

COMPUTATIONAL LEARNING IN DYNAMIC LOGICS

DAY 3: UPDATES ON NEURAL NETWORKS

Nina Gierasimczuk and Caleb Schultz Kisby

@NASSLLI, June 2025

Course Homepage:

<https://sites.google.com/view/nasslli25-learning-in-del>

PLAN FOR TODAY

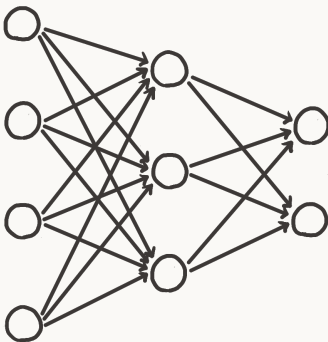
- 1 Overview of Neural Networks
- 2 A Logic for Neural Network Inference
- 3 Neural Network Update in Dynamic Logic

PLAN FOR TODAY

- 1 Overview of Neural Networks
- 2 A Logic for Neural Network Inference
- 3 Neural Network Update in Dynamic Logic

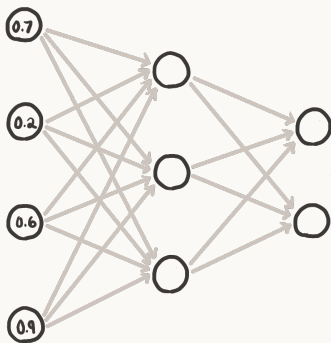
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



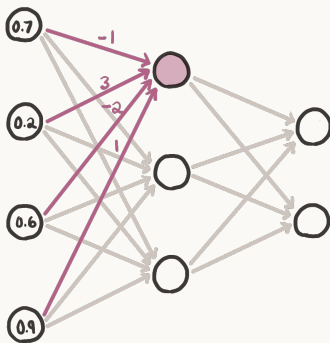
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



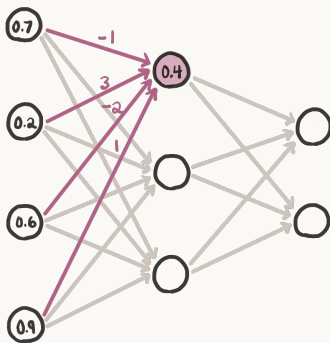
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



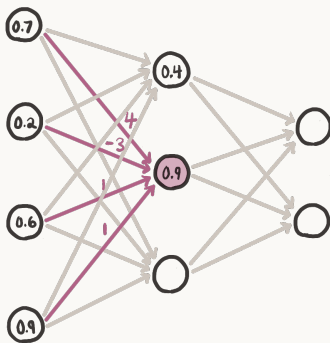
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



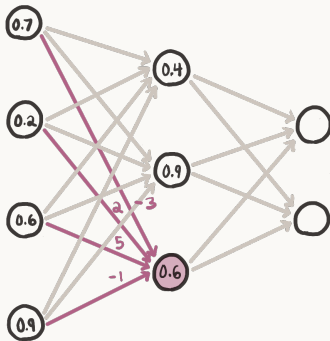
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



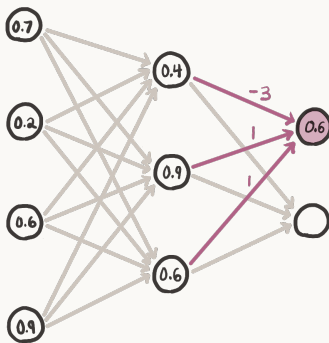
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



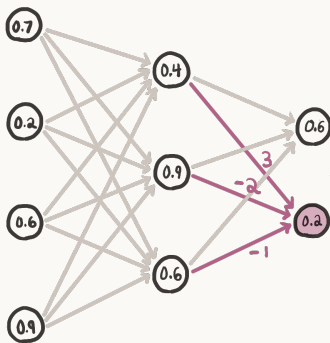
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



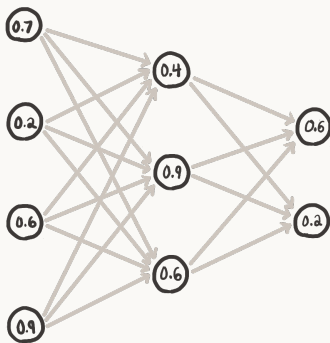
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



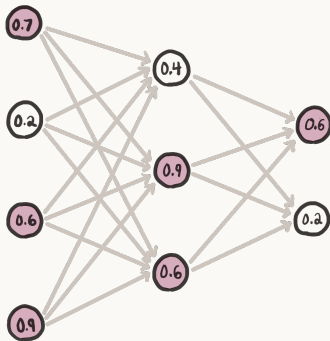
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



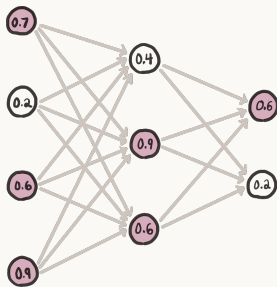
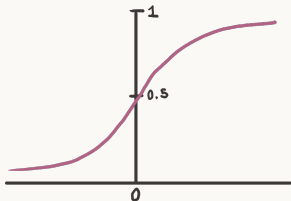
ARTIFICIAL NEURAL NETWORKS

- A neural network is just $\mathcal{N} = (N, E, W, A)$
 - neurons, edges, weights, activation function
- Neurons are successively activated by their predecessors:



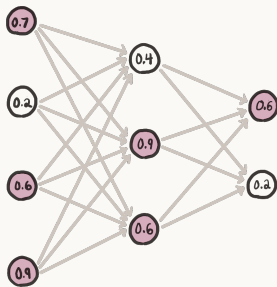
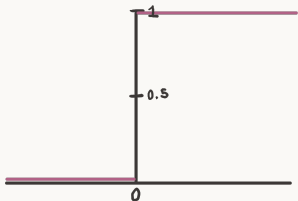
¬FLIES

(BINARY) ARTIFICIAL NEURAL NETWORKS



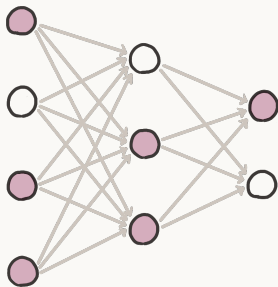
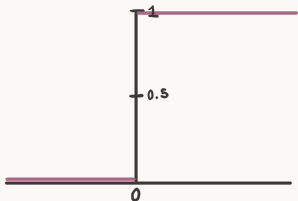
- We take the activation function A to be a binary step function
- This is a useful abstraction for connecting nets with logic, formal languages, and automata
- The net's activation patterns are just **sets of neurons**.

(BINARY) ARTIFICIAL NEURAL NETWORKS



- We take the activation function A to be a binary step function
- This is a useful abstraction for connecting nets with logic, formal languages, and automata
- The net's activation patterns are just **sets of neurons**.

(BINARY) ARTIFICIAL NEURAL NETWORKS



- We take the activation function A to be a binary step function
- This is a useful abstraction for connecting nets with logic, formal languages, and automata
- The net's activation patterns are just **sets of neurons**.

PLAN FOR TODAY

- 1 Overview of Neural Networks
- 2 A Logic for Neural Network Inference
- 3 Neural Network Update in Dynamic Logic

SYNTAX: LANGUAGE OF NEURAL NETWORK INFERENCE

Definition (Language of Epistemic Logic)

Take a countable set of propositions PROP .

$$\varphi := \top \mid p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{A}\varphi \mid \langle \mathbf{C} \rangle \varphi$$

for all $p \in \text{PROP}$. The usual abbreviations are \vee , \rightarrow , and \mathbf{C} (dual to $\langle \mathbf{C} \rangle$)

SYNTAX: LANGUAGE OF NEURAL NETWORK INFERENCE

Definition (Language of Epistemic Logic)

Take a countable set of propositions PROP .

$$\varphi := \top \mid p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{A}\varphi \mid \langle \mathbf{C} \rangle \varphi$$

for all $p \in \text{PROP}$. The usual abbreviations are \vee , \rightarrow , and \mathbf{C} (dual to $\langle \mathbf{C} \rangle$)

- Notice that we're giving the semantics in terms of the \diamond -variant $\langle \mathbf{C} \rangle$

SYNTAX: LANGUAGE OF NEURAL NETWORK INFERENCE

Definition (Language of Epistemic Logic)

Take a countable set of propositions PROP .

$$\varphi := \top \mid p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{A}\varphi \mid \langle \mathbf{C} \rangle \varphi$$

for all $p \in \text{PROP}$. The usual abbreviations are \vee , \rightarrow , and \mathbf{C} (dual to $\langle \mathbf{C} \rangle$)

- Notice that we're giving the semantics in terms of the \diamond -variant $\langle \mathbf{C} \rangle$
- Just like before, we will interpret $\langle \mathbf{C} \rangle \varphi$ in a model, at a world.

SYNTAX: LANGUAGE OF NEURAL NETWORK INFERENCE

Definition (Language of Epistemic Logic)

Take a countable set of propositions PROP .

$$\varphi := \top \mid p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{A}\varphi \mid \langle \mathbf{C} \rangle \varphi$$

for all $p \in \text{PROP}$. The usual abbreviations are \vee , \rightarrow , and \mathbf{C} (dual to $\langle \mathbf{C} \rangle$)

- Notice that we're giving the semantics in terms of the \diamond -variant $\langle \mathbf{C} \rangle$
- Just like before, we will interpret $\langle \mathbf{C} \rangle \varphi$ in a model, at a world.
- **The intended interpretation:**

$\langle \mathbf{C} \rangle \varphi$ holds in a net, at a neuron w if w is activated by input φ .

SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- A binary neural network is just $\mathcal{N} = (N, E, W, A)$

SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- A binary neural network is just $\mathcal{N} = (N, E, W, A)$
- Each choice of E, W, A specifies a transition function from one activation pattern $S \subseteq N$ to the next

SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- A binary neural network is just $\mathcal{N} = (N, E, W, A)$
- Each choice of E, W, A specifies a transition function from one activation pattern $S \subseteq N$ to the next
- Given initial state S_0 , $F_{S_0} : \wp(N) \rightarrow \wp(N)$ is given by

$$F_{S_0}(S) = S_0 \cup \{w \mid A(\sum_{u \in \text{preds}(w)} W(u, w) \cdot \chi_S(u)) = 1\}$$

“the set of all nodes w activated by their immediate predecessors u ”

SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- A binary neural network is just $\mathcal{N} = (N, E, W, A)$
- Each choice of E, W, A specifies a transition function from one activation pattern $S \subseteq N$ to the next
- Given initial state S_0 , $F_{S_0} : \wp(N) \rightarrow \wp(N)$ is given by

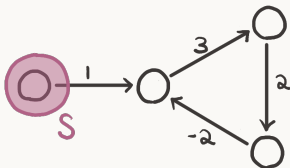
$$F_{S_0}(S) = S_0 \cup \{w \mid A(\sum_{u \in \text{preds}(w)} W(u, w) \cdot \chi_S(u)) = 1\}$$

“the set of all nodes w activated by their immediate predecessors u ”

- $\chi_S(u) = 1$ iff $u \in S$ indicates whether u was activated previously

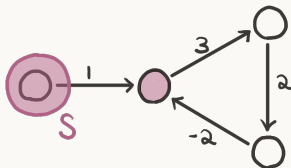
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



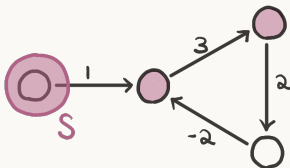
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



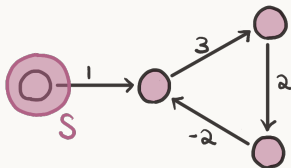
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



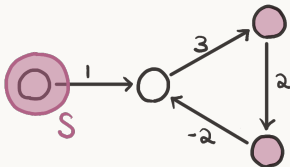
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



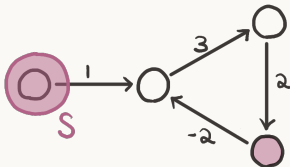
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



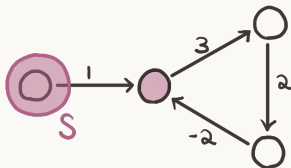
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



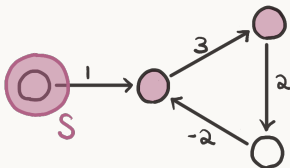
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



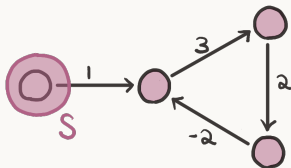
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



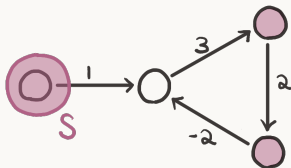
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



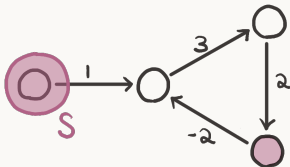
SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

- Notice that the activated nodes could have oscillatory behavior!
- But we only want nets that have a unique “answer” for each input



SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

Postulate

We assume for all $S_0 \subseteq N$, F_{S_0} repeatedly applied to S_0 ,

$$S_0, F_{S_0}(S_0), F_{S_0}(F_{S_0}(S_0)), \dots, F_{S_0}^k(S_0), \dots$$

eventually stabilizes to a unique activation pattern.

SEMANTICS: NEURAL NETWORK CLOSURE OPERATOR

Postulate

We assume for all $S_0 \subseteq N$, F_{S_0} repeatedly applied to S_0 ,

$$S_0, F_{S_0}(S_0), F_{S_0}(F_{S_0}(S_0)), \dots, F_{S_0}^k(S_0), \dots$$

eventually stabilizes to a unique activation pattern.

Definition

Let $\text{Clos} : \wp(N) \rightarrow \wp(N)$ be the function that produces this stable activation pattern.

SEMANTICS: FORMAL DEFINITION

Definition (Neural Network Semantics)

Given a binary neural network model $\mathcal{N} = (N, E, W, A, V)$, where $V : Prop \rightarrow \wp(N)$, and a neuron (“world”) $w \in N$:

$$\mathcal{N}, w \models p \quad \text{iff} \quad w \in V(p) \text{ for each } p \in Prop$$

$$\mathcal{N}, w \models \neg\varphi \quad \text{iff} \quad \text{not } \mathcal{N}, w \models \varphi$$

$$\mathcal{N}, w \models \varphi \wedge \psi \quad \text{iff} \quad \mathcal{N}, w \models \varphi \text{ and } \mathcal{N}, w \models \psi$$

$$\mathcal{N}, w \models A\varphi \quad \text{iff} \quad \text{for all } w \in N \text{ whatsoever, } \mathcal{N}, w \models \varphi$$

$$\mathcal{N}, w \models \langle \mathbf{C} \rangle \varphi \quad \text{iff} \quad w \in \text{Clos}(\llbracket \varphi \rrbracket)$$

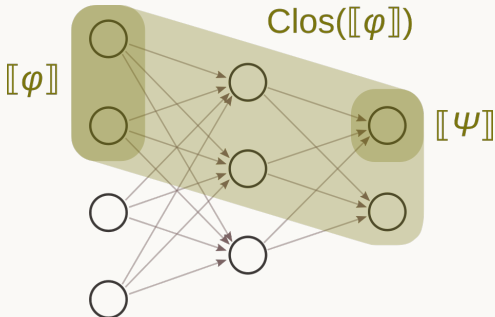
and dually:

$$\mathcal{N}, w \models \mathbf{C}\varphi \quad \text{iff} \quad w \in (\text{Clos}(\llbracket \varphi \rrbracket)^{\mathbf{C}})^{\mathbf{C}}$$

where $\llbracket \varphi \rrbracket = \{u \mid \mathcal{N}, u \models \varphi\}$ is the set of worlds where φ holds (the set of neurons that are active for φ)

EXPRESSING NEURAL NETWORK INFERENCE

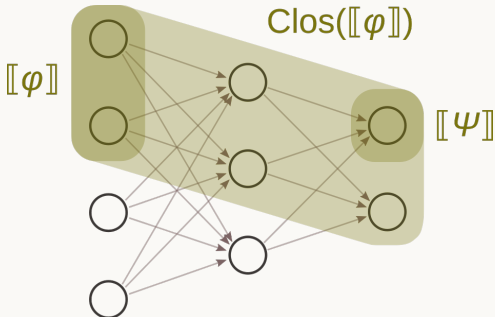
- The **C** modality gives information about the net's answer to an input



The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$

EXPRESSING NEURAL NETWORK INFERENCE

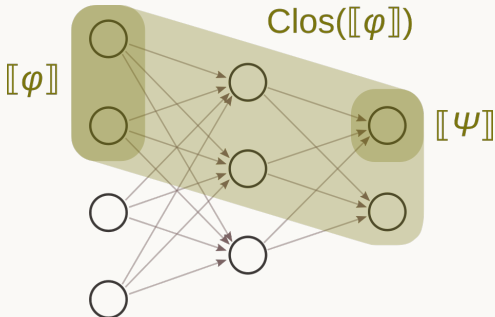
- The **C** modality gives information about the net's answer to an input



The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$

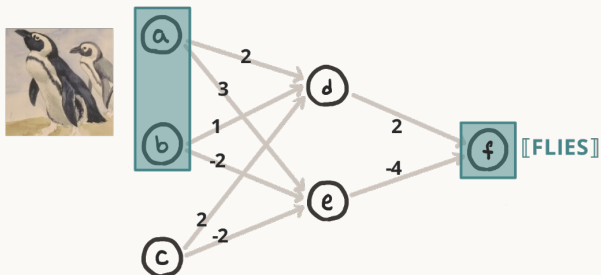
EXPRESSING NEURAL NETWORK INFERENCE

- The **C** modality gives information about the net's answer to an input



The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

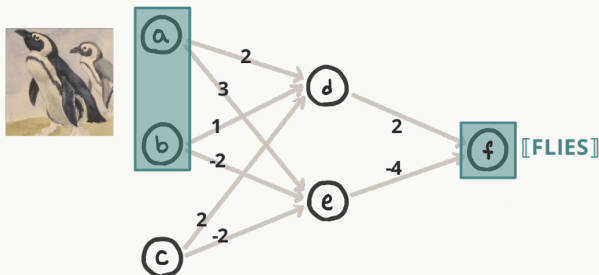
EXAMPLE: EXPRESSING NEURAL NETWORK INFERENCE



- In the exercises, we will ask you will to show

$$\mathcal{N} \not\models \mathbf{A}(\mathbf{C}(\text{PENGUIN}) \rightarrow \text{FLIES})$$

EXAMPLE: EXPRESSING NEURAL NETWORK INFERENCE

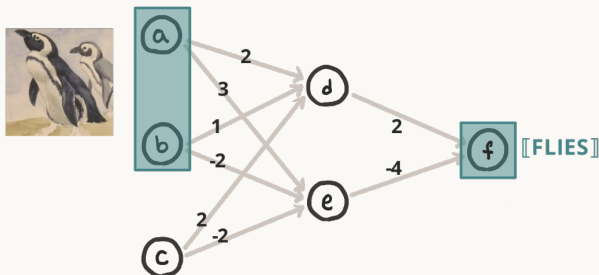


- In the exercises, we will ask you will to show

$$\mathcal{N} \neq \mathbf{A}(\mathbf{C}(\text{PENGUIN}) \rightarrow \text{FLIES})$$

- This means the net does not classify penguins as flying

EXAMPLE: EXPRESSING NEURAL NETWORK INFERENCE



- In the exercises, we will ask you will to show

$$\mathcal{N} \not\models \mathbf{A}(\mathbf{C}(\text{PENGUIN}) \rightarrow \text{FLIES})$$

- This means the net does not classify penguins as flying
- Yet, if we take $\llbracket \text{BIRD} \rrbracket = \{a, b, c\}$,

$$\mathcal{N} \models \mathbf{A}(\mathbf{C}(\text{BIRD}) \rightarrow \text{FLIES})$$

ADDITIONAL COMMENTS ON THIS LOGIC

The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

- What does this remind you of?

ADDITIONAL COMMENTS ON THIS LOGIC

The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

- What does this remind you of?
 - $\text{best}_{\leq}(\llbracket \varphi \rrbracket) \subseteq \llbracket \psi \rrbracket$

ADDITIONAL COMMENTS ON THIS LOGIC

The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

- What does this remind you of?
 - $\text{best}_{\leq}(\llbracket \varphi \rrbracket) \subseteq \llbracket \psi \rrbracket$
 - $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$, taken as a conditional, behaves exactly like $B^{\varphi}\psi$

ADDITIONAL COMMENTS ON THIS LOGIC

The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

- What does this remind you of?
 - $\text{best}_{\leq}(\llbracket \varphi \rrbracket) \subseteq \llbracket \psi \rrbracket$
 - $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$, taken as a conditional, behaves exactly like $B^{\varphi}\psi$
 - You can think of this as the net's conditional belief

ADDITIONAL COMMENTS ON THIS LOGIC

The net satisfies $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$ iff The net satisfies $\mathbf{A}(\psi \rightarrow \langle \mathbf{C} \rangle(\varphi))$
iff $\text{Clos}(\llbracket \varphi \rrbracket) \supseteq \llbracket \psi \rrbracket$
iff **The net classifies φ as ψ**

- What does this remind you of?
 - $\text{best}_{\leq}(\llbracket \varphi \rrbracket) \subseteq \llbracket \psi \rrbracket$
 - $\mathbf{A}(\mathbf{C}(\varphi) \rightarrow \psi)$, taken as a conditional, behaves exactly like $B^{\varphi}\psi$
 - You can think of this as the net's conditional belief
- Interpreting \mathbf{C} on its own is less clear...

PLAN FOR TODAY

- 1 Overview of Neural Networks
- 2 A Logic for Neural Network Inference
- 3 Neural Network Update in Dynamic Logic

UPDATES ON NEURAL NETWORKS

- **Unsupervised Updates**

- The network learns from data that is **unlabeled** (no expected answer or classification)
- Each update softly increases the net's preference for the input
- Hebb's rule, Oja's rule, & competitive learning rule

- **Supervised Updates**

- The network learns from **labeled** data with an **expected answer**
- Each update softly increases the net's accuracy on a function
- Backpropagation rule & delta learning rule

UPDATES ON NEURAL NETWORKS

- **Unsupervised Updates**

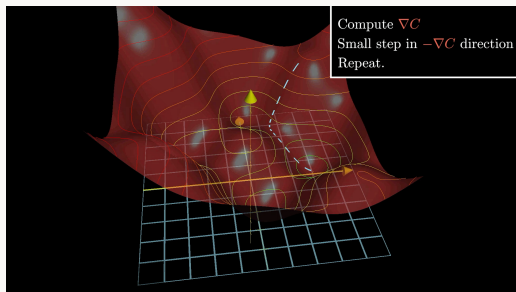
- The network learns from data that is **unlabeled** (no expected answer or classification)
- Each update softly increases the net's preference for the input
- **Hebb's rule**, Oja's rule, & competitive learning rule

- **Supervised Updates**

- The network learns from **labeled** data with an **expected answer**
- Each update softly increases the net's accuracy on a function
- **Backpropagation rule** & delta learning rule

BACKPROPAGATION RULE

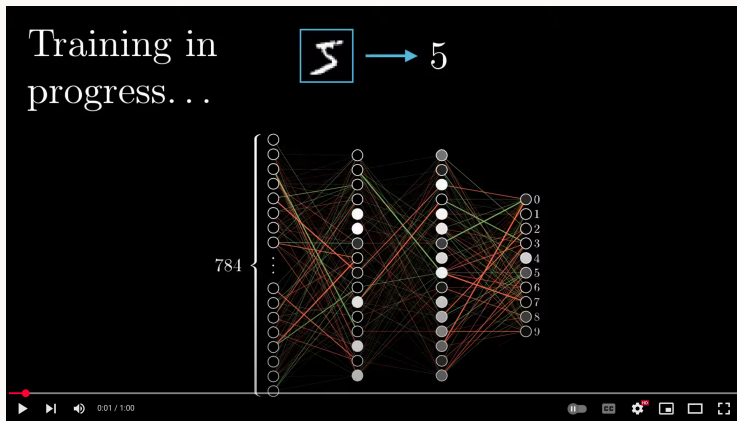
- Backpropagation is the most widely used neural network update rule
- **Main idea:** Backprop implements **gradient descent** on a net's weights



- Given an input \vec{x} with label y , the neural network gives its answer y' to \vec{x} , and each weight of the net is adjusted according to its contribution to the error (difference between y' and y).

BACKPROPAGATION RULE

<https://www.youtube.com/watch?v=cANqroNVdl8>



ADDITIONAL COMMENTS ON NEURAL NETWORK UPDATES

- What if we could have a complete characterization of Backprop, like we did for public announcement, LEX, and MINI?
 - That would be wonderful!

ADDITIONAL COMMENTS ON NEURAL NETWORK UPDATES

- What if we could have a complete characterization of Backprop, like we did for public announcement, LEX, and MINI?
 - That would be wonderful!
 - Unfortunately, this is still an open problem

ADDITIONAL COMMENTS ON NEURAL NETWORK UPDATES

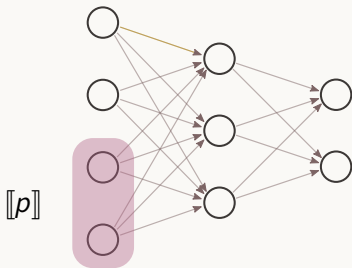
- What if we could have a complete characterization of Backprop, like we did for public announcement, LEX, and MINI?
 - That would be wonderful!
 - Unfortunately, this is still an open problem
- **Proof of concept:** Can we do this for any neural network update at all?

ADDITIONAL COMMENTS ON NEURAL NETWORK UPDATES

- What if we could have a complete characterization of Backprop, like we did for public announcement, LEX, and MINI?
 - That would be wonderful!
 - Unfortunately, this is still an open problem
- **Proof of concept:** Can we do this for any neural network update at all?
 - Let's consider the simplest possible one: **Hebbian learning**

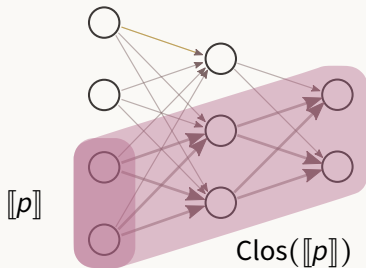
HEBBIAN UPDATE RULE

Neurons that fire together wire together



HEBBIAN UPDATE RULE

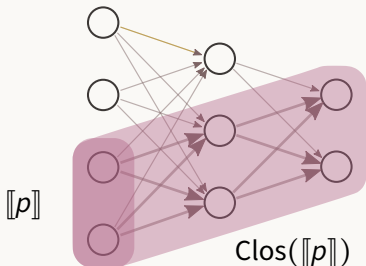
Neurons that fire together wire together



- Each edge involved in the activation is “bumped up” by a fixed learning rate $\eta \geq 0$

HEBBIAN UPDATE RULE

Neurons that fire together wire together



- Each edge involved in the activation is “bumped up” by a fixed learning rate $\eta \geq 0$
- Formally: $\text{HEBB}(\mathcal{N}, \llbracket \varphi \rrbracket) = (N, E, W', A)$, where
$$W'(u, w) = W(u, w) + \eta \cdot \chi_{\text{Clos}(\llbracket \varphi \rrbracket)}(u) \cdot \chi_{\text{Clos}(\llbracket \varphi \rrbracket)}(w)$$

HEBBIAN UPDATE RULE: EXAMPLES

- If nobody ever told you that penguins don't fly, how could you come to believe they don't?

HEBBIAN UPDATE RULE: EXAMPLES

- If nobody ever told you that penguins don't fly, how could you come to believe they don't?
 - Observe animals with similar features that **don't** fly

HEBBIAN UPDATE RULE: EXAMPLES

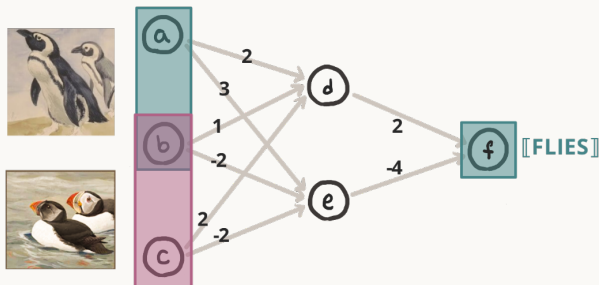
- If nobody ever told you that penguins don't fly, how could you come to believe they don't?
 - Observe animals with similar features that **don't** fly
- Now imagine you believe penguins don't fly. What could cause you to change your mind?

HEBBIAN UPDATE RULE: EXAMPLES

- If nobody ever told you that penguins don't fly, how could you come to believe they don't?
 - Observe animals with similar features that **don't** fly
- Now imagine you believe penguins don't fly. What could cause you to change your mind?
 - Observe animals with similar features that **do** fly

HEBBIAN UPDATE RULE: EXAMPLES

- If nobody ever told you that penguins don't fly, how could you come to believe they don't?
 - Observe animals with similar features that **don't** fly
- Now imagine you believe penguins don't fly. What could cause you to change your mind?
 - Observe animals with similar features that **do** fly



ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate

ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar

ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar
 - **Key idea:** All weights must sum to 1
 - “use it or lose it”

ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar
 - **Key idea:** All weights must sum to 1
 - “use it or lose it”
 - It's an open problem to completely characterize stable Hebbian update rules

ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar
 - **Key idea:** All weights must sum to 1
 - “use it or lose it”
 - It's an open problem to completely characterize stable Hebbian update rules
- HEBB is more gradual than LEX or MINI

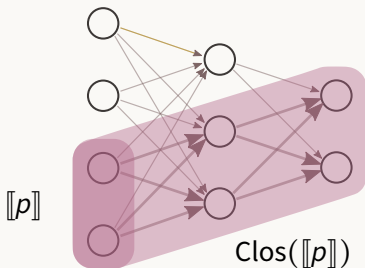
ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar
 - **Key idea:** All weights must sum to 1
 - “use it or lose it”
 - It's an open problem to completely characterize stable Hebbian update rules
- HEBB is more gradual than LEX or MINI
 - The result of LEX and MINI is a change in belief

ADDITIONAL COMMENTS ABOUT HEBBIAN UPDATE

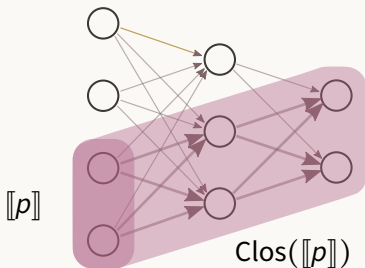
- This kind of Hebbian update is **unstable** — weights will continue to increase until they saturate
 - It's possible to prevent this by using **Oja's rule** or similar
 - **Key idea:** All weights must sum to 1
 - “use it or lose it”
 - It's an open problem to completely characterize stable Hebbian update rules
- HEBB is more gradual than LEX or MINI
 - The result of LEX and MINI is a change in belief
 - HEBB gently nudges us in the direction of a belief

HEBB*: “FIXED-POINT” HEBBIAN UPDATE



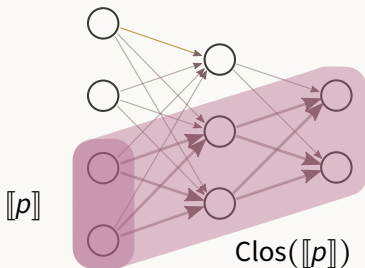
- If we repeatedly apply $\text{HEBB}(\mathcal{N}, S)$, eventually these weights will saturate (they will not inhibit any incoming activations)

HEBB*: “FIXED-POINT” HEBBIAN UPDATE



- If we repeatedly apply $\text{HEBB}(\mathcal{N}, S)$, eventually these weights will saturate (they will not inhibit any incoming activations)
- Let iter be the number of iterations needed to reach this fixed point

HEBB*: “FIXED-POINT” HEBBIAN UPDATE



- If we repeatedly apply $\text{HEBB}(\mathcal{N}, S)$, eventually these weights will saturate (they will not inhibit any incoming activations)
- Let iter be the number of iterations needed to reach this fixed point
- Let $\text{HEBB}^*(\mathcal{N}, S) = (N, E, W', A)$, where

$$W'(u, w) = W(u, w) + \text{iter} \cdot \eta \cdot \chi_{\text{Clos}(\llbracket \varphi \rrbracket)}(u) \cdot \chi_{\text{Clos}(\llbracket \varphi \rrbracket)}(w)$$

NEURAL NETWORK UPDATES IN DYNAMIC LOGIC

- We can use the DEL trick to give semantics using neural network updates

Definition (Neural Network Semantics)

Let \mathcal{N} be a binary neural network model, $w \in N$, and let

$\mathcal{U} : \mathbf{Net} \rightarrow \mathcal{L} \rightarrow \mathbf{Net}$ be any unsupervised update:

$$\mathcal{N}, w \models [\varphi]\psi \quad \text{iff} \quad \text{Update}(\mathcal{N}, \llbracket \varphi \rrbracket), w \models \psi$$

For Hebbian updates in particular:

$$\mathcal{N}, w \models [\varphi]_{\text{HEBB}}\psi \quad \text{iff} \quad \text{HEBB}(\mathcal{N}, \llbracket \varphi \rrbracket), w \models \psi$$

$$\mathcal{N}, w \models [\varphi]_{\text{HEBB}^*}\psi \quad \text{iff} \quad \text{HEBB}^*(\mathcal{N}, \llbracket \varphi \rrbracket), w \models \psi$$

REDUCTION LAWS FOR HEBB^{*}

The following formulas are valid over neural network models:

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

Hebb*-Propositions $[\varphi]_{\text{HEBB}^*} p \leftrightarrow p$

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

Hebb*-Propositions $[\varphi]_{\text{HEBB}^*} p \leftrightarrow p$

Hebb*-Negation $[\varphi]_{\text{HEBB}^*} \neg \psi \leftrightarrow \neg [\varphi]_{\text{HEBB}^*} \psi$

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

Hebb*-Propositions $[\varphi]_{\text{Hebb}^*} p \leftrightarrow p$

Hebb*-Negation $[\varphi]_{\text{Hebb}^*} \neg \psi \leftrightarrow \neg [\varphi]_{\text{Hebb}^*} \psi$

Hebb*-Conjunction $[\varphi]_{\text{Hebb}^*} (\psi \wedge \theta) \leftrightarrow [\varphi]_{\text{Hebb}^*} \psi \wedge [\varphi]_{\text{Hebb}^*} \theta$

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

Hebb*-Propositions $[\varphi]_{\text{Hebb}^*} p \leftrightarrow p$

Hebb*-Negation $[\varphi]_{\text{Hebb}^*} \neg \psi \leftrightarrow \neg [\varphi]_{\text{Hebb}^*} \psi$

Hebb*-Conjunction $[\varphi]_{\text{Hebb}^*} (\psi \wedge \theta) \leftrightarrow [\varphi]_{\text{Hebb}^*} \psi \wedge [\varphi]_{\text{Hebb}^*} \theta$

Hebb*-Diamond $[\varphi]_{\text{Hebb}^*} \Diamond \psi \leftrightarrow \Diamond [\varphi]_{\text{Hebb}^*} \psi$

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

Hebb*-Propositions $[\varphi]_{\text{Hebb}^*} p \leftrightarrow p$

Hebb*-Negation $[\varphi]_{\text{Hebb}^*} \neg \psi \leftrightarrow \neg [\varphi]_{\text{Hebb}^*} \psi$

Hebb*-Conjunction $[\varphi]_{\text{Hebb}^*} (\psi \wedge \theta) \leftrightarrow [\varphi]_{\text{Hebb}^*} \psi \wedge [\varphi]_{\text{Hebb}^*} \theta$

Hebb*-Diamond $[\varphi]_{\text{Hebb}^*} \Diamond \psi \leftrightarrow \Diamond [\varphi]_{\text{Hebb}^*} \psi$

Hebb*-Closure $[\varphi]_{\text{Hebb}^*} \langle \mathbf{C} \rangle \psi \leftrightarrow$
 $\langle \mathbf{C} \rangle ([\varphi]_{\text{Hebb}^*} \psi \vee (\langle \mathbf{C} \rangle \varphi \wedge \Diamond (\langle \mathbf{C} \rangle \varphi \wedge \langle \mathbf{C} \rangle [\varphi]_{\text{Hebb}^*} \psi)))$

REDUCTION LAWS FOR HEBB*

The following formulas are valid over neural network models:

$$\mathbf{Hebb}^*\text{-Propositions } [\varphi]_{\text{HEBB}^*} p \leftrightarrow p$$

$$\mathbf{Hebb}^*\text{-Negation } [\varphi]_{\text{HEBB}^*} \neg \psi \leftrightarrow \neg [\varphi]_{\text{HEBB}^*} \psi$$

$$\mathbf{Hebb}^*\text{-Conjunction } [\varphi]_{\text{HEBB}^*} (\psi \wedge \theta) \leftrightarrow [\varphi]_{\text{HEBB}^*} \psi \wedge [\varphi]_{\text{HEBB}^*} \theta$$

$$\mathbf{Hebb}^*\text{-Diamond } [\varphi]_{\text{HEBB}^*} \Diamond \psi \leftrightarrow \Diamond [\varphi]_{\text{HEBB}^*} \psi$$

$$\mathbf{Hebb}^*\text{-Closure } [\varphi]_{\text{HEBB}^*} \langle \mathbf{C} \rangle \psi \leftrightarrow \\ \langle \mathbf{C} \rangle ([\varphi]_{\text{HEBB}^*} \psi \vee (\langle \mathbf{C} \rangle \varphi \wedge \Diamond (\langle \mathbf{C} \rangle \varphi \wedge \langle \mathbf{C} \rangle [\varphi]_{\text{HEBB}^*} \psi)))$$

Note: We didn't define \Diamond for neural networks; here are its semantics:

$$\mathcal{N}, w \models \Diamond \varphi \quad \text{iff} \quad \text{there is an } E\text{-path from some } u \in \llbracket \varphi \rrbracket \text{ to } w.$$

REDUCTION LAWS FOR HEBB*: INTUITION

Let's look at this last law:

Hebb*-Closure. $[\varphi]_{\text{Hebb}^*} \langle \mathbf{C} \rangle \psi \leftrightarrow$

$$\langle \mathbf{C} \rangle ([\varphi]_{\text{Hebb}^*} \psi \vee (\langle \mathbf{C} \rangle \varphi \wedge \diamond(\langle \mathbf{C} \rangle \varphi \wedge \langle \mathbf{C} \rangle [\varphi]_{\text{Hebb}^*} \psi)))$$

What does it mean? Why does it hold?

REDUCTION LAWS FOR HEBB*: INTUITION

Let's look at this last law:

Hebb*-Closure. $[\varphi]_{\text{Hebb}^*} \langle \mathbf{C} \rangle \psi \leftrightarrow$

$$\langle \mathbf{C} \rangle ([\varphi]_{\text{Hebb}^*} \psi \vee (\langle \mathbf{C} \rangle \varphi \wedge \diamond(\langle \mathbf{C} \rangle \varphi \wedge \langle \mathbf{C} \rangle [\varphi]_{\text{Hebb}^*} \psi)))$$

What does it mean? Why does it hold?

Again, look at what this means for the propositional case:

REDUCTION LAWS FOR HEBB*: INTUITION

Let's look at this last law:

Hebb*-Closure. $[\varphi]_{\text{Hebb}^*} \langle \mathbf{C} \rangle \psi \leftrightarrow$

$$\langle \mathbf{C} \rangle ([\varphi]_{\text{Hebb}^*} \psi \vee (\langle \mathbf{C} \rangle \varphi \wedge \diamond(\langle \mathbf{C} \rangle \varphi \wedge \langle \mathbf{C} \rangle [\varphi]_{\text{Hebb}^*} \psi)))$$

What does it mean? Why does it hold?

Again, look at what this means for the propositional case:

$$[\textcolor{violet}{p}]_{\text{Hebb}^*} \langle \mathbf{C} \rangle \textcolor{teal}{q} \leftrightarrow \langle \mathbf{C} \rangle (\textcolor{teal}{q} \vee (\langle \textcolor{brown}{C} \rangle \textcolor{brown}{p} \wedge \diamond(\langle \textcolor{brown}{C} \rangle \textcolor{brown}{p} \wedge \langle \textcolor{brown}{C} \rangle \textcolor{brown}{q})))$$

REDUCTION LAWS FOR HEBB*: INTUITION

$$[p]_{\text{Hebb}^*} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond(\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

REDUCTION LAWS FOR HEBB^{*}: INTUITION

$$[p]_{\text{Hebb}^*} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb^* has on the closure Clos . This is the description:

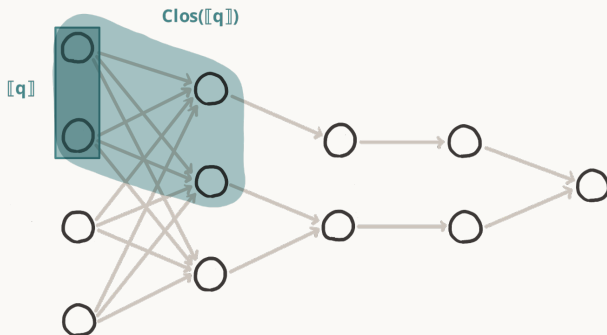
$$\text{Clos}_{\text{Hebb}^*}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$

REDUCTION LAWS FOR HEBB[★]: INTUITION

$$[p]_{\text{Hebb}^\star} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb[★] has on the closure Clos. This is the description:

$$\text{Clos}_{\text{Hebb}^\star}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$

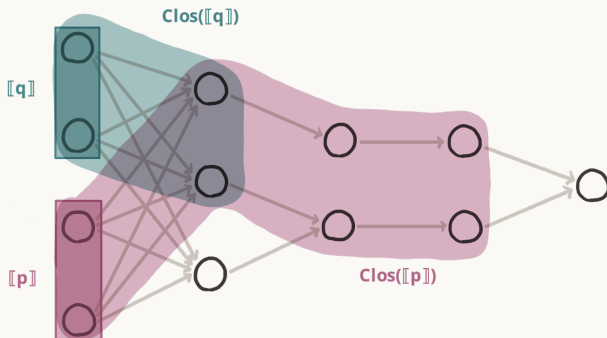


REDUCTION LAWS FOR HEBB[★]: INTUITION

$$[p]_{\text{Hebb}^\star} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb^\star has on the closure Clos . This is the description:

$$\text{Clos}_{\text{Hebb}^\star}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$

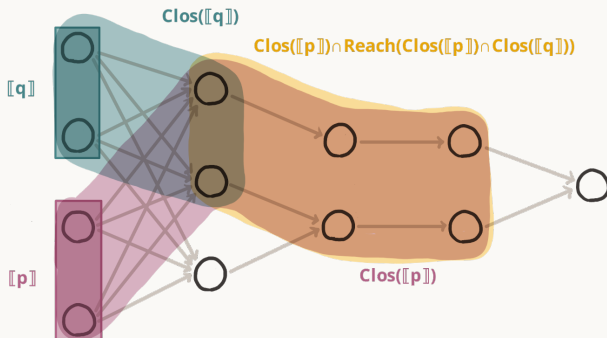


REDUCTION LAWS FOR HEBB*: INTUITION

$$[p]_{\text{Hebb}^*} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb^* has on the closure Clos . This is the description:

$$\text{Clos}_{\text{Hebb}^*}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$

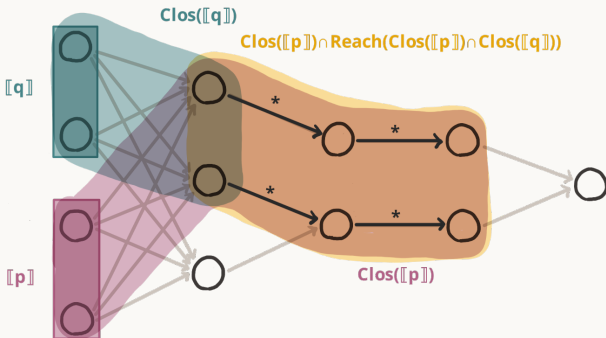


REDUCTION LAWS FOR HEBB*: INTUITION

$$[p]_{\text{Hebb}^*} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb^* has on the closure Clos . This is the description:

$$\text{Clos}_{\text{Hebb}^*}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$

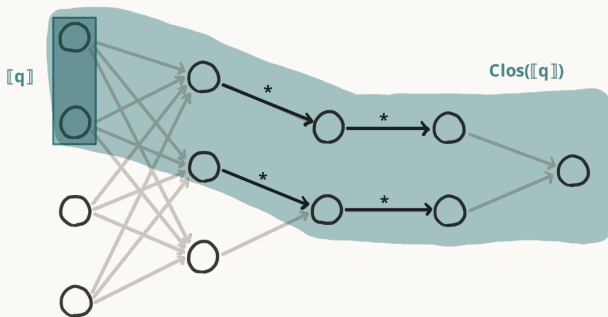


REDUCTION LAWS FOR HEBB*: INTUITION

$$[p]_{\text{Hebb}^*} \langle \mathbf{C} \rangle q \leftrightarrow \langle \mathbf{C} \rangle (q \vee (\langle \mathbf{C} \rangle p \wedge \diamond (\langle \mathbf{C} \rangle p \wedge \langle \mathbf{C} \rangle q)))$$

- This encodes in our logic a complete description of the effect Hebb^* has on the closure Clos . This is the description:

$$\text{Clos}_{\text{Hebb}^*}(\mathcal{N}, [p])([q]) = \text{Clos}([q] \cup (\text{Clos}([p]) \cap \text{Reach}(\text{Clos}([p]) \cap \text{Clos}([q])))$$



TO LEARN MORE CHECK IT...

Neural Network Models of Conditionals: An Introduction

Hannes Leitgeb

Department of Philosophy, University of Bristol,
Hannes.Leitgeb@bristol.ac.uk

Abstract

The “biconditional style” article gives a brief survey of neural network models of conditionals. After short introductions into the studies of neural networks and conditionals, we turn to the notion of an *incomplete classical system* as a working concept in the logical investigation of dynamic systems in general, and of neural networks in particular. The explicit *how* conditionals are represented by temporal dynamical systems, which logical forms these conditionals obey, and what the main open problems in the area are.

Keywords: Conditionals, Neural networks, Dynamical systems, Nonmonotonic logic.

1 Introduction

Neural networks are abstract models of brain structures capable of adapting to new information. The learning abilities of artificial neural networks have given rise to successful computer implementations of various cognitive tasks, from the recognition of facial images to the prediction of currency movement.

Logic deals with formal systems of reasoning in particular, including logic studies formal systems of reasoning towards plausible but uncertain conclusions. As evidence accumulates, the degree to which it supports a hypothesis, is measured by the logic, should tend to indicate that the hypothesis is likely to be true.

Although during a great focus on informatics and reasoning, until recently these two areas developed in opposition to each other: neural networks are quantitative dynamic systems, while logical systems are symbolic systems, networks are described by mathematical equations, whereas logic is subject to normative statements about how we ought to reason, neural networks have been studied by scientists, while the “norms of induction” is regarded as belonging to the humanities.

At present, the interdisciplinary research in the area of conditionals, for example, is increasing and the discovery that these formalisms apply in neural net processes on the representational level give rise to the question that the dynamics of artificial neural networks can be understood in terms of logically valid, and thus rational, rules of inference. As neural networks, connections reasoning, and symbolic

1

The New English-Old Conference on Artificial Intelligence (2024)

What do Hebbian Learners Learn? Reduction Axioms for Iterated Hebbian Learning

Calvin Schultz Kisby¹, Saul A. Blanco², Lawrence S. Moss²

¹Department of Computer Science, Indiana University
²Department of Mathematics, Indiana University
Bloomington, IN 47404, USA
{calvin, olafsson, lawson}@indiana.edu

Abstract

This paper is a contribution to neural network semantics, a formalism for the study of neural networks. It is the first in a series of papers in the logical systems can be mapped to operations on neural network states. In this paper, we do this for a form of network learning operation. We study a dynamic system of iterated Hebbian learning, a single learning policy that updates a neural network by repeatedly applying Hebb's learning rule over the same training data points. The model is a form of “modular logic” for learning, or reduction axioms. This means that complexity for the logic of iterated Hebbian learning follows from complexity of the base logic. These reduction axioms also provide a human interpretable description of iterated Hebbian learning as a kind of abstracting upgrade, and (2) an approach to finding neural networks with guarantees on what they can learn.

Introduction

The two dominant paradigms of AI, connectionist neural networks and symbolic systems, have long seemed irreconcilable. Symbolic systems are well suited for giving explicit instructions to a human-interpretable language, but are brittle and fail to adapt to new situations. On the other hand, neural networks are flexible and robust at learning from unstructured data, but are considered black boxes due to their difficulty to do so without data mining. In response to this dichotomy, the field of neural-symbolic AI has emerged — a community-wide effort to integrate neural and symbolic systems, while retaining the advantages of both. Despite the many different proposals for neural-symbolic AI, few seem to “fit” the field and the field (Bard et al., 2017; Sarker et al., 2022). There is little agreement on what the interface between the two might be. There is a clear need for a working theory that can explain the relationship between neural networks and symbolic systems (Bianco et al., 2022). In fact, there is an open-ended foundational theory for neural-symbolic systems, which we call *neural network axioms*. Its key insight is that neural networks can be taken

The central question this theory aims to answer are: **Satisfaction**. What axioms are needed for neural network operations? Can neural operations be mapped to classical ones in a useful way? Note that checking conditions is equivalent to **formally verifying** properties of nets.

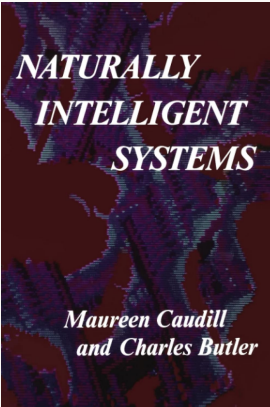
Complexities. What are the complexity axioms for neural network operations? This is captured in **model building**. Can we build a neural network that obeys a set of logical constraints? Can we build a neural network from a classical model?

We write the reader to the landmark survey (Olsson and d'Ávila Garcia 2022), which shows that this framework encompasses a wide class of neural-symbolic systems. We will discuss other work that we consider part of the core theory in the next section.

The standard example is the *forward propagation operator* P_{forward} and \mathcal{L} , which measures to take β sequentially across new neurons and eventually the state of the net \mathcal{N} (Bianco et al., 2022). P_{forward} returns the state at the final point. A classic result from (Bianco 2003) is that, for any conditionals

$\mathcal{L} \models \varphi \iff \mathcal{L} \models P_{\text{forward}}(\varphi)$
i.e., φ is satisfied by input φ iff “the net, when run on φ ,” satisfies the other conditionals \mathcal{L} . This is completely characterized by the long-conditional condition laws of (Kahn, Leifer, and Margalef 1998). The result is robust and can be extended to different choices of conditional axioms and neural network architectures (Gagliardi 2006). This general theory is that formal propagation corresponds to a non-monotonic condition.

A central challenge for this theory is to do the same for *neural network learning operations*, the present work (Schultz, Blanco, and Moss 2022) considers a simple learning policy — iterated Hebbian learning (Hebbian that “fire together wire together”) — on a binary, feed-forward net. Although this work often used axioms for Hebbian learning, the question of completeness is left open.



(a) **Neural Network Models of Conditionals: An Introduction** by Leitgeb

(b) **Reduction Axioms for Iterated Hebbian Learning** by Schultz Kisby, Blanco, & Moss

(c) **Naturally Intelligent Systems** by Caudill & Butler

END OF LECTURE 3

Thank you!