



# Deep Learning

## Fundamentals

### Machine Learning

*Dongsuk Yook*  
*Artificial Intelligence Laboratory*  
*Korea University*

# Contents

- ☐ Perceptrons
- ☐ Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications

# Class Objectives

- ☐ Understanding the fundamentals of deep learning
- ☐ Being able to build a basic deep neural network for a given application

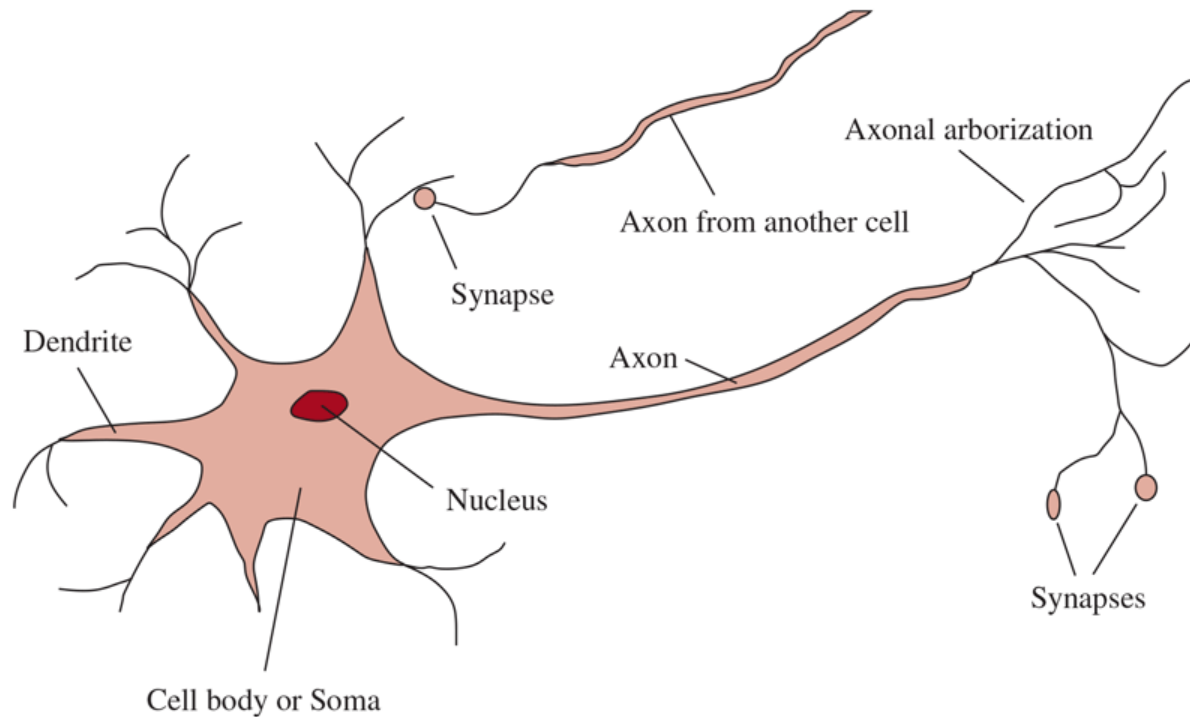
# Contents

- ❑ Perceptrons
  - Neurons
  - Learning for Perceptrons
- ❑ Multilayer Perceptrons
- ❑ Error Backpropagation Algorithm
- ❑ Applications

# Neurons

## □ Biological neural networks

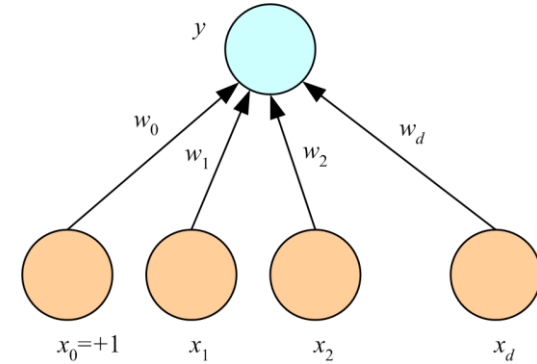
- Approximately  $10^{11}$  interconnected *neurons*, each connected to 10–100,000 others through *synapses*
- More than  $10^{14}$  synaptic connections
- Highly complex and massively parallel



# Perceptrons

## □ Perceptron [McCulloch and Pitts 1943, Rosenblatt 1957]

- Input  $x_j$
- Output  $y$
- Connection *weight* (or synaptic weight)  $w_j$



## □ Linear regression

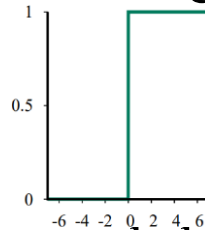
- $y = \sum_j w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$

;  $\mathbf{x} \equiv [1, x_1, \dots, x_d]$   
 $w_0$ : bias

## □ Binary classification

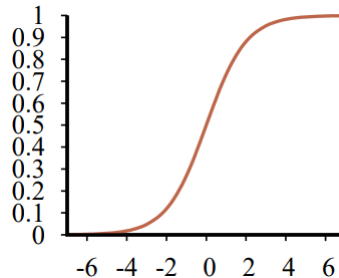
- Linear discriminant function using a *hard threshold* function

- $y = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$



## □ Sigmoid function for posterior probability

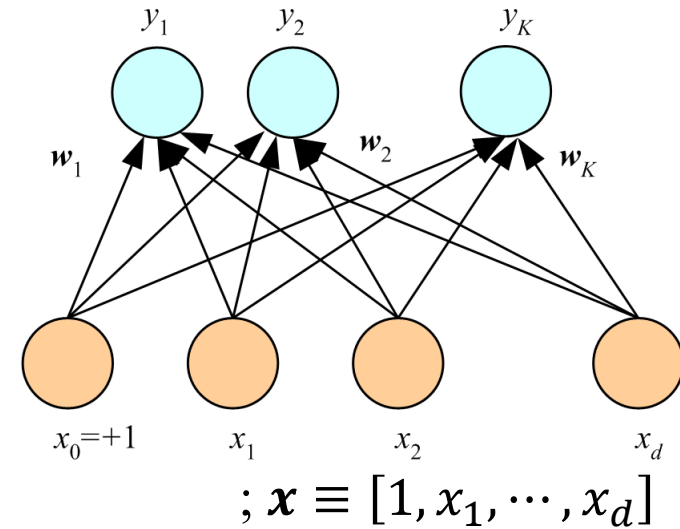
- $y = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$



;  $\exp(-\mathbf{w}^T \mathbf{x}) \equiv e^{-\mathbf{w}^T \mathbf{x}}$

# Multiple Perceptrons

- Multiple perceptrons
  - Input  $x_j$
  - Output  $y_i$
  - Weight  $w_{ij}$  (between input  $x_j$  and output  $y_i$ )



- Multivariate linear regression
  - $y_i = \sum_j w_{ij}x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}$   
 $\mathbf{y} = \mathbf{W}\mathbf{x}$

- Multiclass classification
  - $\mathbf{y} = \mathbf{W}\mathbf{x}$   
 $\mathbf{x} \in c_k$  if  $k = \arg \max_i y_i$

- *Softmax* function for posterior probability
  - $s_i = \sum_j w_{ij}x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}$   
 $y_i = \frac{\exp(s_i)}{\sum_i \exp(s_i)}$

# Contents

- ❑ Perceptrons
  - Neurons
  - Learning for Perceptrons
- ❑ Multilayer Perceptrons
- ❑ Error Backpropagation Algorithm
- ❑ Applications



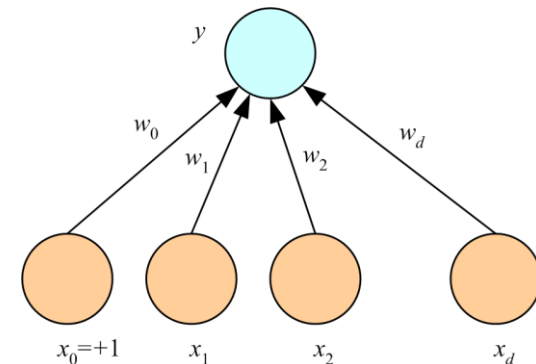
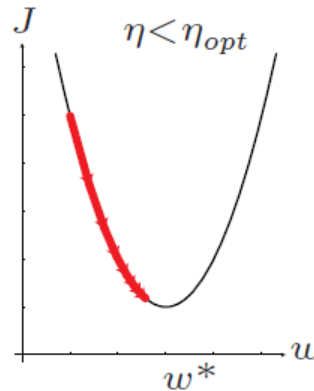
# Learning for a Linear Output Node

□ Stochastic gradient descent (SGD) for a linear output node

- $E(\mathbf{w}|\mathbf{x}, r) \equiv \frac{1}{2}(r - y)^2$  ; mean squared error (MSE)

- $$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial w_j}$$
$$= \underbrace{-(r - y)}_{\frac{\partial E}{\partial y}} \underbrace{x_j}_{\frac{\partial y}{\partial w_j}}$$

- $w_j \leftarrow w_j + \eta(r - y)x_j$  ;  $\eta$ : learning rate or step size



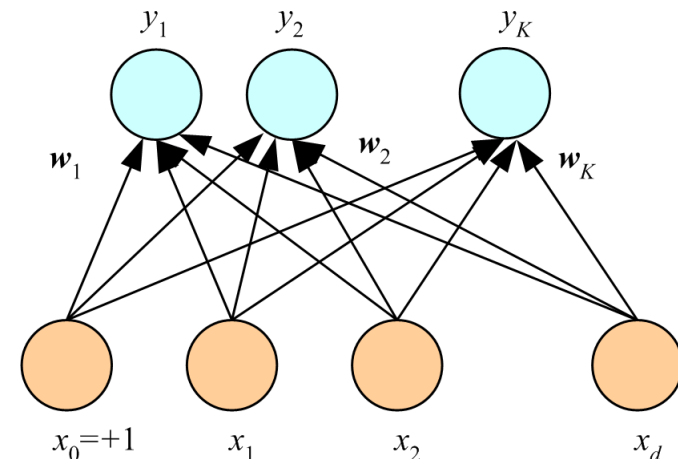
# Learning for Multiple Linear Output Nodes

□ Stochastic gradient descent for multiple linear output nodes

- $E(\mathbf{W}|\mathbf{x}, \mathbf{r}) \equiv \sum_k \underbrace{\frac{1}{2}(r_k - y_k)^2}_{E_k} \quad ; y_i = \mathbf{w}_i^T \mathbf{x} = \sum_j w_{ij} x_j + w_{i0}$

- $$\begin{aligned} \frac{\partial E}{\partial w_{ij}} &= \sum_k \frac{\partial E_k}{\partial y_i} \frac{\partial y_i}{\partial w_{ij}} \\ &= \underbrace{-(r_i - y_i)}_{\frac{\partial E_i}{\partial y_i}} \underbrace{x_j}_{\frac{\partial y_i}{\partial w_{ij}}} \end{aligned}$$

- $w_{ij} \leftarrow w_{ij} + \eta(r_i - y_i)x_j$

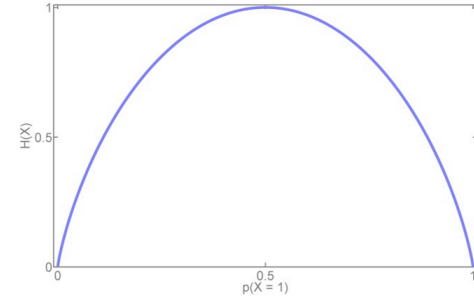


# Entropy

## □ Entropy

- $\mathbb{H}(p) \equiv -\sum_x p(x) \log_2 p(x)$
- e.g.,  $-\left[\theta \log_2 \theta + (1 - \theta) \log_2 (1 - \theta)\right]$

;  $\mathbb{H}(x)$



## □ Cross entropy

- $\mathbb{H}(p, q) \equiv -\sum_x p(x) \log_2 q(x)$

; minimum when  $p = q$   
( $\because$  Jensen's inequality)

## □ Kullback-Leibler (KL) divergence (also called *relative entropy*)

- $\mathbb{KL}(p||q) \equiv \mathbb{H}\left(p(x), \frac{q(x)}{p(x)}\right) = -\sum_x p(x) \log_2 \frac{q(x)}{p(x)}$   
 $= -\sum_x p(x) \log_2 q(x) + \sum_x p(x) \log_2 p(x) = \mathbb{H}(p, q) - \mathbb{H}(p)$

## □ Conditional entropy

- $\mathbb{H}(x|y) \equiv \mathbb{H}(p(x, y), p(x|y)) = -\sum_{x,y} \overbrace{p(x, y)}^{p(y)p(x|y)} \log_2 p(x|y)$   
 $= -\sum_y \sum_x p(y)p(x|y) \log_2 p(x|y) = -\sum_y p(y) \mathbb{H}(p(x|y))$

## □ Mutual information (MI, also called *information gain*)

- $\mathbb{I}(x, y) \equiv \mathbb{KL}(p(x, y)||p(x)p(y)) = \mathbb{H}(x) - \mathbb{H}(x|y) = \mathbb{H}(y) - \mathbb{H}(y|x)$

# Learning for a Sigmoid Output Node

□ Stochastic gradient descent for a sigmoid output node

- $E(\mathbf{w}|\mathbf{x}, r) \equiv -r \log y - (1 - r) \log(1 - y)$

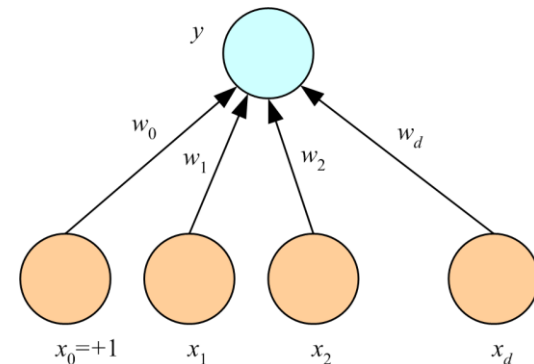
- $$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial E}{\partial y} \frac{\partial y}{\partial s} \frac{\partial s}{\partial w_i} \\ &= \underbrace{\left(-\frac{r}{y} + \frac{1-r}{1-y}\right)}_{\frac{\partial E}{\partial y}} \underbrace{y(1-y)}_{\frac{\partial y}{\partial s}} \underbrace{x_i}_{\frac{\partial s}{\partial w_i}} \\ &= -(r - y)x_i \end{aligned}$$

- $w_i \leftarrow w_i + \eta(r - y)x_i$

; *cross entropy*

;  $y = \frac{1}{1+\exp(-s)}, s = \mathbf{w}^T \mathbf{x}$

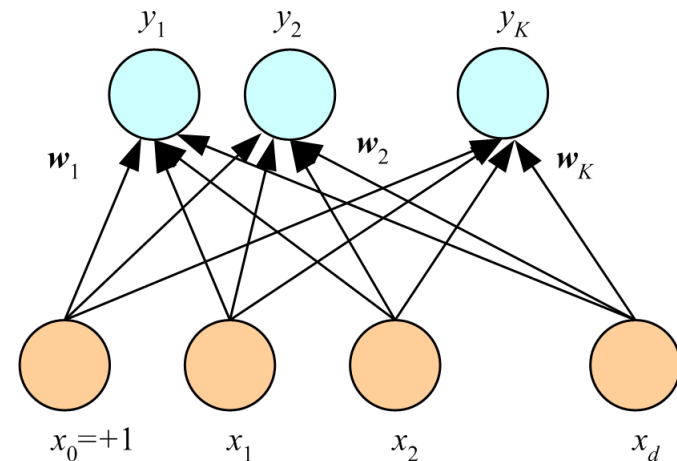
$$\begin{aligned} y' &= -\frac{1}{(1+\exp(-s))^2} \exp(-s) (-1) \\ &= \frac{1}{1+\exp(-s)} \frac{\exp(-s)}{1+\exp(-s)} \\ &= \frac{1}{1+\exp(-s)} \left( \frac{1+\exp(-s)}{1+\exp(-s)} - \frac{1}{1+\exp(-s)} \right) \\ &= y(1 - y) \end{aligned}$$



# Learning for Softmax Output Nodes

□ Stochastic gradient descent for a softmax output node

- $E(\mathbf{W}|\mathbf{x}, \mathbf{r}) \equiv -\sum_k \underbrace{r_k \log y_k}_{E_k}$  ;  $y_k = \frac{\exp(s_k)}{\sum_{\hat{k}} \exp(s_{\hat{k}})}$ ,  $s_i = \mathbf{w}_i^\top \mathbf{x}$
- $$\begin{aligned} \frac{\partial E}{\partial w_{ij}} &= -\sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial s_i} \frac{\partial s_i}{\partial w_{ij}} \\ &= -\left( \sum_{k \neq i} \underbrace{\frac{r_k}{y_k}}_{\frac{\partial E_k}{\partial y_k}} \underbrace{\left( -\frac{\exp(s_k) \exp(s_i)}{(\sum_{\hat{k}} \exp(s_{\hat{k}}))^2} \right)}_{\frac{\partial y_k}{\partial s_i}} + \underbrace{\frac{r_i}{y_i}}_{\frac{\partial E_i}{\partial y_i}} \underbrace{\left( \frac{\exp(s_i)}{\sum_{\hat{i}} \exp(s_{\hat{i}})} - \frac{\exp(s_i) \exp(s_i)}{(\sum_{\hat{i}} \exp(s_{\hat{i}}))^2} \right)}_{\frac{\partial y_i}{\partial s_i}} \right) \underbrace{x_j}_{\frac{\partial s_i}{\partial w_{ij}}} \\ &= -\left( \sum_{k \neq i} \frac{r_k}{y_k} (-y_k y_i) + \frac{r_i}{y_i} (y_i - y_i^2) \right) x_j \\ &= -(\sum_{k \neq i} r_k (-y_i) + r_i (1 - y_i)) x_j \\ &= -(\sum_k r_k (-y_i) + r_i) x_j \\ &= -(r_i - y_i) x_j \end{aligned}$$
- $w_{ij} \leftarrow w_{ij} + \eta(r_i - y_i)x_j$



# Stochastic Gradient Descent for Perceptrons

- Input:  $\mathcal{D} = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$
1. Initialize  $\mathbf{W}$  with small random numbers.
  2. **repeat**
  3.   **for** each  $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{D}$  in random order
  4.     **for** each output node index  $i$
  5.        $s_i \leftarrow \mathbf{w}_i^\top \mathbf{x}^t$
  6.     **end**
  7.     **for** each output node index  $i$
  8.        $y_i \leftarrow \exp(s_i) / \sum_i \exp(s_i)$
  9.     **end**
  10.    **for** each output node index  $i$
  11.     **for** each input node index  $j$
  12.        $w_{ij} \leftarrow w_{ij} + \eta(r_i^t - y_i)x_j^t$
  13.     **end**
  14.    **end**
  15. **end**
  16. **until** convergence

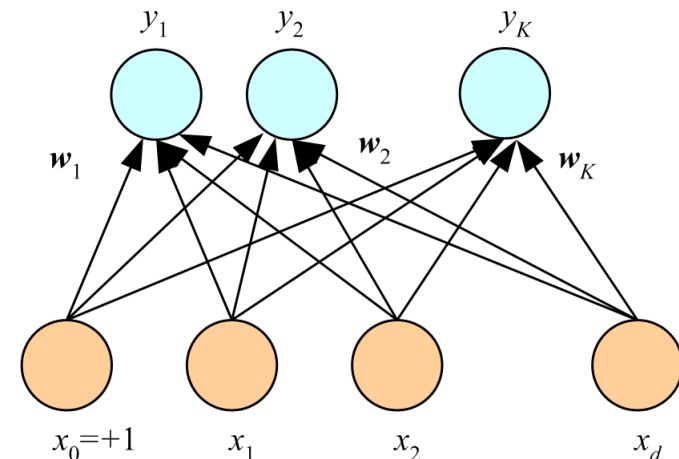
; epoch

; regression:  $y_i \leftarrow s_i$

sigmoid:  $y \leftarrow 1 / (1 + \exp(-\mathbf{w}^\top \mathbf{x}))$

Hard threshold:  $y \leftarrow \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$

- Output:  $\mathbf{W}$



# Contents

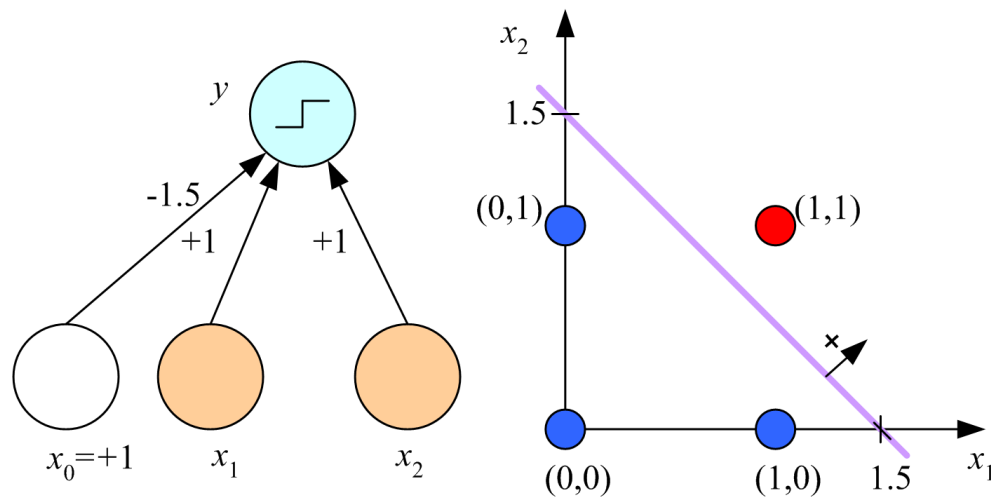
- ☐ Perceptrons
- ☐ Multilayer Perceptrons
  - Boolean Functions
  - Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications

# Boolean AND

## □ Boolean AND operation

- $y = x_1 \text{ AND } x_2$
- $y = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$
- $\mathbf{w}^T \mathbf{x} = w_0 + w_1 x_1 + w_2 x_2 > 0$   
 $\Rightarrow x_2 > -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2}$

$x_1$	$x_2$	$r$
0	0	0
0	1	0
1	0	0
1	1	1



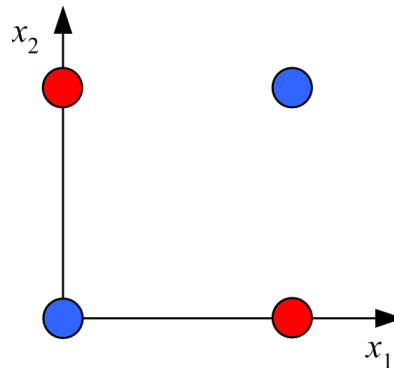


# Boolean XOR

## □ Boolean XOR operation

- $y = x_1 \text{ XOR } x_2$
- $y = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$
- $\mathbf{w}^\top \mathbf{x} = w_0 + w_1 x_1 + w_2 x_2 > 0$   
 $\Rightarrow x_2 > -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2}$

$x_1$	$x_2$	$r$
0	0	0
0	1	1
1	0	1
1	1	0

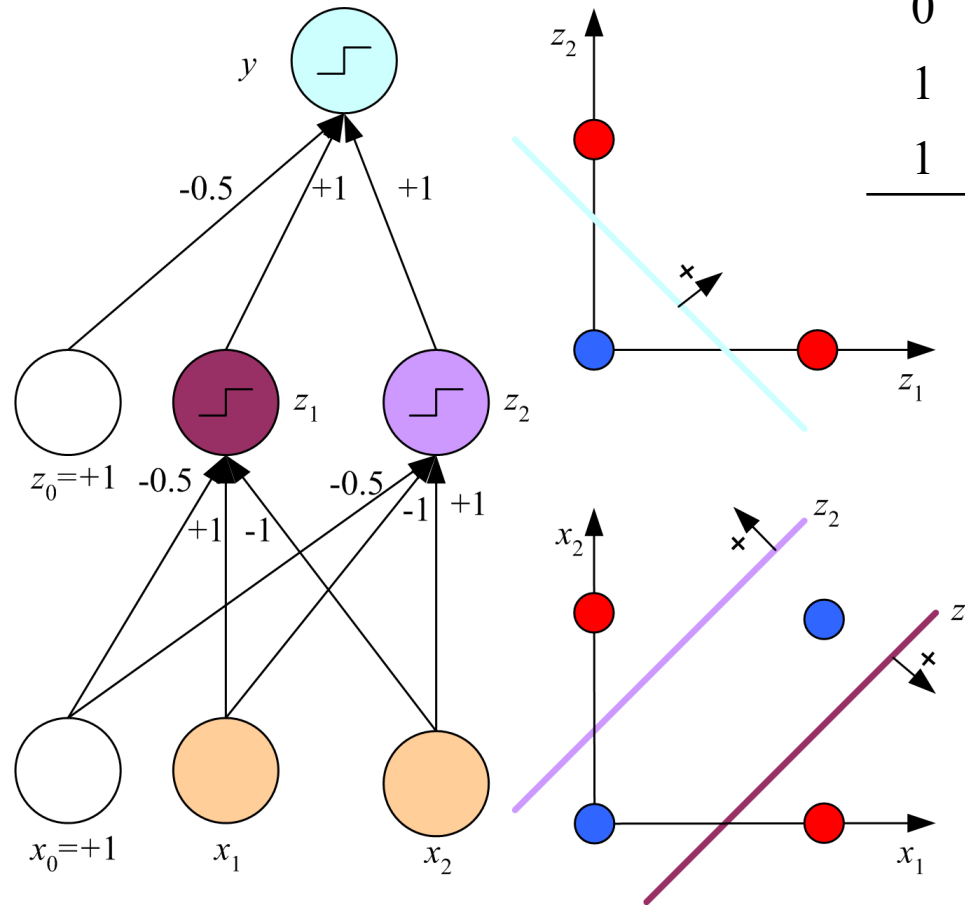


# Boolean XOR

□ Boolean XOR function as a disjunction of conjunctions

▪  $x_1 \text{ XOR } x_2 = (x_1 \text{ AND } \neg x_2) \text{ OR } (\neg x_1 \text{ AND } x_2)$

$x_1$	$x_2$	$z_1$	$z_2$	$r$
0	0	0	0	0
0	1	0	1	1
1	0	1	0	1
1	1	0	0	0



Two-layer perceptrons

# Contents

- ☐ Perceptrons
- ☐ Multilayer Perceptrons
  - Boolean Functions
  - Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications

# Multilayer Perceptrons

## □ Multilayer perceptrons (MLP)

- Input  $x_j$
- Hidden node  $z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$  ; nonlinear activation function
- Weight  $w_{z,ij}$  (between input  $x_j$  and hidden node  $z_i$ )
- Output  $y_i = \begin{cases} \mathbf{w}_{o,i}^\top \mathbf{z} \\ \frac{\exp(s_i)}{\sum_i \exp(s_i)} \end{cases}$  ;  $z_0 \equiv 1$   
;  $s_i = \mathbf{w}_{o,i}^\top \mathbf{z}$
- Weight  $w_{o,ij}$  (between hidden node  $z_j$  and output node  $y_i$ )

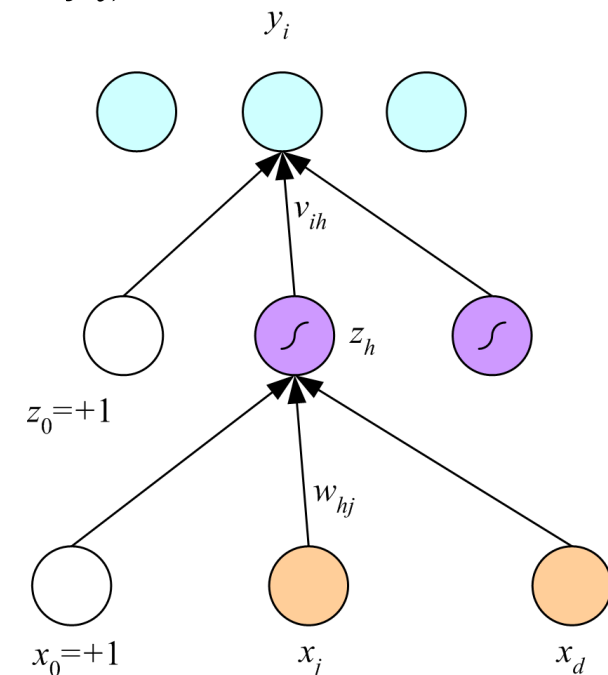
## □ Importance of non-linearity in hidden layers

- $\mathbf{z} = \mathbf{W}\mathbf{x}$   
 $\mathbf{y} = \mathbf{V}\mathbf{z}$   
 $\mathbf{y} = \mathbf{V}\mathbf{W}\mathbf{x} = \mathbf{U}\mathbf{x}$

Output layer

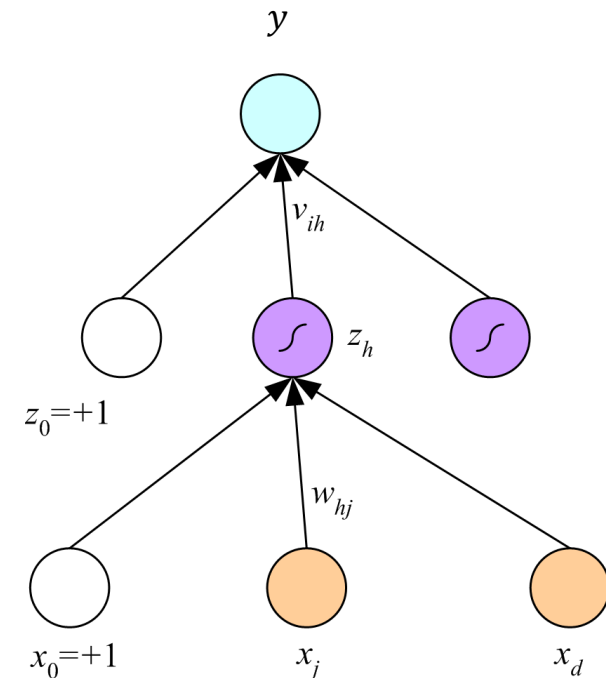
Hidden layer

Input layer


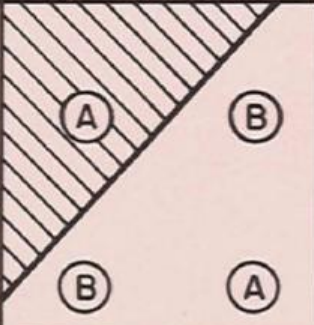
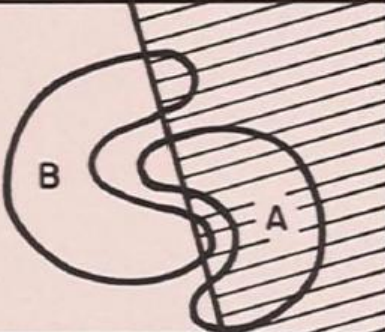

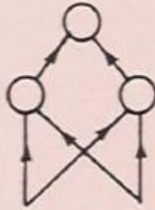
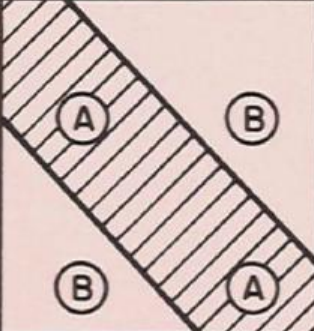
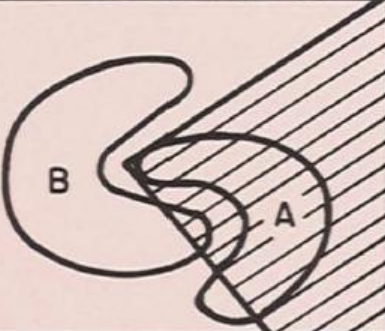
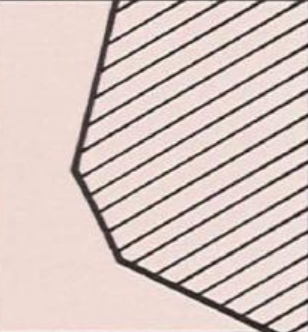
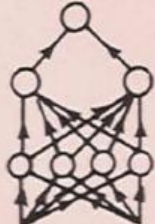
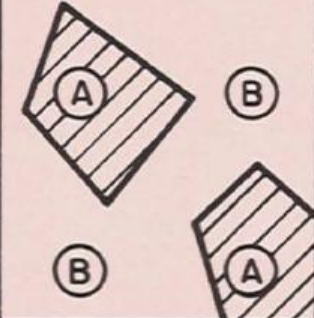
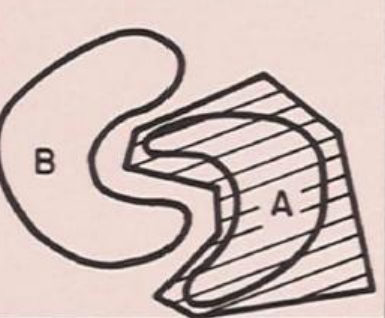



# Universal Approximation Theorem

- Any arbitrary function with continuous input and output can be approximated arbitrarily closely by a two-layer network given sufficient number of hidden units, proper nonlinearities, and weights.
- Proof
  - Fourier theorem



# Decision Regions of MLPs

Structure	Types of Decision Regions	Exclusive-or Problem	Classes with Meshed Regions	Region Shapes
One Layer 	Half-Plane			
Two Layers 	Typically Convex			
Three Layers 	Arbitrary			

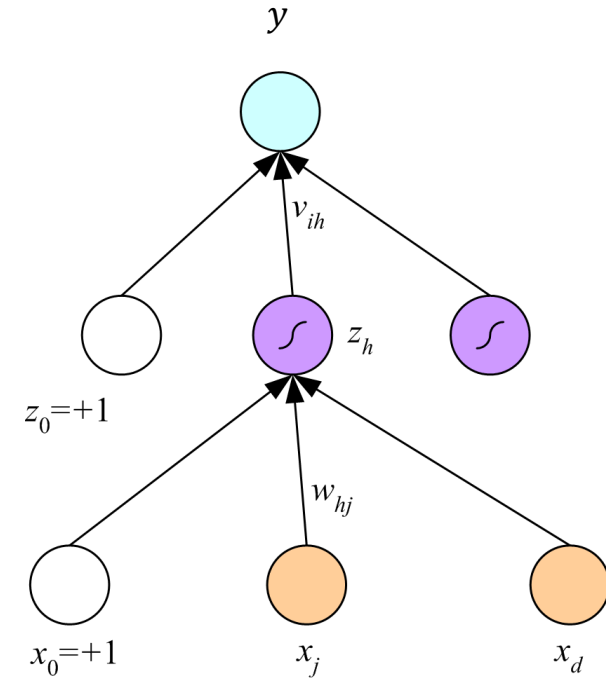
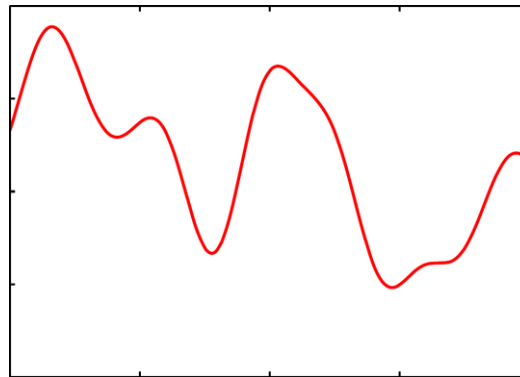
# Contents

- ❑ Perceptrons
- ❑ Multilayer Perceptrons
- ❑ Error Backpropagation Algorithm
  - Regression
  - Classification
  - Deep Neural Networks
- ❑ Applications

# Error Backpropagation

## □ Gradient descent

- $E(\theta|\mathbf{x}, \mathbf{r}) \equiv \frac{1}{2} \sum_i (r_i - y_i)^2$
- $E(\theta|\mathbf{x}, \mathbf{r}) \equiv - \sum_i r_i \log y_i$
- $\theta \leftarrow \theta - \eta \nabla E$





# Learning for Multiple Linear Output Nodes

□ Weight update for multiple linear output nodes

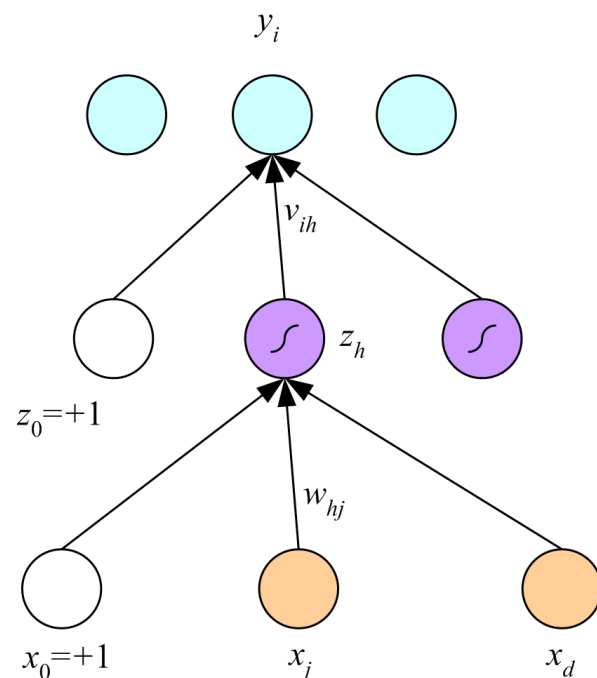
$$\blacksquare E(\mathbf{W}_z, \mathbf{W}_o | \mathbf{x}, \mathbf{r}) \equiv \sum_k \underbrace{\frac{1}{2}(r_k - y_k)^2}_{E_k}$$

$$; y_k = \mathbf{w}_{o,k}^\top \mathbf{z}, z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$$

$$\blacksquare \frac{\partial E}{\partial w_{o,ij}} = \sum_k \frac{\partial E_k}{\partial y_i} \frac{\partial y_i}{\partial w_{o,ij}}$$

$$= \underbrace{-(r_i - y_i)}_{\frac{\partial E_i}{\partial y_i}} \underbrace{z_j}_{\frac{\partial y_i}{\partial w_{o,ij}}}$$

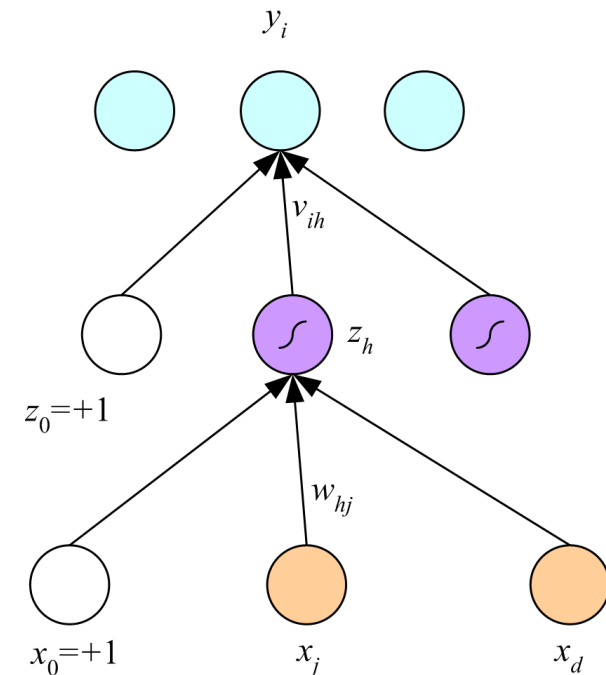
$$\blacksquare w_{o,ij} \leftarrow w_{o,ij} + \eta(r_i - y_i)z_j$$



# Learning for Hidden Nodes

□ Weight update for a hidden node connected to multiple linear output nodes

- $E(\mathbf{W}_z, \mathbf{W}_o | \mathbf{x}, \mathbf{r}) \equiv \sum_k \underbrace{\frac{1}{2}(r_k - y_k)^2}_{E_k} \quad ; y_k = \mathbf{w}_{o,k}^\top \mathbf{z}, z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$
- $$\begin{aligned} \frac{\partial E}{\partial w_{z,ij}} &= \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_i} \frac{\partial z_i}{\partial w_{z,ij}} \\ &= \sum_k \underbrace{-(r_k - y_k)}_{\frac{\partial E_k}{\partial y_k}} \underbrace{w_{o,ki}}_{\frac{\partial y_k}{\partial z_i}} \underbrace{z_i(1 - z_i)x_j}_{\frac{\partial z_i}{\partial w_{z,ij}}} \end{aligned}$$
- $w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k - y_k) w_{o,ki} z_i (1 - z_i) x_j$



# Contents

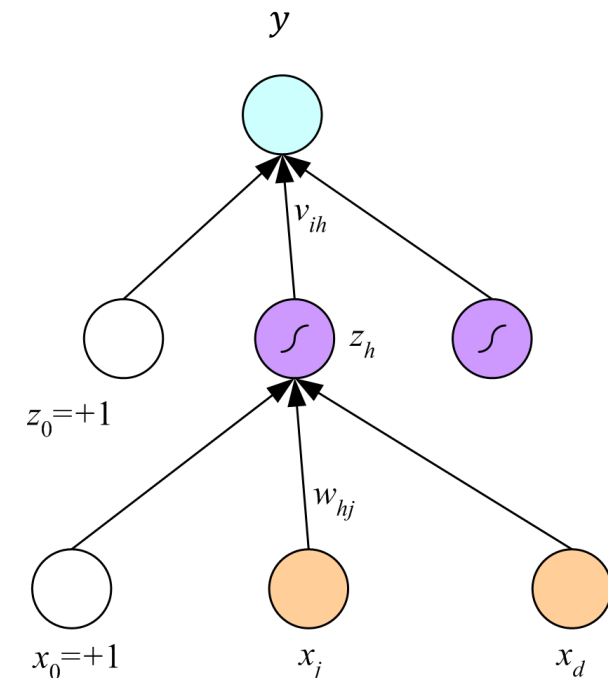
- ❑ Perceptrons
- ❑ Multilayer Perceptrons
- ❑ Error Backpropagation Algorithm
  - Regression
  - Classification
  - Deep Neural Networks
- ❑ Applications

# Learning for a Sigmoid Output Node

□ Weight update for a single sigmoid output node

- $E(\mathbf{W}_z, \mathbf{w}_o | \mathbf{x}, r) \equiv -r \log y - (1 - r) \log(1 - y) \quad ; y = \frac{1}{1 + \exp(-s)}, s = \mathbf{w}_o^T \mathbf{z}$
- $$\begin{aligned} \frac{\partial E}{\partial w_{o,i}} &= \frac{\partial E}{\partial y} \frac{\partial y}{\partial s} \frac{\partial s}{\partial w_{o,i}} \\ &= \underbrace{\left( -\frac{r}{y} + \frac{1-r}{1-y} \right)}_{\frac{\partial E}{\partial y}} \underbrace{y(1-y)}_{\frac{\partial y}{\partial s}} \underbrace{z_i}_{\frac{\partial s}{\partial w_{o,i}}} \\ &= -(r - y)z_i \end{aligned}$$
- $w_{o,i} \leftarrow w_{o,i} + \eta(r - y)z_i$

$$z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^T \mathbf{x})}$$



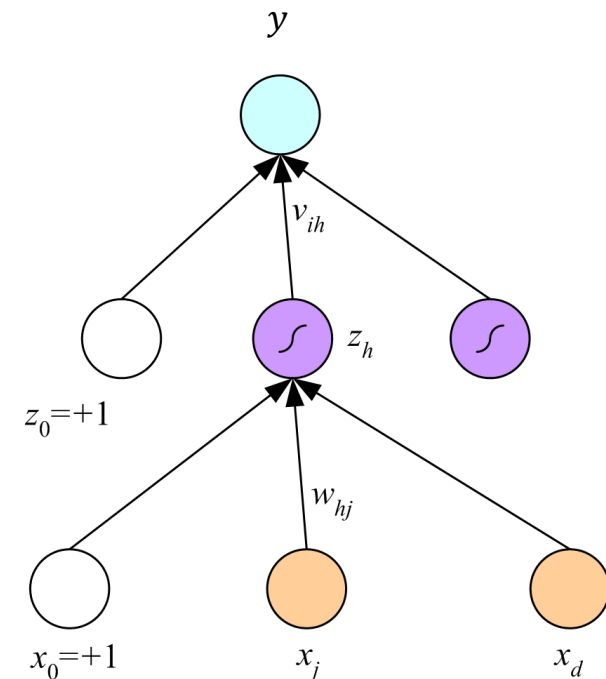
# Learning for Hidden Nodes

□ Weight update for a hidden node connected to a single sigmoid output node

- $E(\mathbf{W}_z, \mathbf{w}_o | \mathbf{x}, r) \equiv -r \log y - (1 - r) \log(1 - y)$  ;  $y = \frac{1}{1 + \exp(-\mathbf{w}_o^\top \mathbf{z})}$
- $\frac{\partial E}{\partial w_{z,ij}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial z_i} \frac{\partial z_i}{\partial w_{z,ij}}$  ;  $z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$ 

$$= \underbrace{\left(-\frac{r}{y} + \frac{1-r}{1-y}\right)}_{\frac{\partial E}{\partial y}} \underbrace{y(1-y)w_{o,i}}_{\frac{\partial y}{\partial z_i}} \underbrace{z_i(1-z_i)x_j}_{\frac{\partial z_i}{\partial w_{z,ij}}}$$

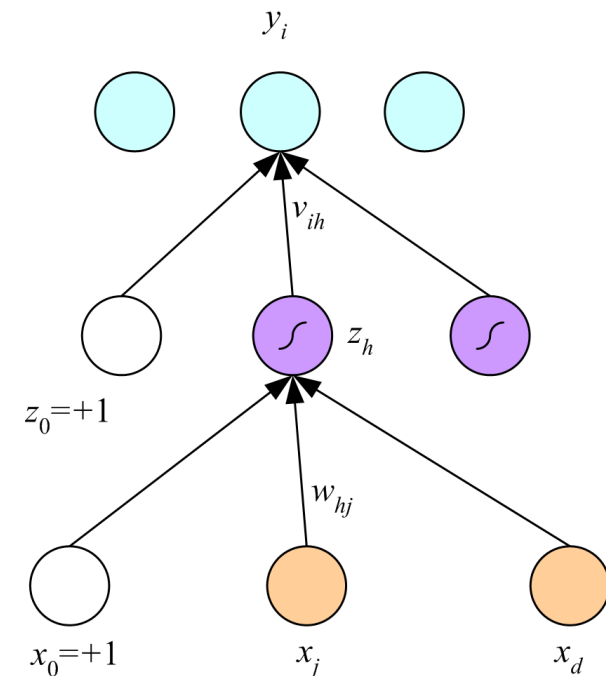
$$= -(r - y)w_{o,i}z_i(1 - z_i)x_j$$
- $w_{z,ij} \leftarrow w_{z,ij} + \eta(r - y)w_{o,i}z_i(1 - z_i)x_j$



# Learning for Softmax Output Nodes

□ Weight update for a softmax output node

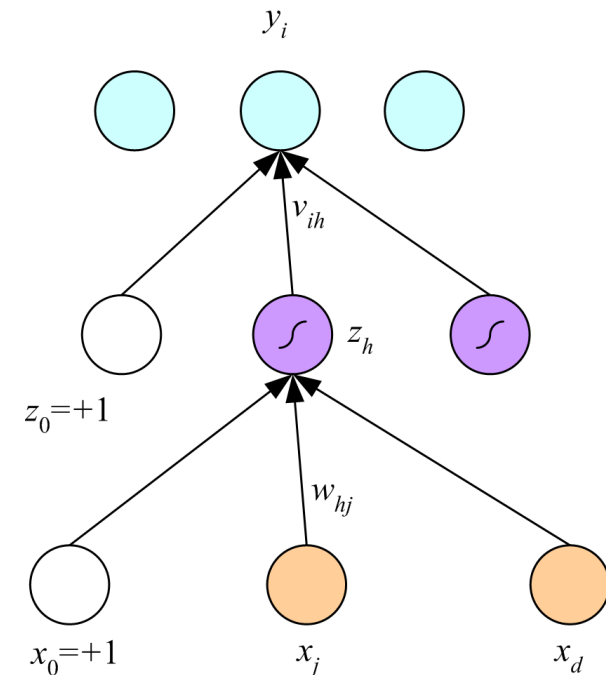
- $E(\mathbf{W}_z, \mathbf{W}_o | \mathbf{x}, \mathbf{r}) \equiv - \sum_k \underbrace{r_k \log y_k}_{E_k}$  ;  $y_k = \frac{\exp(s_k)}{\sum_k \exp(s_k)}$ ,  $s_i = \mathbf{w}_{o,i}^\top \mathbf{z}$
- $\frac{\partial E}{\partial w_{o,ij}} = - \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial s_i} \frac{\partial s_i}{\partial w_{o,ij}}$  ;  $z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$
- $$= - \left( \sum_{k \neq i} \underbrace{\frac{r_k}{y_k}}_{\frac{\partial E_k}{\partial y_k}} \underbrace{\left( - \frac{\exp(s_k) \exp(s_i)}{(\sum_k \exp(s_k))^2} \right)}_{\frac{\partial y_k}{\partial s_i}} + \underbrace{\frac{r_i}{y_i}}_{\frac{\partial E_i}{\partial y_i}} \underbrace{\left( \frac{\exp(s_i)}{\sum_k \exp(s_k)} - \frac{\exp(s_i) \exp(s_i)}{(\sum_k \exp(s_k))^2} \right)}_{\frac{\partial y_i}{\partial s_i}} \right) \underbrace{z_j}_{\frac{\partial s_i}{\partial w_{o,ij}}}$$
- $$= - \left( \sum_{k \neq i} \frac{r_k}{y_k} (-y_k y_i) + \frac{r_i}{y_i} (y_i - y_i^2) \right) z_j$$
- $$= - (\sum_{k \neq i} r_k (-y_i) + r_i (1 - y_i)) z_j$$
- $$= - (\sum_k r_k (-y_i) + r_i) z_j$$
- $$= -(r_i - y_i) z_j$$
- $$w_{o,ij} \leftarrow w_{o,ij} + \eta (r_i - y_i) z_j$$



# Learning for Hidden Nodes

□ Weight update for a hidden node connected to multiple softmax output nodes

- $E(\mathbf{W}_z, \mathbf{W}_o | \mathbf{x}, \mathbf{r}) \equiv - \sum_k r_k \underbrace{\log y_k}_{E_k}$  ;  $y_k = \frac{\exp(s_k)}{\sum_{\hat{k}} \exp(s_{\hat{k}})}$ ,  $s_i = \mathbf{w}_{o,i}^\top \mathbf{z}$
- $\frac{\partial E}{\partial w_{z,ij}} = - \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_i} \frac{\partial z_i}{\partial w_{z,ij}}$  ;  $z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$
- $w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k - y_k) w_{o,ki} z_i (1 - z_i) x_j$



# Learning for Output Nodes

- Weight update for a single sigmoid output node

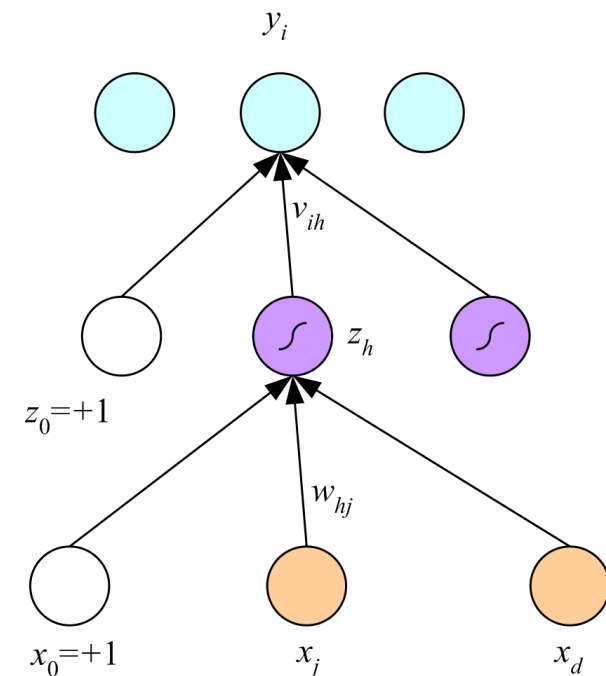
$$\blacksquare w_{o,i} \leftarrow w_{o,i} + \eta(r - y)z_i \quad ; y = \frac{1}{1+\exp(-\mathbf{w}_o^\top \mathbf{z})}, z_i = \frac{1}{1+\exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$$

- Weight update for multiple softmax output nodes

$$\blacksquare w_{o,ij} \leftarrow w_{o,ij} + \eta(r_i - y_i)z_j \quad ; y_i = \frac{\exp(s_i)}{\sum_l \exp(s_l)}, s_i = \mathbf{w}_{o,i}^\top \mathbf{z}$$

- Weight update for multiple linear output nodes

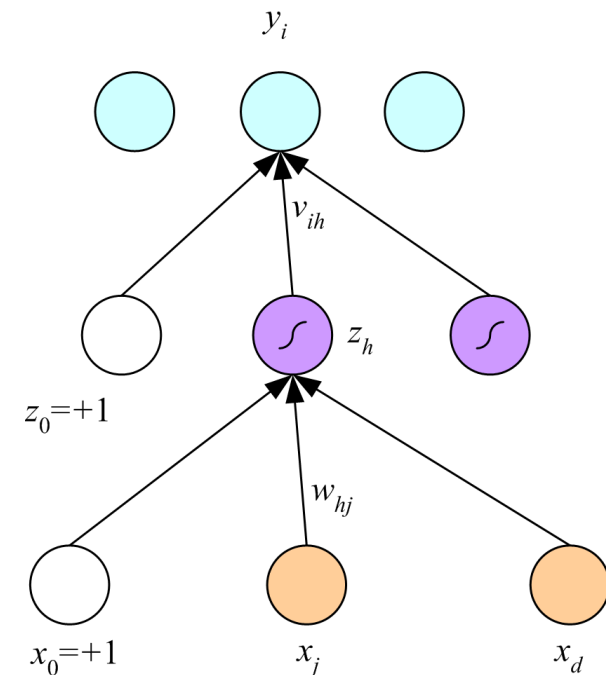
$$\blacksquare w_{o,ij} \leftarrow w_{o,ij} + \eta(r_i - y_i)z_j \quad ; y_i = \mathbf{w}_{o,i}^\top \mathbf{z}$$





# Learning for Hidden Nodes

- Weight update for a hidden node connected to a single sigmoid output node
  - $w_{z,ij} \leftarrow w_{z,ij} + \eta(r - y)w_{o,i}z_i(1 - z_i)x_j$
- Weight update for a hidden node connected to multiple softmax output nodes
  - $w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k - y_k)w_{o,ki} z_i(1 - z_i)x_j$
- Weight update for a hidden node connected to multiple linear output nodes
  - $w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k - y_k)w_{o,ki} z_i(1 - z_i)x_j$

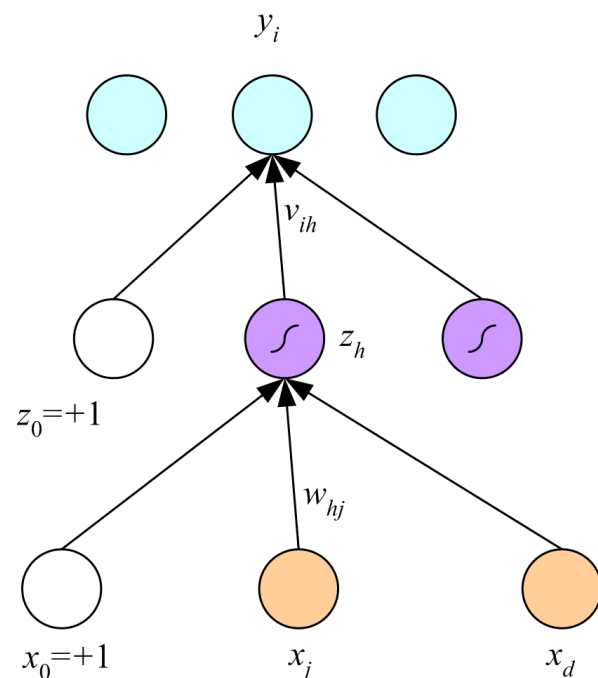


# Stochastic Gradient Descent for MLPs

□ Input:  $\mathcal{D} = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$

1. Initialize  $\mathbf{W}_z$  and  $\mathbf{W}_o$  with small random numbers.
2. **repeat**
3.   **for** each  $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{D}$  in random order
4.     **for** each hidden node index  $i$
5.        $z_i \leftarrow 1/(1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x}^t))$
6.     **end**
7.     **for** each output node index  $i$
8.        $s_i \leftarrow \mathbf{w}_{o,i}^\top \mathbf{z}$
9.     **end**
10.    **for** each output node index  $i$
11.       $y_i \leftarrow \exp(s_i) / \sum_i \exp(s_i)$
12.    **end**
13.    **for** each hidden node index  $i$
14.     **for** each input node index  $j$
15.        $w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k^t - y_k) w_{o,ki} z_i (1 - z_i) x_j^t$
16.     **end**
17.    **end**
18.    **for** each output node index  $i$
19.     **for** each hidden node index  $j$
20.        $w_{o,ij} \leftarrow w_{o,ij} + \eta (r_i^t - y_i) z_j$
21.     **end**
22.    **end**
23. **end**
24. **until** convergence

; sigmoid:  $y \leftarrow 1/(1 + \exp(-\mathbf{w}^\top \mathbf{x}))$



□ Output:  $\mathbf{W}_z$  and  $\mathbf{W}_o$

# Contents

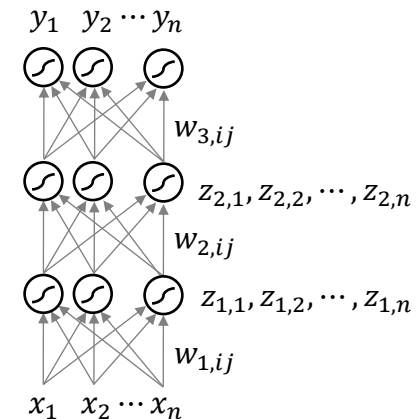
- ☐ Perceptrons
- ☐ Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
  - Regression
  - Classification
  - Deep Neural Networks
- ☐ Applications

# Deep Neural Networks

## ❑ Deep neural network (DNN)

- Multiple hidden layers
  - “deep and narrow” vs. “shallow and wide”
- Representation learning ; more abstract concept
  - e.g., Image recognition: pixels  $\rightarrow$  edges  $\rightarrow$  corners  $\rightarrow$  digits

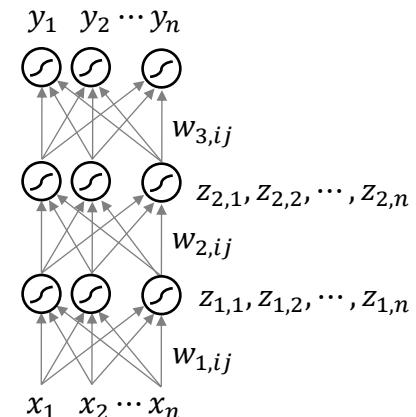
## ❑ Deep learning



# Learning for Multiple Hidden Layers

- Weight update for a hidden node of the first hidden layer in a two hidden layer neural network with multiple linear output nodes

- $E(\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3 | \mathbf{x}, \mathbf{r}) \equiv \sum_k \underbrace{\frac{1}{2}(r_k - y_k)^2}_{E_k} \quad ; y_k = \sum_m w_{3,km} z_{2,m}$
- $\frac{\partial E}{\partial w_{1,ij}} = \sum_k \frac{\partial E_k}{\partial y_k} \underbrace{\frac{\partial y_k}{\partial z_{1,i}}}_{\frac{\partial \sum_m w_{3,km} z_{2,m}}{\partial z_{1,i}}} \frac{\partial z_{1,i}}{\partial w_{1,ij}} \quad ; z_{2,m} = \frac{1}{1 + \exp(-\sum_i w_{2,mi} z_{1,i})}$   
 $\quad \quad \quad \frac{\partial \sum_m w_{3,km} z_{2,m}}{\partial z_{1,i}} = \sum_m \frac{\partial w_{3,km} z_{2,m}}{\partial z_{2,m}} \frac{\partial z_{2,m}}{\partial z_{1,i}} \quad ; z_{1,i} = \frac{1}{1 + \exp(-\sum_j w_{1,ij} x_j)}$   
 $\quad \quad \quad = \sum_k \underbrace{-(r_k - y_k)}_{\frac{\partial E_k}{\partial y_k}} \underbrace{\sum_m w_{3,km} z_{2,m}}_{\frac{\partial w_{3,km} z_{2,m}}{\partial z_{2,m}}} \underbrace{(1 - z_{2,m}) w_{2,mi}}_{\frac{\partial z_{2,m}}{\partial z_{1,i}}} \underbrace{z_{1,i} (1 - z_{1,i}) x_j}_{\frac{\partial z_{1,i}}{\partial w_{1,ij}}}$
- $w_{1,ij} \leftarrow w_{1,ij} + \eta \sum_k (r_k - y_k) \sum_m w_{3,km} z_{2,m} (1 - z_{2,m}) w_{2,mi} z_{1,i} (1 - z_{1,i}) x_j$



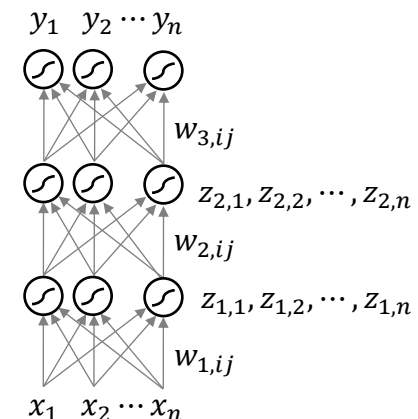
# Learning for Multiple Hidden Layers

- Weight update for a hidden node of the first hidden layer in a two hidden layer neural network with multiple softmax output nodes

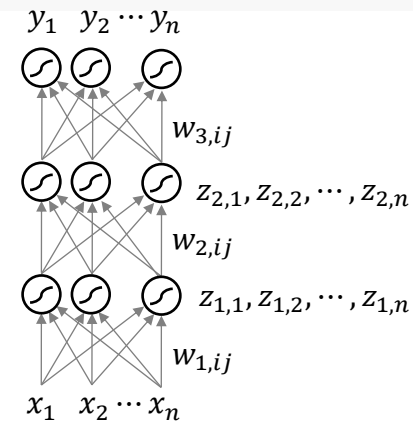
- $E(\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3 | \mathbf{x}, \mathbf{r}) \equiv - \sum_k \underbrace{r_k \log y_k}_{E_k} \quad ; y_k = \frac{\exp(s_k)}{\sum_{\hat{k}} \exp(s_{\hat{k}})}, s_k = \sum_m w_{3,km} z_{2,m}$   
 $z_{2,m} = \frac{1}{1 + \exp(-\sum_i w_{2,mi} z_{1,i})}$   
 $z_{1,i} = \frac{1}{1 + \exp(-\sum_j w_{1,ij} x_j)}$
- $w_{1,ij} \leftarrow w_{1,ij} + \eta \sum_k (r_k - y_k) \sum_m w_{3,km} z_{2,m} (1 - z_{2,m}) w_{2,mi} z_{1,i} (1 - z_{1,i}) x_j$

- Time complexity

- $\mathcal{O}(n^L)$



# Learning for DNNs



## □ Weight update for an output node

- $\frac{\partial E}{\partial w_{3,ij}} = \sum_k \frac{\partial E_k}{\partial y_i} \frac{\partial y_i}{\partial w_{3,ij}} = \sum_k \frac{\partial E_k}{\partial y_i} \frac{\partial y_i}{\partial s_{3,i}} \frac{\partial s_{3,i}}{\partial w_{3,ij}}$
- $w_{3,ij} \leftarrow w_{3,ij} + \eta \underbrace{(r_i - y_i)}_{\delta_{3,i}} z_{2,j} \quad ; s_{l,i} \equiv \sum_j w_{l,ij} z_{l-1,j}$
- $w_{3,ij} \leftarrow w_{3,ij} + \eta \delta_{3,i} z_{2,j} \quad ; \delta_{3,i} \equiv r_i - y_i = \frac{\partial E}{\partial s_{3,i}}$

## □ Weight update for a hidden node connected to multiple output nodes

- $\frac{\partial E}{\partial w_{2,ij}} = \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_{2,i}} \frac{\partial z_{2,i}}{\partial w_{2,ij}} = \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_{2,i}} \frac{\partial z_{2,i}}{\partial s_{2,i}} \frac{\partial s_{2,i}}{\partial w_{2,ij}}$
- $w_{2,ij} \leftarrow w_{2,ij} + \eta \underbrace{\sum_k (r_k - y_k) w_{3,ki}}_{\delta_{2,i}} z'_{2,i} z_{1,j}$
- $w_{2,ij} \leftarrow w_{2,ij} + \eta \delta_{2,i} z_{1,j} \quad ; \delta_{2,i} \equiv \sum_k \delta_{3,k} w_{3,ki} z'_{2,i} = \frac{\partial E}{\partial s_{2,i}}$

## □ Weight update for a hidden node of the first hidden layer in a two hidden layer neural network with multiple output nodes

- $\frac{\partial E}{\partial w_{1,ij}} = \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_{1,i}} \frac{\partial z_{1,i}}{\partial w_{1,ij}} = \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_{1,i}} \frac{\partial z_{1,i}}{\partial s_{1,i}} \frac{\partial s_{1,i}}{\partial w_{1,ij}}$
- $w_{1,ij} \leftarrow w_{1,ij} + \eta \sum_k (r_k - y_k) \sum_m w_{3,km} z'_{2,m} w_{2,mi} z'_{1,i} x_j$
- $w_{1,ij} \leftarrow w_{1,ij} + \eta \underbrace{\sum_m \sum_k (r_k - y_k) w_{3,km} z'_{2,m} w_{2,mi} z'_{1,i}}_{\delta_{1,i}} x_j$
- $w_{1,ij} \leftarrow w_{1,ij} + \eta \delta_{1,i} x_j \quad ; \delta_{1,i} \equiv \sum_m \delta_{2,m} w_{2,mi} z'_{1,i} = \frac{\partial E}{\partial s_{1,i}}$

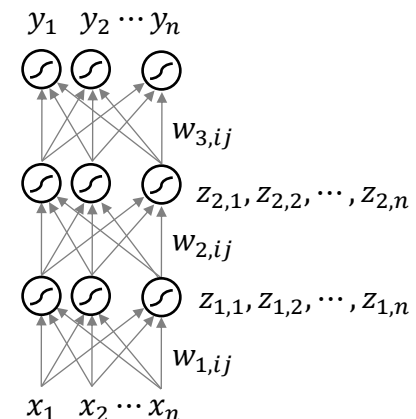
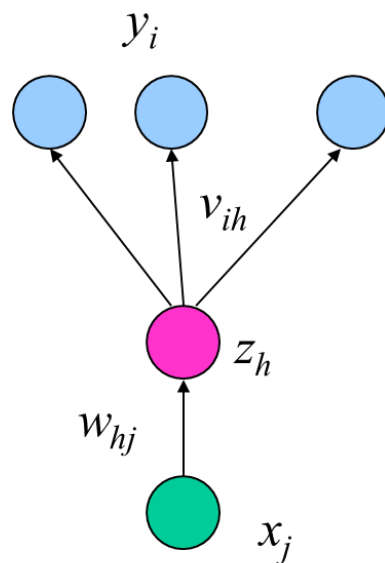
# Learning for DNNs

## □ Weight update for an output node

$$\blacksquare \quad w_{L,ij} \leftarrow w_{L,ij} + \eta \underbrace{\delta_{L,i}}_{\frac{\partial E}{\partial s_{L,i}}} \underbrace{z_{L-1,j}}_{\frac{\partial s_{L,i}}{\partial w_{L,ij}}} \quad ; \delta_{L,i} \equiv r_i - y_i, \partial s_{L,i} \equiv \sum_j w_{L,ij} z_{L-1,j}$$

## □ Weight update for a hidden node of the $l$ -th hidden layer in a multiple hidden layer neural network

$$\blacksquare \quad w_{l,ij} \leftarrow w_{l,ij} + \eta \underbrace{\delta_{l,i}}_{\frac{\partial E}{\partial s_{l,i}}} \underbrace{z_{l-1,j}}_{\frac{\partial s_{l,i}}{\partial w_{l,ij}}} \quad ; \delta_{l,i} \equiv \sum_m \delta_{l+1,m} w_{l+1,mi} z'_{l,i}$$





# Stochastic Gradient Descent for DNNs

□ Input:  $\mathcal{D} = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$

1. Initialize  $\{\mathbf{W}_l\}_{l=1}^L$  with small random numbers.

2. **repeat**

3.     **for** each  $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{D}$  in **random order**

4.         **for** each hidden layer index  $l$  from 1 to  $L - 1$

5.             **for** each node index  $i$  in the  $l$ -th layer

6.                  $z_{l,i} \leftarrow 1/(1 + \exp(-\mathbf{w}_{l,i}^\top \mathbf{z}_{l-1}))$

7.             **end**

8.         **end**

9.         **for** each output node index  $i$

10.              $y_i \leftarrow \exp(\mathbf{w}_{L,i}^\top \mathbf{z}_{L-1}) / \sum_i \exp(\mathbf{w}_{L,i}^\top \mathbf{z}_{L-1})$

11.         **end**

12.         **for** each output node index  $i$

13.              $\delta_{L,i} \leftarrow r_i^t - y_i$

14.         **end**

15.         **for** each layer index  $l$  from  $L$  to 1

16.             **for** each node index  $j$  in the  $(l - 1)$ -th layer

17.                 **if**  $l \neq 1$  **then**  $\delta_{l-1,j} \leftarrow \sum_i \delta_{l,i} w_{l,i,j} z'_{l-1,j}$

18.                 **for** each node index  $i$  in the  $l$ -th layer

19.                      $w_{l,i,j} \leftarrow w_{l,i,j} + \eta \delta_{l,i} z_{l-1,j}$

20.                 **end**

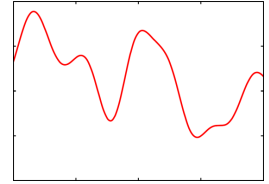
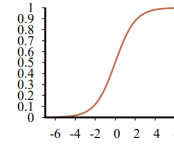
21.             **end**

22.         **end**

23.     **end**

24. **until** **convergence**

□ Output: all  $\{\mathbf{W}_l\}_{l=1}^L$



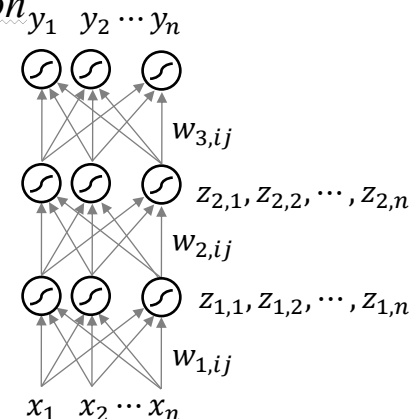
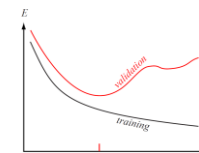
; forward process: feed forward

;  $\mathbf{z}_0 \equiv \mathbf{x}^t$

}  $y_i \leftarrow \mathbf{w}_{L,i}^\top \mathbf{z}_{L-1}$   
 $y \leftarrow 1/(1 + \exp(-\mathbf{w}_L^\top \mathbf{z}_{L-1}))$

; backward process

; error back-propagation

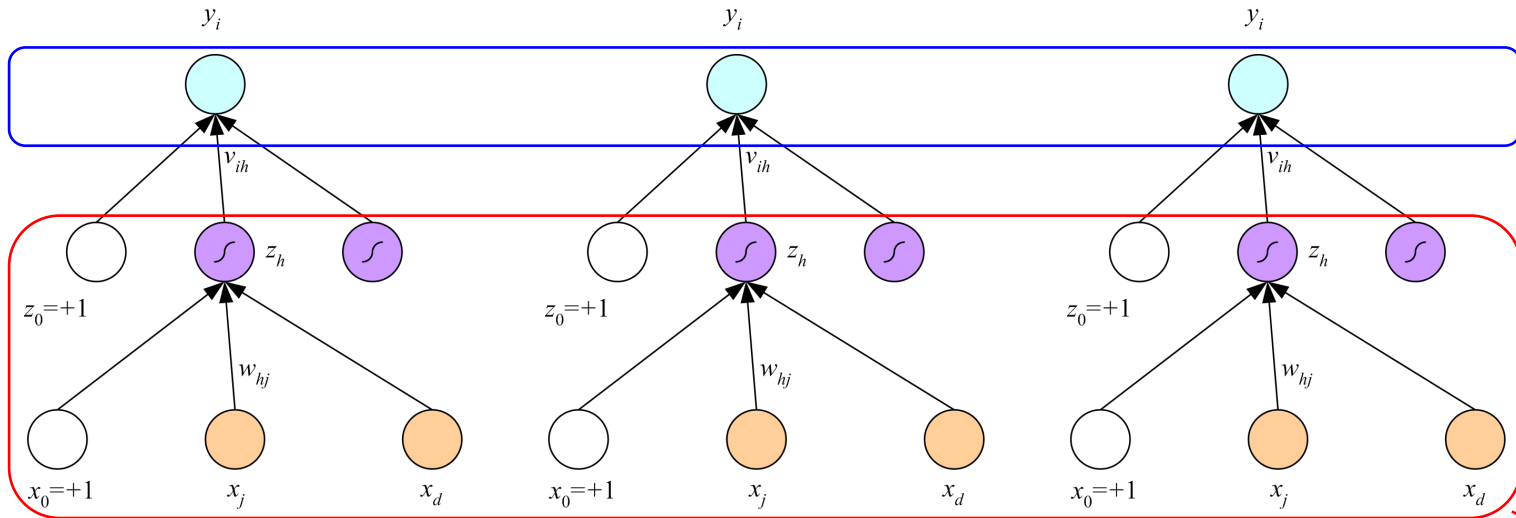


# Contents

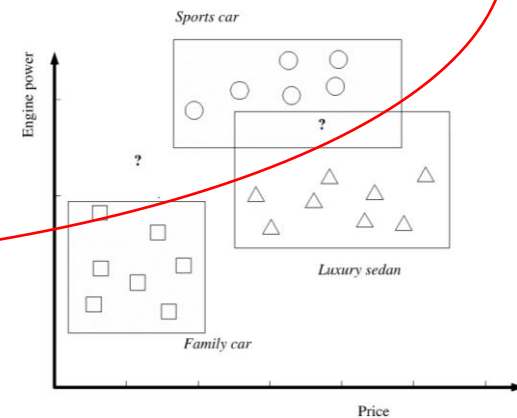
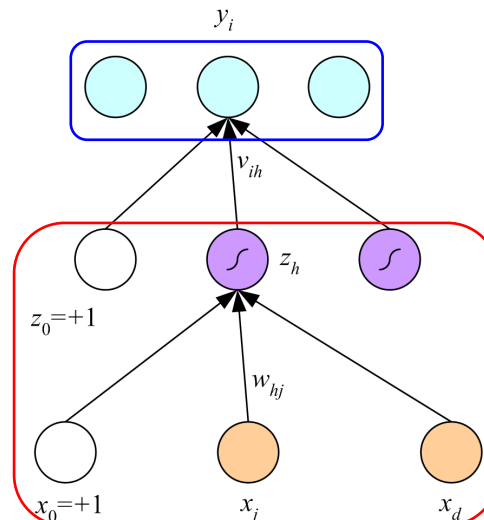
- ☐ Perceptrons
- ☐ Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications
  - Multilabel Classification
  - Dimensionality Reduction
  - Representation Learning

# Multilabel Classification

- $K$  independent two-class classification



- Multitask learning



# Learning for Multilabel Classification

## □ $K$ dependent two-class classification

- $E(\mathbf{W}_z, \mathbf{W}_o | \mathbf{x}, \mathbf{r}) \equiv \sum_k \underbrace{[-r_k \log y_k - (1 - r_k) \log(1 - y_k)]}_{E_k}$
- Weight update for  $K$  sigmoid output nodes

$$\begin{aligned} \bullet \quad \frac{\partial E}{\partial w_{o,ij}} &= \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial w_{o,ij}} \\ &= \underbrace{\left(-\frac{r_i}{y_i} + \frac{1-r_i}{1-y_i}\right)}_{\frac{\partial E_i}{\partial y_i}} \underbrace{y_i(1-y_i)z_j}_{\frac{\partial y_i}{\partial w_{o,ij}}} \\ &= -(r_i - y_i)z_j \end{aligned}$$

$$\bullet \quad w_{o,ij} \leftarrow w_{o,ij} + \eta(r_i - y_i)z_j$$

- Weight update for a hidden node connected to  $K$  sigmoid output nodes

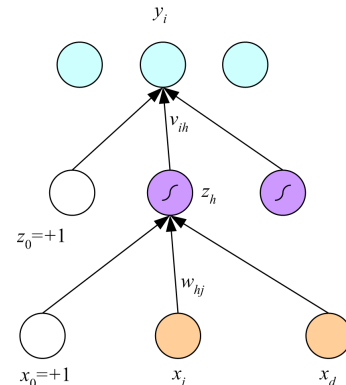
$$\begin{aligned} \bullet \quad \frac{\partial E}{\partial w_{z,ij}} &= \sum_k \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial z_i} \frac{\partial z_i}{\partial w_{z,ij}} \\ &= \sum_k \underbrace{\left(-\frac{r_k}{y_k} + \frac{1-r_k}{1-y_k}\right)}_{\frac{\partial E_k}{\partial y_k}} \underbrace{y_k(1-y_k)w_{o,ki}}_{\frac{\partial y_k}{\partial z_i}} \underbrace{z_i(1-z_i)x_j}_{\frac{\partial z_i}{\partial w_{z,ij}}} \\ &= -\sum_k (r_k - y_k)w_{o,ki} z_i(1-z_i)x_j \end{aligned}$$

$$\bullet \quad w_{z,ij} \leftarrow w_{z,ij} + \eta \sum_k (r_k - y_k)w_{o,ki} z_i(1-z_i)x_j$$

$$r_k = \begin{cases} 1 & \text{if } x \text{ has label } k \\ 0 & \text{otherwise} \end{cases}$$

$$y_k = \frac{1}{1 + \exp(-\mathbf{w}_{o,k}^\top \mathbf{z})}$$

$$z_i = \frac{1}{1 + \exp(-\mathbf{w}_{z,i}^\top \mathbf{x})}$$



# Contents

- ☐ Perceptrons
- ☐ Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications
  - Multilabel Classification
  - Dimensionality Reduction
  - Representation Learning

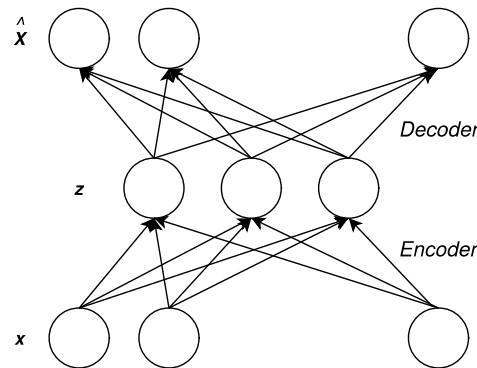
# Autoencoders

## □ Autoencoders [Cottrell, Munro, and Zipser 1987]

- $\mathbf{z} = \text{ENCODE}(\mathbf{x}|\mathbf{W})$
- $\tilde{\mathbf{x}} = \text{DECODE}(\mathbf{z}|\mathbf{V})$

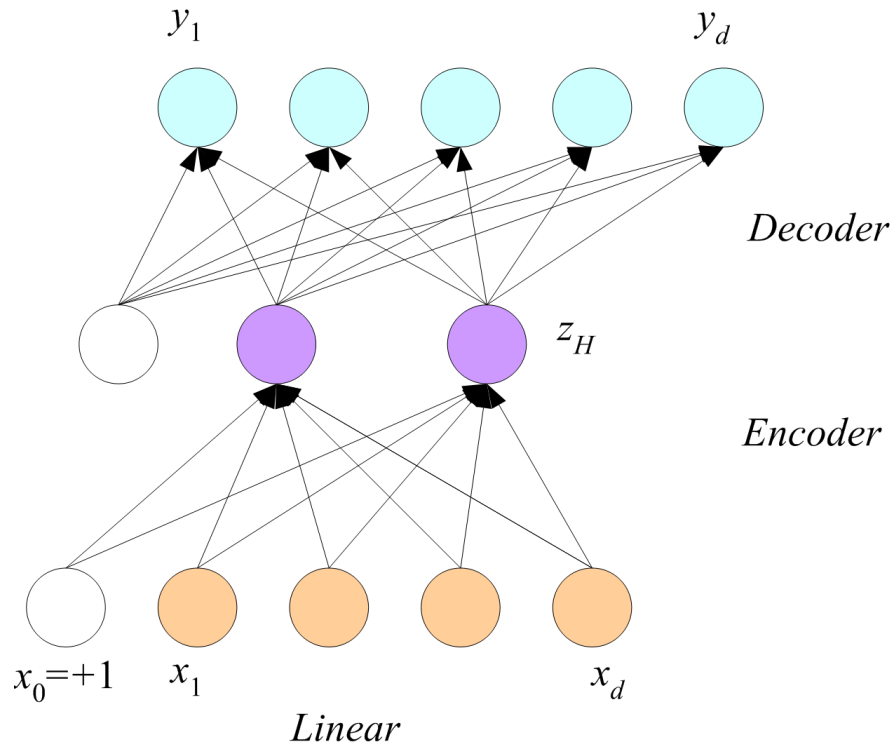
## □ Reconstruction error of an autoencoder

- $$E(\mathbf{W}, \mathbf{V}|\mathcal{D}) = \frac{1}{N} \sum_t \|\mathbf{x}^t - \tilde{\mathbf{x}}^t\|^2 \quad ; \quad \|\mathbf{x}\|_2 \equiv \sqrt{\sum_i x_i^2} \text{ (Euclidean norm or } \ell^2 \text{ norm)}$$
$$= \frac{1}{N} \sum_t \|\mathbf{x}^t - \text{DECODE}(\text{ENCODE}(\mathbf{x}^t|\mathbf{W})|\mathbf{V})\|^2$$

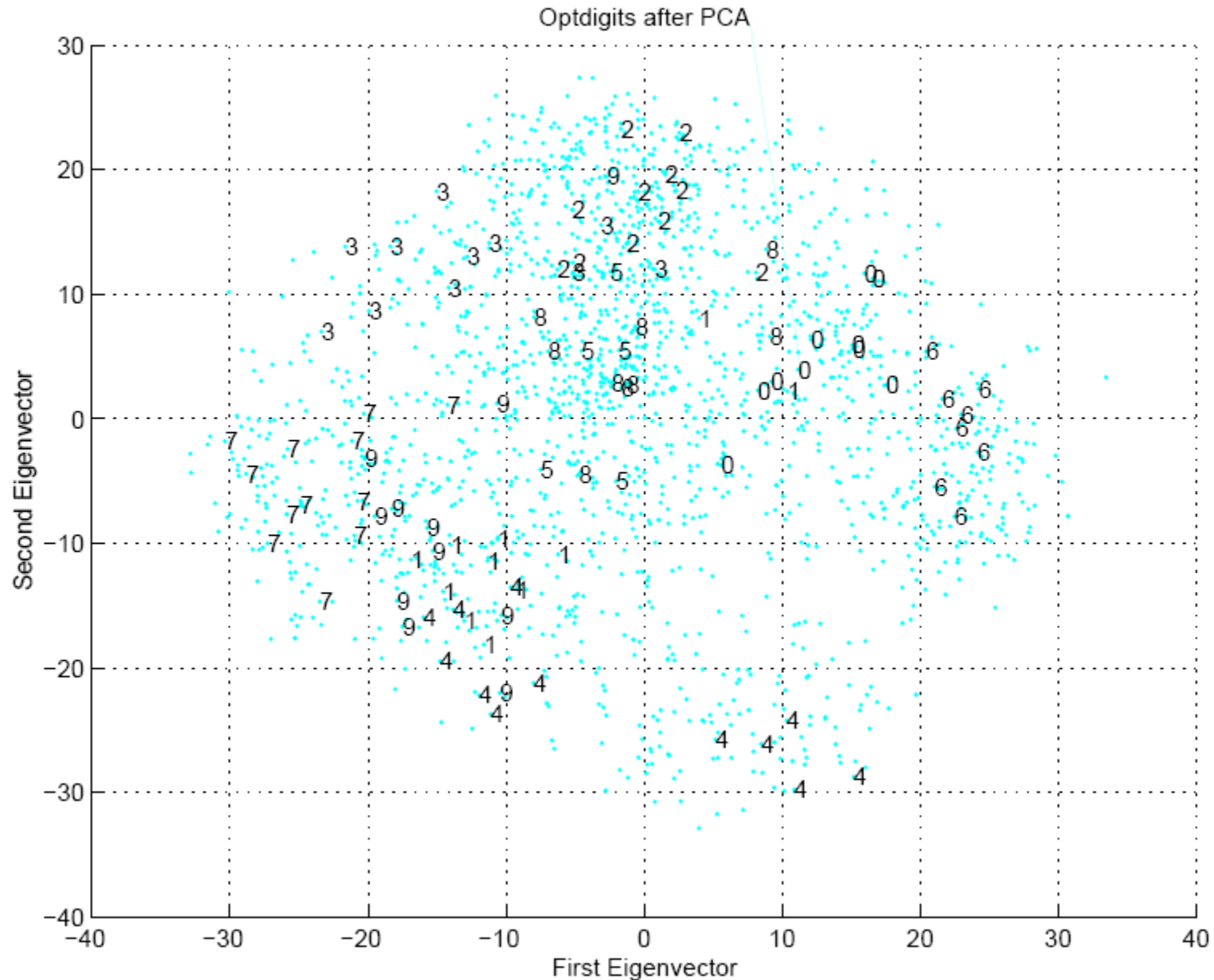


# Dimensionality Reduction Using Autoencoders

- Dimensionality reduction using linear autoencoders



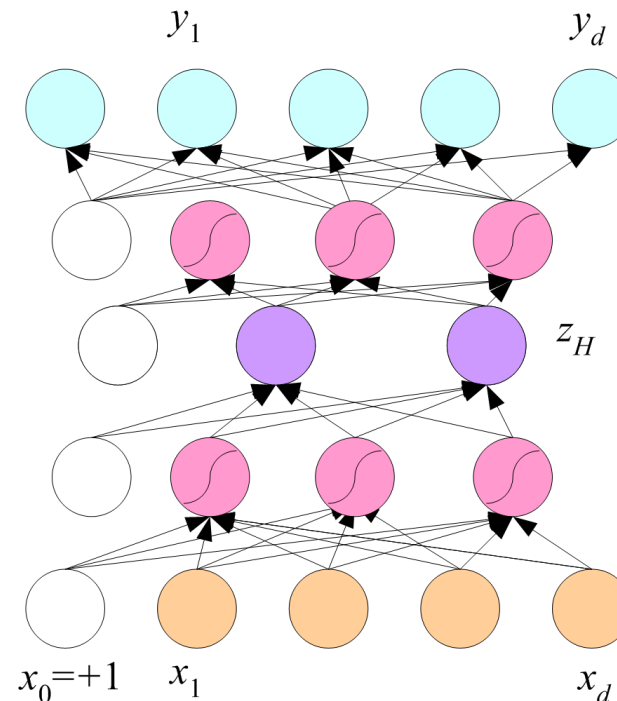
# Dimensionality Reduction Example



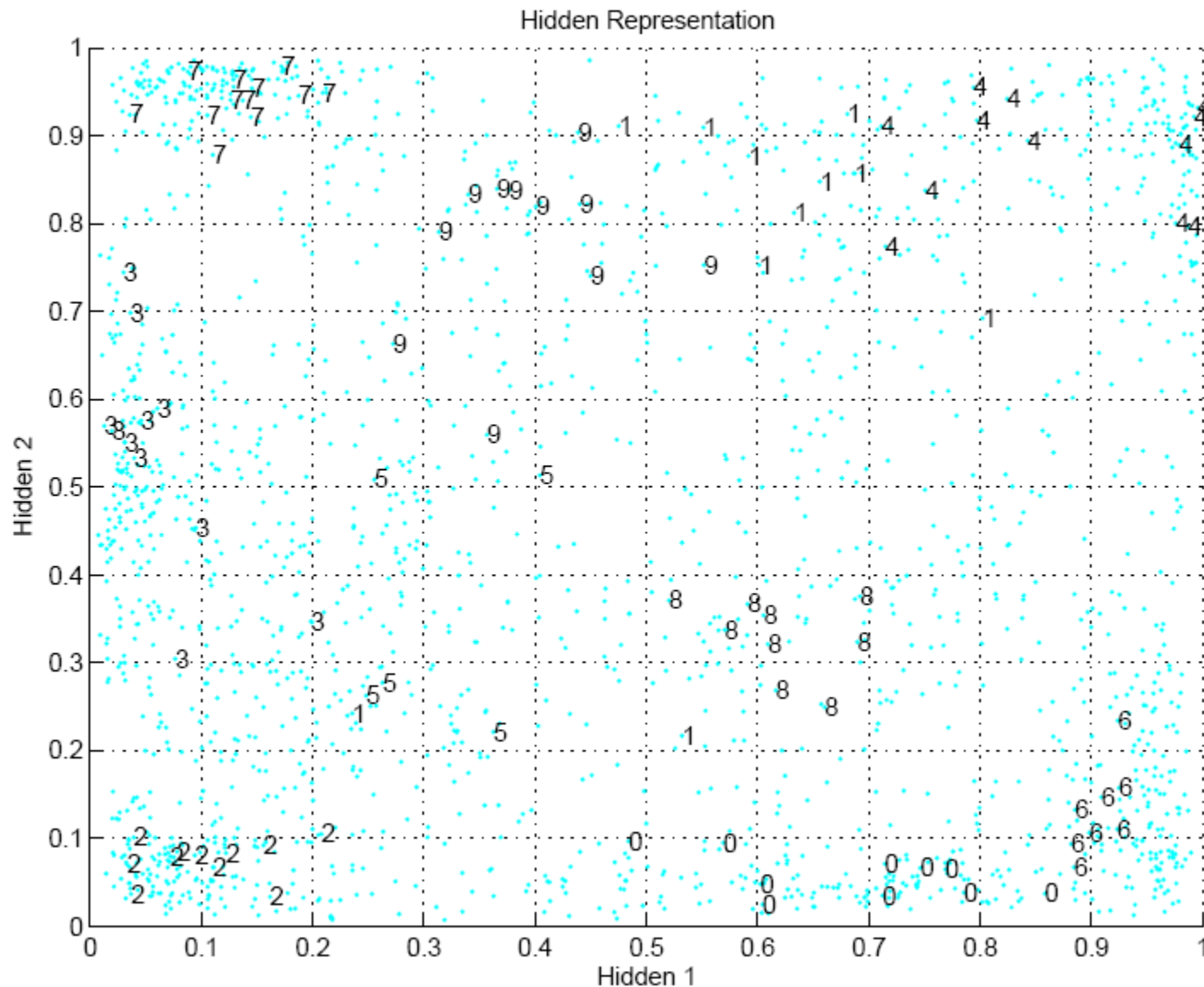


# Variants of Autoencoders

- ❑ Denoising autoencoder [Vincent *et al.* 2008]
  - Noiseless and noisy versions of the same instance with the same desired output
- ❑ Sparse autoencoder [Ranzato *et al.* 2007]
  - dimensionality of input < dimensionality of hidden layer
- ❑ Deep autoencoder and stacked autoencoder



# Dimensionality Reduction Example



# Multidimensional Scaling

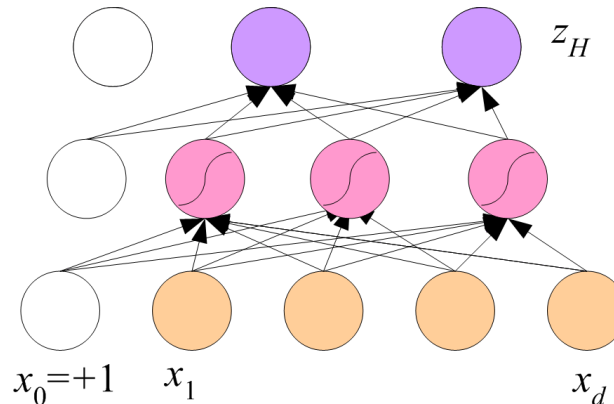
## □ Multidimensional scaling (MDS)

- high-dimensional space  $\Rightarrow$  low-dimensional space
- Trying to keep the inter-point distances unchanged
- Dimensionality reduction

## □ Sammon mapping

- $\mathbf{z} = \mathbf{g}(\mathbf{x}|\theta)$  where  $\mathbf{x} \in \mathcal{R}^d$ ,  $\mathbf{z} \in \mathcal{R}^k$ , and  $k < d$
- Sammon stress: 
$$\sum_{t,s} \frac{(\|\mathbf{g}(\mathbf{x}^t|\theta) - \mathbf{g}(\mathbf{x}^s|\theta)\| - \|\mathbf{x}^t - \mathbf{x}^s\|)^2}{\|\mathbf{x}^t - \mathbf{x}^s\|^2} = \sum_{t,s} \frac{(\|\mathbf{z}^t - \mathbf{z}^s\| - \|\mathbf{x}^t - \mathbf{x}^s\|)^2}{\|\mathbf{x}^t - \mathbf{x}^s\|^2}$$

## □ An MLP with $d$ inputs, hidden units, and $k < d$ output units



# Contents

- ☐ Perceptrons
- ☐ Multilayer Perceptrons
- ☐ Error Backpropagation Algorithm
- ☐ Applications
  - Multilabel Classification
  - Dimensionality Reduction
  - Representation Learning

# Representation Learning

## □ Linear model and nonlinear model

- $y = \sum_i v_i x_i + v_0$

- $y = \sum_i v_i \phi_i(\mathbf{x})$

;  $\phi_i$ : basis function

e.g., polynomial basis function

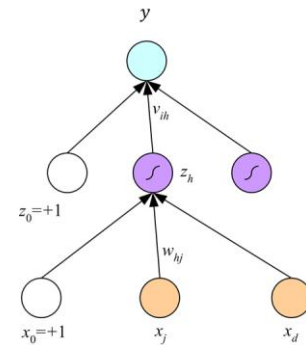
## □ Embedding (also called code or *representation*)

- $y = \sum_i v_i \phi(\mathbf{x}|\mathbf{w}_i)$  ;  $\phi(\mathbf{x}|\mathbf{w}_i) \equiv \frac{1}{1 + \exp(-\mathbf{w}_i^T \mathbf{x})}$

- Learnable basis function

- A new space where the problem becomes linear, e.g., XOR

- Unsupervised learning, e.g., autoencoder



## □ Dimensionality reduction: dimensionality of embedding < dimensionality of input

## □ Transfer learning: using the embedding trained on one task as input to another task

## □ Multitask learning: learning similar tasks with shared hidden layers

## □ Self-supervised learning and finetuning: using the embedding trained on a large unlabeled data set as the input in supervised learning with a small labeled data set

# Word Embedding

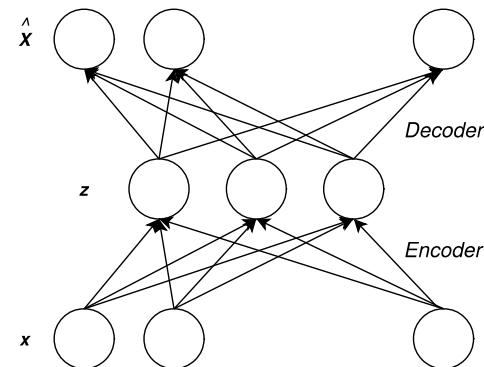
## ❑ Binary vector representation of words

- $d$ -dimensional one-hot vector
- High dimensional
- No syntactic/semantic relationship

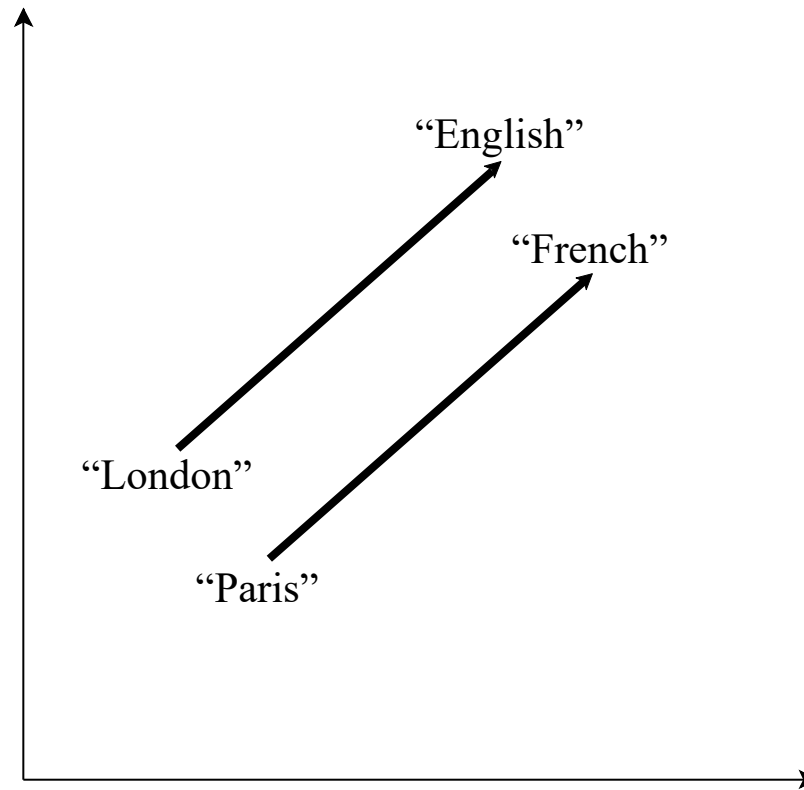
; bag of words

## ❑ Continuous vector representation of words (WORD2VEC, Mikolov *et al.* 2013)

- Input: context (windows of words)
  - Each word in the window at a time (SKIP-GRAM)  
e.g., “Protesters in Paris clash with the French police.”
- Hidden layer: less than  $d$  hidden units
- Output: center word ( $d$ -dimensional one-hot vector)



# Word Embedding Example



$$\text{VEC}(\text{"French"}) - \text{VEC}(\text{"Paris"}) + \text{VEC}(\text{"London"}) = \text{VEC}(\text{"English"})$$

# Summary and Preview

- ❑ Perceptrons
  - Introduction
  - Learning for Perceptrons
- ❑ Multilayer Perceptrons
  - Boolean Functions
  - Multilayer Perceptrons
- ❑ Error Backpropagation Algorithm
  - Regression
  - Classification
  - Deep Neural Networks
- ❑ Applications
  - Multilabel Classification
  - Dimensionality Reduction
  - Representation Learning
- ❑ Architecture of DNNs