

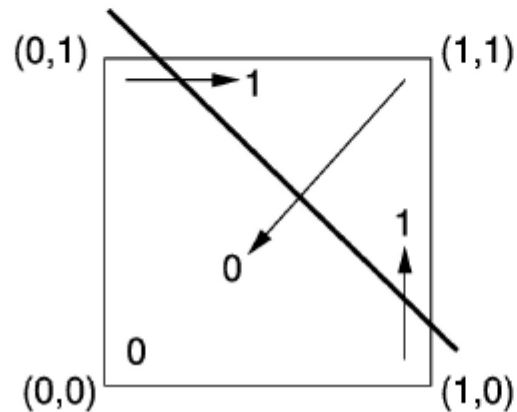
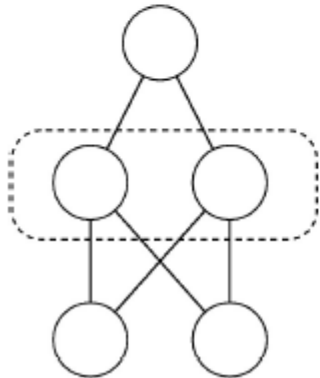
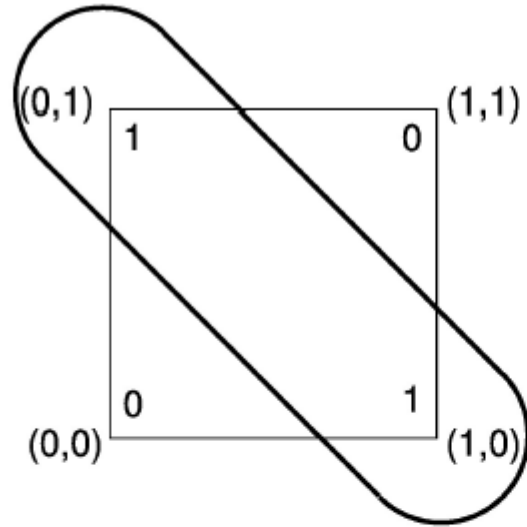
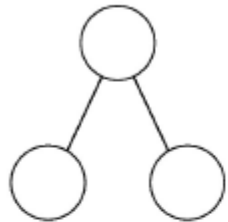
Introduction to Parallel Distributed Processing models

Deon T. Benton
004-Backpropagation

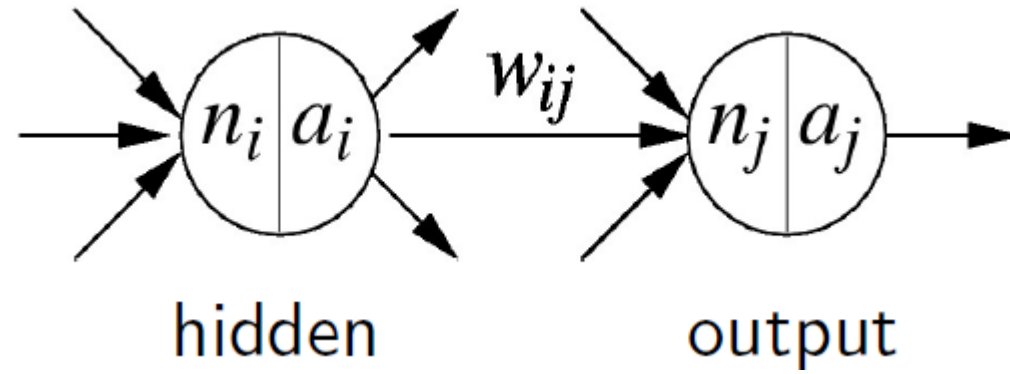
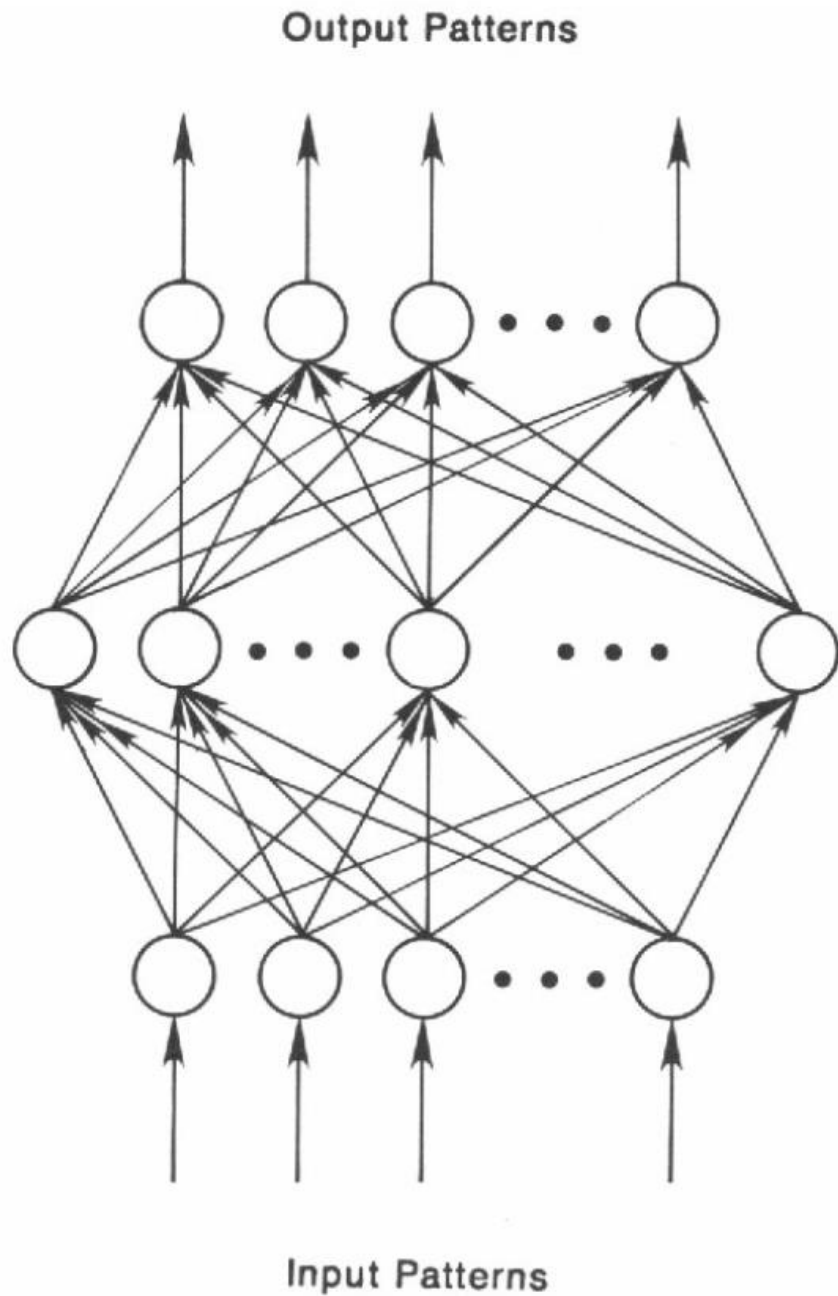
Pre-lecture quiz

- The Delta rule is guaranteed to succeed under which conditions?
- **True or False?** Linear independence and orthogonality mean the same thing; that is, they are identical.
- Under what conditions are the Delta rule and the Hebb rule equivalent?

XOR with intermediate (“hidden”) units

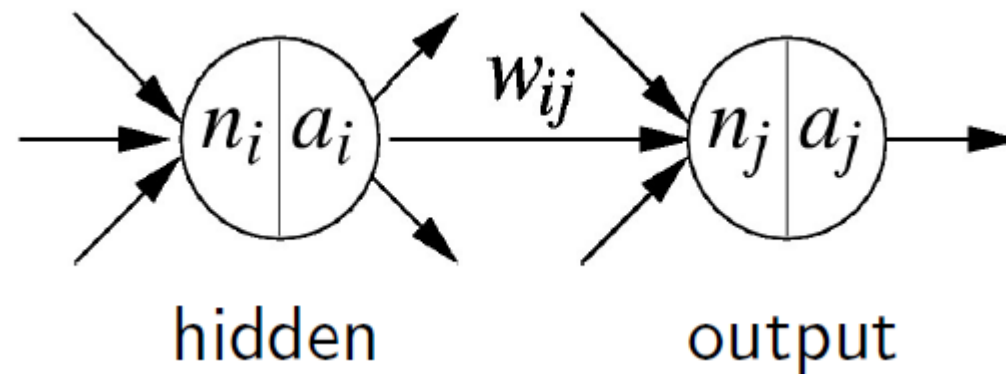
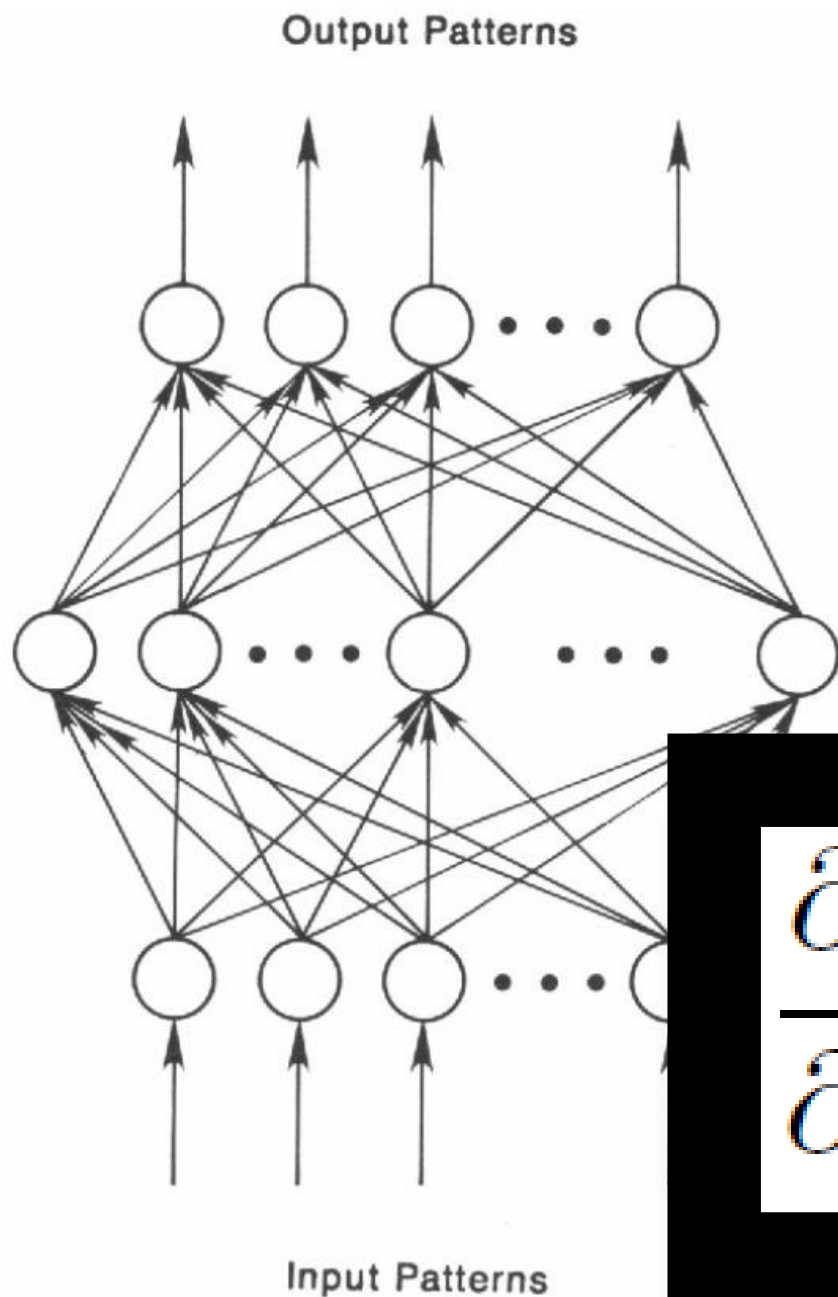


- Intermediate units can re-represent input patterns as new patterns with **altered similarities**
- Targets which are not linearly separable in the input space can be linearly separable in the intermediate representational space
- Intermediate units are called “hidden” because their activations are not determined directly by the training environment (inputs and targets)



- Hidden-to-output weights can be trained with the Delta rule

Why??



- Hidden-to-output weights can be trained with the Delta rule

$$\frac{\partial E}{\partial a_i}$$

can we train input-to-hidden weights?

Hidden units do not have targets (for determining error)

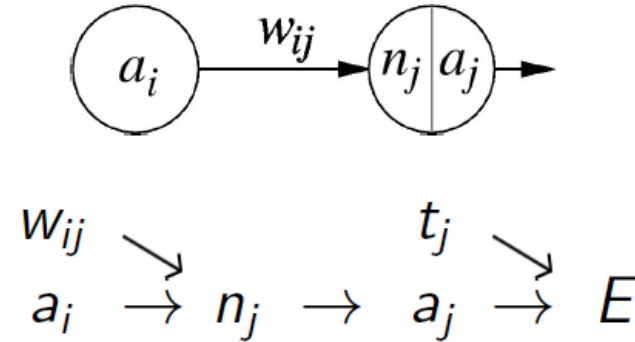
Trick: We don't need targets, we just need to know how hidden activations affect error (i.e., error derivatives)

Delta rule as gradient descent in error (sigmoid units)

$$n_j = \sum_i a_i w_{ij}$$

$$a_j = \frac{1}{1 + \exp(-n_j)}$$

$$\text{Error } E = \frac{1}{2} \sum_j (t_j - a_j)^2$$



Gradient descent: $\Delta w_{ij} = -\epsilon \frac{\partial E}{\partial w_{ij}}$

$$\begin{aligned} \frac{\partial E}{\partial w_{ij}} &= \frac{\partial E}{\partial a_j} \frac{da_j}{dn_j} \frac{\partial n_j}{\partial w_{ij}} \\ &= -(t_j - a_j) \quad a_j (1 - a_j) \quad a_i \end{aligned}$$

$$\Delta w_{ij} = -\epsilon \frac{\partial E}{\partial w_{ij}} = \epsilon (t_j - a_j) a_j (1 - a_j) a_i$$

Quick aside: The derivative of the sigmoid function

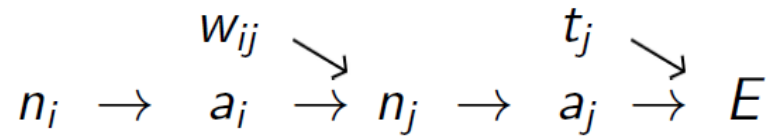
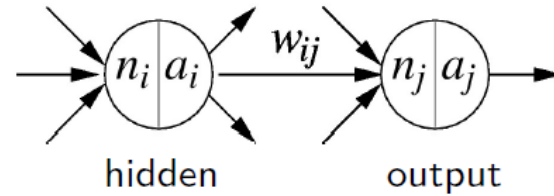
$$\begin{aligned}\frac{d}{dn_j}(a_j) &= \frac{d}{dn_j} \left(\frac{1}{1 + e^{-n_j}} \right) = \frac{1}{(1 + e^{-n_j})} \cdot \frac{(1 + e^{-n_j}) - 1}{(1 + e^{-n_j})} \\ &= \frac{d}{dn_j} (1 + e^{-n_j})^{-1} = \frac{1}{(1 + e^{-n_j})} \cdot \left(\frac{1 + e^{-n_j}}{1 + e^{-n_j}} + \frac{1}{1 + e^{-n_j}} \right) \\ &= -(1 + e^{-n_j})^{-2} (-e^{-n_j}) = a_j(1 - a_j) \\ &= \frac{e^{-n_j}}{(1 + e^{-n_j})^2}\end{aligned}$$

Generalized Delta rule (“backpropagation”)

$$n_j = \sum_i a_i w_{ij}$$

$$a_j = \frac{1}{1 + \exp(-n_j)}$$

$$\text{Error } E = \frac{1}{2} \sum_j (t_j - a_j)^2$$



Gradient descent: $\Delta w_{ij} = -\epsilon \frac{\partial E}{\partial w_{ij}}$

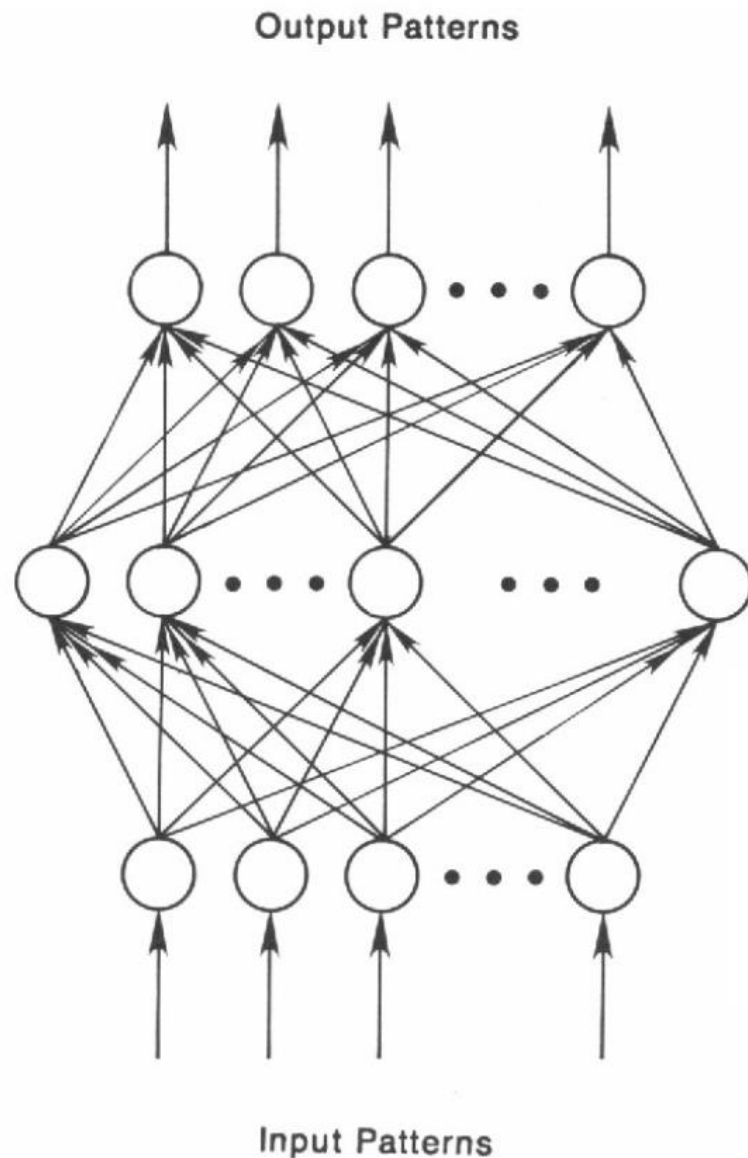
Intermediate notation
("input derivatives" in Lens)

$$\frac{\partial E}{\partial n_j} = \boxed{\frac{\partial E}{\partial a_j}} \frac{da_j}{dn_j} = -(t_j - a_j) a_j (1 - a_j)$$

$$\frac{\partial E}{\partial w_{ij}} = \overline{\frac{\partial E}{\partial n_j}} \frac{\partial n_j}{\partial w_{ij}} = \frac{\partial E}{\partial n_j} a_i$$

$$\frac{\partial E}{\partial a_i} = \sum_j \frac{\partial E}{\partial n_j} \frac{\partial n_j}{\partial a_i} = \sum_j \frac{\partial E}{\partial n_j} w_{ij}$$

Generalized Delta rule (“backpropagation”)



= computed from the top (hidden-to-output) layer
 = computed from the bottom (input-to-hidden) layer

$$\frac{\partial E}{\partial w_{ijL}} = \frac{\frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial n_j} \frac{\partial n_j}{\partial w_{ijL}}}{\frac{\partial a_j}{\partial n_j} \frac{\partial n_j}{\partial w_{ijL}}} = \frac{\frac{\partial E}{\partial n_j} \frac{\partial n_j}{\partial w_{ijL}}}{\frac{\partial n_j}{\partial w_{ijL}}} = \frac{\partial E}{\partial n_j} a_i$$

$$\frac{\partial E}{\partial w_{ijL-1}} = \frac{\frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial n_j} \frac{\partial n_j}{\partial a_i} \frac{\partial a_i}{\partial n_i} \frac{\partial n_i}{\partial w_{ijL-1}}}{\frac{\partial a_j}{\partial n_j} \frac{\partial n_j}{\partial a_i} \frac{\partial a_i}{\partial n_i} \frac{\partial n_i}{\partial w_{ijL-1}}} = \frac{\frac{\partial E}{\partial n_j} \frac{\partial n_j}{\partial a_i} \frac{\partial a_i}{\partial n_i} \frac{\partial n_i}{\partial w_{ijL-1}}}{\frac{\partial n_j}{\partial a_i} \frac{\partial a_i}{\partial n_i} \frac{\partial n_i}{\partial w_{ijL-1}}} = \frac{\frac{\partial E}{\partial n_j} \frac{\partial n_j}{\partial w_{ijL-1}}}{\frac{\partial n_j}{\partial w_{ijL-1}}}$$

Sigmoid derivative₁ points to $\frac{\partial a_j}{\partial n_j}$
 Sigmoid derivative₂ points to $\frac{\partial n_j}{\partial a_i}$
 Sigmoid derivative₂ points to $\frac{\partial a_i}{\partial n_i}$

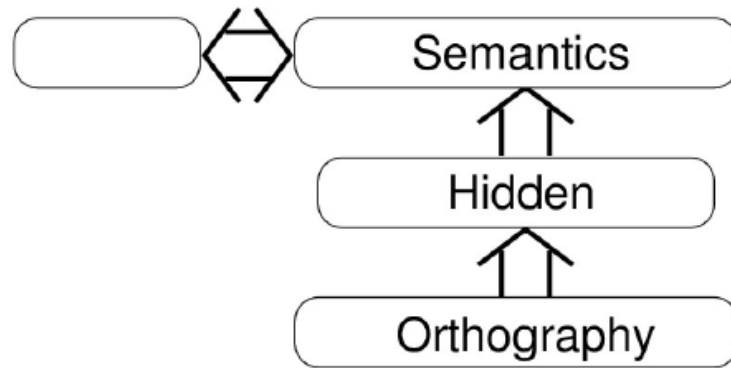
Backpropagation *in words*

Backpropagation is an algorithm that involves:

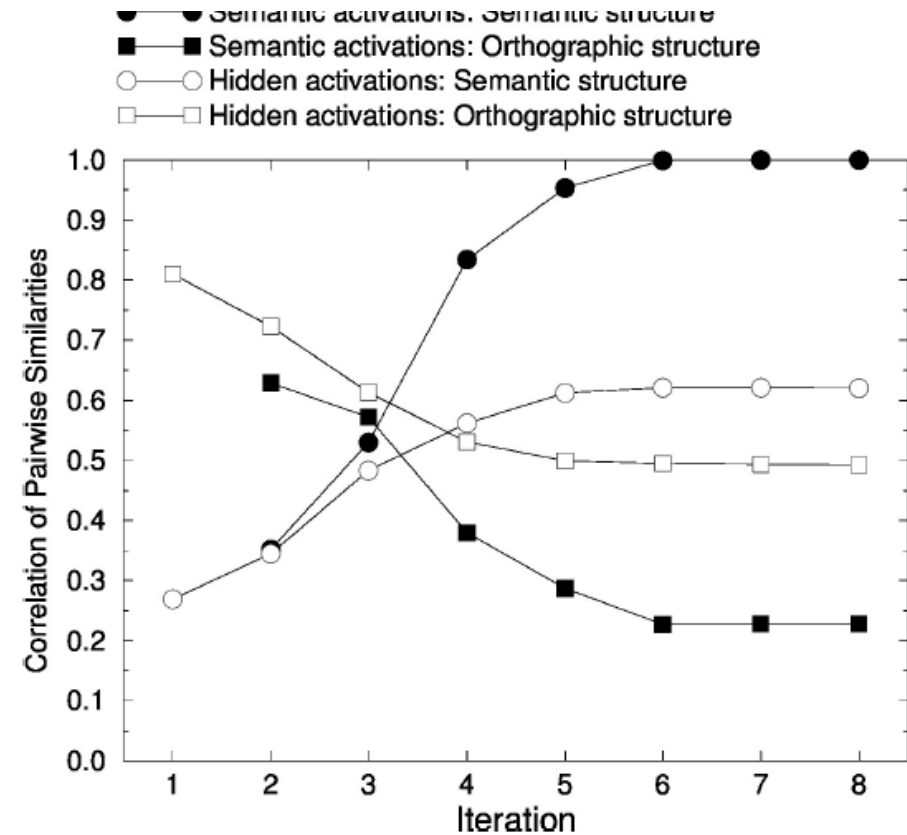
1. Computing net inputs and activations *in a forward pass*.
2. Computing total network error (over all training examples).
3. Adjusting network weights based on the total error *in the backward pass*.
4. Repeat steps 1-3 until the network converges.

What do hidden representations learn

Plaut and Shallice (1993)

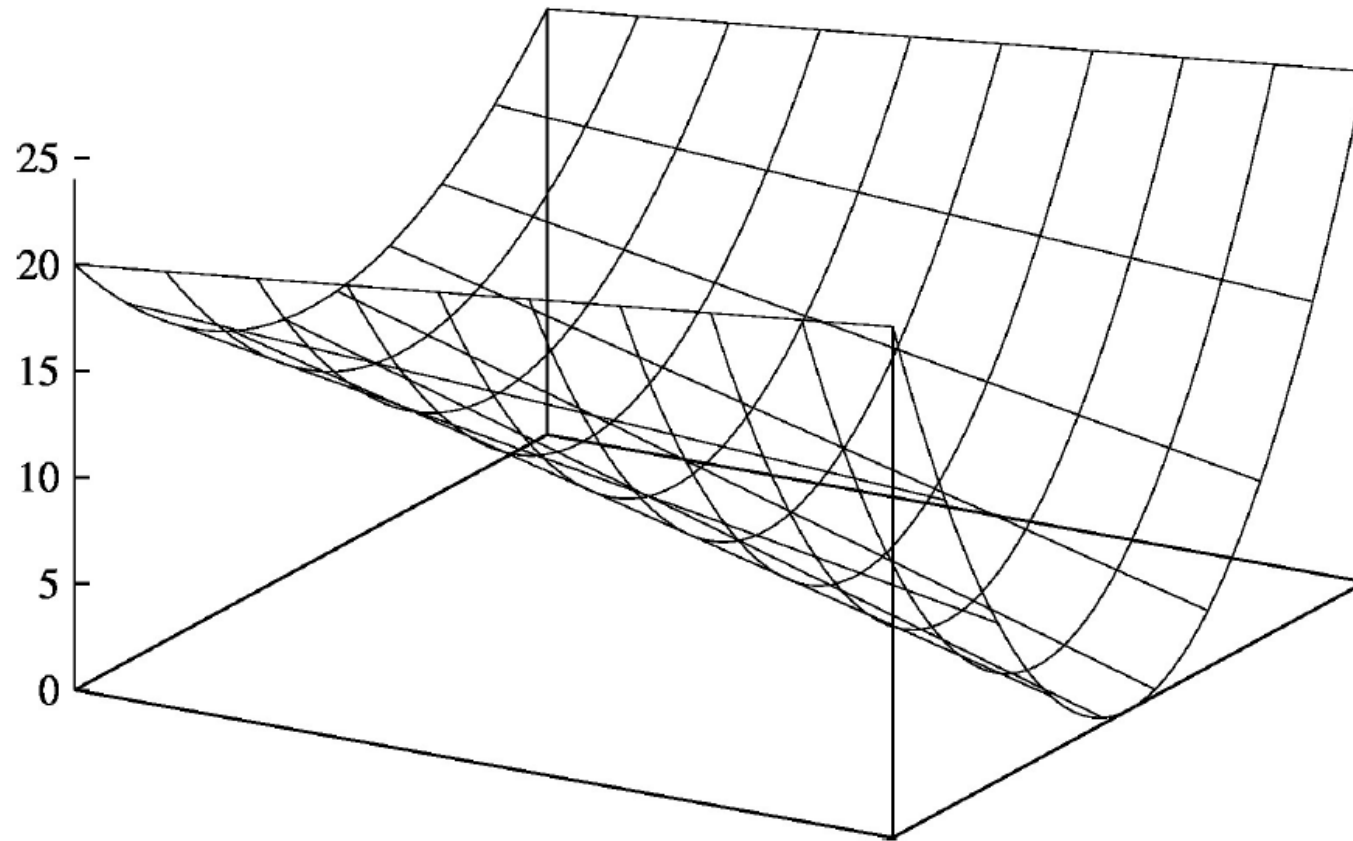


- Mapped orthography to semantics (unrelated similarities)
- Compared similarities among hidden representations to those among orthographic and semantic representations (over settling)



- Hidden representations “**split the difference**” between input and output similarity

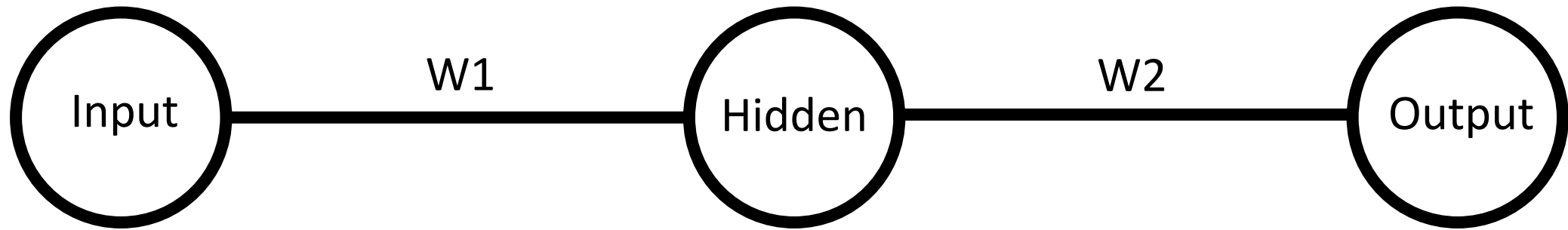
Accelerating learning: Momentum descent



$$\Delta w_{ij}^{[t]} = -\epsilon \frac{\partial E}{\partial w_{ij}} + \alpha \left(\Delta w_{ij}^{[t-1]} \right)$$

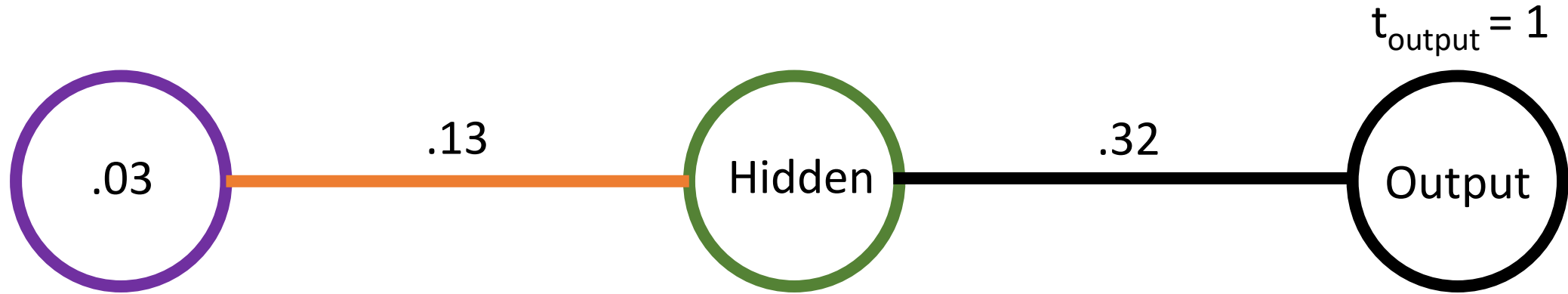
Backpropagation: A really simple example

*The **forward** pass*



Backpropagation: A really simple example

The forward pass

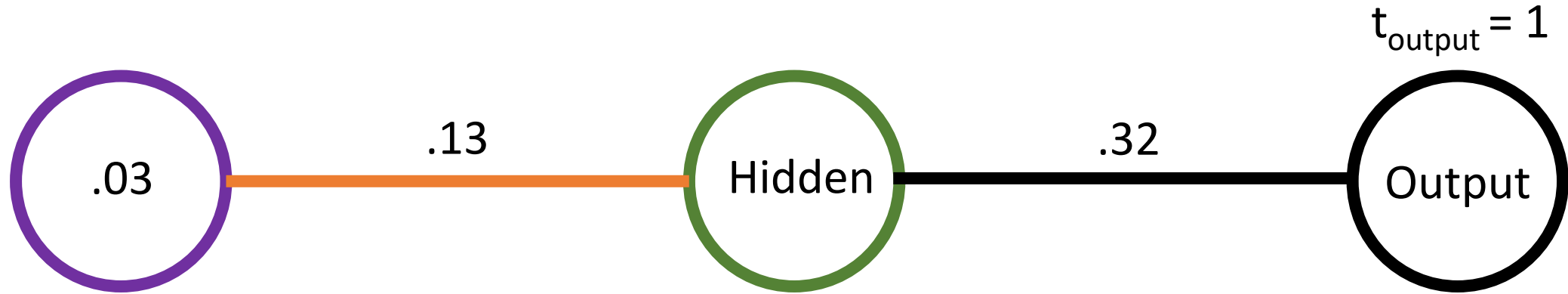


$$\text{net}_{\text{hidden}} = a_{\text{input}} * w_{\text{input} * \text{hidden}}$$

$$a_{\text{hidden}} = \frac{1}{1 + e^{-\text{net}_{\text{hidden}}}}$$

Backpropagation: A really simple example

The forward pass

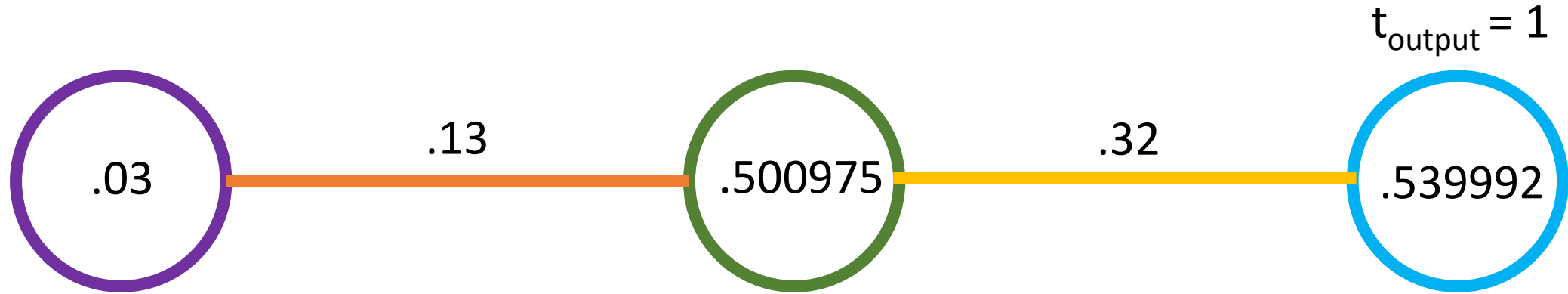


$$\text{net}_{\text{hidden}} = .03 * .13 = .0039$$

$$a_{\text{hidden}} = \frac{1}{1 + e^{-.0039}} \sim .500975$$

Backpropagation: A really simple example

The forward pass

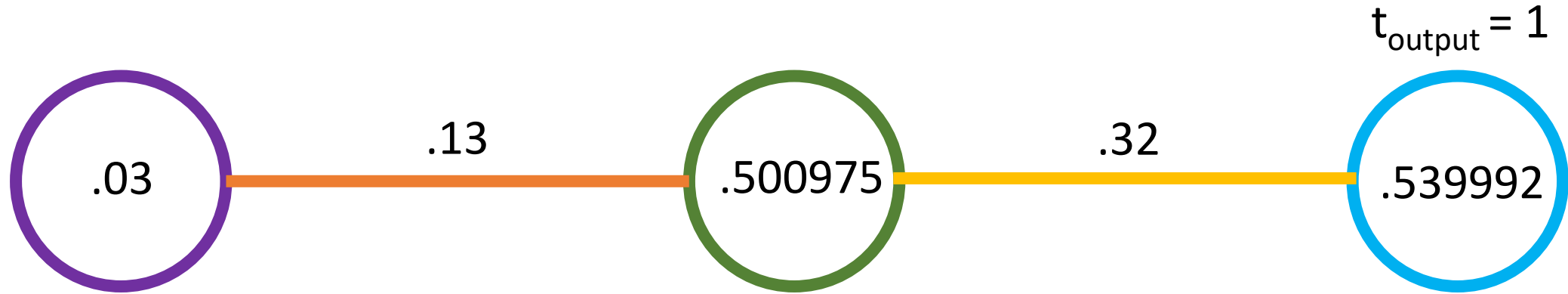


$$\text{net}_{\text{output}} = .500975 * .32 = .160312$$

$$a_{\text{output}} = \frac{1}{1 + e^{-.160312}} \sim .539992$$

Backpropagation: A really simple example

The forward pass

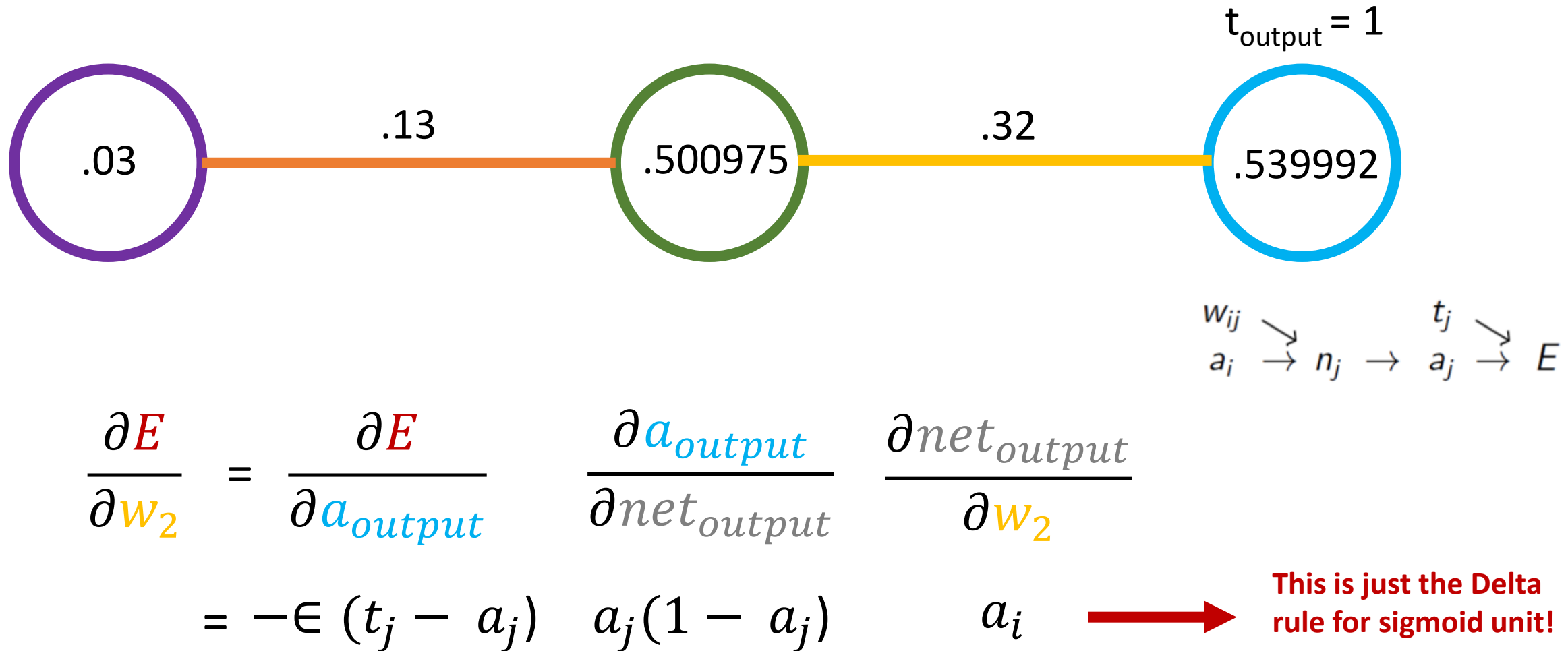


Network error $= \frac{1}{2}(t_j - a_j)^2$

Network error $= \frac{1}{2}(1 - .539992)^2 = .105804$

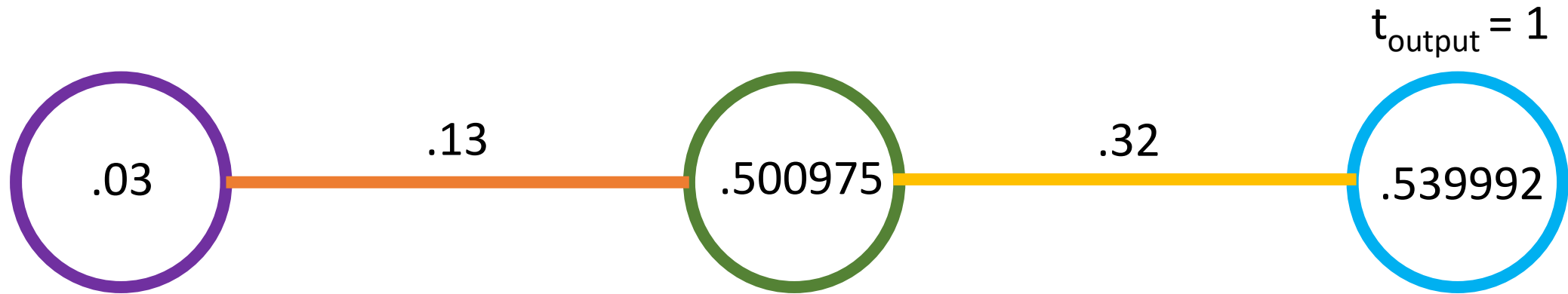
Backpropagation: A really simple example

The backwards pass



Backpropagation: A really simple example

The backwards pass



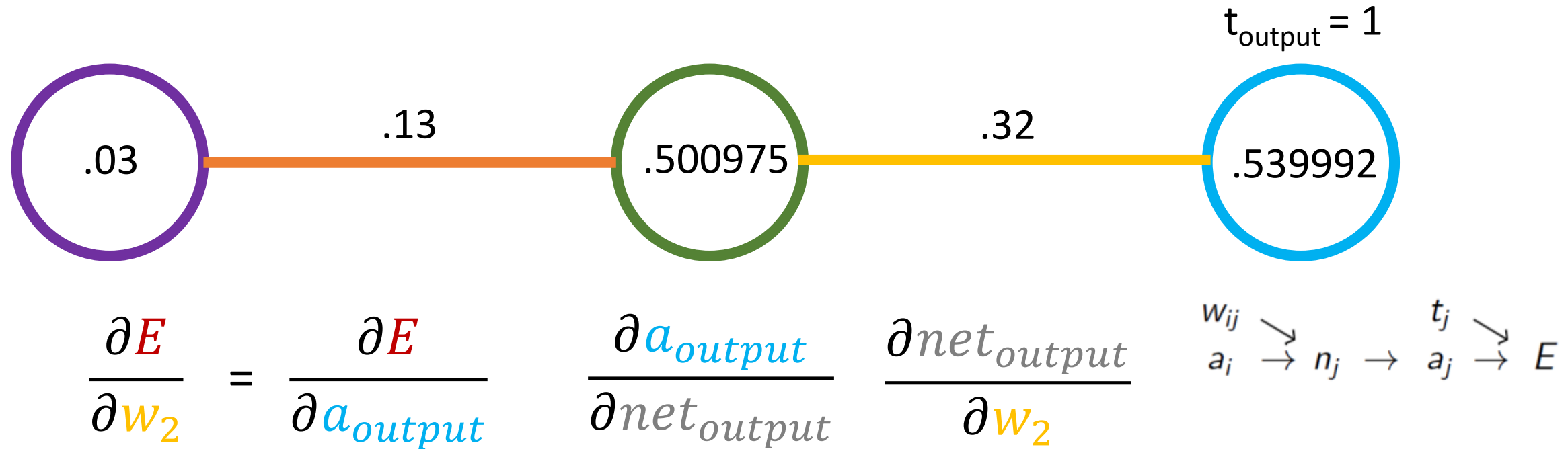
$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial a_{\text{output}}} \frac{\partial a_{\text{output}}}{\partial \text{net}_{\text{output}}} \frac{\partial \text{net}_{\text{output}}}{\partial w_2}$$

$\begin{matrix} w_{ij} & \searrow \\ a_i & \rightarrow n_j \end{matrix} \rightarrow \begin{matrix} t_j & \searrow \\ a_j & \rightarrow E \end{matrix}$

$$\begin{aligned}
 &= (1 - .539992) .539992 (1 - .539992) .500975 \\
 &= .460008 .248401 .500975 \\
 &= .057245
 \end{aligned}$$

Backpropagation: A really simple example

The backwards pass

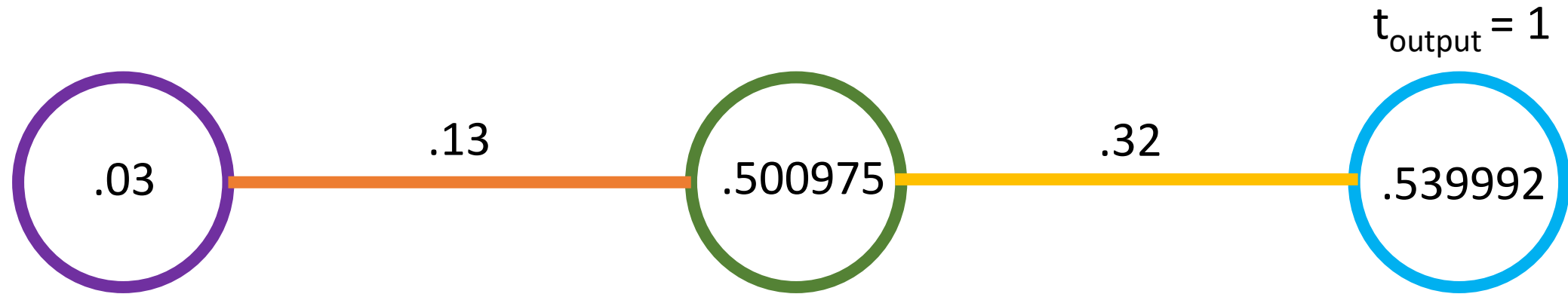


$$\Delta w_2 = .32 - 0.5 * .057245$$

$$\text{New_}w_2 = .3486225$$

Backpropagation: A really simple example

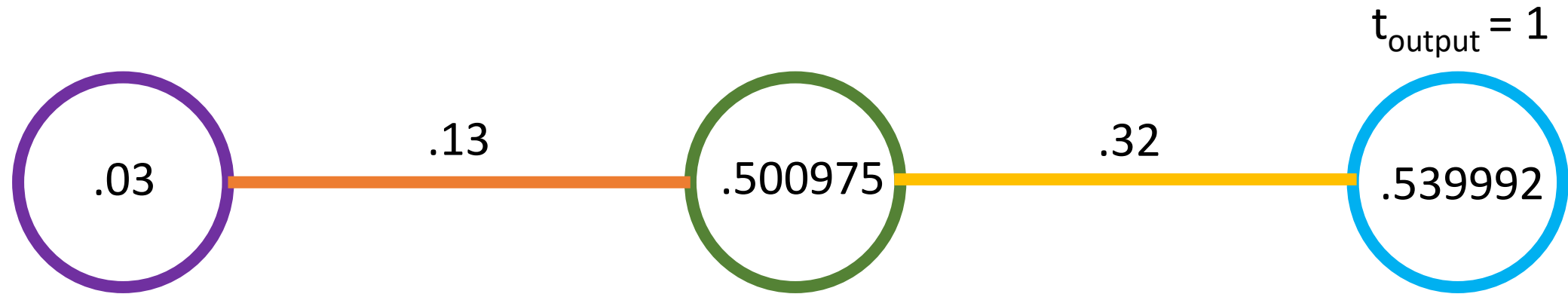
The backwards pass



$$\begin{aligned}
 \frac{\partial E}{\partial w_1} &= \frac{\partial E}{\partial a_{\text{output}}} \frac{\partial a_{\text{output}}}{\partial \text{net}_{\text{output}}} \frac{\partial \text{net}_{\text{output}}}{\partial a_{\text{hidden}}} \frac{\partial a_{\text{hidden}}}{\partial \text{net}_{\text{hidden}}} \frac{\partial \text{net}_{\text{hidden}}}{\partial w_1} \\
 &= -\epsilon (t_j - a_j) a_j (1 - a_j) w_{ij} a_k (1 - a_k) a_{ki} \\
 &= (1 - .539992) .539992 (1 - .539992) .32 .500975 (1 - .500975) .03 \\
 &= .460008 * .248401 * .32 * .249999 * .03 = .000274
 \end{aligned}$$

Backpropagation: A really simple example

The backwards pass



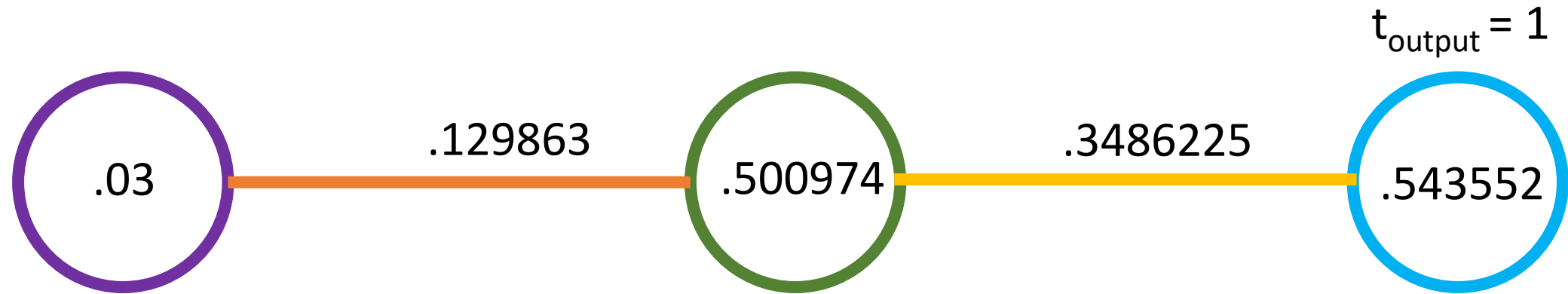
$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial a_{\text{output}}} \frac{\partial a_{\text{output}}}{\partial \text{net}_{\text{output}}} \frac{\partial \text{net}_{\text{output}}}{\partial a_{\text{hidden}}} \frac{\partial a_{\text{hidden}}}{\partial \text{net}_{\text{hidden}}} \frac{\partial \text{net}_{\text{hidden}}}{\partial w_1}$$

$$\Delta w_1 = .13 - 0.5 * .000274$$

$$\text{New_}w_1 = .129863$$

Backpropagation: A really simple example

Update the weights!



We're doing better, albeit slightly! Keep going!