)$ .

Using this rule, when the agent is choosing between  $\sim Learn$  and  $Learn$  at node 0, she is choosing between option sets  $\{a_H\}$  and  $\{a_H, a_0\}$ . Clearly, when  $P(H) = (0, 1)$ , both  $\{a_H, a_0\}$  and  $\{a_H\}$  are admissible option sets, and thus both  $\sim Learn$  and  $Learn$  are admissible options at node 0. The same holds if we attach a small enough price  $\epsilon$  to  $\sim Learn$ . Hence, although paying to avoid free evidence is not mandatory, it is admissible.

There have been different proposals for extending various IP decision rules to sequential decision problems [2, 10]. Yet they all show that it is sometimes admissible for imprecise agents to pay to avoid learning free evidence. That is, we can find a sequential learning problem  $D(\mathcal{E}, \mathcal{A})$  and a credal set  $P$  such that, letting  $Learn$  and  $\sim Learn$  be options defined as above, paying some small price  $\epsilon > 0$  to choose  $\sim Learn$  rather than  $Learn$  is admissible for an agent with imprecise credence  $P$ . In this sense, imprecise probability and decision theory fail to capture our starting intuition that rational agents value the evidence, as expressed by (VE-SC).

Bradley and Steele [2] have sought to soften the blow to IP by showing that, for a sequential version of Maximality, although learning free evidence (as opposed to paying not to learn) is not mandated, it is always admissible. This shows a *weaker* sense in which imprecise agents value the evidence: it is always admissible for them to pursue it when it's free, even though paying to avoid it might be also admissible. Hence, although imprecise agents do not value the evidence in the same way as precise ones, they still value it in this weaker sense.

In this essay I will follow a similar response strategy. I will present a different characterisation of our intuition that rational agent value the evidence, not in terms of sequential decision behaviour, but rather in terms of *deference*. Namely, that an agent values the evidence when she defers to it, treating it like an epistemic expert. In the precise case, both characterisations of the value of evidence are equivalent, in the sense that precise agents value the evidence according to both. But in the imprecise case, the deference characterisation is weaker than the one based on sequential choice. In fact, I will show that agents with imprecise credences always defer to the evidence, even though, as we have seen in this section, they are sometimes allowed to pay to avoid free evidence in sequential learning problems. Thus agents with imprecise credences value the evidence in this weaker sense.

#### 4. The Value of Evidence: Deference Characterisation

Our beliefs are sometimes rationally required to align with those of an expert. If your doctor believes a certain drug will treat your condition, you should also believe this, and if a trusted meteorologist predicts a hurricane is likely to hit your town tomorrow, you should also find this likely. To make this more precise, we need to specify what it means for an agent to match someone's beliefs, and also what kind of beliefs can serve as experts for a given agent. We can do this by specifying a deference principle.

It will help here to introduce some more notation. I will denote by  $\Pi$  the credal set of a deferring agent, writing  $\pi$  as shorthand to denote a singleton credal set  $\Pi = \{\pi\}$ . I will assume throughout that  $\Pi$  is regular. I use  $R$  as a *definite description* of the expert's credal set. This means that  $R$  may denote a different credal set  $R_i$  depending on which  $\omega_i \in \Omega$  is the case (you can think of  $R$  as a function from  $\Omega$  to  $2^{\mathcal{P}\Omega}$ ). For example, let  $\Omega = \{\omega_1, \omega_2\}$ , where  $\omega_1$  is the possibility that the killer was Mr. Green and  $\omega_2$  is the possibility that the killer was Mr. White. Then we could denote by  $R$  the killer's credal set, so that  $R_1$  and  $R_2$  denote the credal sets of Mr. Green and Mr. White, respectively. I will write  $p$  as shorthand for the definite description of a singleton credal set  $R = \{p\}$ . For any random variable

$X : \Omega \rightarrow \mathbb{R}$  and subset  $S \subseteq \mathbb{R}$ , I will write  $[R(X) = S]$  for the event  $\{\omega_i \in \Omega : R_i(X) = S\}$ .

A *deference principle* specifies the relationship that must hold between  $\Pi$  and  $R$  in order for the former to treat the latter as an expert. We can read such a principle in two ways: if we assume a certain  $R$  is worthy of deference, then the principle imposes rationality constraints on the agent's credence  $\Pi$ ; on the other hand, given the agent's credence  $\Pi$ , we can use the principle to determine whether the agent regards a certain  $R$  as an expert. Deference principles have been used to express how rational agents ought to defer to the objective chances [11, 9], to the evidential probabilities (that is, the probabilities that are rational in light of one's total evidence) [3, 7], to their future selves [18, 19], and to other agents [6, 13].

For example, the following is a classical deference principle for precise credences:

- **Reflection Principle (RP):**<sup>5</sup> Let  $\pi$  be an agent's precise credence function, and let  $p$  be the definite description of a precise credence function defined on the same domain. Then  $\pi$  defers to  $p$  iff, for every event  $A \subseteq \Omega$ :

$$\pi(A \mid [p(A) = s]) = s \quad (5)$$

whenever this conditional probability is defined. If this is the case, we say that  $\pi$  Reflects  $p$ .

Reflection says that, conditional on the expert assigning a certain probability to an event, the agent must match them by assigning the same (conditional) probability to that event. This is in line with our intuitions about deference: you defer to someone when, conditional on them having certain credences, you also have those credences.

In the previous section I have looked at a way to capture our intuition that rational agents value the evidence in terms of sequential decision-making (VE-SC). Here I want to propose a different characterisation in terms of deference. The idea is simple: an agent values the evidence when she defers to it. More precisely, since deference is a relationship between credences, we have:<sup>6</sup>

- **Value of Evidence - Deference (VE-D)** A an agent with credal set  $P$  values the evidence when, for any partition  $\mathcal{E} = \{E_1, \dots, E_k\}$ , she defers to the credal set obtained by updating  $P$  on whichever  $E_s \in \mathcal{E}$  is true.

If we take Reflection as our deference principle, and focus on agents with precise credences, then it's easy to show that rational agents value the evidence as specified by (VE-D).

**Proposition 1** *Let  $\pi$  be a regular probability function over  $\Omega$  and  $\mathcal{E} = \{E_1, \dots, E_k\}$  a partition. Let  $p$  be the probability function obtained by updating  $\pi$  on whichever  $E_s \in \mathcal{E}$  is true. Then  $\pi$  Reflects  $p$ .*

<sup>5</sup>This principle was introduced by Van Fraassen [18].

<sup>6</sup>See Dorst [4] for some variants of this view.



Hence agents with precise credences value the evidence both in terms of sequential choice (VE-SC) and terms of deference (VE-D).

The main advantage of (VE-D) is that it does not depend on our theory of sequential action and decision-making. In the imprecise case, this allows us to study the value of evidence without having to worry about extending imprecise decision theory to sequential decision problems. The main obstacle will be to specify what it means for a credal set  $\Pi$  to defer to (the definite description of) a credal set  $R$ . Indeed, some authors have argued that imprecise probability clashes with rational deference principles, as I will discuss in the next section. So my aim will be to specify a notion of deference that does not clash with imprecise probability, and use it to show that imprecise agents value the evidence in the sense of (VE-D).

## 5. Imprecision and Deference

Imprecise analogues of the Reflection Principle have been shown to clash with the phenomenon of credal dilation, which we observed in the Coin Puzzle example. This has led some to argue against the rationality of imprecise probabilities [16, 21].

Here is a natural generalisation of the precise Reflection Principle to the imprecise setting:

- **Value Reflection:** Let  $\Pi$  be an agent's credal set and  $R$  the definite description of a credal set defined on the same domain. Then  $\Pi$  defers to  $R$  iff for every event  $A \subseteq \Omega$  and value set  $S \subseteq \mathbb{R}$ :

$$\Pi(A | [R(A) = S]) = S \quad (6)$$

whenever this conditional credal set is defined.

White [21] has shown that this principle clashes with dilation. To show this, he gives the Coin Toss Puzzle introduced earlier (Example 1). Recall that in this example, you start with  $\Pi(H) = \{1/2\}$  and  $\Pi(A) = (0, 1)$ . Furthermore, you judge the coin toss to be independent of  $A$ . So letting  $E_A$  be the event that the painted coin lands with the “A” face up, for every  $p \in \Pi$  it should be  $p(E_A | A) = p(E_A)$ .

Instead of building a sequential decision problem from this scenario, suppose that Jack is about to toss the coin in front of you. As shown above, regardless of whether you learn  $E_A$  or  $\neg E_A$ , your updated credal set after observing the coin toss will be maximally imprecise about  $H$ . This conflicts with Value Reflection. Before the coin toss, you know your credence in  $H$  will dilate, because it will do so whether you observe  $E_A$  or  $\neg E_A$ . Denoting your future credal set by  $R$ , this means that  $[R(H) = (0, 1)]$  is just  $\Omega$ . So assuming you should defer to your updated credal set  $R$ , Value Reflection requires:

$$\Pi(H) = \Pi(H | [R(H) = (0, 1)]) = (0, 1). \quad (7)$$

This means you should not have  $\Pi(H) = \{1/2\}$  before the coin toss, even though you know that the coin is fair.

To sum things up, the following four conditions are jointly inconsistent:

- (i) Your starting credal set for the Coin Toss Puzzle has  $\Pi(H) = \{1/2\}$ .
- (ii) After observing the coin toss, regardless of how it lands, your updated credal set  $R$  will have  $R(H) = (0, 1)$  (Dilation).
- (iii) You should defer to your updated credal set.
- (iv)  $\Pi$  defers to  $R$  iff for every  $A \subseteq \Omega$  and  $S \subseteq \mathbb{R}$ ,  $\Pi(A | [R(A) = S]) = S$  (Value Reflection).

Most supporters of imprecise probability agree with (i) and (ii), so they must reject either (iii) or (iv).<sup>7</sup> In fact, (iii) and (iv) are deeply connected. As discussed earlier, by specifying a deference principle we specify, for any given credal set, which credal sets it defers to. Therefore, whether it's true that (iii) one should defer to one's informed self (either generally or in this specific example) will depend on the deference principle (iv) we endorse. Note also that claim (iii), that one should defer to one's informed self, is what I presented in the previous section as a way to show that evidence is valuable for imprecise agents. Hence in this essay I will reject (iv) Value Reflection, and my aim will be to specify, and justify as best I can, an IP deference principle that is consistent with (i) - (iii).

## 6. Justifying a Deference Principle

Before presenting a deference principle for agents with imprecise credences, I should say a bit more about how such a principle can be justified. One way to do so is to list desiderata which the principle should satisfy, and then check whether and to what extent a candidate principle satisfies them.

The first desideratum is that the principle should capture some intuitions about what it means to defer to someone.

- **(D1) Captures deference intuitions:** the principle should capture some of our intuitions about what it means to defer to someone. For example, it's natural to think that, if one defers to an expert, then conditional on the expert having a certain opinion one should also have the same opinion.

The second desideratum is related to our discussion in Section 5: the candidate deference principle should establish that imprecise agents defer to their updated selves. After all, if anyone is worthy of being considered an expert, then

<sup>7</sup>Although see [1] for a discussion of update rules which avoid dilation.

surely someone who started from the same prior beliefs as you but has updated these beliefs on more evidence should be considered such.<sup>8</sup>

- **(D2) Defer to informed self:** Let  $\Pi$  be a regular credal set defined over  $\Omega$ ,  $\mathcal{E} = \{E_1, \dots, E_k\}$  be a partition of events, and denote by  $R$  the credal set obtained by updating  $\Pi$  on whichever  $E_s \in \mathcal{E}$  is true. Then  $\Pi$  should defer to  $R$ .

If our deference principle satisfies (D2), then this shows evidence is valuable for imprecise agents in the sense of (VE-D). It also shows that the principle can accommodate credal dilation in cases like the Coin Toss Puzzle, unlike Value Reflection. Since the dilated credal set is obtained by updating the initial credal set, if (D2) is satisfied, then the initial credal set will defer to the updated one.

The third and final desideratum requires that, when only precise credences are involved, our imprecise deference principle should collapse to a reasonable precise deference principle. Often the precise case is easier to study, and if we have good reasons to reject a deference principle in the precise case, these may overpower the reasons we have for supporting it in its more general, imprecise form.

- **(D3) Non-revisionist:** The principle should collapse to a reasonable precise deference principle when both the agent and the expert credal sets are singletons.

## 7. Two IP Deference Principles

I will define imprecise deference principles in terms of gambles which an agent finds desirable. So it will be useful to introduce some notation to talk about these gambles. A gamble is a function  $X : \Omega \rightarrow \mathbb{R}$ , where  $X(\omega_i)$  denotes the value paid by the gamble when  $\omega_i$  is the case. Given an option  $a_j$  and a utility function  $U$ , we can define a corresponding gamble  $X_j = U(a_j, \cdot)$  whose payout at  $\omega_i$  is just the utility of option  $a_j$  if  $\omega_i$  is the case. I denote by  $\mathcal{L}(\Omega)$  the set of all gambles on  $\Omega$ . When  $X \in \mathcal{L}(\Omega)$  is a gamble,  $A \subseteq \Omega$  is an event, and  $p \in \mathcal{P}_\Omega$  is a probability function, I write  $p(X)$  as shorthand for  $\text{Exp}_p(X)$  and  $p(X|A)$  as shorthand for  $\text{Exp}_{p(\cdot|A)}(X)$ .

If  $P$  is an agent's credal set, denote by  $D_P \subseteq \mathcal{L}(\Omega)$  the agent's *set of (strictly) desirable gambles*, defined by:

$$D_P = \{X : \text{for every } p \in P, p(X) > 0\} \quad (8)$$

Intuitively, these are just the gambles that a rational agent with credal set  $P$  would be disposed to accept. Every set

of desirable gambles generated in this way from a regular credal set respects the following *coherence constraints*:

$$0 \notin D \quad (C1)$$

$$X \geq 0, X \neq 0 \implies X \in D \quad (C2)$$

$$X \in D, \lambda > 0 \implies \lambda X \in D \quad (C3)$$

$$X, Y \in D \implies (X + Y) \in D \quad (C4)$$

A set of desirable gambles that respects the above constraints is said to be *coherent*.<sup>9</sup>

We are finally ready to state the two IP deference principles I want to put forward:<sup>10</sup>

- **Strong Total Trust (STT):** Let  $\Pi$  be a regular credal set, and  $R$  be the definite description of a credal set defined on the same domain.  $\Pi$  defers to  $R$  iff for every gamble  $X : \Omega \rightarrow \mathbb{R}$ , we have:

$$X \in D_{\Pi(\cdot|[X \in D_R])} \quad (9)$$

whenever this conditional credal set is defined, and where  $[X \in D_R] = \{\omega_i \in \Omega : X \in D_{R_i}\}$ . If this is the case, we say that  $\Pi$  *S-Trusts*  $R$ .

- **Weak Total Trust (WTT):** Let  $\Pi$  be a regular credal set, and  $R$  be the definite description of a credal set defined on the same domain.  $\Pi$  defers to  $R$  iff for every gamble  $X : \Omega \rightarrow \mathbb{R}$ :

$$-X \notin D_{\Pi(\cdot|[X \in D_R])} \quad (10)$$

whenever this conditional credal set is defined. If this is the case, we say that  $\Pi$  *W-Trusts*  $R$ .

Note that, if  $\Pi$  S-Trusts  $R$ , then  $\Pi$  also W-Trusts  $R$ . This is because, if  $X \in D$ , then  $-X \notin D$ , whenever  $D$  is coherent. It's also worth nothing that, since  $\Pi$  is assumed to be regular on  $\Omega$ , we can rewrite both principles as pointwise constraints on the elements of  $\Pi$ . For example, (STT) is equivalent to the requirement that for every  $X$  such that  $[X \in D_R] \neq \emptyset$ :

$$\pi(X|[X \in D_R]) > 0 \text{ for all } \pi \in \Pi \quad (11)$$

<sup>9</sup>If  $P$  is not regular, then defining  $D_P$  as above does not guarantee that the coherence axiom (C2) is respected. A common way to address this is to define  $D_P = \{X : \text{for every } p \in P, p(X) > 0\} \cup \{X \in \mathcal{L}(\Omega) : X \geq 0, X \neq 0\}$ . However, this definition complicates the relationship between the set  $D_{P(\cdot|A)}$  of desirable gambles according to the conditional credal set  $P(\cdot|A)$ , and the set  $\{X : XA \in D_P\}$  of gambles which are desirable for  $\Pi$  conditional on  $A$ . This relationship is central to the results presented in this essay. I will leave the problem of extending these results to accommodate non-regular credal sets to another time.

<sup>10</sup>The form and name of these principles were inspired by Dorst et al. [5]. Later we will see that, when all credences involved are precise, both STT and WTT are equivalent to the principle given by Dorst et al., called "Total Trust".

<sup>8</sup>At least this is the case when the evidence forms a partition, and one learns whichever element of the partition is true, as in the setup of Good's theorem. Here I will restrict myself to these cases.

If we then write  $[r(X) > 0 \text{ for all } r \in R]$  to denote the event  $\{\omega_i : r(X) > 0 \text{ for all } r \in R_i\}$ , we can further rewrite condition (11) so that the conditioning event is also expressed in pointwise terms:

$$\pi(X|[r(X) > 0 \text{ for all } r \in R]) > 0 \text{ for all } \pi \in \Pi \quad (12)$$

The weaker principle (WTT) can be similarly rewritten as a pointwise constraint.

Let's see how these principles fare against the desiderata (D1-D3) listed above.

### 7.1. D1: Capture Deference Intuitions

Starting from (D1), the intuition behind both STT and WTT is similar to that behind Value Reflection. In both cases, conditional on the expert's imprecise credence having a certain property, the agent's imprecise credence should match it in some way. In the case of Value Reflection the property in question is the set of prevision values assigned to a random variable, whereas in the case of STT it is the disposition to accept a gamble corresponding to that variable. So we can give the following informal definition of STT: an agent defers to an expert when, conditional on the expert finding a gamble desirable, the agent finds that gamble desirable. WTT requires the agent to match the expert in a weaker sense: an agent defers to an expert when, conditional on the expert finding a gamble desirable, the agent does not find it desirable to sell that gamble.

Both STT and WTT have an interesting alternative formulation. Before introducing it, it is useful to define the notion of strict preference between options.

**Definition 2 (Strict preference)** *Let  $P$  be a credal set,  $U$  be a utility function, and  $a_1, a_2$  be two options. We say that  $P$  strictly prefers  $a_1$  to  $a_2$  under  $U$ , when:*

$$(X_1 - X_2) \in D_P \quad (13)$$

where  $X_1 = U(a_1, \cdot)$  is the gamble corresponding to option  $a_1$ , and  $X_2 = U(a_2, \cdot)$  is the gamble corresponding to option  $a_2$ .

Let  $\Pi$  be an agent's credal set, and  $U$  the agent's utility function, and let  $R$  be the definite description of another credal set. Assume as usual that  $P$  is regular. Consider a binary decision problem  $\mathcal{A} = \{a_1, a_2\}$ , and assume that learning  $R$ 's preferences (under  $U$ ) on  $\mathcal{A}$  does not affect the agent's utility function on  $\mathcal{A}$ . Define a new option  $s_1$  which is equal to  $a_2$  whenever  $R$  strictly prefers  $a_2$  to  $a_1$  under  $U$ , and equal to  $a_1$  otherwise. Define  $s_2$  analogously.

$$s_1 = \begin{cases} a_2 & \text{if } R \text{ strictly prefers } a_2 \text{ to } a_1 \text{ under } U \\ a_1 & \text{otherwise.} \end{cases} \quad (14)$$

$$s_2 = \begin{cases} a_1 & \text{if } R \text{ strictly prefers } a_1 \text{ to } a_2 \text{ under } U \\ a_2 & \text{otherwise.} \end{cases} \quad (15)$$

You can think of  $s_1$  as a "black box" option which, if  $R$  has a definite preference in  $\mathcal{A}$ , contains  $R$ 's preferred option, and otherwise contains  $a_1$ . Denote by  $[s_1 \neq a_1]$  the event  $\{\omega_i \in \Omega : R_i \text{ strictly prefers } a_2 \text{ to } a_1\}$ , and similarly for  $[s_2 \neq a_2]$ . We can characterise STT in terms of these black box options as follows:

**Proposition 3**  *$\Pi$  S-Trusts  $R$  iff for every binary decision problem  $\mathcal{A} = \{a_1, a_2\}$ , the following hold:*

1. *If  $[s_1 \neq a_1] \neq \emptyset$ , then  $\Pi$  strictly prefers  $s_1$  to  $a_1$ ,*
2. *if  $[s_2 \neq a_2] \neq \emptyset$ , then  $\Pi$  strictly prefers  $s_2$  to  $a_2$ .*

**Proof** Assuming the utility of the options in  $\mathcal{A}$  is not affected by learning facts about the expert's preferences, we can treat options as gambles. That is, to each option  $a_j$  corresponds a gamble  $X_j$  which pays  $U(a_j, w_i)$  when  $w_i$  is the case. Similarly, write  $S_j$  for the gamble which pays  $U(s_j, w_i)$  when  $w_i$  is the case, where  $s_j$  is defined as above. Then we have:

$$\begin{aligned} S_1 - X_1 &= \\ &= X_1 [(X_2 - X_1) \notin D_R] + X_2 [(X_2 - X_1) \in D_R] - X_1 \\ &= X_2 [(X_2 - X_1) \in D_R] - X_1 [(X_2 - X_1) \in D_R] \\ &= (X_2 - X_1) [(X_2 - X_1) \in D_R] \end{aligned}$$

We say that  $\Pi$  strictly prefers  $s_1$  to  $a_1$  iff  $(S_1 - X_1) \in D_\Pi$ , and by the above this is the case iff:

$$(X_2 - X_1) [(X_2 - X_1) \in D_R] \in D_\Pi \quad (16)$$

which in turn holds iff, for every  $\pi \in \Pi$ :

$$\pi((X_2 - X_1) [(X_2 - X_1) \in D_R]) > 0 \quad (17)$$

If  $[(X_2 - X_1) \in D_R] \neq \emptyset$ , then condition (17) is equivalent to:

$$\pi((X_2 - X_1) | [(X_2 - X_1) \in D_R]) > 0 \quad (18)$$

Condition (18) holds for all  $\pi \in \Pi$  iff, for all  $\pi \in \Pi(\cdot | [(X_2 - X_1) \in D_R])$ :

$$\pi(X_2 - X_1) > 0 \quad (19)$$

Clearly  $\Pi$  S-Trusts  $R$  iff condition (19) holds for every pair of gambles such that  $[(X_2 - X_1) \in D_R] \neq \emptyset$ . And since by definition  $[s_1 \neq a_1] = [(X_2 - X_1) \in D_R]$ , this gives the result. ■

A similar characterisation can be given for W-Trust:

**Proposition 4**  *$\Pi$  W-Trusts  $R$  iff for every binary problem  $\mathcal{A} = \{a_1, a_2\}$ , the following hold:*

1.  $\Pi$  does not strictly prefer  $a_1$  to  $s_1$ ,
2.  $\Pi$  does not strictly prefer  $a_2$  to  $s_2$ .

**Proof** Analogous to Proposition 3. ■

This gives us another intuitive way to think about deference.  $\Pi$  defers to  $R$  when it values  $R$ 's preferences. For STT, this means that a black box containing  $R$ 's preferred option when  $R$  has a definite preference, and containing  $a_j$  otherwise, is at least as good as  $a_j$  according to  $\Pi$ . For WTT, it means that this black box is not definitely worse than  $a_j$  according to  $\Pi$ .

## 7.2. D2: Defer to Informed Self

To capture the intuition that evidence is valuable, we want to show that a rational agent should defer to their updated credences, both in the case of STT and in the case of WTT. This will also ensure that both principles are consistent with credal dilation in cases like the Coin Toss Puzzle (Example 1).

Since STT is the stronger constraint, it suffices to show that agents S-Trust their updated credences.

**Proposition 5** *Let  $\Pi$  be a regular credal set,  $\mathcal{E} = \{E_1, \dots, E_k\}$  be a partition such that  $\Pi(\cdot|E_s)$  is defined for every  $E_s \in \mathcal{E}$ , and denote by  $R$  the credal set obtained by updating  $\Pi$  on whichever  $E_s \in \mathcal{E}$  is true. Then  $\Pi$  S-Trusts  $R$ .*

**Proof** Assume by way of contradiction that  $\Pi$  does not S-Trust  $R$ . Then there is some gamble  $X$  such that  $X \notin D_{\Pi(\cdot|[X \in D_R])}$ , where  $[X \in D_R] \neq \emptyset$ . Under the assumption that  $\Pi$  is regular, this is equivalent to:

$$\pi(X) \leq 0 \text{ for some } \pi \in \Pi(\cdot|[X \in D_R]) \quad (20)$$

$$\Leftrightarrow \pi(X|[X \in D_R]) \leq 0 \text{ for some } \pi \in \Pi \quad (21)$$

$$\Leftrightarrow \pi(X[X \in D_R]) \leq 0 \text{ for some } \pi \in \Pi \quad (22)$$

$$\Leftrightarrow X[X \in D_R] \notin D_{\Pi} \quad (23)$$

We know  $R$  is obtained by updating  $\Pi$  on whichever  $E_s \in \mathcal{E}$  is true, so we can rewrite this as:

$$X \bigcup_{s: X \in D_{\Pi(\cdot|E_s)}} E_s \notin D_{\Pi} \quad (24)$$

And because the members of  $\mathcal{E}$  are mutually exclusive, this is the same as:

$$\sum_{s: X \in D_{\Pi(\cdot|E_s)}} XE_s \notin D_{\Pi} \quad (25)$$

So there must be some  $E_s$  such that  $XE_s \notin D_{\Pi}$ , while also  $X \in D_{\Pi(\cdot|E_s)}$ . But as above,  $X \in D_{\Pi(\cdot|E_s)}$  is equivalent to  $XE_s \in D_{\Pi}$ , contradiction. ■

The fact that imprecise agents defer to their informed selves shows that they value the evidence in the sense of (VE-D). Yet Good's theorem fails for these agents when we extend imprecise decision theory to sequential problems, and so imprecise agents do not value the evidence in the sense of (VE-SC). Thus it's easier to value the evidence in the sense of (VE-D) than (VE-SC) for agents with imprecise credences. In fact, note how the above result does not require us to settle on a specific decision rule for imprecise agents, nor does it require to extend this rule to the sequential case. All that is needed to define our deference principles and to prove Proposition 5 is the notion of desirability of gambles, which is fairly uncontroversial.<sup>11</sup>

The above result also ensures that STT and WTT do not clash with credal dilation in the same way as Value Reflection did. In the Coin Toss Puzzle, as discussed in Section 5,  $R$  is your credal set updated on whichever element of the partition  $\{E_A, \neg E_A\}$  is true. So by Proposition 5 your initial credal set S-Trusts and W-Trusts  $R$ .

## 7.3. D3: Non-Revisionist

If both  $\Pi$  and all candidate experts  $R_i$  are singleton sets containing a single regular probability function (call their elements  $\pi$  and  $p_i$ ), then both STT and WTT are equivalent to the following deference principle, introduced by Dorst et al. [5].

- **Total Trust**  $\pi$  defers to  $p$  iff for every gamble  $X$ :

$$\pi(X|[p(X) \geq 0]) \geq 0 \quad (26)$$

whenever this conditional prevision is defined. If this is the case, we say that  $\pi$  *Totally Trusts*  $p$ .

A thorough defense of this principle can be found in Dorst et al. [5]. Here I will just mention one reason to prefer Total Trust to Precise Reflection, since it will be relevant for the comparison of imprecise deference principles in the next section. The reason is that Precise Reflection is known to be problematic in cases where *modest*, i.e. where some candidate expert credence  $p_i$ , is such that  $p_i([p = p_i]) < 1$ . Here is an example:

**Example 2** *You are in a room with three scientists, who have precise credences  $p_1, p_2$ , and  $p_3$  defined over the same finite possibility space  $\Omega = \{w_1, w_2, w_3\}$ . Their credences*

<sup>11</sup>This notion of desirability arguably does impose some constraints on the decision rule. For example, an imprecise agent with credal set  $\Pi$  who uses the  $\Gamma$ -maximin decision rule may choose the constant gamble 0 over some other gamble  $X$ , even though  $X \in D_{\Pi}$ . Hence supporters of  $\Gamma$ -maximin may find our notion of desirability inadequate. But there are independent reasons to reject  $\Gamma$ -maximin [15], and the two most popular IP decision rules, E-admissibility and Maximality, coincide on binary decision problems, and are consistent with our notion of desirability [17].



are defined as follows:

$$p_i(\{\omega_j\}) = \begin{cases} 0.8 & \text{if } i = j; \\ 0.1 & \text{otherwise.} \end{cases} \quad (27)$$

so that  $p$  is a definite description of the credence of the most accurate scientist in the room.

The key feature of this example is that the candidate experts are modest: if  $\{\omega_i\}$  is the case, the most accurate scientist in the room, who has credence  $p_i$ , is not certain that they are the most accurate scientist in the room, since  $p_i([p = p_i]) = p_i(\{\omega_i\}) = 0.8$ . We also have that  $\{\omega_1\} \equiv [p(\{\omega_1\}) = 0.8]$ , because from the fact that the most accurate scientist assigns probability 0.8 to  $\omega_1$  we can infer that their credence is  $p_1$ . So Precise Reflection imposes the following constraint on your credence  $\pi$ :

$$\begin{aligned} \pi(\{\omega_1\}|\{\omega_1\}) &= \pi(\{\omega_1\} | [p(\{\omega_1\}) = 0.8]) \\ &= 0.8 \quad (\text{by Reflection}) \\ &< 1 \end{aligned}$$

which violates the ratio formula. Hence, if you are coherent, you cannot defer to the most accurate scientist in Example 2 according to Reflection. On the other hand, it can be shown that  $\pi$  Totally Trusts (and therefore also S/W-Trust)  $p$  in Example 2 iff  $\pi$  is a convex combination of  $p_1$ ,  $p_2$ , and  $p_3$ .<sup>12</sup> If we think that it should be possible to defer to modest experts in cases like Example 2, then Total Trust does better than Reflection as a precise deference principle.

## 8. An Alternative IP Deference Principle

This final section compares STT/WTT with an alternative IP deference principle given in the literature, which I will call *Identity Reflection*.<sup>13</sup>

- **Identity Reflection:** Let  $\Pi$  be an agent's credal set and  $R$  the definite description of a credal set defined on the same domain. Then  $\Pi$  defers to  $R$  iff for any credal set  $M$ :

$$\Pi(\cdot | [R = M]) = M \quad (28)$$

whenever this conditional credal set is defined. If this is the case, we say  $\Pi$  *I-Reflects*  $R$ .

Much like Value Reflection, this principle is intuitively appealing. If you defer to  $R$  then, given a full specification  $M$  of a credal set, your credal set conditional on  $R = M$  should be  $M$ . So Identity Reflection satisfies desideratum (D1).

<sup>12</sup>This follows from Theorem B.14 in Dorst et al. [5]

<sup>13</sup>Similar principles are mentioned in [16, 14]. The name and formulation given here is due to Moss [12].

Identity Reflection also addresses the clash with credal dilation discussed in Section 5. Indeed, it's easy to show that the following analogue of Proposition 5 holds for Identity Reflection:

**Proposition 6** *Let  $\Pi$  be a regular credal set,  $\mathcal{E} = \{E_1, \dots, E_k\}$  be a partition such that  $\Pi(\cdot | E_s)$  is defined for every  $E_s \in \mathcal{E}$ , and denote by  $R$  the credal set obtained by updating  $\Pi$  on whichever  $E_s \in \mathcal{E}$  is true. Then  $\Pi$  I-Reflects  $R$ .*

So coherent credal sets always defer to their updated selves according to Identity Reflection, showing that Identity Reflection satisfies desideratum (D2). This implies that imprecise agents value the evidence in the sense of (VE-D), and also that Identity Reflection is consistent with credal dilation in the Coin Toss Puzzle example.

A potential problem for Identity Reflection is that it constrains the agent's opinions conditional on a full specification of the expert's credal set  $R$ . It is natural to wonder whether and how this principle constrains the agent's opinions conditional on hypotheses about more local features of the expert's credal set. Moss has expressed this worry as follows (focusing on the case where  $R$  is your future credal set):

Fans of imprecise credences should value Reflection principles that are easy to operationalize. Identity Reflection constrains your credences in light of your opinions about extremely strong hypotheses. A more valuable Reflection principle would constrain your current credences in light of more targeted opinions about your future credences—for instance, constraining your current imprecise credence in  $p$  in light of your estimates of your future imprecise credence in that same proposition. [12][p. 633]

The next result responds to this worry by showing that Identity Reflection *does* constrain your opinions about a proposition (or more generally, about a random variable) conditional on hypotheses detailing  $R$ 's opinions about that proposition (random variable). In particular, in order to I-Reflect  $R$  you must S-Trust  $R$ .

**Proposition 7** *Let  $\Pi$  be a regular credal set, and  $R$  the definite description of a credal set defined on the same domain. If  $\Pi$  I-Reflects  $R$ , then  $\Pi$  S-Trusts  $R$ .*

**Proof** Assume  $\Pi$  I-Reflects  $R$ . Then let  $X : \Omega \rightarrow \mathbb{R}$  such that  $[X \in D_R] \neq \emptyset$ . Then we have:

$$X[X \in D_R] = X \bigcup_{R_i: X \in D_{R_i}} [R = R_i] \quad (29)$$

$$= \sum_{R_i: X \in D_{R_i}} X[R = R_i]. \quad (30)$$

But for every  $R_i$  such that  $X \in D_{R_i}$ , we know by Identity Reflection that  $\Pi(\cdot | [R = R_i]) = R_i$ , and therefore  $X \in D_{\Pi(\cdot | [R=R_i])}$ . Since  $\Pi$  is regular, this in turn implies  $X[R = R_i] \in D_{\Pi}$ . So each gamble in the sum above belongs to  $D_{\Pi}$ , and by coherence their sum belongs to  $D_{\Pi}$ . Thus we have:

$$X[X \in D_R] \in D_{\Pi} \quad (31)$$

which, since  $[X \in D_R] \neq \emptyset$  and  $\Pi$  is regular, implies  $X \in D_{\Pi(\cdot | [X \in D_R])}$ . ■

Hence, despite its formulation, Identity Reflection does manage to impose substantive constraints on the agent's credences conditional on hypotheses about local features of the expert's credence. It requires that, conditional on the expert finding a gamble desirable, the agent must find that gamble desirable.

Proposition 7 shows that Identity Reflection is no weaker than STT. In fact, we can show Identity Reflection is strictly stronger than STT. That is, we can find  $\Pi$  and  $R$  such that  $\Pi$  S-Trusts  $R$ , but  $\Pi$  does not I-Reflect  $R$ . To see this, note that in the setup of Example 2, Identity Reflection fails in the same way as Precise Reflection does, since conditioning on  $[p_1(\{\omega_1\}) = 0.8]$  is the same as conditioning on  $[p = p_1]$ . Hence, in Example 2, no coherent  $\pi$  I-Reflects  $p$ . Indeed, we can prove the following more general result:

**Proposition 8** *Let  $\Pi$  be a regular credal set and  $R$  the rigid designator of a credal set. If there is some  $\omega_i \in \Omega$  such that  $R_i([R = R_i]) \neq \{1\}$ , then  $\Pi$  does not I-Reflect  $R$ .*

This shows that, according to Identity Reflection, an agent cannot defer to an expert if it is possible that this expert is modest.

But we have seen that many coherent credences Totally Trust  $p$  in Example 2. And because STT/WTT are both equivalent to Total Trust in the precise case, these credences also S/W-Trust  $p$ . So if we think it should be possible to defer to modest experts, this gives us a reason to prefer STT/WTT over Identity Reflection as a deference principle for imprecise credences.

## 9. Conclusion

Rational agents value the evidence. This intuition can be captured in terms of sequential decision-making, as the claim that rational agents never pay to avoid free evidence (VE-SC). Good's theorem shows agents with precise credences value the evidence in this sense. But the theorem fails for agents with imprecise credences, raising a worry that imprecise probabilities are inadequate for a theory of rationality.

This essay looks at a different way to capture our intuitions about the value of evidence, which does not rely on a theory

of sequential decision-making, according to which an agent values the evidence when they defer to their informed selves (VE-D). According to popular notions of deference for precise probabilities, agents with precise credences also value the evidence in this sense. To extend this result to imprecise probabilities, I introduce two imprecise deference principles, STT and WTT. These principles have a natural characterisation in terms of the desirability of black-box options, they are consistent with credal dilation, and collapse to a reasonable deference principle in the precise case. Using these deference principles, I show that coherent imprecise agents defer to their informed selves, and thus value the evidence according to (VE-D). Finally, I show that Identity Reflection, an alternative IP deference principle discussed in the literature, is strictly stronger than STT and WTT. In particular, it's impossible for an agent to defer to a modest expert under Identity Reflection, whereas modest experts can be deferred to under STT/WTT.

An open question is whether STT/WTT can be modified to produce interesting constraints, of the kind expressed by Propositions 3 and 4, for arbitrary decision problems, instead of being limited to cases where the option set is binary. This limitation is related to the fact that whether  $\Pi$  defers to  $R$  under STT/WTT depends only on the corresponding sets of desirable gambles  $D_{\Pi}$  and  $D_R$ . But different credal sets may produce the same set of desirable gambles even when they make importantly different probabilistic judgements. In particular, the E-admissible choices for a credal set in a decision problem are not generally determined by the set of desirable gambles associated to that credal set, which only captures binary preference. A future goal would be to define an IP deference principle that is sensitive to these differences.

## Acknowledgments

Thanks to Jason Konek, Arthur Van Camp, Kevin Blackwell, and the ISIPTA referees for their helpful comments. I was supported by funds from the ERC Starting Grant "Epistemic Utility for Imprecise Probability" during my work on this paper.

## References

- [1] Seamus Bradley and Katie Steele. Uncertainty, learning, and the "problem" of dilation. *Erkenntnis*, 79: 1287–1303, 2014.
- [2] Seamus Bradley and Katie Steele. Can free evidence be bad? Value of information for the imprecise probabilist. *Philosophy of Science*, 83(1):1–28, 2016.
- [3] David Christensen. Rational reflection. *Philosophical Perspectives*, 24:121–140, 2010.

- [4] Kevin Dorst. Evidence: A guide for the uncertain. *Philosophy and Phenomenological Research*, 100(3): 586–632, 2020.
- [5] Kevin Dorst, Benjamin A Levinstein, Bernhard Salow, Brooke E Husic, and Branden Fitelson. Deference done better. *Philosophical Perspectives*, 35(1):99–150, 2021.
- [6] Adam Elga. Reflection and disagreement. *Noûs*, 41(3):478–502, 2007.
- [7] Adam Elga. The puzzle of the unmarked clock and the new rational reflection principle. *Philosophical Studies*, 164:127–139, 2013.
- [8] I. J. Good. On the principle of total evidence. *The British Journal for the Philosophy of Science*, 17(4): 319–321, 1967.
- [9] Ned Hall. Correcting the guide to objective chance. *Mind*, 103(412):505–517, 1994.
- [10] Joseph B Kadane, Mark Schervish, and Teddy Seidenfeld. Is ignorance bliss? *The Journal of Philosophy*, 105(1):5–36, 2008.
- [11] David Lewis. A subjectivist’s guide to objective chance. In Richard C. Jeffrey, editor, *Studies in Inductive Logic and Probability*, volume 2, pages 267–297. University of California Press, 1980.
- [12] Sarah Moss. Global constraints on imprecise credences: Solving reflection violations, belief inertia, and other puzzles. *Philosophy and Phenomenological Research*, 103(3):620–638, 2021.
- [13] Richard Pettigrew and Michael G. Titelbaum. Deference done right. *Philosophers’ Imprint*, 14:1–19, 2014.
- [14] Miriam Schoenfield. Chilling out on epistemic rationality: A defense of imprecise credences (and other imprecise doxastic attitudes). *Philosophical Studies*, 158(2):197–219, 2012.
- [15] Teddy Seidenfeld. A contrast between two decision rules for use with (convex) sets of probabilities:  $\gamma$ -maximin versus e-admissibility. *Synthese*, 140(1/2): 69–88, 2004.
- [16] Brett Topey. Coin flips, credences and the reflection principle. *Analysis*, 72(3):478–488, 2012.
- [17] Matthias CM Troffaes. Decision making under uncertainty using imprecise probabilities. *International journal of approximate reasoning*, 45(1):17–29, 2007.
- [18] Bas C Van Fraassen. Belief and the will. *The Journal of Philosophy*, 81(5):235–256, 1984.
- [19] Bas C Van Fraassen. Belief and the problem of Ulysses and the sirens. *Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition*, 77(1):7–37, 1995.
- [20] Peter Walley. *Statistical reasoning with imprecise probabilities*, volume 42. Springer, 1991.
- [21] Roger White. Evidential symmetry and mushy credence. *Oxford studies in epistemology*, 3:161–186, 2010.