

# Predictable polarization

Kevin Dorst

24.223 Rationality

## I. Confirmation bias as reflection failures

We saw from Kelly:

- The choice to scrutinize selectively can certainly be Bayesian.
- So long as *failing* to find a flaw *lowers* your credence, the update can be perfectly rational.

Why think this is a form of bias? Because often we can *predict* (or *have expectations for*) how it will shift our beliefs.

What confirmation bias is not:

- Not about being *likely* to raise your credence.
- Not about someone *who knows more than you* being able to predict how your beliefs will shift.

*Proposal:* Your inquiry exhibits **confirmation bias** toward  $q$  iff your expectation of your updated credence in  $q$  is higher than your prior:

It's definitely possible to search for evidence for  $q$  without exhibiting confirmation bias toward  $q$ .

Example: word searches.

Possibilities =  $(n, c, f)$ : no word, completable but don't find, and find.

$$P = \begin{pmatrix} n & c & f \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}, \text{ while } P^+ = \begin{pmatrix} 2/3 & 1/3 & 0 \\ 2/3 & 1/3 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Salow: selective scrutiny is rational only if it satisfies Reflection.

Dorst: the problem of polarization is that we violate Reflection.

*Examples:* Googling symptoms; biased or one-sided sources; Pascal's Wager; lawyer's argument; going to college.

## II. Is confirmation bias irrational?

An update  $(P, P^+)$  satisfies the **value of evidence** iff, for all decisions<sup>1</sup>,  $P$  expects the option that  $P^+$  recommends<sup>2</sup> to be better than the option  $P$  recommends.

→  $P$  wants to give power of attorney to  $P^+$ .

*Fact:* there are updates that satisfy the value of evidence that exhibit confirmation bias. *Example:*

$$P = \left(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right), \text{ while } P^+ = \begin{pmatrix} 2/3 & 1/3 & 0 \\ 1/3 & 2/3 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Lottery with 10 tickets; credence you lost?

Me knowing you'll raise your credence that my Dad's birthday is in June.

$\mathbb{E}_P(P^+(q)) > P(q)$ . I.e. violate the expectation-version of Reflection.

$P$  is stationary wrt  $P^+$ :  $\mathbb{E}_P(P^+) = P$ .

$\mathbb{E}_P(P^+(q)) = P(q)$ .

<sup>1</sup> Set of options  $\{X_1, \dots, X_n\}$ , fns from worlds  $w$  to utilities  $X_i(w) \in \mathbb{R}$

<sup>2</sup> The  $X_i$  that maximizes  $\mathbb{E}_{P^+}(X_i)$

⇒  $P$  expects  $P^+$  to be more accurate than  $P$

$\mathbb{E}_P(P^+) = \left(\frac{5}{12}, \frac{4}{12}, \frac{3}{12}\right)$ , so  
 $\mathbb{E}_P(P^+(\text{word})) = \frac{7}{12} \approx 0.58 > 0.5 = P(\text{word})$

How? This is possible because  $P^+$  has *higher-order uncertainty*: it is unsure of its own values.

→ an *ambiguity-asymmetry*.

**Claim 1:** Since this update satisfies the value of evidence, it can be rational despite exhibiting confirmation bias.

**Claim 2:** Repeating this (rational) process can lead to predictable, profound polarization.

The rational response to *not-finding* is to be unsure what the rational response is.

## References

- Henderson, Leah and Gebharter, Alexander, 2021. 'The role of source reliability in belief polarisation'. *Synthese*, 1–23.
- Jern, Alan, Chang, Kai Min K., and Kemp, Charles, 2014. 'Belief polarization is not always irrational'. *Psychological Review*, 121(2):206–224.
- Nickerson, Raymond S., 1998. 'Confirmation bias: A ubiquitous phenomenon in many guises.' *Review of General Psychology*, 2(2):175–220.
- Salow, Bernhard, 2018. 'The Externalist's Guide to Fishing for Compliments'. *Mind*, 127(507):691–728.