
THE LOGIC OF HEBBIAN LEARNING

Caleb Kisby¹, Saúl A. Blanco¹, Lawrence S. Moss²

¹Department of Computer Science, Indiana University

²Department of Mathematics, Indiana University

Bloomington, IN 47408, USA

{cckisby, sbblancor, lmoss}@indiana.edu

Note: This arXiv print is a fuller version of our FLAIRS paper [Kisby et al., 2022], complete with proofs we could not include in the conference format. We have also corrected a couple of errors, as well as certain aesthetic choices that made the proof system somewhat awkward. Otherwise, we have kept as close to the original as possible; for a more up-to-date version of the system presented here, see our follow-up work [Kisby et al., 2024].

ABSTRACT

We present the logic of Hebbian learning, a dynamic logic whose semantics¹ are expressed in terms of a layered neural network learning via Hebb’s associative learning rule. Its language consists of modality $T\varphi$ (read “typically φ ,” formalized as forward propagation), conditionals $\varphi \Rightarrow \psi$ (read “typically φ are ψ ”), as well as dynamic modalities $[\varphi^+]\psi$ (read “evaluate ψ after performing Hebbian update on φ ”). We give axioms and inference rules that are sound with respect to the neural semantics; these axioms characterize Hebbian learning and its interaction with propagation. The upshot is that this logic describes a neuro-symbolic agent that both learns from experience and also reasons about what it has learned.

Introduction

Artificial intelligence has long been marked by a schism between two of its major paradigms: symbolic reasoning and connectionist learning. Neural systems have had wild success with learning from unstructured data, whereas symbolic reasoners are notorious for their rigidity. On the other hand, symbolic systems excel at sophisticated (static) reasoning tasks that neural systems cannot readily learn. Symbolic systems also tend to have more explainable reasoning, thanks to their use of explicit inferences in an intuitive language. Moreover, due to their connection with logic, it is straightforward to compare the relative power and complexity of different symbolic reasoners.

But as Valiant famously put it, intelligent cognitive agents must have *both* “the ability to learn from experience, and the ability to reason from what has been learned” [Valiant, 2003]. *Neuro-symbolic artificial intelligence* has emerged in the last few decades to address this challenge — a monumental effort to integrate neural and symbolic systems, while retaining the advantages of both (see [Bader and Hitzler, 2005] and [Sarker et al., 2021], two surveys that span the decades). Despite the cornucopia of neuro-symbolic proposals, the field has not yet agreed on an interface between the two that satisfyingly preserves both flexible learning and expressive reasoning.

Following the path set out by [Balkenius and Gärdenfors, 1991] and [Leitgeb, 2001, 2003], we advance the following proposal for the neuro-symbolic interface. Rather than viewing the neural and symbolic as two different systems to be combined, we view them as two ways of interpreting the same agent. More precisely, we view the dynamics of neural networks as the semantics to a formal logic. This logic serves as a bridge between the neural network model and formal inference.

¹A Python implementation of our semantics, using Tensorflow & Keras [Abadi et al., 2015]), is available at

<https://github.com/ais-climber/neural-semantics>

Previous work, particularly [Leitgeb, 2001], has considered how forward propagation in binary feed-forward nets forms a sound and complete semantics for the (static) conditional logic **CL** (*loop-cumulative*). The novelty of our paper is that we extend this logic by viewing a simple learning policy — Hebbian update (“neurons that fire together wire together”) — as a dynamic modality. By doing so, we demonstrate that the dynamics of Hebbian learning (in binary feed-forward nets) directly corresponds to a particular dynamic multimodal logic that we call *the logic of Hebbian learning*. This logic meets Valiant’s challenge: It characterizes a cognitive agent that can learn from experience and also reason about what it has learned.

Our main result is the soundness of axioms and inference rules that characterize Hebbian learning. The most interesting axioms involve the interaction between Hebbian update and forward propagation. We also demonstrate how our logic models the learning of a concrete neural network. And although we leave the question of completeness open, we close by considering the importance of completeness for logics of this kind.

Related Work

Logics with Neural Semantics.

The idea that we can view neural networks as the semantics for symbolic reasoning dates back to [McCulloch and Pitts, 1943]. Our work builds on a recent reimagining of this à la [Balkenius and Gärdenfors, 1991], [Leitgeb, 2001, 2003, 2018], which formally characterize the dynamics of inhibitory neural networks as conditional logics. Similarly, [Blutner, 2004] demonstrates that Hopfield networks correspond to the logic of what he calls “weight-annotated Poole systems.” More recently, [Giordano et al., 2021] describe multilayer perceptrons and self-organizing maps in terms of defeasible description logics. Yet no neural semantics to date has tackled the issue of learning — doing this for Hebbian learning is precisely the contribution of our paper.

Neuro-Symbolic AI.

Across the neuro-symbolic literature, an ubiquitous premise is that integration involves combining or composing otherwise distinct neural and symbolic modules. In contrast, this paper presents the neural and symbolic as two perspectives we can have about the same agent.

To our knowledge, the combined work of [Garcez et al., 2001] and [Garcez et al., 2008] is the only neuro-symbolic proposal (besides neural semantics, see above) that exhibits this intimate interface between the two. The former gives a formally sound method for extracting conditionals from a network and the latter gives a method for build neural network models from rules (in a variety of different logics). When combined, we can freely translate between a neural network and its beliefs. But unlike our work, this framework does not offer a logical account of the neural network’s learning.

Dynamic Logics for Learning.

Two recent papers, [Baltag et al., 2019a] and [Baltag et al., 2019b], also present dynamic multimodal logics that characterize learning. The former models an individual’s learning in the limit, whereas the latter models supervised learning as a game played between student and teacher. But it is unclear how learning policies expressed in these logics might relate to specific neural implementations of learning such as Hebbian update and backpropagation.

Furthermore, the syntax and inferences of our logic do not resemble either of these in a meaningful way. Perhaps the closest logics to ours are dynamic logics of *preference upgrade*, in the sense of [Van Benthem and Liu, 2007]. In particular, consider the modalities $[\uparrow\varphi]$ (lexicographic upgrade) and $[\uparrow\varphi]$ (elite change) [Van Benthem, 2007]. Both of these operators implement policies for modifying an agent’s preference relation $<$ over possible worlds. As with our logic, the key axioms characterizing these policies deal with their interaction with conditionals $\varphi \Rightarrow \psi$. But the semantics of our logic are very different; we leave the issue of how our neural semantics relate to classical preference relations to future work. In addition, both $[\uparrow\varphi]$ and $[\uparrow\varphi]$ are reducible to the static language of conditionals, whereas it is presently unclear how our $[\varphi^+]$ might reduce to its base language.

Background

Neural Network Models

A model of the logic of Hebbian learning is just a special type of artificial neural network that we call a *binary feedforward neural network* (BFNN).

Definition 1. A BFNN is a pointed directed graph $\mathcal{N} = \langle N, E, W, A, O, \eta \rangle$, where

- N is a finite nonempty set (the set of neurons)
- $E \subseteq N \times N$ (the set of excitatory connections)
- $W : N \times N \rightarrow \mathbb{R}$ (the weight of a given connection)
- A is a function which maps each $n \in N$ to $A^{(n)} : \mathbb{R}^k \rightarrow \mathbb{R}$ (the activation function for n , where k is the indegree of n)
- O is a function which maps each $n \in N$ to $O^{(n)} : \mathbb{R} \rightarrow \{0, 1\}$ (the output function for n)
- $\eta \in \mathbb{R}, \eta \geq 0$ (the learning rate)

As shorthand, we sometimes write W_{ij} to mean $W(i, j)$, for $(i, j) \in E$. Moreover, BFNNs are *feed-forward*, i.e. they do not contain cycles of edges with all nonzero weights. BFNNs are also *binary*, i.e. the output of each neuron is in $\{0, 1\}$. This binary assumption is unrealistic in practice, although letting it go is just a matter of extending our two-valued logic towards a fuzzy-valued logic (left to future work).

We further require that each composition of activation and output functions $O^{(n)} \circ A^{(n)}$ is *strictly* monotonically increasing, i.e. for all $\vec{x}, \vec{y} \in \mathbb{R}^k$ if $\vec{x} < \vec{y}$ then $O^{(n)}(A^{(n)}(\vec{x})) < O^{(n)}(A^{(n)}(\vec{y}))$. We will more often refer to the equivalent condition:

$$\vec{x} \leq \vec{y} \quad \text{iff} \quad O^{(n)}(A^{(n)}(\vec{x})) \leq O^{(n)}(A^{(n)}(\vec{y})) \quad (*)$$

Our activation functions include in particular those sigmoid functions commonly used for neural networks in practice.

The Dynamics of Propagation

Of course, BFNNs are not merely static directed graphs, but are dynamic in nature. When a BFNN receives a signal (which we model as the initial state), it propagates that signal forward until the state of the net stabilizes. This stable state of the net is considered to contain the net's response (answer) to the given signal (question). We model forward propagation as follows, drawing heavily from the approach proposed by [Leitgeb, 2001].

We consider a neuron n active if its activation $A^{(n)}$ is high enough to trigger an output $O^{(n)}$ of 1 (intuitively, if the neuron fires). Since our BFNNs are binary, either a given neuron is active (1) or it is not (0). So we can identify the state of \mathcal{N} with the set of neurons that are active. For a given BFNN \mathcal{N} , let its set of states be

$$\text{Set} = \{S \mid S \subseteq N\}$$

We can get the activation value of a particular neuron in a state S using the following characteristic function:

Definition 2. For $S \in \text{Set}$, let $\chi_S : N \rightarrow \{0, 1\}$ be given by $\chi_S(n) = 1$ iff $n \in S$

Neurons in a state $S \in \text{Set}$ can subsequently activate new neurons, which activate yet more neurons, until eventually the state of \mathcal{N} stabilizes. We call this final state of affairs $\text{Prop}(S)$, the *propagation* of S .

Definition 3. Let $\text{Prop} : \text{Set} \rightarrow \text{Set}$ be defined recursively as follows: $n \in \text{Prop}(S)$ iff either Let $\text{Prop} : \text{Set} \rightarrow \text{Set}$ be defined recursively as follows: $n \in \text{Prop}(S)$ iff either

(Base Case) $n \in S$, or

(Constructor) The weighted sum of predecessors of n subsequently activate it; for those m_1, \dots, m_k such that $m_i E n$ we have

$$O^{(n)}(A^{(n)}(\sum_{m_i E n} W(m_i, n) \cdot \chi_{\text{Prop}(S)}(m_i))) = 1$$

Alternatively, consider a finite automaton with state space Set and transition function $F_{S^*} : \text{Set} \rightarrow \text{Set}$ tracking the propagation of an initial state S^* through \mathcal{N} . We can view $\text{Prop}(S^*)$ as a fixed point of F_{S^*} [Leitgeb, 2001].

The following theorem, due to [Leitgeb, 2001], says that we can neatly characterize the algebraic structure of Prop as a closure operator. Note that Leitgeb proves this for *inhibition nets*, i.e. weightless BFNNs with both excitatory and inhibitory connections. But inhibition nets and our BFNNs are equivalent with respect to their propagation structure — we prove this result again for BFNNs as a kind of “sanity check” that our definitions are correct.

Theorem 1. Let $\mathcal{N} \in \text{Net}$. For all $S, S_1, S_2 \in \text{Set}$, Prop satisfies

(Inclusion) $S \subseteq \text{Prop}(S)$

(Idempotence) $\text{Prop}(S) = \text{Prop}(\text{Prop}(S))$

(Cumulative) If $S_1 \subseteq S_2 \subseteq \text{Prop}(S_1)$ then $\text{Prop}(S_1) = \text{Prop}(S_2)$

(Loop) If $S_1 \subseteq \text{Prop}(S_0), \dots, S_k \subseteq \text{Prop}(S_{k-1})$ and $S_0 \subseteq \text{Prop}(S_k)$, then $\text{Prop}(S_i) = \text{Prop}(S_j)$ for all $i, j \in \{0, \dots, k\}$

Proof. We prove each in turn:

(Inclusion) If $n \in S$, then $n \in \text{Prop}(S)$ by the base case of **Prop**.

(Idempotence) The (\subseteq) direction is just Inclusion. As for (\supseteq) , let $n \in \text{Prop}(\text{Prop}(S))$, and proceed by induction on $\text{Prop}(\text{Prop}(S))$.

Base Step: $n \in \text{Prop}(S)$, and so we are done.

Inductive Step: Let m_1, \dots, m_k be the predecessors of n . We have

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(\text{Prop}(S))}(m_i))) = 1$$

By inductive hypothesis, each $m_i \in \text{Prop}(\text{Prop}(S))$ iff $m_i \in \text{Prop}(S)$, and so

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(S)}(m_i))) = 1$$

By definition, this means that $n \in \text{Prop}(S)$.

(Cumulative) For the (\subseteq) direction, let $n \in \text{Prop}(S_1)$. We proceed by induction on $\text{Prop}(S_1)$.

Base Step: $n \in S_1$. Well, $S_1 \subseteq S_2 \subseteq \text{Prop}(S_2)$, so $n \in \text{Prop}(S_2)$.

Inductive Step: Let m_1, \dots, m_k be the predecessors of n . We have

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(S_1)}(m_i))) = 1$$

By inductive hypothesis, each $m_i \in \text{Prop}(S_1)$ iff $m_i \in \text{Prop}(S_1)$, and so

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(S_2)}(m_i))) = 1$$

By definition, this means that $n \in \text{Prop}(S_2)$.

Now consider the (\supseteq) direction. The Inductive Step holds similarly (just swap S_1 and S_2). As for the Base Step, if $n \in S_2$ then since $S_2 \subseteq \text{Prop}(S_1)$, $n \in S_1$.

(Loop) Let $k \geq 0$ and suppose the hypothesis. Our goal is to show that for each i , $\text{Prop}(S_i) \subseteq \text{Prop}(S_{i-1})$, and additionally $\text{Prop}(S_0) \subseteq \text{Prop}(S_k)$. This will show that all $\text{Prop}(S_i)$ contain each other, and so are equal. Let $i \in \{0, \dots, k\}$ (if $i = 0$ then $i - 1$ refers to k), and let $n \in \text{Prop}(S_i)$. We proceed by induction on $\text{Prop}(S_i)$.

Base Step: $n \in S_i$, and since $S_i \subseteq \text{Prop}(S_{i-1})$ by assumption, $n \in \text{Prop}(S_{i-1})$.

Inductive Step: Let m_1, \dots, m_k be the predecessors of n . We have

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(S_i)}(m_i))) = 1$$

By inductive hypothesis, each $m_i \in \text{Prop}(S_i)$ iff $m_i \in \text{Prop}(S_{i-1})$, and so

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W(m_i, n) \cdot \chi_{\text{Prop}(S_{i-1})}(m_i))) = 1$$

By definition, this means that $n \in \text{Prop}(S_{i-1})$. □

In the terminology of [Kraus et al., 1990], **Prop** is *loop-cumulative* — it satisfies both the cumulative and loop properties above. In fact, **Prop** is *not* a fully monotonic closure operator, as the following fact shows.

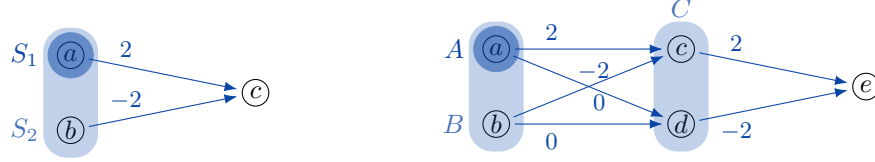


Figure 1: TODO

Proposition 1. It is not the case that for all $S_1, S_2 \in \text{Set}$, if $S_1 \subseteq S_2$, then $\text{Prop}(S_1) \subseteq \text{Prop}(S_2)$.

Proof. Consider the BFNN \mathcal{N} in Figure 1. Let $O^{(n)} \circ A^{(n)}$ be determined by a step function with threshold 0, i.e. $O^{(n)}(A^{(n)}(x)) = 1$ iff $x > 0$. We have $S_1 = \{a\} \subseteq \{a, b\} = S_2$, and so the hypothesis holds. But $\text{Prop}(S_1) = \{a, c\} \not\subseteq \{a, b\} = \text{Prop}(S_2)$. (Observe that c does not get activated in $\text{Prop}(S_2)$ because the weights cancel each other.) \square

From Hebbian Learning to Logic

The Dynamics of Hebbian Learning

The plan from here is to extend this logic of propagation by providing an account of Hebbian learning. Our goal is to cast Hebbian update as a dynamic modality, so that we can explore its interactions with **Prop** in symbolic language. As with **Prop**, we start by outlining the algebraic structure of Hebbian update.

Hebb’s classic learning rule [Hebb, 1949] states that when two adjacent neurons are simultaneously and persistently active, the connection between them strengthens. In contrast with, e.g. backpropagation, Hebbian learning is errorless and unsupervised. Another key difference is that Hebbian update is local — the change in a weight ΔW_{ij} depends only on the activation of the immediately adjacent neurons. For this reason, the Hebbian family of learning policies has traditionally been considered more biologically plausible than backpropagation. There are many variations of Hebbian learning, but we only consider the most basic (unstable, no weight decay) form of Hebb’s rule: $\Delta W_{ij} = \eta x_i x_j$, where η is the learning rate and x_i, x_j are the outputs of adjacent neurons i and j , respectively.

In order to incorporate Hebb’s rule into our framework, we introduce a function **Hebb** (“Hebbian update”) to strengthen those edges in a BFNN \mathcal{N} whose neurons are active when we feed \mathcal{N} a signal $S \in \text{Set}$.

Definition 4. Let $\text{Hebb} : \text{Net} \times \text{Set} \rightarrow \text{Net}$ be given by $\text{Hebb}(\langle N, E, W, A, O, \eta \rangle, S) = \langle N, E, W^*, A, O, \eta \rangle$, where

$$W_{ij}^* = W_{ij} + \eta \cdot \chi_{\text{Prop}(S)}(i) \cdot \chi_{\text{Prop}(S)}(j)$$

Notice that we propagate S before getting the active status of neurons. This is because otherwise we would never strengthen connections beyond the input layer. Now let’s consider the algebraic properties of **Hebb**. We were able to formulate the algebraic properties of **Prop** in terms of **Set** containment. Similarly, we can express certain properties of **Hebb** in terms of **Net** containment.

Definition 5. Let $\mathcal{N}_1, \mathcal{N}_2 \in \text{Net}$ differ only in their weights. We write

$$\mathcal{N}_1 \preceq \mathcal{N}_2$$

to mean that for all $S \in \text{Set}$, $\text{Prop}_{\mathcal{N}_1}(S) \subseteq \text{Prop}_{\mathcal{N}_2}(S)$. We use $\mathcal{N}_1 \cong \mathcal{N}_2$ to express that $\mathcal{N}_1 \preceq \mathcal{N}_2$ and $\mathcal{N}_2 \preceq \mathcal{N}_1$.

For example, consider the least upper bound \mathcal{N}^{lub} of \preceq . \mathcal{N}^{lub} is that net whose weights have been “maximally” strengthened — that is, increased to the point that every propagation $\text{Prop}(S)$ results in all neurons graph-reachable from S . By construction, \mathcal{N}^{lub} is a supernet of every other net.

We have the following test to determine if $\mathcal{N}_1 \preceq \mathcal{N}_2$.

Lemma 2. Suppose \mathcal{N}_1 and \mathcal{N}_2 are the same except for their weights, and let $S \in \text{Set}$. Then $\text{Prop}_{\mathcal{N}_1}(S) \subseteq \text{Prop}_{\mathcal{N}_2}(S)$ iff for all n and all m_1, \dots, m_k such that $(m_i, n) \in E$,

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_1}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_1}(S)}(m_i))) = 1$$

implies

(**)

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_2}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_2}(S)}(m_i))) = 1$$

In words: For all S, n , if n is activated by its predecessors in \mathcal{N}_1 , then it is activated by its predecessors in \mathcal{N}_2 as well.

Proof. (\rightarrow) Consider the contrapositive; suppose that there are m_1, \dots, m_k, n such that $(m_i, n) \in E$ with

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_1}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_1}(S)}(m_i))) = 1$$

and yet

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_2}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_2}(S)}(m_i))) = 0$$

But by definition of **Prop**, this means $n \in \text{Prop}_{\mathcal{N}_1}(S)$, but $n \notin \text{Prop}_{\mathcal{N}_2}(S)$. So $\text{Prop}_{\mathcal{N}_1}(S) \not\subseteq \text{Prop}_{\mathcal{N}_2}(S)$.

(\leftarrow) We need to show that $\text{Prop}_{\mathcal{N}_1}(S) \subseteq \text{Prop}_{\mathcal{N}_2}(S)$. Let $n \in \text{Prop}_{\mathcal{N}_1}(S)$. We show $n \in \text{Prop}_{\mathcal{N}_2}(S)$ by structural induction on $\text{Prop}_{\mathcal{N}_1}(S)$:

Base Step: $n \in S$. But then $n \in \text{Prop}_{\mathcal{N}_2}(S)$ by the base case of $\text{Prop}_{\mathcal{N}_2}(S)$.

Inductive Step: We have $n \in \text{Prop}_{\mathcal{N}_1}(S)$ from the constructor case, i.e. because for those m_1, \dots, m_k such that $(m_i, n) \in E$,

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_1}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_1}(S)}(m_i))) = 1$$

By assumption, for these n and predecessors m_1, \dots, m_k in particular,

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\mathcal{N}_2}(m_i, n) \cdot \chi_{\text{Prop}_{\mathcal{N}_2}(S)}(m_i))) = 1$$

Notice that our hypothesis is so strong that there is no need to apply our inductive hypothesis. We immediately get, by definition of **Prop**, $n \in \text{Prop}_{\mathcal{N}_2}(S)$. \square

Corollary 1. Let $\mathcal{N}_1, \mathcal{N}_2$ be the same except for their weights. Then $\mathcal{N}_1 \preceq \mathcal{N}_2$ iff for all $n \in N$ and for those $m_1, \dots, m_k \in N$ such that $(m_i, n) \in E$, $(**)$ holds.

Proof. This follows straightforwardly from the previous lemma:

$$\begin{aligned} \mathcal{N}_1 \preceq \mathcal{N}_2 & \text{ iff } \text{for all } S, \text{Prop}_{\mathcal{N}_1}(S) \subseteq \text{Prop}_{\mathcal{N}_2}(S) & (\text{By defn of } \preceq) \\ & \text{ iff } (**) \text{ holds for all } n \text{ and its predecessors } m_1, \dots, m_k & (\text{By Lemma 2}) \end{aligned} \quad \square$$

This test is a convenient way to show $\mathcal{N}_1 \preceq \mathcal{N}_2$: prove that for all n , if n activates in \mathcal{N}_1 then n activates in \mathcal{N}_2 . But the test is not always applicable — often we still need to do a full proof by induction and leverage the inductive hypothesis to prove net containment properties. But the proof is still worth understanding, since it serves as an example for how to do these (somewhat nasty) inductive proofs.

We are now ready to state and prove the following algebraic characterization of **Hebb**. Note that $\text{Prop}(S)$ abbreviates $\text{Prop}_{\mathcal{N}}(S)$, the propagation in the net before update.

Theorem 3. For all $\mathcal{N}, \mathcal{N}_1, \mathcal{N}_2 \in \text{Net}$ and $A, B, B_1, \dots, B_k \in \text{Set}$, **Hebb** satisfies

(Monotonicity in \mathcal{N}) if $\mathcal{N}_1 \preceq \mathcal{N}_2$ then $\text{Hebb}(\mathcal{N}_1, A) \preceq \text{Hebb}(\mathcal{N}_2, A)$

(Absorption) $\text{Hebb}(\mathcal{N}, \text{Prop}(A)) \cong \text{Hebb}(\mathcal{N}, A)$

(Local) $\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B) \subseteq \text{Prop}(A) \cup \text{Prop}(B)$

(Cumulative) If $\text{Prop}(B_1) \subseteq \text{Prop}(B_2)$ and $\text{Prop}(B_2) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_1)$, then $\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_1) = \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_2)$

(Loop) If $\text{Prop}(B_1) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_0), \dots, \text{Prop}(B_n) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_{n-1})$, and $\text{Prop}(B_0) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_n)$, then $\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_i) = \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B_j)$ for all $i, j \in \{0, \dots, n\}$

Proof. Unfortunately, we are not able to use Lemma 2, so we resort to doing full proofs by induction on **Prop**.

(Monotonicity in \mathcal{N}) Suppose $\mathcal{N}_1 \preceq \mathcal{N}_2$, i.e. for all B , $\text{Prop}_{\mathcal{N}_1}(B) \subseteq \text{Prop}_{\mathcal{N}_2}(B)$. Now let $A, B \in \text{Set}$. We want to show that $\text{Hebb}(\mathcal{N}_1, A) \preceq \text{Hebb}(\mathcal{N}_2, A)$; in other words, for all A, B ,

$$\text{Prop}_{\text{Hebb}(\mathcal{N}_1, A)}(B) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}_2, A)}(B)$$

Let $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}_1, A)}(B)$, and proceed by induction on the structure of this Prop .

Base Step: $n \in B$. So $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}_2, A)}(B)$ by the base case of Prop .

Inductive Step: We have $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}_1, A)}(B)$ from the constructor case, i.e. because for those m_1, \dots, m_k such that $(m_i, n) \in E$,

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\text{Hebb}(\mathcal{N}_1, A)}(m_i, n) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}_1, A)}(B)}(m_i))) = 1$$

TODO

(Absorption) To prove $\text{Hebb}(\mathcal{N}, \text{Prop}(A)) \cong \text{Hebb}(\mathcal{N}, A)$, we need to show that for all sets A, B ,

$$\text{Prop}_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(B) = \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$$

For the (\supseteq) direction, let $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(B)$. We proceed by induction on this Prop .

Base Step: $n \in B$, and so $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$ by the base case of Prop .

Inductive Step: Let m_1, \dots, m_k be the predecessors of n . We have

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(m_i, n) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(B)}(m_i))) = 1$$

First, by inductive hypothesis, each $m_i \in \text{Prop}_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(B)$ iff $m_i \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$, and so

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\text{Hebb}(\mathcal{N}, \text{Prop}(A))}(m_i, n) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)}(m_i))) = 1$$

From here, we expand the weights in the updated net $\text{Hebb}(\mathcal{N}, \text{Prop}(A))$:

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} (W_{\mathcal{N}}(m_i, n) \cdot \chi_{\text{Prop}(\text{Prop}(A))}(m_i) \cdot \chi_{\text{Prop}(\text{Prop}(A))}(n)) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)}(m_i))) = 1$$

By idempotence of Prop , we can collapse each $\text{Prop}(\text{Prop}(A)) = \text{Prop}(A)$.

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} (W_{\mathcal{N}}(m_i, n) \cdot \chi_{\text{Prop}(A)}(m_i) \cdot \chi_{\text{Prop}(A)}(n)) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)}(m_i))) = 1$$

But these is exactly the weights for $\text{Hebb}(\mathcal{N}, A)$, and so

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\text{Hebb}(\mathcal{N}, A)}(m_i, n) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)}(m_i))) = 1$$

We conclude that $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$, by the constructor of Prop .

As for the (\supseteq) direction, both the base step and inductive step hold similarly (this time, expanding the nested Prop using idempotence in reverse).

(Local) Rather than showing $\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B) \subseteq \text{Prop}(A) \cup \text{Prop}(B)$ directly, we instead prove a stronger claim:

$$\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B) - \text{Prop}(A) = \text{Prop}(B) - \text{Prop}(A) \quad (***)$$

First, let's explain why this is sufficient. Suppose $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$. There are two cases:

- $n \in \text{Prop}(A)$. So clearly $n \in \text{Prop}(A) \cup \text{Prop}(B)$.
- $n \notin \text{Prop}(A)$. So $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B) - \text{Prop}(A)$. By $(***)$, this means $n \in \text{Prop}(B) - \text{Prop}(A)$. But then $n \in \text{Prop}(B)$, so $n \in \text{Prop}(A) \cup \text{Prop}(B)$.

This means that $(***)$ implies the Local property. We have left to prove $(***)$. For the (\subseteq) direction, suppose $n \in \text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B) - \text{Prop}(A)$, and proceed by structural induction on $\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)$.

Base Step: $n \in B$. By the base case of **Prop**, $n \in \text{Prop}(B)$. So $n \in \text{Prop}(B) - \text{Prop}(A)$.

Inductive Step: Let m_1, \dots, m_k be the predecessors of n . We have

$$O^{(n)}(A^{(n)}(\sum_{m_i \in n} W_{\text{Hebb}(\mathcal{N}, A)}(m_i, n) \cdot \chi_{\text{Prop}_{\text{Hebb}(\mathcal{N}, A)}(B)}(m_i))) = 1$$

(Cumulative & Loop) Finally, the Cumulative and Loop properties for **Hebb** follow from the Cumulative and Loop properties for **Prop** in Theorem 1 by substituting (for each $i \in \{0, \dots, n\}$)

$$\begin{array}{lcl} \mathcal{N} & \rightsquigarrow & \text{Hebb}(\mathcal{N}, A) \\ B_i & \rightsquigarrow & \text{Prop}_{\mathcal{N}}(B_i) \end{array}$$

□

Proposition 2. **Hebb** is not monotonic in S .

Proof. **FILL IN**

□

Syntax and Semantics

We can now introduce the logic of Hebbian learning. Let p, q, \dots be finitely many propositional variables. These represent fixed, ‘ontic’ states, i.e. established choices of neurons that correspond to features in the external world. For example, p might be the set of neurons that encapsulates the color *pink*. We presume that we already agree on these states, although we acknowledge that this is a major unresolved empirical issue. As for more complex formulas:

Definition 6. Formulas of our language \mathcal{L} are given by

$$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \Rightarrow \varphi \mid \mathbf{T}\varphi \mid [\varphi^+]\varphi$$

where p is any propositional variable. We define $\top, \perp, \vee, \rightarrow, \leftrightarrow, \Leftrightarrow$, and the dual modalities $\langle \mathbf{T} \rangle, \langle \varphi^+ \rangle$ in the usual way.

The modalities \mathbf{T} and $[\varphi^+]$ reflect our two operations **Prop** and **Hebb**, respectively. We intend for $\mathbf{T}\varphi$ to denote “the propagation of signal φ ,” and for $[\varphi^+]\psi$ to denote “after performing Hebbian update on φ , evaluate ψ .” We import $\varphi \Rightarrow \psi$ from [Leitgeb, 2001], read “the propagation of signal φ contains ψ ”. Note that $\varphi \Rightarrow \psi$ is redundant (equivalent to $\mathbf{T}\varphi \rightarrow \psi$ using the semantics below), though we keep it in our syntax because it conveniently expresses “the net classifies φ as ψ ” (if φ is interpreted as an input and ψ as a classification).

Our formulas also have more classical alternative readings, divorced from the dynamics of neural networks. Following [Leitgeb, 2001], we will define $\varphi \Rightarrow \psi$ such that it has the conditional reading “typically φ are ψ ” (where φ and ψ are read as generics, e.g. “typically birds fly”). This gives us a natural preferential reading for $\mathbf{T}\varphi$ as “typically φ ” or “the typical φ .”² Finally, Hebbian learning $[\varphi^+]\psi$ has a dual reading as *preference upgrade* [Van Benthem and Liu, 2007]. As mentioned in the Related Work section, we leave the question concerning how $[\varphi^+]$ can be viewed classically as updating a preference relation to future work.

A model of our logic is just a BFNN \mathcal{N} equipped with an interpretation function $\llbracket \cdot \rrbracket : \mathcal{L} \rightarrow \text{Set}_{\mathcal{N}}$.

Definition 7. Let $\mathcal{N} \in \text{Net}$. Our semantics are defined recursively as follows:

$\llbracket p \rrbracket$	$\in \text{Set}$ is fixed, nonempty
$\llbracket \top \rrbracket$	$= \emptyset$
$\llbracket \neg\varphi \rrbracket$	$= \overline{\llbracket \varphi \rrbracket}$
$\llbracket \varphi \wedge \psi \rrbracket$	$= \llbracket \varphi \rrbracket \cap \llbracket \psi \rrbracket$
$\llbracket \varphi \Rightarrow \psi \rrbracket$	$= \llbracket \mathbf{T}\varphi \rightarrow \psi \rrbracket$
$\llbracket \mathbf{T}\varphi \rrbracket$	$= \text{Prop}(\llbracket \varphi \rrbracket)$
$\llbracket [\varphi^+]\psi \rrbracket$	$= \llbracket \psi \rrbracket_{\text{Hebb}(\mathcal{N}, \llbracket \varphi \rrbracket)}$

²Our notation takes inspiration from [Giordano et al., 2021], which formalizes the dynamics of a net via a concept constructor \mathbf{T} in the description logic \mathcal{ALC} . Note the subtle difference between their typicality inclusions $\mathbf{T}(\varphi) \sqsubseteq \psi$ and our $\mathbf{T}\varphi \rightarrow \psi$: Ours flips the direction of containment.

Notice that these semantics are “flipped” in the sense that \wedge is interpreted as union (instead of intersection), and consequently \rightarrow is interpreted as superset (instead of subset). This choice may seem odd, but it reflects the intuition that neurons act as “elementary-feature-detectors” [Leitgeb, 2001]. For example, say $\llbracket \varphi \rrbracket$ represents those neurons that are *necessary* for detecting an apple, and $\llbracket \psi \rrbracket$ represents those neurons that are *necessary* for detecting the color red. If the net observes a red apple ($\varphi \wedge \psi$), both the neurons detecting red-features $\llbracket \varphi \rrbracket$ and the neurons detecting apple-features $\llbracket \psi \rrbracket$ necessarily activate, i.e. $\llbracket \varphi \rrbracket \cup \llbracket \psi \rrbracket$ activates. As for implication, “every apple is red” ($\varphi \rightarrow \psi$) holds for a net iff whenever the neurons detecting apple-features $\llbracket \varphi \rrbracket$ necessarily activate, so do the neurons detecting red-features $\llbracket \psi \rrbracket$. But this is only true if $\llbracket \varphi \rrbracket \supseteq \llbracket \psi \rrbracket$. This justifies us reading propositional connectives classically, despite the backwards flavor of the semantics.

Our interpretation of formulas is completely algebraic, in the sense that formulas denote sets rather than truth-values. But we can consider formulas to have truth-values as follows.

Definition 8. $\mathcal{N} \models \varphi$ iff $\llbracket \varphi \rrbracket_{\mathcal{N}} = \emptyset$.

This choice also appears to be strange at its surface. But it is a natural one in light of the fact that we defined $\llbracket \top \rrbracket := \emptyset$. For example, consider implication: $\mathcal{N} \models \varphi \rightarrow \psi$ holds iff $\llbracket \varphi \rightarrow \psi \rrbracket = \emptyset = \llbracket \top \rrbracket$, which holds iff $\llbracket \varphi \rrbracket \supseteq \llbracket \psi \rrbracket$ by our semantics.

A curious consequence is that if $\mathcal{N} \models \varphi$ and φ cannot be written to contain an implication \rightarrow , then φ must be a tautology. But we do not consider this troubling, since it only makes sense to consider a neural network’s judgment of φ when given a state $\llbracket \psi \rrbracket$ the net is in.

Basic Axioms and Inference Rules		T Axioms	
(PC)	All propositional tautologies	(LOOP)	$(\mathbf{T}\varphi_0 \rightarrow \varphi_1) \wedge \dots \wedge (\mathbf{T}\varphi_k \rightarrow \varphi_0)$ $\rightarrow (\mathbf{T}\varphi_0 \leftrightarrow \mathbf{T}\varphi_k)$
(MP)	$\frac{\varphi \quad \varphi \rightarrow \psi}{\psi}$	(DUAL)	$\langle \mathbf{T} \rangle \varphi \leftrightarrow \neg \mathbf{T} \neg \varphi$
(NECT)	$\frac{\varphi}{\mathbf{T}\varphi}$	(T)	$\mathbf{T}\varphi \rightarrow \varphi$
(NEC+)	$\frac{\psi}{[\varphi^+] \psi}$	(4)	$\mathbf{T}\varphi \rightarrow \mathbf{T}\mathbf{T}\varphi$
Reduction Axioms		Interaction Axioms	
(R _p)	$[\varphi^+] p \leftrightarrow p$	(NEST _T)	$[\mathbf{T}\varphi^+] \psi \leftrightarrow [\varphi^+] \psi$
(R _¬)	$[\varphi^+] \neg \psi \leftrightarrow \neg [\varphi^+] \psi$	(NS)	$[\varphi^+] \mathbf{T}\psi \rightarrow \mathbf{T}[\varphi^+] \psi$
(R _∧)	$[\varphi^+] (\psi \wedge \rho) \leftrightarrow ([\varphi^+] \psi \wedge [\varphi^+] \rho)$	(TP)	$\mathbf{T}[\varphi^+] \psi \wedge \mathbf{T}\varphi \rightarrow [\varphi^+] \mathbf{T}\psi$

Figure 2: A list of sound rules and axioms of the logic of Hebbian learning. We leave the question of completeness to future work.

Inference and Axioms

The proof system for our logic is as follows. We have $\vdash \varphi$ iff either φ is an axiom, or φ follows from previously obtained formulas by one of the inference rules. If $\Gamma \subseteq \mathcal{L}$ is a set of formulas, we consider $\Gamma \vdash \varphi$ to hold whenever there exist finitely many $\psi_1, \dots, \psi_k \in \Gamma$ such that $\vdash \psi_1 \wedge \dots \wedge \psi_k \rightarrow \varphi$.

We list the axioms and inference rules for our logic in Figure 2. Our main result is the soundness of these axioms and rules — we do not claim that this list forms a complete axiomatization (we revisit the question of completeness in the Conclusion).

The static base of our logic is the modal logic characterized by **T**. If we translate $\varphi \Rightarrow \psi$ via $\mathbf{T}\varphi \rightarrow \psi$ as before, we see that this modal logic contains the conditional logic **CL** (loop-cumulative). As a modality, **T** is neither normal, regular, nor monotonic, but it is classical. Note for instance that the normal modal property (K) (expressed in terms of **T**) is equivalent to

$$(K) \quad \mathbf{T}(\varphi \wedge \psi) \leftrightarrow (\mathbf{T}\varphi \wedge \mathbf{T}\psi)$$

neither direction of which is sound in our logic. Instead, we have the (LOOP) axiom expressing the loop-cumulativity of Prop. In fact, (LOOP) is all that is needed for this weakening of monotonicity — note that we have dropped cumulativity (C_{\Rightarrow}), as well as (LOOP₊) and (C_+) for $[\varphi^+]$, since they all follow from the present axioms (compare against [Kisby et al., 2022]).

Since Hebbian update only affects the propagation of states, we have reduction axioms (R_p), (R_¬), (R_∧), as well as the axiom (NEST_T) for terms that nest **T** within $[\varphi^+]$.

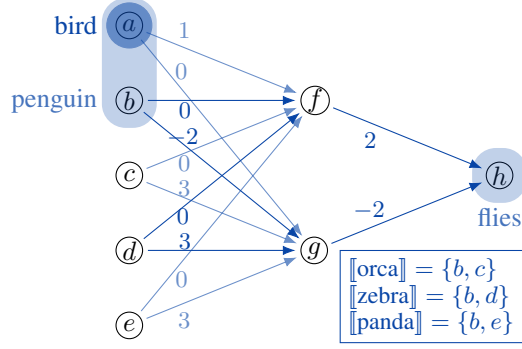


Figure 3: A BFNN \mathcal{N} , equipped with the ReLU activation function, $T = 1$, and $\eta = 1$. After observing the dataset $\langle \text{orca}, \text{zebra}, \text{panda} \rangle$, \mathcal{N} learns that penguins do not fly, while preserving the fact that birds typically fly. **TODO**

In lieu of a full reduction for $[\varphi^+]\mathbf{T}\psi$, we instead have the weaker axioms (NS) and (TP). These two axioms capture key cognitive biases of a Hebbian agent. Consider the axiom (TP), i.e. (Typicality Preservation)

$$(TP) \quad \mathbf{T}[\varphi^+]\psi \wedge \mathbf{T}\varphi \rightarrow [\varphi^+]\mathbf{T}\psi$$

This says that if our agent expects ψ is normally true after learning φ , but she also happens to expect φ , then after learning φ the typicality of ψ will be preserved. This is a peculiar kind of cognitive bias whereby a Hebbian agent maintains her prior attitudes when presented with news she already expects.

The axiom (NS), i.e. (No Surprises)

$$(NS) \quad [\varphi^+]\mathbf{T}\psi \rightarrow \mathbf{T}[\varphi^+]\psi$$

says that if after learning φ , our agent thinks normally ψ , then she would have expected ψ to be true after learning φ in the first place. Loosely: She will never be surprised.

Soundness of these axioms is largely a matter of matching each axiom with its corresponding property of Hebb.

Theorem 4. The rules and axioms above are sound, i.e. hold for all $\mathcal{N} \in \text{Net}$.

Proof. **FILL IN** □

Applying the Logic: A Concrete Example

We now demonstrate our neuro-symbolic interface by way of an example neural network in a machine learning context. The task: Given an image of an animal, classify it as flying or non-flying. Suppose we have the partially pre-trained BFNN \mathcal{N} in Figure 3.

For simplification's sake, let's suppose that our animal images can be reduced to 5-dimensional vectors in order to be fed into the input layer of \mathcal{N} . Say:

penguin	$\langle 11000 \rangle$	orca	$\langle 01100 \rangle$
zebra	$\langle 01010 \rangle$	panda	$\langle 01001 \rangle$

In addition, suppose an image activates the first node if and only if it depicts a bird.

We can identify each animal with the set of nodes it activates in the input layer. This gives us the sets shown in Figure 3. We can also identify the class of things that fly with the output node, i.e. $\llbracket \text{flies} \rrbracket = \{h\}$. In principle we can identify propositions with sets containing hidden nodes as well, although in practice the meaning of hidden nodes is often unclear.

With this interpretation in mind, we see that $\mathcal{N} \models \text{bird} \Rightarrow \text{flies}$, but also $\mathcal{N} \models \text{penguin} \Rightarrow \text{flies}$ (which is incorrect). Our hope is that \mathcal{N} corrects this mistake via Hebbian learning.

Say we expose \mathcal{N} to non-flying animals that share the black-and-white color of penguins, e.g. we train \mathcal{N} on the dataset $\langle \text{orca}, \text{zebra}, \text{panda} \rangle$. The propagations of each instance will increase W_{bg} . Once we have given \mathcal{N} the entire dataset ($W_{bg} = 1$), $\text{Prop}(\llbracket \text{penguin} \rrbracket)$ will contain g , which will cancel the signal given by $f \rightarrow h$. Our logic successfully models this behavior:

$$\begin{aligned} \mathcal{N} &\models [\text{orca}^+][\text{zebra}^+][\text{panda}^+](\text{bird} \Rightarrow \text{flies}), \text{ yet} \\ \mathcal{N} &\not\models [\text{orca}^+][\text{zebra}^+][\text{panda}^+](\text{penguin} \Rightarrow \text{flies}) \end{aligned}$$

i.e. \mathcal{N} learns that penguins do not fly while preserving the fact that birds typically fly.

Conclusion and Future Work

In this paper, we gave sound axioms and rules characterizing the logic of Hebbian learning. This logic interfaces the neuro-symbolic divide by characterizing conditionals \Rightarrow and modalities \mathbf{T} , $[\varphi^+]$ in terms of the propagation and Hebbian update of signals in a neural network. The upshot of all this is that this logic describes a neuro-symbolic agent that learns associatively and also reasons about what it has learned.

We leave open the question of whether the axioms and rules we list are complete. But we take this opportunity to stress the importance of having strong completeness for logics of this kind. Strong completeness for a *static* neural semantics provides a bridge across which we can extract a set of rules Γ from an interpreted network, and also build an interpreted neural network implementing Γ . But once the neural network updates, we lose the interpretations of neurons that allow for these translations. If we had strong completeness for the *dynamic* logic, we could fully track the interpretations while the net learns and preserve this neuro-symbolic correspondence.

Beyond the logic of Hebbian learning, we believe that this framework will be a fruitful way to explore the neuro-symbolic interface for a variety of neural networks and learning policies. Exciting future directions include:

1. Mapping more expressive syntax to neural activity
2. Generalizing to a broader class of neural networks
3. Generalizing to a broader class of activation functions
4. Characterizing other learning policies in logical terms

The holy grail of this line of work is to completely axiomatize the (1) first-order logic of (2) nonbinary (fuzzy-valued) neural networks with (3) more varied (e.g. ReLU and GELU) activation functions that (4) learn via backpropagation.

Acknowledgements

We thank the anonymous reviewers for their careful reviews and helpful comments. C. Kisby was supported in part by the US Department of Defense [Contract No. W52P1J2093009].

Corrections to the FLAIRS Paper

The original FLAIRS paper [Kisby et al., 2022] contained an error, starting with the definition of propagation (Definition 3) and affecting the results of Lemma 2, Corollary 1, Theorem 3, and Theorem 4. We only noticed the error after trying to formally verify some of these results in Lean. See our follow-up paper [Kisby et al., 2024], where the major results have been checked in Lean.

In this arXiv version, we have presented corrected definitions and results. In this section we will state which statements in the FLAIRS paper were false, and explain how they have been corrected.

In the FLAIRS paper, propagation was defined as follows:

Definition 3. Let $\text{Prop} : \text{Set} \rightarrow \text{Set}$ be defined recursively as follows: $n \in \text{Prop}(S)$ iff either

(Base Case) $n \in S$, or

(Constructor) For those $m_1, \dots, m_k \in \text{Prop}(S)$ such that $(m_i, n) \in E$ we have

$$O^{(n)}(A^{(n)}(\vec{W}(m_i, n))) = 1$$

The issue here is surprisingly subtle: We only apply the activation function to the weights of those predecessors that were already active in the previous step ($m_1, \dots, m_k \in \text{Prop}(S)$). This means that the vector $\vec{W}(m_i, n)$ doesn't comply with the arity that $A^{(n)}$ is expecting; since the arity of $A^{(n)}$ is the indegree of n , we must give $A^{(n)}$ the vector of weights over *all* predecessors. But a more serious problem is that this tricks us into mistakenly ignoring those predecessors which were *not* already active.

In machine learning work, the standard way of doing this is to feed the activation function a weighted sum of all predecessor activations. This simple fix resolves the issue. So in this draft we instead write:

Definition 3. Let $\text{Prop} : \text{Set} \rightarrow \text{Set}$ be defined recursively as follows: $n \in \text{Prop}(S)$ iff either

(Base Case) $n \in S$, or

(Constructor) For those m_1, \dots, m_k such that $m_i E n$ we have

$$O^{(n)}(A^{(n)}(\sum_{m_i E n} W(m_i, n) \cdot \chi_{\text{Prop}(S)}(m_i))) = 1$$

This inner sum term reflects the weighted sum of predecessor activations, and $\chi_{\text{Prop}(S)}(m_i)$ tells us whether each m_i is already active (in $\text{Prop}(S)$).

Lemma 2. Suppose \mathcal{N}_1 and \mathcal{N}_2 are the same except for their weights, and let $S \in \text{Set}$. Then $\text{Prop}_{\mathcal{N}_1}(S) \subseteq \text{Prop}_{\mathcal{N}_2}(S)$ iff for all $n \in \text{Prop}_{\mathcal{N}_1}(S)$ and for those $m_1, \dots, m_k \in \text{Prop}_{\mathcal{N}_1}(S)$ such that $(m_i, n) \in E$,

$$\begin{aligned} O^{(n)}(A^{(n)}(\vec{W}_{\mathcal{N}_1}(m_i, n))) &= 1 \\ \text{implies} & \\ O^{(n)}(A^{(n)}(\vec{W}_{\mathcal{N}_2}(m_i, n))) &= 1 \end{aligned} \quad (**)$$

Corollary 1. Let $\mathcal{N}_1, \mathcal{N}_2$ be the same except for their weights. Then $\mathcal{N}_1 \preceq \mathcal{N}_2$ iff for all $n \in N$ and for those $m_1, \dots, m_k \in N$ such that $(m_i, n) \in E$, $(**)$ holds.

Theorem 3. For all $\mathcal{N}, \mathcal{N}_1, \mathcal{N}_2 \in \text{Net}$ and $S, S_1, S_2 \in \text{Set}$, Hebb satisfies

(Inclusion) $\mathcal{N} \preceq \text{Hebb}(\mathcal{N}, S)$

(Absorption) $\text{Hebb}(\mathcal{N}, \text{Prop}(S)) \cong \text{Hebb}(\mathcal{N}, S)$

(Monotonicity in \mathcal{N}) if $\mathcal{N}_1 \preceq \mathcal{N}_2$ then $\text{Hebb}(\mathcal{N}_1, S) \preceq \text{Hebb}(\mathcal{N}_2, S)$

(Local) $\text{Prop}_{\text{Hebb}(\mathcal{N}, S_2)}(S_1) \subseteq \text{Prop}_{\mathcal{N}}(S_1) \cup \text{Prop}_{\mathcal{N}}(S_2)$

(Cumulative) If $\text{Prop}_{\mathcal{N}}(S_1) \subseteq \text{Prop}_{\mathcal{N}}(S_2)$ and $\text{Prop}_{\mathcal{N}}(S_2) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_1)$, then $\text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_1) = \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_2)$

(Loop) If $\text{Prop}_{\mathcal{N}}(S_1) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_0), \dots, \text{Prop}_{\mathcal{N}}(S_n) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_{n-1})$, and $\text{Prop}_{\mathcal{N}}(S_0) \subseteq \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_n)$, then $\text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_i) = \text{Prop}_{\text{Hebb}(\mathcal{N}, S)}(S_j)$ for all $i, j \in \{0, \dots, n\}$

(State the axioms that were wrong, and the revised form of them — really, this is about No Surprises and Typicality Preservation!)

References

- Martín Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration-a structured survey. *arXiv preprint cs/0511042*, 2005.
- Christian Balkenius and Peter Gärdenfors. Nonmonotonic Inferences in Neural Networks. In *KR*, pages 32–39, 1991.
- Alexandru Baltag, Nina Gierasimczuk, Aybükce Özgün, Ana Lucia Vargas Sandoval, and Sonja Smets. A dynamic logic for learning theory. *Journal of Logical and Algebraic Methods in Programming*, 109:100485, 2019a.
- Alexandru Baltag, Dazhu Li, and Mina Young Pedersen. On the right path: a modal logic for supervised learning. In *International Workshop on Logic, Rationality and Interaction*, pages 1–14. Springer, 2019b.
- Reinhard Blutner. Nonmonotonic inferences and neural networks. In *Information, Interaction and Agency*, pages 203–234. Springer, 2004.
- Artur SD’Avila Garcez, Luis C Lamb, and Dov M Gabbay. *Neural-symbolic cognitive reasoning*. Springer Science & Business Media, 2008.
- AS d’Avila Garcez, Krysia Broda, and Dov M Gabbay. Symbolic knowledge extraction from trained neural networks: A sound approach. *Artificial Intelligence*, 125(1-2):155–207, 2001.
- Laura Giordano, Valentina Gliozzi, and Daniele Theseider Dupré. From common sense reasoning to neural network models through multiple preferences: An overview. *CoRR*, abs/2107.04870, 2021. URL <https://arxiv.org/abs/2107.04870>.

- Donald Hebb. *The organization of behavior: A neuropsychological theory*. Wiley, New York, 1949.
- Caleb Kisby, Saúl Blanco, and Lawrence Moss. The logic of hebbian learning. In *The International FLAIRS Conference Proceedings*, volume 35, 2022.
- Caleb Schultz Kisby, Saúl A Blanco, and Lawrence S Moss. What do hebbian learners learn? reduction axioms for iterated hebbian learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 14894–14901, 2024.
- Sarit Kraus, Daniel Lehmann, and Menachem Magidor. Nonmonotonic reasoning, preferential models and cumulative logics. *Artificial intelligence*, 44(1-2):167–207, 1990.
- Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets. *Artificial Intelligence*, 128(1-2):161–201, 2001.
- Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets II. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 11(supp02):105–135, 2003.
- Hannes Leitgeb. Neural Network Models of Conditionals. In *Introduction to Formal Philosophy*, pages 147–176. Springer, 2018.
- Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, 1943.
- Md Kamruzzaman Sarker, Lu Zhou, Aaron Eberhart, and Pascal Hitzler. Neuro-symbolic artificial intelligence: Current trends. *arXiv preprint arXiv:2105.05330*, 2021.
- Leslie G Valiant. Three problems in computer science. *Journal of the ACM (JACM)*, 50(1):96–99, 2003.
- Johan Van Benthem. Dynamic logic for belief revision. *Journal of applied non-classical logics*, 17(2):129–155, 2007.
- Johan Van Benthem and Fenrong Liu. Dynamic logic of preference upgrade. *Journal of Applied Non-Classical Logics*, 17(2):157–182, 2007.