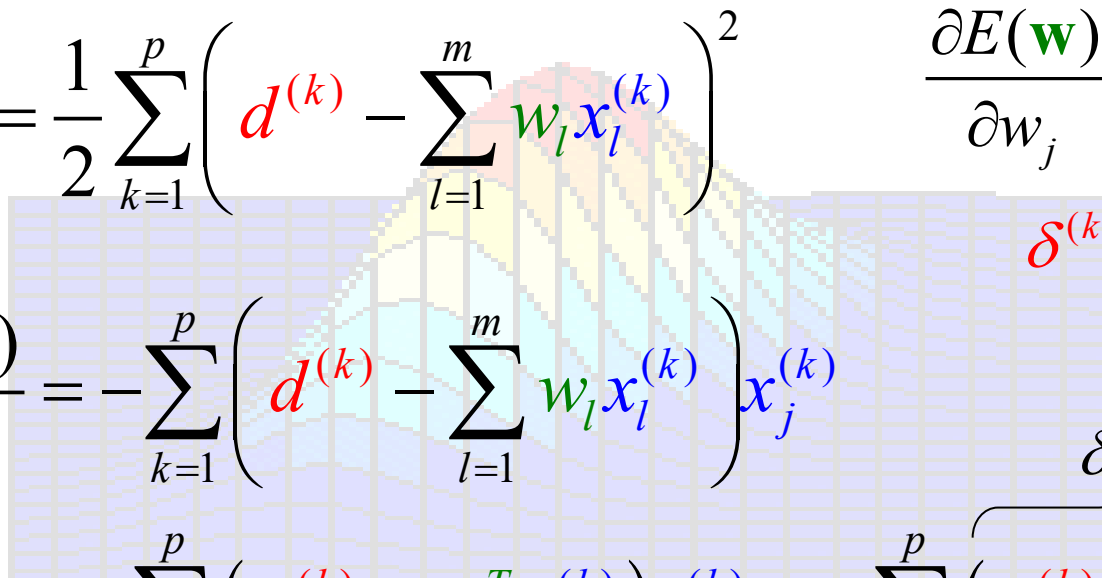


# LMS (Least Mean Square)

Minimize the **cost** function (**error** function):


$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p \left( d^{(k)} - \sum_{l=1}^m w_l x_l^{(k)} \right)^2$$
$$\frac{\partial E(\mathbf{w})}{\partial w_j} = - \sum_{k=1}^p \delta^{(k)} x_j^{(k)}$$
$$\delta^{(k)} = d^{(k)} - y^{(k)}$$
$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial w_j} &= - \sum_{k=1}^p \left( d^{(k)} - \sum_{l=1}^m w_l x_l^{(k)} \right) x_j^{(k)} \\ &= - \sum_{k=1}^p \left( d^{(k)} - \mathbf{w}^T \mathbf{x}^{(k)} \right) x_j^{(k)} = - \sum_{k=1}^p \left( \overbrace{d^{(k)} - y^{(k)}}^{\delta^{(k)}} \right) x_j^{(k)} \end{aligned}$$

# Adaline Learning Rule

Minimize the **cost** function (**error** function):

$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p \left( d^{(k)} - \sum_{l=1}^m w_l x_l^{(k)} \right)^2$$
$$\frac{\partial E(\mathbf{w})}{\partial w_j} = - \sum_{k=1}^p \delta^{(k)} x_j^{(k)}$$
$$\delta^{(k)} = d^{(k)} - y^{(k)}$$

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \left( \frac{\partial E(\mathbf{w})}{\partial w_1}, \frac{\partial E(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial E(\mathbf{w})}{\partial w_m} \right)^T$$

$$\Delta \mathbf{w} = -\eta \nabla_{\mathbf{w}} E(\mathbf{w}) \quad \text{--- Weight Modification Rule}$$

# Learning Modes

- Batch Learning Mode:

$$\Delta w_j = \eta \sum_{k=1}^p \delta^{(k)} x_j^{(k)}$$

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = - \sum_{k=1}^p \delta^{(k)} x_j^{(k)}$$

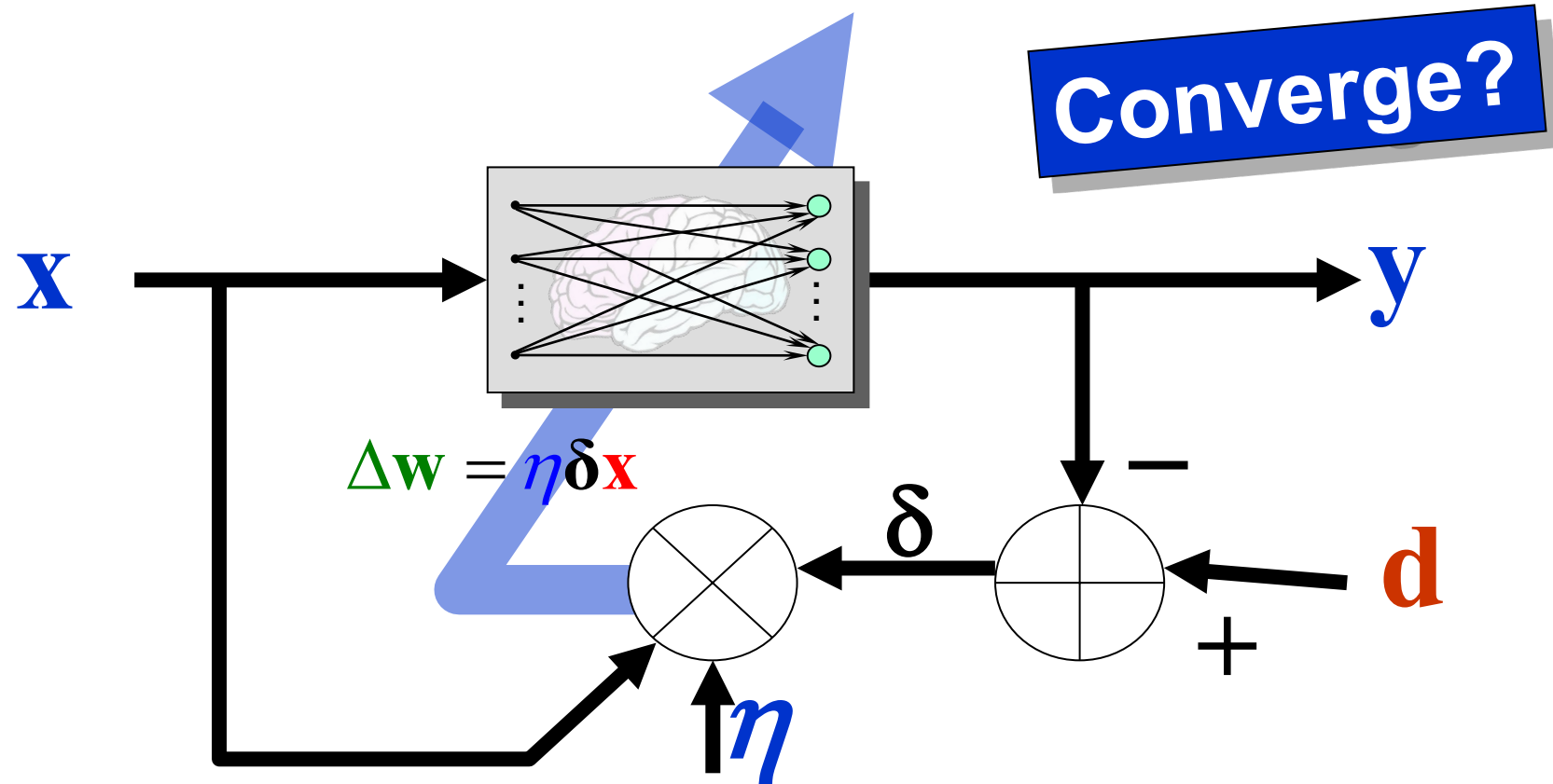
- Incremental Learning Mode:

$$\delta^{(k)} = d^{(k)} - y^{(k)}$$

$$\Delta w_j = \eta \delta^{(k)} x_j^{(k)}$$

# Summary – Adaline Learning Rule

$\delta$ -Learning Rule  
LMS Algorithm  
Widrow-Hoff Learning Rule



# LMS Convergence

Based on the **independence theory** (Widrow, 1976).

1. The **successive input** vectors are statistically independent.
2. At time  $t$ , the **input vector  $\mathbf{x}(t)$**  is statistically independent of all **previous samples of the desired response**, namely  $d(1)$ ,  $d(2)$ , ...,  $d(t-1)$ .
3. At time  $t$ , the **desired response  $d(t)$**  is dependent on  $\mathbf{x}(t)$ , but statistically independent of all **previous values of the desired response**.
4. The input vector  $\mathbf{x}(t)$  and desired response  $d(t)$  are drawn from **Gaussian distributed populations**.

# LMS Convergence

It can be shown that LMS is convergent if

$$0 < \eta < \frac{2}{\lambda_{\max}}$$

where  $\lambda_{\max}$  is the largest eigenvalue of the correlation matrix  $\mathbf{R}_x$  for the inputs.

$$\mathbf{R}_x = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^{\infty} \mathbf{x}_i \mathbf{x}_i^T$$

# LMS Convergence

It can be shown that LMS is convergent if

$$0 < \eta < \frac{2}{\lambda_{\max}}$$

Where  $\lambda_{\max}$  is the largest eigenvalue of the correlation matrix  $\mathbf{R}_x$  for the inputs.

$$\mathbf{R}_x = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^{\infty} \mathbf{x}_i \mathbf{x}_i^T$$

Since  $\lambda_{\max}$  is hardly available, we commonly use

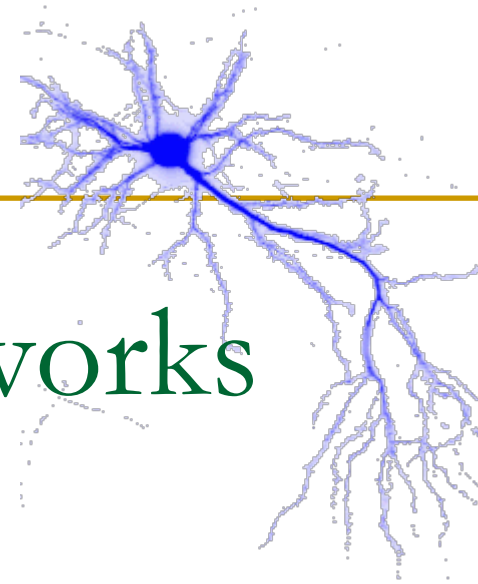
$$0 < \eta < \frac{2}{\text{tr}(\mathbf{R}_x)}$$

# Comparisons

	Perceptron Learning Rule	Adaline Learning Rule (Widrow-Hoff)
Fundamental	Hebbian Assumption	Gradient Descent
Convergence	In finite steps	Converge Asymptotically
Constraint	Linearly Separable	Linear Independence



# Feed-Forward Neural Networks

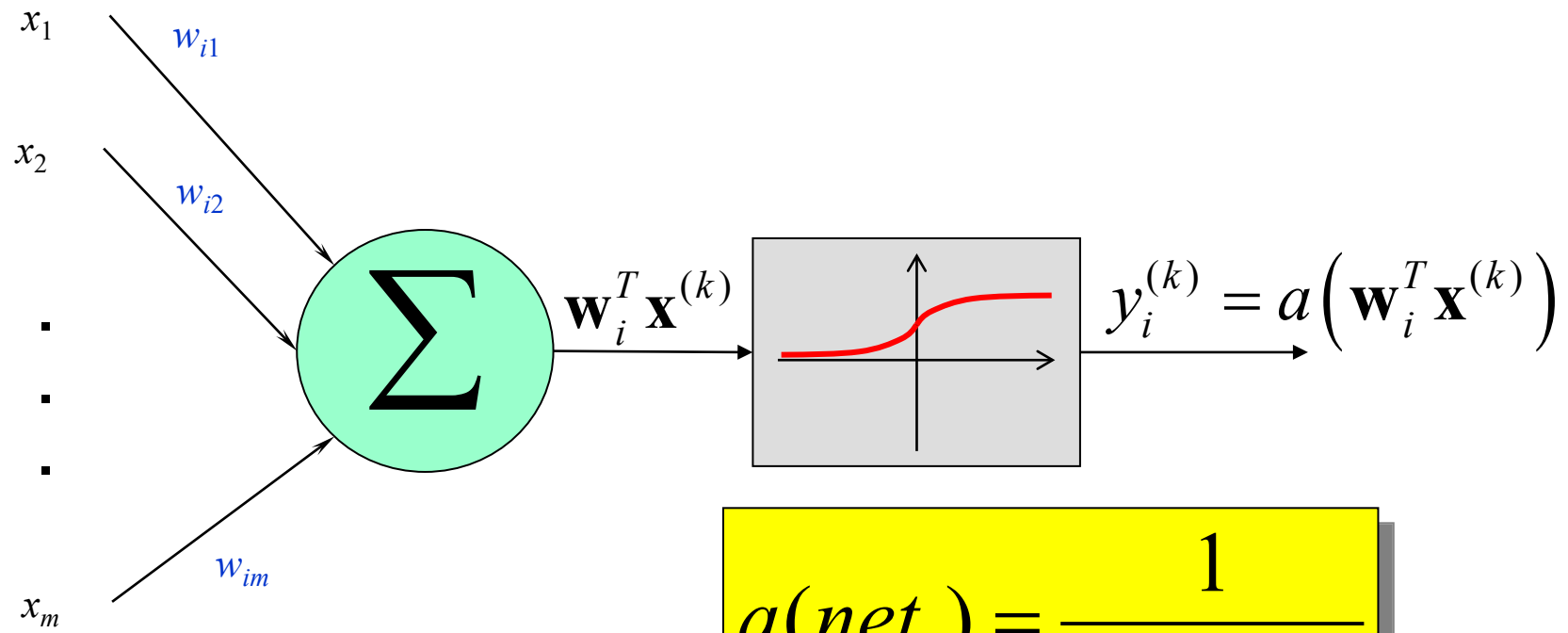


---

## Learning Rules for Single-Layered Perceptron Networks

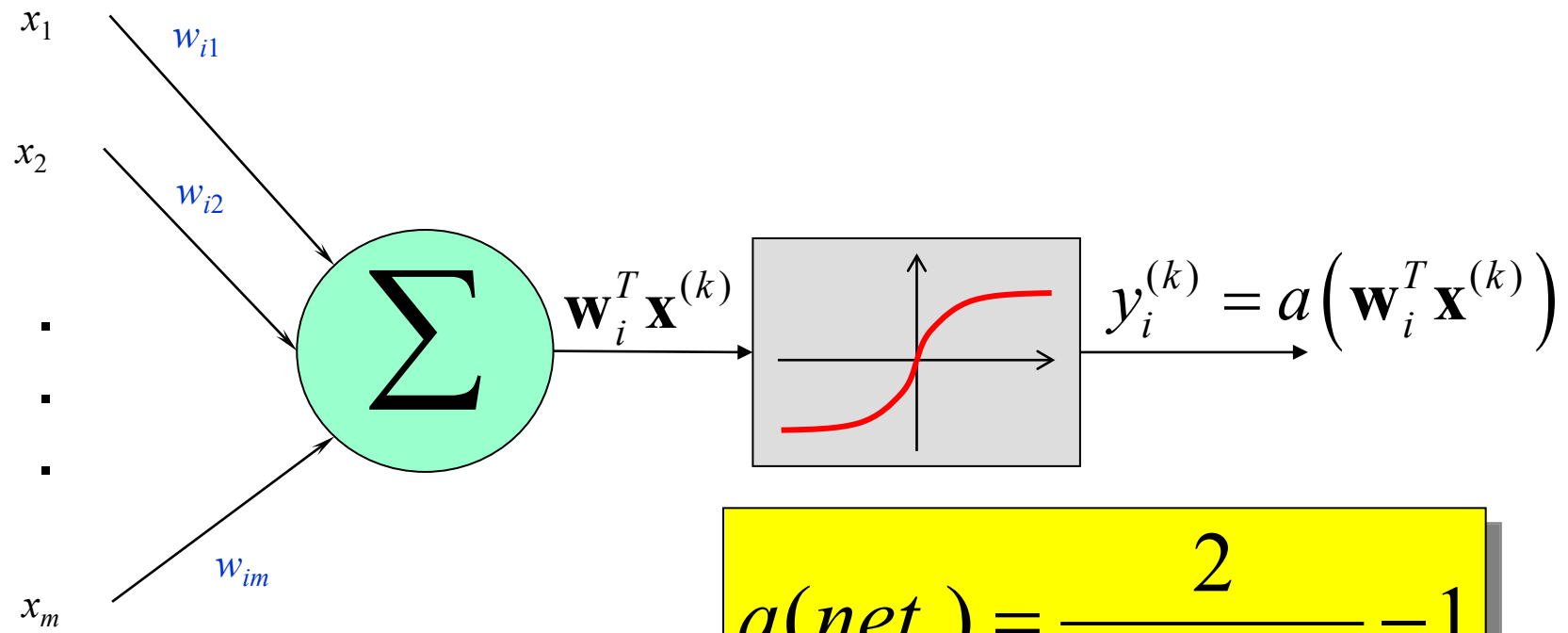
- Perceptron Learning Rule
- Adaline Learning Rule
- $\delta$ -Learning Rule

# Unipolar Sigmoid



$$a(net_i) = \frac{1}{1 + e^{-\lambda net_i}}$$

# Bipolar Sigmoid



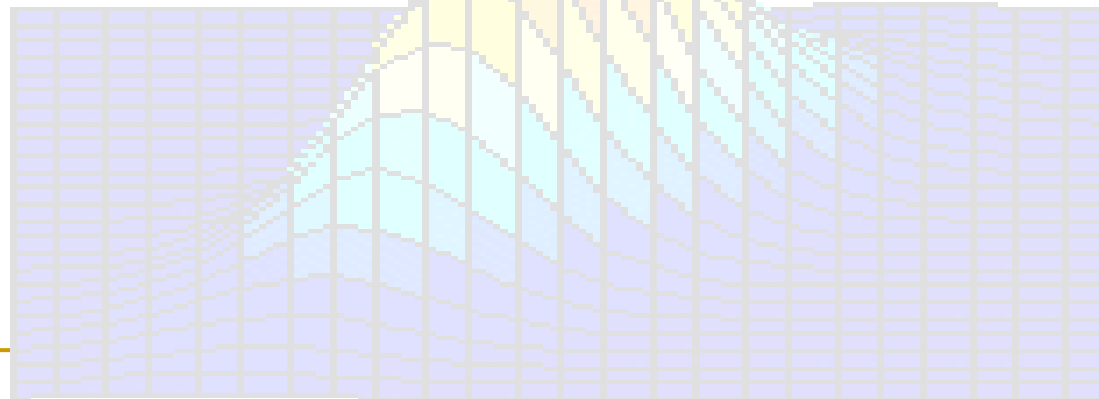
$$a(net_i) = \frac{2}{1 + e^{-\lambda net_i}} - 1$$

# Goal

Minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p (\mathbf{d}^{(k)} - \mathbf{y}^{(k)})^2$$

$$= \frac{1}{2} \sum_{k=1}^p [\mathbf{d}^{(k)} - a(\mathbf{w}^T \mathbf{x}^{(k)})]^2$$



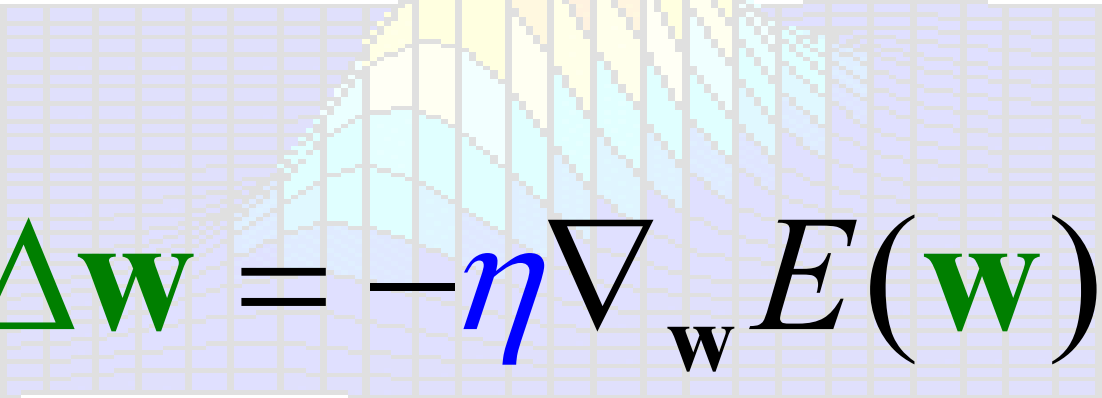
# Gradient Descent Algorithm

$$\nabla_w E(\mathbf{w}) = \left( \frac{\partial E(\mathbf{w})}{\partial w_1}, \frac{\partial E(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial E(\mathbf{w})}{\partial w_m} \right)^T$$

Minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p (d^{(k)} - y^{(k)})^2$$

$$= \frac{1}{2} \sum_{k=1}^p [d^{(k)} - a(\mathbf{w}^T \mathbf{x}^{(k)})]^2$$


$$\Delta \mathbf{w} = -\eta \nabla_w E(\mathbf{w})$$

# The Gradient

$$\nabla_w E(\mathbf{w}) = \left( \frac{\partial E(\mathbf{w})}{\partial w_1}, \frac{\partial E(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial E(\mathbf{w})}{\partial w_m} \right)^T$$

$$y^{(k)} = a(\mathbf{w}^T \mathbf{x}^{(k)})$$

Minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p (d^{(k)} - y^{(k)})^2$$

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial w_j} &= - \sum_{k=1}^p (d^{(k)} - y^{(k)}) \frac{\partial y^{(k)}}{\partial w_j} \\ &= - \sum_{k=1}^p (d^{(k)} - y^{(k)}) \underbrace{\frac{\partial a(\text{net}^{(k)})}{\partial \text{net}^{(k)}}}_{?} \underbrace{\frac{\partial \text{net}^{(k)}}{\partial w_j}}_{?} \end{aligned}$$

Depends on the activation function used.

$$\text{net}^{(k)} = \mathbf{w}^T \mathbf{x}^{(k)} = \sum_{i=1}^m w_i x_i^{(k)} \Rightarrow \frac{\partial \text{net}^{(k)}}{\partial w_j} = x_j^{(k)}$$

# Weight Modification Rule

$$\nabla_w E(\mathbf{w}) = \left( \frac{\partial E(\mathbf{w})}{\partial w_1}, \frac{\partial E(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial E(\mathbf{w})}{\partial w_m} \right)^T$$

$$y^{(k)} = a(\text{net}^{(k)})$$

Minimize  $E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p (d^{(k)} - y^{(k)})^2$

$$\delta^{(k)} = d^{(k)} - y^{(k)}$$

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = - \sum_{k=1}^p (d^{(k)} - y^{(k)}) x_j^{(k)} \frac{\partial a(\text{net}^{(k)})}{\partial \text{net}^{(k)}}$$

Batch

$$\Delta w_j = \eta \sum_{k=1}^p \delta^{(k)} x_j^{(k)} \frac{\partial a(\text{net}^{(k)})}{\partial \text{net}^{(k)}}$$

Learning  
Rule

Incremental

$$\Delta w_j = \eta \delta^{(k)} x_j^{(k)} \frac{\partial a(\text{net}^{(k)})}{\partial \text{net}^{(k)}}$$

# The Learning Efficacy

$$\nabla_w E(\mathbf{w}) = \left( \frac{\partial E(\mathbf{w})}{\partial w_1}, \frac{\partial E(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial E(\mathbf{w})}{\partial w_m} \right)^T$$

$$y^{(k)} = a(\text{net}^{(k)})$$

Minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p (d^{(k)} - y^{(k)})^2$$

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = - \sum_{k=1}^p (d^{(k)} - y^{(k)}) x_j^{(k)} \frac{\partial a(\text{net}^{(k)})}{\partial \text{net}^{(k)}}$$

Sigmoid

Adaline

$$a(\text{net}) = \text{net}$$

$$\frac{\partial a(\text{net})}{\partial \text{net}} = 1$$

Unipolar

$$a(\text{net}) = \frac{1}{1 + e^{-\lambda \text{net}}}$$

$$\frac{\partial a(\text{net})}{\partial \text{net}} = \lambda y^{(k)} (1 - y^{(k)})$$

Bipolar

$$a(\text{net}) = \frac{2}{1 + e^{-\lambda \text{net}}} - 1$$

**Exercise**



# Comparisons

Adaline

Batch

$$\Delta w_j = \eta \sum_{k=1}^p \delta^{(k)} x_j^{(k)}$$

Incremental

$$\Delta w_j = \eta \delta^{(k)} x_j^{(k)}$$

Sigmoid

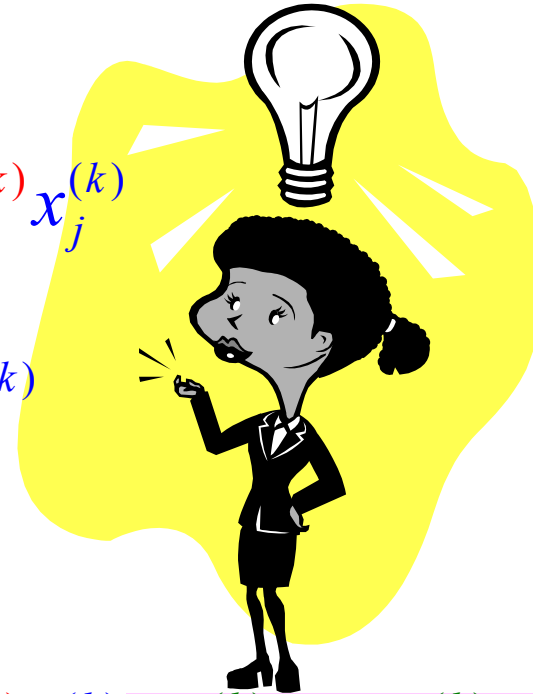
Batch

$$\Delta w_j = \eta \sum_{k=1}^p \delta^{(k)} x_j^{(k)} \lambda y^{(k)} (1 - y^{(k)})$$

Incremental

$$\Delta w_j = \eta \delta^{(k)} x_j^{(k)} \lambda y^{(k)} (1 - y^{(k)})$$

$$\lambda y^{(k)} (1 - y^{(k)})$$



# Training of Artificial Neurons

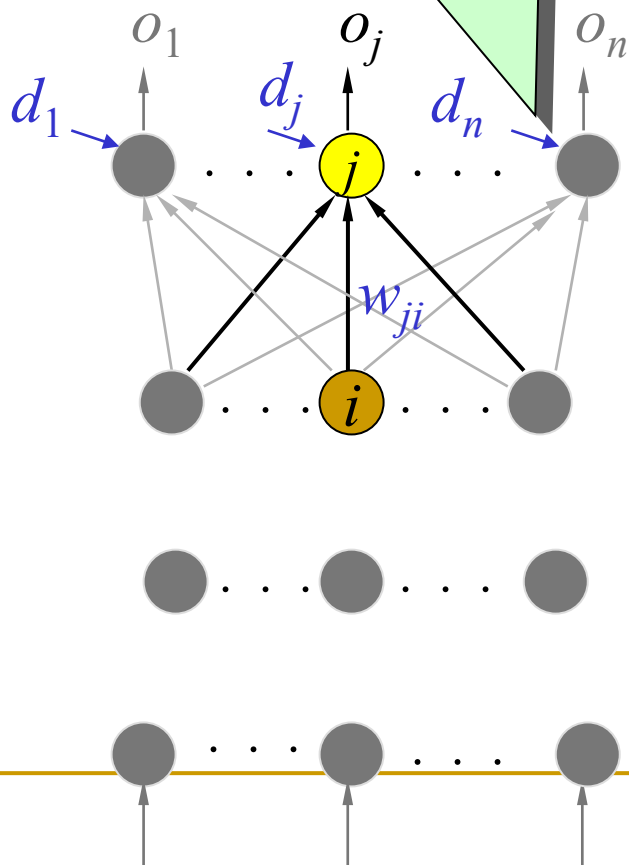
$$\delta_j^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} = -(d_j^{(l)} - o_j^{(l)}) \lambda o_j^{(l)} (1 - o_j^{(l)})$$

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$

$$o_j^{(l)} = a(\text{net}_j^{(l)})$$

$$\text{net}_j^{(l)} = \sum w_{ji} o_i^{(l)}$$



$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ji}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ji}} = \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} \frac{\partial \text{net}_j^{(l)}}{\partial w_{ji}}$$

$$\frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} = \frac{\partial E^{(l)}}{\partial o_j^{(l)}} \frac{\partial o_j^{(l)}}{\partial \text{net}_j^{(l)}}$$

$$-(d_j^{(l)} - o_j^{(l)})$$

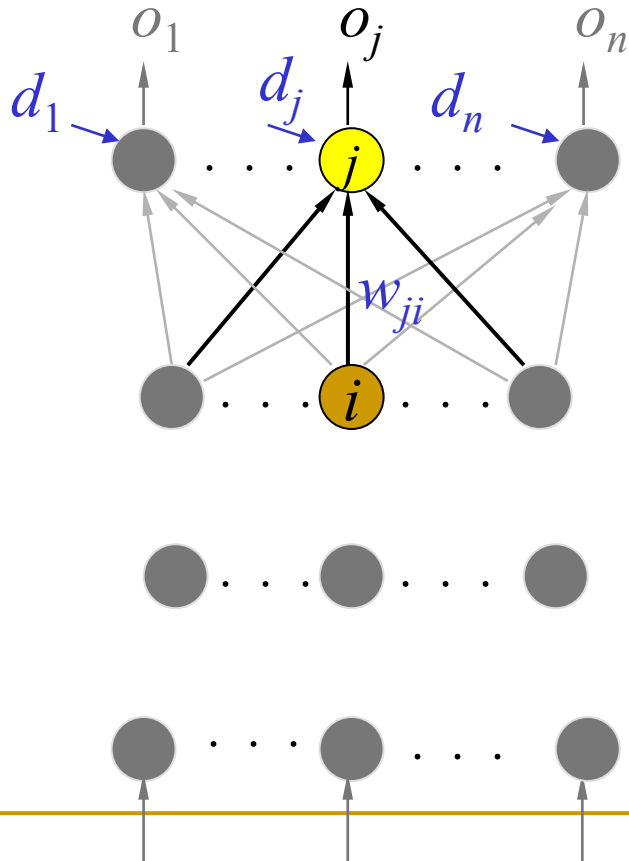
Using sigmoid,

$$\lambda o_j^{(l)} (1 - o_j^{(l)})$$

# Learning on Output Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$



$$o_j^{(l)} = a(\text{net}_j^{(l)})$$

$$\text{net}_j^{(l)} = \sum w_{ji} o_i^{(l)}$$

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ji}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ji}} = \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} \boxed{\frac{\partial \text{net}_j^{(l)}}{\partial w_{ji}}} \rightarrow o_i^{(l)}$$

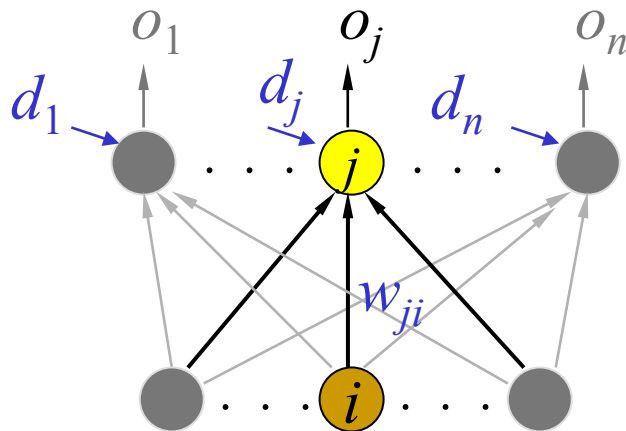
$$\frac{\partial E^{(l)}}{\partial w_{ji}} = \delta_j^{(l)} o_i^{(l)}$$

$$= -(d_j^{(l)} - o_j^{(l)}) \lambda o_j^{(l)} (1 - o_j^{(l)}) o_i^{(l)}$$

# Learning on Output Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$



$$o_j^{(l)} = a(\text{net}_j^{(l)})$$

$$\text{net}_j^{(l)} = \sum w_{ji} o_i^{(l)}$$

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial}{\partial \text{net}_j^{(l)}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}}$$

*How to train the weights connecting to output neurons?*

$$\frac{\partial E}{\partial w_{ji}} = \sum_{l=1}^p \delta_j^{(l)} o_i^{(l)}$$

$$\frac{\partial E^{(l)}}{\partial w_{ji}} = \delta_j^{(l)} o_i^{(l)}$$

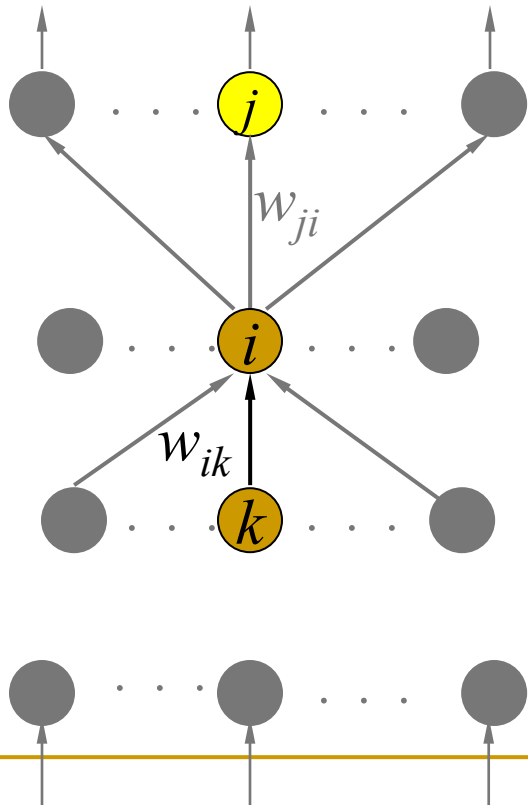
$$\Delta w_{ji} = -\eta \sum_{l=1}^p \delta_j^{(l)} o_i^{(l)}$$

$$= -(d_j^{(l)} - o_j^{(l)}) \lambda o_j^{(l)} (1 - o_j^{(l)}) o_i^{(l)}$$

# Learning on Hidden Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$



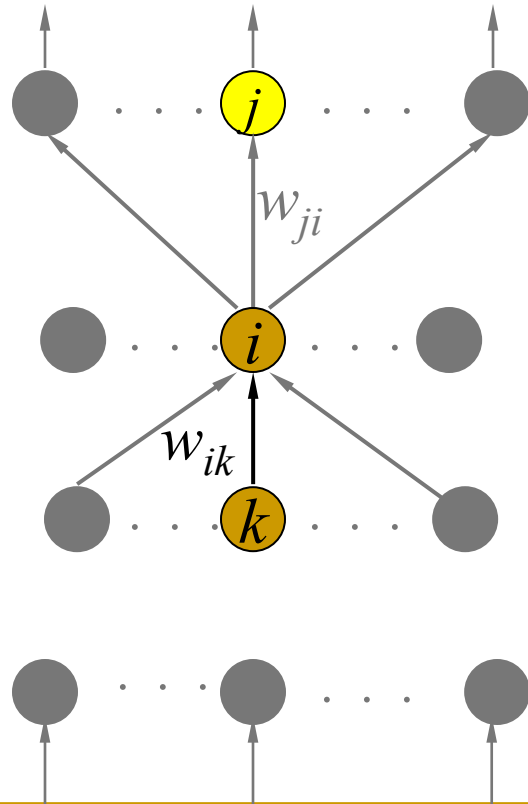
$$\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ik}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ik}} = \underbrace{\frac{\partial E^{(l)}}{\partial net_i^{(l)}}}_{?} \underbrace{\frac{\partial net_i^{(l)}}{\partial w_{ik}}}_{?}$$

# Learning on Hidden Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$



$$\delta_i^{(l)}$$

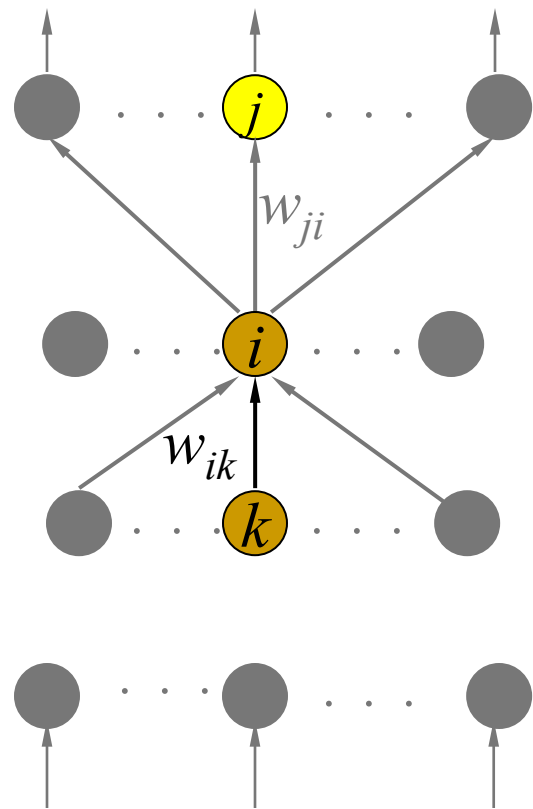
$$\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ik}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ik}} = \frac{\partial E^{(l)}}{\partial net_i^{(l)}} \boxed{\frac{\partial net_i^{(l)}}{\partial w_{ik}}} \rightarrow o_k^{(l)}$$

# Learning on Hidden Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$



$$\delta_i^{(l)}$$

$$\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ik}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ik}} = \frac{\partial E^{(l)}}{\partial net_i^{(l)}} \frac{\partial net_i^{(l)}}{\partial w_{ik}} \rightarrow o_k^{(l)}$$

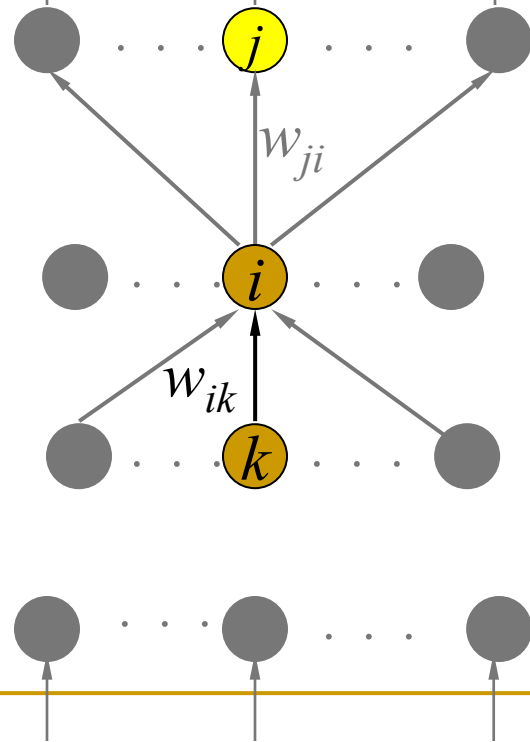
$$\frac{\partial E^{(l)}}{\partial net_i^{(l)}} = \underbrace{\frac{\partial E^{(l)}}{\partial o_i^{(l)}}}_{\delta_i^{(l)}} \frac{\partial o_i^{(l)}}{\partial net_i^{(l)}} \rightarrow \lambda o_i^{(l)} (1 - o_i^{(l)})$$

# Learning on Hidden Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$

$$\delta_i^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}} = \lambda o_i^{(l)} (1 - o_i^{(l)}) \sum_j w_{ji} \delta_j^{(l)}$$

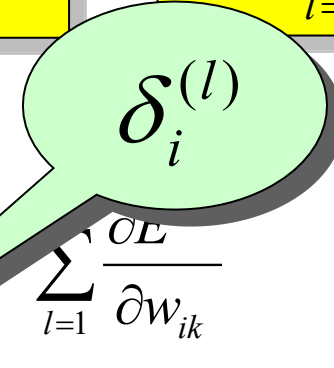


$$\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E}{\partial w_{ik}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ik}} = \frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}} \frac{\partial \text{net}_i^{(l)}}{\partial w_{ik}} \rightarrow o_k^{(l)}$$

$$\frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}} = \frac{\partial E^{(l)}}{\partial o_i^{(l)}} \frac{\partial o_i^{(l)}}{\partial \text{net}_i^{(l)}} \rightarrow \lambda o_i^{(l)} (1 - o_i^{(l)})$$

$$\frac{\partial E^{(l)}}{\partial o_i^{(l)}} = \sum_j \underbrace{\frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}}}_{\delta_j^{(l)}} \underbrace{\frac{\partial \text{net}_j^{(l)}}{\partial o_i^{(l)}}}_{w_{ji}}$$



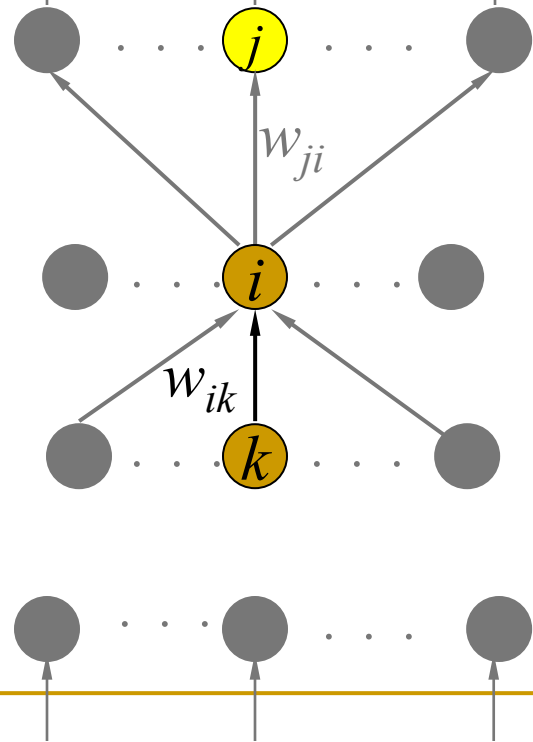


# Learning on Hidden Neurons

$$E^{(l)} = \frac{1}{2} \sum_{j=1}^n [d_j^{(l)} - o_j^{(l)}]^2$$

$$E = \sum_{l=1}^p E^{(l)}$$

$$\delta_i^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}} = \lambda o_i^{(l)} (1 - o_i^{(l)}) \sum_j w_{ji} \delta_j^{(l)}$$



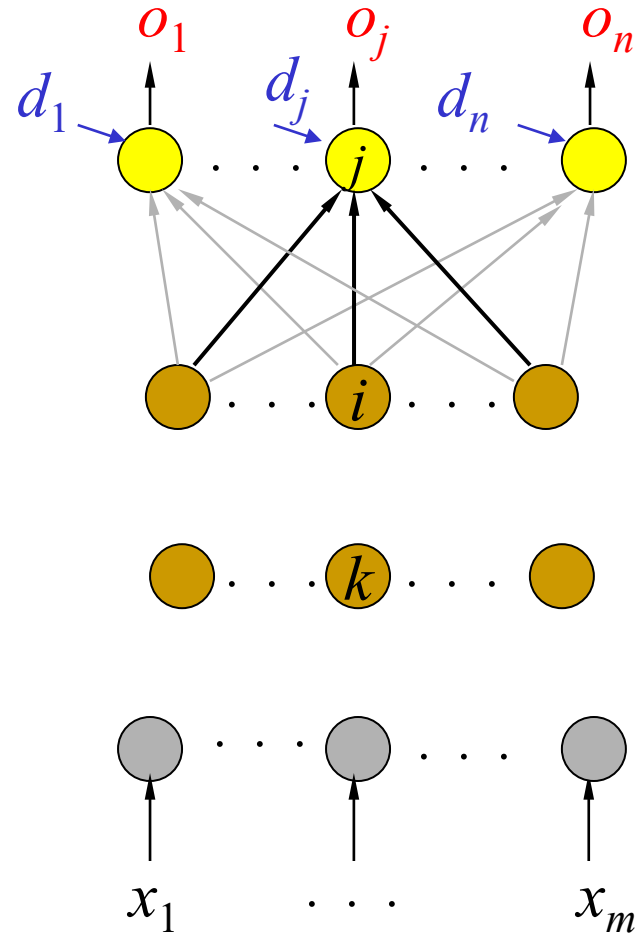
$$\frac{\partial E}{\partial w_{ik}} = \frac{\partial}{\partial w_{ik}} \sum_{l=1}^p E^{(l)} = \sum_{l=1}^p \frac{\partial E^{(l)}}{\partial w_{ik}}$$

$$\frac{\partial E^{(l)}}{\partial w_{ik}} = \boxed{\frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}}} \frac{\partial \text{net}_i^{(l)}}{\partial w_{ik}} \rightarrow o_k^{(l)}$$

$$\frac{\partial E}{\partial w_{ik}} = \sum_{l=1}^p \delta_i^{(l)} o_k^{(l)}$$

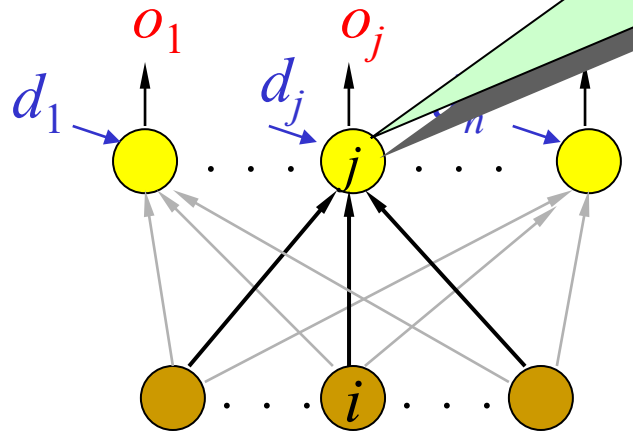
$$\Delta w_{ik} = -\eta \sum_{l=1}^p \delta_i^{(l)} o_k^{(l)}$$

# Back Propagation

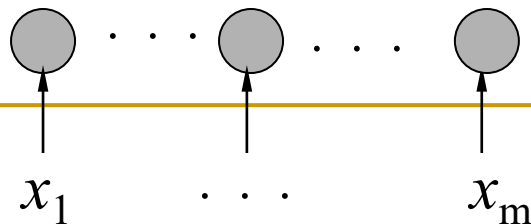


# Back Propagation

$$\delta_j^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} = -\lambda(d_j^{(l)} - o_j^{(l)})o_j^{(l)}(1 - o_j^{(l)})$$

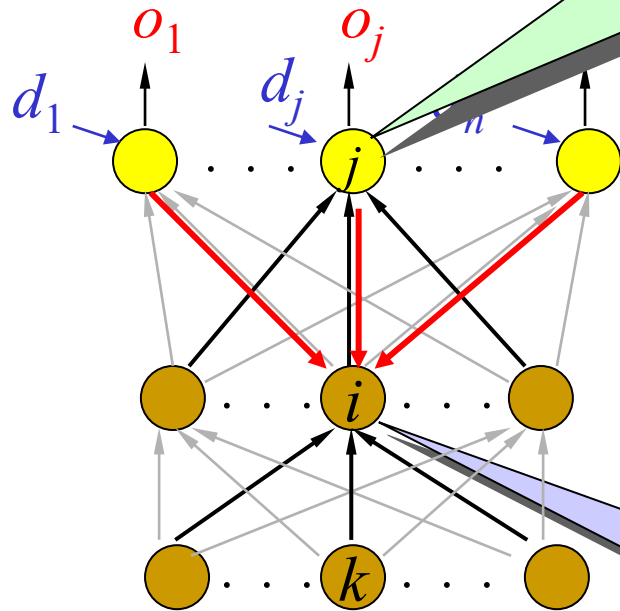


$$\Delta w_{ji} = -\eta \sum_{l=1}^p \delta_j^{(l)} o_i^{(l)}$$



# Back Propagation

$$\delta_j^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_j^{(l)}} = -\lambda(d_j^{(l)} - o_j^{(l)})o_j^{(l)}(1 - o_j^{(l)})$$



$$\Delta w_{ji} = -\eta \sum_{l=1}^p \delta_j^{(l)} o_i^{(l)}$$

$$\Delta w_{ik} = -\eta \sum_{l=1}^p \delta_i^{(l)} o_k^{(l)}$$

$$\delta_i^{(l)} = \frac{\partial E^{(l)}}{\partial \text{net}_i^{(l)}} = \lambda o_i^{(l)} (1 - o_i^{(l)}) \sum_j w_{ji} \delta_j^{(l)}$$

# Learning Factors

- Initial Weights
- Learning Constant ( $\eta$ )
- Cost Functions
- Momentum
- Update Rules
- Training Data and Generalization
- Number of Layers
- Number of Hidden Nodes

# Reading Assignments

- Shi Zhong and Vladimir Cherkassky, “Factors Controlling Generalization Ability of MLP Networks.” In Proc. IEEE Int. Joint Conf. on Neural Networks, vol. 1, pp. 625-630, Washington DC. July 1999. (<http://www.cse.fau.edu/~zhong/pubs.htm>)
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986b). "Learning Internal Representations by Error Propagation," in Parallel Distributed Processing: Explorations in the Microstructure of Cognition, vol. I, D. E. Rumelhart, J. L. McClelland, and the PDP Research Group. MIT Press, Cambridge (1986). (<http://www.cnbc.cmu.edu/~plaut/85-419/papers/RumelhartETAL86.backprop.pdf>).