

# **Backpropagation Learning, Internal Representations and Semantics**

# Learning

---

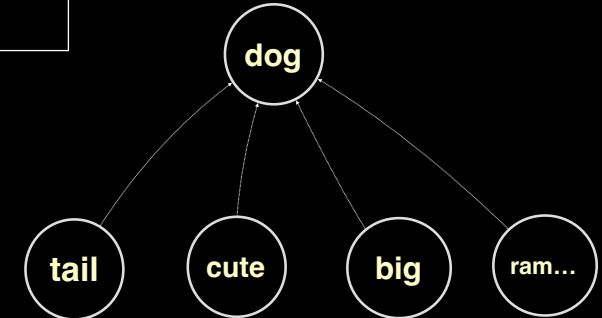
- **Unsupervised Learning**
  - Hebbian Learning Rule
  - Pattern associator
  - Self-organized maps
  - Topographic structure
  - Pattern detectors
- **Supervised Learning**
  - **Scalar Learning**
    - Classical and Instrumental Conditioning
    - Sequential learning and Prediction
  - **Vector-Based Learning**
    - **Generalized Delta Rule**
    - **Backpropagation**
    - **Semantics (Deep Learning)**



# Simple Pattern Associator

Pure Correlational learning:

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-



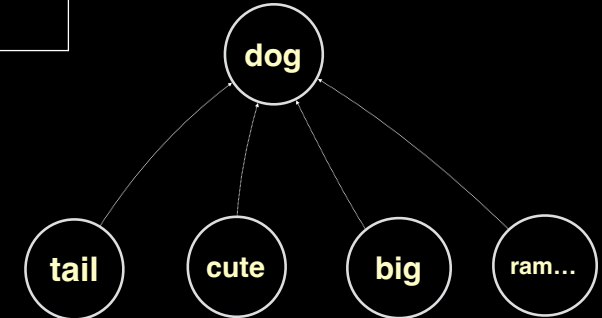
# Simple Pattern Associator

Pure Correlational learning:

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

- **Observe:**

- *dog* is uncorrelated with *tail*, *cute*, and *rambunctious*



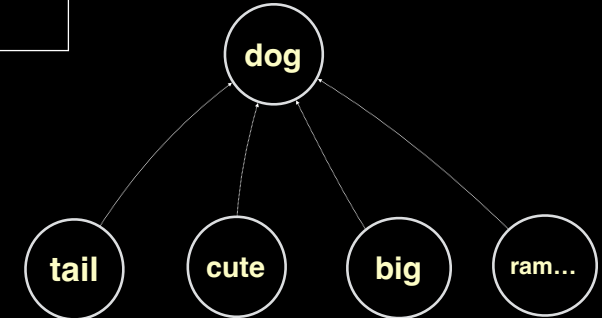
# Simple Pattern Associator

Pure Correlational learning:

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

- **Observe:**

- *dog* is uncorrelated with *tail*, *cute*, and *rambunctious*
- $\therefore$  weights from these units to *dog* will be 0



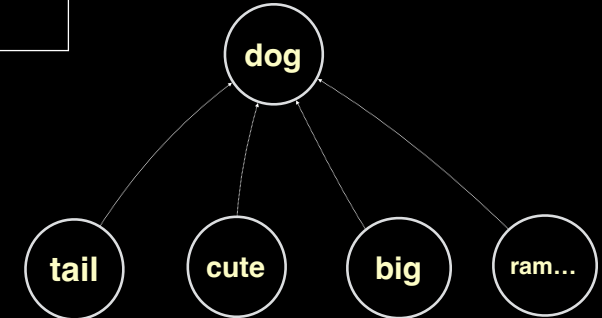
# Simple Pattern Associator

Pure Correlational learning:

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

- **Observe:**

- *dog* is uncorrelated with *tail*, *cute*, and *rambunctious*
- $\therefore$  weights from these units to *dog* will be 0
- *dog* is positively correlated with *big*



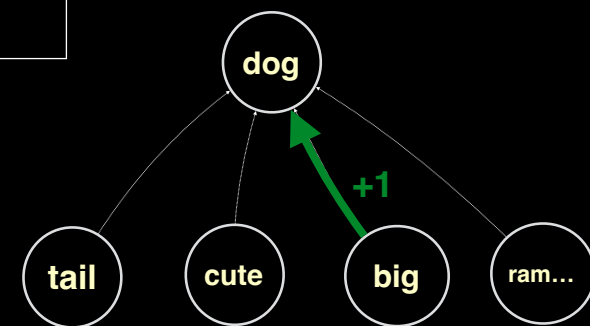
# Simple Pattern Associator

Pure Correlational learning:

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

- **Observe:**

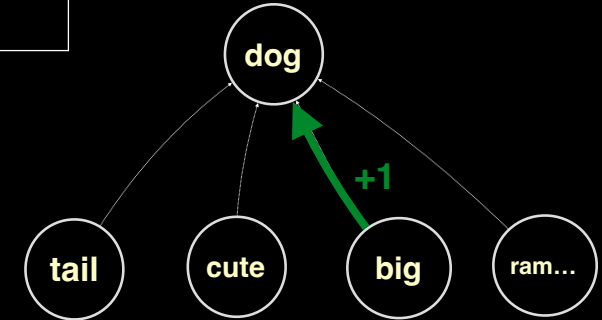
- *dog* is uncorrelated with *tail*, *cute*, and *rambunctious*
- $\therefore$  weights from these units to *dog* will be 0
- *dog* is positively correlated with *big*
- $\therefore$  a positive weight will develop between *big* and *dog*, and *dog* will be positive for *Felix* when it should be negative





# Simple Pattern Associator

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

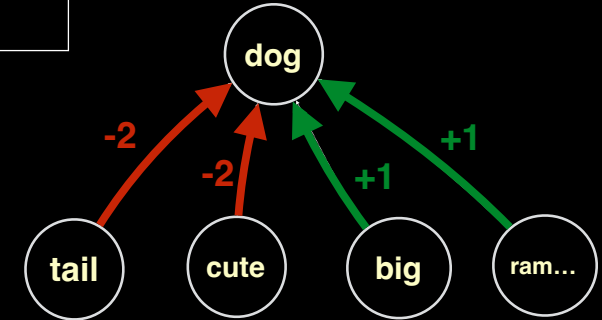
- There is a set of weights that will do the trick





# Simple Pattern Associator

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

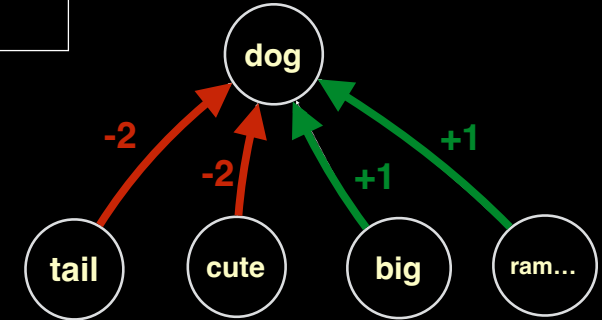
- There is a set of weights that will do the trick



# Simple Pattern Associator

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

- There is a set of weights that will do the trick
  - how can these be learned?

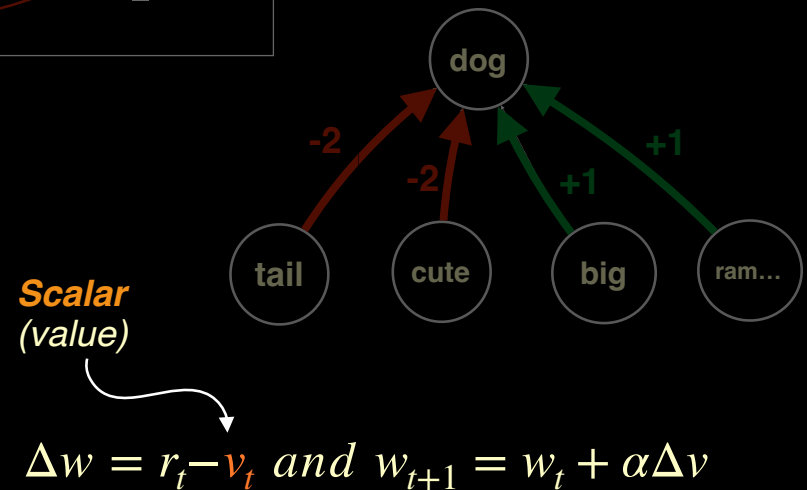


# Simple Pattern Associator

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-

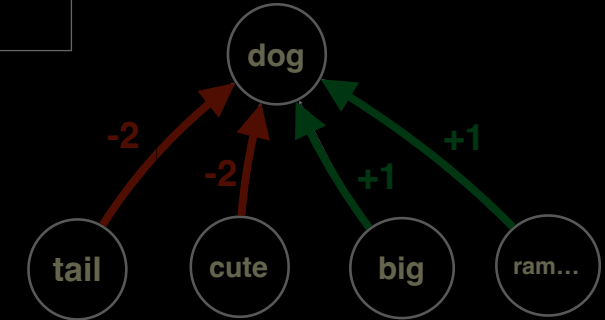
*Note: In the original image, the 'dog' column for Felix and Natasha, and the 'big' column for Felix, are circled. Green arrows point from these circles to the 'dog' column for Felix. Red arrows point from the 'tail' and 'cute' columns for Felix to the 'dog' column for Natasha.*

- There is a set of weights that will do the trick
  - how can these be learned?
- Prediction (error-driven) learning:
  - Reinforcement learning (*Rescorla-Wagner*):



# Simple Pattern Associator

		<i>tail</i>	<i>cute</i>	<i>big</i>	<i>rambunctious</i>	<i>dog</i>
	<i>Fido</i>	+	-	+	-	+
	<i>Rover</i>	+	+	+	+	+
	<i>Felix</i>	+	+	+	-	-
	<i>Natasha</i>	+	-	-	+	-



- There is a set of weights that will do the trick
  - how can these be learned?
- Prediction (error-driven) learning:
  - Reinforcement learning (*Rescorla-Wagner*):
  - Delta rule (*Widrow-Hoff*):

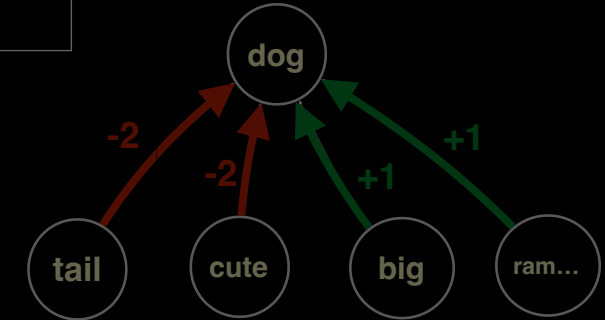
$$\Delta w = r_t - v_t \text{ and } w_{t+1} = w_t + \alpha \Delta v$$

$$\Delta w = r_t - a_t \text{ and } w_{t+1} = w_t + \alpha \Delta v$$

**Vector**  
(activity)

# Simple Pattern Associator

		tail	cute	big	rambunctious	dog
	Fido	+	-	+	-	+
	Rover	+	+	+	+	+
	Felix	+	+	+	-	-
	Natasha	+	-	-	+	-



- There is a set of weights that will do the trick
  - how can these be learned?

- Prediction (error-driven) learning:

- Reinforcement learning (*Rescorla-Wagner*):  $\Delta w = r_t - v_t$  and  $w_{t+1} = w_t + \alpha \Delta v$
- Delta rule (*Widrow-Hoff*):  $\Delta w = r_t - a_t$  and  $w_{t+1} = w_t + \alpha \Delta v$

- Can be extended to train many output units...

**Vector**  
(activity)

# Delta Rule

---

*(also known as the Widrow Hoff or  
Least Mean Square (LMS) learning rule)*

# Delta Rule

---

*(also known as the Widrow Hoff or  
Least Mean Square (LMS) learning rule)*

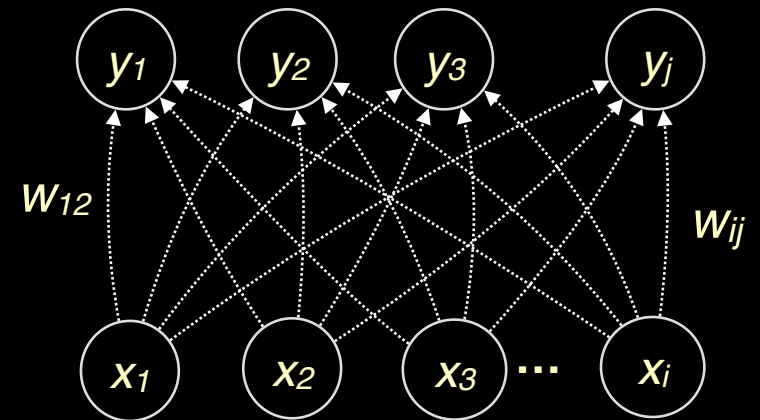
- **Generalization of Rescorla-Wagner to many output units:**

# Delta Rule

*(also known as the Widrow Hoff or  
Least Mean Square (LMS) learning rule)*

- Generalization of Rescorla-Wagner to many output units:

$$\Delta w_{ij} = \epsilon e_j x_j$$



# Delta Rule

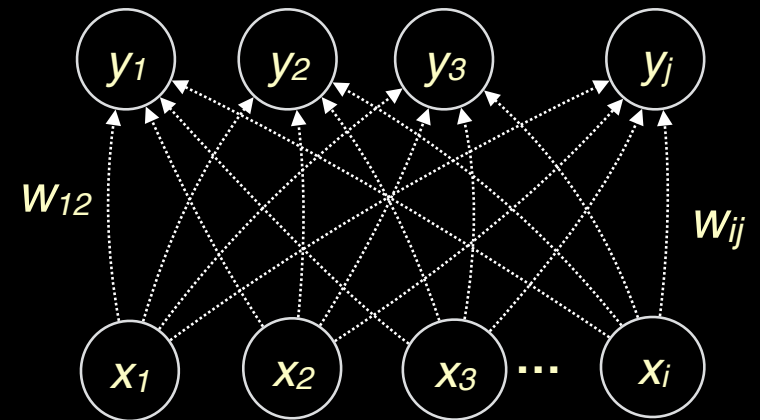
*(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)*

- **Generalization of Rescorla-Wagner to many output units:**

$$\Delta w_{ij} = \epsilon e_j x_j$$

where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$



# Delta Rule

*(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)*

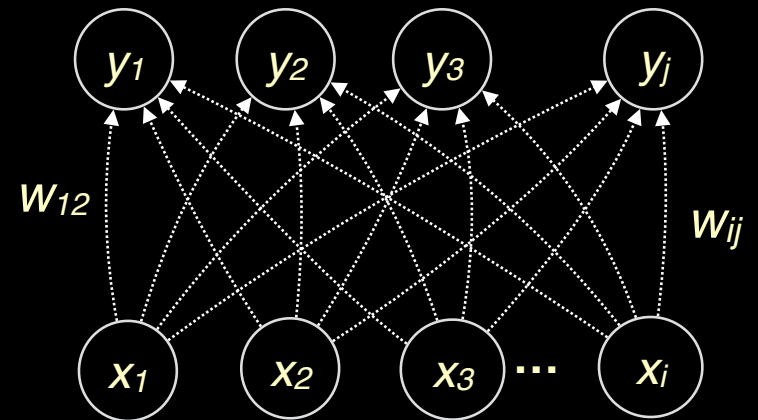
- **Generalization of Rescorla-Wagner to many output units:**

$$\Delta w_{ij} = \epsilon e_j x_j$$

where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$

$t_j$  is the target activation for  $y_j$ ,  
(in the training pattern)



# Delta Rule

*(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)*

- **Generalization of Rescorla-Wagner to many output units:**

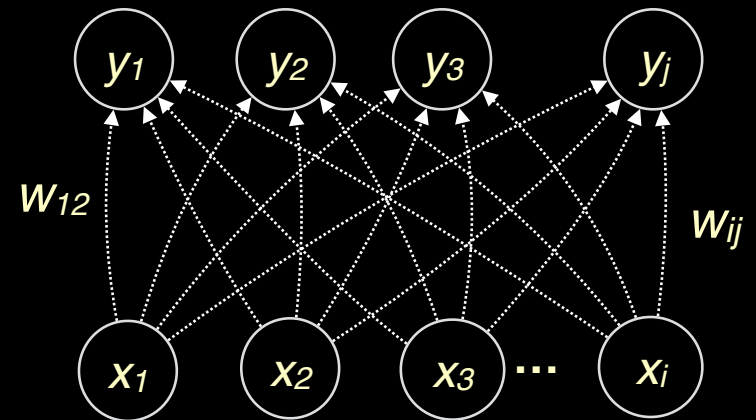
$$\Delta w_{ij} = \epsilon e_j x_j$$

where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$

$t_j$  is the target activation for  $y_j$ ,  
(in the training pattern)

and  $y_j$  is its actual activity



# Delta Rule

(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)

- Generalization of Rescorla-Wagner to many output units:

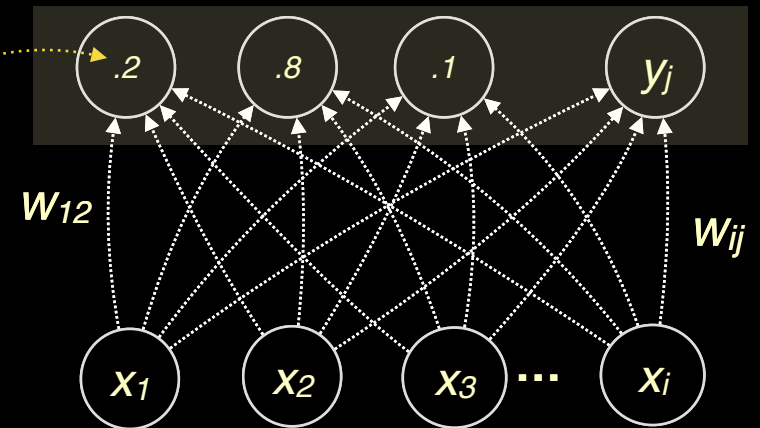
$$\Delta w_{ij} = \epsilon e_j x_j$$

where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$

$t_j$  is the target activation for  $y_j$ ,  
(in the training pattern)

and  $y_j$  is its actual activity



# Delta Rule

(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)

- Generalization of Rescorla-Wagner to many output units:

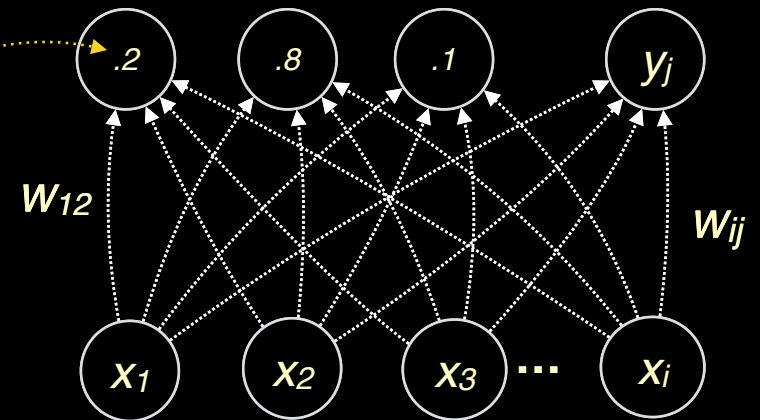
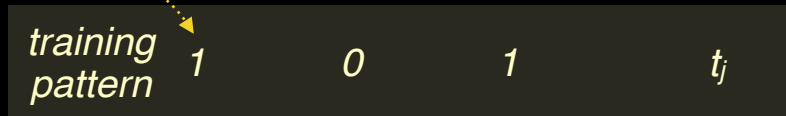
$$\Delta w_{ij} = \epsilon e_j x_j$$

where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$

$t_j$  is the target activation for  $y_j$ ,  
(in the training pattern)

and  $y_j$  is its actual activity



# Delta Rule

(also known as the Widrow Hoff or Least Mean Square (LMS) learning rule)

- Generalization of Rescorla-Wagner to many output units:

$$\Delta w_{ij} = \epsilon e_j x_j$$

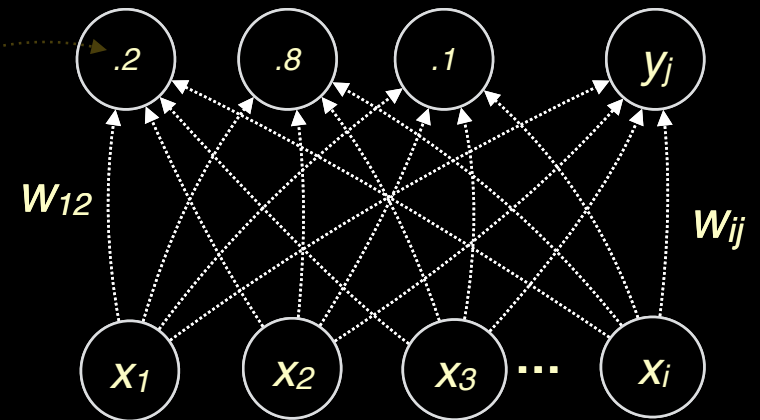
where  $e_j$  is the error for output unit  $j$ :

$$e_j = t_j - y_j$$

$t_j$  is the target activation for  $y_j$ ,  
(in the training pattern)

and  $y_j$  is its actual activity

training pattern	1	0	1	$t_j$
error	.8	-.8	.9	$t_j$



# **Training Procedure**

# **Training Procedure**

---

- **Two phases:**
  1. **Activation**
  2. **Error computation and weight adjustment**

# Training Procedure

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units

# Training Procedure

---

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units
  - Determine the resulting output

# Training Procedure

---

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units
  - Determine the resulting output
  - **Compute error: difference between target and actual output**

# Training Procedure

---

- Two phases:
  1. Activation
  2. Error computation and weight adjustment
- For each pattern:
  - Turn on input units
  - Determine the resulting output
  - Compute error: difference between target and actual output
  - Adjust strength of each weight according to the delta rule

# Training Procedure

---

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units
  - Determine the resulting output
  - Compute error: difference between target and actual output
  - Adjust strength of each weight according to the delta rule
- **Epoch:**
  - One pass of training through all patterns

# Training Procedure

---

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units
  - Determine the resulting output
  - Compute error: difference between target and actual output
  - Adjust strength of each weight according to the delta rule
- **Epoch:**
  - One pass of training through all patterns
- **Performance measures:**
  - *Sum squared error (SSE) over output units  $y_j$  for pattern  $p$  (pss):  $\sum_j (t_j - y_j)^2$*

# Training Procedure

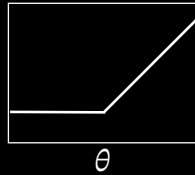
---

- **Two phases:**
  1. Activation
  2. Error computation and weight adjustment
- **For each pattern:**
  - Turn on input units
  - Determine the resulting output
  - Compute error: difference between target and actual output
  - Adjust strength of each weight according to the delta rule
- **Epoch:**
  - One pass of training through all patterns
- **Performance measures:**
  - *Sum squared error (SSE) over output units  $y_j$  for pattern  $p$  (pss):*  $\sum_j (t_j - y_j)^2$
  - *Total sum of squares (tss) over a set of patterns:*  $\sum_p SSE_p$

# The Perceptron

*Frank Rosenblatt, 1957*

---



# The Perceptron

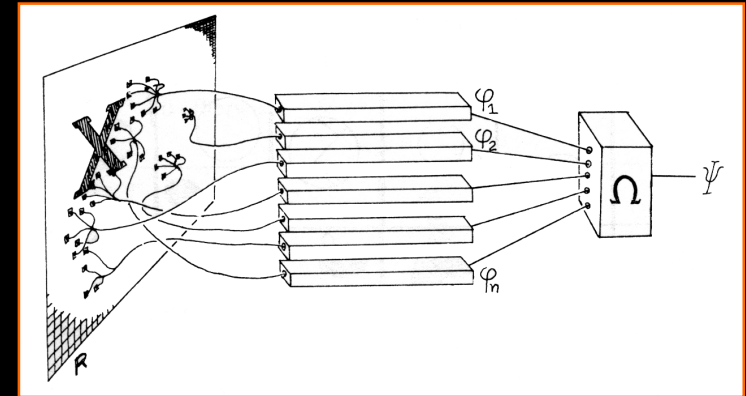
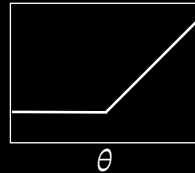
Frank Rosenblatt, 1957

- **Perceptron**

- **binary classifier**

- **linearly thresholded units:**  
*aka Rectified Linear Units (ReLU)*

$$\text{Output} = \begin{cases} 1 & \text{if net} = \sum_i w_i i_i > \theta \\ 0 & \text{otherwise} \end{cases}$$



# The Perceptron

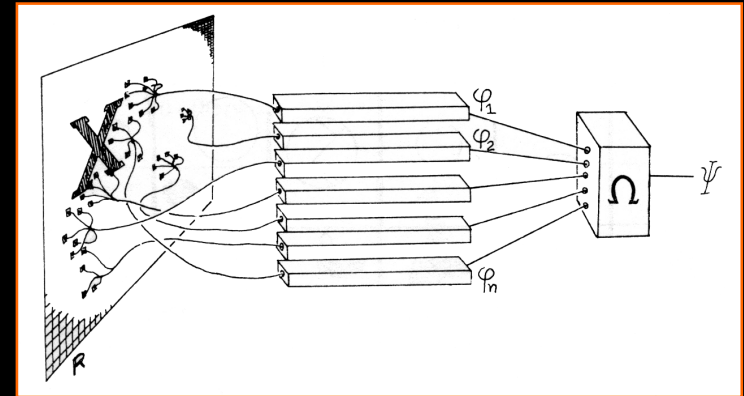
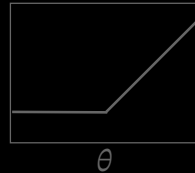
Frank Rosenblatt, 1957

- **Perceptron**

- binary classifier

- linearly thresholded units:  
aka Rectified Linear Units (ReLU)

Output =  $\begin{cases} 1 & \text{if net} = \sum_i w_i i_i > \theta \\ 0 & \text{otherwise} \end{cases}$



- learning rule (similar to the Delta rule):

$$\Delta\theta = -(t_p - o_p) = -\delta_p \text{ (threshold)}$$

$$\Delta w_i = (t_p - o_p) i_{pi} = \delta_p i_p \text{ (weights)}$$

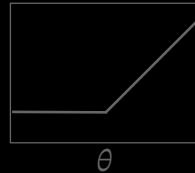
# The Perceptron

Frank Rosenblatt, 1957

- **Perceptron**

- binary classifier
- linearly thresholded units:  
aka Rectified Linear Units (ReLU)

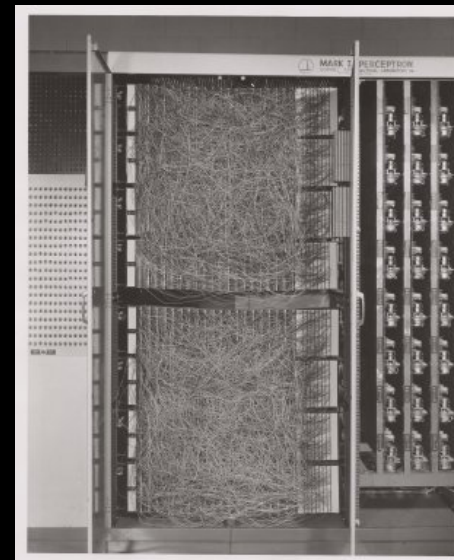
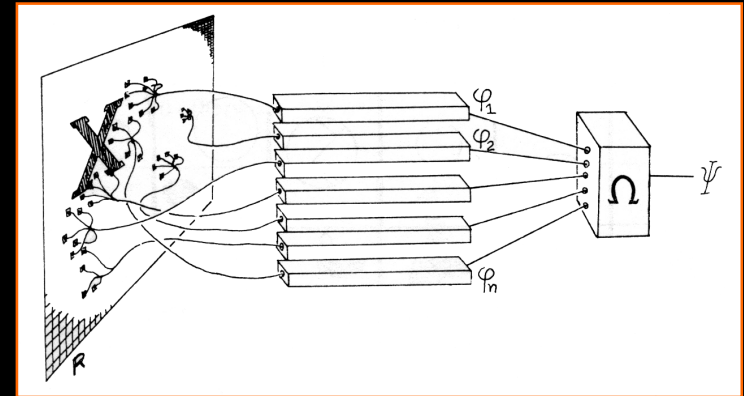
$$\text{Output} = \begin{cases} 1 & \text{if net} = \sum_i w_i i_i > \theta \\ 0 & \text{otherwise} \end{cases}$$



- learning rule (similar to the Delta rule):

$$\Delta\theta = -(t_p - o_p) = -\delta_p \text{ (threshold)}$$

$$\Delta w_i = (t_p - o_p) i_{pi} = \delta_p i_p \text{ (weights)}$$



## Mark 1 Perceptron

Input:  
400 photocells

weights:  
potentiometers

weight updates:  
electric motors

# The Perceptron

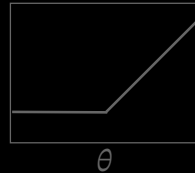
Frank Rosenblatt, 1957

- **Perceptron**

- binary classifier

- linearly thresholded units:  
*aka Rectified Linear Units (ReLU)*

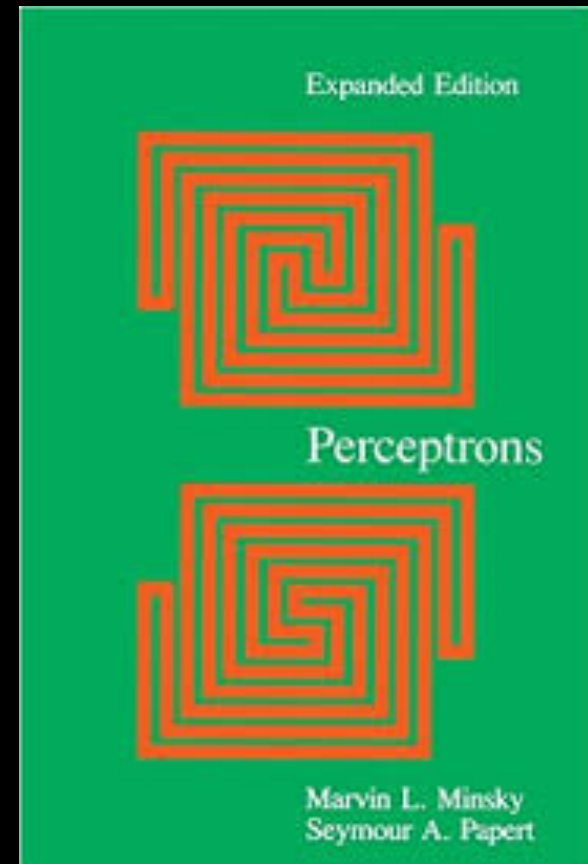
Output =  $\begin{cases} 1 & \text{if net} = \sum_i w_i i_i > \theta \\ 0 & \text{otherwise} \end{cases}$



- learning rule (similar to the Delta rule):

$$\Delta \theta = -(t_p - o_p) = -\delta_p \text{ (threshold)}$$

$$\Delta w_i = (t_p - o_p) i_{pi} = \delta_p i_p \text{ (weights)}$$



# The Perceptron

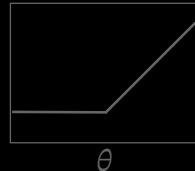
Frank Rosenblatt, 1957

- **Perceptron**

- binary classifier

- linearly thresholded units:  
*aka Rectified Linear Units (ReLU)*

$$\text{Output} = \begin{cases} 1 & \text{if net} = \sum_i w_i i_i > \theta \\ 0 & \text{otherwise} \end{cases}$$



- learning rule (similar to the Delta rule):

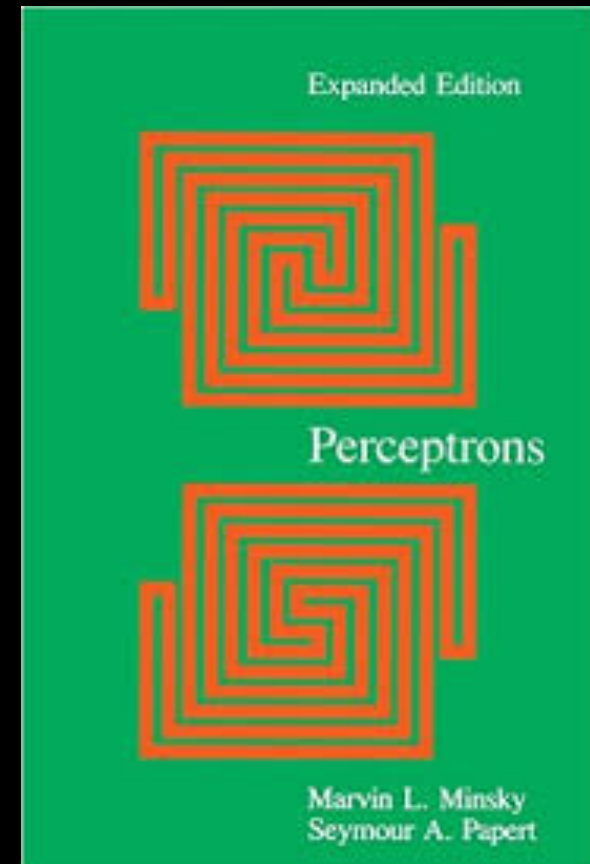
$$\Delta \theta = - (t_p - o_p) = - \delta_p \text{ (threshold)}$$

$$\Delta w_i = (t_p - o_p) i_{pi} = \delta_p i_p \text{ (weights)}$$

- **Minsky & Papert (1969)**

- Perceptrons can't solve the XOR problem

- ∴ not computationally general



# The Problem

---

# The Problem

---

Actually, two problems:

## 1) Additive

- can't solve problems that are not "*linearly separable*" (e.g., XOR)...

# The Problem

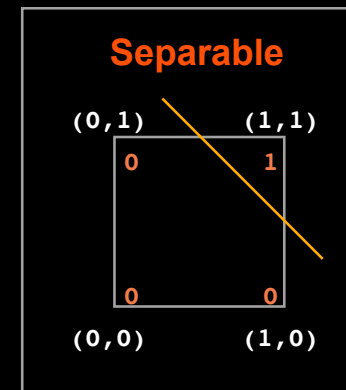
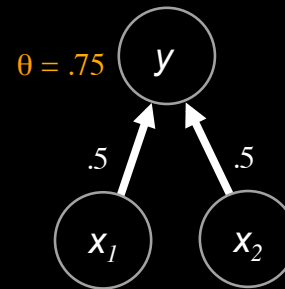
---

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>AND</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	0
1	1	1



# The Problem

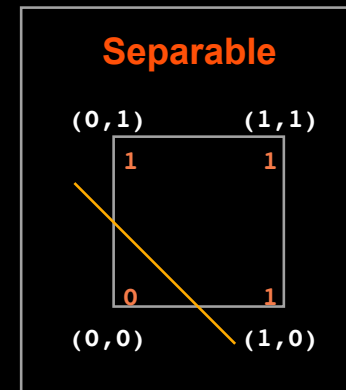
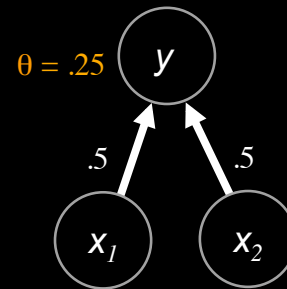
---

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1



# The Problem

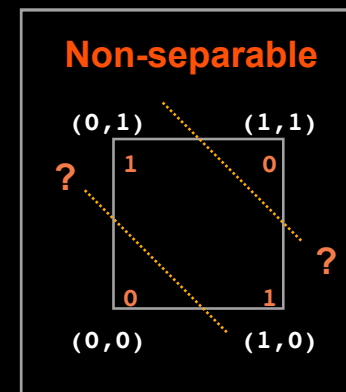
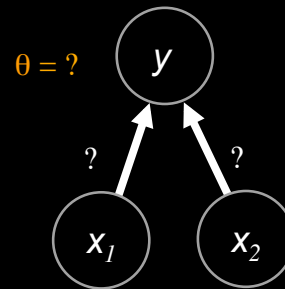
---

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



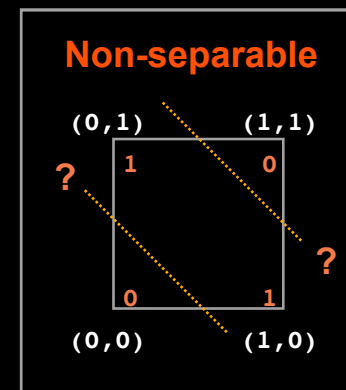
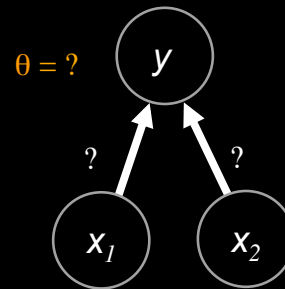
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



**Example:**

no food

⇒ don't like

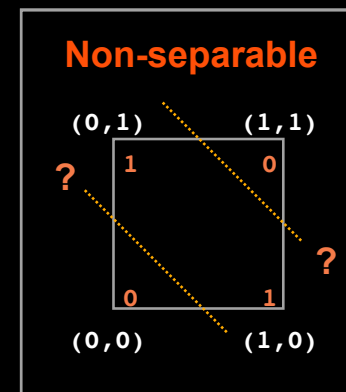
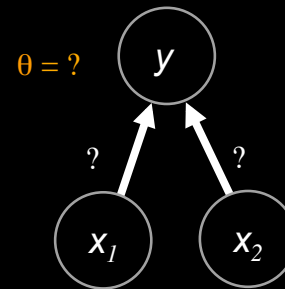
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



**Example:**

no food  
pickles

$\Rightarrow$  don't like  
 $\Rightarrow$  like

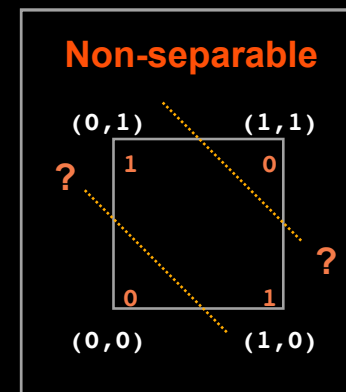
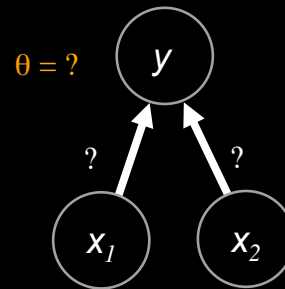
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



**Example:**

no food	⇒ don't like
pickles	⇒ like
ice cream	⇒ like

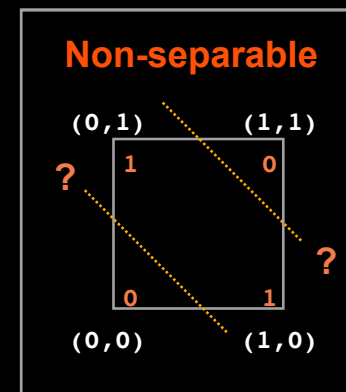
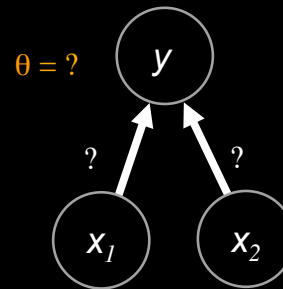
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



### Example:

no food	$\Rightarrow$	don't like
pickles	$\Rightarrow$	like
ice cream	$\Rightarrow$	like
pickles + ice cream	$\Rightarrow$	don't like

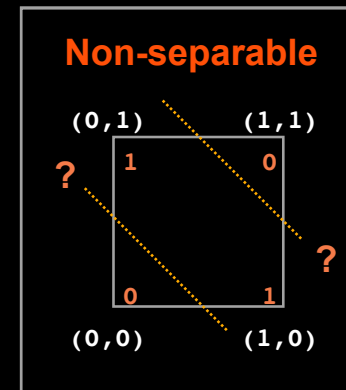
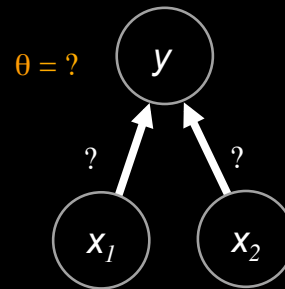
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not "*linearly separable*" (e.g., XOR)...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



### Example:

no food  $\Rightarrow$  don't like  
pickles  $\Rightarrow$  like  
ice cream  $\Rightarrow$  like  
pickles + ice cream  $\Rightarrow$  don't like (unless its Bent Spoon)

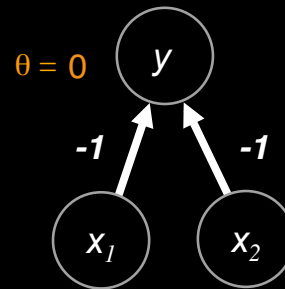
# The Problem

---

Actually, two problems:

- 1) Multiplicative would solve the problem:  
(if inputs can range from *negative* to positive)

<u>X-OR</u>		
$x_1$	$x_2$	$y$
-1	-1	1
-1	1	0
1	-1	0
1	1	1



Separable

(-1,1)      (1,1)

(-1,-1)      (1,-1)

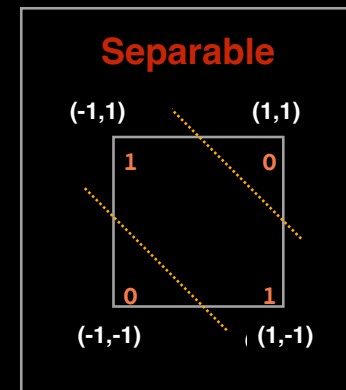
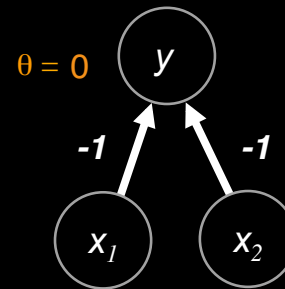
# The Problem

---

Actually, two problems:

- 1) Multiplicative would solve the problem:  
(if inputs can range from *negative* to positive)

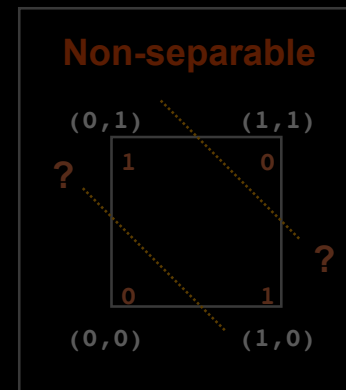
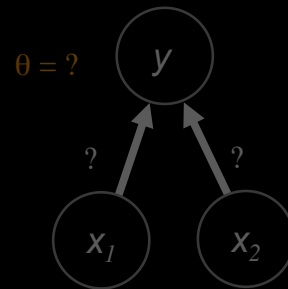
<u>X-OR</u>		
$x_1$	$x_2$	$y$
-1	-1	1
-1	1	0
1	-1	0
1	1	1



# The Problem

---

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



# The Problem

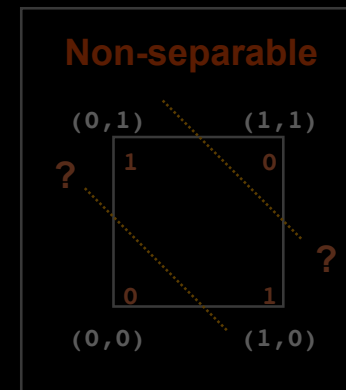
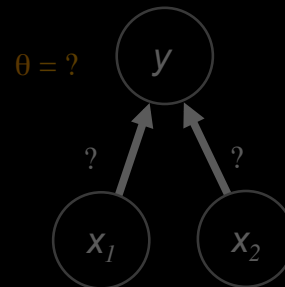
---

Actually, two problems:

## 1) Additive

– can't solve problems that are not “*linearly separable*”...

<u>X-OR</u>		
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



## 2) Single layered

– “*hidden units*” solve the problem...

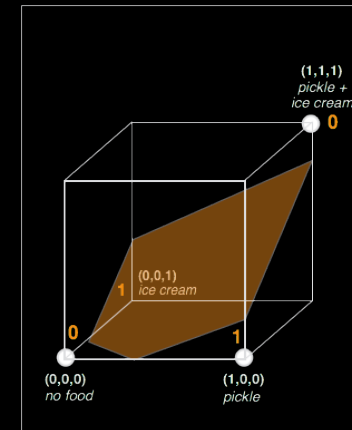
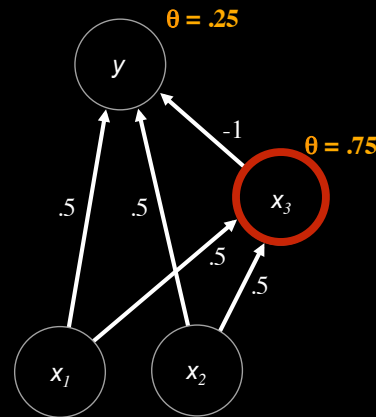
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not “linearly separable”...

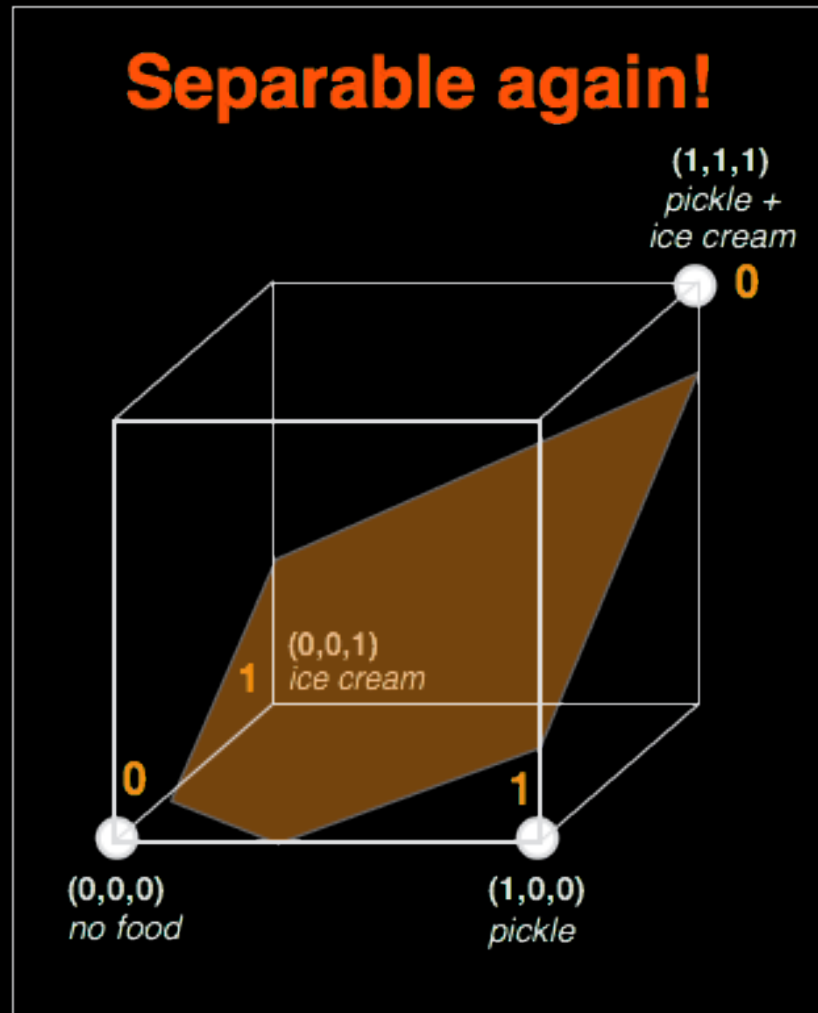
$x_1$	$x_2$	$x_3$	$y$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



## 2) Multi-layered

– “hidden units” solve the problem

# Add a 3rd Dimension



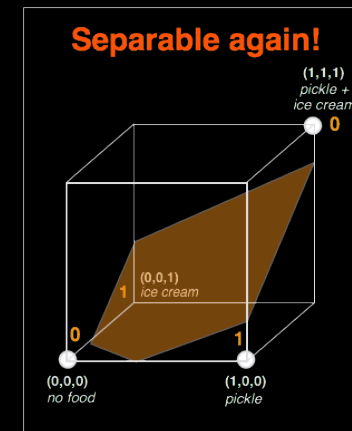
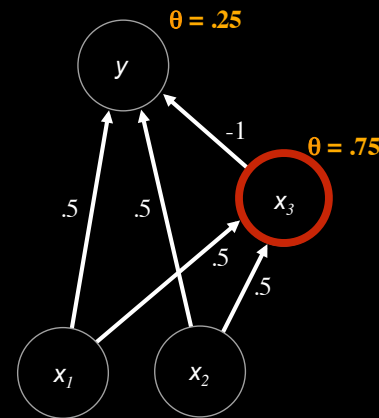
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not “linearly separable”...

$x_1$	$x_2$	$x_3$	$y$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



## 2) Multi-layered

– “hidden units” solve the problem (Minsky & Papert knew that...)

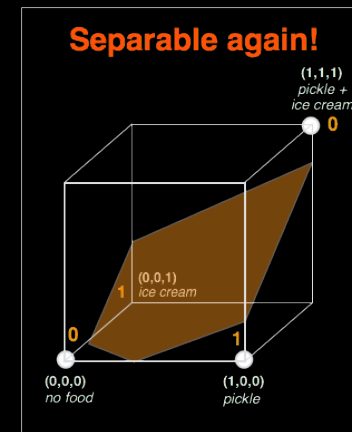
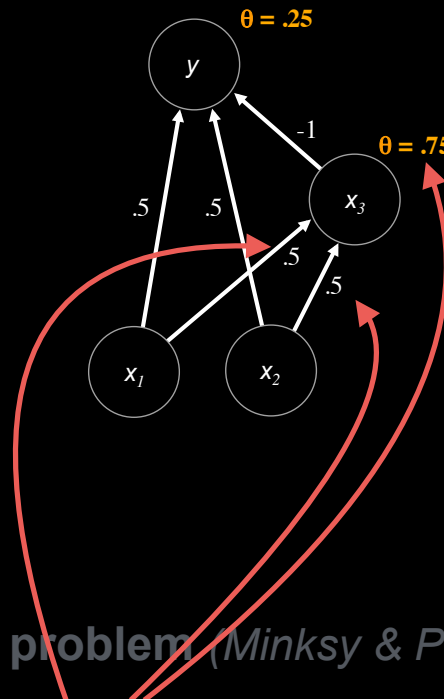
# The Problem

Actually, two problems:

## 1) Additive

– can't solve problems that are not “linearly separable”...

$x_1$	$x_2$	$x_3$	$y$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



## 2) Multi-layered

– “hidden units” solve the problem (Minsky & Papert knew that)...

– **but how do you train hidden units? where does the error come from?**

# The Solution

---

# The Solution

---

- **Backpropagation Algorithm** (*Generalized Delta Rule*)  
Generalizes the LMS procedure (Delta Rule) to apply to:
  - Non-linear units
  - Hidden units

# The Solution

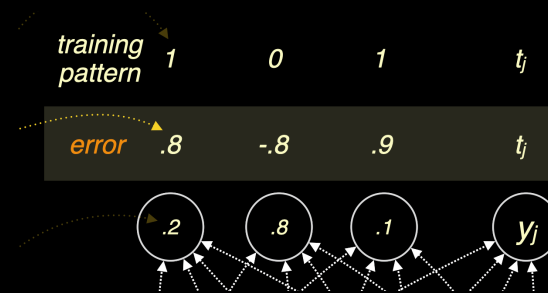
- **Backpropagation Algorithm** (*Generalized Delta Rule*)

Generalizes the LMS procedure (Delta Rule) to apply to:

- Non-linear units
- Hidden units

- **Start with Delta Rule:**

$$\Delta w_{ij} = \epsilon \delta_j x_i$$



# The Solution

- **Backpropagation Algorithm** (*Generalized Delta Rule*)

Generalizes the LMS procedure (Delta Rule) to apply to:

- Non-linear units
- Hidden units

- **Start with Delta Rule:**

$$\Delta w_{ij} = \epsilon \delta_j x_i$$

- **For the output units:**

$$\delta_j = (t_j - y_j) \cdot f'_j(\text{net}_j) \quad \text{where} \quad f'_j(\text{net}_j) = \frac{d(x_j)}{d(\text{net}_j)}$$



# The Solution

- **Backpropagation Algorithm** (*Generalized Delta Rule*)

Generalizes the LMS procedure (Delta Rule) to apply to:

- Non-linear units
- Hidden units

- **Start with Delta Rule:**

$$\Delta w_{ij} = \epsilon \delta_j x_i$$

- **For the output units:**

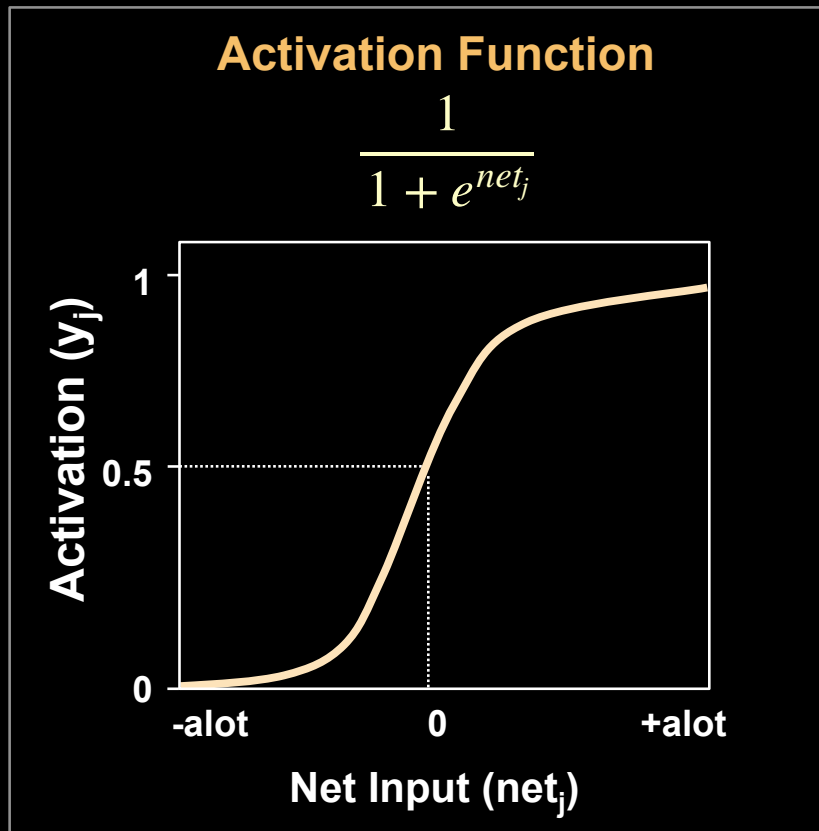
$$\delta_j = (t_j - y_j) \cdot f'_j(\text{net}_j) \quad \text{where} \quad f'_j(\text{net}_j) = \frac{d(x_j)}{d(\text{net}_j)}$$

**Note:**

- this is similar to the Delta Rule, except that the error term is now scaled by the **derivative** of the **activation function**
- this requires that the activation function be **differentiable**, and therefore **continuous**



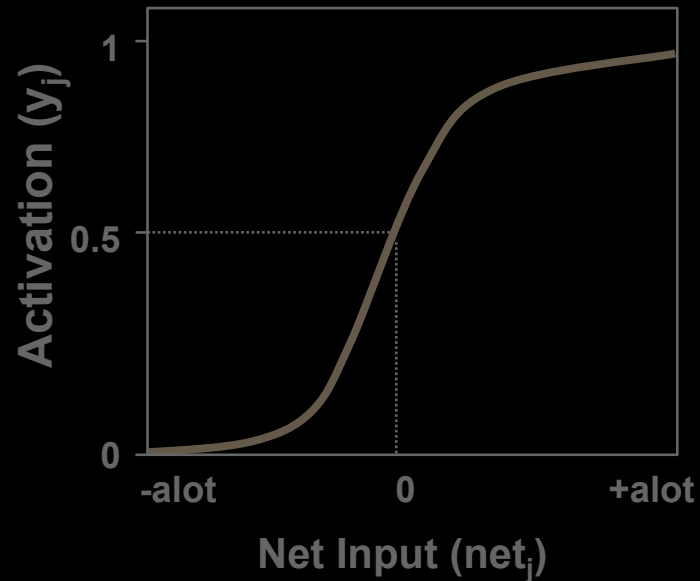
# Logistic Activation Function



# Logistic Activation Function

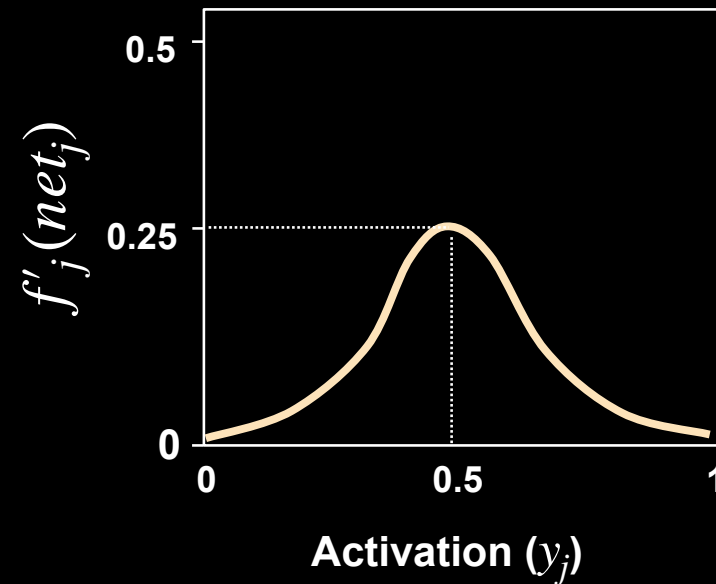
Activation Function

$$\frac{1}{1 + e^{-net_j}}$$

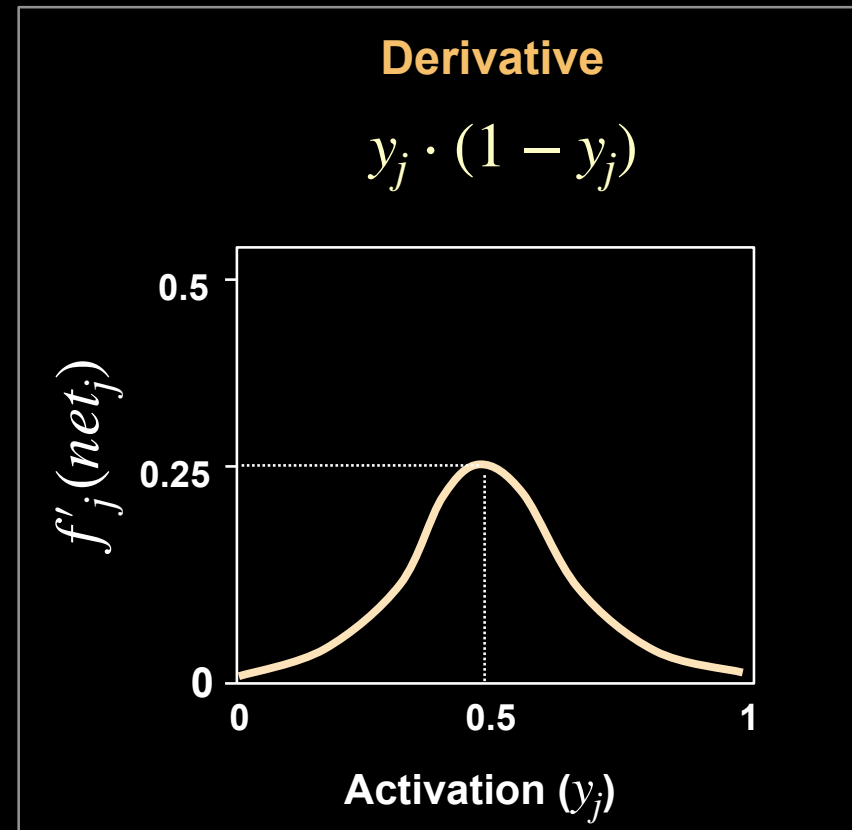
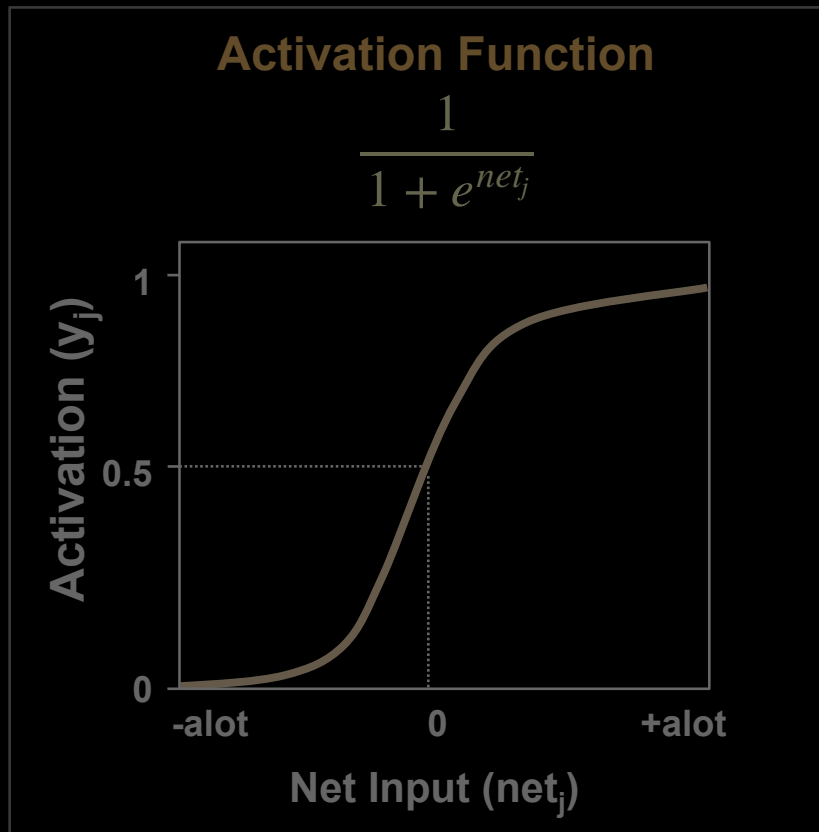


Derivative

$$y_j \cdot (1 - y_j)$$



# Logistic Activation Function

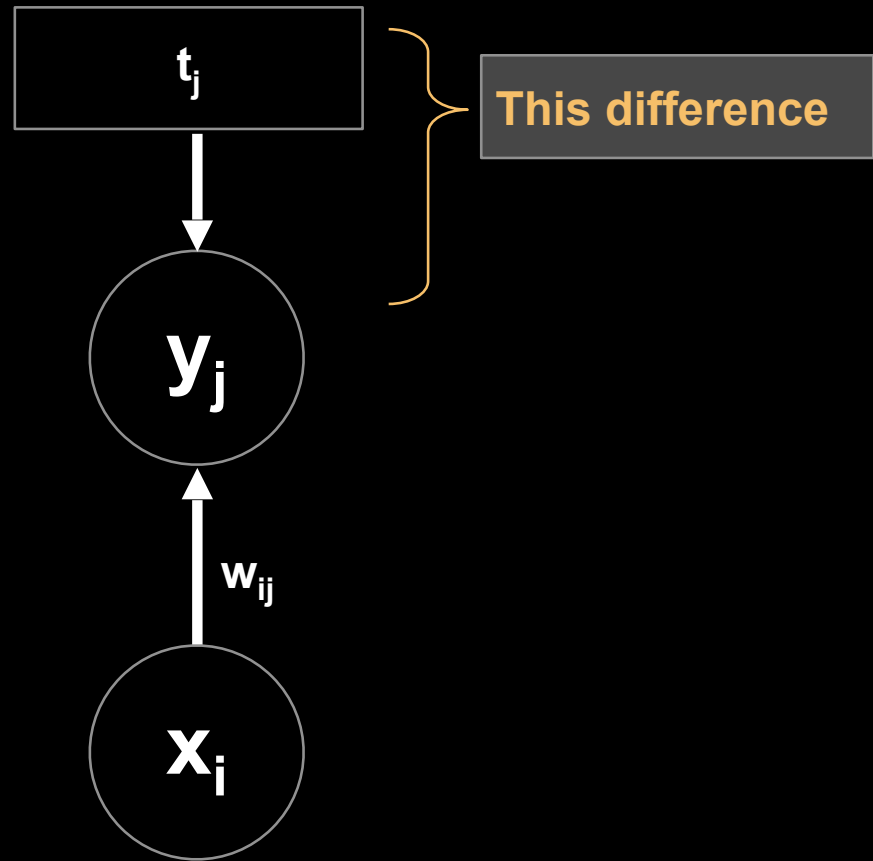


## Note:

- Derivative (slope) is maximal when activation is 0.5 (i.e., unit is uncommitted)

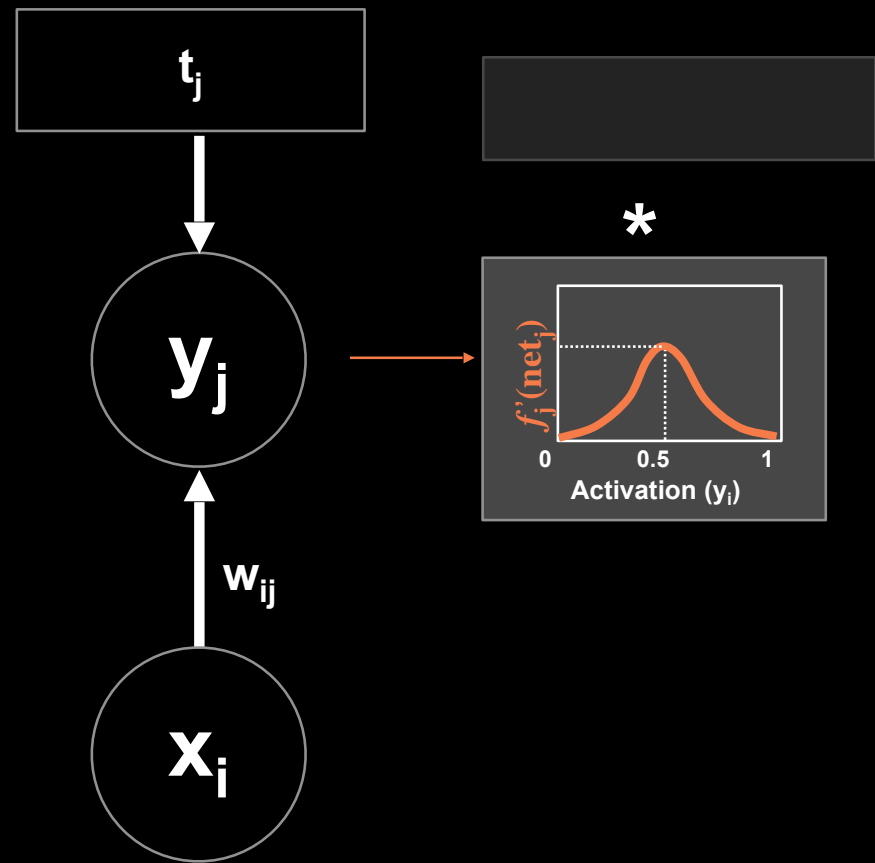
# Effect of Non-Linearity

$$\Delta W_{ij} = (t_j - y_j)$$



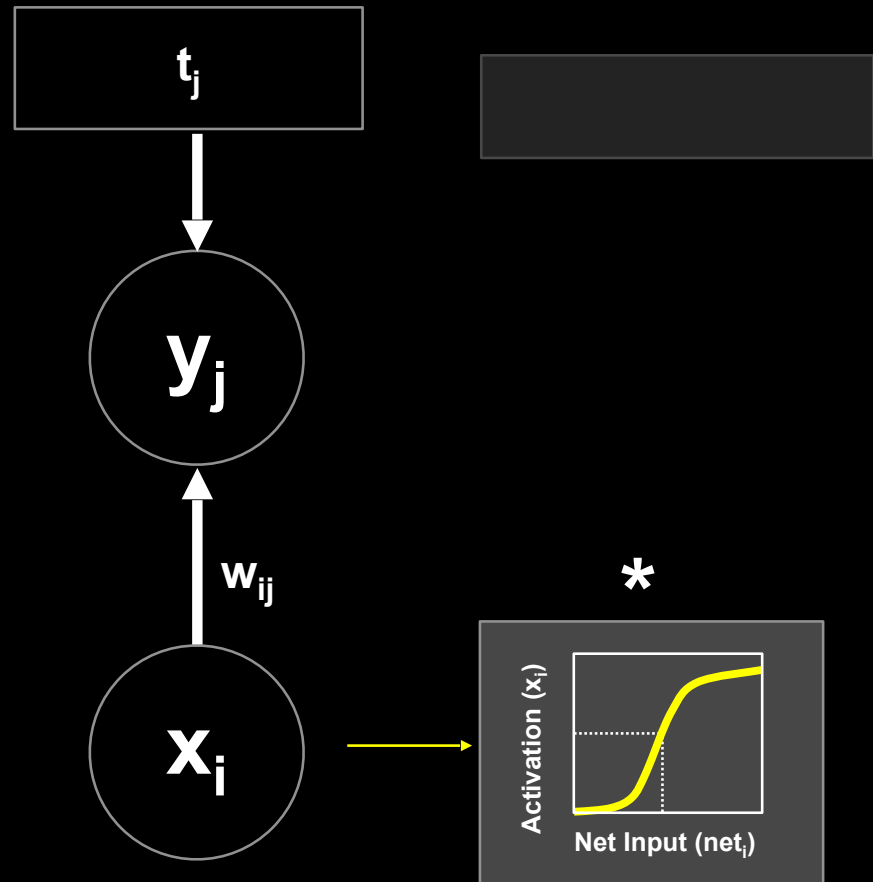
# Effect of Non-Linearity

$$\Delta W_{ij} = (t_j - y_j) \cdot f'_j(\text{net}_j)$$



# Effect of Non-Linearity

$$\Delta W_{ij} = (t_j - y_j) \cdot f'_j(\text{net}_j) \cdot x_j$$

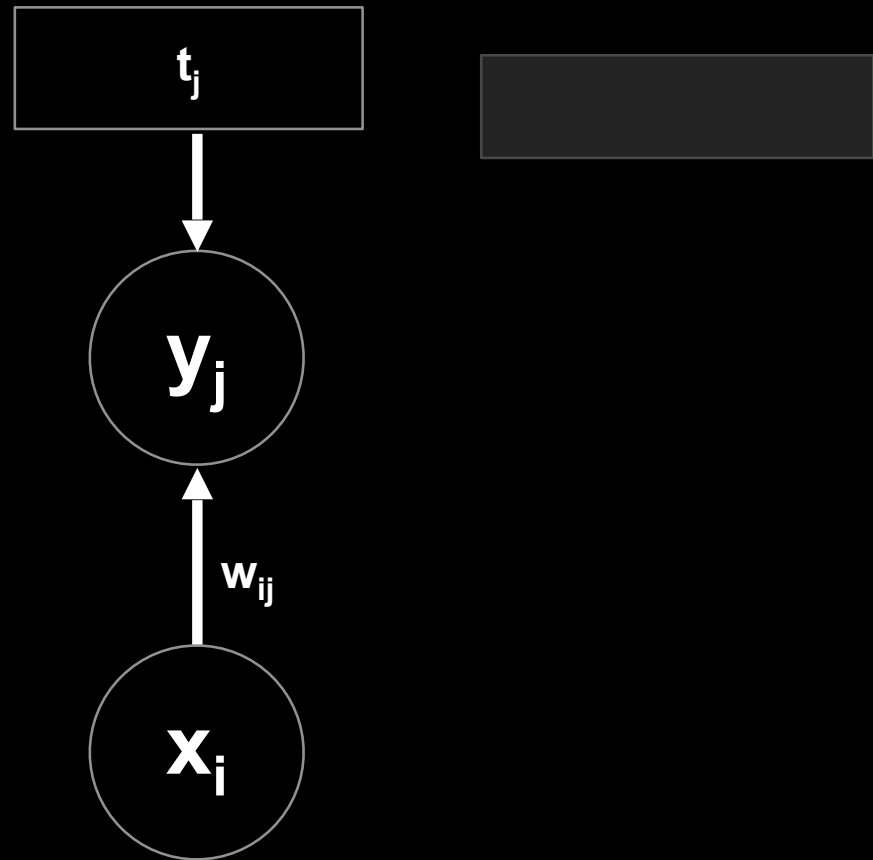


# Effect of Non-Linearity

$$\Delta W_{ij} = (t_j - y_j) \cdot f'_j(\text{net}_j) \cdot x_j$$

- **Note:**

- Biggest weights changes are made to “uncommitted” receiving units ( $y$ ) (i.e., for which activity  $\sim 0.5$ )
- Weight change is scaled by activity of sending unit ( $x$ )



# Backpropagation Algorithm

- **Backpropagation Learning Algorithm**

Generalizes the LMS procedure to apply to:

- Hidden units
- Non-linear units

- **Start with LMS (Delta) rule:**

$$\Delta w_{ij} = \epsilon \delta_j x_j$$

- **For the output units:**

$$\delta_j = (t_j - y_j) \cdot f'_j(\text{net}_j) \quad \text{where} \quad f'_j(\text{net}_j) = \frac{d(x_j)}{d(\text{net}_j)}$$

- **For the hidden units:**

$$\delta_k = (\sum_j \delta_j w_{jk}) \cdot f'_k(\text{net}_k) \quad \text{where} \quad f'_k(\text{net}_k) = f'_j(\text{net}_j) = \frac{d(x_k)}{d(\text{net}_k)}$$

# Backpropagation Algorithm

- Backpropagation Learning Algorithm

Generalizes the LMS procedure to apply to:

- Hidden units
- Non-linear units

- Start with LMS (Delta) rule:

$$\Delta w_{ij} = \epsilon \delta_j x_j$$

- For the output units:

$$\delta_j = (t_j - y_j) \cdot f'_j(\text{net}_j) \quad \text{where} \quad f'_j(\text{net}_j) = \frac{d(x_j)}{d(\text{net}_j)}$$

- For the hidden units:

$$\delta_k = (\sum_j \delta_j w_{jk}) \cdot f'_k(\text{net}_k) \quad \text{where} \quad f'_k(\text{net}_k) = f'_j(\text{net}_j) = \frac{d(x_k)}{d(\text{net}_k)}$$

**KEY OBSERVATION:**

- The “error”  $\delta_k$  for hidden unit  $k$  is the sum over the **errors** (deltas  $\delta_j$ ) **of the units; to which it projects, weighted by its connections  $w_{jk}$  to them.**

# Backpropagation Algorithm

- **Backpropagation Learning Algorithm**

Generalizes the LMS procedure to apply to:

- Hidden units
- Non-linear units

- **Start with LMS (Delta) rule:**

$$\Delta w_{ij} = \epsilon \delta_j x_j$$

- **For the output units:**

$$\delta_j = (t_j - y_j) \cdot f'_j(\text{net}_j) \quad \text{where} \quad f'_j(\text{net}_j) = \frac{d(x_j)}{d(\text{net}_j)}$$

- **For the hidden units:**

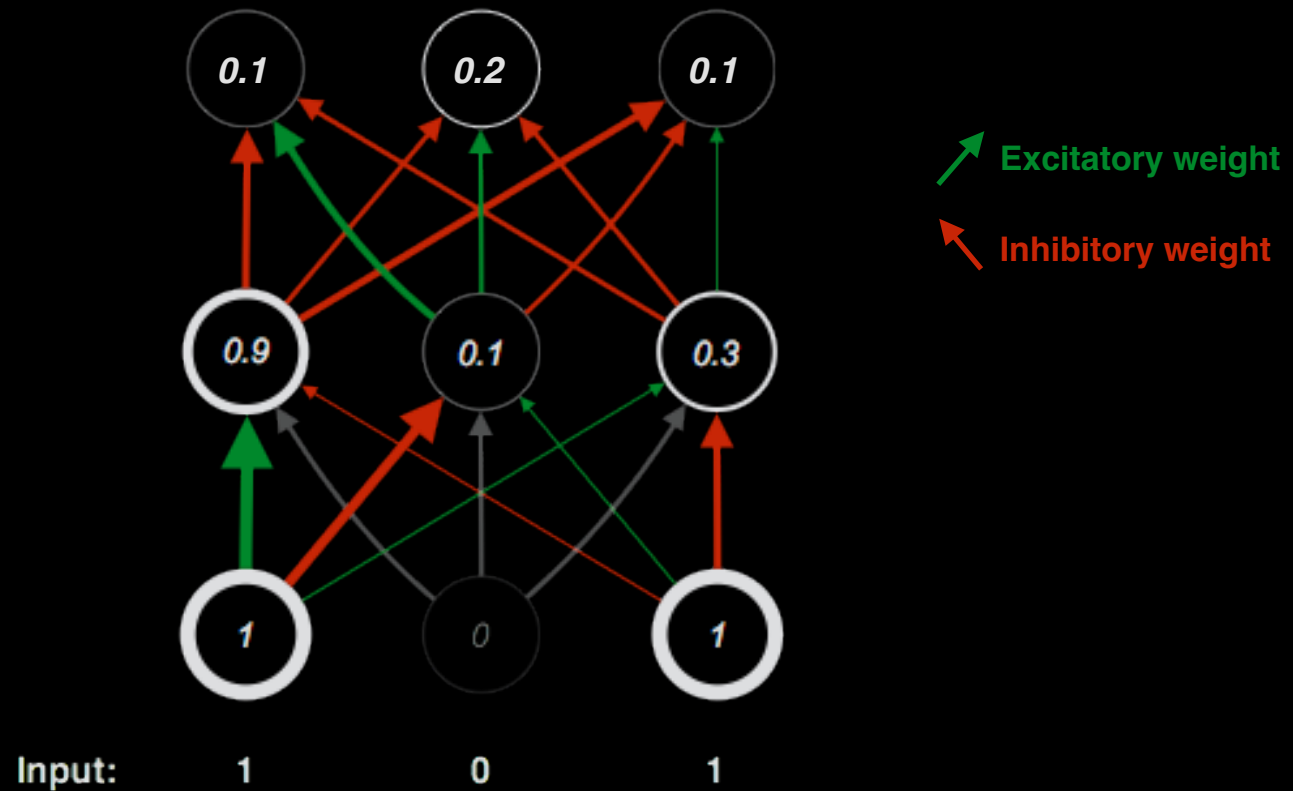
$$\delta_k = (\sum_j \delta_j w_{jk}) \cdot f'_k(\text{net}_k) \quad \text{where} \quad f'_k(\text{net}_k) = f'_j(\text{net}_j) = \frac{d(x_k)}{d(\text{net}_k)}$$

**KEY OBSERVATION:**

- The “error”  $\delta_k$  for hidden unit  $k$  is the sum over the *errors* (deltas  $\delta_j$ ) of the units  $j$  to which it projects, weighted by its connections  $w_{jk}$  to them.
- Minimizing this reduces the average *impact each unit  $k$  has on the error at the layer above it*

# Error Calculation

---



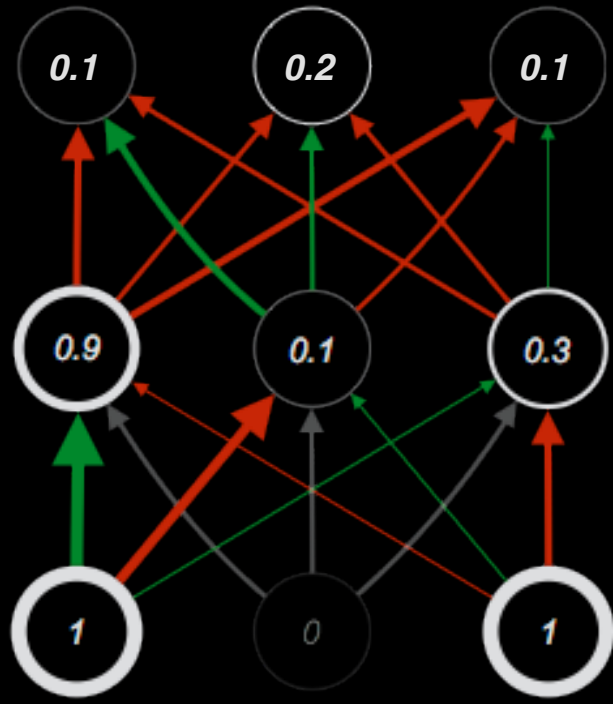
# Error Calculation

Target:

0

1

0



Input:

1

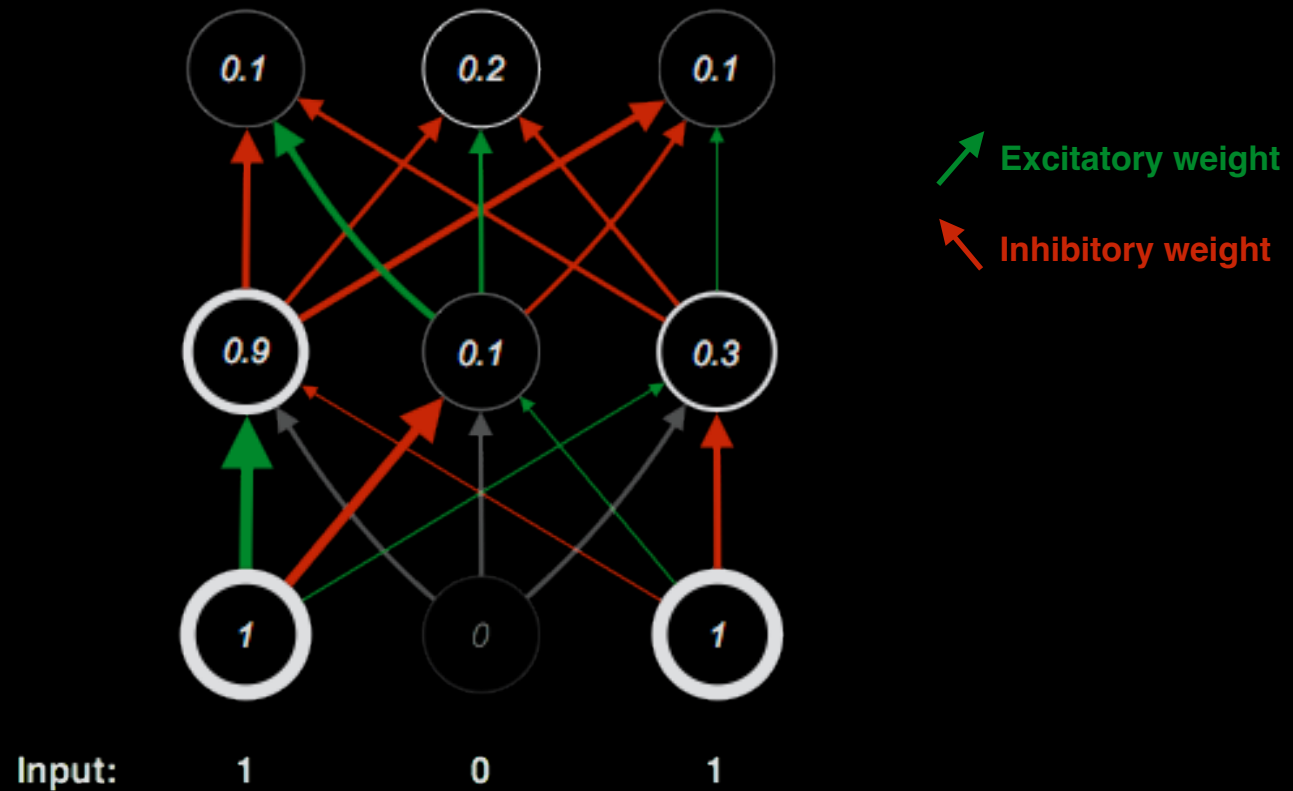
0

1

Excitatory weight  
Inhibitory weight

# Error Calculation

Target:      0            1            0  
Output:      0.1          0.2          0.1

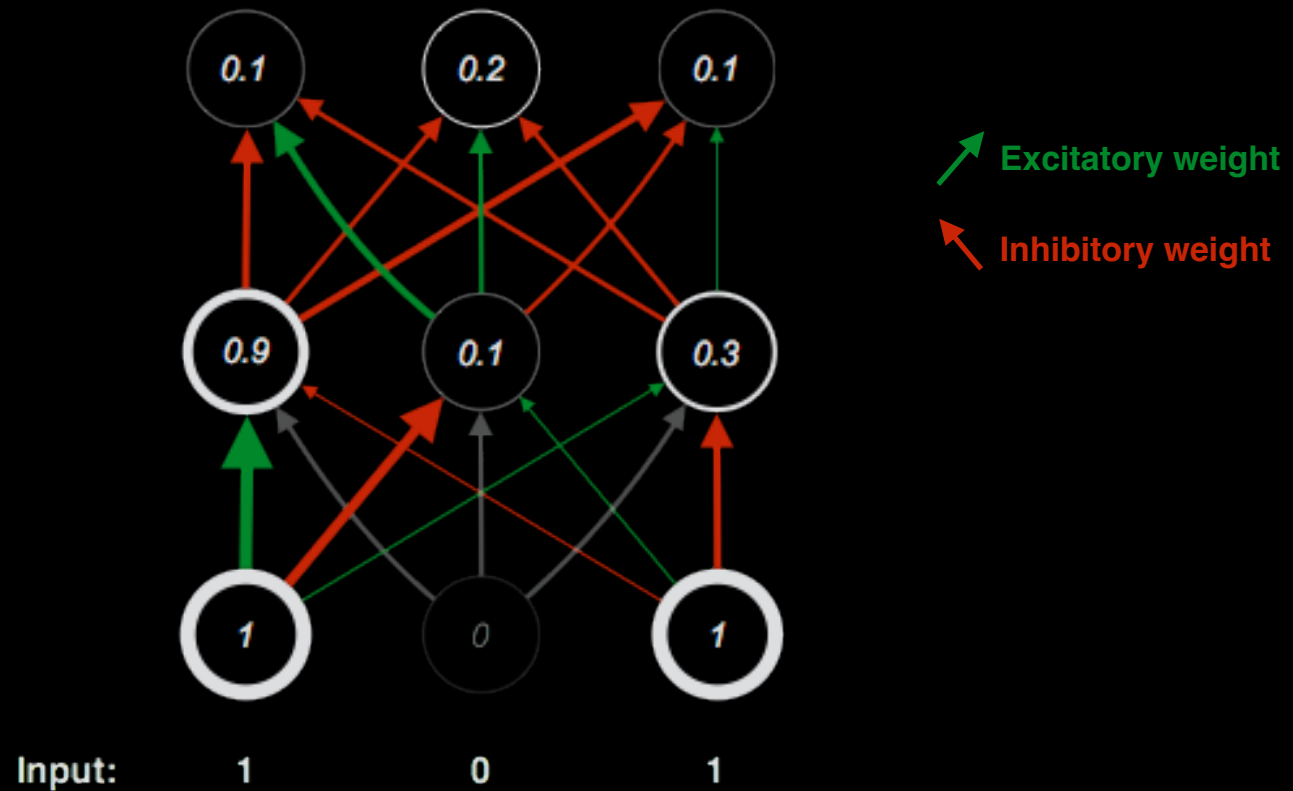


# Error Calculation

---

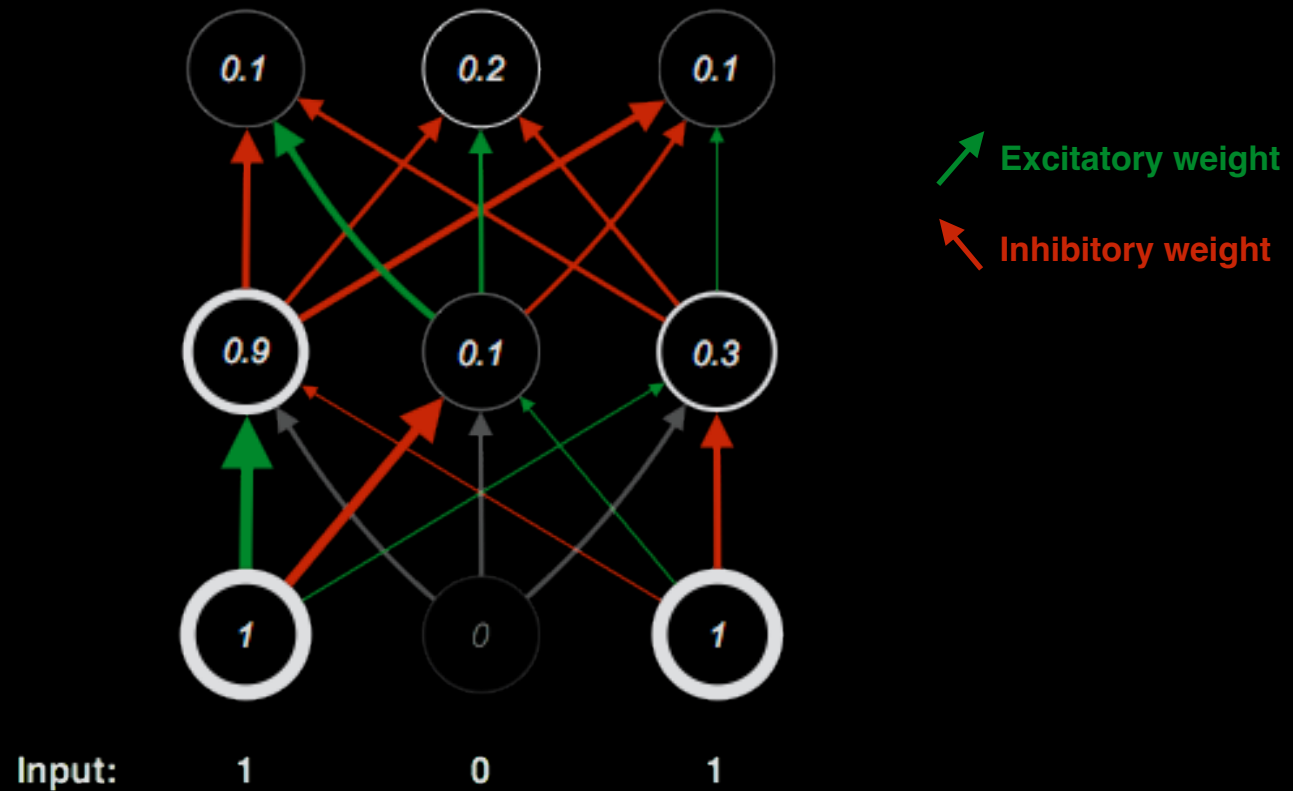
Target:	0	1	0
- Output:	0.1	0.2	0.1

---



# Error Calculation

Target:	0	1	0
- Output:	0.1	0.2	0.1
= Error:	-0.1	+0.8	-0.1

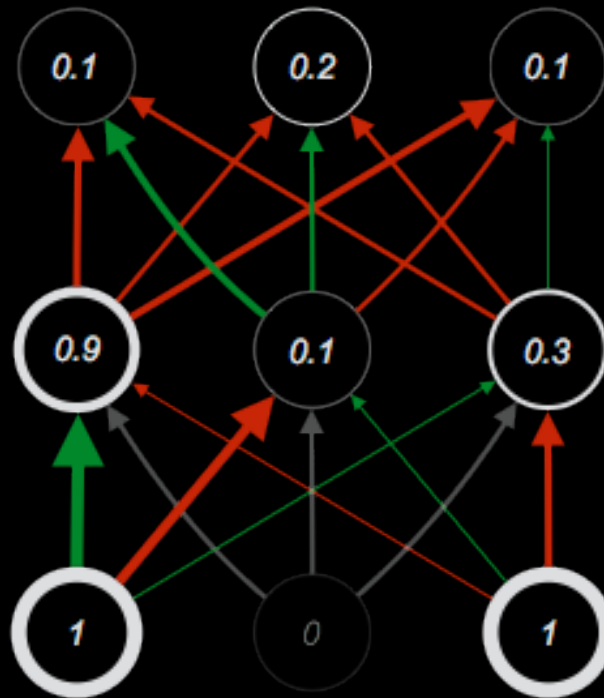




# Error Calculation

Target:	0	1	0
- Output:	0.1	0.2	0.1
= Error:	-0.1	+0.8	-0.1

- **Hint:**

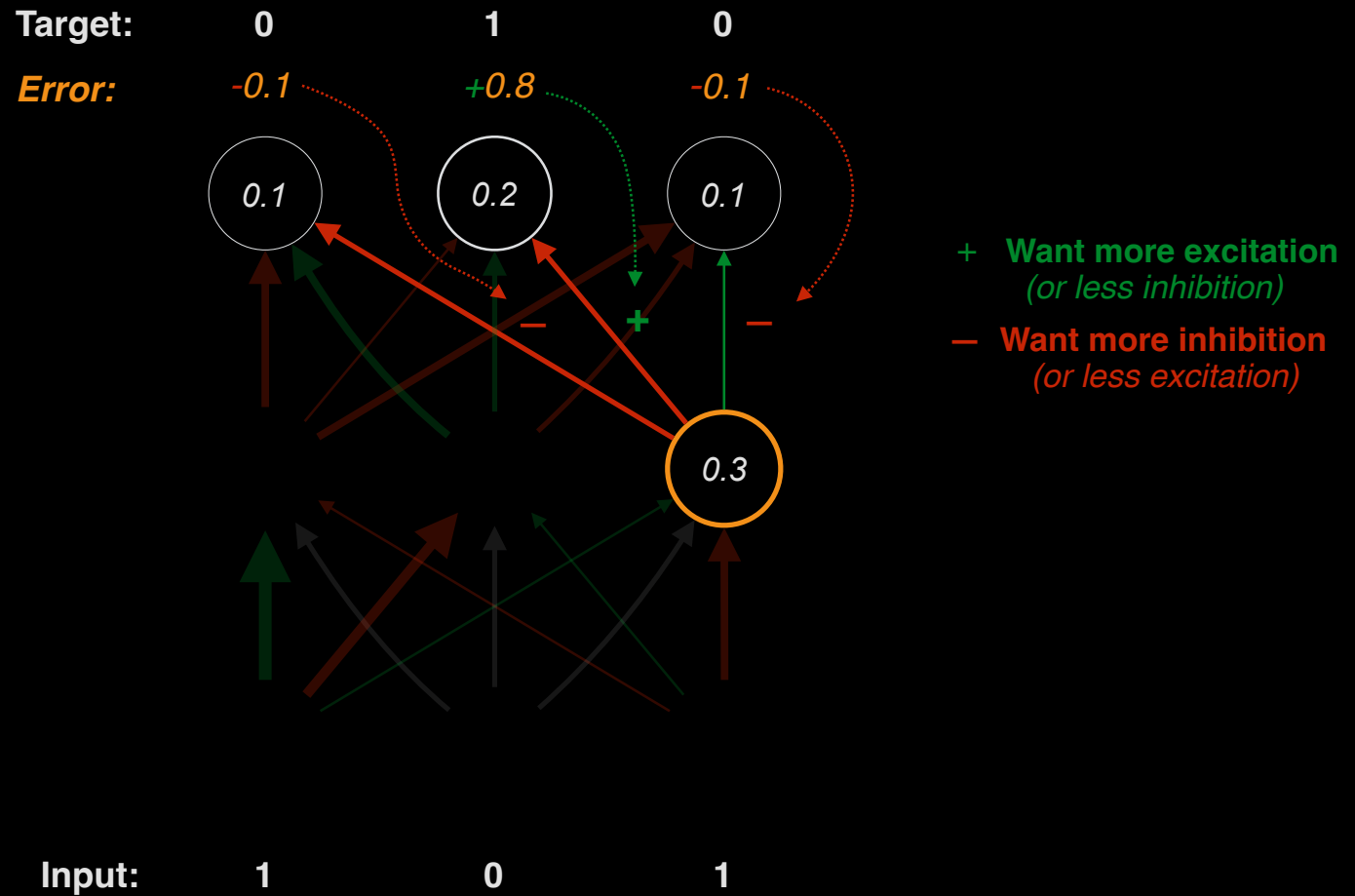
Think of "error" as desired amount of change



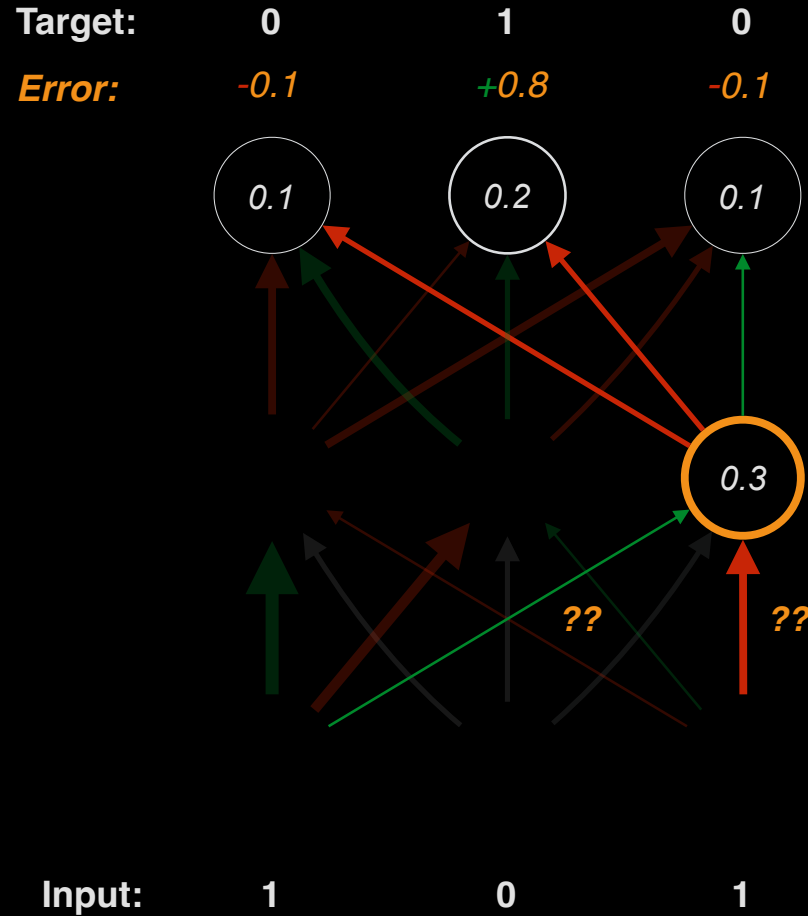
 Excitatory weight  
 Inhibitory weight

Input: 1 0 1

# Output Layer: Delta Rule

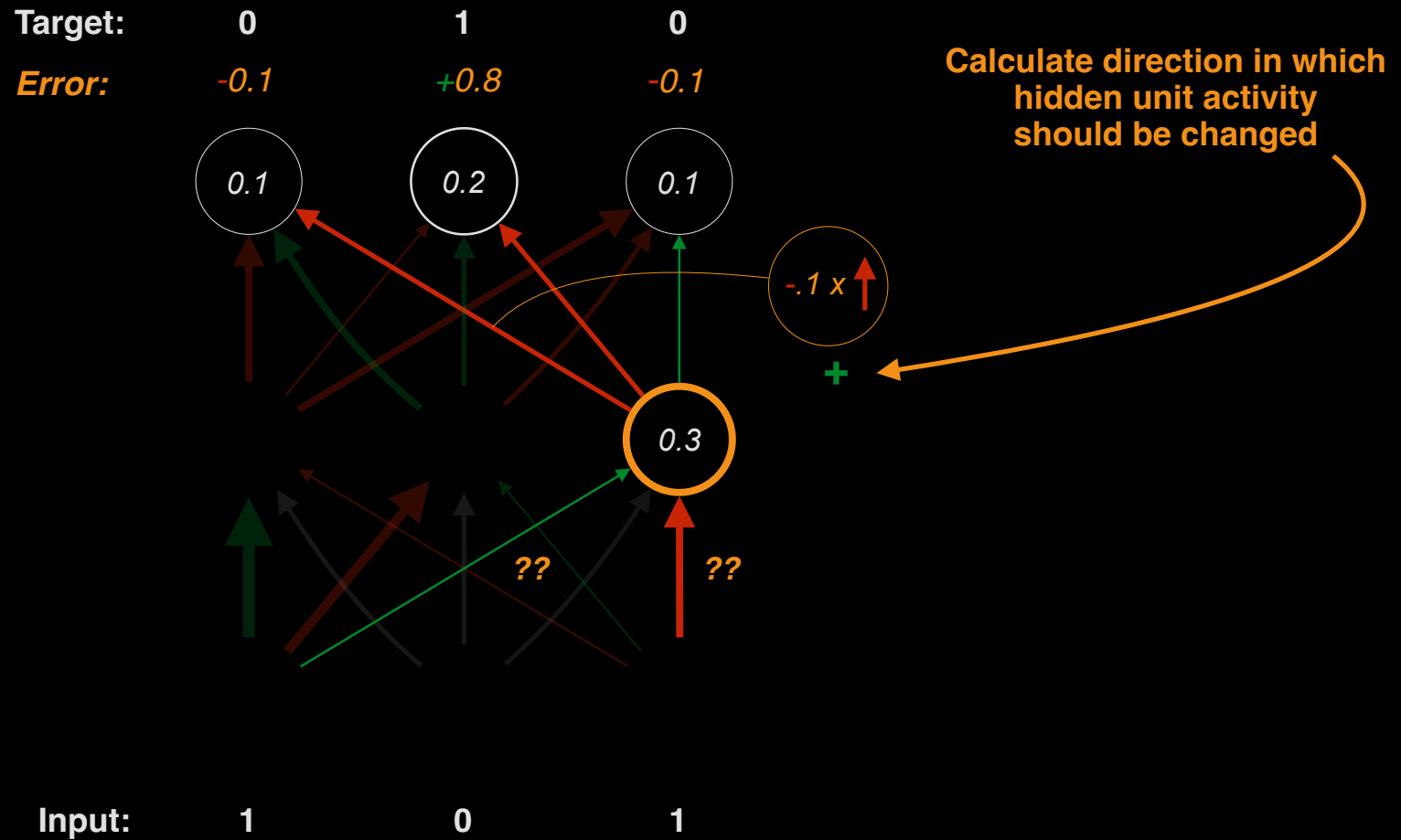


# Hidden Layer: *Generalized Delta Rule*

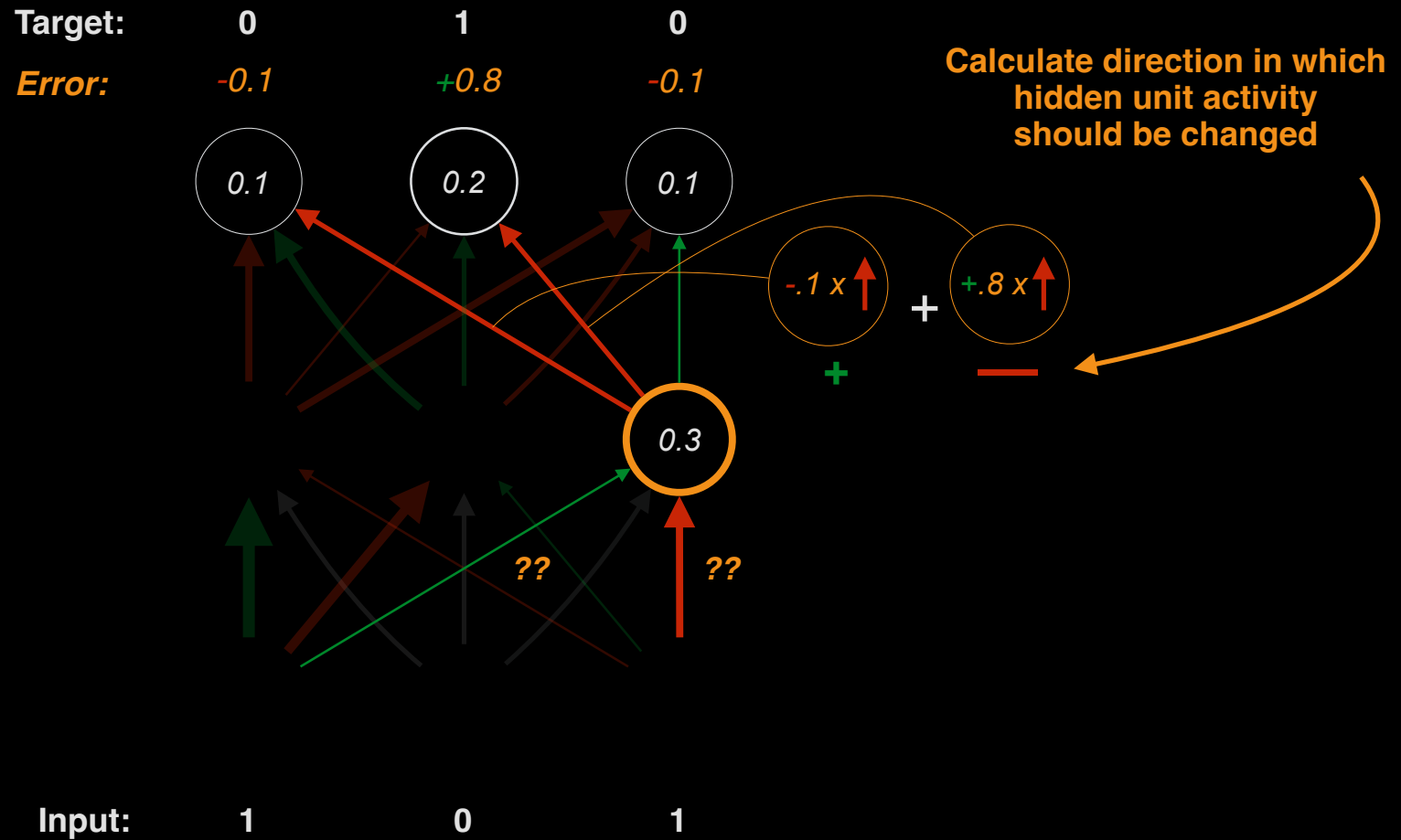


Calculate direction in which hidden unit activity should be changed

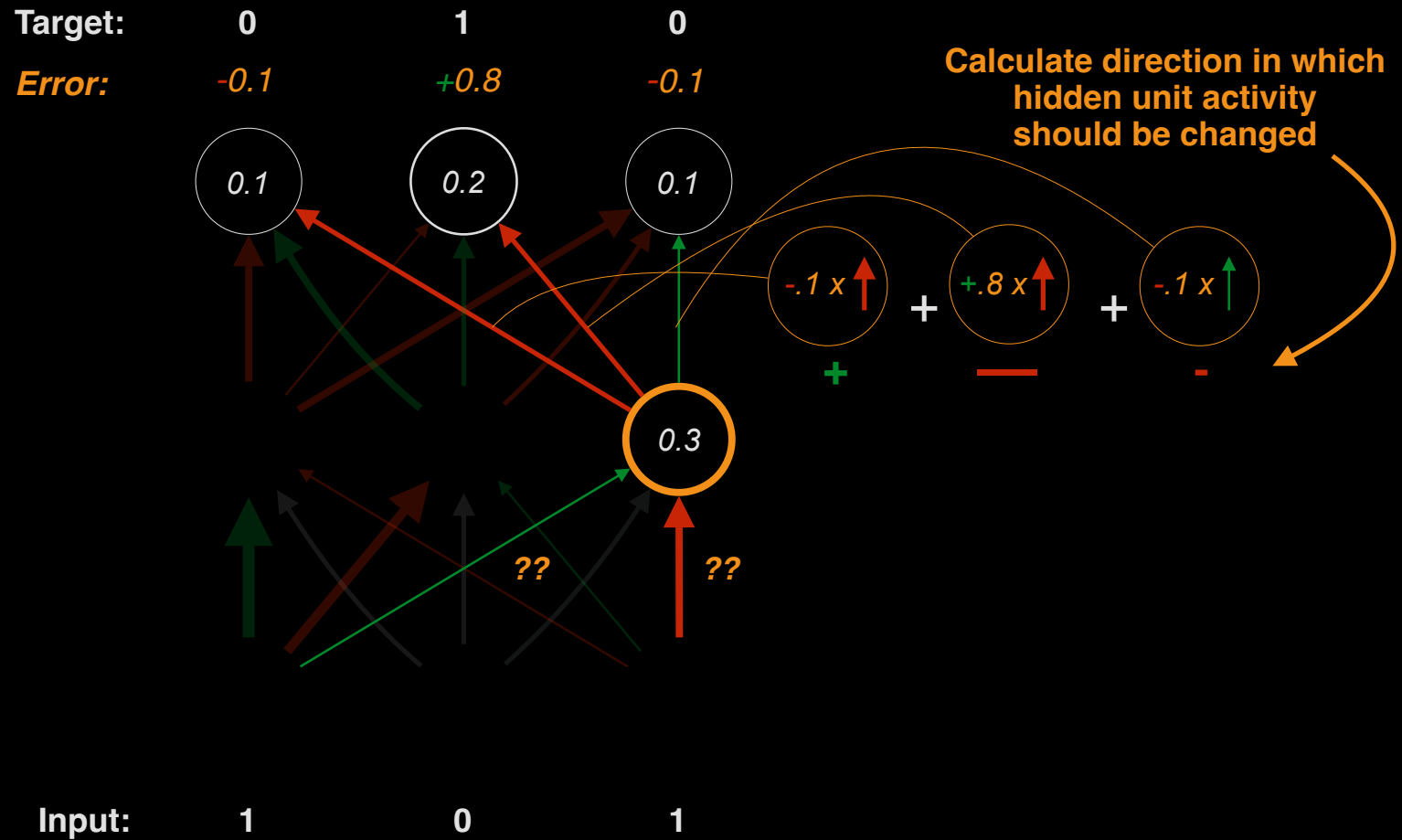
# Hidden Layer: *Generalized Delta Rule*



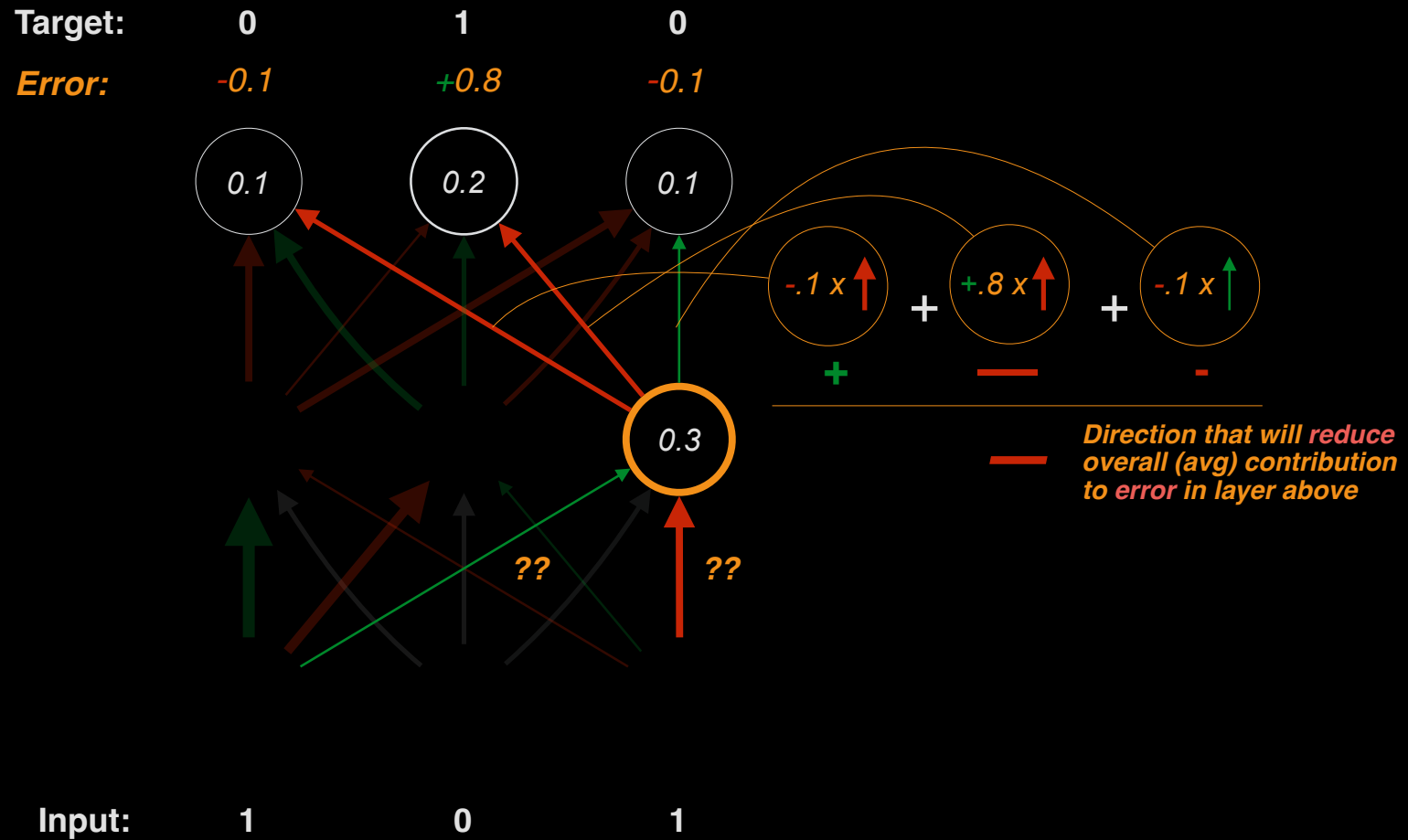
# Hidden Layer: *Generalized Delta Rule*



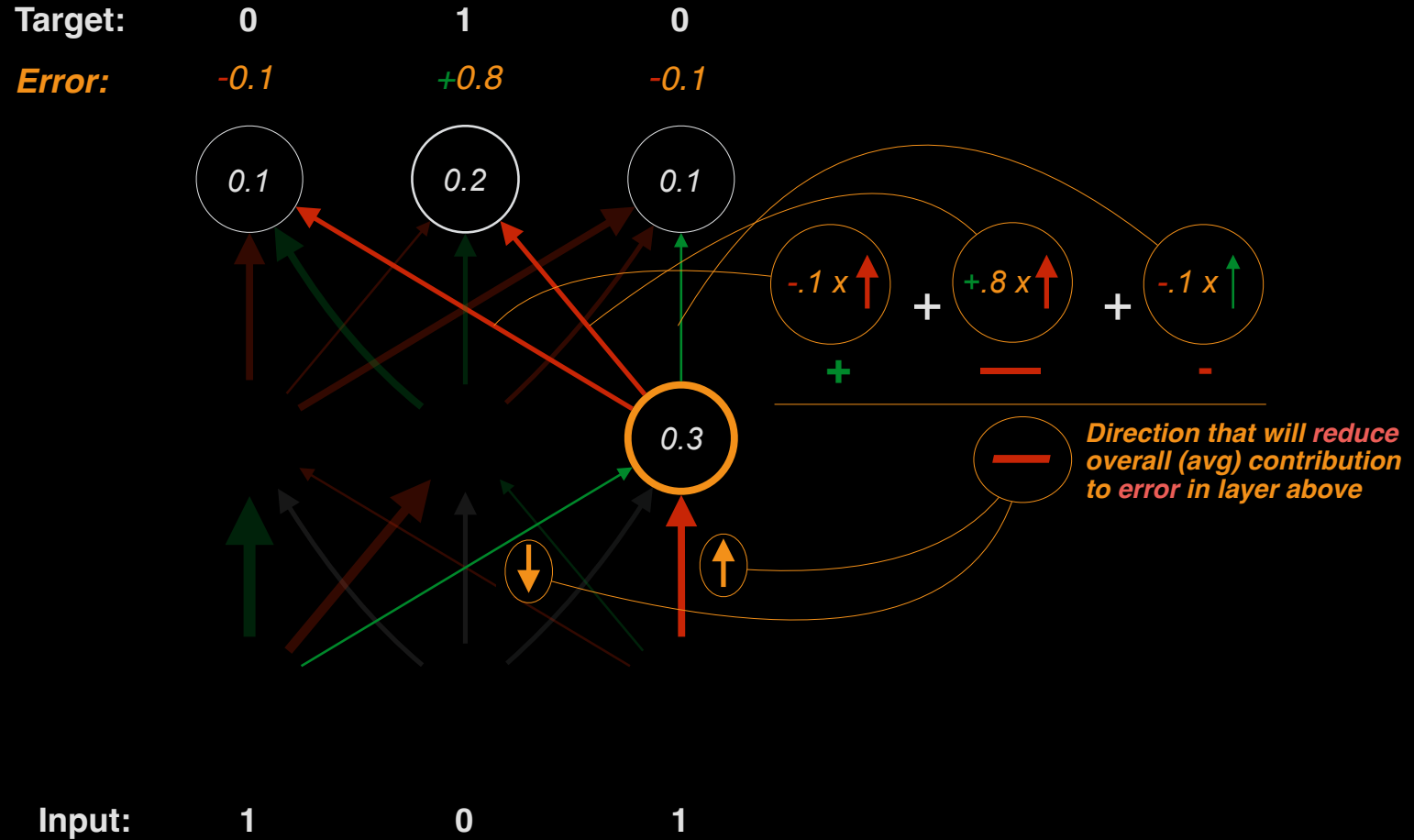
# Hidden Layer: *Generalized Delta Rule*



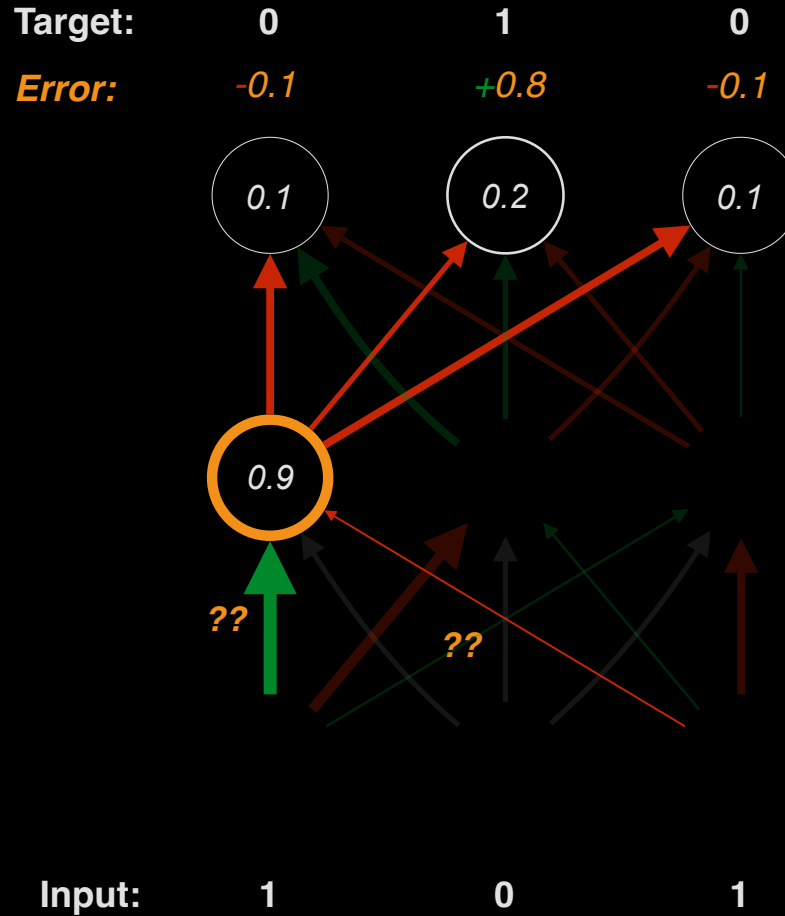
# Hidden Layer: *Generalized Delta Rule*



# Hidden Layer: *Generalized Delta Rule*



# Hidden Layer: *Generalized Delta Rule*



Calculate direction in which hidden unit activity should be changed

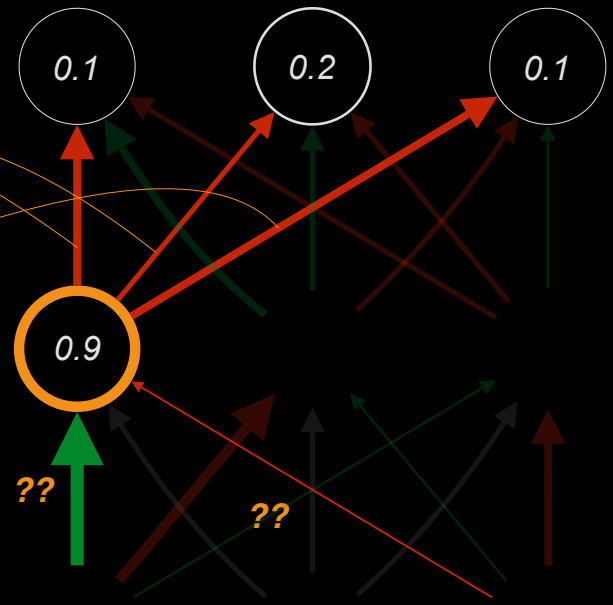
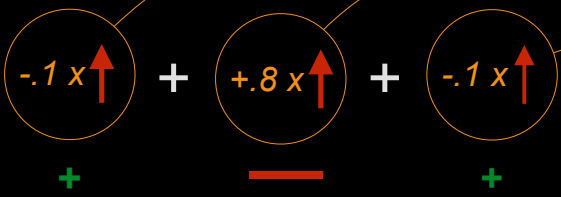




# Hidden Layer: *Generalized* Delta Rule

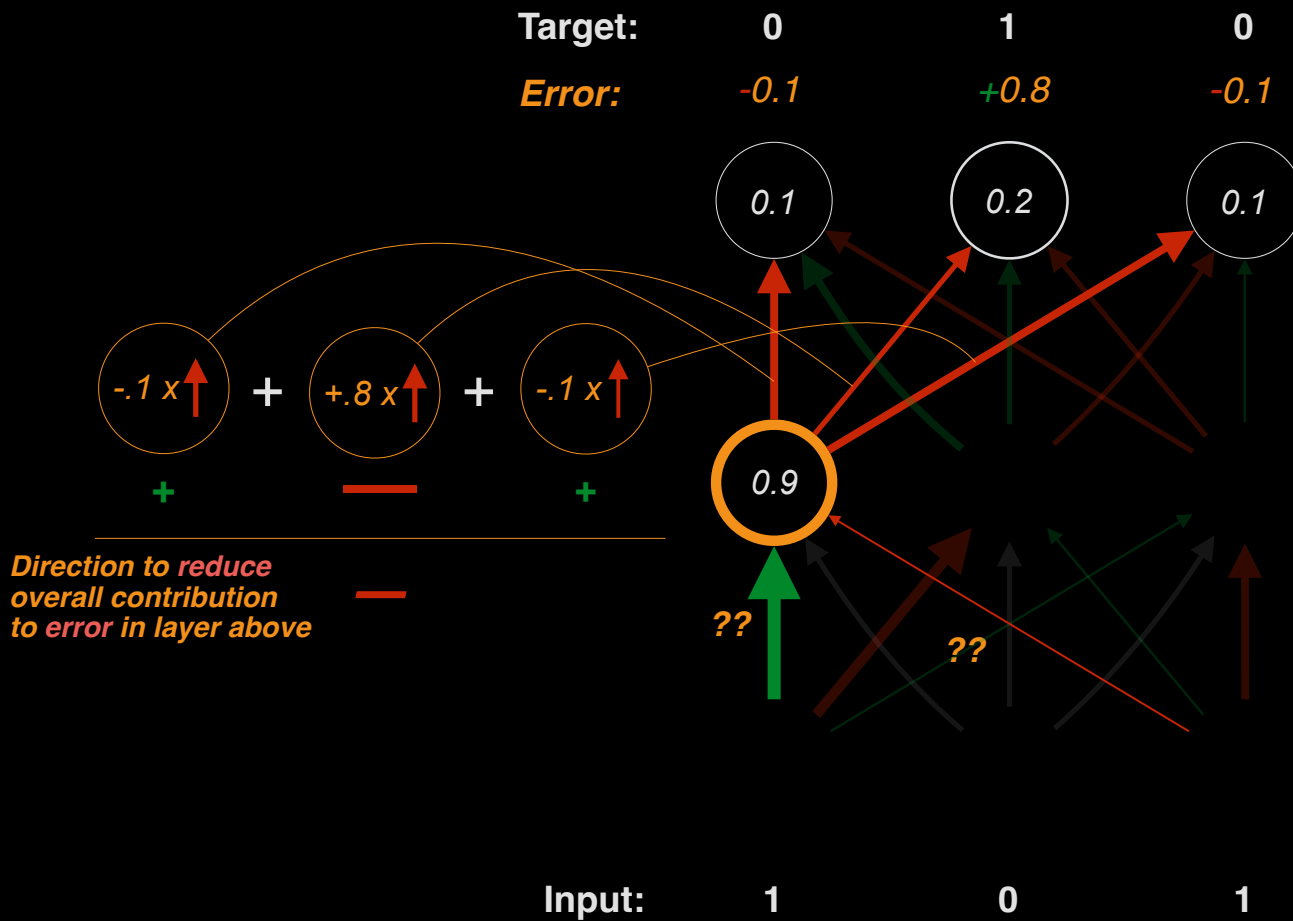
Target:      0            1            0  
Error:        -0.1        +0.8        -0.1

Calculate direction in which hidden unit activity should be changed



Input:      1            0            1

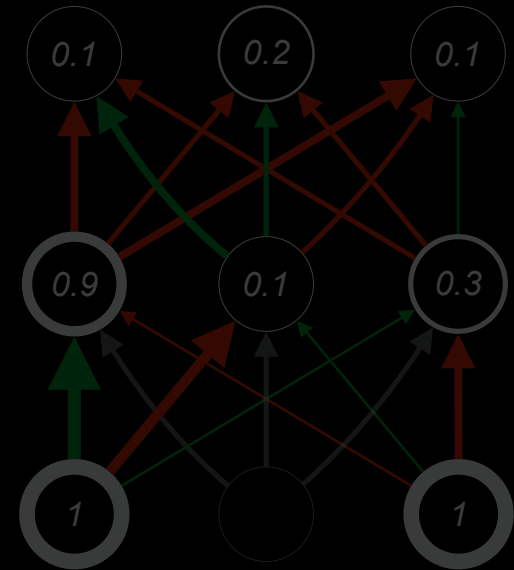
# Hidden Layer: *Generalized* Delta Rule





# The Solution

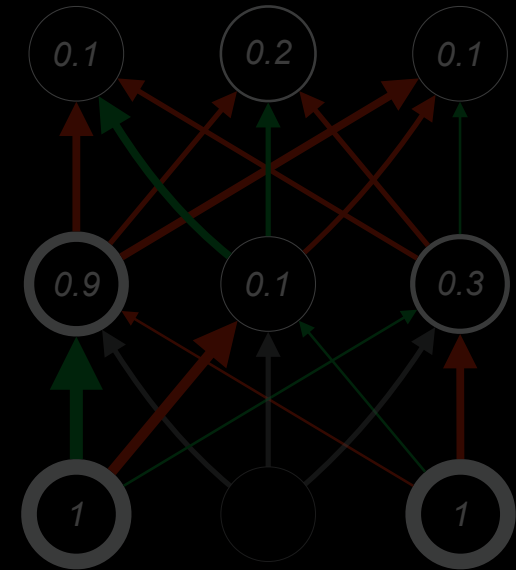
---



# The Solution

---

- The “Generalized Delta Rule:”
  - change each weight to the hidden unit so that it minimizes its contribution to the error at the layer above



- **Version that handles non-linear activation functions:**  
(e.g., the logistic)
  - **Backpropagation Learning Algorithm** (Rumelhart, Hinton & Williams, 1986)  
(aka “deep learning”):  
$$\delta_j = (\sum_k \delta_k w_{jk}) * f'_j(\text{net}_j)$$
 where  $f'_j(\text{net}_j) = d(x_j) / d(\text{net}_j)$  (derivate of activation fct)

# Backpropagation

---

# Backpropagation

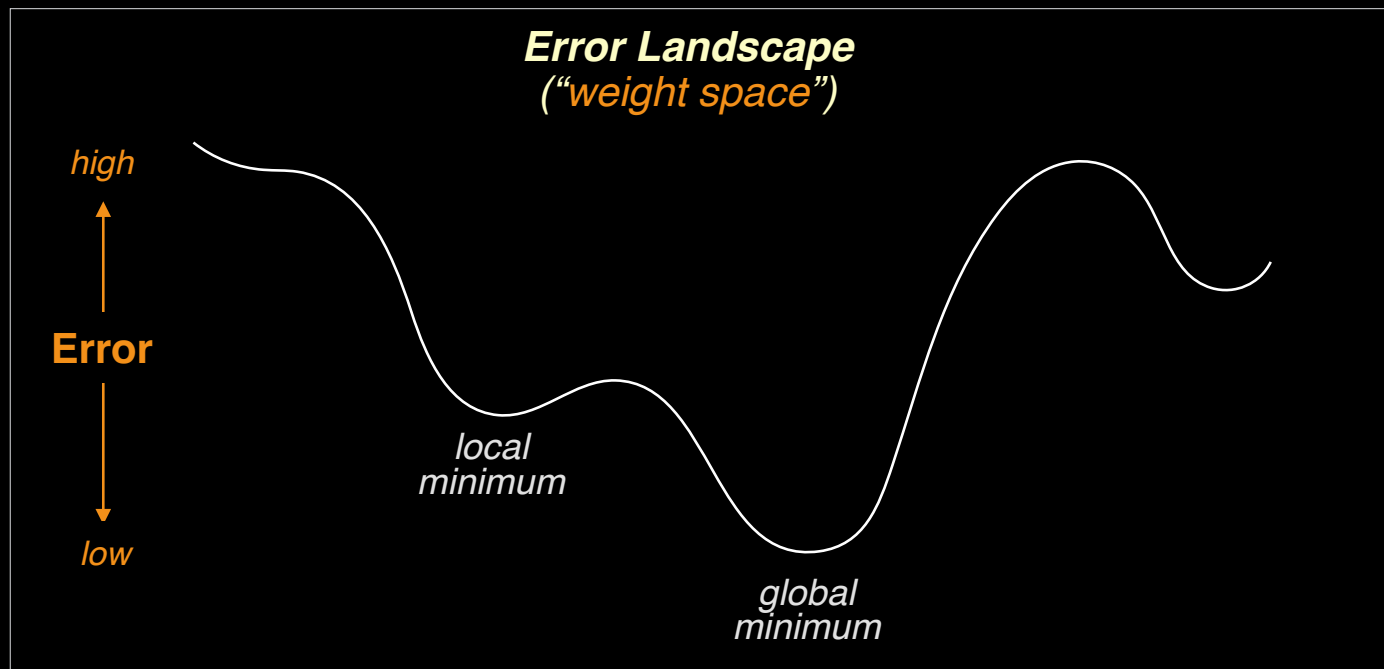
---

- “Gradient descent” algorithm
  - subject to “*local minima*”

# Backpropagation

---

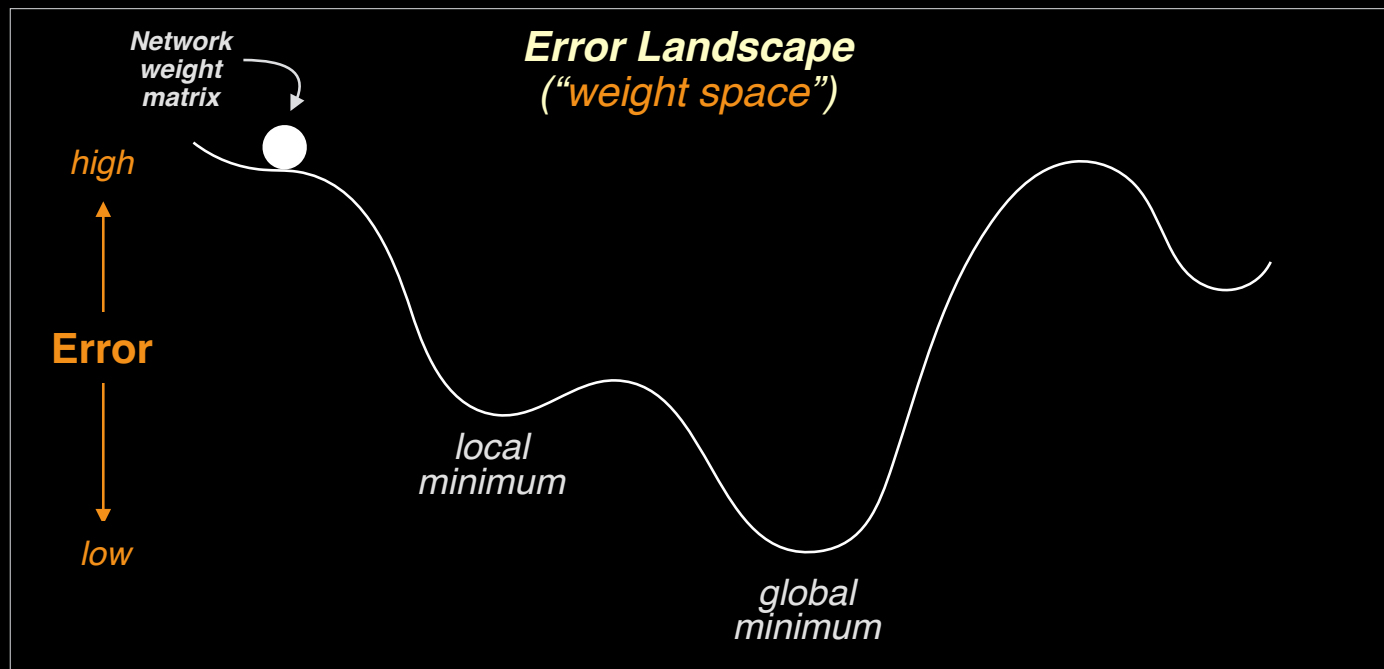
- “Gradient descent” algorithm
  - subject to “*local minima*”



# Backpropagation

---

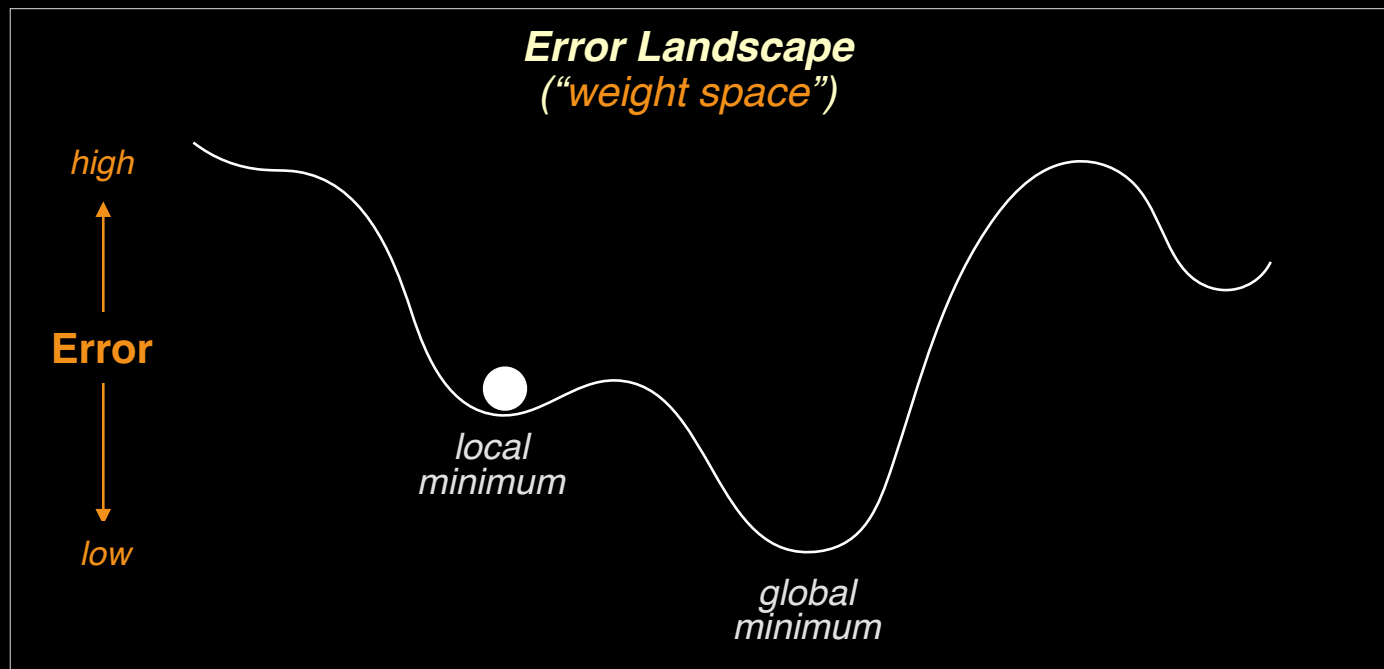
- “Gradient descent” algorithm
  - subject to “*local minima*”



# Backpropagation

---

- “Gradient descent” algorithm
  - subject to “*local minima*”



# Backpropagation

---

# Backpropagation

---

- **Learning rule:**

$$\delta_j = (\sum_k \delta_k w_{jk}) * f'_j(\text{net}_j) \text{ where } f'_j(\text{net}_j) = d(x_j) / d(\text{net}_j) \text{ (derivate of activation fct)}$$

- **“Gradient descent” algorithm**

- subject to *“local minima”*

- **Nevertheless, works almost all of the time**

- **example of how models / simulation:**

- **are an antidote to self-deception (*mental “muddles”*)**

- **help develop intuitions about behavior of complex (*e.g., nonlinear*) systems**

# Backpropagation

---

- **Learning rule:**  
$$\delta_j = (\sum_k \delta_k w_{jk}) * f'_j(\text{net}_j) \text{ where } f'_j(\text{net}_j) = d(x_j) / d(\text{net}_j) \text{ (derivate of activation fct)}$$
- **“Gradient descent” algorithm**
  - subject to *“local minima”*
- **Nevertheless, works almost all of the time**
  - example of how models / simulation:
    - are an antidote to self-deception (*mental “muddles”*)
    - help develop intuitions about behavior of complex (*e.g., nonlinear*) systems
- **Learns internal representations!**  
**(hidden units)**

# Backpropagation

---

- Learning rule:

$$\delta_j = (\sum_k \delta_k w_{jk}) * f'_j(\text{net}_j) \text{ where } f'_j(\text{net}_j) = d(x_j) / d(\text{net}_j) \text{ (derivate of activation fct)}$$

- “Gradient descent” algorithm

- subject to “local minima”

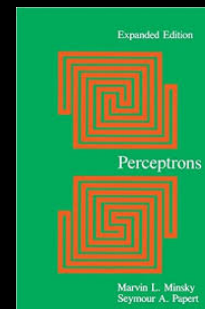
- Nevertheless, works almost all of the time

- example of how models / simulation:

- are an antidote to self-deception (*mental “muddles”*)

- help develop intuitions about behavior of complex (*e.g., nonlinear*) systems

- **Learns internal representations!**  
(*hidden units*)



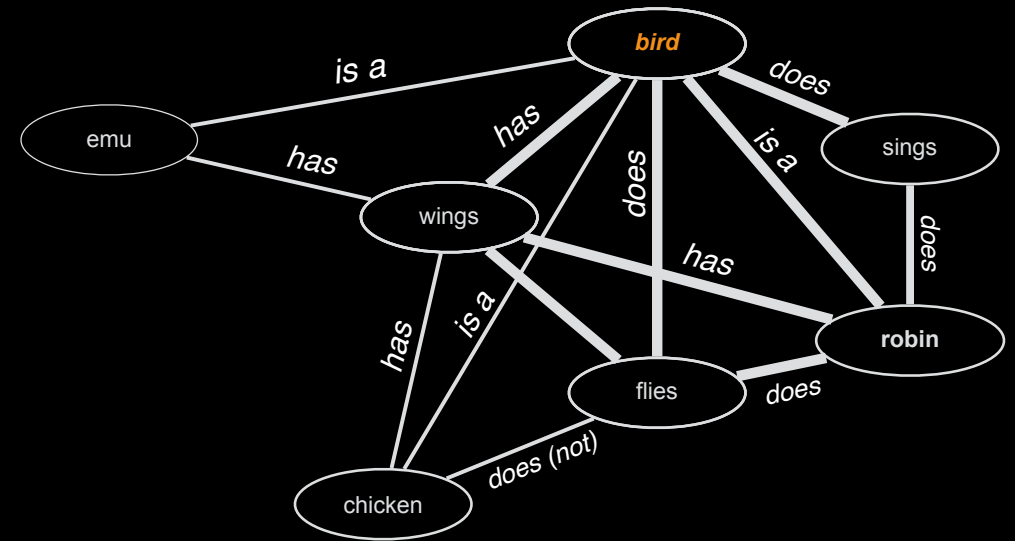
# Backpropagation

---



# Semantic Knowledge

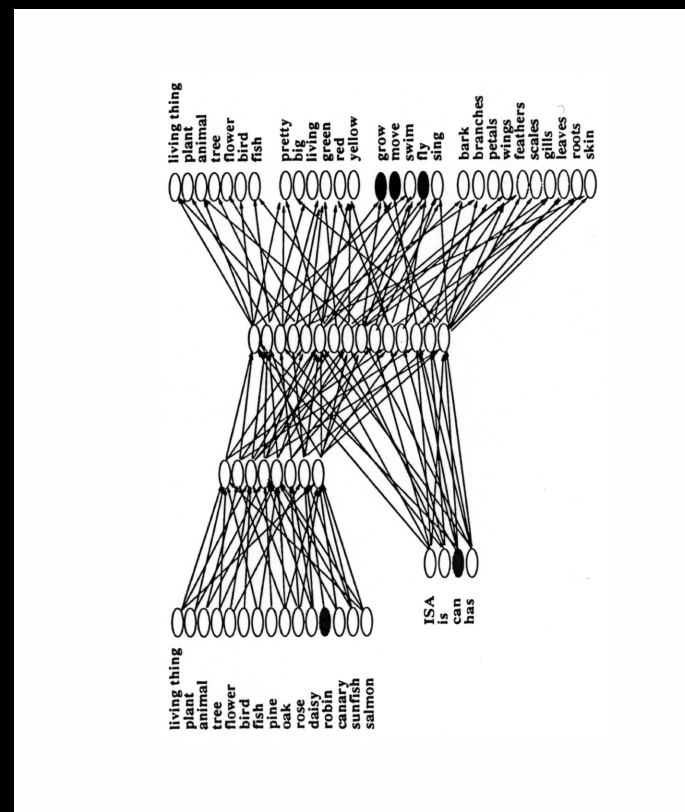
---



# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)



# Semantic Knowledge

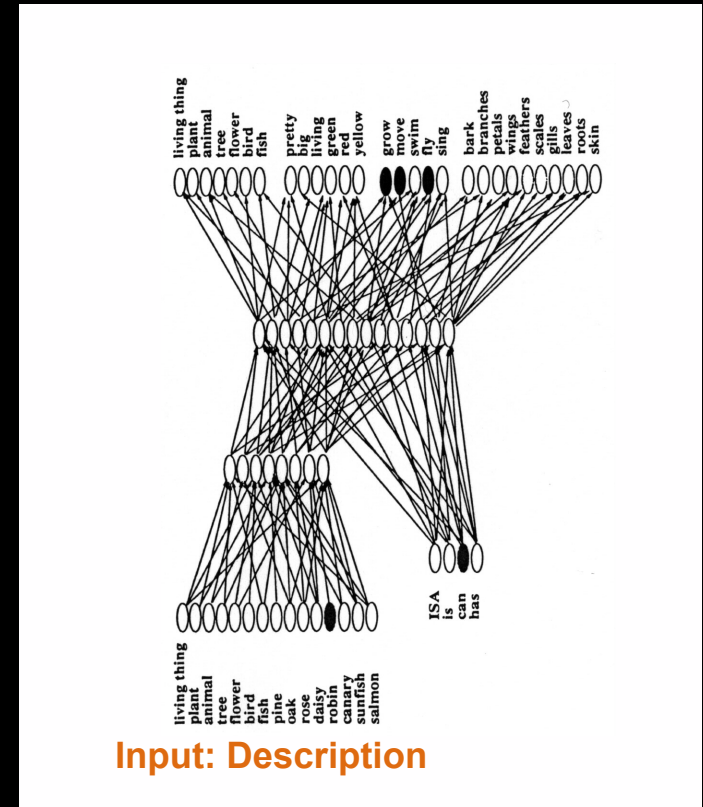
## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- Inputs:

- Descriptions

- ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)



# Semantic Knowledge

## Rumelhart's Semantic Network

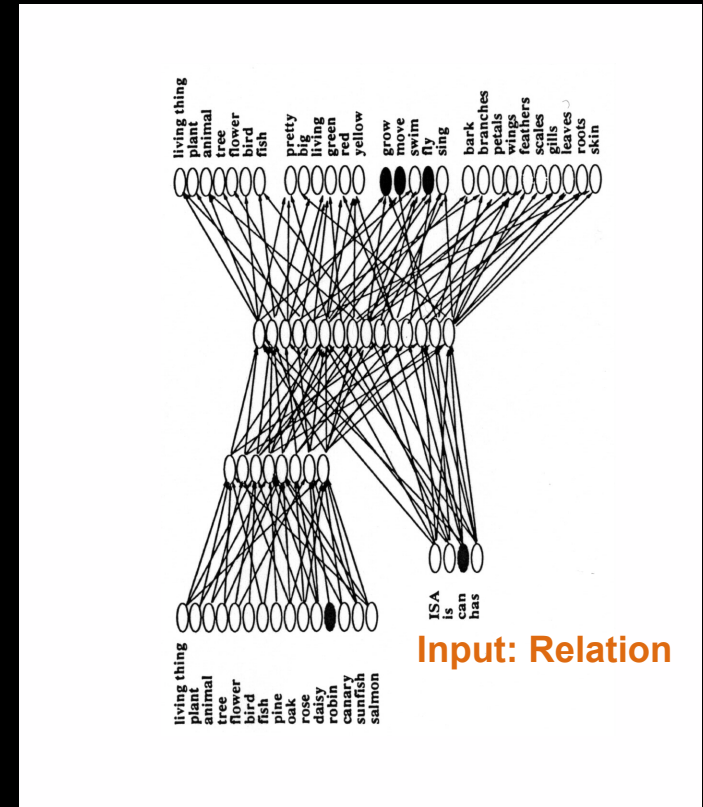
(Rumelhart & Todd, 1987)

- **Inputs:**

- **Descriptions**

- ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)

- **Relations (ISA, is, can, has)**



# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

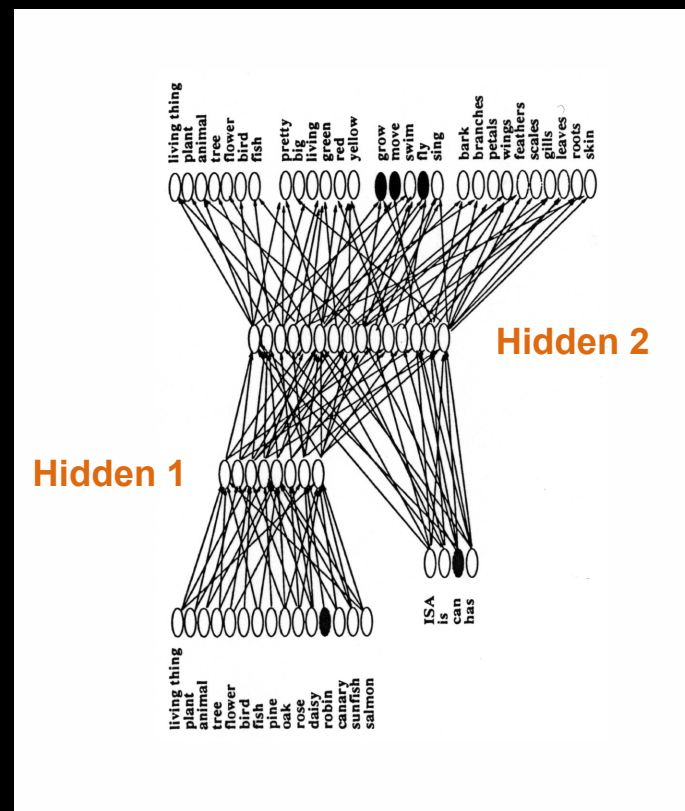
- Inputs:

- Descriptions

- ◆ features (mouth, wings, leaves...)
- ◆ category (bird, tree...)
- ◆ actions (moves, lives, flies, grows...)
- ◆ names (robin, pine, emu..)

- Relations (ISA, is, can, has)

- Hidden Layers: associative

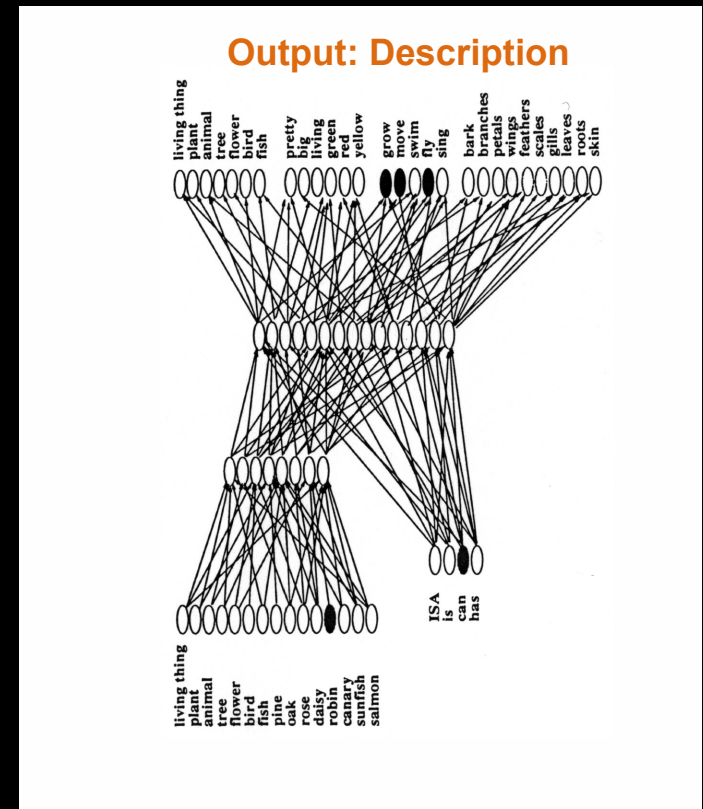


# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- **Inputs:**
  - **Descriptions**
    - ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)
  - **Relations (ISA, is, can, has)**
- **Hidden Layers: associative**
- **Output:**
  - **Descriptions**  
(similar to set used for input)

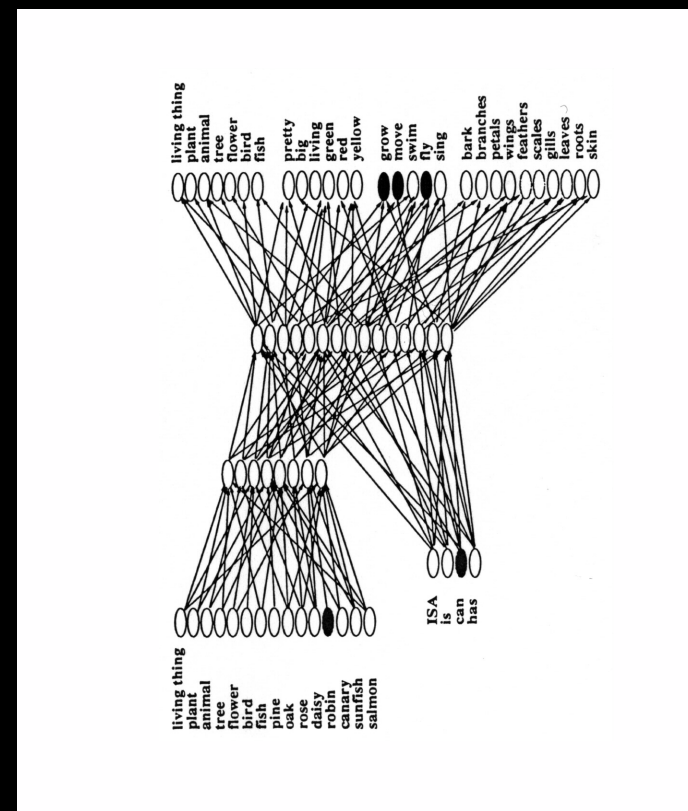


# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- **Inputs:**
  - **Descriptions**
    - ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)
  - **Relations (ISA, is, can, has)**
- **Hidden Layers: associative**
- **Output:**
  - **Descriptions**  
(similar to set used for input)
- **Training:**
  - **Exemplars of birds and trees**

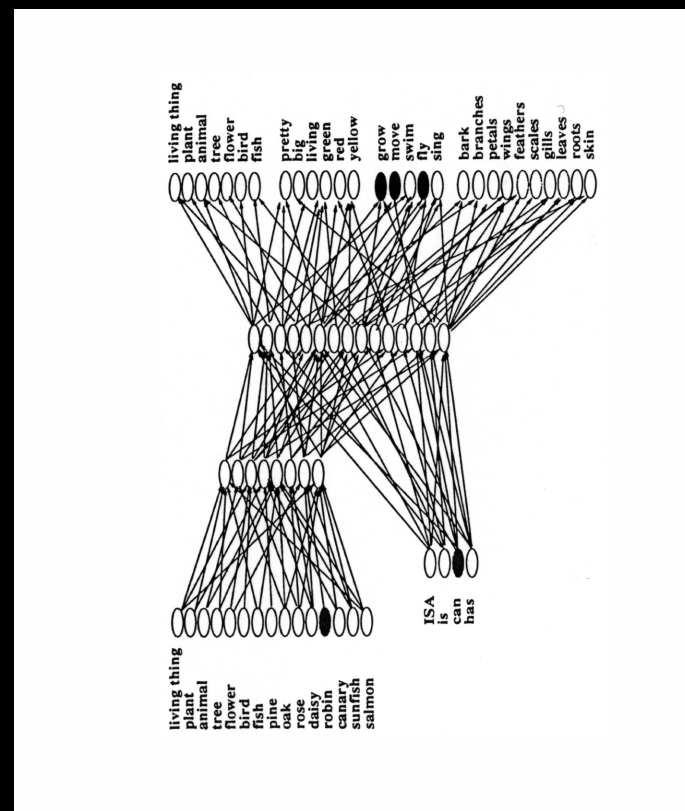


# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- **Inputs:**
  - **Descriptions**
    - ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)
  - **Relations (ISA, is, can, has)**
- **Hidden Layers: associative**
- **Output:**
  - **Descriptions**  
(similar to set used for input)
- **Training:**
  - Exemplars of birds and trees
- **Results: developed a set of internal representations that:**
  - Had distinct recognizable patterns for different categories (e.g., birds vs. trees)

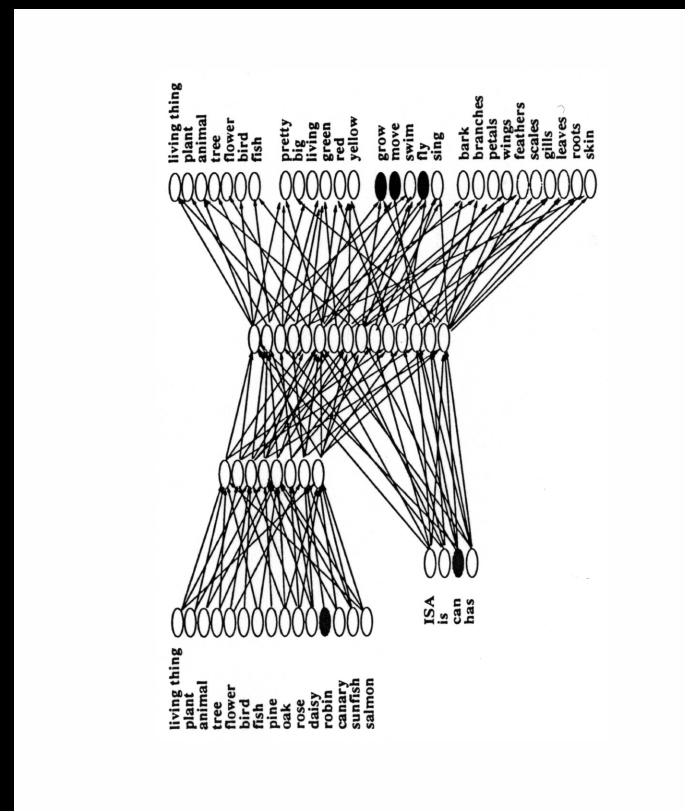


# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- **Inputs:**
  - **Descriptions**
    - ◆ features (mouth, wings, leaves...)
    - ◆ category (bird, tree...)
    - ◆ actions (moves, lives, flies, grows...)
    - ◆ names (robin, pine, emu..)
  - **Relations (ISA, is, can, has)**
- **Hidden Layers: associative**
- **Output:**
  - **Descriptions**  
(similar to set used for input)
- **Training:**
  - Exemplars of birds and trees
- **Results: developed a set of internal representations that:**
  - Had distinct recognizable patterns for different categories (e.g., birds vs. trees)
  - Assigned dedicated representations to exceptional cases



# Semantic Knowledge

## Rumelhart's Semantic Network

(Rumelhart & Todd, 1987)

- **Inputs:**

- **Descriptions**

- ◆ features (mouth, wings, leaves...)
- ◆ category (bird, tree...)
- ◆ actions (moves, lives, flies, grows...)
- ◆ names (robin, pine, emu..)

- **Relations (ISA, is, can, has)**

- **Hidden Layers: associative**

- **Output:**

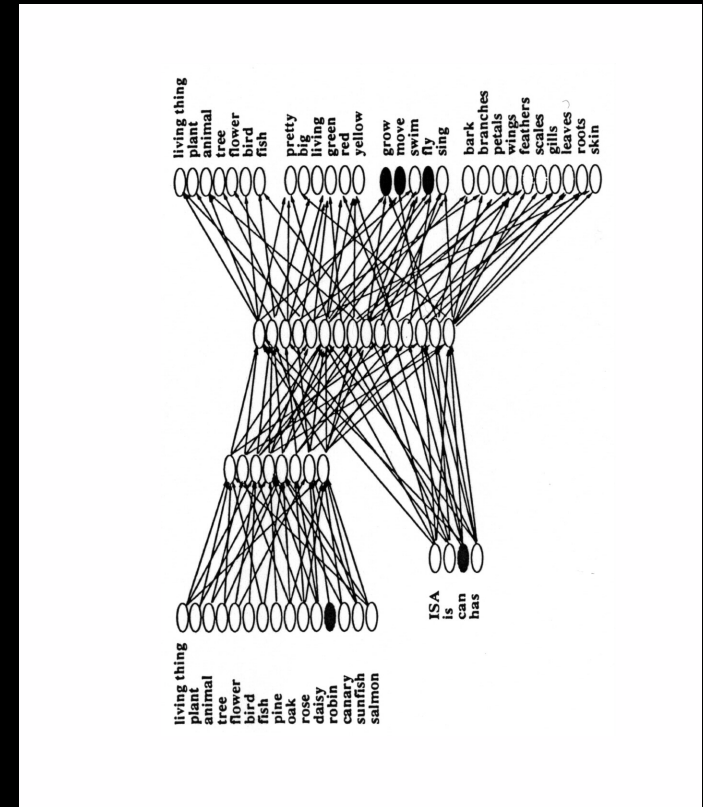
- **Descriptions**  
(similar to set used for input)

- **Training:**

- Exemplars of birds and trees

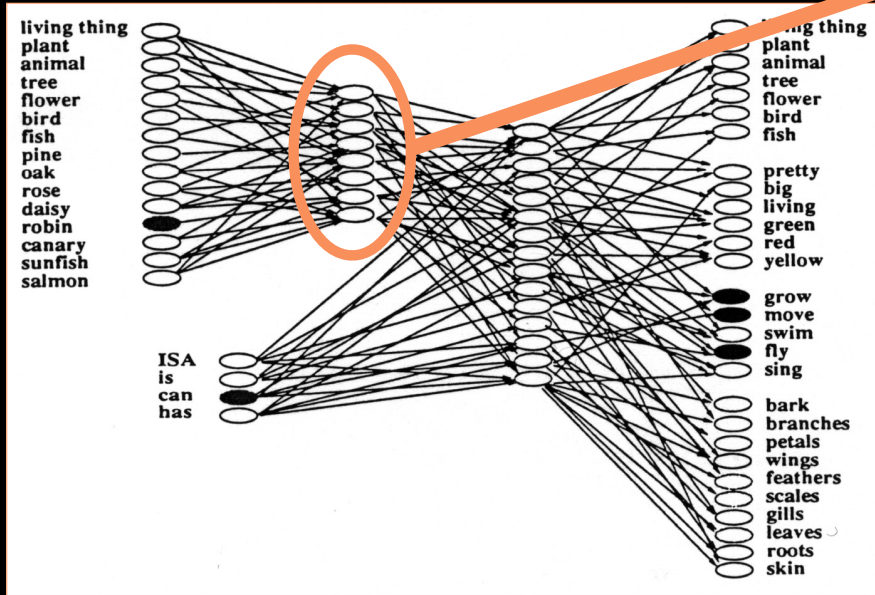
- **Results: developed a set of internal representations that:**

- Had distinct recognizable patterns for different categories (e.g., birds vs. trees)
- Assigned dedicated representations to exceptional cases
- Could properly categorize novel exemplars if they were “typical”

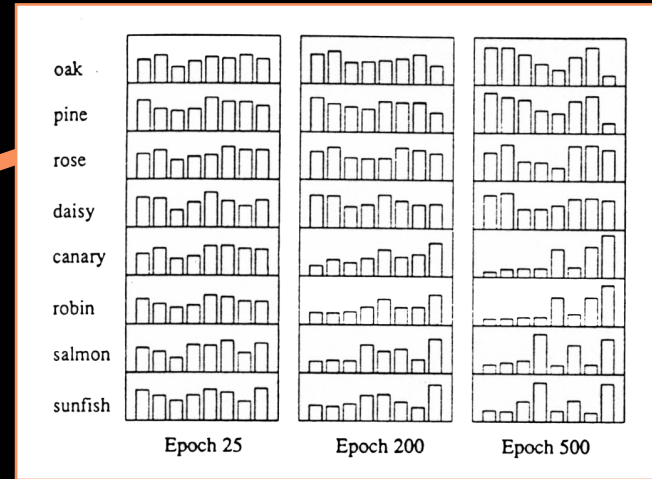


# Rumelhart's Semantic Network

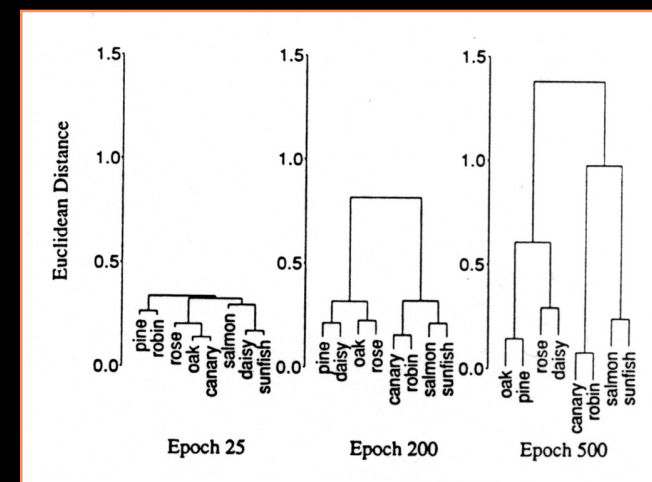
## Rumelhart's Semantic Network



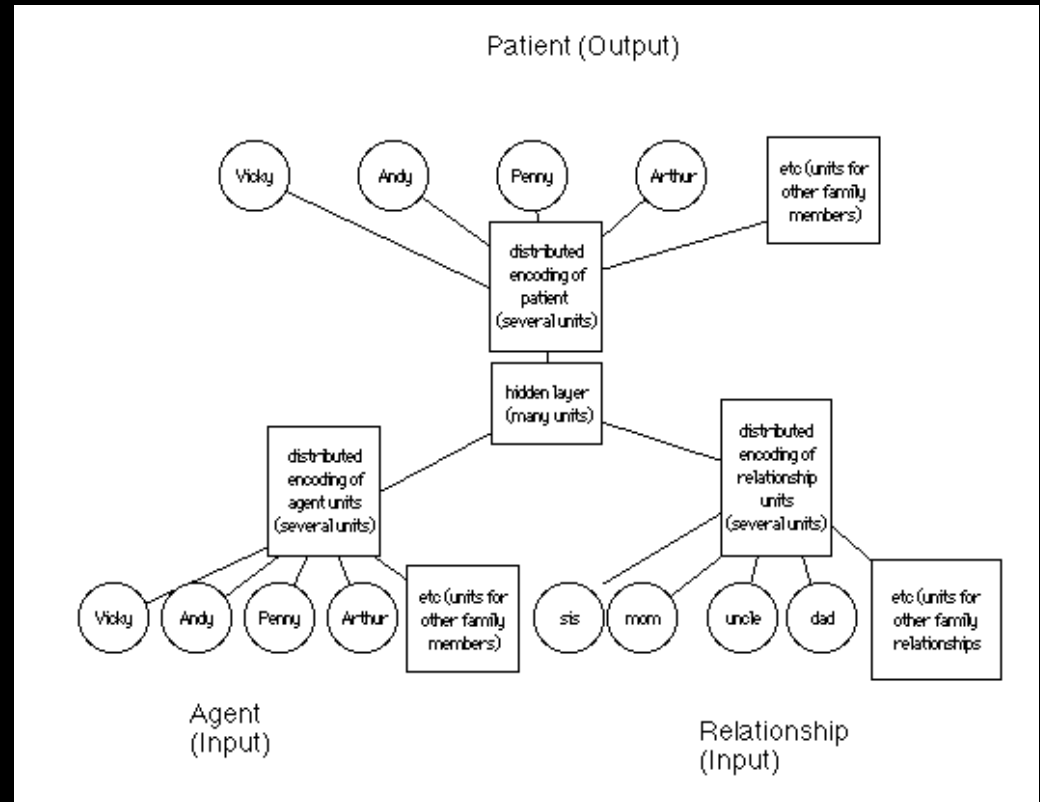
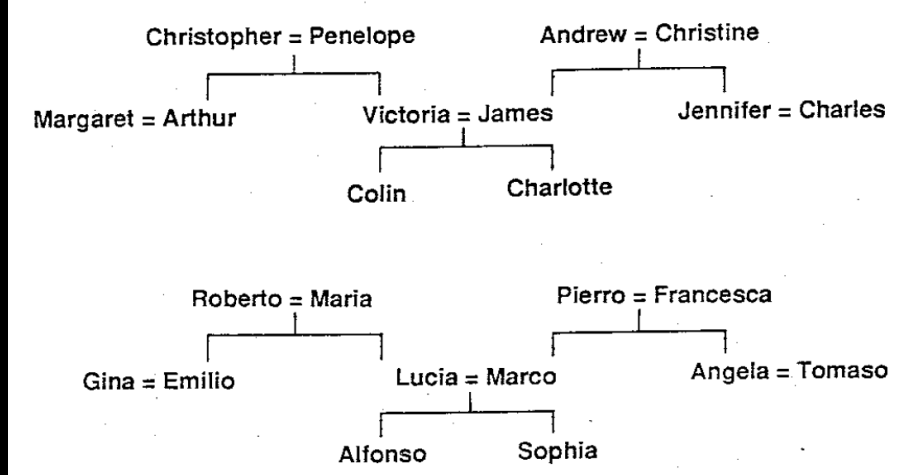
With sufficient training...



acquired rich semantic structure



# Hinton's Family Trees Model



# “Deep Learning”



**Dave Rumelhart**  
(UCSD / Stanford)  
**Backpropagation**



**Yann LeCun**  
(NYU)  
**Convolutional Networks**



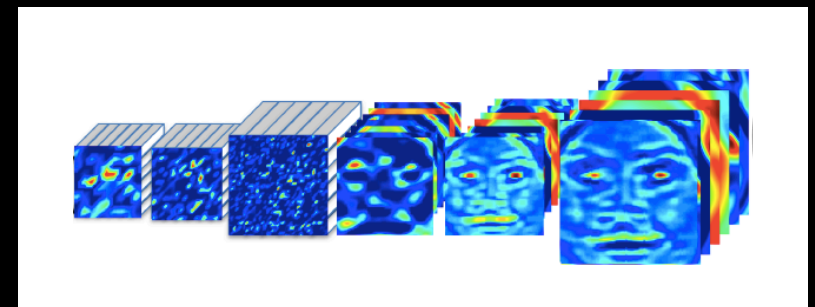
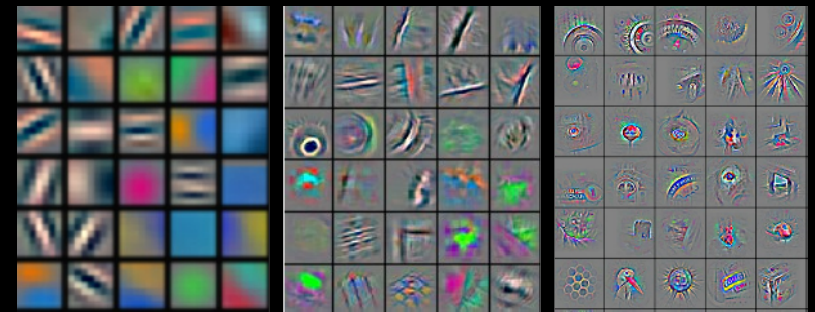
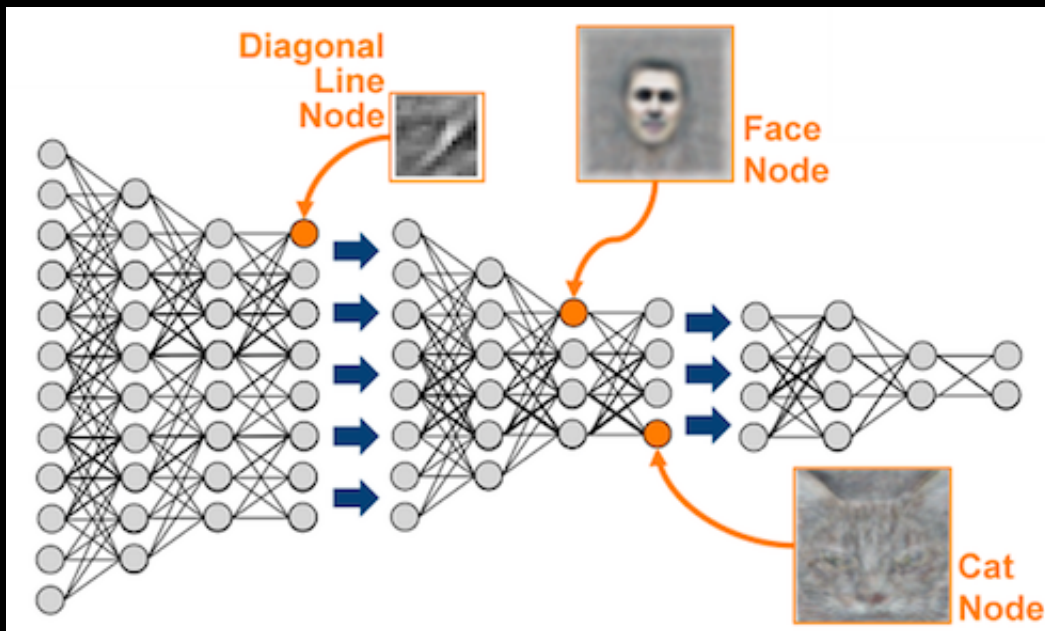
**Geoff Hinton**  
(Toronto / Google)  
**Alex Net**



**Bruno Olshausen**  
(Redwood Institute)  
**Sparse Coding**



**Fei Fei Li**  
(Princeton/Stanford)  
**Image Net**



# “Deep Learning”

## Internal representations

