

1 An Improved Argument for  
2 Superconditionalization

3 Julia Staffel

Glauber De Bona

4 (draft of May 2023, please cite published version in *Erkenntnis*)

5 **Abstract**

6 Standard arguments for Bayesian conditionalizing rely on assumptions  
7 that many epistemologists have criticized as being too strong: (i) that con-  
8 ditionalizers must be logically infallible, which rules out the possibility of  
9 rational logical learning, and (ii) that what is learned with certainty must  
10 be true (factivity). In this paper, we give a new factivity-free argument  
11 for the superconditionalization norm in a personal possibility framework  
12 that allows agents to learn empirical and logical falsehoods. We then  
13 discuss how the resulting framework should be interpreted. Does it still  
14 model norms of rationality, or something else, or nothing useful at all?  
15 We discuss five ways of interpreting our results, three that embrace them  
16 and two that reject them. We find one of each kind wanting, and leave  
17 readers to choose among the remaining three.

18 **1 Introduction**

19 Standard arguments for Bayesian conditionalization rely on assumptions that  
20 many epistemologists have criticized as being too strong, in particular: (i) that

21 conditionalizers must be logically infallible, which precludes the possibility of  
22 rational logical learning, and (ii) that what is learned with certainty must be  
23 true (factivity), which disregards the possibility of rationally updating on a  
24 falsehood. For each of these assumptions, it has been shown that we can drop it  
25 by modifying the arguments for conditionalization, resulting in a more general  
26 version of the rule.

27 Pettigrew (2021a) and Rescorla (2020) have shown that we can remove the  
28 factivity assumption, resulting in arguments that show that agents who become  
29 certain of empirical falsehoods should conditionalize on them. These arguments'  
30 novelty is showing that non-conditionalizers are exposed to both accuracy dom-  
31 inance and Dutch book sure loss at all possible worlds, not only at worlds  
32 compatible with the learned evidence.

33 In separate work, Pettigrew (2021b) has demonstrated that by using a frame-  
34 work of personal possibilities instead of logical possibilities, the requirement of  
35 logical infallibility is removed, and agents can be modeled as updating by su-  
36 perconditionalization when they learn logical facts. Even though factivity is not  
37 explicitly assumed, his framework models learning via discarding the set of per-  
38 sonally possible worlds that are incompatible with the evidence. His arguments  
39 for superconditionalization ignore accuracy and betting outcomes in worlds that  
40 are inconsistent with what the agent is certain of. We call this *pseudo-factivity*.  
41 The resulting arguments show that an agent who doesn't (super)conditionalize  
42 is *internally* irrational, since they fully believe they are in a world where they  
43 are exposed both to a Dutch book and to accuracy domination. Nevertheless,  
44 this pseudo-factive argument does not demonstrate guaranteed sure loss and  
45 accuracy domination in all possible worlds: if the agent learns something false,  
46 the actual world is not included in the pseudo-factive argument, and hence it  
47 doesn't show what happens there.

48 Our paper has two aims: The first is to strengthen the argument for (su-  
49 per)conditionalization in a personal possibility setting that allows agents to  
50 learn falsehoods. Unlike Pettigrew’s argument, which assumes pseudo-factivity,  
51 our argument considers accuracy and sure loss at all worlds, including those in-  
52 compatible with what the agent considers certain. This adds robustness to the  
53 agent’s decision to (super)conditionalize – we show that it is the only rational  
54 update even if certainty has been misplaced. Our second aim is to explore how  
55 the resulting framework is best interpreted. Does it still model norms of ratio-  
56 nality, or something else, or nothing useful at all? For example, it demands that  
57 agents “conditionalize” when they become certain of logical falsehoods. This  
58 might seem palatable in cases of highly complex logical reasoning, in which even  
59 a skilled reasoner could easily make an error. Yet, this version of conditional-  
60 ization also requires updating when agents become certain of obvious logical  
61 falsehoods. This raises the question of whether we’ve gone too far - a Bayesian  
62 framework that allows agents to “learn” logical falsehoods and update on them  
63 might at best seem too soft, and at worst seriously misguided.<sup>1</sup> We discuss five  
64 ways of interpreting our results, three that embrace them and two that reject  
65 them. We find one of each kind wanting, and leave readers to choose among the  
66 remaining three.

## 67 **2 Factivity-Free Arguments for Conditioning**

68 Informal glosses of the conditionalization norm tend to go roughly like this: *If*  
69 *an agent becomes certain of some piece of evidence  $E$ , then they should up-*  
70 *date (or plan to update) their credences by making their new unconditional cre-*  
71 *dences equal to their old credences that were conditional on  $E$ .* This formulation  
72 mentions the agent’s attitude towards  $E$ , but omits an important detail that  
73 is usually assumed in arguments for conditionalization: that  $E$  must also be

74 true. Rescorla (2020) has recently drawn attention to this assumption, and  
75 argues that it is desirable to provide proofs of the theorems that underlie the  
76 arguments for conditionalization (in particular, the Dutch book theorems) that  
77 don't rely on it. Pettigrew (2021a) concurs and proves a more general version  
78 of Rescorla's factivity-free Dutch book theorem, as well as a factivity-free ver-  
79 sion of the accuracy-dominance argument for conditionalization (see Briggs and  
80 Pettigrew (2020)).

81 The philosophical motivation for removing the factivity assumption is easy  
82 to see: conditionalization is a norm that tells rational agents how to update  
83 their credences. Rationality is commonly understood as an internalist notion,  
84 hence, agents can have high credences or beliefs in false propositions without  
85 committing a rational error (Comesana 2020). For example, a brain in a vat,  
86 or someone who is deceived by an evil demon, might have evidence that seems  
87 impeccable from their perspective, plausibly making it rational for them to  
88 conditionalize on it. The falsity of their evidence is not attributable to a failure  
89 to be a rational learner, but to having the bad fortune of being placed in an  
90 unreliable learning environment. We can think of more realistic cases as well in  
91 which learners come to acquire false information through no fault of their own.  
92 Rescorla points to instances of scientific reasoning that involve rational updates  
93 on incorrect data, among other examples.<sup>2</sup>

94 Rescorla (2020) gives a non-factive argument for conditionalization that is  
95 based on an improved version of the standard Dutch book theorem for condi-  
96 tionalization. His setting allows an agent to become certain of a proposition  $E$   
97 that is not true, thus abandoning factivity. The impact of this adaptation can  
98 be better understood if we first look at an application of the standard Dutch  
99 book theorem for conditionalization.

100 **Example 1.** *Suppose the agent has the following initially coherent credences at*

101  $t_1$ :  $c(E) = 0.8$ ,  $c(X \& E) = 0.4$ . The agent is considering two updating rules:

102 • U1:  $c_E(X) = 0.4$  (Don't Conditionalize)

103 • U2:  $c_E(X) = 0.5$  (Conditionalize)

104 Since U1 violates conditionalization, there is a factive Dutch Book against  
105 it:

106 • A: at  $t_1$ , the agent buys a bet for 0.40 that returns 1 iff  $X \& E$  is true;

107 • B: at  $t_1$ , the agent buys a bet for 0.08 that returns 0.40 iff  $E$  is false;

108 • C: at  $t_2$ , if the agent becomes certain of  $E$ , they sell a bet for 0.40 that  
109 returns 1 iff  $X$  is true.

110 Suppose that bets A, B and C take place, since the agent becomes certain of  
111  $E$  between  $t_1$  and  $t_2$ . The agent spent 0.48 on bets A and B at  $t_1$ , and received  
112 0.40 back at  $t_2$  by selling C, with a current net loss of 0.08. The agent is certain  
113 that B will not pay back and that A and C cancel each other; for the agent,  $X$   
114 is true iff  $X \& E$  is true. Thus, the agent is certain that they are losing 0.08 for  
115 sure, which is indeed the case if the evidence is true (factivity).

116 But what happens if  $E$  is false, unbeknownst to the agent and the bookie?  
117 The agent has again spent 0.48 on bets A and B. Since  $E$  is false, A returns  
118 nothing, and B returns 0.40, leaving the agent with a net loss of 0.08 from  
119 those two bets. Since the agent and the bookie become certain of  $E$ , despite its  
120 falsity, bet C is also placed, being sold by the agent for 0.40. The outcome then  
121 depends on whether  $X$  is true or false. If  $X$  is true, the agent must pay out 1  
122 on bet C, leading to an overall net loss of  $0.08 + 0.60 = 0.68$ . But if  $X$  ends up  
123 being false, the agent keeps the selling price from bet C, leading to an overall  
124 net gain of  $0.40 - 0.08 = 0.32$ . Hence, failing to conditionalize does not imply  
125 that the agent loses money via a factive Dutch book.<sup>3</sup> Thus, factivity cannot

126 *be simply discarded from an argument for Conditionalization, while keeping the*  
127 *same (factive) Dutch book theorem, or there might be other permissible updates.*

128 A quick and dirty way to patch up the standard argument would be to  
129 replace factivity by what we call *pseudo-factivity*: We narrow the possibilities,  
130 after becoming certain of some (true or false) evidence  $E$ , to the set of possible  
131 worlds consistent with  $E$ , say  $W_E$ . As only worlds  $w \in W_E$  after updating are  
132 considered, the standard Dutch-book and accuracy-dominance theorems of the  
133 classical, factive arguments for conditionalization would be applicable.

134 Consider the situation in Example 1. Assuming pseudo-factivity, after be-  
135 coming certain of  $E$ , the agent and the bookie rule out every world where  $E$   
136 is false. Being aware of the factive Dutch book above, the agent knows U1  
137 (but not U2) makes them vulnerable to sure loss in every world they are still  
138 considering as possible at  $t_2$ . Being certain of  $E$ , the agent simply ignores the  
139 possibility of  $E$  being false, where bets A, B and C could actually give them  
140 profit. Therefore, the agent, from their point of view, is compelled to adopt U2  
141 (conditionalize).

142 The resulting arguments show that a non-conditionalizer should view their  
143 update as irrational, for, from their point of view, they are accuracy dominated  
144 and exposed to Dutch books. Nonetheless, these arguments would not show the  
145 pragmatic or epistemic problems of not conditionalizing from an impartial point  
146 of view, which includes worlds  $w \notin W_E$  where  $E$  is false. This is undesirable,  
147 since it makes the arguments rather weak. Take one of the real-life examples  
148 that motivates Rescorla: some scientists receive data  $E$  that is, unbeknownst  
149 to them, false. How should the scientists update, given that they have become  
150 certain of  $E$ ? It seems plausible that their best option is to conditionalize on  
151  $E$ . Rescorla concurs, arguing that “even if the scientist should not have become  
152 certain of  $E$ , we can still assess how well she reallocates her other credences in

153 light of her faulty certainty.” If we hold the agent’s certainty in  $E$  fixed, it’s not  
154 only true from the agent’s internal perspective that they should conditionalize,  
155 rather, a rational evaluation from a third-personal perspective intuitively agrees.  
156 But this third-personal perspective is left out if we assume pseudo-factivity, and  
157 just consider worlds not ruled out by the agent in the argument.

158       However, as an anonymous reviewer points, out, it’s not obvious that on an  
159 internalist view of rationality, it matters whether there is support for condition-  
160 ing from this impartial perspective in addition to the agent’s own point of view.  
161 We maintain, however, that even from the agent’s perspective, conditionaliza-  
162 tion stands on a stronger footing if it can be supported by an argument that  
163 doesn’t assume pseudo-factivity, but considers all possible worlds. Here’s why:  
164 as is widely acknowledged in discussions of preface-paradoxical cases, rational  
165 agents realize that they sometimes make mistakes, even if they are unable to  
166 spot them. Similarly, an agent who always conditionalizes on the claims they  
167 become certain of realizes that they occasionally update on false things, unbe-  
168 knownst to them. They might thus wonder if always conditioning is the most  
169 desirable strategy for them to pursue in light of this. An argument that as-  
170 sumes pseudo-factivity only tells them that they will *think* (with credence 1)  
171 that conditioning is best in each case. By contrast, an argument that shows that  
172 conditioning is the best strategy in all possible worlds, even in those that the  
173 agent has ruled out, shows them that conditioning is in fact the most desirable  
174 updating strategy for them to implement (assuming what they become certain  
175 of is held fixed). Hence, even from the perspective of the agent, an argument  
176 for conditionalization that doesn’t rely on pseudo-factivity is stronger than one  
177 that does.

178       While Rescorla doesn’t explicitly consider a pseudo-factive modification of  
179 the standard Dutch strategy argument, his own solution cleverly avoids it, and

180 is thus stronger. In his version of the non-factive Dutch strategy argument,  
 181 when  $c_E(X) \neq c(X|E)$ , Rescorla suggests that the bookie make the bet at  $t_2$   
 182 conditional on  $E$ , so that its fair relative price to the agent, who becomes certain  
 183 of  $E^*$ , would be  $c_{E^*}(X|E)$ . The agent sees that bet as fair since either  $E = E^*$ ,  
 184 and  $c_{E^*}(X|E) = c_E(X)$ , or  $c_{E^*}(E) = 0$ , when the agent is certain that the bet  
 185 will be called off. When  $E$  is not the case and the bet at  $t_2$  is called off, the  
 186 situation is analogous to the standard Dutch book for conditionalization when  
 187  $E$  is false, no bet takes place at  $t_2$  and a suitable bet at  $t_1$  on  $E$  guarantees the  
 188 loss to the agent. When  $E$  is the case, the bet at  $t_2$  is not called off and the  
 189 difference between  $c_E(X)$  and  $c(X|E)$  ensures the agent's net loss, no matter  
 190 whether or not  $E = E^*$ . In Example 1, for instance, had the bookie made the  
 191 bet C on  $X$ , at  $t_2$ , conditional on  $E$ , it would have been called off in case  $E$  is  
 192 false, and the agent would still have lost  $0.08 = 0.48 - 0.40$  for sure, as bets A  
 193 and B cost together  $0.48 = 0.40 + 0.08$  and bet B would have paid 0.40 back.  
 194 Rescorla's converse non-factive Dutch book theorem for conditionalization is a  
 195 direct consequence of the standard version, for a sure loss in every outcome in  
 196 Rescorla's scenario implies a sure loss in every scenario where the agent learns  
 197 a true  $E$ .

198 Pettigrew (2021a) formulates different non-factive arguments for condition-  
 199 alization, but he also avoids the problematic assumption of pseudo-factivity. He  
 200 argues for the General Reflection Principle (GRP), a stronger norm than con-  
 201 ditionalization, employing both Dutch book and accuracy considerations. GRP  
 202 is a generalization of Van Frassen's Reflection Principle and in its weaker form  
 203 demands that the current credence function, at time  $t_1$ , be a convex combina-  
 204 tion of the possible future credence functions at time  $t_2$ . Formally, the principle  
 205 reads:

**Weak General Reflection Principle (wGRP)**(Pettigrew 2021a) Sup-



pose  $c$  is the agent's credence function at  $t_1$  and  $c' = \langle c'_1, \dots, c'_n \rangle$  is a tuple of credence functions they might have at  $t_2$ . Then rationality requires that there is, for each  $c'_i$  in  $c'$ , a weight  $\lambda_i$  such that  $\sum_{i=1}^n \lambda_i = 1$  and

$$c(-) = \sum_{i=1}^n \lambda_i c'_i(-)$$

206

207 Pettigrew shows how conditionalization can be directly derived from wGRP  
 208 without assuming factivity. Suppose the agent will become certain of exactly  
 209 one member of a partition  $\{E_1, \dots, E_n\}$  and has a planned (coherent) credence  
 210 function  $c'_i$  to update to when becoming certain of each  $E_i$ , such that  $c_i(E_i) = 1$ .  
 211 Now, if the current credences  $c$  together with  $c' = \langle c'_1, \dots, c'_n \rangle$  satisfy wGRP,  
 212 then  $c(X \& E_j) = c(E_j)c'_j(X)$  for any  $X$ .

213 The first of Pettigrew's arguments for wGRP employs Dutch strategies  
 214 formed by a set of acts, which are a general form of bet. Formally, an act  
 215  $A : W \rightarrow \mathbb{R}$  is a function associating a utility  $A(w)$  with each possible world  
 216  $w \in W$ . From a probabilistic credence function  $c$  whose domain contains a  
 217 proposition  $w$  representing each possible world<sup>4</sup>, one can compute the expected  
 218 utility of an act  $A$ , defined as  $\sum_w c(w)A(w)$ . A credence function  $c$  is said to  
 219 prefer one act out of a pair if it has the higher expected utility. A pair  $\langle c, c' \rangle$ ,  
 220 formed by the prior  $c$  and a tuple  $c'$  of possible posteriors, is said to be vulnera-  
 221 ble to a Strong Dutch strategy if there are acts  $A, B, A', B'$  such that:  $c$  prefers  
 222  $A$  to  $B$ , each  $c'_i$  in  $c'$  prefers  $A'$  to  $B'$  and  $B(w) + B'(w) > A(w) + A'(w)$  for  
 223 every possible word  $w$ . Now a theorem uses Dutch strategies to characterizes  
 224 those  $\langle c, c' \rangle$  satisfying wGRP:

225 **Theorem 1** (Pettigrew (2021a)). *Let  $c$  be a probabilistic credence function and*  
 226  *$c' = \langle c'_1, \dots, c'_n \rangle$  be the possible future probabilistic credence functions defined*

227 over a set of credal objects  $\mathcal{F}$  where each possible world  $w$  is represented.

228 (a) If  $\langle c, c' \rangle$  violates wGRP, then it is vulnerable to a Strong Dutch Strategy.

229 (b) If  $\langle c, c' \rangle$  satisfies wGRP, then it is not vulnerable to a Strong Dutch Strategy.

230 In the accuracy-based argument for wGRP, Pettigrew adopts an additive,  
231 continuous, strictly proper inaccuracy measure  $\mathfrak{I}$  which assigns values to cre-  
232 dences at each possible world. This means there is a continuous strictly proper<sup>5</sup>  
233 scoring rule  $s : [0, 1] \times [0, 1] \rightarrow [0, \infty]$  such that  $\mathfrak{I}(c, w) = \sum_X s(v_w(X), c(X))$ . A  
234 pair  $\langle c, c' \rangle$  is then said to be accuracy dominated if there is an alternative pair,  
235  $\langle c^*, c'^* \rangle$  such that  $\mathfrak{I}(c^*, w) + \mathfrak{I}(c'^*_i, w) < \mathfrak{I}(c, w) + \mathfrak{I}(c'_i, w)$ , for all possible worlds  
236  $w$  and any  $1 \leq i \leq n$ . A theorem then states that any pair  $\langle c, c' \rangle$ , formed by the  
237 current credence function  $c$  and the set of possible posteriors  $c' = \langle c'_1, \dots, c'_n \rangle$ ,  
238 is accuracy dominated if wGRP is violated, and satisfying wGRP avoids such  
239 domination:

240 **Theorem 2** (Pettigrew (2021a)). *Let  $c$  be a probabilistic credence function and*  
241  *$c' = \langle c'_1, \dots, c'_n \rangle$  be the possible future probabilistic credence functions defined*  
242 *over a set of credal objects  $\mathcal{F}$  where each possible world  $w$  is represented.*

243 (a) *If  $\langle c, c' \rangle$  violates wGRP, then it is accuracy dominated.*

244 (b) *If  $\langle c, c' \rangle$  satisfies wGRP, then it is not accuracy dominated.*

245 Pettigrew thus shows us two more routes towards arguing for non-factive  
246 conditionalization, both of which show that conditioning is the only rational  
247 update rule in all possible worlds, regardless of whether the agent learns a truth  
248 or a falsehood.

249 In the next section, we will turn to another one of Pettigrew's arguments for  
250 conditioning, which he has offered within a personal possibility framework. We  
251 will argue that it is inferior to the arguments just discussed, because it relies on  
252 pseudo-factivity.

### 253 **3 Arguments for Conditioning in a Personal Pos-** 254 **sibility Framework**

255 Dropping factivity is not the only modification to arguments for conditionaliza-  
256 tion that people have made in order to better model realistic learning scenarios.  
257 In another, unrelated strand of the literature, it has been debated how logical  
258 learning can be modeled in a Bayesian framework. Standard Bayesian models  
259 that are based on classical probabilities assume that rational agents are logi-  
260 cally infallible. While this doesn't mean that an agent needs to know every  
261 possible logical truth, it still requires that, insofar they have any attitude at all  
262 towards a proposition, they assign credence 1 to it if it is a logical truth, and  
263 credence 0 if it is a logical falsehood. Being uncertain about, i.e., assigning mid-  
264 dling credences to, logical truths and falsehoods is not permitted by standard  
265 Bayesian models. Also, a rational agent's credences have to correctly reflect  
266 other logical relations between the contents of their attitudes, for example, if  
267 they have credences towards two propositions  $X$  and  $Y$ , and the former entails  
268 the latter, then their credence in  $X$  can't be higher than their credence in  $Y$ .  
269 This precludes Bayesian models from representing learning experiences in which  
270 agents come to be aware of logical facts and relations that they were previously  
271 ignorant of. Yet, this kind of logical (and also mathematical) learning is com-  
272 mon for human reasoners. Being uncertain about a logical or mathematical fact  
273 might be a failure of ideal rationality, but is not necessarily a rational defect  
274 given standards of human rationality.

275 A common suggestion for incorporating logical learning into a Bayesian  
276 framework is to replace logical possibilities with what is possible from the  
277 agent's perspective in formulating norms of probabilistic coherence and updat-  
278 ing. First proposed by Hacking (1967), this idea has recently been developed

279 further by Pettigrew (2021b). Pettigrew proposes to model an agent’s growing  
280 logical awareness by replacing logical with personal possibilities as the contents  
281 of the agent’s attitudes. This replacement does not preclude us from formulat-  
282 ing Dutch book or accuracy arguments for coherence or conditionalization. The  
283 main difference is that the agent is now required to be coherent with regard to  
284 what is possible from their perspective, rather than what is logically possible.  
285 Further, Pettigrew argues that we should argue for a slightly modified version  
286 of conditionalization called “superconditionalization.” Superconditionalization  
287 is slightly more general than conditionalization in the following sense: standard  
288 conditionalization assumes that there is always a proposition that the agent  
289 learns with certainty and assigns credence 1 to. Superconditionalization does  
290 not require this. Instead, an agent can directly rule out possibilities, without  
291 there being a proposition that corresponds to those possibilities, and to which  
292 the agent had assigned a credence. Pettigrew’s argument for superconditioning  
293 in the personal possibility framework does not assume factivity, hence, agents  
294 can learn things that are false. Unfortunately, however, it assumes pseudo-  
295 factivity, which means that it doesn’t show that conditionalizing is the only  
296 rational updating rule in all the worlds regarded as possible before the learn-  
297 ing experience occurs. The argument only takes into account the worlds the  
298 agent considers live after  $E_i$  has been learned. We will explain how the argu-  
299 ment works, and then motivate the need to reformulate the argument without  
300 pseudo-factivity.

301 Formally, Pettigrew’s framework, which we mainly follow from here on,  
302 employs a set  $W$  of personally possible worlds at which each credal object  
303 from a set  $\mathcal{F}$  is either true or false. Each  $w \in W$  corresponds to a valua-  
304 tion  $v_w : \mathcal{F} \rightarrow \{0, 1\}$ , with  $v_w(X) = 0$  if  $X$  is false at  $w$  and  $v_w(X) = 1$  if  
305  $X$  is true at  $w$ , for any  $X \in \mathcal{F}$ . The set of these valuations is denoted by

306  $W_{\mathcal{F}} = \{v_w | w \in W\}$ . Note that a contradiction can be true in a given person-  
307 ally possible world, or a tautology false. Also, there might be worlds where  $X$   
308 and  $\neg X$  are both true, or both false, for some  $X \in \mathcal{F}$ . A credence function  
309  $c : \mathcal{F} \rightarrow [0, 1]$  represents the agent's numerical credences on the credal objects.  
310 In the learning scenario, the agent is about to become certain of, between  $t_1$   
311 and  $t_2$ , exactly one  $E_i$  from a partition  $\mathcal{E} = \{E_1, \dots, E_n\}$  of  $W$ .<sup>6</sup>

312 Note again that each  $E_i$  need not correspond to a credal object  $X \in \mathcal{F}$  that  
313 is true only in worlds  $w \in E_i$ ; that is,  $E_i$  need not be represented in  $\mathcal{F}$ . An  
314 updating rule  $c'$  is a function that takes each  $E_i \in \mathcal{E}$  and returns a credence  
315 function  $c'_i$ , the posterior at  $t_2$  endorsed by the rule when the agent becomes  
316 certain of  $E_i$ . Given a fixed partition  $\mathcal{E} = \{E_1, \dots, E_n\}$ , we can denote an  
317 updating rule by the tuple  $c' = \langle c'_1, \dots, c'_n \rangle$ , hence it can be seen as a set of  
318 possible future credences. We call a pair  $\langle c, c' \rangle$ , formed by a credence function  
319 at  $t_1$  and an updating rule, a credal strategy.

320 Using personally possible worlds  $W$ , the (synchronic) incoherence of an  
321 agent's credence function  $c$  is defined as the existence of a set of bets it en-  
322 dors that, taken together, causes loss to the agent at every world in  $W$ . A  
323 theorem by de Finetti (1974) characterizes the coherent  $c$  as those inside the  
324 convex hull of  $W_{\mathcal{F}}$ , denoted by  $W_{\mathcal{F}}^+$ . This motivates the following version of the  
325 Probabilism norm, parametrized by  $W$ :

326 **Personal Probabilism**(Pettigrew 2021b) Suppose  $c$  is the agent's credence  
327 function and  $W$  is the set of their personally possible worlds. Then  $c$  ought to  
328 be in  $W_{\mathcal{F}}^+$ .

329 A personally probabilistic  $c$  must be some weighted average of the valuations  
330  $v_w : \mathcal{F} \rightarrow [0, 1]$ . That is, there must be weights  $p : W \rightarrow [0, 1]$  such that  
331  $c(X) = \sum_w p(w)v_w(X)$  for all  $X \in \mathcal{F}$ . The function  $p : W \rightarrow [0, 1]$  can be seen  
332 as a way to coherently extend  $c$  to (credal objects representing) each personally

333 possible world, meaning that  $c$  together with  $p$  remains personally probabilistic  
 334 and immune to Dutch books.

335 Assuming again a set  $W$  of personally possible worlds, the diachronic inco-  
 336 herence of a credal strategy  $\langle c, c' \rangle$  is analogously defined: there is a set of bets  
 337  $B$  endorsed by  $c$ , a set  $B_i$ , for each  $i$ , endorsed by  $c'_i$ , and, for any  $E_i \in \mathcal{E}$  and  
 338 personally possible world  $w \in E_i$ ,  $B$  together with  $B_i$  lead to a loss of money.  
 339 To characterize the coherent credal strategies, over  $W$ , Pettigrews generalizes  
 340 conditionalization to consider cases where  $E_i$  is not represented in  $\mathcal{F}$ :

**Definition 1.**  $\langle c, c' \rangle$  is *superconditionalizing* if there is a function  $p : W \rightarrow [0, 1]$ ,  
 with  $\sum_{w \in W} p(w) = 1$ , such that, for all  $X \in \mathcal{F}$ ,  $c(X) = \sum_{w \in W} p(w)v_w(X)$  and for  
 each  $E_i \in \mathcal{E}$  with  $\sum_{w \in E_i} p(w) > 0$ :

$$c'_i(X) = \frac{\sum_{w \in E_i} p(w)v_w(X)}{\sum_{w \in E_i} p(w)}$$

341 Pettigrew proceeds to prove that a credal strategy  $\langle c, c' \rangle$  is diachronically  
 342 coherent if, and only if,  $\langle c, c' \rangle$  is superconditionalizing. This yields his first  
 343 argument for the following norm:

344 **Superconditionalization**(Pettigrew 2021b) Suppose  $c$  is the agent's cre-  
 345 dence function and  $c'$  is their updating rule. Then  $\langle c, c' \rangle$  is superconditionalizing.

346  
 347 The second argument for superconditionalization is based on accuracy dom-  
 348 inance. Assuming a continuous, additive, strictly proper inaccuracy measure  $\mathfrak{I}$ ,  
 349 Pettigrew defines that  $\langle c^*, c'^* \rangle$  accuracy dominates  $\langle c, c' \rangle$  if, for any  $E_i \in \mathcal{E}$  and  
 350 any  $w \in E_i$ ,  $\mathfrak{I}(c^*, w) + \mathfrak{I}(c'^*_i, w) < \mathfrak{I}(c, w) + \mathfrak{I}(c'_i, w)$ . Pettigrew then proves  
 351 that a credal strategy  $\langle c, c' \rangle$  satisfies superconditionalization if, and only if, it is  
 352 not accuracy dominated.

353 The similarities between the accuracy dominance argument employed in The-  
 354 orem 2 and the one defined above hide a crucial difference: the worlds consid-  
 355 ered. In the argument for superconditionalization, Pettigrew assumes pseudo-  
 356 factivity, evaluating sure losses and accuracy after updating only for worlds  
 357 consistent with the evidence, as the agent discards the other possibilities while  
 358 learning. This has a similar effect as assuming factivity, since the set of pos-  
 359 sible worlds is the learned  $E$ . His proof shows that a superconditioning credal  
 360 strategy is not accuracy dominated, and this holds even if we drop factivity or  
 361 consider the initial set of worlds  $W$ . Nevertheless, if we consider all initially  
 362 possible worlds  $w \in W$ , his proof does not ensure that *only* superconditioning  
 363 credal strategies will not be accuracy dominated.

364 Formally, the accuracy dominance mentioned in Theorem 2 holds for every  
 365 pair  $\langle c'_i, w \rangle$ . In Pettigrew's accuracy-based argument for superconditionaliza-  
 366 tion, the dominance is defined for every  $\langle c'_i, w \rangle$  such that  $w \in E_i$ , thus ignoring,  
 367 for each  $E_i$ , all worlds  $w \notin E_i$ . That is, the credal strategies  $\langle c, c' \rangle$  and  $\langle c^*, c'^* \rangle$   
 368 are compared at a world  $w \in E_i$  only via  $c'_i$  and  $c_i^*$ . And in fact this detail is  
 369 used in Pettigrew's proof. That is, for a non-superconditioning  $\langle c, c' \rangle$ , Pettigrew  
 370 does not show a pair  $\langle c^*, c'^* \rangle$  with  $\mathfrak{J}(c^*, w) + \mathfrak{J}(c'^*_i, w) < \mathfrak{J}(c, w) + \mathfrak{J}(c'_i, w)$  for  
 371 all  $i$  and  $w$ , including those  $w \notin E_i$ .

372 Something similar occurs in his Dutch book argument for superconditional-  
 373 ization. In the definition of diachronic incoherence, after the agent learns  $E_i$ ,  
 374 updates, and the bets  $B_i$  endorsed by  $c'_i$  take place, only net gains at worlds  
 375  $w \in E_i$  are considered, due to pseudo-factivity. However, if the agent becomes  
 376 certain of a false  $E_j$ , updating to  $c'_j$ , they will not necessarily engage in the bets  
 377  $B_i$  that would cause them sure loss at worlds  $w \in E_i$ . Again, if pseudo-factivity  
 378 were dropped and we considered all possible combinations of worlds  $w \in W$   
 379 and pieces of evidence  $E_i \in \mathcal{E}$ , superconditionalization would still avoid Dutch

380 books, but there is no proof that only superconditioning credal strategies would  
381 do so. Indeed, a Dutch book relying on pseudo-factivity could give profit to  
382 a non-superconditionalizer at a world ignored for being incompatible with the  
383 learned evidence, as Example 1 shows.

384 Pettigrew's arguments are pseudo-factive: even though they do not assume  
385 the evidence  $E_i$  is true, the live possibilities while updating are reduced to the  
386 worlds consistent with  $E_i$  – as factivity would imply. Above, we argued that  
387 pseudo-factive arguments should be avoided when supporting conditionalization  
388 in cases of learning false evidence. While the pseudo-factive arguments prove  
389 that non-conditionalizers are irrational in the worlds compatible with  $E_i$ , they  
390 don't show that they are irrational in the worlds that the agent has ruled out  
391 after learning  $E_i$ . But in cases of empirical learning in a standard Bayesian  
392 framework, we wanted an argument that shows that when false evidence is  
393 learned, conditionalization is the best updating strategy not just from the per-  
394 spective of the agent (and their misplaced certainty), but also from an impartial  
395 perspective that has all possible worlds in view. Both Rescorla and Pettigrew  
396 deliver such arguments in that context.

397 This raises the question of whether there is something special about the  
398 framework of personal possibilities that makes pseudo-factive arguments for  
399 conditionalization more appropriate. One might point out, for example, that  
400 we're only trying to model the agent's perspective, making it superfluous to  
401 attend to possibilities the agent has ruled out. We don't find this reasoning very  
402 persuasive. Even in a personal possibility framework, an agent's certainty can be  
403 misplaced. This is the case for both empirical and logical certainties. Just like in  
404 the cases discussed before, we don't just want to know whether conditioning is  
405 the only rational updating strategy from within the agent's current perspective,  
406 we also want to know if, holding the agent's certainties fixed, conditioning is



407 the only way to go from a third-personal perspective that need not share the  
408 agent’s assessment of which worlds are live possibilities. As explained above,  
409 this also reassures the agent that they should always conditionalize, even if they  
410 realize they are sometimes wrong via preface-paradox-style reasoning. If our  
411 aim is to show that conditionalization is robustly applicable even in non-ideal  
412 conditions, then it is desirable to show that it is the uniquely rational updating  
413 strategy not only in the absence of logical omniscience, but also in the presence  
414 of misplaced certainty. Avoiding pseudo-factivity is thus especially desirable in  
415 a personal possibility framework.

## 416 **4 A Non-Factive Argument for Conditioning in** 417 **a Personal Possibility Framework**

418 In this section, we will show how both modifications to the standard arguments  
419 for conditioning can be combined - we can have a personal-possibility argu-  
420 ment for superconditioning that is properly non-factive, i.e., it avoids assuming  
421 pseudo-factivity. We will show that superconditionalization can be derived from  
422 the Weak General Reflection Principle extended to personally possible worlds.  
423 The arguments for superconditionalization thus depend on those for wGRP,  
424 which we first need to adapt to the personally possible worlds framework.

425 First, we must reinterpret and refine the Weak General Reflection Principle  
426 in light of our framework. It suffices to assume a fixed set of credal objects  $\mathcal{F}$   
427 over which the credences are assigned. Note that wGRP does not mention a  
428 set of worlds, but only credence functions, which in principle may even violate  
429 probabilism. We can explicitly add a set of personally possible worlds  $W$  to the  
430 definition though, to refer to in the following arguments:

**Weak General Reflection Principle (wGRP)** Consider a set of person-

ally possible worlds  $W$  and a set of credal objects  $\mathcal{F}$ , each of which is either true or false at a given world  $w \in W$ . Suppose  $c : \mathcal{F} \rightarrow [0, 1]$  is the agent's credence function at  $t_1$  and  $c' = \langle c'_1, \dots, c'_n \rangle$  is a tuple of credence functions  $c'_i : \mathcal{F} \rightarrow [0, 1]$  they might have at  $t_2$ . Then rationality requires that there is, for each  $c'_i$  in  $c'$ , a weight  $\lambda_i$  such that  $\sum_{i=1}^n \lambda_i = 1$  and, for all  $X \in \mathcal{F}$

$$c(X) = \sum_{i=1}^n \lambda_i c'_i(X)$$

431

432 Both arguments for wGRP by Pettigrew (2021a), employing Dutch strategy  
 433 or accuracy considerations, assume probabilistic credences that are also deter-  
 434 mined for propositions representing each possible world. This brings about a  
 435 problem for our framework as the set  $\mathcal{F}$  of credal objects does not necessarily  
 436 contain those propositions (to allow for logical learning). In the Dutch strategy  
 437 argument, such assumptions are employed to determine the preference of a cre-  
 438 dence function  $c$  over a pair of acts. To address this issue, we can redefine this  
 439 preference using the credence functions  $p : W \rightarrow [0, 1]$  that coherently extend  $c$   
 440 to all the personally possible worlds:

441 **Definition 2.** Given two acts  $A : W \rightarrow \mathbb{R}$  and  $B : W \rightarrow \mathbb{R}$ , a personally proba-  
 442 bilistic credence function  $c : \mathcal{F} \rightarrow [0, 1]$  *prefers*  $A$  to  $B$  if, for every credence func-  
 443 tion  $p : W \rightarrow [0, 1]$  that coherently extends  $c$ ,  $\sum_{w \in W} p(w)A(w) > \sum_{w \in W} p(w)B(w)$ .

444 The Strong Dutch Strategy definition can then be applied to credence func-  
 445 tions defined over an arbitrary set  $\mathcal{F}$  of credal objects. The theorem that charac-  
 446 terizes the pairs  $\langle c, c' \rangle$  vulnerable to a Strong Dutch Strategy as those violating  
 447 wGRP can now be reworked to consider an arbitrary  $\mathcal{F}$ .<sup>7</sup>

448 **Theorem 3.** *Let  $c : \mathcal{F} \rightarrow [0, 1]$  be a personally probabilistic credence function*  
 449 *and  $c' = \langle c'_1, \dots, c'_n \rangle$  be the set of possible future personally probabilistic credence*

450 functions defined over  $\mathcal{F}$ .

451 (a) If  $\langle c, c' \rangle$  violates wGRP, then it is vulnerable to a Strong Dutch Strategy.

452 (b) If  $\langle c, c' \rangle$  satisfies wGRP, then it is not vulnerable to a Strong Dutch Strategy.

453 The accuracy-based argument for wGRP put forward by Pettigrew relies  
454 on Theorem 2, which also requires some adaptation to the personally possible  
455 worlds framework, as the arbitrary set  $\mathcal{F}$  of credal objects need not contain a  
456 proposition for each possible world.

457 **Theorem 4.** Let  $c : \mathcal{F} \rightarrow [0, 1]$  be a personally probabilistic credence function  
458 and  $c' = \langle c'_1, \dots, c'_n \rangle$  be the set of possible future personally probabilistic credence  
459 functions defined over  $\mathcal{F}$ .

460 (a) If  $\langle c, c' \rangle$  violates wGRP, then it is accuracy dominated.

461 (b) If  $\langle c, c' \rangle$  satisfies wGRP, then it is not accuracy dominated.

462 Now that we have two non-factive (and non-pseudo-factive) arguments for  
463 wGRP, considering personally possible worlds and an arbitrary set  $\mathcal{F}$  of credal  
464 objects, we need to derive superconditionalization from it. The idea is that an  
465 updating rule  $c'$ , with a planned  $c'_j$  for the case of becoming certain of each  $E_j$   
466 from a given partition  $\{E_1, \dots, E_n\}$  of  $W$ , is a set of possible future (person-  
467 ally probabilistic) credence functions, which should satisfy wGPR in order to  
468 avoid Dutch strategies and accuracy domination. When we simply drop factiv-  
469 ity, without replacing it by something else, any credal strategy  $\langle c, c' \rangle$  satisfying  
470 wGRP would not be vulnerable to Dutch books or accuracy domination. For  
471 instance, the agent could plan to hold their credences fixed regardless of which  
472 evidence they become certain of, and they would still seem rational if factivity  
473 is not replaced by a suitable property<sup>8</sup>. But, of course, becoming certain of  
474  $E_i$  implies some restrictions on the updated credence function, and, actually,

475 assuming  $E_1, \dots, E_n$  are among the considered credal objects, then imposing  
 476  $c'_i(E_j) = 1$  whenever  $i = j$ <sup>9</sup>, given personal probabilism, suffices for wGRP to  
 477 imply Conditionalization without factivity, as Pettigrew (2021a) shows<sup>10</sup>. How-  
 478 ever, as  $E_j$  might not be in  $\mathcal{F}$  in our framework, this assumption has to be  
 479 slightly modified: each  $c'_j$  must be coherently (according to personal probabil-  
 480 ism) extendable to a credence function  $c^*$  with  $c^*(E_j) = 1$ . If each  $c'_j$  satisfies  
 481 that property, captured by the following definition, a credal strategy  $\langle c, c' \rangle$  sat-  
 482 ifying wGRP will be superconditioning.

483 **Definition 3.** A set of credence functions  $c' = \langle c'_1, \dots, c'_n \rangle$  respects a partition  
 484  $\{E_1, \dots, E_n\}$  of  $W$  if, for each  $1 \leq i \leq n$ , there is a function  $p_i : W \rightarrow [0, 1]$ ,  
 485 with  $\sum_{w \in W} p_i(w) = 1$ , such that  $c'_i(X) = \sum_{w \in W} p_i(w)v_w(X)$  for each  $X \in \mathcal{F}$  and  
 486  $\sum_{w \in E_i} p_i(w) = 1$ .

487 Note that respecting a partition implies that all credence functions in the  
 488 set  $c'$  are personally probabilistic. When each  $E_i$  is in  $\mathcal{F}$ , a set of credence  
 489 functions  $c' = \langle c'_1, \dots, c'_n \rangle$  respects a partition  $\{E_1, \dots, E_n\}$  if, for all  $i$  and  $j$ ,  
 490  $c'_i(E_j) = 1$  whenever  $i = j$ ; and personal probabilism then implies  $c'_i(E_j) = 0$  for  
 491  $j \neq i$ . When some  $E_i$  is not in  $\mathcal{F}$ , respecting the partition means the agent can  
 492 extend the credence functions' range to  $\mathcal{F} \cup \{E_i\}$  and assign  $c'_j(E_i) = 1$  for  $j = i$   
 493 without violating personal probabilism. If  $c'$  is an update rule for the partition  
 494 it respects, the agent plans to adopt a credence function when becoming certain  
 495 of an  $E_i$  that is coherent with assigning credence 1 to  $E_i$  and credence 0 to the  
 496 other  $E_j \neq E_i$ .

497 The next result derives superconditionalization for a credal strategy satisfy-  
 498 ing wGRP whose updating rule respects the partition for which it is defined:

499 **Theorem 5.** *Let  $c$  be a credence function. Let  $c' = \langle c'_1, \dots, c'_n \rangle$  be an updating*  
 500 *rule for a partition  $\{E_1, \dots, E_n\}$  of  $W$ , respecting it. If the pair  $\langle c, c' \rangle$  satisfies*  
 501 *the Weak General Reflection Principle, then  $\langle c, c' \rangle$  is superconditioning.*

502 When factivity does not hold, and the agent might become certain of some  
503 false  $E_j$ , adopting the credence function  $c'_j$ , according to their updating rule,  
504 Theorem 5 shows that not superconditioning on  $E_j$  implies Dutch book vulnera-  
505 bility and accuracy dominance. In fact, in the theorem we could have defined  $c'$   
506 simply as a set of possible future credence functions instead of an updating rule  
507 for  $\{E_1, \dots, E_n\}$ . In that case, the agent would not need to commit to adopt  
508 specifically  $c'_j$  when becoming certain of  $E_j$ . As long as those future credence  
509 functions respect a partition, wGRP requires them to be superconditioning on  
510 that partition.

511 Putting it all together, we have provided two arguments for a stronger version  
512 of wGRP, which holds for an arbitrary set of credal objects. Furthermore, we  
513 proved that if an agent has an updating rule  $c' = \langle c'_1, \dots, c'_n \rangle$  for a partition  
514  $\{E_1, \dots, E_n\}$ , and each  $c'_i$  is personally probabilistic implying  $c'_i(E_i) = 1^{11}$ , then  
515 wGRP entails superconditionalization.

## 516 5 Consequences and Responses

517 In the previous section, we generated an argument for superconditionalization  
518 that drops both the factivity assumption (without assuming pseudo-factivity)  
519 and swaps logical for personal possibilities. Our argument treats the cases of  
520 learning a truth and becoming certain of a falsehood in a parallel way. In  
521 both cases, the uniquely rational response to becoming certain of some  $E$  is to  
522 superconditionalize on it, regardless of whether the learned claim is empirical or  
523 logical, and regardless of whether we focus on all the worlds, including ones the  
524 agent no longer considers live, or on just the ones not ruled out by the agent.

525 In what follows, we will discuss how the resulting framework is best in-  
526 terpreted. While Rescorla argues in some detail for dropping factivity, and  
527 Pettigrew motivates the need to represent logical learning with personal possi-

528 bilities, there has so far been no discussion of a model that combines both, even  
529 though this possibility is already implicit in Pettigrew's theory, as we explained.  
530 Our discussion is independent of our previous argument, in the sense that noth-  
531 ing we say in this section depends on accepting that Pettigrew's pseudo-factive  
532 arguments for superconditionalization should be replaced by ours.

533 While the two ways of modifying standard Bayesianism might seem individ-  
534 ually compelling, one might worry that once we combine them, the resulting  
535 version of superconditionalization goes too far. For example, suppose an agent,  
536 call him Bob, is deliberating about installing a tree swing for his children. He  
537 is currently not sure if this can be done safely, so he needs to calculate whether  
538 the tree is strong enough to withstand the force generated by the swing. Sup-  
539 pose he does the calculation, which is well within his mathematical capabilities,  
540 but he makes an error. His result suggests the swing is safe, even though it's  
541 not. Still, our argument for superconditionalization recommends conditioning  
542 on the faulty result, which would then lead the agent to further conclude that  
543 building the swing is safe. Even by the standards of non-ideal norms of human  
544 rationality, our argument's verdict might seem overly permissive.

545 We will now discuss different possible responses to our results. To keep the  
546 discussion manageable, we will assume that readers are generally sympathetic  
547 to Bayesian theories of norms of rationality, and the standard arguments for  
548 supporting them, such as accuracy and Dutch book arguments. The question  
549 we're interested in is whether in fully dropping factivity and embracing personal  
550 possibilities, we've relaxed the standard framework too far. We will first discuss  
551 what can be said in favor of the results we've generated, and after that, discuss  
552 ways of pushing back on them. There are different ways in which one might  
553 embrace the results, which we will call (i) embrace completely, (ii) embrace and  
554 supplement, and (iii) embrace and reinterpret.

555 (i) Embrace Completely

556 One possible reaction is to think that we're getting things exactly right. On  
557 this view, the framework operates correctly in constraining the agent's credences  
558 only in light of the possibilities that are distinguished by the agent, regardless  
559 of how they map onto the logical possibilities. Further, the only relevant con-  
560 sideration in licensing an update is whether the agent has become certain of  
561 any of the possibilities (or ruled out any of them ). What is actual and whether  
562 the update is based on logical or empirical information is irrelevant. Hence,  
563 this view essentially formalizes the idea that rational norms of coherence and  
564 reasoning should be entirely dependent on the agent's perspective, regardless  
565 of how empirically and logically (in)accurate their take on the world is. If we  
566 embrace this interpretation, the example is not taken to be worrisome: Bob is  
567 correct in thinking that he should decide whether to build the tree swing based  
568 on a calculation of the strength of the tree. And if his calculation shows him  
569 that the tree is strong enough, then, *from his perspective*, he should update his  
570 credences and decide accordingly. This is true even if his math is in fact mis-  
571 taken, and the tree would break under the load. If we take seriously the idea  
572 that we're modeling what follows from the agent's *actual* point of view, then  
573 our framework *should* say that he ought to decide to build the swing.

574 One might further explain the motivation behind this response by pointing  
575 out that the agent's perspective can also be incorrect due to empirical factors.  
576 For example, suppose the agent knows that swings are safe to install in trees  
577 of species A but not species B, but he is unsure what species his tree belongs  
578 to. He hires a tree expert to advise him, but due to a mixup, the expert tells  
579 him species A, which is the wrong answer. Yet, having no reason to distrust the  
580 expert, the agent comes to think that his tree belongs to species A. Again, the  
581 agent would be advised by our framework to conditionalize on this information,

582 and it would capture that, *from his perspective*, this is the sensible thing to do.  
583 On this view, which takes our framework to capture the rational way to reason  
584 given whatever input the agent has, the parallel between the two versions of the  
585 tree swing example is the correct result.

586 (ii) Embrace and Supplement

587 Even if we think that models that abandon factivity and embrace personal  
588 possibilities capture correctly how an agent should reason *given their perspective*,  
589 we might still want to be able to critically assess how they arrived at their  
590 perspective. It's one thing to think that *if* you believe you are the pope, you  
591 should infer from that that you're catholic. It's another to think that it's rational  
592 to begin with for you to think you are the pope (unless you are, which, gentle  
593 reader, is unlikely).

594 On this view, the personal probability model is an important ingredient  
595 in explaining what makes an agent's attitudes rational, but its verdicts are  
596 conditional in nature. If the agent's attitudes that serve as input for the model  
597 are rational, then the model tells the agent how to keep their attitudes coherent  
598 and update them. But the rational norms that govern inputs are external to the  
599 model. An account along these lines is defended by McHugh and Way (2018).

600 It's important to note that even standard Bayesian models that assume  
601 factivity and a classical logical possibility framework need to depend on such  
602 external norms to some extent. Suppose an agent becomes certain of an empir-  
603 ical truth simply by guessing correctly. If the agent then conditionalizes their  
604 credences on this truth, there is nothing in the standard Bayesian framework  
605 that would rule against the rationality of their attitude or reasoning. But we  
606 still want to say that it's irrational to become certain of something based on a  
607 pure guess, even if it the agent got lucky and guessed correctly. This verdict  
608 can only be delivered by a norm of rationality that is external to the Bayesian



609 framework.

610 In our current setup, the role of these external norms has to be significantly  
611 expanded, since false logical and mathematical beliefs are no longer constituting  
612 a violation of the rules of the model. Hence, we need a set of rational norms  
613 to supplement our model that judge which of the agent's attitudes have been  
614 rationally formed and which ones have not. This gives the overall theory a  
615 much greater degree of flexibility than the standard Bayesian framework, be-  
616 cause depending on the demandingness of one's views on rationality, verdicts  
617 about which empirical and logical judgments were rationally formed might vary  
618 considerably. For example, if our tree swing builder from before made a rather  
619 subtle error in his calculation, some theories of rational a priori belief might  
620 count his update as rationally permissible, while stricter theories might rule  
621 even subtle errors to be irrational. Similarly, depending on how sketchy the  
622 supposed tree expert appeared to be and what our standard for rational trust  
623 in testimony is, Bob's resulting high credence that his tree belongs to species A  
624 may or may not be considered rational.

625 But does the *embrace and supplement* strategy really alleviate the worry  
626 that the factivity-free personal probability framework is too permissive? We  
627 think one's answer to this question depends on how much of a contribution to  
628 a theory of rationality one expects from a normative formal model of rational  
629 credence. If one's expectations are fairly minimal, one might not worry about  
630 external norms doing too much of the heavy lifting. But for those who think  
631 that the formal model should be the central part of a theory of rational belief  
632 and updating, putting in so many constraints "by hand" won't be a satisfactory  
633 strategy.

634 (iii) Embrace and Reinterpret

635 Another possibility is to deny that these models are still normative.<sup>12</sup> Once

636 we move to personal possibilities, and the model just captures what the agent  
637 takes to be appropriate inputs to their reasoning, these models are better inter-  
638 preted as formal representations or *descriptions* of the agent’s actual reasoning.

639 How compelling one finds this suggestion partly depends on how one inter-  
640 prets the normative force of these models. For example, Dogramaci (2018b) is  
641 worried that once we move to personal possibilities, the constraints of the frame-  
642 work become essentially impossible to violate. He says that “any case where it  
643 would initially *appear* someone is violating it [the additivity principle] will be  
644 ultimately better described, and *correctly* described, as a case where they are  
645 not violating it. Suppose I initially appear to violate [the additivity principle]  
646 by saying there’s half a chance of rain tomorrow and half a chance of snow, and  
647 I think there’s three quarters of a chance of rain or snow. Any such case will  
648 be better described as one where I turn out to think it might both rain and  
649 snow, and thus there are simply more doxastic possibilities (dreamt of in my  
650 philosophy) than it first appeared [...]” Pettigrew pushes back, claiming that as  
651 long as the cognitive processes by which we rule out personally possible worlds  
652 are not identical to the processes by which we assign credences, violations of  
653 personal probabilism are possible.

654 We think that, while Dogramaci’s argument presents a serious challenge at  
655 the synchronic level, it is far less clear that the same reinterpretation strategy  
656 can be used to argue that the constraints imposed by superconditonalization  
657 are toothless. Take the agent from Dogramaci’s example, whose credences have  
658 been charitably reinterpreted to take seriously the possibility that it might both  
659 snow and rain, so that  $Cr(snow) = 0.5$ ,  $Cr(rain) = 0.5$ , and  $Cr(snow \vee rain) =$   
660  $0.75$ . Suppose the agent learns that it’s snowing. As a result, superconditioning  
661 provides some substantive constraints on the person’s updated credences, for  
662 example that  $Cr'(snow) = 1$ , and that  $Cr'(snow \vee rain) \geq 0.5$ . What if the

663 agent's updated credences diverge from this, so that, for instance,  $Cr'(snow) =$   
664  $1$  and  $Cr'(snow \vee rain) = 0.25$ ? If we wanted to reinterpret the agent's personal  
665 probabilities to try to make them look coherent, we would have to go back to  
666 adjust our initial interpretation of the agent's credences, and we might have to  
667 engage in some serious gerrymandering of personally possible worlds to achieve  
668 this. While this is certainly a strategy for immunizing agents from violating  
669 norms of rationality, we're not sure whether it's *always* possible to reinterpret  
670 the agent's starting credences to make their updates seem rational. But even if it  
671 is often possible to do so, the idea that this could be a legitimate way of ascribing  
672 mental states to the agent is not very plausible. When an agent updates their  
673 credences in a way that appears to violate superconditioning, it's not clear why  
674 we shouldn't take this observation at face value, rather than conclude that they  
675 had some rather bizzare set of initial personally possible worlds and resulting  
676 credences that rationalize this update.

677     One's take on this matter will likely depend in part on one's view of how  
678 to attribute mental states to agents. We won't decide this here. But suppose  
679 that you find yourself siding with Dogramaci's argument that these models are  
680 lacking in normative force. If personal probability models can't be violated, this  
681 does not mean that these models are automatically well suited to be descriptive  
682 of the agent's reasoning. There are many lively debates in cognitive science  
683 and psychology about the exact heuristics and strategies that generate our per-  
684 formance on various reasoning tasks. Descriptive theories of human reasoning  
685 usually make specific predictions about how humans will think about particular  
686 problems or approach cognitive tasks. These theories are not just supposed to  
687 accommodate the data after the fact. If personal probability models are really  
688 as malleable as Dogramaci claims, then they are too malleable to make those  
689 substantive predictions. But if they can make substantive predictions about

690 reasoning patterns, especially in diachronic cases, then Dogramaci's argument  
691 that the models have no normative force is unconvincing, since in that case,  
692 the model's prediction can be interpreted as a normative constraint on updat-  
693 ing one's credences. Hence, either these models put substantive constraints on  
694 credences, especially in diachronic contexts, or they don't. If they do, then a  
695 normative interpretation is feasible. If they don't, then those models can't make  
696 interesting descriptive predictions about how agents will reason. Those who are  
697 unconvinced by the normative interpretation of personal probability models are  
698 left to conclude that these models live in the no man's land of pointless formal  
699 constructions that lack an interesting philosophical application.

700 We said above that we take our audience to be those who are generally  
701 sympathetic to Bayesian models of rational belief and updating. So those who  
702 are dissatisfied with the three strategies just discussed must think that we took  
703 our modifications of the standard Bayesian framework too far. We will thus now  
704 discuss reasons for (iv) rejecting the switch from logical to personal possibilities,  
705 and for (v) keeping factivity.

#### 706 (iv) Reject Personal Possibilities

707 Swapping in personal for logical possibilities was supposed to modify Bayesian  
708 models to make them suitable for representing logical learning. But one might  
709 worry that Bayesian models are just the wrong tool for the job, and that even  
710 with the modification, we're trying to stretch them past their reasonable domain  
711 of application. How might one defend this position?

712 Take a standard Bayesian model and consider how it represents empirical  
713 learning. It takes a change in the agent's credences as input (this applies both in  
714 cases of standard and Jeffrey conditioning), and it outputs which new credence  
715 assignment the agent should adopt. In doing so, the model makes no reference  
716 to any sorts of reasoning steps the agent might undergo. Hence, it captures

717 neither how the agent might arrive at their new credences nor whether their  
718 new credences are properly based on their previous attitudes. This makes sense  
719 - firstly, there is plausibly more than one permissible cognitive path towards  
720 arriving at the correctly updated attitudes. Secondly, a probabilistic model  
721 doesn't have sufficient structure to capture the nature of the basing relations  
722 between the agent's attitudes. As a result, various authors have interpreted  
723 the Bayesian framework as representing relations of *propositional* rationality,  
724 rather than relations of *doxastic* rationality. This means that the framework  
725 shows us what the rational attitudes are for the agent *to adopt*, in light of  
726 their evidence and prior credences, but it doesn't show us whether the agent's  
727 credences are rationally held in the doxastic sense (Smithies 2015; Wedgwood  
728 2017; Dogramaci 2018a; Staffel 2019; Titelbaum 2019).

729 If we interpret the Bayesian framework in this way, then the standard appeal  
730 to logical possibilities makes a lot of sense. Assuming that logical and mathe-  
731 matical facts are knowable a priori, the agent is in a sense already in possession  
732 of the needed evidence that rationalizes the relevant credences in the standard  
733 framework. On this interpretation, it doesn't make sense to try to represent  
734 states of temporary logical ignorance in the framework if its real purpose is  
735 to show which attitudes are rationalized by the agent's evidence, regardless of  
736 whether the agent has worked this out already at the current moment. On this  
737 view, the use of personal probabilities to represent steps in logical reasoning is  
738 simply a confused repurposing of what the framework is supposed to model.

739 A common objection to this interpretation is that what an agent's evidence  
740 indicates to them should be somehow dependent on their cognitive abilities or  
741 recognitional capacities (Lord 2018; Turri 2010). Perhaps I have, in some sense,  
742 entailing evidence for or against Goldbach's conjecture, but it is still a stretch  
743 to say I have propositional justification for/against it, since it is completely

744 beyond me to figure this out. This view might propose to relativize Bayesian  
745 norms to what is within the agent’s cognitive reach to figure out. Still, this  
746 view can preserve the idea that the framework models what is propositionally  
747 rational for the agent (see Dogramaci (2018a, 2018b) for this kind of view). It  
748 endorses taking some steps towards de-idealizing standard Bayesianism, without  
749 endorsing the idea that it is suitable for modeling logical learning.

750 If, in light of these arguments, we resist the switch from logical to personal  
751 possibilities, then we avoid sanctioning positive credence assignments to logical  
752 and mathematical falsehoods as rational.<sup>13</sup> We can thus resist calling the mis-  
753 calculating tree swing builder rational. This is the case even if we get rid of  
754 factivity and allow rational agents to update on *empirical* falsehoods, such as  
755 the misleading testimony about the tree species.

756 (v) Keep Factivity

757 Some philosophers, especially those who favor the knowledge-first program  
758 in epistemology, might balk at the idea that agents can rationally update on  
759 false evidence. Knowledge-firsters tend to argue that our evidence must be  
760 known, and that it is a violation of rational norms to become certain of, and  
761 conditionalize on a falsehood. Cases in which agents become certain of and  
762 update on falsehoods despite trying to “do everything right”, like the case of  
763 the brain in the vat, or the example of Bob being misled by the arborist, are  
764 handled by saying that these norm violations are excused.<sup>14</sup>

765 Proponents of such a view might be tempted to think that Pettigrew’s  
766 pseudo-factive argument discussed in section 3 is better than our argument.  
767 But the sense in which Pettigrew’s argument assumes factivity doesn’t capture  
768 what knowledge-firsters want. The pseudo-factive argument doesn’t claim that  
769 it’s irrational to supercondition on a falsehood, quite the opposite. This does  
770 not capture the knowledge-firster’s claim that becoming certain of and condi-

771 tioning on a falsehood is always irrational. Hence, the knowledge-firster doesn't  
772 gain anything from endorsing Pettigrew's argument as opposed to ours.

773 In fact, the knowledge-firster might even prefer our argument to Pettigrew's,  
774 for the following reason: knowledge-firsters tend to argue that in cases in which  
775 an agent doesn't have knowledge, but in which the agent has done everything  
776 in their power to be a knower, the best course of action for them is to reason  
777 as if they had knowledge. In this type of case, the agent would be excused for  
778 violating knowledge norms, since it is only due to external forces not in their  
779 control that they failed to follow the knowledge norms. In those cases, agents  
780 who become certain of a falsehood should update by superconditionalization,  
781 because that's what would be uniquely rational for someone who has knowledge.  
782 Our argument delivers this verdict, but Pettigrew's doesn't.<sup>15</sup>

783 Knowledge-firsters thus need to appeal to norms external to the model in  
784 order to impose a rational prohibition on becoming certain of, and condition-  
785 alizing on, falsehoods. But if they are interested in formalizing their idea that  
786 unlucky non-knowers are excused for their norm-violations if they update like  
787 knowers, they might very well prefer our more robust arguments for supercon-  
788 ditionalization to Pettigrew's pseudo-factive arguments. The resulting position  
789 would ultimately turn out to be a version of the "embrace and supplement"  
790 strategy discussed above.

## 791 **6 Conclusion**

792 We have offered an argument for superconditionalization in a personal possibility  
793 framework, which shows that if an agent becomes certain of an empirical or  
794 logical claim, the uniquely rational updating strategy is superconditionalization,  
795 regardless of whether the learned claim is true or false. This means that we're  
796 greatly expanding the applicability of the superconditionalization norm. By

797 using personal possibilities instead of logical ones, the norm applies to cases of  
798 logical learning, which it doesn't cover in standard Bayesian models. Further,  
799 since our argument avoids assuming pseudo-factivity, it more robustly supports  
800 superconditionalization as the *uniquely* rational updating rule than Pettigrew's  
801 argument.

802 Yet, one might have mixed feelings about such a far-reaching version of su-  
803 perconditionalization. As we saw, it applies even in cases like Bob's, who makes  
804 an avoidable mathematical error when calculating whether his tree can support  
805 a swing. Three possible reactions to this result stood out as most attractive in  
806 our discussion in section 5. Readers are invited to pick their favorite.

807 We think that the most promising way of *resisting* our argument is to say  
808 that the Bayesian framework is unsuitable for modeling logical learning. On this  
809 view, Bayesian models are best seen as modeling relations of propositional jus-  
810 tification, which hold independently of whether the agent has recognized them  
811 through reasoning. Such a view might still embrace Rescorla's and Pettigrew's  
812 arguments that agents should update by (super-)conditionalization when they  
813 learn empirical falsehoods, but it would resist the switch to personal possibili-  
814 ties.

815 If we accept the switch to personal possibilities, there are two plausible in-  
816 terpretations of our results. The first one, "embrace completely", welcomes  
817 our expansion of superconditionalization, and interprets the resulting models  
818 as showing us what is rational from the agent's own perspective. It takes the  
819 agent's attitudes as a fixed input without passing judgment on them, and shows  
820 which reasoning moves seem rational from the agent's perspective. This inter-  
821 pretation is quite radical, as it doesn't make room for the idea that for certain  
822 irrational inputs, agents should not reason with them, but instead try to correct  
823 them.



824 A less radical interpretation suggests to “embrace and supplement” our ar-  
 825 gument. The idea here is that we supplement our models with separate, exter-  
 826 nal norms for evaluating the attitudes that serve as inputs to the model, and  
 827 then the formalism shows how the agent should update. This allows us to say  
 828 that for an agent to be rational, their input attitudes must be rational, and  
 829 they must be personally probabilistic and update by superconditionalization.  
 830 This proposal is a way of spelling out McHugh and Way’s account of good rea-  
 831 soning which says that a pattern of reasoning is good if it leads agents from  
 832 fitting input attitudes to fitting output attitudes (McHugh and Way 2018). On  
 833 this interpretation of our view, external norms make a significant contribution  
 834 to determining whether an agent has rational credences. But even standard  
 835 Bayesian views must rely on such external norms to some degree, which means  
 836 that we would be merely expanding our reliance on them. One advantage of the  
 837 “embrace and supplement” strategy is that it can formulate model-independent  
 838 norms for rational input attitudes for both empirical and logical claims, whereas  
 839 the standard Bayesian framework comes with fixed norms for rational attitudes  
 840 towards logical claims.

## 841 Appendix

**Lemma 1.** *Let  $\mathcal{F} = \{X_1, \dots, X_m\}$  be a set of credal objects over a set of worlds  $W$ . Given real numbers  $a_1, \dots, a_m \in \mathbb{R}$ , let  $A : W \rightarrow [0, 1]$  be an act defined as  $A(w) = \sum_{i=1}^m a_i v_w(X_i)$  for all  $w \in W$ . If  $p : W \rightarrow [0, 1]$  coherently extends a personally probabilistic credence function  $c : \mathcal{F} \rightarrow [0, 1]$ , then:*

$$\sum_{w \in W} p(w)A(w) = \sum_{i=1}^m c(X_i)a_i$$

842 *Proof.* Consider a  $p : W \rightarrow [0, 1]$  that coherently extends a personally proba-

843 bilistic credence function  $c : \mathcal{F} \rightarrow [0, 1]$ . For each  $1 \leq i \leq m$ , define an act  
 844  $A_i : W \rightarrow [0, 1]$  such that  $A_i(w) = a_i v_w(X_i)$  for all  $w \in W$ . So we have that:

$$\begin{aligned} \sum_{w \in W} p(w)A(w) &= \sum_{w \in W} p(w) \sum_{i=1}^m A_i(w) \\ &= \sum_{i=1}^m \sum_{w \in W} p(w)A_i(w) \end{aligned}$$

845 Splitting the inner sum according to the truth value of  $X_i$ , for  $1 \leq i \leq m$  it  
 846 holds that:

$$\begin{aligned} \sum_{w \in W} p(w)A_i(w) &= \sum_{w, v_w(X_i)=1} p(w)A_i(w) + \sum_{w, v_w(X_i)=0} p(w)A_i(w) \\ &= \sum_{w, v_w(X_i)=1} p(w)a_i + \sum_{w, v_w(X_i)=0} p(w)0 \\ &= a_i \sum_{w, v_w(X_i)=1} p(w) \end{aligned}$$

847 As  $p$  coherently extends  $c$ ,  $\sum_{w, v_w(X_i)=1} p(w) = c(X_i)$ . Finally, we can conclude  
 848 that:

$$\sum_{w \in W} p(w)A(w) = \sum_{i=1}^m c(X_i)a_i$$

849 □

850 **Theorem 3 .** Let  $c : \mathcal{F} \rightarrow [0, 1]$  be a personally probabilistic credence function  
 851 and  $c' = \langle c'_1, \dots, c'_n \rangle$  be the set of possible future personally probabilistic credence  
 852 functions defined over  $\mathcal{F}$ .

- 853 (a) If  $\langle c, c' \rangle$  violates *wGRP*, then it is vulnerable to a Strong Dutch Strategy.  
 854 (b) If  $\langle c, c' \rangle$  satisfies *wGRP*, then it is not vulnerable to a Strong Dutch Strategy.

*Proof.* (a) Supposing  $\mathcal{F} = \{X_1, \dots, X_m\}$ , any credence function  $\hat{c} : \mathcal{F} \rightarrow [0, 1]$

can be viewed as a vector  $\langle \hat{c}(X_1), \dots, \hat{c}(X_m) \rangle \in \mathbb{R}^m$ . If  $\langle c, c' \rangle$  violates wGRP, then  $c$  is not in the convex hull of the vectors  $c'_1, \dots, c'_n$ . Thus, by the Separating Hyperplane Theorem, there are numbers  $a, a' \in \mathbb{R}$  and a vector  $b \in \mathbb{R}^m$  such that, for each  $c'_j \in c'$ :

$$\sum_{i=1}^m c(X_i)b_i < a < a' < \sum_{i=1}^m c'_j(X_i)b_i$$

855 Now consider the constant acts  $A : W \rightarrow [0, 1]$  and  $A' : W \rightarrow [0, 1]$  such that  
856  $A(w) = a$  and  $A'(w) = -a'$  for all  $w \in W$ . Furthermore, consider the act  
857  $B : W \rightarrow [0, 1]$  defined in the following way:  $B(w) = \sum_{i=1}^m b_i v_w(X_i)$ . That is, for  
858 each world  $w \in W$ , determine those  $X_i$  that are true and sum the corresponding  
859  $b_i$  to obtain  $B(w)$ . If all  $X_i \in \mathcal{F}$  are false in  $w$ , then  $B(w) = 0$ . Now, define the  
860 act  $B' : W \rightarrow [0, 1]$  via  $B'(w) = -B(w)$  for all  $w \in W$ . By Lemma 1, for any  $p :$   
861  $W \rightarrow [0, 1]$  that coherently extends  $c$  we have that  $\sum_w p(w)B(w) = \sum_{i=1}^m c(X_i)b_i$ .  
862 As  $\sum_w p(w)A(w) = a$ ,  $c$  prefers  $A$  to  $B$ . Analogously, by Lemma 1, each  $c'_j$   
863 prefers  $A'$  to  $B'$ . Finally, note that  $B(w) + B'(w) = 0 > a - a' = A(w) + A'(w)$   
864 for any  $w \in W$ , therefore  $\langle c, c' \rangle$  is vulnerable to a Strong Dutch Strategy.

865 (b) Assume there are weights  $\lambda_j \in [0, 1]$ , summing up to one, such that  
866  $c(X) = \sum_j \lambda_j c'_j(X)$  for all  $X \in \mathcal{F}$ . To prove by contradiction, suppose there  
867 are acts  $A, A', B, B'$  such that  $c$  prefers  $A$  to  $B$ , each  $c_j$  prefers  $A'$  to  $B'$ , but  
868  $A(w) + A'(w) < B(w) + B'(w)$  at any  $w \in W$ . For each  $c'_j$ , consider a credence  
869 function  $p'_j : W \rightarrow [0, 1]$  that coherently extends it, thus preferring  $A'$  to  $B'$ .  
870 Defining  $p : W \rightarrow [0, 1]$  via  $p(w) = \sum_j \lambda_j p'_j(w)$  for all  $w \in W$ , we have that  
871  $\langle p, p' \rangle$  satisfies wGRP. Note that  $p$  coherently extends  $c$ , hence preferring  $A$   
872 to  $B$ . Consequently,  $\langle p, p' \rangle$  is vulnerable to a Strong Dutch Strategy, which  
873 contradicts Theorem 1(b).  $\square$

874 **Theorem 4 .** *Let  $c : \mathcal{F} \rightarrow [0, 1]$  be a personally probabilistic credence function*

875 and  $c' = \langle c'_1, \dots, c'_n \rangle$  be the set of possible future personally probabilistic credence  
 876 functions defined over  $\mathcal{F}$ .

877 (a) If  $\langle c, c' \rangle$  violates wGRP, then it is accuracy dominated.

878 (b) If  $\langle c, c' \rangle$  satisfies wGRP, then it is not accuracy dominated.

879 *Proof.* (a) See the ( $\rightarrow$ )-part of the proof of Theorem 2 (Pettigrew 2021a), just  
 880 interpreting  $W$  as a set of personally possible worlds.

881 (b) Suppose  $\langle c, c' \rangle$  satisfies wGRP, so that there are  $\lambda_1, \dots, \lambda_n \in [0, 1]$  such  
 882 that  $\sum_j \lambda_j = 1$  and  $c(X) = \sum_{j=1}^n \lambda_j c_j(X)$  for any  $X \in \mathcal{F}$ . To prove by con-  
 883 tradiction, assume  $\langle c, c' \rangle$  is accuracy dominated: there are credence functions  
 884  $c^*, c_1^*, \dots, c_n^*$ , defined on  $\mathcal{F}$ , such that  $\mathfrak{J}(c, w) + \mathfrak{J}(c'_j, w) > \mathfrak{J}(c^*, w) + \mathfrak{J}(c_j^*, w)$  for  
 885 all  $w \in W$  and  $1 \leq j \leq n$ . Since accuracy is measured with an additive strictly  
 886 proper  $\mathfrak{J}$ , there is a strictly proper scoring rule  $s$  such that, for any credence  
 887 function  $\hat{c}$ ,  $\mathfrak{J}(\hat{c}, w) = \sum_{X \in \mathcal{F}} s(v_w(X), \hat{c}(X))$ . For  $s$  is strictly proper, we have that,  
 888 for any  $X \in \mathcal{F}$  and any  $1 \leq j \leq n$ :

$$c(X)s(1, c(X)) + (1 - c(X))s(0, c(X)) \leq c(X)s(1, c^*(X)) + (1 - c(X))s(0, c^*(X)) \quad (1)$$

$$c'_j(X)s(1, c'_j(X)) + (1 - c'_j(X))s(0, c'_j(X)) \leq c'_j(X)s(1, c_j^*(X)) + (1 - c'_j(X))s(0, c_j^*(X)) \quad (2)$$

889 If we replace those  $c(X)$  out of the scope of  $s(\cdot)$  by  $\sum_j \lambda_j c_j(X)$  in Expression  
 890 (1) (note also that  $\sum_j \lambda_j = 1$ ) and, in Expression (2), multiply both sides by  
 891  $\lambda_j$  before summing for all  $j$ , we obtain, respectively:

$$\begin{aligned} & \sum_j \lambda_j c'_j(X)s(1, c(X)) + \sum_j \lambda_j (1 - c'_j(X))s(0, c(X)) \leq \\ & \sum_j \lambda_j c'_j(X)s(1, c^*(X)) + \sum_j \lambda_j (1 - c'_j(X))s(0, c^*(X)) \quad (3) \\ & \sum_j \lambda_j [c'_j(X)s(1, c'_j(X)) + (1 - c'_j(X))s(0, c'_j(X))] \leq \end{aligned}$$

$$\sum_j \lambda_j [c'_j(X)s(1, c_j^*(X)) + (1 - c'_j(X))s(0, c_j^*(X))] \quad (4)$$

892 Grouping the summations in  $j$  in each side of Expression (3), it can be added  
893 to Expression (4) to obtain:

$$\begin{aligned} & \sum_j \lambda_j [c'_j(X)(s(1, c(X)) + s(1, c'_j(X))) + (1 - c'_j(X))(s(0, c(X)) + s(0, c'_j(X)))] \leq \\ & \sum_j \lambda_j [c'_j(X)(s(1, c^*(X)) + s(1, c_j^*(X))) + (1 - c'_j(X))(s(0, c^*(X)) + s(0, c_j^*(X)))] \quad (5) \end{aligned}$$

894 Since  $\langle c, c' \rangle$  is accuracy dominated by  $c^*$  and  $\langle c_1^*, \dots, c_n^* \rangle$ ,  $\mathfrak{J}(c, w) + \mathfrak{J}(c'_j, w) >$   
895  $\mathfrak{J}(c^*, w) + \mathfrak{J}(c_j^*, w)$  for all  $w \in W$  and  $1 \leq j \leq n$ . Multiplying each side of this  
896 inequality by  $\lambda_j p_j(w)$ , for a  $p_j$  that coherently extends  $c'_j$ , and summing for all  
897  $w \in W$  and all  $1 \leq j \leq n$ , we obtain:

$$\sum_j \lambda_j \sum_w p_j(w) [\mathfrak{J}(c, w) + \mathfrak{J}(c'_j, w)] > \sum_j \lambda_j \sum_w p_j(w) [\mathfrak{J}(c^*, w) + \mathfrak{J}(c_j^*, w)] \quad (6)$$

898 Recall that, for any  $\hat{c} : \mathcal{F} \rightarrow [0, 1]$ ,  $\mathfrak{J}(\hat{c}, w) = \sum_{X \in \mathcal{F}} s(v_w(X), \hat{c}(X))$ . Thus, for any  
899  $\hat{c}$ ,  $\mathfrak{J}(\hat{c}, w)$  can be rewritten as:

$$\begin{aligned} \mathfrak{J}(\hat{c}, w) &= \sum_{X \in \mathcal{F}} v_w(X) s(1, \hat{c}(X)) + \sum_{X \in \mathcal{F}} (1 - v_w(X)) s(0, \hat{c}(X)) \\ &= \sum_{X \in \mathcal{F}} v_w(X) s(1, \hat{c}(X)) - \sum_{X \in \mathcal{F}} v_w(X) s(0, \hat{c}(X)) + \sum_{X \in \mathcal{F}} s(0, \hat{c}(X)) \quad (7) \end{aligned}$$

900 If  $\hat{c}$  is fixed, the first two summations in Expression (7) can be viewed as acts in  
901 the format  $\sum_i a_i v_w(X_i)$ . Hence, applying Lemma 1 to  $\sum_w p_j(w) \mathfrak{J}(c, w)$  yields:

$$\sum_{X \in \mathcal{F}} c'_j(X) s(1, c(X)) - \sum_{X \in \mathcal{F}} c'_j(X) s(0, c(X)) + \sum_w p_j(w) \sum_{X \in \mathcal{F}} s(0, c(X))$$

902 As  $\sum_w p_j(w) = 1$ , a bit of algebraic manipulation results in:

$$\sum_{X \in \mathcal{F}} [c'_j(X)s(1, c(X)) + (1 - c'_j(X))s(0, c(X))]$$

903 Analogously, we can apply Lemma 1 to each  $\sum_w p_j(w)\mathcal{J}(\cdot, w)$  resulting from  
 904 Expression (6), obtaining:

$$\begin{aligned} & \sum_j \sum_{X \in \mathcal{F}} \lambda_j [c'_j(X)(s(1, c(X)) + s(1, c'_j(X))) + (1 - c'_j(X))(s(0, c(X)) + s(0, c'_j(X)))] > \\ & \sum_j \sum_{X \in \mathcal{F}} \lambda_j [c'_j(X)(s(1, c^*(X)) + s(1, c'_j(X))) + (1 - c'_j(X))(s(0, c^*(X)) + s(0, c'_j(X)))] \quad (8) \end{aligned}$$

905 But note that this is just the negation of Expression (5) summed for all  $X \in \mathcal{F}$ ,  
 906 which is a contradiction, completing the proof.  $\square$

907 **Theorem 5 .** *Let  $c$  be a credence function. Let  $c' = \langle c'_1, \dots, c'_n \rangle$  be an updating*  
 908 *rule for a partition  $\{E_1, \dots, E_n\}$  of  $W$ , respecting it. If the pair  $\langle c, c' \rangle$  satisfies*  
 909 *the Weak General Reflection Principle, then  $\langle c, c' \rangle$  is superconditioning.*

910 *Proof.* As  $c'$  respects the partition  $\{E_1, \dots, E_n\}$ , for each  $1 \leq i \leq n$  there is a  
 911 function  $p_i : W \rightarrow [0, 1]$ , with  $\sum_{w \in W} p_i(w) = 1$ , such that  $c'_i(X) = \sum_{w \in W} p_i(w)v_w(X)$   
 912 for all  $X \in \mathcal{F}$  and  $\sum_{w \in E_i} p_i(w) = 1$ . wGRP implies that, for some  $\lambda_1, \dots, \lambda_n \in$   
 913  $[0, 1]$  with  $\sum_{i=1}^n \lambda_i = 1$ , we have that  $c(X) = \sum_{i=1}^n \lambda_i c'_i(X)$  for all  $X \in \mathcal{F}$ . Thus,  
 914  $c(X) = \sum_{i=1}^n \lambda_i \sum_{w \in W} p_i(w)v_w(X) = \sum_{w \in W} \sum_{i=1}^n \lambda_i p_i(w)v_w(X)$  for all  $X \in \mathcal{F}$ . Let the  
 915 function  $p : W \rightarrow [0, 1]$  be such that  $p(w) = \sum_{i=1}^n \lambda_i p_i(w)$  for all  $w \in W$ . Note  
 916 that  $\sum_{w \in W} p(w) = 1$ .

Now consider an element  $E_j$  of the partition. We have, for all  $X \in \mathcal{F}$ ,  
 that  $\sum_{w \in E_j} p(w)v_w(X) = \sum_{w \in E_j} v_w(X) \sum_{i=1}^n \lambda_i p_i(w)$ . Given that, for any  $w \in$   
 $E_j$ ,  $p_i(w) = 0$  whenever  $i \neq j$ , as the partition is respected, it follows that  
 $\sum_{w \in E_j} v_w(X) \sum_{i=1}^n \lambda_i p_i(w) = \sum_{w \in E_j} v_w(X) \lambda_j p_j(w) = \lambda_j \sum_{w \in W} v_w(X) p_j(w)$ . For all

$w \in E_j$ , we also have that  $\sum_{w \in E_j} p(w) = \sum_{w \in E_j} \sum_{i=1}^n \lambda_i p_i(w) = \sum_{w \in E_j} \lambda_j p_j(w)$ ,  
 thus  $\sum_{w \in E_j} p(w) = \lambda_j \sum_{w \in E_j} p_j(w) = \lambda_j$ . Finally, for each  $1 \leq j \leq n$  with  
 $\sum_{w \in E_j} p(w) = \lambda_j > 0$ , we obtain, for all  $X \in \mathcal{F}$ :

$$c'_j(X) = \sum_{w \in W} p_j(w) v_w(X) = \frac{\lambda_j \sum_{w \in W} p_j(w) v_w(X)}{\lambda_j} = \frac{\sum_{w \in E_j} p(w) v_w(X)}{\sum_{w \in E_j} p(w)}$$

917 Consequently,  $\langle c, c' \rangle$  is superconditioning. □

## 918 Notes

920 <sup>1</sup>Henceforth, we will use “learn” and “become certain of” interchangeably. Hence, when  
 921 we say that an agent learns something, what they learn can be true or false.

922 <sup>2</sup>One objection is worth mentioning here, although we won’t discuss it in detail, as it would  
 923 lead us away from our main train of thought. It says that becoming *certain* of empirical  
 924 evidence is always irrational. Proponents of the regularity principle argue that we should at  
 925 most invest high credence in any empirical propositions, but not credence 1. This makes Jeffrey  
 926 conditionalization the only relevant update rule, thus omitting the need for a factivity-free  
 927 version of conditioning. We think that the case for regularity is unconvincing. For discussion,  
 928 see Rescorla (2020) and especially Hájek (2012). Also, proponents of contextualist versions of  
 929 Bayesianism give good reasons to avoid regularity as a requirement (Greco 2017; Salow 2019).

930 <sup>3</sup>A similar problem would occur if  $E$  was true but the agent became certain of a  $E' \neq E$ .  
 931 In that case, the bet at  $t_2$  would not occur, for the bookie would also be certain of  $E'$ , but  
 932 the standard Dutch book relies on it, for  $E$  is true.

933 <sup>4</sup>Abusing the notation, we use  $w$  for both a possible world and the proposition that is true  
 934 only there.

935 <sup>5</sup> $s$  is strictly proper if only  $x = p$  minimizes  $ps(1, x) + (1 - p)s(0, x)$  for any  $p \in [0, 1]$ .

936 <sup>6</sup>We assume in our argument below that  $W$  is kept fixed after learning, to avoid pseudo-  
 937 factivity, while in Pettigrew’s argument this set is narrowed to the learned  $E_i$ , which makes  
 938 his arguments pseudo-factive.

939 <sup>7</sup>Proofs can be found in the Appendix.

940 <sup>8</sup>Pseudo-factivity would do the job, but we are trying to avoid it.

941 <sup>9</sup>Note this property is, given (Personal) Probabilism, implied by but weaker than restricting  
942 the set of worlds  $W$  to those compatible with the learned  $E_j$ .

943 <sup>10</sup>Also Rescorla (2020) assumes a similar assumption, for an agent that became certain of  
944 an  $E_i$  to accept any bet conditional on a  $E_j \neq E_i$ .

945 <sup>11</sup>If  $E_i \notin \mathcal{F}$ ,  $c'_i$  implies  $c'_i(E_i) = 1$  if it is the only credence that personally probabilistically  
946 extends  $c'_i$  to  $E_i$ .

947 <sup>12</sup>Thanks to Kenny Easwaran for suggesting that we discuss this option.

948 <sup>13</sup>Except perhaps in cases in which it is beyond an agent's cognitive capacities to assign the  
949 correct credences.

950 <sup>14</sup>This strategy is defended by Williamson (forthcoming). For a critical discussion, see  
951 Greco (forthcoming).

952 <sup>15</sup>Notice that moving to the standard arguments for conditioning would not help much here.  
953 In the standard framework, it's irrational for the agent to become certain of a logical falsehood.  
954 But when an agent becomes certain of an empirical falsehood, the arguments give the same  
955 verdict as Pettigrew's version from section 3: conditionalizing on the falsehood avoids Dutch-  
956 bookability and accuracy dominance, but it's not shown that it's the unique strategy with  
957 these properties.

## 958 References

959 Briggs, R. A. and R. Pettigrew (2020). An accuracy-dominance argument for  
960 conditionalization. *Noûs* 54(1), 162–181.

961 Comesana, J. (2020). *Being Rational and Being Right*. Oxford University  
962 Press.

963 de Finetti, B. (1974). *Theory of probability*. John Wiley and Sons, Chichester.

964 Dogramaci, S. (2018a). Rational credence through reasoning. *Philosophers'*  
965 *Imprint* 18.

966 Dogramaci, S. (2018b). Solving the problem of logical omniscience. *Philosophical*  
967 *Issues* 28(1), 107–128.

968 Greco, D. (2017). Cognitive mobile homes. *Mind* 126(501), 93–121.



- 969 Greco, D. (2021). Justifications and excuses in epistemology. *Noûs* 55(3),  
970 517–537.
- 971 Hacking, I. (1967). Slightly more realistic personal probability. *Philosophy of*  
972 *Science* 4, 311–325.
- 973 Hájek, A. (2012). Is strict coherence coherent? *Dialectica* 66(3), 411–424.
- 974 Lord, E. (2018). *The Importance of Being Rational*. Oxford, UK: Oxford  
975 University Press.
- 976 McHugh, C. and J. Way (2018). What is good reasoning? *Philosophy and*  
977 *Phenomenological Research*, 153–174.
- 978 Pettigrew, R. (2021a). Bayesian updating when what you learn might be false.  
979 *Erkenntnis*, 1–16.
- 980 Pettigrew, R. (2021b). Logical ignorance and logical learning. *Syn-*  
981 *these* 198(10), 9991–10020.
- 982 Rescorla, M. (2020). An improved dutch book theorem for conditionalization.  
983 *Erkenntnis.*, 1–29.
- 984 Salow, B. (2019). Elusive externalism. *Mind* 128(510), 397–427.
- 985 Smithies, D. (2015). Ideal rationality and logical omniscience. *Syn-*  
986 *these* 192(9), 2769–2793.
- 987 Staffel, J. (2019). *Unsettled Thoughts: A Theory of Degrees of Rationality*.  
988 Oxford University Press.
- 989 Titelbaum, M. (2019). Return to reason. In A. S.-P. M. Skipper (Ed.), *Higher-*  
990 *Order Evidence: New Essays*. Oxford: Oxford University Press.
- 991 Turri, J. (2010). On the relationship between propositional and doxastic jus-  
992 tification. *Philosophy and Phenomenological Research* 80(2), 312–326.
- 993 Wedgwood, R. (2017). *The Value of Rationality*. Oxford University Press.

994 Williamson, T. (forthcoming). Justifications, excuses, and sceptical scenarios.  
995 In F. Dorsch and J. Dutant (Eds.), *The New Evil Demon*. Oxford: Oxford  
996 University Press.