

## Problem Statement

One big open problem in belief revision, dynamic epistemic logic, and machine learning is to develop the “model theory of learning,” i.e.

**Question 1.** What axioms capture important properties of learning (e.g. induction, no forgetting, no hallucinations, robustness)? What do models of these properties look like? Are there “correspondences” between axioms and properties of models, as there often are in modal logic?

A lot of work has already been done on this — see the two very relevant papers [1, 2]. In particular, [1] shows that, given a stream of true data, announcement and lexicographic upgrade converge to the full truth (whereas conservative upgrade does not). And if the stream has finitely many false data, only lexicographic upgrade converges to the full truth (of the three).

I was reading the second volume of [3] recently. In it, George Pólya makes some empirical observations about the role of guessing and induction in mathematics, and then tries to formalize these as logical rules. I got interested in this because I noticed that his rules of induction capture axioms for *plausibility upgrade* / *belief revision* that I believe have been overlooked by work on DEL and belief revision.

For the rest of this note, I want to formalize Pólya’s inductive rules as DEL axioms and explore the models that satisfy them. This is a bit of a detour from the question above, but I’ve been having fun thinking about it (and I hope you do too)!

## Pólya’s Stepwise Induction

Pólya charts about 16 different rules in his book, but for now I’m interested in the two that relate to induction. First, let’s look at his Stepwise Induction rule

$$\frac{P \rightarrow Q \quad Q \text{ true}}{P \text{ more plausible}}$$

Pólya gives this rule the following reading: If we discover that some consequence  $Q$  of  $P$  is true, then  $P$  becomes more plausible. He refers to this as just “induction”; I prefer to call it *stepwise* induction because a modern understanding of induction is *convergence to the truth*, and this rule says what ought to happen in a single discovery step rather than at some point of convergence.

Pólya wrote this 30 years before logics for plausibility upgrade were developed. But with the hindsight we have today, it’s clear that his rule expresses a DEL constraint for plausibility upgrade.

Suppose we have a plausibility model  $\mathcal{M} = \langle W, R, \leq, V \rangle$ . Let  $[P]$  be the action “make  $P$  more plausible,” which updates  $\mathcal{M}$  to  $\mathcal{M}_P^*$ . Let  $[\!P]$  be the action “discover that  $P$  is true, which updates  $\mathcal{M}$  to  $\mathcal{M}_P^!$  (we could do this by conditionalization / public announcement).

To make things concrete, let’s consider the full language

$$p \mid \neg\varphi \mid \varphi \wedge \psi \mid \Diamond\varphi \mid \mathbf{K}\varphi \mid \mathbf{B}^p\psi \mid [\!P]\varphi \mid [P]\varphi$$

with the semantics

$\mathcal{M}, w \Vdash p$	iff	$w \in V(p)$
$\mathcal{M}, w \Vdash \neg\varphi$	iff	$\mathcal{M}, w \not\Vdash \varphi$
$\mathcal{M}, w \Vdash \varphi \wedge \psi$	iff	$\mathcal{M}, w \Vdash \varphi$ and $\mathcal{M}, w \Vdash \psi$
$\mathcal{M}, w \Vdash \Diamond\varphi$	iff	there exists $u$ with $\mathcal{M}, u \Vdash \varphi$
$\mathcal{M}, w \Vdash \mathbf{K}\varphi$	iff	for all $u$ with $wRu$ , $\mathcal{M}, u \Vdash \varphi$
$\mathcal{M}, w \Vdash \mathbf{B}^\varphi\psi$	iff	$\text{best}_\leq(\llbracket\varphi\rrbracket) \subseteq \llbracket\psi\rrbracket$
$\mathcal{M}, w \Vdash [!P]\varphi$	iff	$\mathcal{M}_P^!, w \models \varphi$
$\mathcal{M}, w \Vdash [P]\varphi$	iff	$\mathcal{M}_P^*, w \models \varphi$

where  $\llbracket\varphi\rrbracket = \{u \mid \mathcal{M}, u \Vdash \varphi\}$ , and  $\text{best}_\leq(S) = \{u \in S \mid u \text{ is } \leq\text{-minimal over } S\}$ .

We can't express the rule in DEL as-stated, since  $[P]$  and  $[!P]$  can't occur by themselves. But we can rephrase the rule in terms of the *effects* that the updates have: The rule says that any effect of making  $P$  more plausible is also an effect of discovering that  $Q$  is true. More formally, if  $\mathcal{M}, w \Vdash P \rightarrow Q$  then

$$\mathcal{M}_P^*, w \Vdash \varphi \text{ implies } \mathcal{M}_Q^!, w \Vdash \varphi$$

And this holds iff  $\mathcal{M}, w$  satisfies

$$(P \rightarrow Q) \rightarrow ([P]\varphi \rightarrow [!Q]\varphi) \quad (\text{STEP-IND})$$

The second rule I'd like to explore is Shaded Stepwise Induction ("shaded" is Pólya's terminology)

$$\frac{P \rightarrow Q \quad Q \text{ more plausible}}{P \text{ more plausible}}$$

We can do a similar analysis on this, and express it in DEL as

$$(P \rightarrow Q) \rightarrow ([P]\varphi \rightarrow [Q]\varphi) \quad (\text{SHADED-STEP-IND})$$

I want to ask the same questions from before, but now zoomed-in at these two axioms. Specifically,

**Question 2.** What do satisfying models of (STEP-IND) and (SHADED-STEP-IND) look like? I'll be happy to find just one for each, but I'd like to completely axiomatize them if possible.

**Question 3.** What is the relationship between (STEP-IND) and (SHADED-STEP-IND)? Is one 'stronger' than the other, and in what sense?

**Question 4.** What is the relationship between *stepwise* induction and induction understood as convergence to the truth? Does stepwise induction imply regular induction (as the name suggests)? Or is this notion of stepwise induction irrelevant to regular induction?

## Progress So Far

### A Model for Shaded Stepwise Induction

**Definition 1.** Given  $\mathcal{M} = \langle W, R, \leq, V \rangle$ , let the update  $\mathcal{M}_P^* = \langle W, R, \leq_P^*, V \rangle$ , where  $\leq_P^*$  is reordered according to the rule:

Put all states where *some* (non- $\top$ ) consequence of  $P$  holds in front of all the other states; the order within the two groups is kept the same.

Note that this is a variation on Lexicographic upgrade.

**Lemma 1.** For all formulas  $P, Q$  and all sets  $S$ ,  $\text{best}_{\leq_Q^*}(S) \subseteq \text{best}_{\leq_P^*}(S)$ .

*Proof.* We have two cases:

- No state  $w \in S$  satisfies a consequence of  $Q$ . Note that this also implies that no state satisfies a consequence of  $P$ , since TODO

□

**Theorem 2.** (SHADED-STEP-IND) is sound for this update.

*Proof.* By induction on  $\varphi$ . The interesting case is conditional belief  $\mathbf{B}^\varphi\psi$ , since it is the only case affected by the upgrade. Suppose  $\mathcal{M}, w \Vdash P \rightarrow Q$ , and  $\mathcal{M}, w \Vdash [P]\mathbf{B}^\varphi\psi$ . So  $\text{best}_{\leq_P^*}(\llbracket\varphi\rrbracket) \subseteq \llbracket\psi\rrbracket$ . By Lemma 1 we have

$$\text{best}_{\leq_Q^*}(\llbracket\varphi\rrbracket) \subseteq \text{best}_{\leq_P^*}(\llbracket\varphi\rrbracket) \subseteq \llbracket\psi\rrbracket$$

And so  $\mathcal{M}, w \Vdash [Q]\mathbf{B}^\varphi\psi$ .

□

So we have an update that models (SHADED-STEP-IND). My goal will now be to completely axiomatize this update.

**Proposition 3.** We can completely characterize the effect this update has on best. For all  $S \subseteq W$ ,

$$\text{best}_{\leq_P^*}(S) = \begin{cases} \text{best}_{\leq}(S \cap \{w \mid \exists X \neq \top \text{ such that } \mathcal{M}, w \Vdash P \rightarrow X \text{ and } \mathcal{M}, w \Vdash X\}) & \text{if } S \cap \{w \mid \exists X \neq \top \text{ such that } \mathcal{M}, w \Vdash P \rightarrow X \text{ and } \mathcal{M}, w \Vdash X\} \neq \emptyset \\ \text{best}_{\leq}(S) & \text{otherwise} \end{cases}$$

*Proof.* Notice that this expression is exactly the characterization of lexicographic upgrade

$$\text{best}_{\leq_P^\uparrow}(S) = \begin{cases} \text{best}_{\leq}(S \cap \llbracket P \rrbracket) & \text{if } S \cap \llbracket P \rrbracket \neq \emptyset \\ \text{best}_{\leq}(S) & \text{otherwise} \end{cases}$$

except instead of upgrading the set

$$\llbracket P \rrbracket = \{w \mid \mathcal{M}, w \Vdash P\}$$

we upgrade

$$\{w \mid \exists X \neq \top \text{ such that } \mathcal{M}, w \Vdash P \rightarrow X \text{ and } \mathcal{M}, w \Vdash X\}$$

From here, the proof is identical to the proof for lexicographic upgrade.

□

**Discussion.** Consider the expression

$$\exists X \neq \top \text{ such that } \mathcal{M}, w \Vdash P \rightarrow X \text{ and } \mathcal{M}, w \Vdash X$$

If we could express this in DEL, we would have a complete axiomatization for the upgrade (just modify the reduction axioms for lexicographic upgrade). This has an undeniably modal flavor to it; think of  $P \rightarrow X$  as an accessibility relation over formulas. It doesn't seem to be expressible in our language so far, so let's go ahead and define it.

**Definition 2.** Let  $\odot$  be a new operator in our language, with semantics given by

$$\mathcal{M}, w \Vdash \odot P \text{ iff } \exists \text{ formula } X \neq \top \text{ such that } \mathcal{M}, w \Vdash P \rightarrow X \text{ and } \mathcal{M}, w \Vdash X$$

We read  $\odot P$  as “some consequence of  $P$  holds.”  $\odot$  is an existential modality; we define the dual  $[\odot]P \leftrightarrow \neg \odot \neg P$  (“all consequences of  $P$  hold”). (The  $\odot$  symbol evokes an image of the winding paths we may take in finding consequences of  $P$ .)

**Proposition 4.** The following are sound:

1. If  $\vdash P$  then  $\vdash [\odot]P$  (NECESS)
2.  $(P \rightarrow Q) \rightarrow (\odot Q \rightarrow \odot P)$  (ANTITONE)
3.  $P \rightarrow \odot P$  (REFL)

*Proof.* (1) is easy: If  $P$  holds in all  $\mathcal{M}, w$ , then so do all of its (nontrivial) consequences. (3) is also easy, since  $P$  is a consequence of itself.

For (2), suppose  $\mathcal{M}, w \Vdash P \rightarrow Q$  and  $\mathcal{M}, w \Vdash \odot Q$ . Then  $\mathcal{M}, w \Vdash X$  for some  $X \neq \top$  such that  $\mathcal{M}, w \Vdash Q \rightarrow X$ . But then by transitivity of consequence ( $X$  is a consequence of  $Q$ , which is a consequence of  $P$ ),  $\mathcal{M}, w \Vdash P \rightarrow X$ , which gives us  $\mathcal{M}, w \Vdash \odot P$ .  $\square$

**Corollary 5.** Here are some interesting consequences of (ANTITONE) and (REFL). Notice that they are flipped from their usual  $\square$  presentation.

1.  $\odot P \wedge \odot Q \rightarrow \odot(P \wedge Q)$  (M)
2.  $\odot(P \vee Q) \rightarrow \odot P \vee \odot Q$  (OR)
3.  $\odot \odot P \leftrightarrow \odot P$  (TRANS)

*Proof.* For (1), since  $\vdash P \wedge Q \rightarrow P$ , by (ANTITONE) we have  $\vdash \odot P \rightarrow \odot(P \wedge Q)$ . But then we have  $\vdash \odot P \wedge \odot Q \rightarrow \odot(P \wedge Q)$ . (2) is similar: Since  $\vdash P \rightarrow P \vee Q$ , by (ANTITONE) we have  $\vdash \odot(P \vee Q) \rightarrow \odot P$ . Consequently,  $\vdash \odot(P \vee Q) \rightarrow \odot P \vee \odot Q$ .

Now consider (3). The backwards direction is just (REFL). As for the other direction, (REFL) gives us  $\vdash P \rightarrow \odot P$ . Applying (ANTITONE) we get  $\vdash \odot \odot P \rightarrow \odot P$ .  $\square$

Importantly,  $\odot$  is *not* a normal modality — the converse of (M),  $\odot(P \wedge Q) \rightarrow (\odot P \wedge \odot Q)$ , is not sound. Surprisingly,  $\odot$  is completely axiomatized by (NECESS), (ANTITONE), and (REFL).

**Lemma 6.** (Existence Lemma). Suppose  $\Delta$  is a maximally consistent set with  $\odot \varphi \in \Delta$ . Then there is a formula  $X \neq \top$  of complexity less than or equal to  $\varphi$  with  $\varphi \rightarrow X \in \Delta$  and  $X \in \Delta$ .

*Proof.* Suppose for contradiction that there is no such  $X$ , i.e. for all  $X$  of complexity less than or equal to  $\varphi$ , if  $\varphi \rightarrow X \in \Delta$  then  $X \notin \Delta$ . Let's now construct the formula

$$\psi = \bigvee \{X \neq \top \mid X \text{ has complexity } \leq \varphi \text{ and } \varphi \rightarrow \odot X \in \Delta\}$$

Intuitively, this formula says “some consequence of some consequence of  $\varphi$  holds.” Is it well-defined? Well, by (REFL)  $\varphi \rightarrow \odot\varphi \in \Delta$ , so the set is nonempty. Moreover, since the complexity of these formulas is bounded there are only finitely many of them. So this formula is well-defined.

First, since  $\odot\varphi \in \Delta$  and  $\Delta$  is maximally consistent,  $\psi \rightarrow \odot\varphi \in \Delta$  (by  $\rightarrow$ -introduction: if we suppose  $\psi \in \Delta$  then trivially  $\odot\varphi \in \Delta$ ).

Second, I claim that  $\neg\odot\psi \in \Delta$ . Why? Let  $\psi = X_1 \vee \dots \vee X_k$ . For each  $X_i$  occurring in  $\psi$ ,  $(\varphi \rightarrow \odot X_i) \in \Delta$ . By our very first assumption, every  $\odot X_i \notin \Delta$ . So  $\odot X_1 \vee \dots \vee \odot X_k \notin \Delta$ . But then by (OR),

$$\odot\psi = \odot(X_1 \vee \dots \vee X_k) \notin \Delta$$

And since  $\Delta$  is maximally consistent, we have  $\neg\odot\psi \in \Delta$ .

Let's put everything together. Since  $\psi \rightarrow \odot\varphi \in \Delta$ , by (ANTITONE) we have  $\odot\odot\varphi \rightarrow \odot\psi \in \Delta$ . (REFL) gives us  $\odot\varphi \rightarrow \odot\odot\varphi$ , and then by modus ponens we have  $\odot\varphi \rightarrow \odot\psi \in \Delta$ . But since  $\neg\odot\psi \in \Delta$ , by contraposition  $\neg\odot\varphi \in \Delta$ . But  $\odot\varphi \in \Delta$ , so this contradicts the fact that  $\Delta$  is consistent.  $\square$

**Proposition 7.** Suppose we build the canonical model  $\mathcal{M}^c$  using the usual maximally-consistent set construction for the base logic of  $\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi$ . The  $\odot$ -case of the truth lemma holds, i.e. for all maximally consistent  $\Delta$ ,

$$\mathcal{M}^c, \Delta \Vdash \odot\varphi \quad \text{iff} \quad \odot\varphi \in \Delta$$

*Proof.* ( $\rightarrow$ ) Suppose  $\mathcal{M}^c, \Delta \Vdash \odot\varphi$ . So  $\mathcal{M}^c, \Delta \Vdash X$  for some  $X \neq \top$  such that  $\mathcal{M}^c, \Delta \Vdash \varphi \rightarrow X$ . Now apply the inductive hypothesis to both  $\varphi \rightarrow X$  and  $X$  to get  $\varphi \rightarrow X \in \Delta$  and  $X \in \Delta$ . From  $X \in \Delta$  and (REFL) we have  $\odot X \in \Delta$ . And from  $\varphi \rightarrow X \in \Delta$ ,  $\odot X \in \Delta$ , and (ANTITONE), we have  $\odot\varphi \in \Delta$ .

( $\leftarrow$ ) Now suppose  $\odot\varphi \in \Delta$ . By Lemma 6 there is a formula  $X \neq \top$  of complexity less than or equal to  $\varphi$  with  $\varphi \rightarrow X \in \Delta$  and  $X \in \Delta$ . Since its complexity is less than or equal to  $\varphi$ , we can apply the inductive hypothesis to  $\varphi \rightarrow X$  and  $X$ . So  $\mathcal{M}^c, \Delta \Vdash \varphi \rightarrow X$  and  $\mathcal{M}^c, \Delta \Vdash X$ . By the semantics of  $\odot$ ,  $\mathcal{M}^c, \Delta \Vdash \odot\varphi$ .  $\square$

**Theorem 8.** Assuming model building for the logic  $\{\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi\}$ , given any consistent  $\Gamma$  over  $\{\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi, \odot\varphi\}$  we can build a model and point  $\mathcal{M}, w \models \Gamma$ .

*Proof.* This proof is totally standard:

First, extend  $\Gamma$  to maximally consistent  $\Gamma \supseteq \Gamma$  in the usual way, and let  $\mathcal{M}^c$  be the canonical model for the logic  $\{\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi\}$  (no need to modify it for  $\odot$ ). By the Truth Lemma, for all  $\varphi \in \Delta$ ,  $\mathcal{M}^c, \Delta \Vdash \varphi$ . So in particular, for all  $\varphi \in \Gamma$ ,  $\mathcal{M}^c, \Delta \Vdash \varphi$ .  $\square$

**Corollary 9.** Assuming completeness for the logic  $\{\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi\}$ , the logic  $\{\Diamond\varphi, \mathbf{K}\varphi, \mathbf{B}^\varphi\psi, \odot\varphi\}$  is also complete.

*Proof.* This proof is also standard:

Suppose contrapositively that  $\Gamma \not\models \varphi$ . Then  $\Gamma \cup \{\neg\varphi\}$  is consistent. We can then build a model  $\mathcal{M}$  and point  $w$  such that  $\mathcal{M}, w \models \Gamma \cup \{\neg\varphi\}$ . But then  $\mathcal{M}, w \models \Gamma$  yet  $\mathcal{M}, w \not\models \varphi$ , which is what we wanted to show.  $\square$

**Theorem 10.** The upgrade defined above is completely axiomatized by the reduction axioms (over DEL +  $\odot$ ):

$$\begin{array}{ll}
[P]p & \leftrightarrow p \\
[P]\neg\varphi & \leftrightarrow \neg[P]\varphi \\
[P](\varphi \wedge \psi) & \leftrightarrow [P]\varphi \wedge [P]\psi \\
[P]\Diamond\varphi & \leftrightarrow \Diamond[P]\varphi \\
[P]\odot\varphi & \leftrightarrow \odot[P]\varphi \\
[P]\mathbf{K}\varphi & \leftrightarrow \mathbf{K}[P]\varphi \\
[P]\mathbf{B}^\varphi\psi & \leftrightarrow (\Diamond(\odot P \wedge [P]\varphi) \wedge \mathbf{B}^{\odot P \wedge [P]\varphi}[P]\psi) \\
& \vee (\neg\Diamond(\odot P \wedge [P]\varphi) \wedge \mathbf{B}^{[P]\varphi}[P]\psi)
\end{array}$$

*Proof.* TODO  $\square$

## Things To Try

1. go through mathreviews and google scholar to make sure this hasn't been done yet
2. Prove that  $\odot$  is not expressible in the base language (it is a genuine language extension).
3. Find a model satisfying (STEP-IND) (do a similar analysis of this upgrade).
4. Think about Question 3: What is the relationship between (STEP-IND) and (SHADED-STEP-IND)?
5. Think about Question 4: Do either of these axioms guarantee convergence to the truth, under reasonable assumptions?

## References

- [1] Alexandru Baltag, Nina Gierasimczuk, and Sonja Smets. "Truth-tracking by belief revision". In: *Studia Logica* 107 (2019), pp. 917–947.
- [2] Alexandru Baltag et al. "A dynamic logic for learning theory". In: *Journal of Logical and Algebraic Methods in Programming* 109 (2019), p. 100485.
- [3] George Polya. *Mathematics and plausible reasoning: Induction and analogy in mathematics*. Vol. 2. Princeton University Press, 1954.