

# Bayesianism & Explanationism

Branden Fitelson



<http://fitelson.org/B&E.pdf>

There are two main kinds of epistemic requirements. First, let's examine these in the case of full (all-or-nothing) belief.

- **Correctness Requirements** (for full belief)
  - If  $p$  is false, then believing that  $p$  is incorrect.
  - Correctness requirements are *factive* — they involve relations between one's attitudes and *the facts* [29, 23, 30].
- **Rational Requirements** (for full belief)
  - If one believes both  $p$  and  $\neg p$ , then one's beliefs are *structurally* irrational (*viz.*, *incoherent*).
  - If one's total evidence  $K$  counter-supports  $p$ , then believing that  $p$  is *substantively* irrational.
  - Rational requirements are *non-factive* — they involve **either** (a) relations (of coherence) among one's attitudes (*structural* rationality) **or** (b) relations between one's total evidence  $K$  and one's attitudes (*substantive* rationality) [31].

This distinction also applies to Bayesian epistemology...

Let us suppose that our agent is equipped with a *confidence ordering*  $\geq$  over propositions. So that  $p \geq q$  iff  $S$  is *at least as confident* in the truth of  $p$  as in the truth of  $q$ .

And,  $p > q$  iff  $S$  is *strictly more confident* in  $p$  than in  $q$ .

- **Correctness Requirements** (for comparative confidence)
  - If  $p$  is false and  $q$  is true, then  $p > q$  is incorrect.
    - Heuristically, correctness requirements involve *agreement with the attitudes of an omniscient agent* [16, 14, 20].
- **Rational Requirements** (for comparative confidence)
  - Intransitive  $\geq$  relations are *structurally* irrational. *E.g.*, the combination of attitudes  $\{p > q, q > r, r > p\}$  is *incoherent*.
  - If one's total evidence  $K$  supports  $p$  strictly more strongly than  $K$  supports  $q$ , then  $q \geq p$  is *substantively* irrational.

Bayesians also like to talk about *numerical degrees of confidence* (credences). Indeed, this will be our main focus today...

Let us suppose that our agent has (at each time  $t$ ) *degrees of confidence* (*credences*) represented by a function  $cr_t(\cdot)$ .

The function  $cr_t(\cdot | E)$  will reflect the agent's degrees of confidence (at  $t$ ) — *on the indicative supposition that  $E$  is true*.

Finally, suppose our agent learns (exactly)  $E$  (with certainty) between  $t_0$  and  $t_1$ . So, her transition from  $cr_{t_0}(\cdot)$  to  $cr_{t_1}(\cdot)$  reflects the upshot of learning (precisely) the content  $E$ .

**Structural Bayesianism** involves a Trinity of Requirements [24].

1. **Synchronic (Non-Suppositional) Probabilism.** At each time  $t$ , the agent's (non-suppositional) credence function  $cr_t(\cdot)$  should obey the (Kolmogorov) probability axioms.
2. **Synchronic (Suppositional) Ratio Formula.** At each time  $t$ , the agent's (suppositional) credence function  $cr_t(\cdot | E)$  should obey the ratio formula  $cr_t(\cdot | E) = \frac{cr_t(\cdot \& E)}{cr_t(E)}$ .
3. **Diachronic Conditionalization.**  $cr_{t_1}(\cdot)$  should equal  $\frac{cr_{t_0}(\cdot \& E)}{cr_{t_0}(E)}$ .

The Fourth Pillar of Structural Bayesianism *follows from* (1)–(3).

4. **Learning & Supposing.**  $cr_{t_1}(\cdot)$  should equal  $cr_{t_0}(\cdot | E)$ .

de Finetti [3] gave Dutch Book Arguments for both (1) and (2).

Accuracy-dominance arguments for (1) are now popular [20]. One can also give an accuracy-dominance argument for (2) [12].

Lewis and others have aimed to adapt de Finetti’s argument for (2) into a “diachronic Dutch Book” argument for (3) [21].

Like de Finetti, I view (1) and (2) as the *fundamental* Bayesian principles governing the (structural) rationality of credences.

I have always been more skeptical about the existence and nature of “diachronic coherence requirements.” Indeed, there seem to be *many* reasons to worry about (3) & (4) [13, 32, 25].

I will focus, primarily, on the synchronic requirements (1) & (2), and how they (allegedly) interact with Explanationism.

There are also *substantive* Bayesian requirements. In general, these are of the following generic form: If one’s total evidence at time  $t$  ( $K_t$ ) is such and so, then  $cr_t(\cdot)$  should be thus and such.

Let  $\{H_1, \dots, H_n\}$  be some partition of alternative hypotheses (putative explanations of  $E$ ) entertained by an agent (at time  $t$ ).

**Substantive Bayesianism** (some example requirements)

- **The Principle of Indifference** [8]. If  $K_t$  does not favor any  $H_i$  over any  $H_j$ , then  $cr_t(H_i)$  should equal  $cr_t(H_j)$ ,  $\forall i, j$ .
- **The Principal Principle** [19]. If  $K_t$  entails that the objective chance of  $H_i$  is  $c$  (and  $K_t$  doesn’t contain/imply any inadmissible evidence), then  $cr_t(H_i)$  should equal  $c$ .
- **The Requirement of Total Evidence** [17, 1, 30].  $cr_t(H_i)$  should be equal to the evidential probability  $Pr(H_i | K_t)$ .

Next up: Explanationism. I will follow Douven’s [5] recent discussion of (the various explications of) Explanationism.

As Douven [5] explains, there have been various historical views regarding the proper formulation of Explanationism.

Rather than rehearsing Douven’s list of historical explications in detail, I will remain at a higher level of abstraction.

☞ The idea behind *Explanationism* is that some “epistemic credit” should accrue to a hypothesis  $H$  in virtue of its being the best (or only) explanation of  $E$  (among the available alternatives  $\{H_k\}$ ).

The question in which I am interested is: What is the best way to accommodate this basic Explanationist idea — of “credit” accruing to  $E$ ’s best explanation — *within a Bayesian framework?*

van Fraassen [27] and Douven [5] maintain that a Bayesian should incorporate Explanationism *by revising some of the basic requirements of Structural Bayesianism: (1)–(3)*.

I will focus on Douven’s proposal, since it is more precise, and it can be couched in purely synchronic terms [as a revision of (2)].

Douven recommends that Bayesians *revise the Ratio Formula* (2) in such a way that the following alternative to Bayes’s Theorem is adopted (for hypotheses  $H_i \in \{H_1, \dots, H_n\}$  and evidence  $E$ ).

$$\text{EXPL. } cr_t(H_i | E) = \frac{cr_t(H_i) \cdot cr_t(E | H_i) + c(H_i, E)}{\sum_{k=1}^n [cr_t(H_k) \cdot cr_t(E | H_k) + c(H_k, E)]'}$$

where  $c(H, E) \in [0, 1)$  is  $H$ ’s “ $E$ -abductive credit score.” And,  $c(H, E) > 0$  iff  $H$  best explains  $E$  (o.w.  $c(H, E) = 0$ ).<sup>1</sup>

Note that *if there is no best explanation* of  $E$  among the  $\{H_k\}$ , then **EXPL** reduces to Bayes’s Theorem (since all of the credit scores  $c(H_k, E)$  will be equal to zero in such a case).

Because **EXPL** leads to violations of Structural Bayesianism [*viz.*, (2)], Douven discusses various ways a defender of **EXPL** might respond to de Finetti’s [3] Dutch Book argument for (2).

<sup>1</sup>This definition can’t be quite right, since sometimes  $c(H, E)$  will need to be  $\geq 1$  in order to emulate some Bayesian abductive updates. See Extras.

☞ There is a more elegant way to capture the idea of “abductive credit,” which avoids the probabilistic incoherencies of [5, 27].

Note that there is something odd about thinking of **EXPL** as a *structural* requirement in the first place. **EXPL** has this form:

**EXPL.**  $cr_t(H_i | E)$  should receive a boost — over and above the value prescribed by Bayes’s Theorem — just in case  $(\mathbb{A}_i) H_i$  is the best explanation of  $E$  (among the  $\{H_k\}$ ).

Read literally, then, **EXPL** relates  $cr_t(H_i | E)$  to *the fact* that  $H_i$  best explains  $E$ , which makes **EXPL** a *correctness* requirement.

In order for **EXPL** to be a *structural* requirement, it would have to be stated *as a relation among the agent’s credences*. To wit:

**EXPL<sub>1</sub>.**  $cr_t(H_i | E)$  should receive a boost — over & above its Bayes’s Thm value — iff *the agent is certain at  $t$  that*  $\mathbb{A}_i$ .<sup>2</sup>

<sup>2</sup>We also have to be careful to allow this “boost” to occur *only once* — presumably, when  $\mathbb{A}_i$  is *first learned*. Otherwise,  $H_i$  will be “over-boosted.”

Even **EXPL<sub>1</sub>** is arguably not a (pure) structural requirement, since it only constrains agents who *entertain*  $\mathbb{A}_i$ , for some  $H_i$  and  $E$ .

Pure structural requirements (*e.g.*, probabilism) do not presuppose anything about the *contents* of an agent’s attitudes.

☞ If **EXPL** is going to be a (non-vacuous) *rational* requirement, then it must presuppose that the agent *entertains*  $\mathbb{A}_i$ , for some  $H_i$  and  $E$ . And, in that case, I think a preferable explication exists.

**Evidential Relevance of Abduction (ERA).** Let  $\mathbb{A}_j$  assert that  $H_j$  is the best (or, better still, the *only* [4, 9]) explanation of  $E$  (among the available  $\{H_k\}$ ). Then, it is — in *some cases* — *rationally required* that an agent’s credence function at  $t_0$  be such that

$$cr_{t_0}(H_j | E \ \& \ \mathbb{A}_j) > cr_{t_0}(H_j | E \ \& \ \neg \mathbb{A}_j).$$

Alternatively, let  $\mathbb{A}(H_i, H_j, E) \stackrel{\text{def}}{=} H_j$  is a *better explanation* of  $E$  than  $H_i$  is. And, then explicate **ERA** in terms of *favoring* [10].

Supposing  $E$  (at  $t_0$ ),  $\mathbb{A}(H_i, H_j, E)$  *favors*  $H_j$  over  $H_i$ .

I’m not the first one to propose something like **ERA**. Climenhaga, Hartmann *et. al.*, Lange, and Weisberg have all argued convincingly for similar principles [2, 18, 28, 4, 15, 9].

Roche & Sober [22] agree that **ERA** is the right *formulation* of Explanationism; but, they argue that **ERA** is *false* — *viz.*, that  $\mathbb{A}_j$  is *never* relevant to  $H_j$  (on the supposition of the evidence  $E$ ).

I think Climenhaga [2] and Lange [18] do a pretty good job of responding to Roche & Sober’s skeptical argument [22].

The main thing I would add to Climenhaga’s and Lange’s trenchant responses Roche & Sober is the following point.

☞ In order to refute R&S’s skepticism, all that is required is a *single example* in which  $\mathbb{A}_j$  is relevant to  $H_j$  (given  $E$ ).

I will close by discussing just such an example (involving Newton, Einstein, and the motion of Mercury), which is a well-known instance of The Problem of Old Evidence [6].

Let  $E \stackrel{\text{def}}{=} \text{what was (collectively) known about the precession of the perihelion of Mercury in } t_1 = 1914$ .<sup>3</sup>

Let  $H_1 \stackrel{\text{def}}{=} \text{Newton’s theory of planetary motion}$ , and  $H_2 \stackrel{\text{def}}{=} \text{Einstein’s theory of general relativity}$ . This yields the following 3-element partition of hypotheses:  $\{H_1, H_2, H_3 = \neg H_1 \ \& \ \neg H_2\}$ .

If **ERA** applies in this case, then our agent should be such that

$$cr_{t_0}(H_2 | E \ \& \ \mathbb{A}_2) > cr_{t_0}(H_2 | E \ \& \ \neg \mathbb{A}_2).$$

Thus, after  $E$  was learned (between  $t_0$  and  $t_1$ ), we should have

$$cr_{t_1}(H_2 | \mathbb{A}_2) > cr_{t_1}(H_2 | \neg \mathbb{A}_2).$$

When  $\mathbb{A}_2$  was subsequently learned in  $t_2 = 1915$  [26, 7], it — and *not  $E$* , which was *old evidence* — provided a boost to  $H_2$  [15, 9].

<sup>3</sup> $E$  had been known long before 1914. It is generally agreed that  $t_0$  is no later than 1882, when Newcomb calculated the Newtonian discrepancy in Mercury’s precession precisely. See [26] for a detailed history of this case.

How would EXPL<sub>1</sub> handle this Old Evidence Problem?

Let us suppose that our agent became certain of A<sub>2</sub> in t<sub>2</sub> = 1915.

Then, according to EXPL<sub>1</sub>, H<sub>2</sub> receives a “positive credit score” [c(H<sub>2</sub>, E) > 0] and a boost in credence (at t<sub>2</sub>) — over and above what Bayes’s Thm would recommend. Hence, EXPL<sub>1</sub> yields

$$cr_{t_2}(H_2 | E) > cr_{t_2}(H_2).$$

In this sense, according to EXPL<sub>1</sub>, E confirms H<sub>2</sub> in 1915.

☞ I see (at least) two problems with this application of EXPL<sub>1</sub>.

- It is *misleading* to say that E is what confirms H<sub>2</sub> in 1915. It is A<sub>2</sub> — by ensuring c(H<sub>2</sub>, E) > 0 — that is doing the work.
- In addition to the machinery of probability theory [modulo (2)], Douven also needs a “theory of credit scores.” In this sense, the (Bayesian) ERA approach is *more parsimonious*. [See Extras for a detailed formal example, which brings this out.]

The simplest formal example of Douven v. Bayes involves a 2-element partition {H<sub>1</sub> = H, H<sub>2</sub> = ¬H}, evidence E, and an abductive claim A asserting that H explains E better than ¬H.

Our Bayesian agent will begin with a prior cr<sub>t<sub>0</sub></sub>(·) over the algebra generated by the three atoms H, E, A. Between t<sub>0</sub> and t<sub>1</sub>, they will learn E, and between t<sub>1</sub> and t<sub>2</sub> they will learn A.

I will assume cr<sub>t<sub>0</sub></sub>(·) satisfies the following six (6) constraints.

- cr<sub>t<sub>0</sub></sub>(·) is regular.
- E (initially) confirms H [cr<sub>t<sub>0</sub></sub>(H | E) > cr<sub>t<sub>0</sub></sub>(H | ¬E)].
- ERA applies [cr<sub>t<sub>0</sub></sub>(H | E & A) > cr<sub>t<sub>0</sub></sub>(H | E & ¬A)].
- E is (a priori) irrelevant to A. [cr<sub>t<sub>0</sub></sub>(A | E) = cr<sub>t<sub>0</sub></sub>(A | ¬E)].
- A is (a priori) irrelevant to H. [cr<sub>t<sub>0</sub></sub>(H | A) = cr<sub>t<sub>0</sub></sub>(H | ¬A)].
- cr<sub>t<sub>0</sub></sub>(H) = cr<sub>t<sub>0</sub></sub>(E) = cr<sub>t<sub>0</sub></sub>(A) = 1/2.

There are many priors that satisfy (i)–(vi). I will choose a relatively simple one for illustrative purposes.

A	E	H	cr <sub>t<sub>0</sub></sub> (·)	cr <sub>t<sub>1</sub></sub> (·) = cr <sub>t<sub>0</sub></sub> (·   E)	cr <sub>t<sub>2</sub></sub> (·) = cr <sub>t<sub>0</sub></sub> (·   E & A)
T	T	T	7/32	7/16	7/8
T	T	F	1/32	1/16	1/8
T	F	T	1/32	0	0
T	F	F	7/32	0	0
F	T	T	3/32	3/16	0
F	T	F	5/32	5/16	0
F	F	T	5/32	0	0
F	F	F	3/32	0	0

Our Bayesian agent is such that cr<sub>t<sub>2</sub></sub>(H) = cr<sub>t<sub>0</sub></sub>(H | E & A) = 7/8.

Douven’s rule for computing cr<sub>t<sub>2</sub></sub>(H | E) will yield the Bayesian answer of 7/8 here *only if* H’s credit score is c(H, E) = 1 ∉ [0, 1).

In fact, we’ll need to allow c(H, E) ∈ [0, ∞). See my PrSAT [11] notebook fitelson.org/douven.pdf (.nb) for technical details.

☞ I would argue that this is how the values of c(H, E) should be “reverse engineered.” But, then, why *not* just go *fully Bayesian*?

- R. Carnap, *Logical Foundations of Probability*, 1950.
- N. Climenhaga, *How explanation guides confirmation*, 2017.
- B. de Finetti, *Foresight: Its Logical Laws, Its Subjective Sources*, 1937.
- R. Dawid, S. Hartmann, and J. Sprenger, *The No Alternatives Argument*, 2015.
- I. Douven, *The Art of Abduction*, 2022.
- E. Eells, *Bayesian Problems of Old Evidence*, 1990.
- A. Einstein, *Explanation of the Perihelion Motion of Mercury from General Relativity Theory*, 1915.
- B. Eva, *Principles of Indifference*, 2019.
- B. Eva and S. Hartmann, *On the Origins of Old Evidence*, 2020.
- B. Fitelson, *Contrastive Bayesianism*, 2012.
- B. Fitelson, *A Decision Procedure for Probability Calculus with Applications*, 2008.
- J.D. Gallow, *Learning and value change*, 2019.
- J.D. Gallow, *Updating for externalists*, 2021.
- A. Gibbard, *Rational Credence and the Value of Truth*, 2006.
- S. Hartmann and B. Fitelson, *A new Garber-style solution to the problem of old evidence*, 2015.
- J. Joyce, *Prospects for an Alethic Epistemology of Partial Belief*, 2009.
- J.M. Keynes, *A Treatise on Probability*, 1921.
- M. Lange, *Putting explanation back into “inference to the best explanation,”* 2020.
- D. Lewis, *A Subjectivist’s Guide to Objective Chance*, 1980.
- R. Pettigrew, *Accuracy and the Laws of Credence*, 2016.
- R. Pettigrew, *The Dutch Book Arguments*, 2020.
- W. Roche and E. Sober, *Explanatoriness is evidentially irrelevant...*, 2013.
- J. Thomson, *Normativity*, 2008.
- M. Titelbaum, *Fundamentals of Bayesian Epistemology*, 2022.
- M. Titelbaum, *Quitting Certainties*, 2013.
- K.J. Treschman, *Early Astronomical Tests of General Relativity...*, 2014.
- B. van Fraassen, *Laws and symmetry*, 1989.
- J. Weisberg, *Locating IBE in the Bayesian Framework*, 2009.
- B. Williams, *Internal and External Reasons*, 1981.
- T. Williamson, *Knowledge and its Limits*, 2000.
- A. Worsnip, *Fitting Things Together: Coherence and the Demands of Structural Rationality*, 2021.
- J. Zhao, V. Crupi, K. Tentori, B. Fitelson, and D. Osherson, *Updating: Learning vs. supposing*, 2012.