



Artificial Neural Network

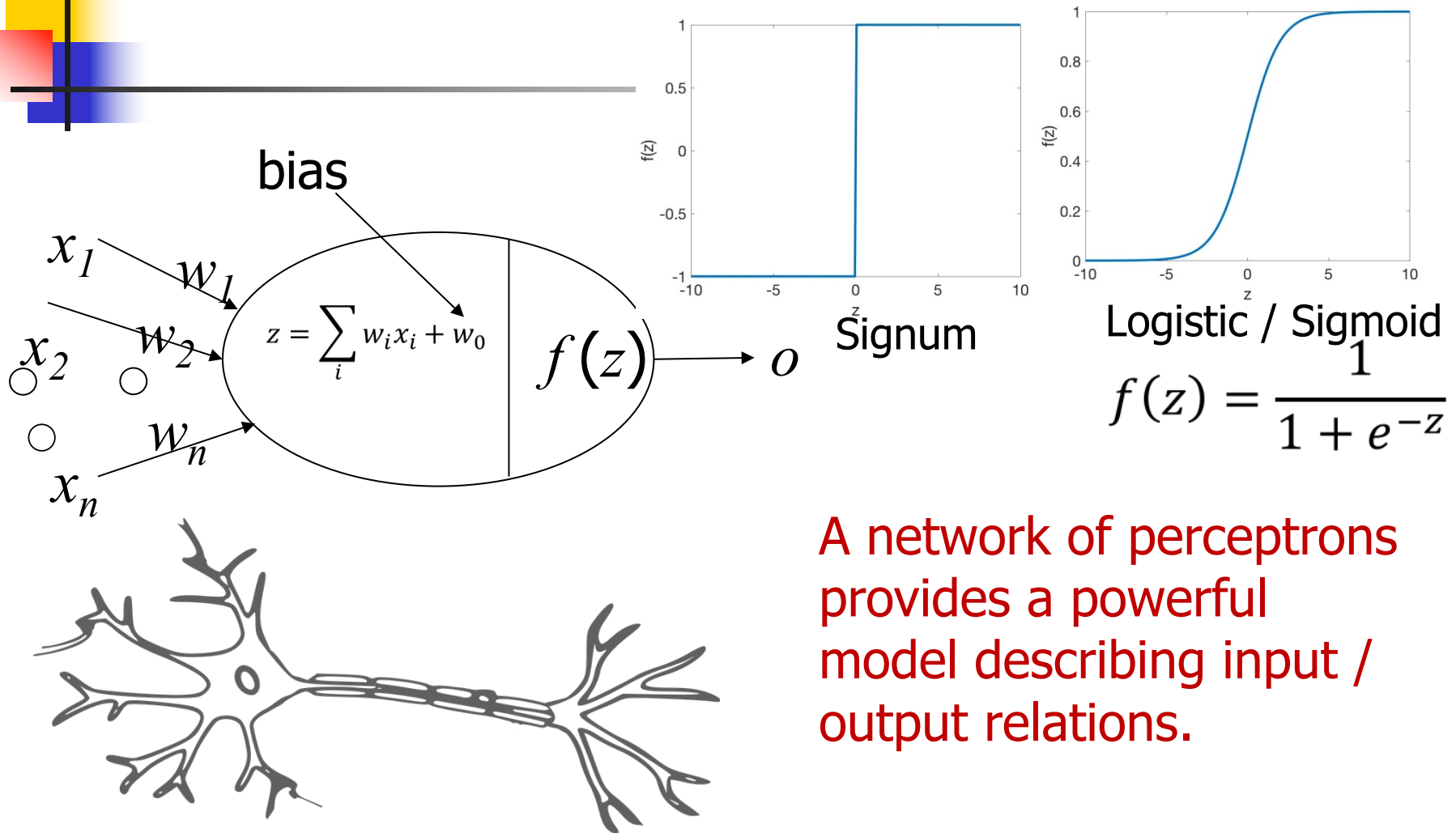
Jayanta Mukhopadhyay
Dept. of Computer Science and Engg.



Books

- Chapter 6 of “Pattern Classification” by R.O. Duda, P. E. Hart and D. G. Stork
- Chapter 11 of “Introduction to Machine Learning” by Ethem Alpaydin.
- Chapter 4 of “Machine learning” by Tom M. Mitchel.

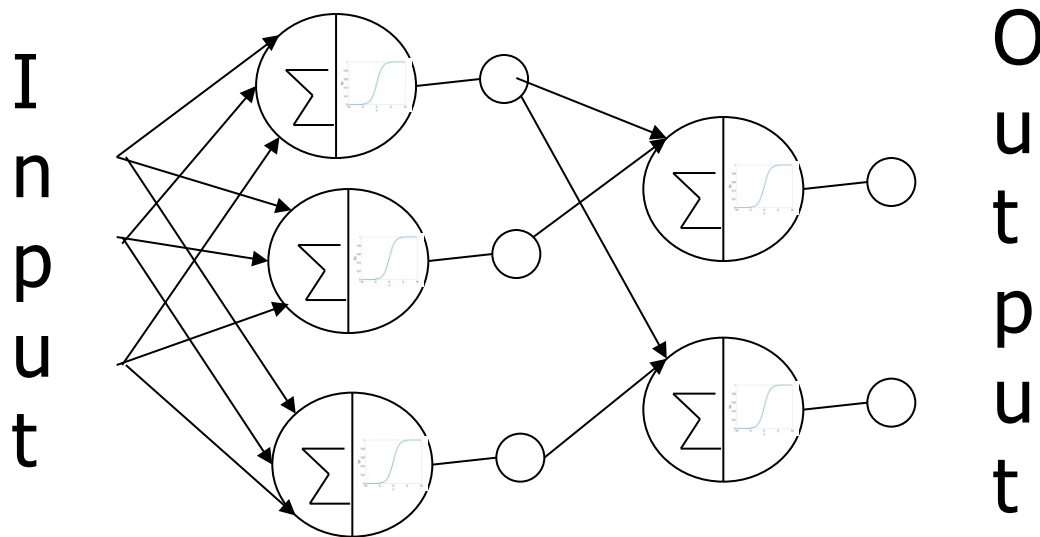
Perceptron modelling a neuron



A network of perceptrons provides a powerful model describing input / output relations.

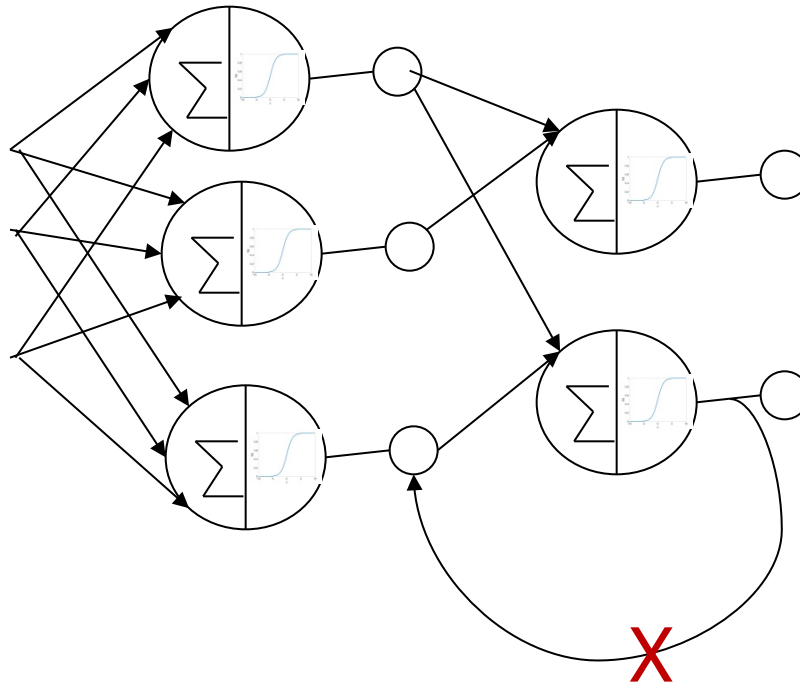
Artificial Neural Network

- A network of perceptrons.
 - Input: A vector
 - Output: A vector / A scalar

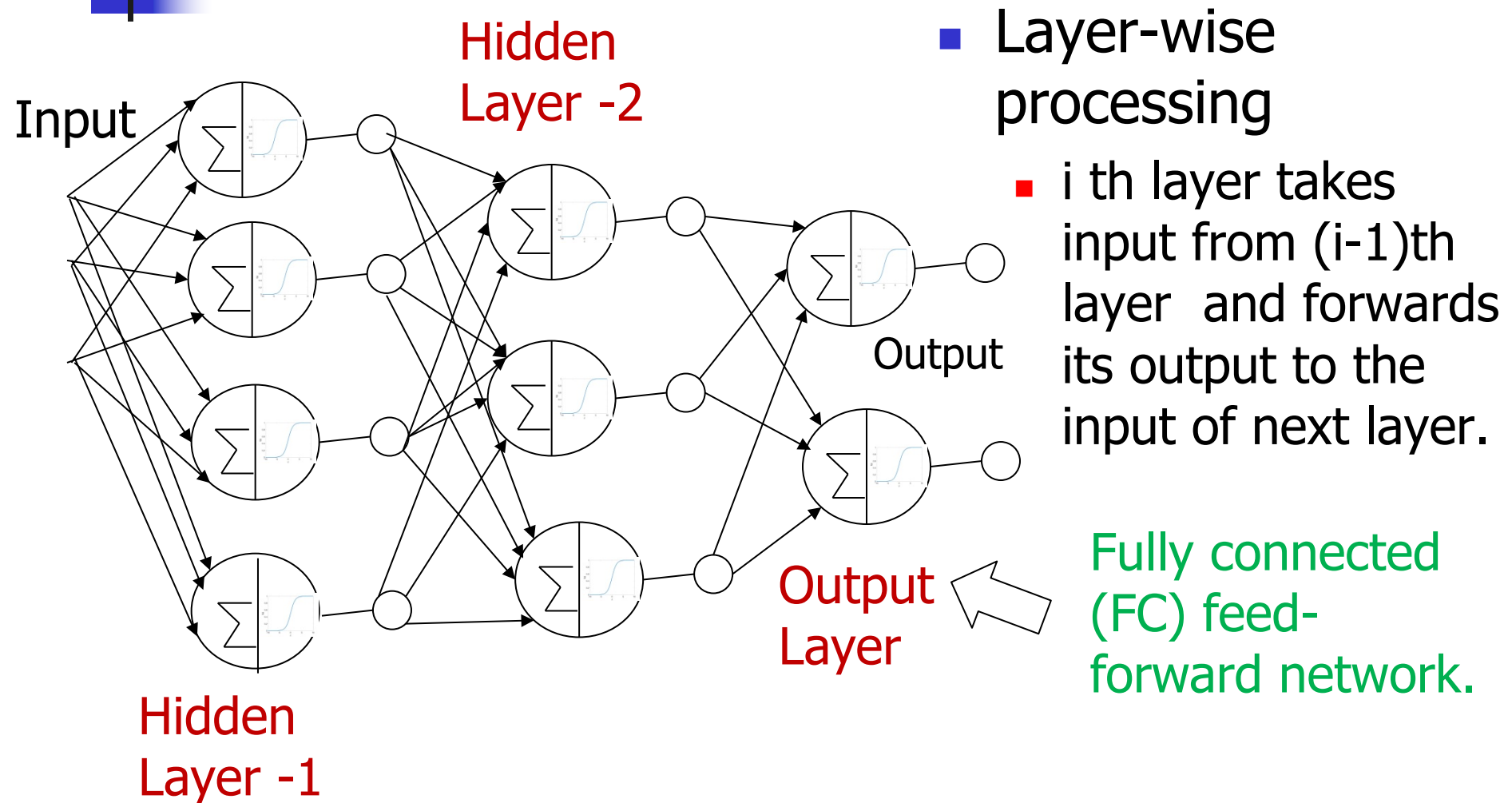


Feed-forward Network

- No feed back or loop in the network.



Multilayered feed-forward Network





Mathematical description of the model

- Let j th neuron of i th layer be $ne_j^{(i)}$.
- Its corresponding weights
 - $W_j^{(i)} = (w_{j1}^{(i)}, w_{j2}^{(i)}, \dots, w_{jn_{(i-1)}}^{(i)})$
 - Bias: $w_{j0}^{(i)}$
 - $n_{(i-1)}$: Dimension of input to the neuron
 - n_i : Dimension of output at i th layer

- Output of the neuron:

$$y_j^{(i)} = f\left(W_j^{(i)T} X^{(i-1)} + w_{j0}^{(i)}\right)$$

Mathematical description of the model

- Output of j th neuron in i th layer:

$$y_j^{(i)} = f\left(W_j^{(i)T} X^{(i-1)} + w_{j0}^{(i)}\right)$$

- Input output relation in i th layer

$$\mathbf{Z}^{(i)} = \begin{bmatrix} W_1^{(i)T} \\ W_2^{(i)T} \\ \vdots \\ W_{n_i}^{(i)T} \end{bmatrix} X^{(i-1)} + \begin{bmatrix} w_{10}^{(i)} \\ w_{20}^{(i)} \\ \vdots \\ w_{n_i0}^{(i)} \end{bmatrix}$$

$\mathbf{W}^{(i)}$ points to the weight matrix $\begin{bmatrix} W_1^{(i)T} \\ W_2^{(i)T} \\ \vdots \\ W_{n_i}^{(i)T} \end{bmatrix}$. $\mathbf{b}^{(i)}$ points to the bias vector $\begin{bmatrix} w_{10}^{(i)} \\ w_{20}^{(i)} \\ \vdots \\ w_{n_i0}^{(i)} \end{bmatrix}$.

Input output relation

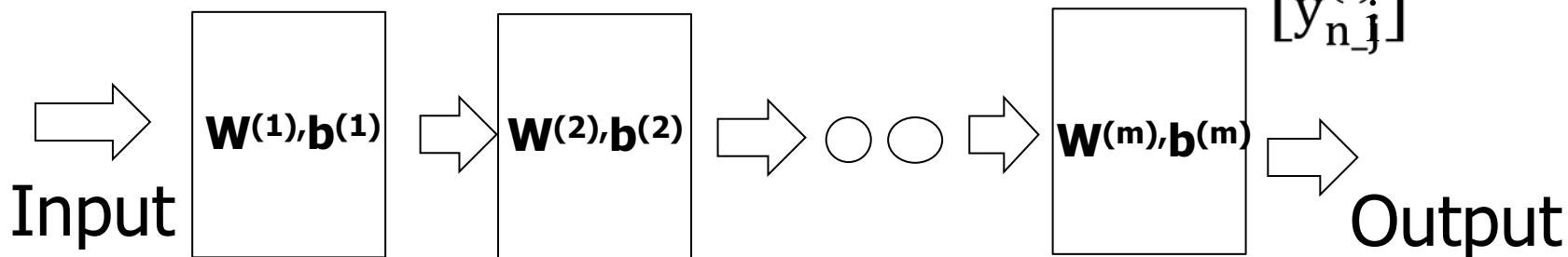
- Output of j th neuron in i th layer:

$$y_j^{(i)} = f\left(W_j^{(i)T} X^{(i-1)} + w_{j0}^{(i)}\right)$$

- Input output relation in i th layer

$$Z^{(i)} = \begin{bmatrix} W_1^{(i)T} \\ W_2^{(i)T} \\ \vdots \\ W_{n_j}^{(i)T} \end{bmatrix} X^{(i-1)} + \begin{bmatrix} w_{10}^{(i)} \\ w_{20}^{(i)} \\ \vdots \\ w_{n_j 0}^{(i)} \end{bmatrix}$$

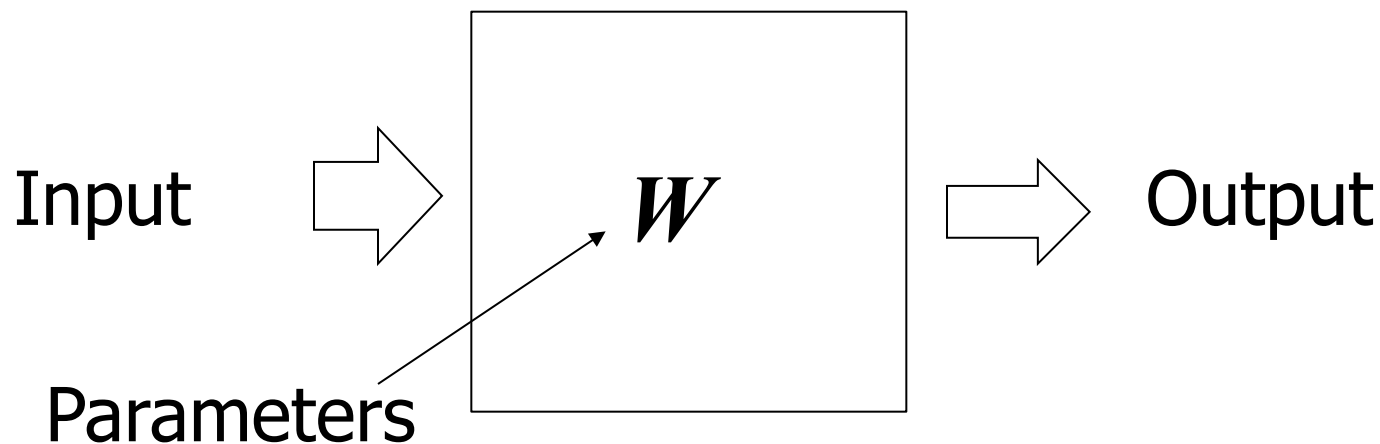
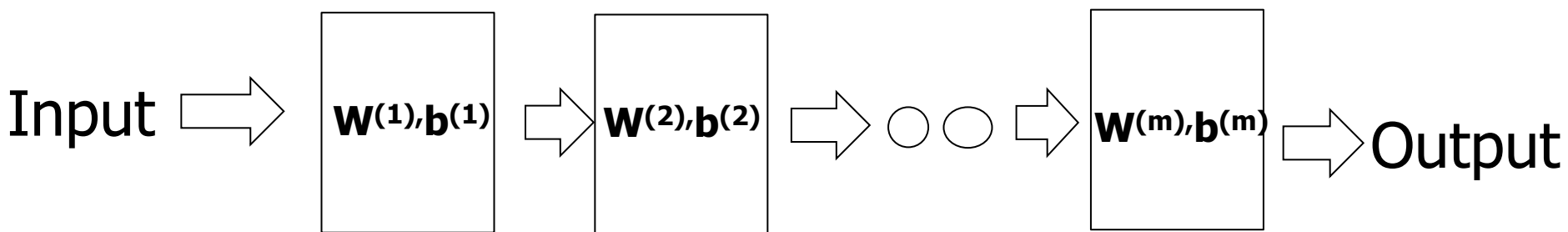
$$Y^{(i)} = f(W^{(i)} X^{(i-1)} + b^{(i)}) \equiv \begin{bmatrix} y_1^{(i)} \\ y_2^{(i)} \\ \vdots \\ y_{n_j}^{(i)} \end{bmatrix}$$





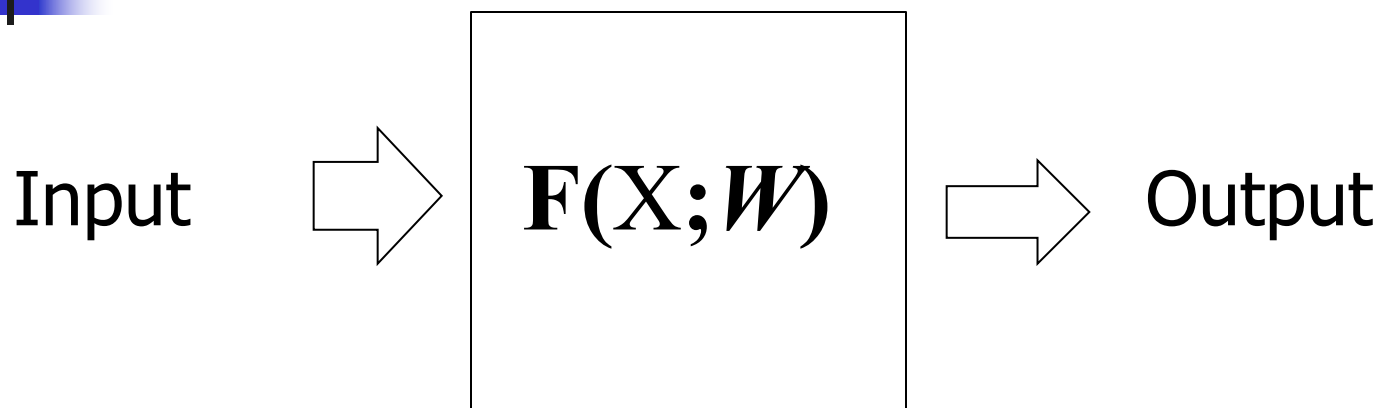
Input output relation

$$Y^{(i-1)} = \mathbf{f}(\mathbf{W}^{(i)}X^{(i-1)} + \mathbf{b}^{(i)})$$





Optimization problem

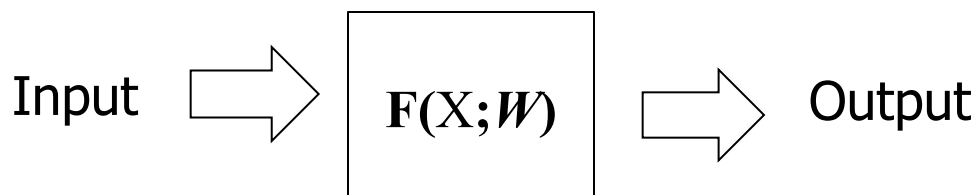


Given $\{(X_i, O_i)\}$, $i=1,2,\dots,N$, find \mathbf{W} such that it produces O_i given input X_i for all i .

Minimize:
$$J_n(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \|O_i - F(X_i; \mathbf{W})\|^2$$

Apply the same gradient descent procedure to obtain the solution.

Optimization problem



Training samples:
 $\{(X_i, O_i)\}, i=1,2,\dots,N$

Minimize:

$$J_n(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \|O_i - F(X_i; \mathbf{W})\|^2$$

Apply the same
gradient descent
procedure to obtain
the solution.

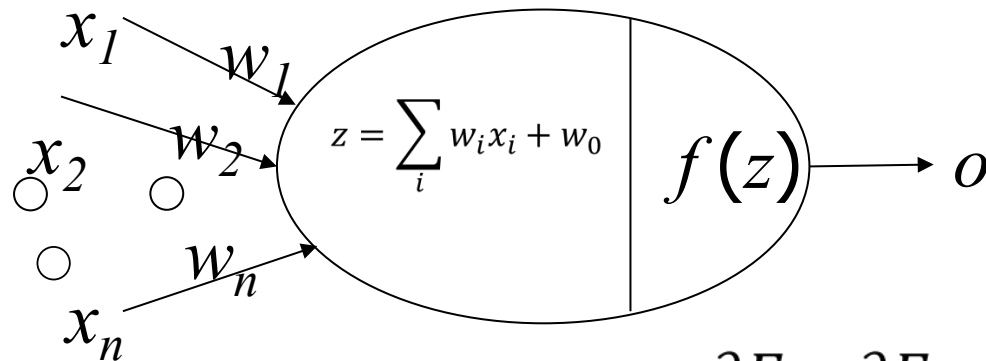
1. Start with an initial \mathbf{W}_0 .
2. Update \mathbf{W} iteratively.

$$\mathbf{W}_i = \mathbf{W}_{i-1} + \eta(i) \sum_k (O_k - F(X_k; \mathbf{W}_{i-1})) \nabla F(X_k; \mathbf{W}_{i-1})$$

Stochastic gradient descent:

$$\mathbf{W}_i = \mathbf{W}_{i-1} + \eta(i)(O_k - F(X_k; \mathbf{W}_{i-1})) \nabla F(X_k; \mathbf{W}_{i-1})$$

Chain rule of computing gradient of a single neuron



Target response: t

Error:

$$E = (t - o)^2$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial E}{\partial o} \frac{\partial o}{\partial z} \frac{\partial z}{\partial w_i}$$

Annotations for the chain rule above:

- $\frac{\partial E}{\partial o} \rightarrow -2(t - o)$
- $\frac{\partial o}{\partial z} \rightarrow f'(z)$
- $\frac{\partial z}{\partial w_i} \rightarrow x_i$

$$\nabla(W) = \left(\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right)$$

$$f(z) = \frac{1}{1 + e^{-z}}$$

$$f'(z) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$\frac{\partial E}{\partial w_i} = -2(t - o)f'(z)x_i$$

Analytical method!
Computed given the functional values.

$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial o} \frac{\partial o}{\partial z} \frac{\partial z}{\partial x_i}$$

$$\frac{\partial E}{\partial x_i} = -2(t - o)f'(z)w_i$$

Annotation for the chain rule above:

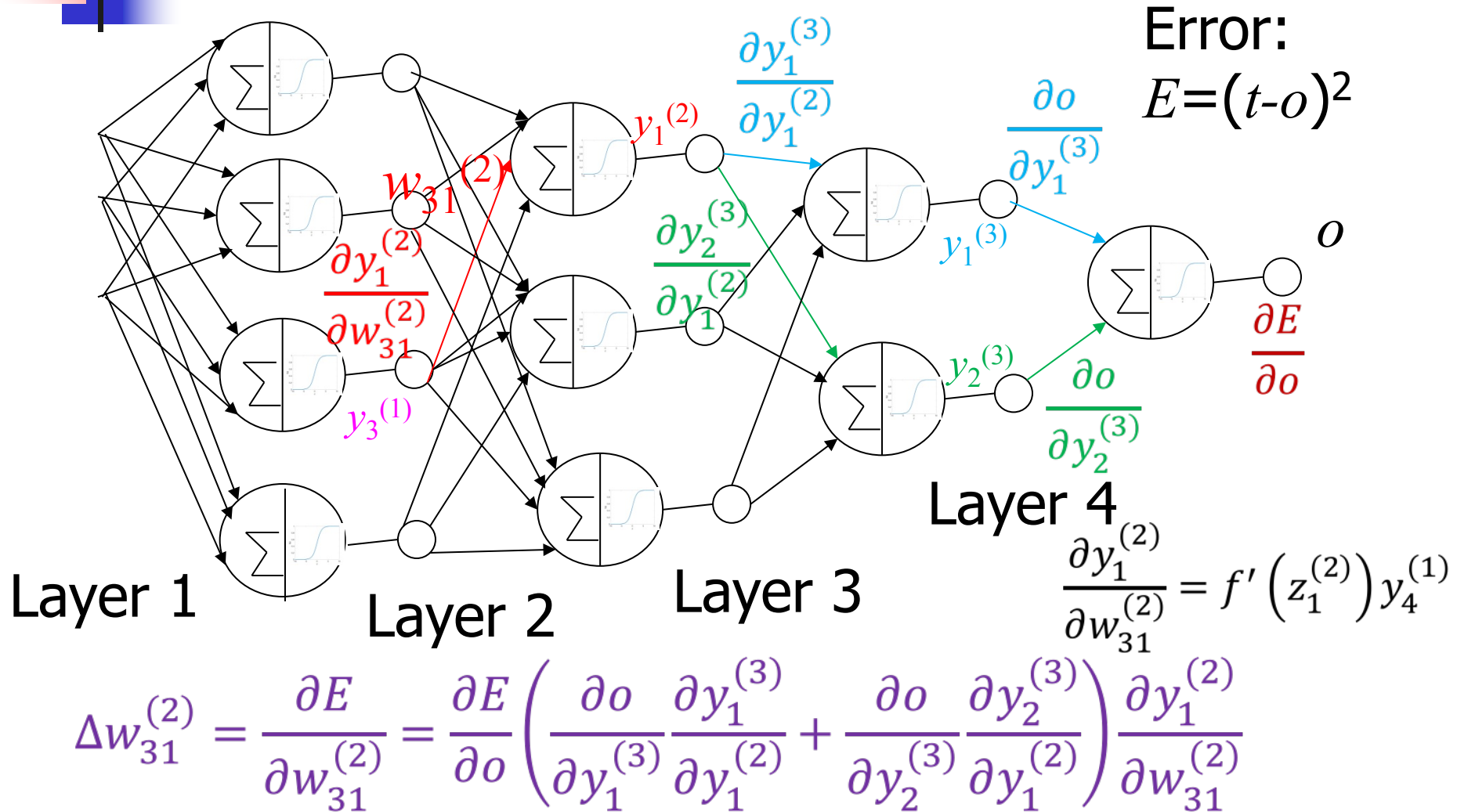
- $\frac{\partial o}{\partial z} \rightarrow f(z)(1 - f(z))$ (via the identity $f'(z) = f(z)(1 - f(z))$)



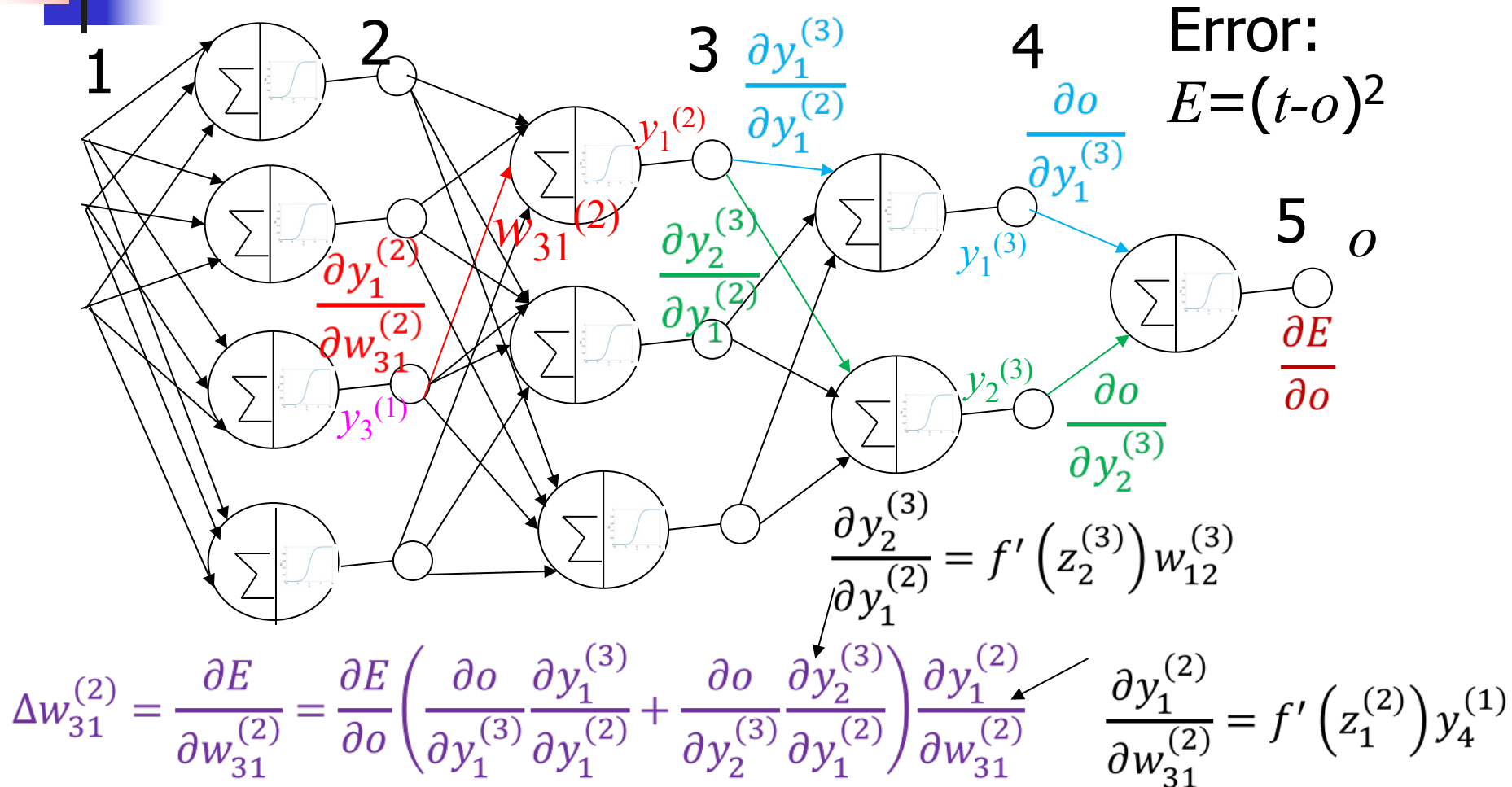
Computing gradient: Back propagation method

- For multi-layered feed forward network.
- Apply chain rule.
 - From output to toward input.
 - From output layer to toward input layer.
 - Compute partial derivatives of weights at $(i-1)$ th layer from the i th layer.

Back propagation: Concept



Back propagation: Concept



Back propagation: Delta rule

Error:

$$E = (t - o)^2$$

$$\frac{\partial y_2^{(3)}}{\partial y_1^{(2)}} = f'(z_2^{(3)}) w_{12}^{(3)}$$

$$\frac{\partial y_1^{(3)}}{\partial y_1^{(2)}} = f'(z_1^{(3)}) w_{11}^{(3)}$$

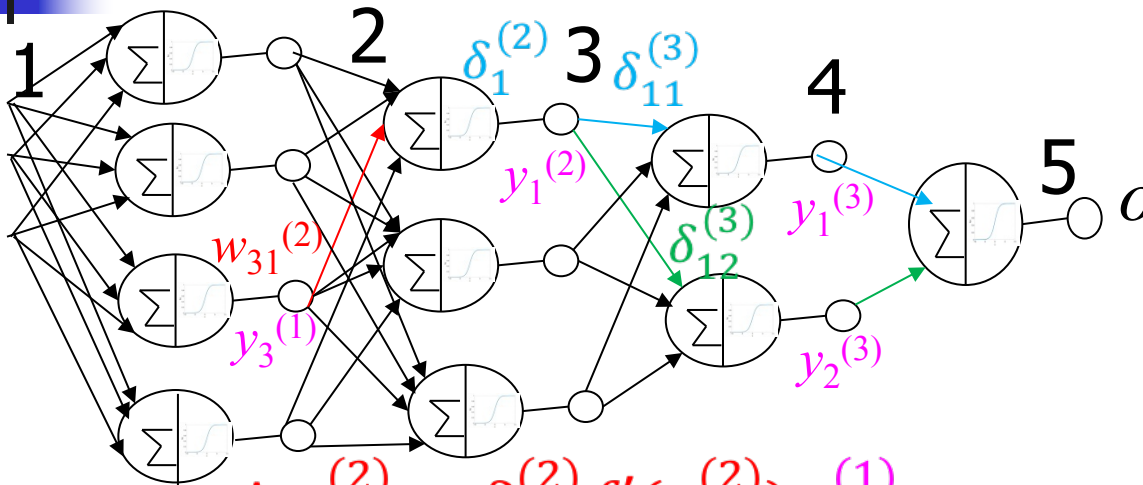
$$\Delta w_{31}^{(2)} = \delta_1^{(2)} f'(z_1^{(2)}) y_3^{(1)}$$

$$\Delta w_{31}^{(2)} = \frac{\partial E}{\partial w_{31}^{(2)}} = \frac{\partial E}{\partial o} \left(\frac{\partial o}{\partial y_1^{(3)}} \frac{\partial y_1^{(3)}}{\partial y_1^{(2)}} + \frac{\partial o}{\partial y_2^{(3)}} \frac{\partial y_2^{(3)}}{\partial y_1^{(2)}} \right) \frac{\partial y_1^{(2)}}{\partial w_{31}^{(2)}}$$

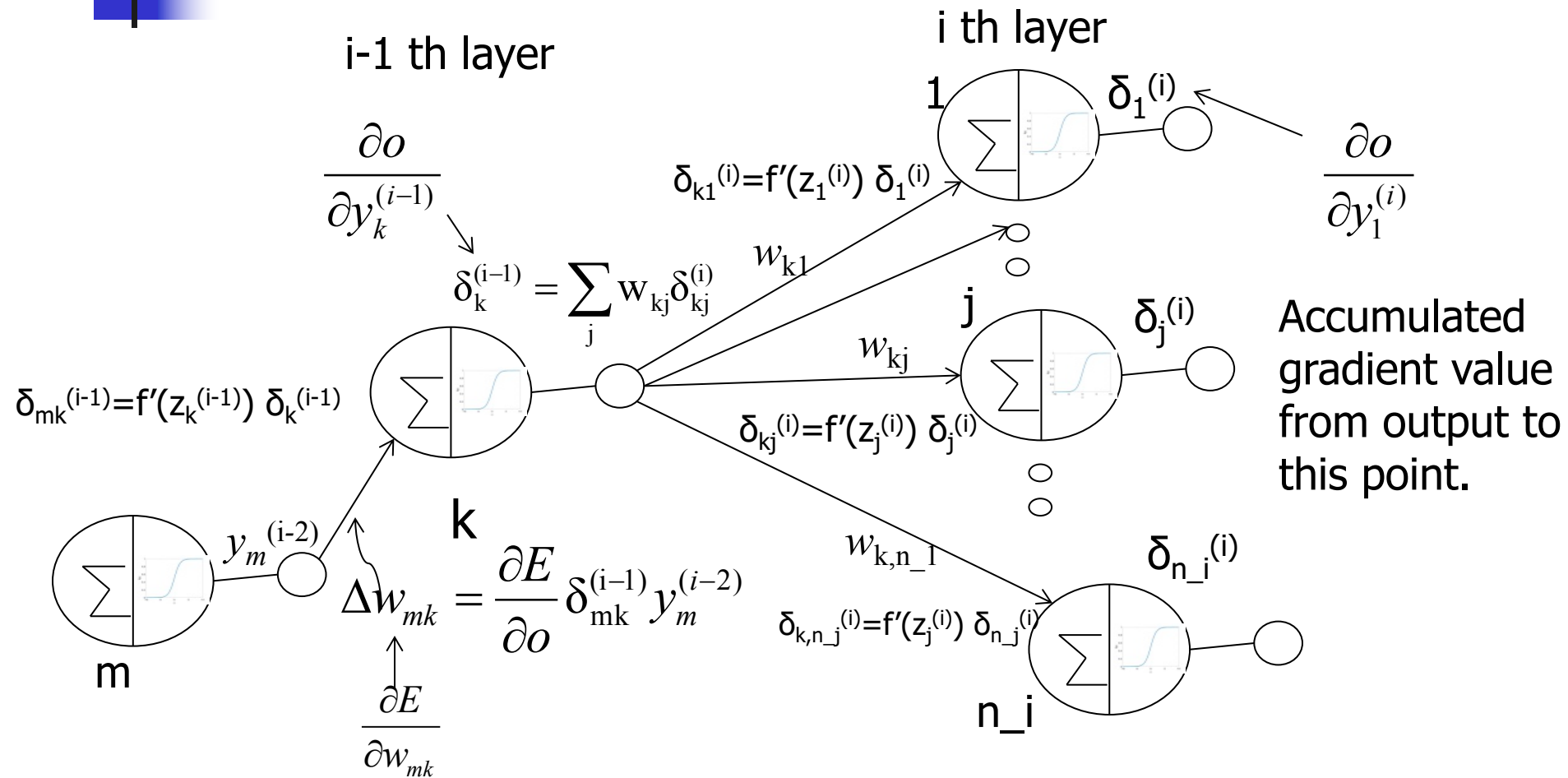
$$\Delta w_{31}^{(2)} = \frac{\partial E}{\partial w_{31}^{(2)}} = \frac{\partial E}{\partial o} \left(\underbrace{\frac{\partial o}{\partial y_1^{(3)}} f'(z_1^{(3)}) w_{11}^{(3)}}_{\delta_{11}^{(3)}} + \underbrace{\frac{\partial o}{\partial y_2^{(3)}} f'(z_2^{(3)}) w_{12}^{(3)}}_{\delta_{12}^{(3)}} \right) \frac{\partial y_1^{(2)}}{\partial w_{31}^{(2)}}$$

Delta rule:

$$\delta_1^{(2)} = \delta_{11}^{(3)} w_{11}^{(3)} + \delta_{12}^{(3)} w_{12}^{(3)}$$



Back propagation: Delta rule





ANN training

- Initialize $W^{(0)}$.
- For each training sample (x_i, o_i) do
 - Compute functional values of each neuron in the forward pass.
 - Update weights of each link starting from the output layer using back propagation.
 - Continue till it converges.



Improving convergence

- Momentum

- Gradients may change abruptly in consecutive iteration. To avoid we may use running average of weight updates to be added with the gradient.

$$\Delta w_i^{(t)} = \alpha \Delta w_i^{(t-1)} - \eta \frac{\partial E^{(t)}}{\partial w_i}$$

Usually ranges
between 0.5 to 1.0

- Adaptive learning rate:

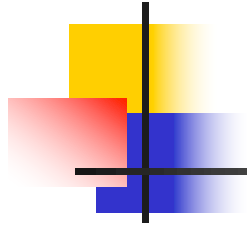
$$\Delta \eta = \begin{cases} +a & \text{if } E^{t+T} < E^t \\ -b\eta & \text{Otherwise} \end{cases}$$

- We increase learning rate at constant steps if error decreases, else decrease it geometrically.



Summary

- A perceptron models a neuron.
- A network of perceptron nodes provide a powerful model for classifying linearly non-separable classes as well.
- A multilayer feed forward network of neurons can be trained using back propagation algorithm (to compute gradients of error w.r.t weights).
- Can be used for regression as well.



Thank you!