

# CH. 11: Multilayer Perceptrons

## Ups and Downs of NN Study:

- 1958: Perceptron (linear model)
- 1969: Perceptron has limitation
- 1980: **Multi-Layer Perceptron** -- Do not have significant difference from DNN today
- 1986: **Backpropagation** -- Efficient for training multi-layer neural networks
- 2006: **Encoder-Decoder, Restricted Boltzmann Machine** -- For pre-training multi-layer neural networks
- 2009: **GPU**
- 2011: Start to be popular in speech and image processing

## 11.1 Introduction

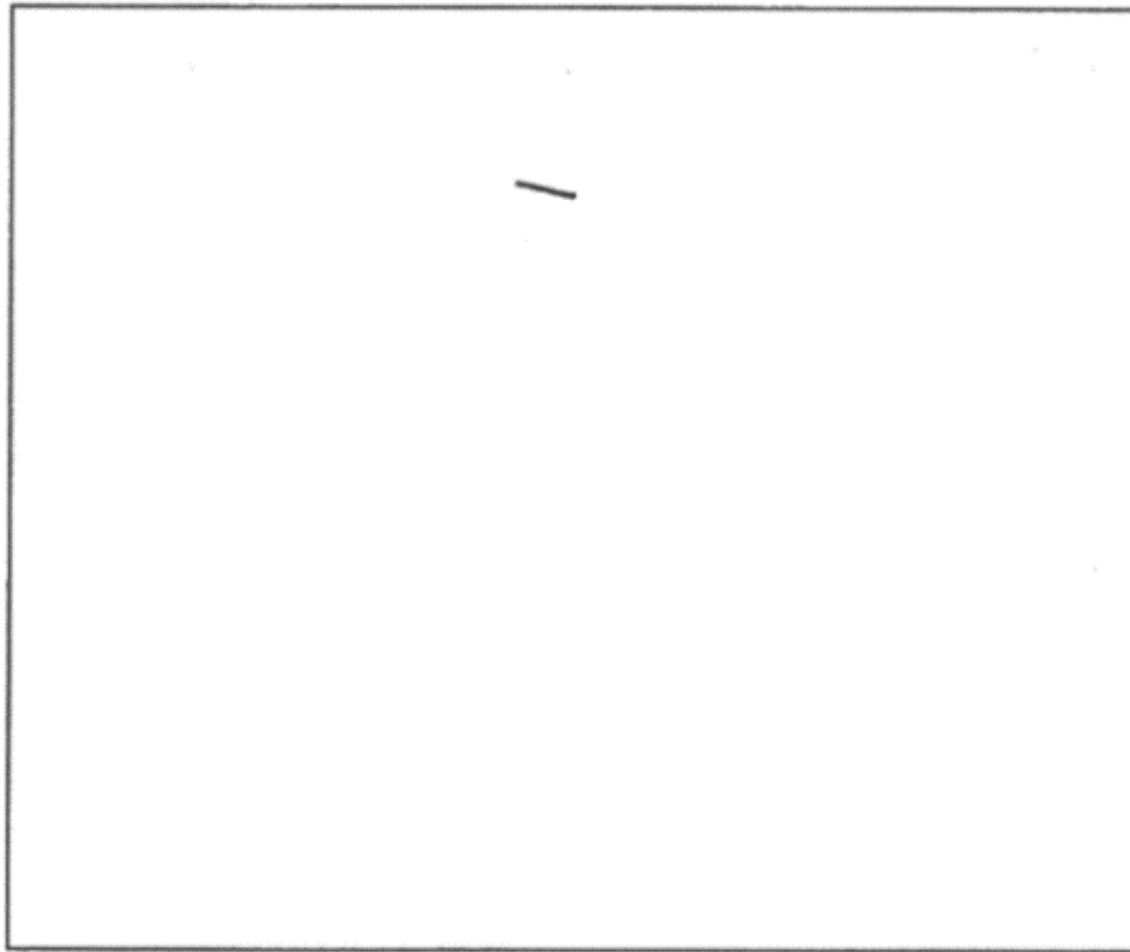
### □ ANN technology

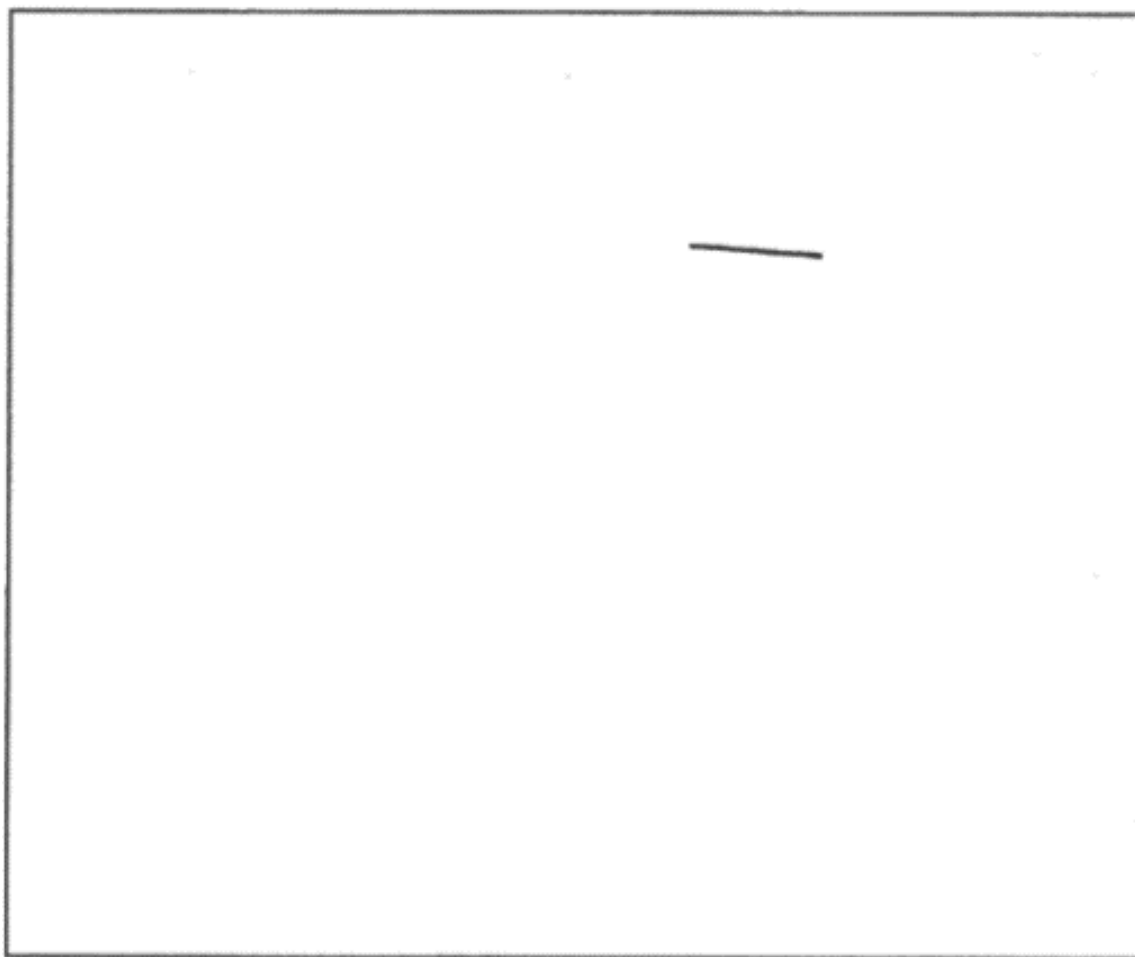
- Human brain is the most complex computing device that we have ever known.
- Conventional computers are good in scientific and mathematical computations, creation and manipulation of databases, control functions, graphics, word processing, .....
- How do computers learn, analyze, organize, adapt comprehend, associate, recognize, plan, decide, .....

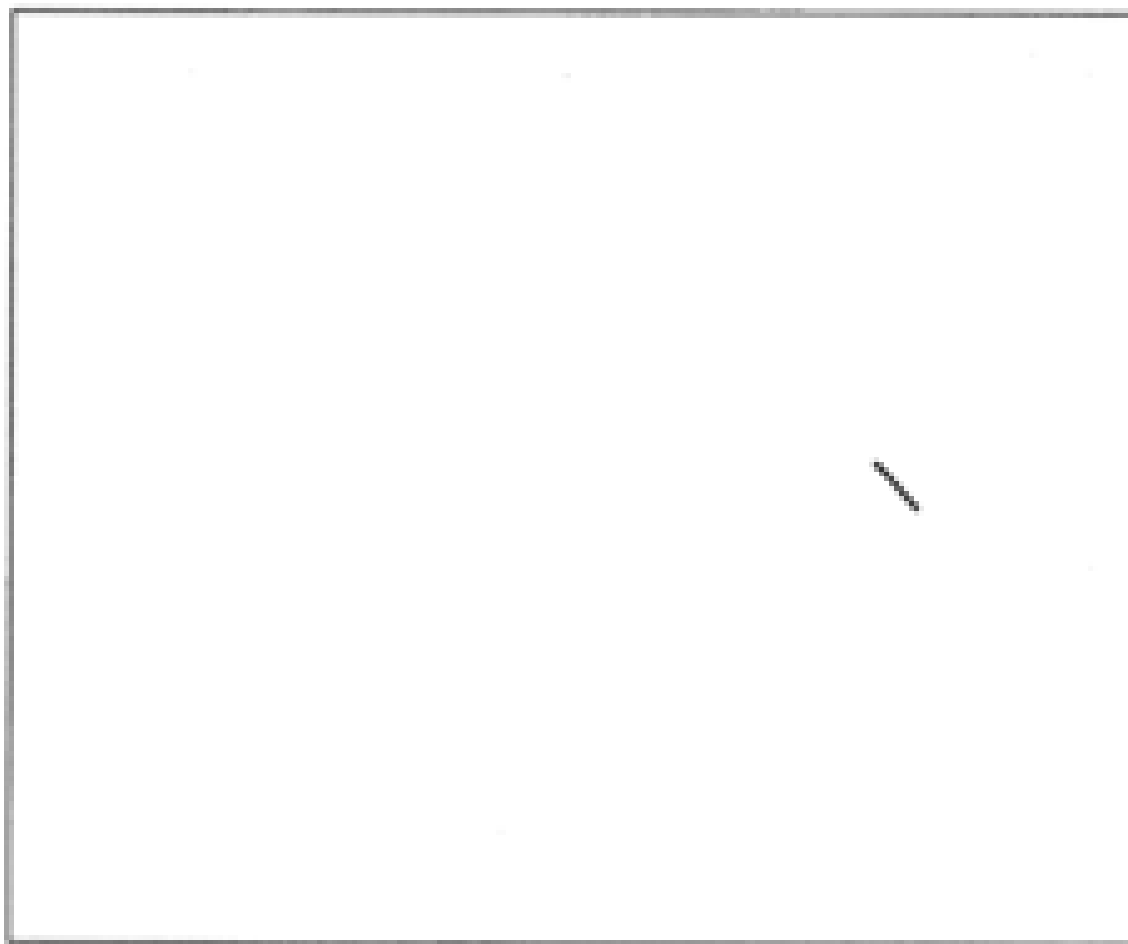
- Many problems are not suitable to be solved by sequential machines and algorithms.

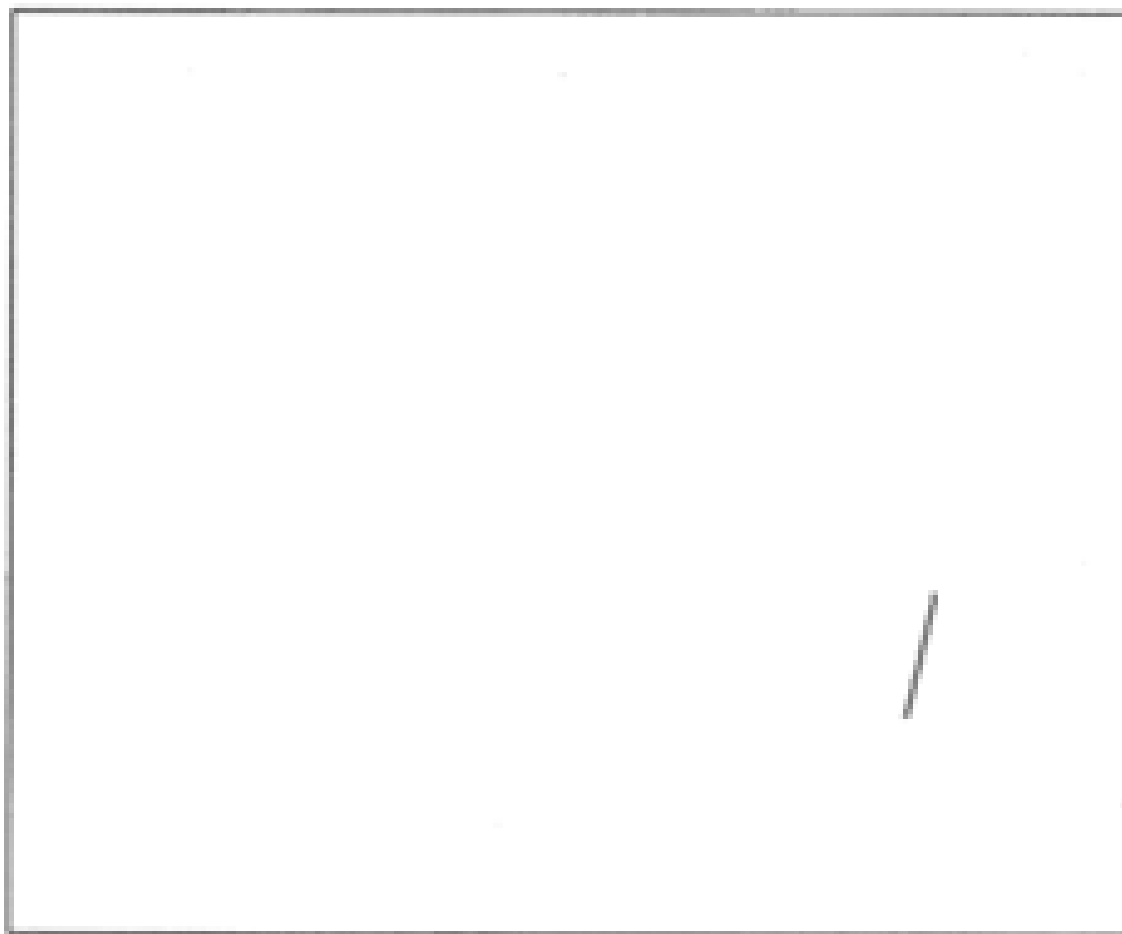
Parallel-processing architectures may be good tools for solving difficult-to-solve, or unsolved problems.

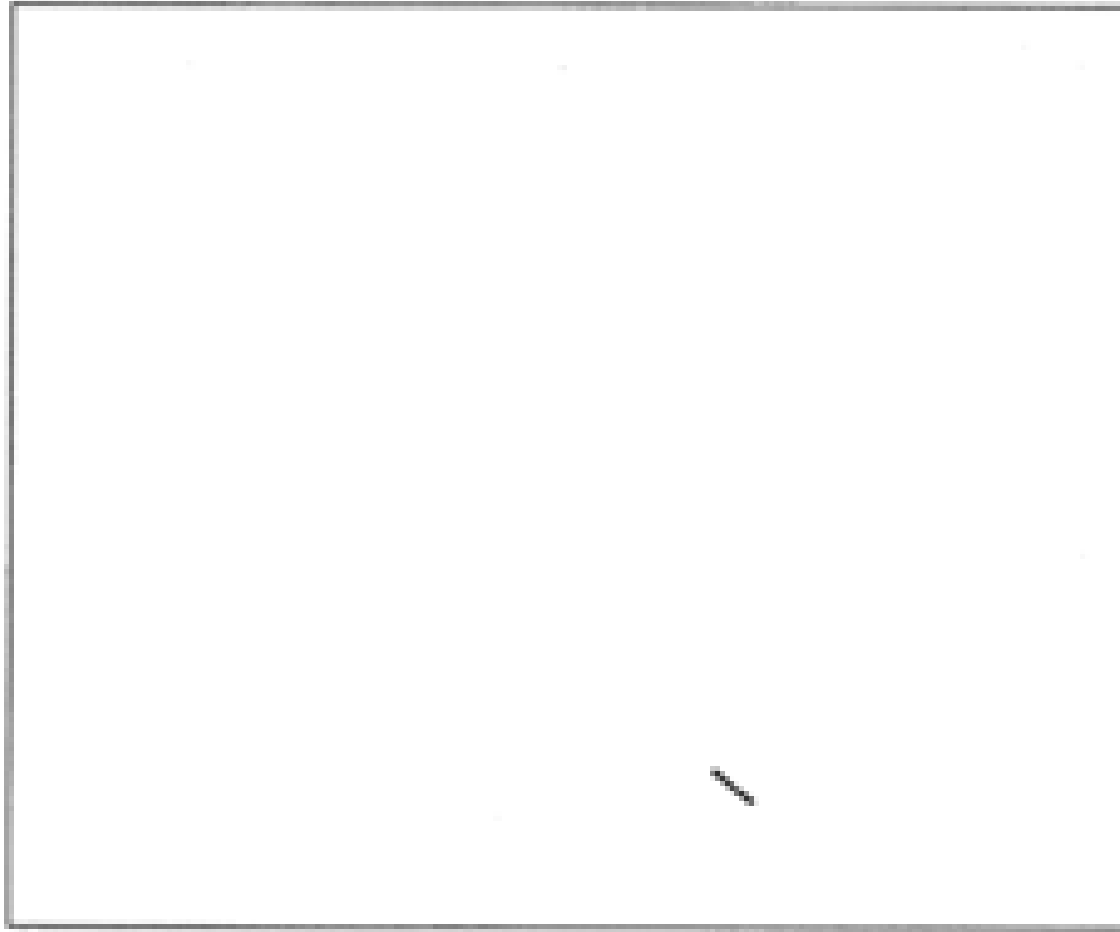
## Example 1: Perceptual-Organization Problem (23 line drawings)



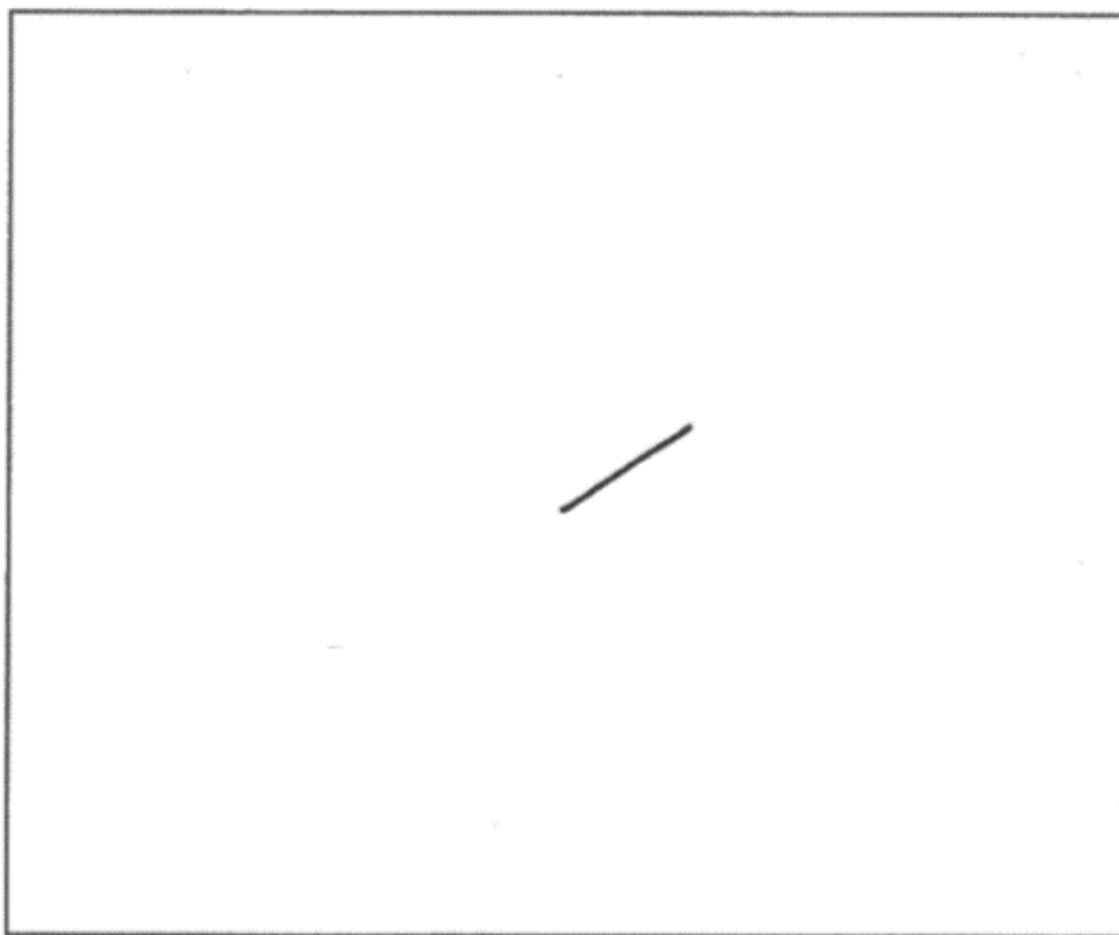


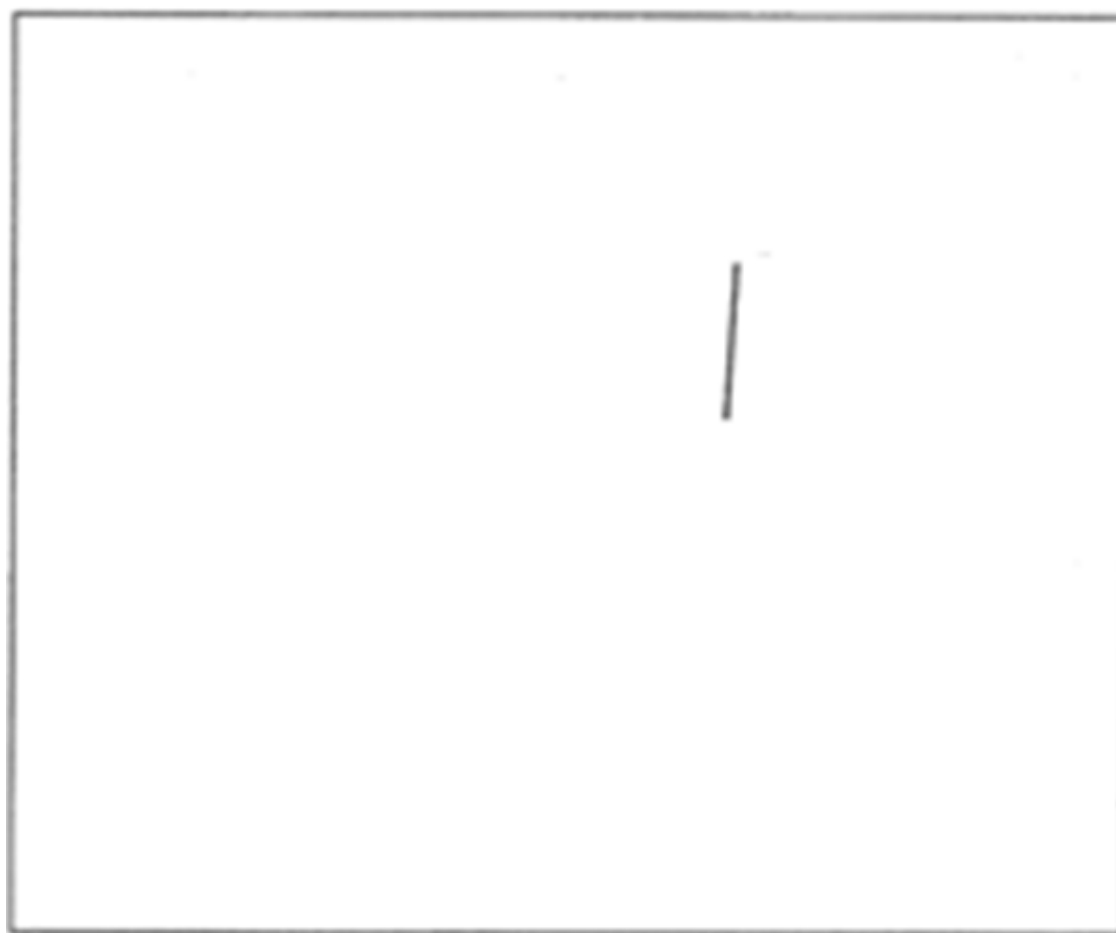


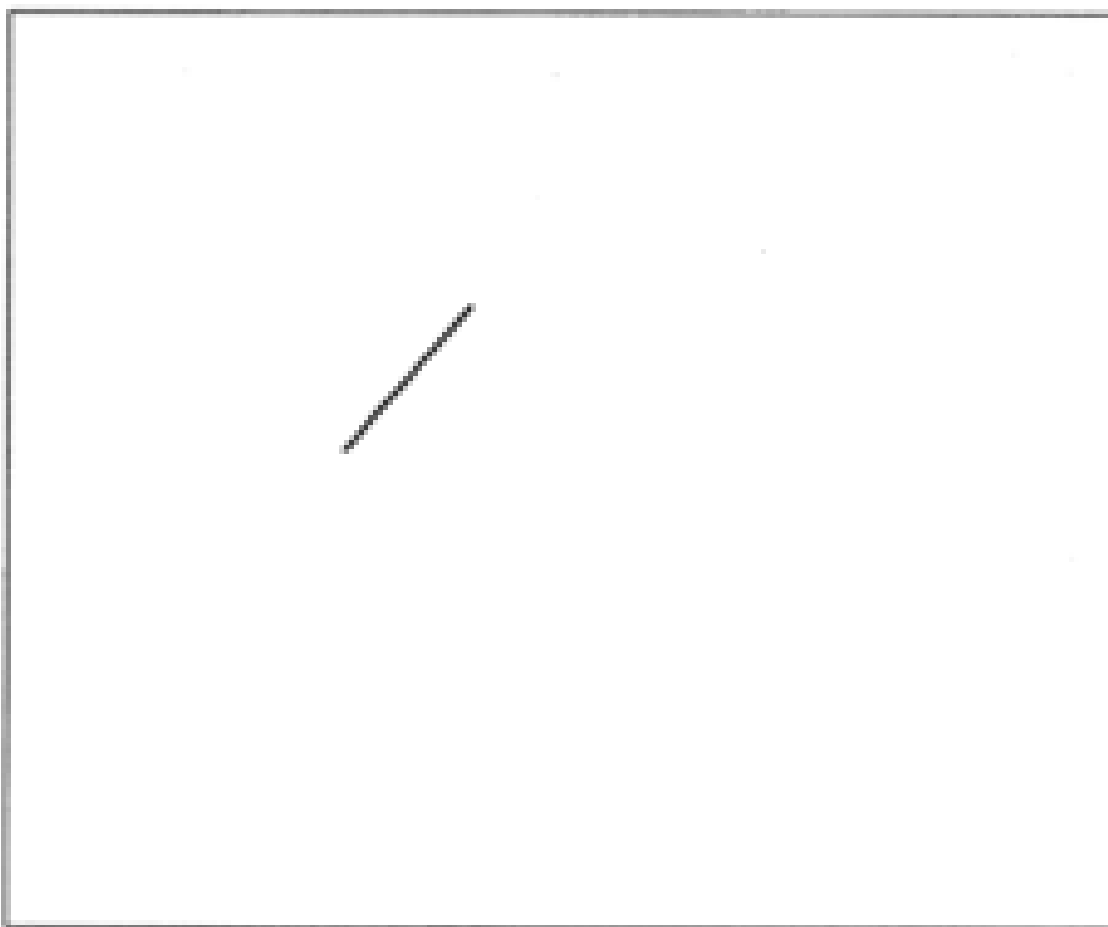


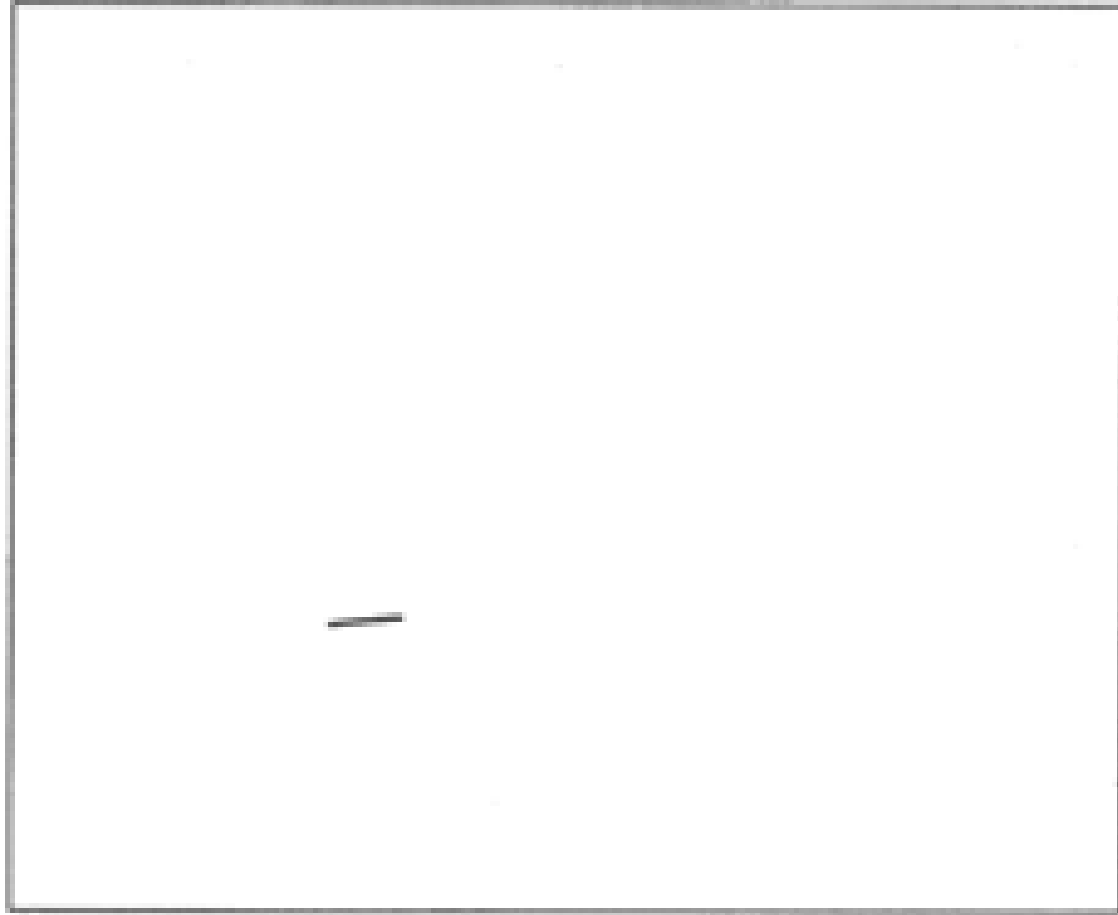








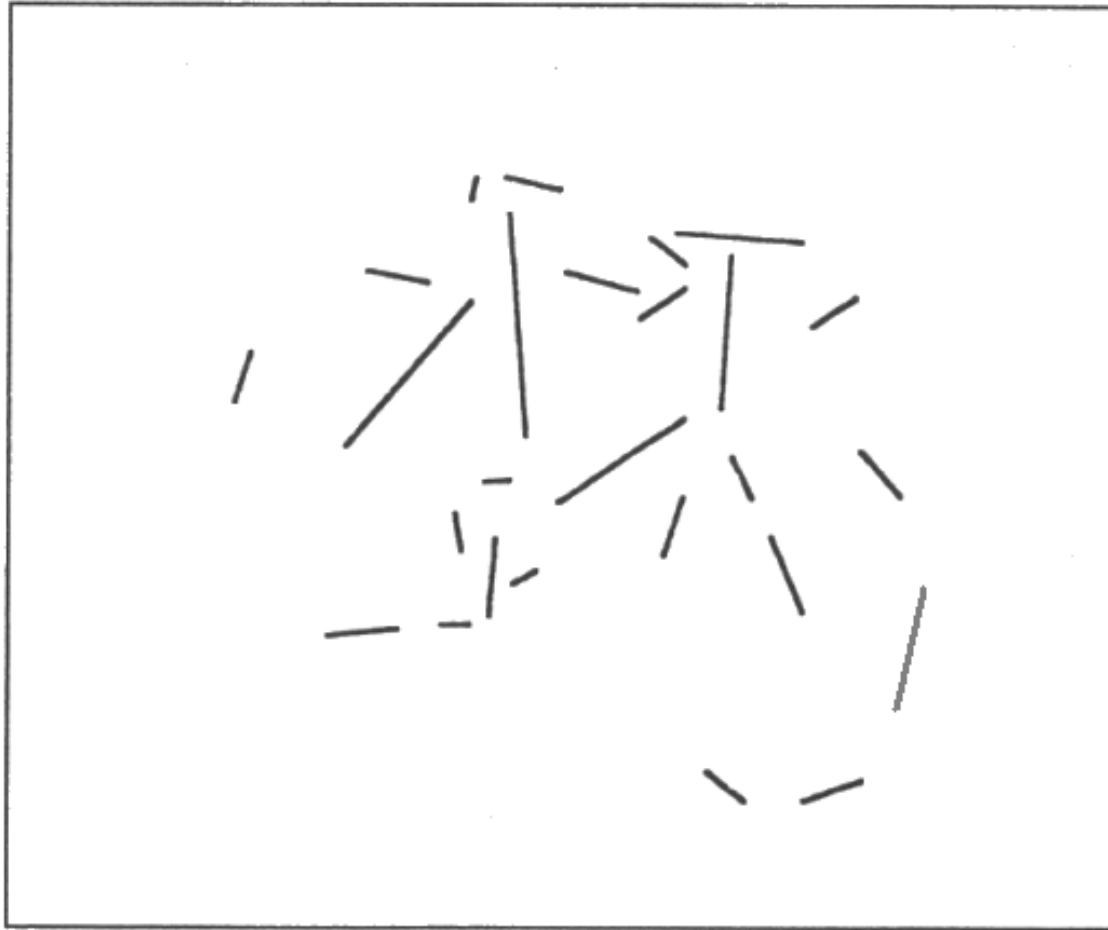




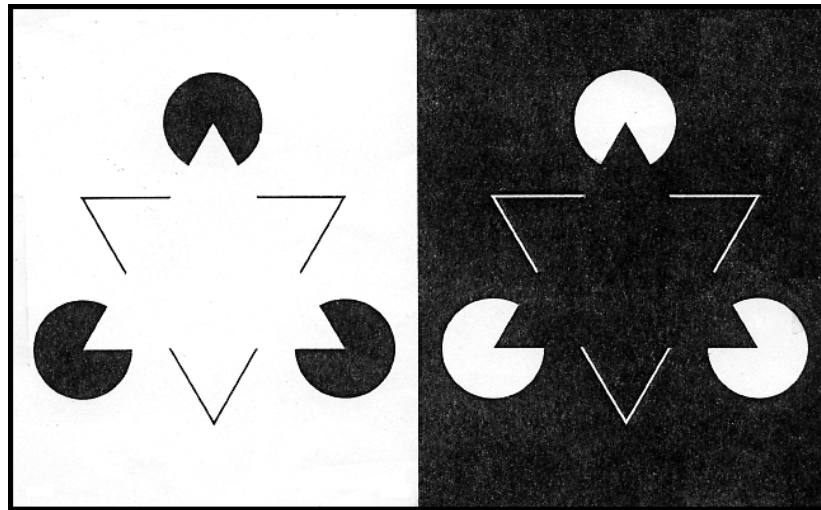




Perceptual organization is well processed in a parallel manner.



Example 2: Subjective inference (reasoning)



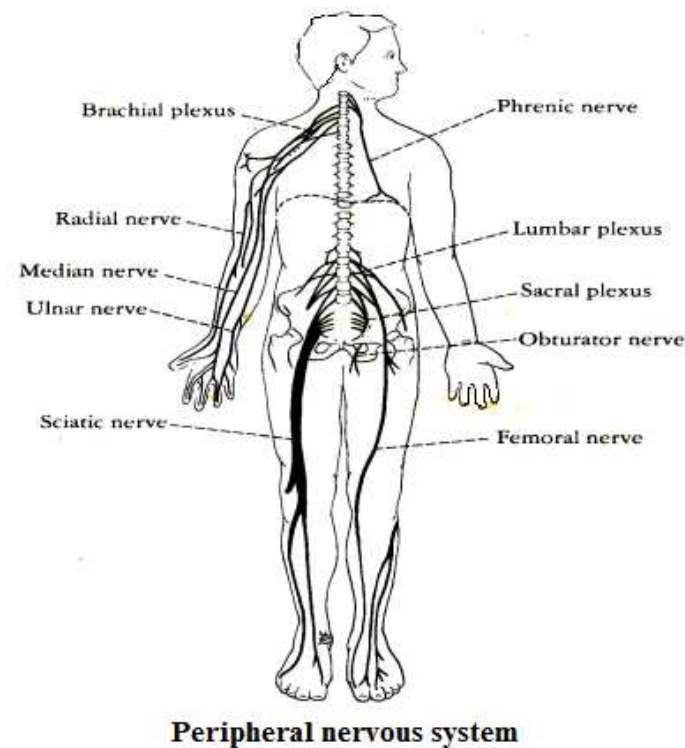
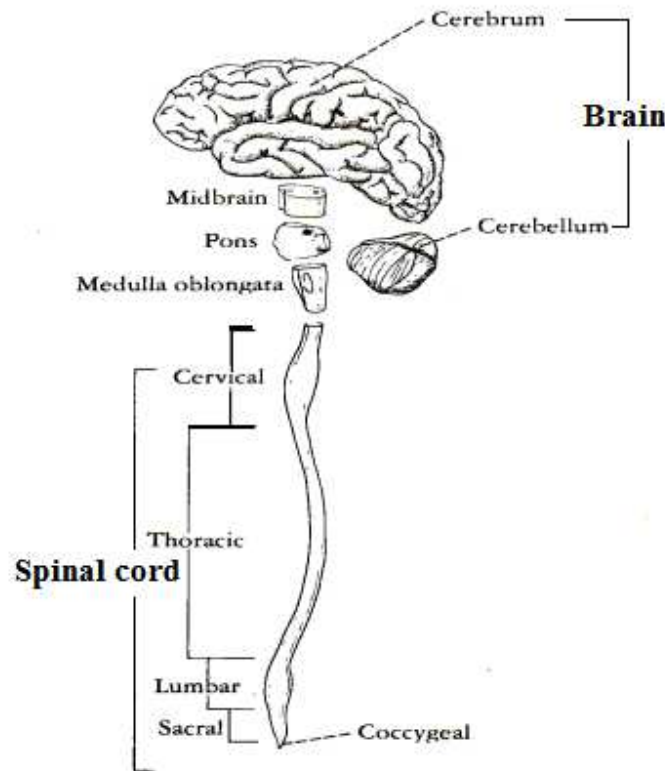


### Example 3: Visual Pattern Recognition

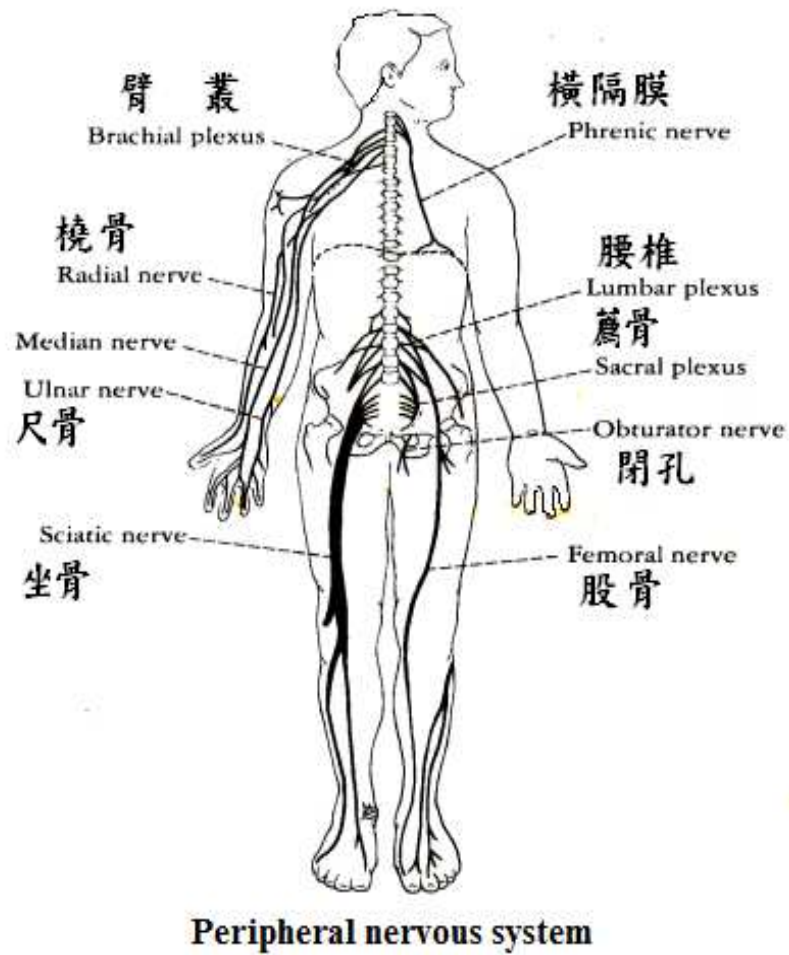


## 11.2 Neurophysiology

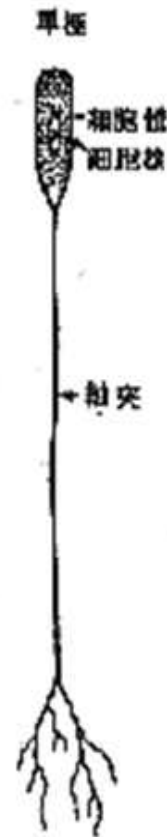
Three major components constructing the human nervous system: **brain**, **spinal cord**, **periphery**



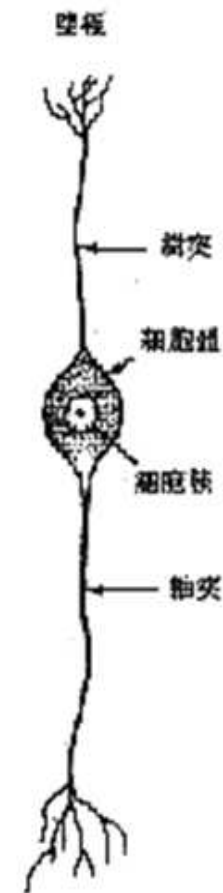
- Periphery



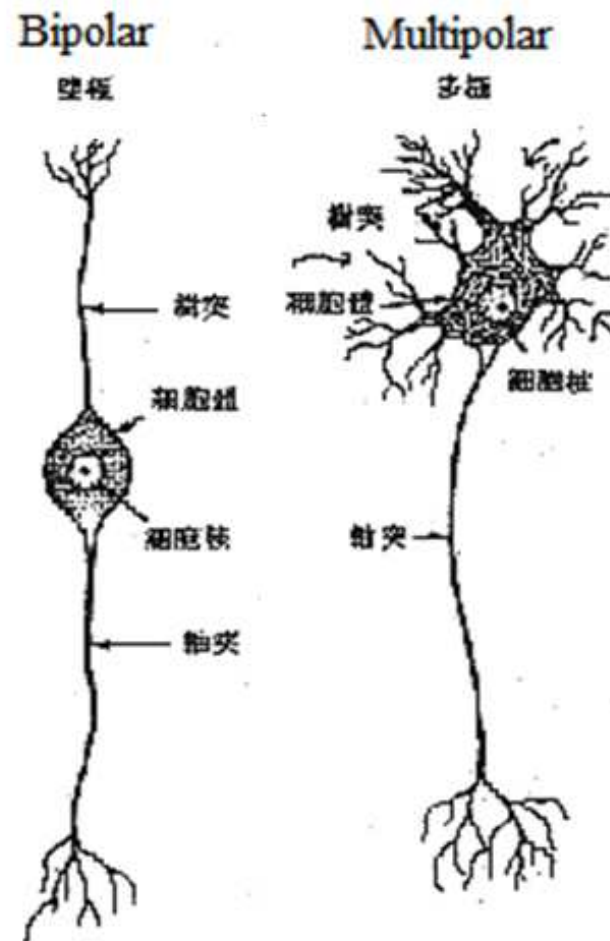
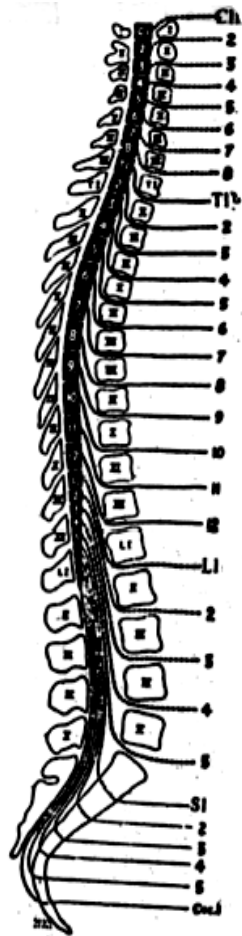
Unipolar



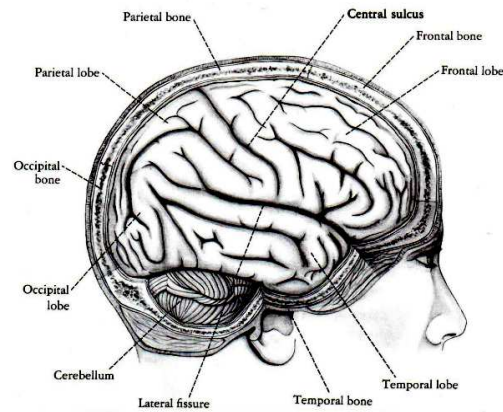
Bipolar



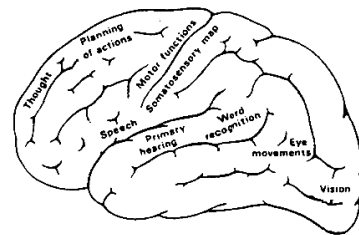
- Spinal Nerves



- Brain



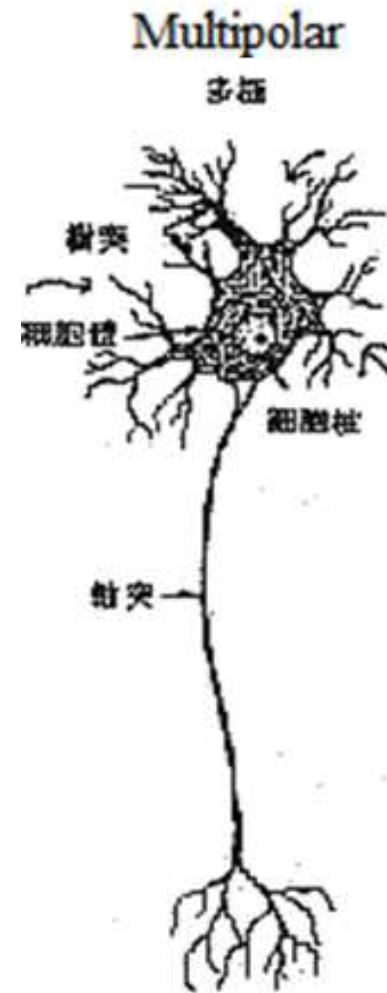
## Cerebral cortex



Size:  $1 \text{ m}^2$

Thick: 2-4 mm

Layers: 6



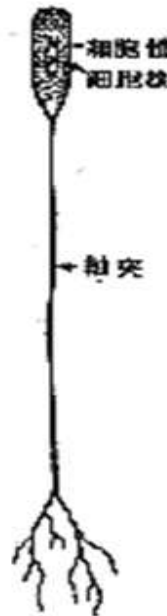
## ◎ Single-Neurons Physiology

Three types of neurons:

1. unipolar
2. bipolar
3. multipolar

Unipolar

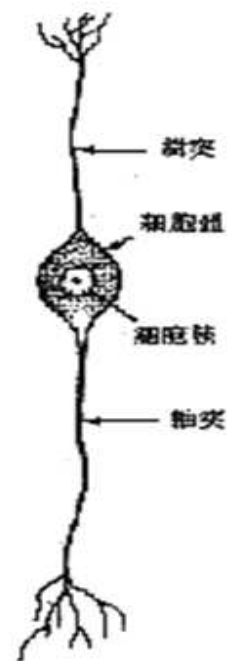
單極



Terminal  
末梢神經  
感覺器官

Bipolar

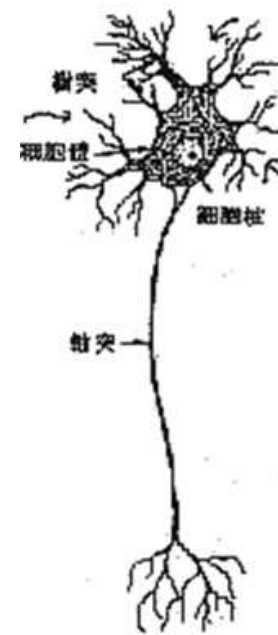
雙極



Pathway  
連絡神經  
脊髓

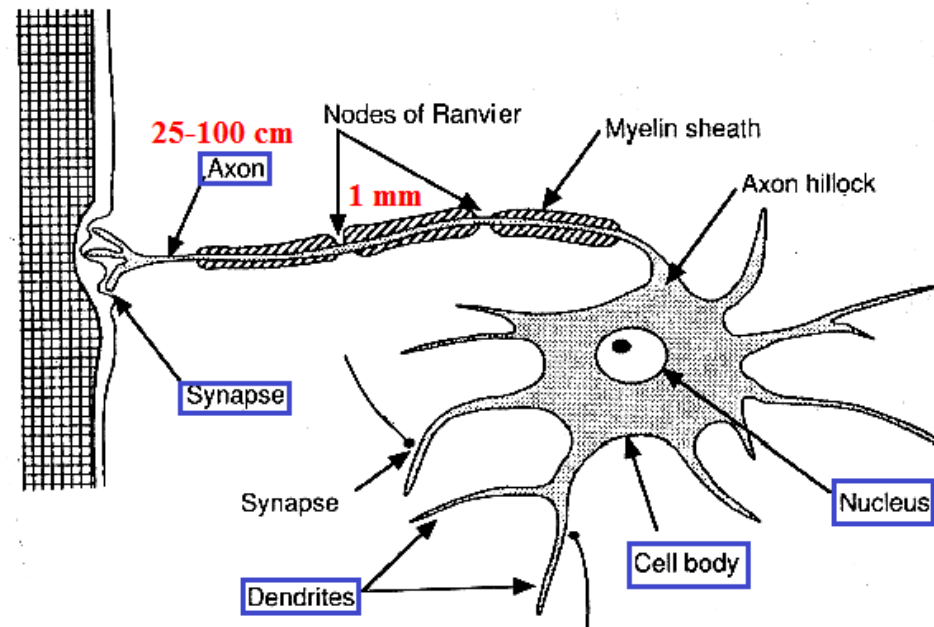
Multipolar

多極



Central  
中樞神經  
腦

## ◎ Multi-Polar Neuron

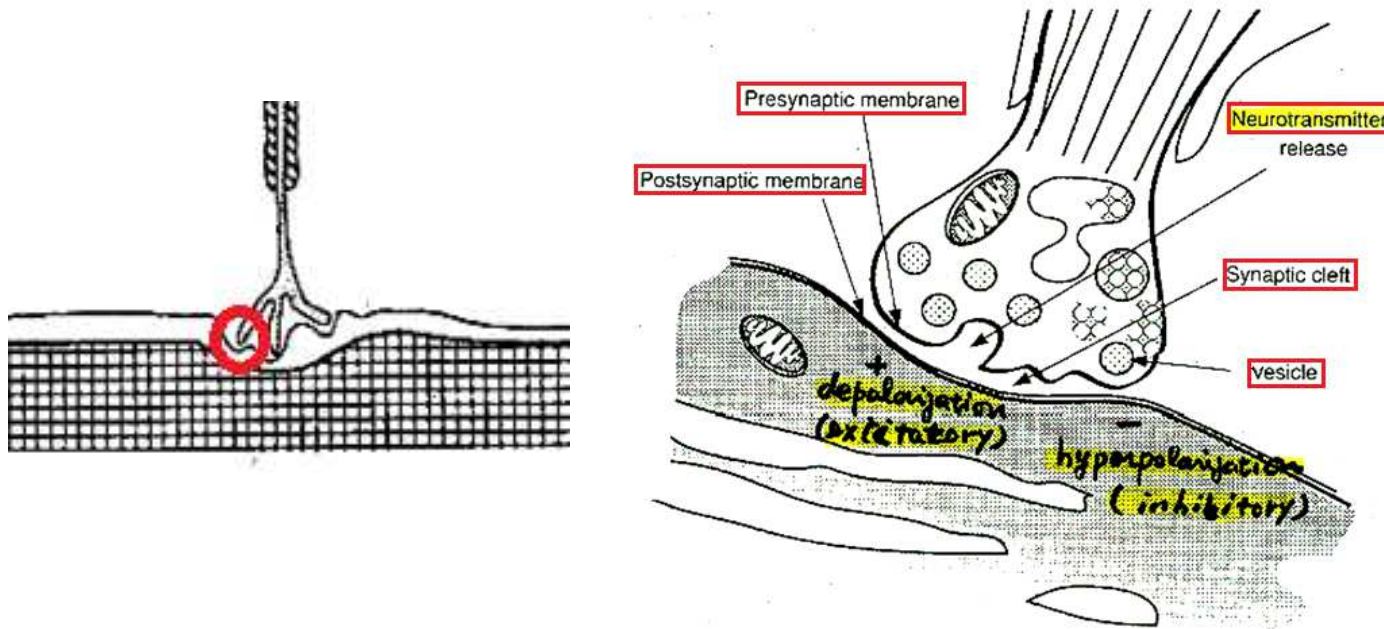


The input impulses can be **excitatory** or **inhibitory** and are summed at the **axon hillock**. If the summed potential is sufficient, an **action potential** is generated.



## ◎ Synaptic Function

Neurotransmitters held in vesicles and are released near presynaptic membrane into synaptic cleft and absorbed by postsynaptic membrane.





**Signal** : frequency of pulses

**Learning**: adjusting synaptic gaps

**Memory** : strength of synaptic connections

Knowledge is acquired through a learning process.

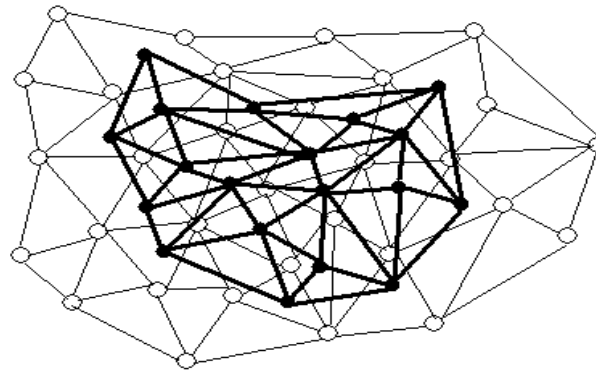
Acquired knowledge is stored in interneuron links in terms of strengths.

© Comparison between nervous and computer systems

	Neurons	Processing units
Switch time	10e-3 sec	10e-9 sec
Synapses	10e+3	10
Number	10e+11	~

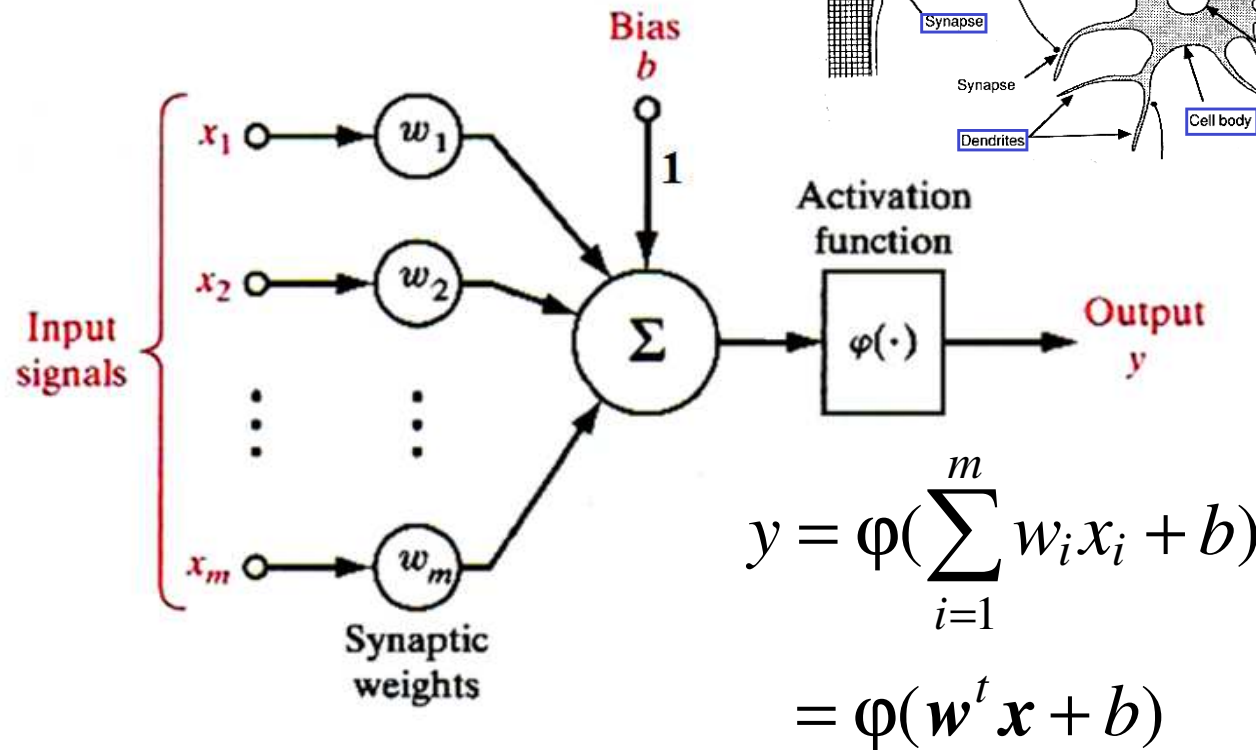
## 11.3 Artificial Neural Networks

- A collection of parallel processing units (neurons) interconnected together
- Various neural systems with different characteristics and functions result from
  - i) different functions of neurons
  - ii) different ways of connections
  - iii) different flows of information



- Characteristics: nonlinearity, non-locality, non-algorithm, dynamics, adaptivity, fault-tolerance, input-output mapping, evidential response, self-organization
- Functions: learn, analyze, associate, organize, comprehend, recognize, plan, decide

## ◎ Neuron Model

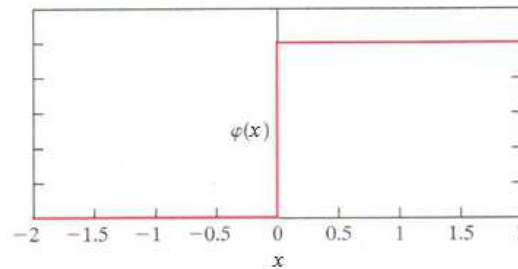


**Bias term  $b$ :** feedback, error, gain, adjustment

Types of activation function  $\varphi(\cdot)$ :

1. **Threshold function** (Heaviside function)

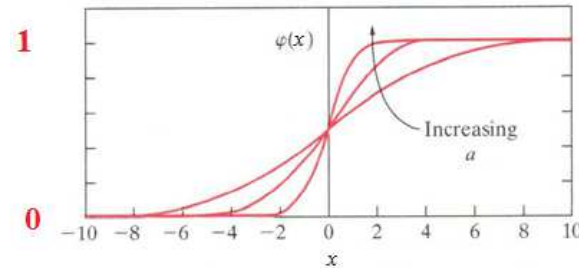
$$\varphi(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



2. **Sigmoid function**

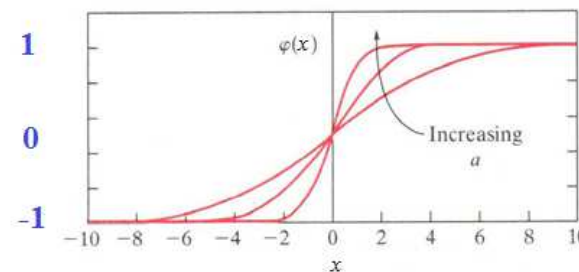
Logistic function

$$\varphi(x) = \frac{1}{1 + \exp(-ax)}$$



Hyperbolic function

$$\varphi(x) = \tanh(ax)$$



### 3. Softsign function

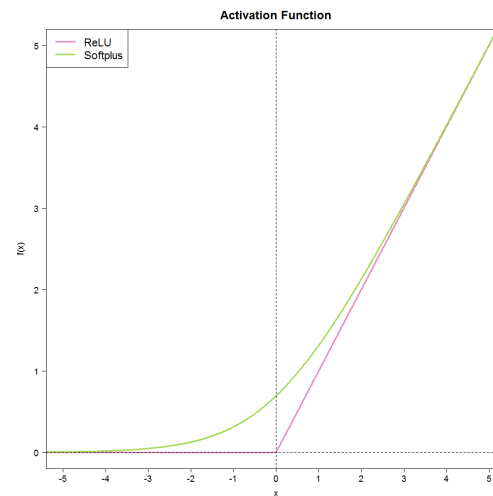
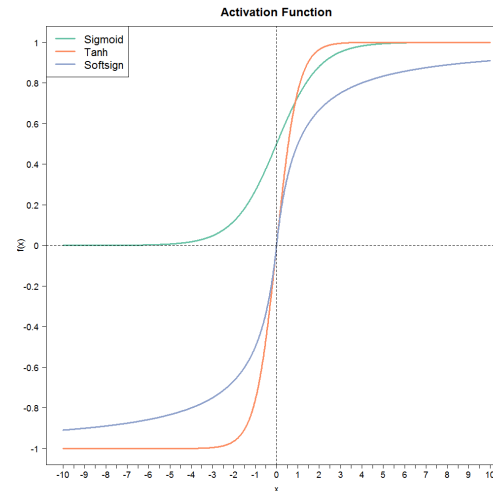
$$\varphi(x) = \frac{x}{1 + |x|}$$

### 4. Softplus function

$$\varphi(x) = \ln(1 + e^x)$$

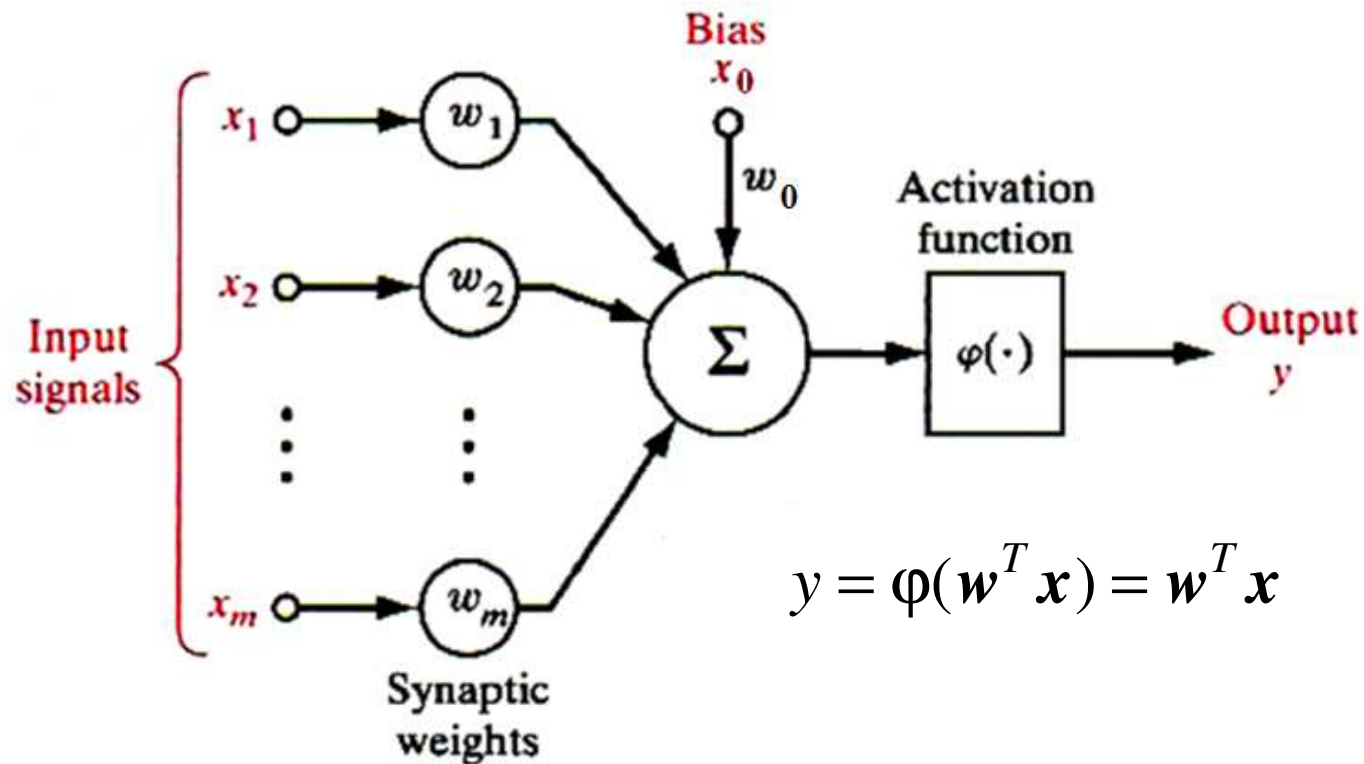
### 5. ReLU function

$$\varphi(x) = \max(0, x)$$



**Perceptron:** Neuron model with linear active function

$$\varphi(x) = x \quad \therefore y = \varphi(\mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \mathbf{x}$$



## 11.4 Training a Perceptron

### 11.4.1 Least Mean Square (LMS) Learning

⊙ Input vectors :  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L\}$

Ideal outputs :  $\{d_1, d_2, \dots, d_L\}$

Actual outputs :  $\{y_1, y_2, \dots, y_L\}$

Mean square error:

$$\begin{aligned}\langle \varepsilon_k^2 \rangle &= \frac{1}{L} \sum_{k=1}^L \varepsilon_k^2 = \langle (d_k - y_k)^2 \rangle = \langle (d_k - \mathbf{w}^T \mathbf{x}_k)^2 \rangle \\ &= \langle d_k^2 \rangle + \mathbf{w}^T \langle \mathbf{x}_k \mathbf{x}_k^T \rangle \mathbf{w} - 2 \langle d_k \mathbf{x}_k^T \rangle \mathbf{w} \quad \text{-- (2.4)}\end{aligned}$$

Let  $\xi = \langle \varepsilon_k^2 \rangle$ ,  $\mathbf{p} = \langle d_k \mathbf{x}_k \rangle$ ,  $\mathbf{R} = \langle \mathbf{x}_k \mathbf{x}_k^T \rangle$ : correlation matrix



$$(2.4) \Rightarrow \xi = \langle d_k^2 \rangle + \mathbf{w}^T R \mathbf{w} - 2 \mathbf{p}^T \mathbf{w}$$

Idea:  $\mathbf{w}^* = \arg \min_{\mathbf{w}} \xi(\mathbf{w})$

Let  $\frac{d\xi(\mathbf{w})}{d\mathbf{w}} = 2R\mathbf{w} - 2\mathbf{p} = 0$ . Obtain  $\mathbf{w}^* = R^{-1}\mathbf{p}$ .

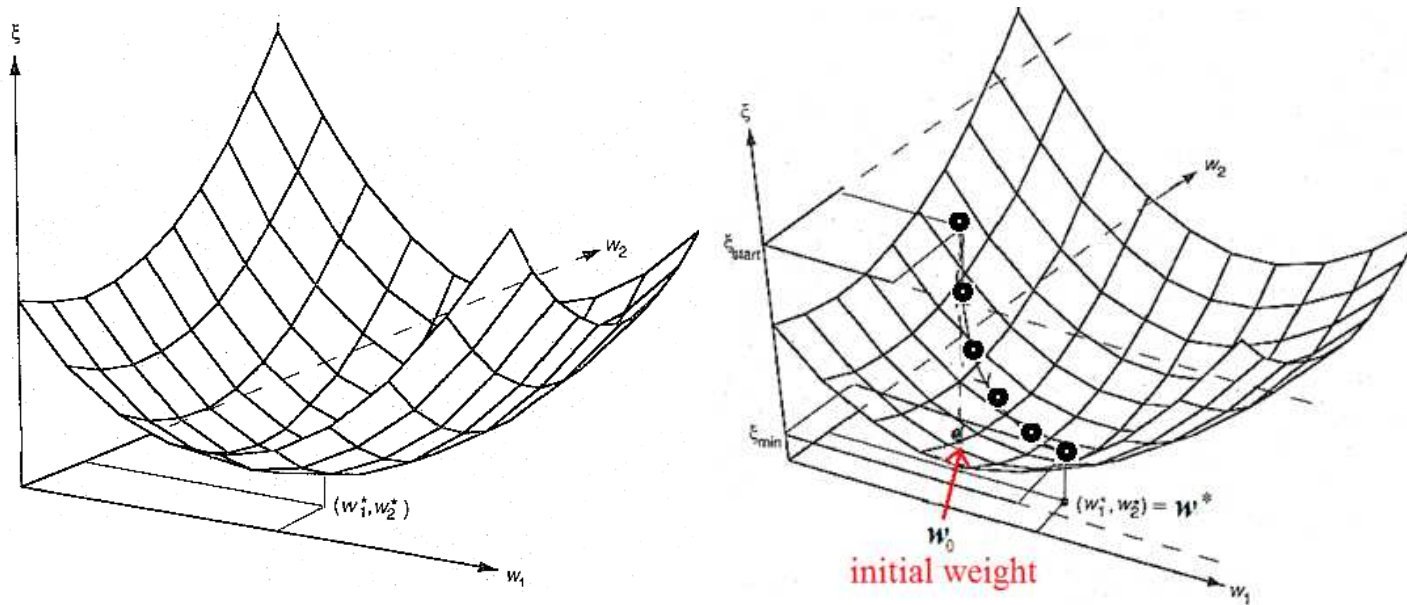
Practical difficulties of analytical formula :

1. Large dimensions –  $R^{-1}$  difficult to calculate
2. Need to calculate the expected values of

$\mathbf{p} = \langle d_k \mathbf{x}_k \rangle$ ,  $R = \langle \mathbf{x}_k \mathbf{x}_k^T \rangle$  – lack of the knowledge of probabilities

## 11.4.2 Gradient Descent (GD) Learning

The graph of  $\xi(\mathbf{w}) = \langle d_k^2 \rangle + \mathbf{w}^T R \mathbf{w} - 2 \mathbf{p}^T \mathbf{w}$  is a paraboloid.



- Steps:**
1. Initialize weight values  $\mathbf{w}(t_0)$
  2. Determine the steepest descent direction

$$-\nabla_{\mathbf{w}} \xi(\mathbf{w}(t)) = -\frac{d\xi(\mathbf{w}(t))}{d\mathbf{w}(t)} = 2(\mathbf{p} - R\mathbf{w}(t))$$

$$\text{Let } \Delta\mathbf{w}(t) = -\nabla_{\mathbf{w}} \xi(\mathbf{w}(t)) = 2(\mathbf{p} - R\mathbf{w}(t))$$

3. Modify weight values

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \mu\Delta\mathbf{w}(t), \quad \mu: \text{step size}$$

4. Repeat 2~3.

◆ No calculation of  $R^{-1}$

Drawbacks:

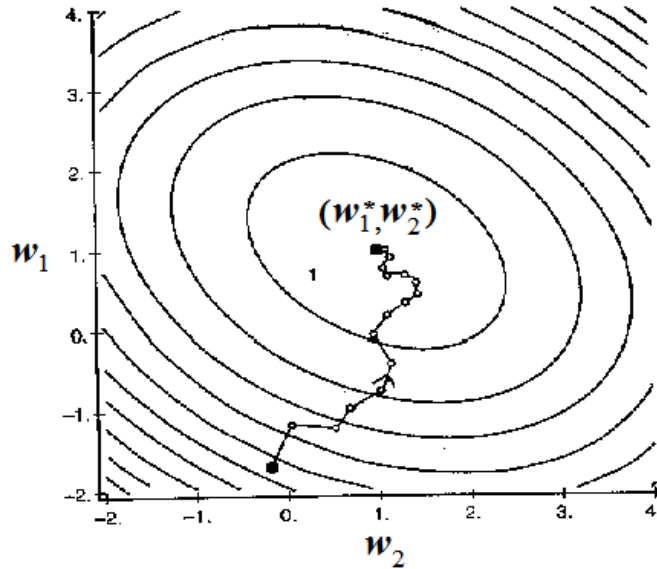
- i) Need to calculate  $\mathbf{p}$  and  $R$ ,
- ii) Steepest descent is a batch training method.

### 11.4.3 Stochastic Gradient Descent (SGD) Learning

Approximate  $-\nabla_{\mathbf{w}} \xi(\mathbf{w}(t)) = 2(\mathbf{p} - R\mathbf{w}(t))$  by randomly selecting one training example at a time

1. Apply an input vector  $\mathbf{x}_k$
2.  $\varepsilon_k^2(t) = (d_k - y_k)^2 = (d_k - \mathbf{w}^T(t) \cdot \mathbf{x}_k)^2$
3.  $-\nabla_{\mathbf{w}} \xi(\mathbf{w}(t)) = -\nabla_{\mathbf{w}} \langle \varepsilon_k^2(t) \rangle \approx -\nabla_{\mathbf{w}} \varepsilon_k^2(t)$   
$$= 2(d_k - \mathbf{w}^T(t) \cdot \mathbf{x}_k) \mathbf{x}_k = 2\varepsilon_k(t) \mathbf{x}_k$$
4.  $\mathbf{w}(t+1) = \mathbf{w}(t) + \mu \Delta \mathbf{w}(t) = \mathbf{w}(t) + 2\mu \varepsilon_k(t) \mathbf{x}_k$
5. Repeat 1~4 with the next input vector

◆ No calculation of  $\mathbf{p}$  and  $R$



Drawback: time consuming.

Improvement: mini-batch training method.

- Practical Considerations:
  - (a) No. of training vectors, (b) Stopping criteria
  - (c) Initial weights, (d) Step size

### 11.4.4 Conjugate Gradient Descent (CGD) Learning

-- **Drawback**: can only minimize quadratic functions,

$$\text{e.g., } f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T A \mathbf{w} - \mathbf{b}^T \mathbf{w} + c \quad A = 2R$$

- Adequate for our error function  $\mathbf{b}^T = 2\mathbf{p}^T$

$$\xi(\mathbf{w}) = \mathbf{w}^T R \mathbf{w} - 2\mathbf{p}^T \mathbf{w} + \langle d_k^2 \rangle \quad c = \langle d_k^2 \rangle$$

**Advantage**: guarantees to find the optimum solution  
in at most  $n$  iterations, where  $n$  is the  
size of matrix  $A$ .

- The size of correlation matrix  $R$  is the dimension  
of input vectors  $\mathbf{x}$ .

## A-Conjugate Vectors:

Let  $A_{n \times n}$  : square, symmetric, positive-definite matrix.

$\{s(0), s(1), \dots, s(n-1)\}$  are  $A$ -conjugate vectors if

$$s^T(i)As(j) = 0, \quad \forall i \neq j$$

\* If  $A = I$  (identity matrix), conjugacy = **orthogonality**.

• CGD finding the  $\mathbf{w}$  to minimize  $f(\mathbf{w})$  is through

$$\mathbf{w}(i+1) = \mathbf{w}(i) + \eta(i)s(i), \quad \text{where } i = 0, \dots, n-1,$$

$$(\text{GD: } \mathbf{w}(i+1) = \mathbf{w}(i) - \eta(i) \nabla_{\mathbf{w}} f(\mathbf{w}(i)) )$$

$$\eta(i) \text{ is determined by } \min_{\eta} f(\mathbf{w}(i) + \eta s(i)).$$

$$\left( \begin{aligned} \because f(\mathbf{w}) &= \frac{1}{2} \mathbf{w}^T A \mathbf{w} - \mathbf{b}^T \mathbf{w} + c \\ f(\mathbf{w}(i) + \eta \mathbf{s}(i)) &= \frac{1}{2} (\mathbf{w}(i) + \eta \mathbf{s}(i))^T A (\mathbf{w}(i) + \eta \mathbf{s}(i)) \\ &\quad - \mathbf{b}^T (\mathbf{w}(i) + \eta \mathbf{s}(i)) + c \end{aligned} \right.$$

$$\text{Let } \frac{d}{d\eta} [f(\mathbf{w}(i) + \eta \mathbf{s}(i))] = 0$$

$$\begin{aligned} \frac{d}{d\eta} [f(\mathbf{w}(i) + \eta \mathbf{s}(i))] &= \frac{d}{d\eta} \left[ \mathbf{w}(i)^T A \mathbf{w}(i) + \eta \mathbf{s}(i)^T A \mathbf{w}(i) \right. \\ &\quad \left. + \eta \mathbf{w}(i)^T A \mathbf{s}(i) + \eta^2 \mathbf{s}(i)^T A \mathbf{s}(i) - \mathbf{b}^T \mathbf{w}(i) - \eta \mathbf{b}^T \mathbf{s}(i) \right] \\ &= \mathbf{s}(i)^T A \mathbf{w}(i) + \mathbf{w}(i)^T A \mathbf{s}(i) + 2\eta \mathbf{s}(i)^T A \mathbf{s}(i) - \mathbf{b}^T \mathbf{s}(i) = 0 \end{aligned}$$

$$\eta = \frac{\mathbf{b}^T \mathbf{s}(i) - (\mathbf{s}(i)^T A \mathbf{w}(i) + \mathbf{w}(i)^T A \mathbf{s}(i))}{2\mathbf{s}(i)^T A \mathbf{s}(i)} \quad \text{---- (A)}$$



How to determine  $s(i)$ ?

$$f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T A \mathbf{w} - \mathbf{b}^T \mathbf{w} + c \Rightarrow -\nabla_{\mathbf{w}} f(\mathbf{w}) = \mathbf{b} - A \mathbf{w}$$

$$\text{Define } \mathbf{r}(i) = -\nabla_{\mathbf{w}} f(\mathbf{w}) = \mathbf{b} - A \mathbf{w}(i) \text{ --- (B)}$$

Let  $s(0) = \mathbf{r}(0) = \mathbf{b} - A \mathbf{w}(0)$ ,  $\mathbf{w}(0)$ : arbitrary vector

$$s(i) = \mathbf{r}(i) + \alpha(i)s(i-1), \quad i = 1, \dots, n-1 \text{ --- (C)}$$

To determine  $\alpha(i)$ , multiply (C) by  $s(i-1)^T A$ ,

$$s^T(i-1) A s(i) = s^T(i-1) A (\mathbf{r}(i) + \alpha(i)s(i-1)) \text{ --- (D)}$$

In order to be  $A$ -conjugate:  $s^T(i) A s(j) = 0, \quad \forall i \neq j$

$$(D) \Rightarrow 0 = s^T(i-1) A \mathbf{r}(i) + \alpha(i)s^T(i-1) A s(i-1).$$

$$\alpha(i) = -\frac{s^T(i-1)Ar(i)}{s^T(i-1)As(i-1)} \quad \text{--- (D)}$$

**Consider the error function:**

$$\xi(\mathbf{w}) = \mathbf{w}^T R \mathbf{w} - 2 \mathbf{p}^T \mathbf{w} + \langle d_k^2 \rangle$$

Comparing with  $f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T A \mathbf{w} - \mathbf{b}^T \mathbf{w} + c,$

$$A = 2R, \mathbf{b}^T = 2\mathbf{p}^T, c = \langle d_k^2 \rangle$$

From (A)  $\eta(i) = \frac{\mathbf{b}^T s(i) - (s(i)^T A \mathbf{w}(i) + \mathbf{w}(i)^T A s(i))}{2s(i)^T A s(i)},$

$$\eta(i) = \frac{\mathbf{p}^T s(i) - (s(i)^T R \mathbf{w}(i) + \mathbf{w}(i)^T R s(i))}{2s(i)^T R s(i)}.$$

From (B)  $\mathbf{r}(i) = \mathbf{b} - A\mathbf{w}(i)$ ,

$$\mathbf{r}(i) = 2(\mathbf{p} - R\mathbf{w}(i)).$$

From (C)  $\mathbf{s}(0) = \mathbf{r}(0) = 2(\mathbf{p} - R\mathbf{w}(0))$ ,

$$\mathbf{s}(i) = \mathbf{r}(i) + \alpha(i)\mathbf{s}(i-1).$$

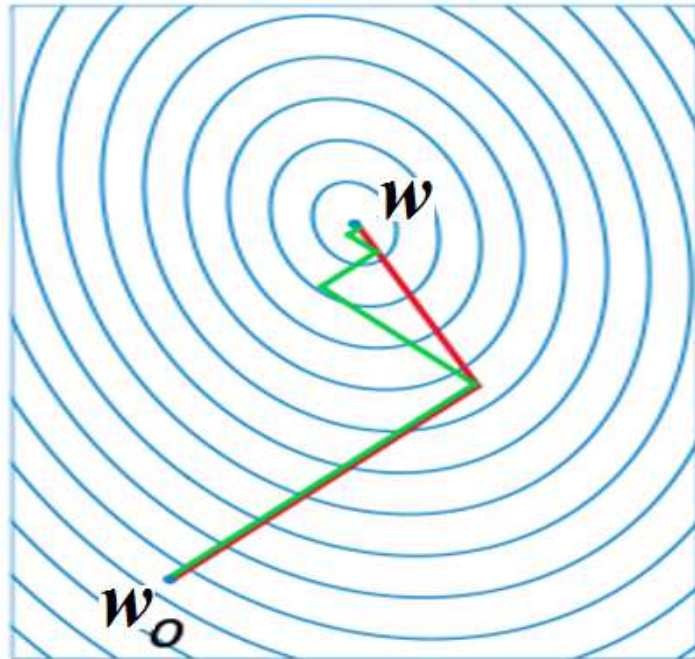
From (D)  $\alpha(i) = -\frac{\mathbf{s}^T(i-1)A\mathbf{r}(i)}{\mathbf{s}^T(i-1)A\mathbf{s}(i-1)},$

$$\alpha(i) = -\frac{\mathbf{s}^T(i-1)R\mathbf{r}(i)}{\mathbf{s}^T(i-1)R\mathbf{s}(i-1)}.$$

Update rule:  $\mathbf{w}(i+1) = \mathbf{w}(i) + \eta(i)\mathbf{s}(i)$ ,  $i = 0, 1, \dots, n-1$

$\mathbf{w}(0)$ : arbitrary starting vector

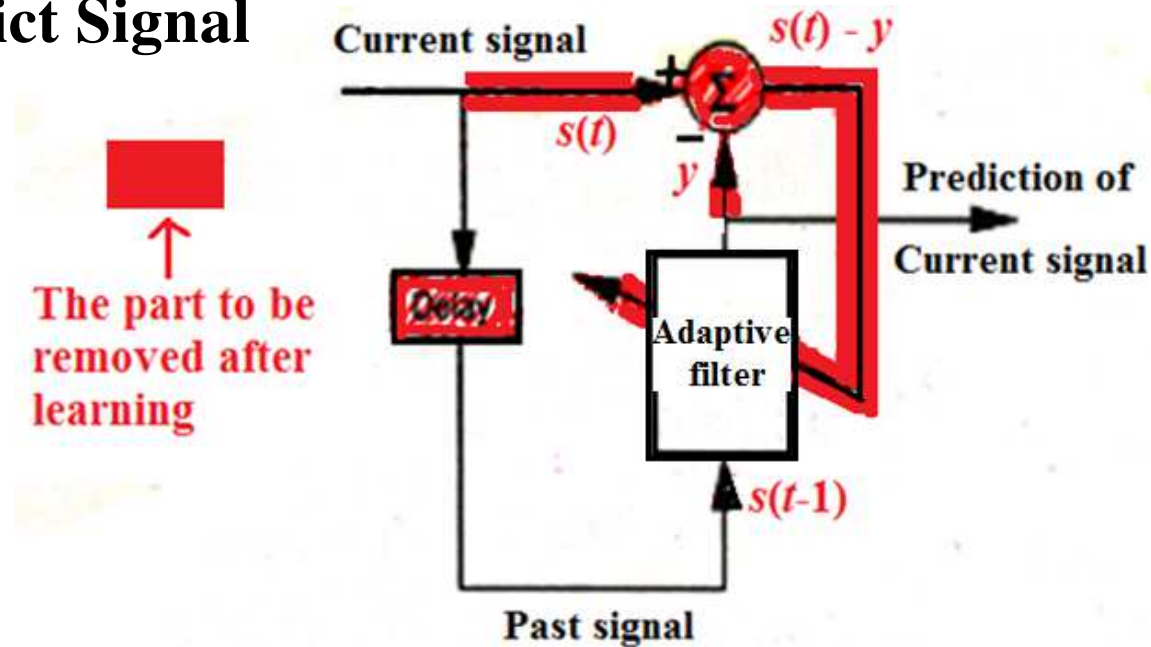
**Example:** A comparison of the convergences of gradient descent (green) and conjugate gradient (red) for minimizing a quadratic function.



Conjugate gradient converges in at most  $n$  steps where  $n$  is the size of the matrix of the system (here  $n=2$ ).

## 11.5. Applications

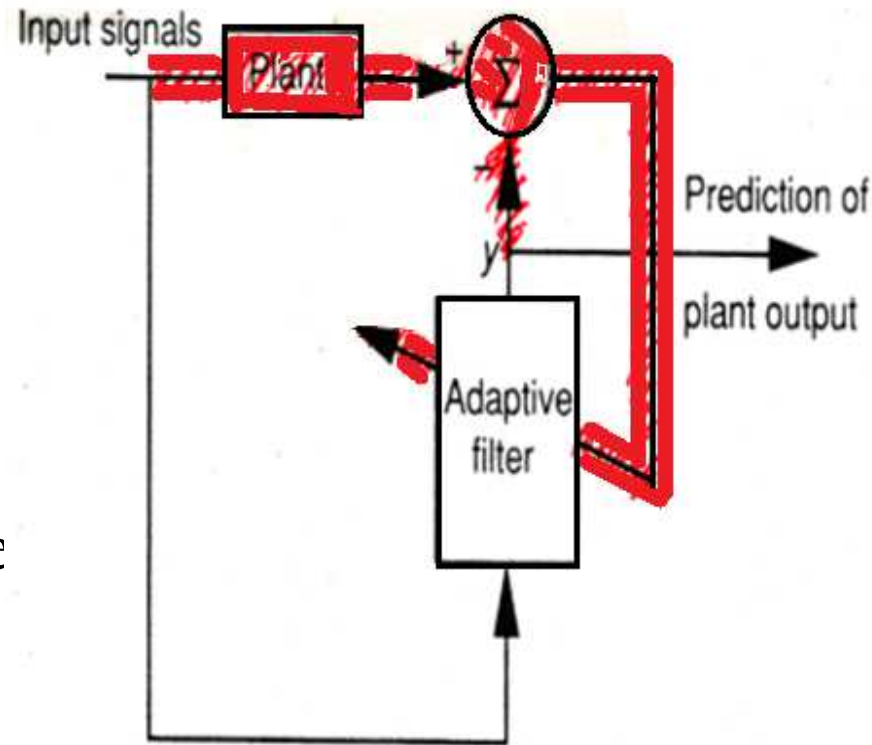
### 11.5.1 Predict Signal



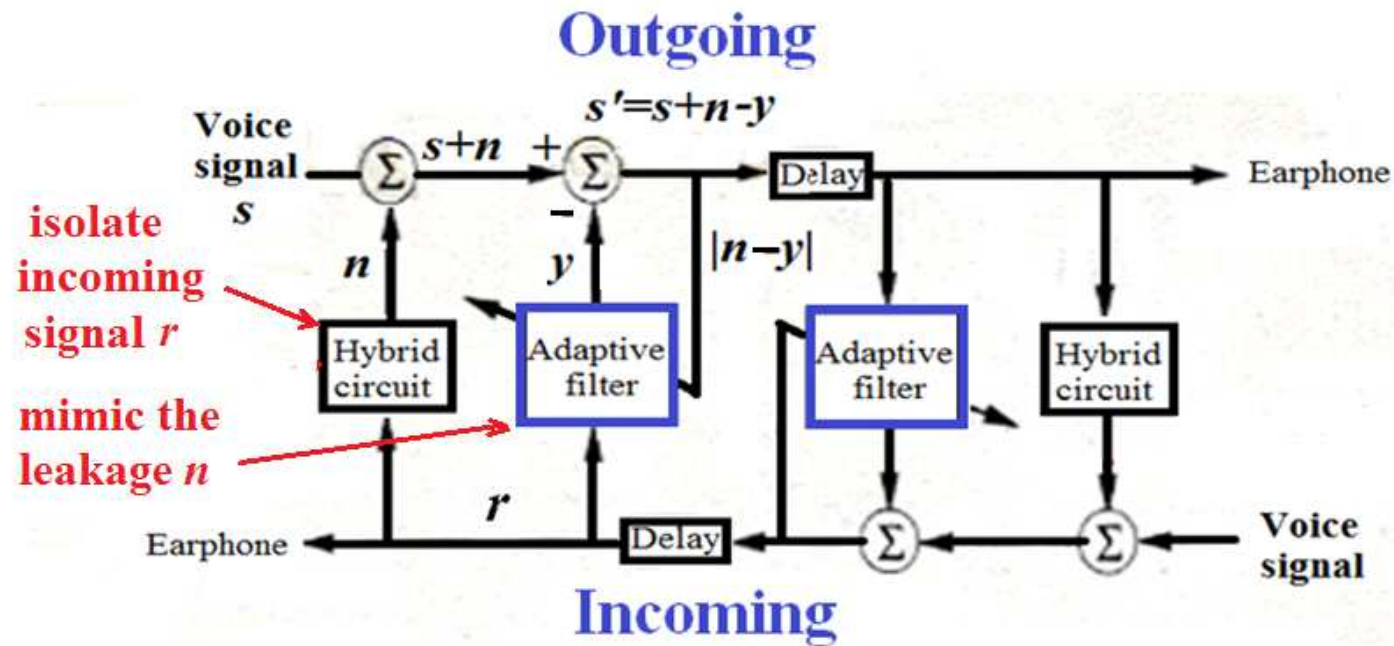
An adaptive filter is trained to predict signal. The Signal used to train the filter is a delayed actual signal. The expected output is the current signal.

## 11.5.2 Reproduce Signal

An adaptive filter is used to model a plant. Inputs to the filter are the same as those to the plant. The filter adjusts its weights based on the difference between its output and the output of the plant.



### 11.5.3. Echo Cancellation in Telephone Circuits



$s$  : outgoing voice,  $r$  : incoming voice

$n$  : noise (leakage of  $r$ )

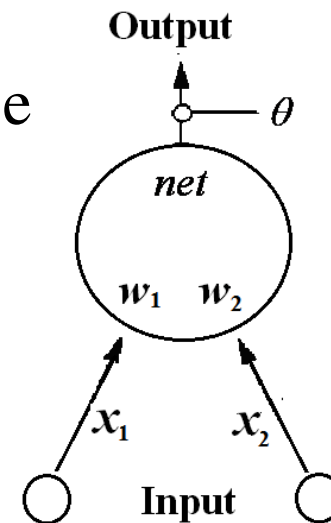
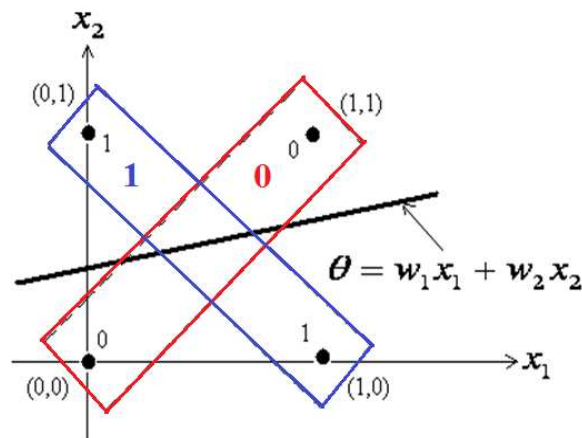
$y$  : the output of the filter mimics  $n$

## 11.5 Multilayer Perceptrons (MLP)

### ○ XOR function

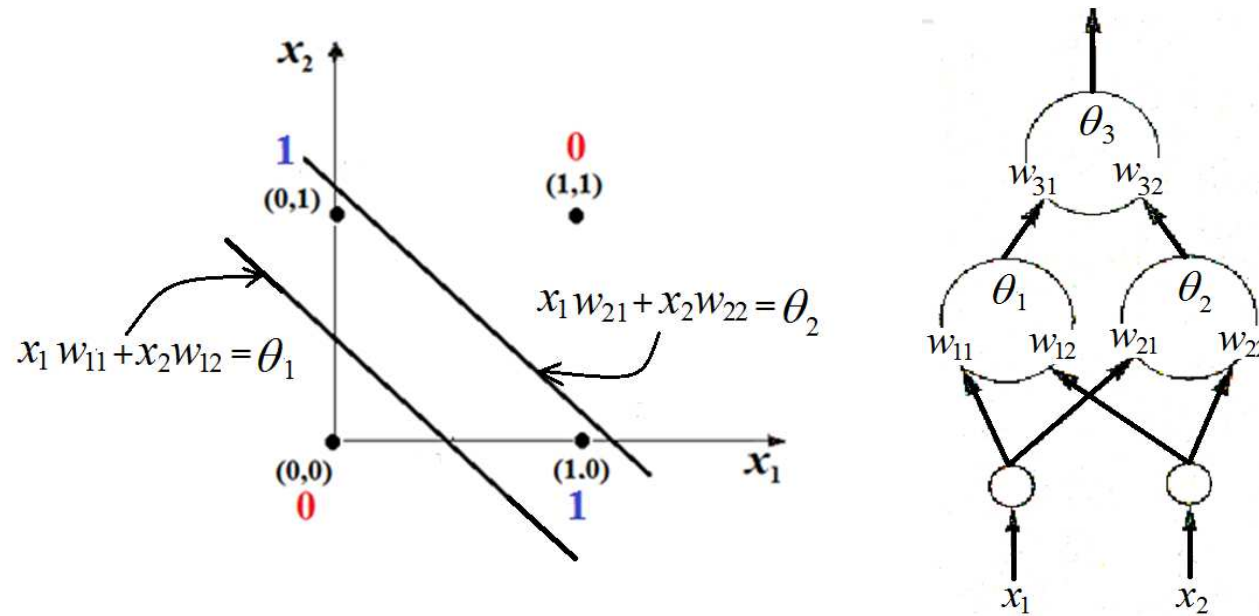
$x_1$	$x_2$	XOR
0	0	0
0	1	1
1	0	1
1	1	0

This problem cannot be solved by an adaline.



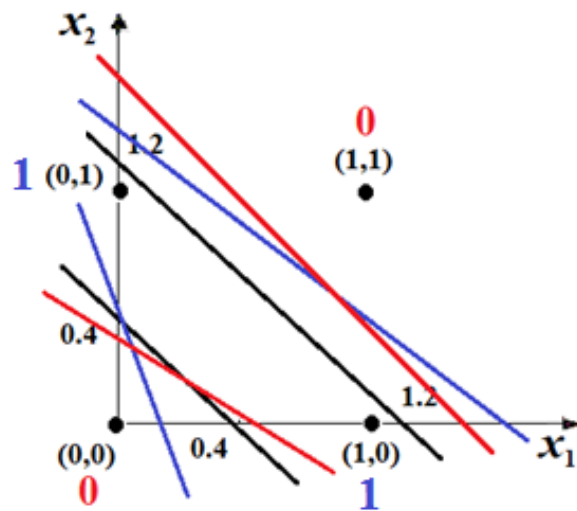
Reason:  $w_1x_1 + w_2x_2 = \theta$  specifies is a line in the  $(x_1, x_2)$  plane.



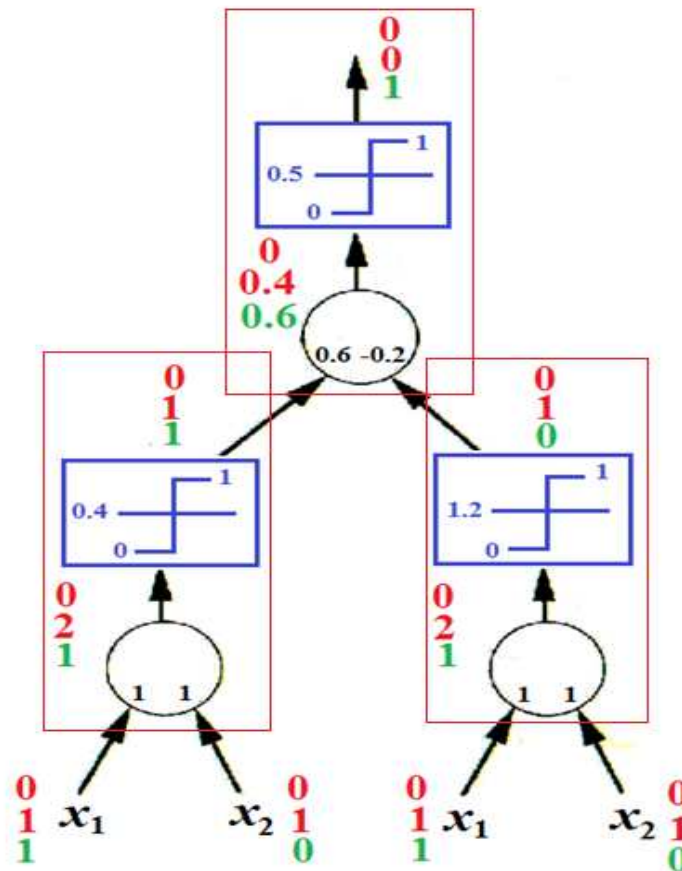


The two neurons in the hidden layer provides two lines that can separate the plane into three regions. The two regions containing  $(0,0)$  and  $(1,1)$  are associated with the network output of 0. The central region is associated with the network output of 1.

There are many solution pairs of lines.

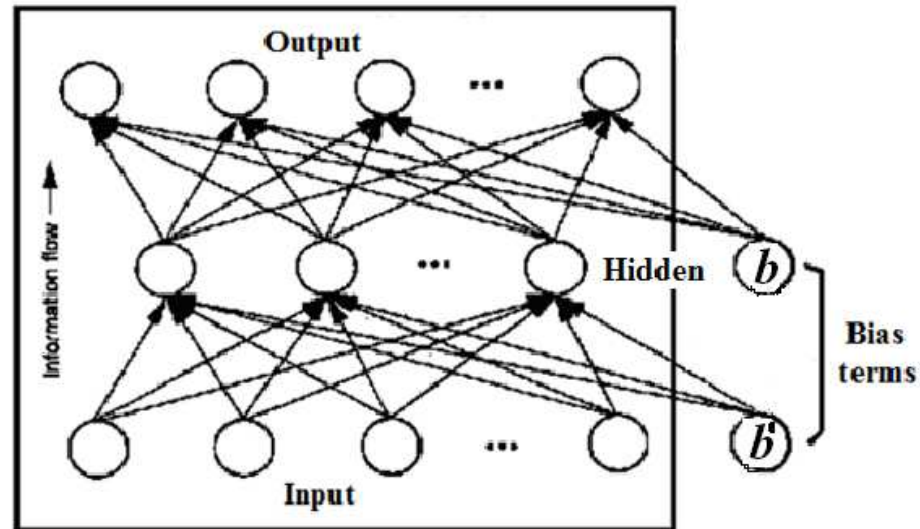


One solution:



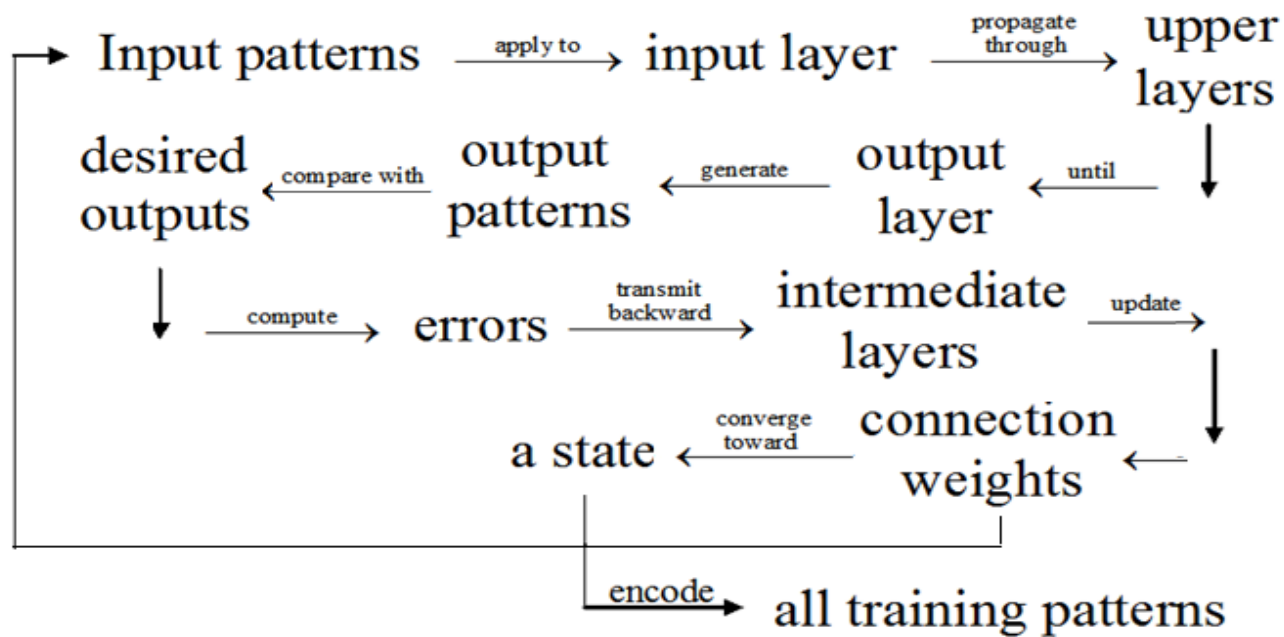
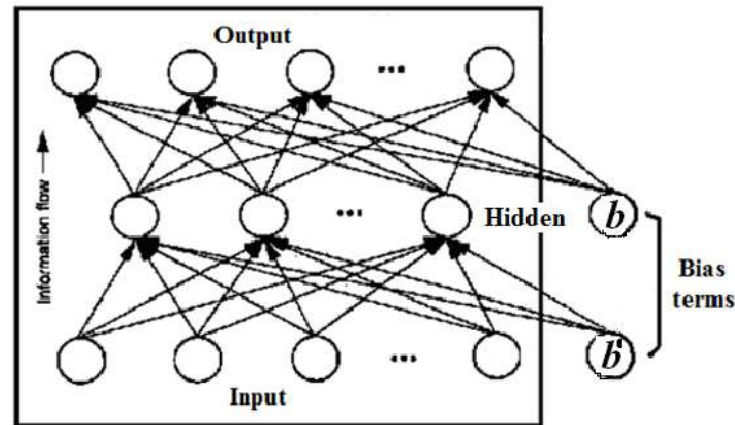
## 11.6 Train MLP by Backpropagation

- Architecture of MLP



During training, self-organization of nodes on the intermediate layers s.t. different nodes recognize different **features** or their **relationships**. Noisy and incomplete patterns can thus be handled.

- **Learning cycle**

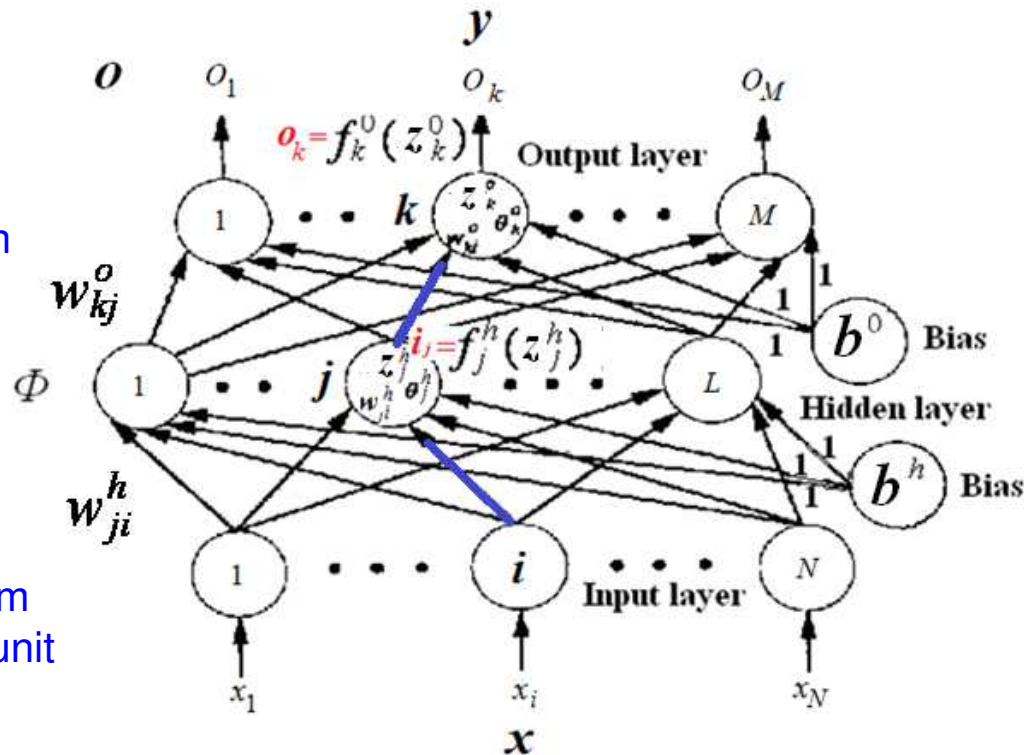


Given training examples:  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_p, \mathbf{y}_p), \dots, (\mathbf{x}_P, \mathbf{y}_P)$   
 where  $\mathbf{y}_p = \Phi(\mathbf{x}_p)$ , find an  $\text{NN}(\theta)$  approximating  $\Phi$   
 where  $\theta = (w_{ij}^l, b_i^l)_{i,j,l}$ . Consider input  $\mathbf{x} = (x_1, x_2, \dots, x_N)$

Net input to the  
 $j$ th hidden unit:

$$z_j^h = \sum_{i=1}^N w_{ji}^h x_i + b_j^h$$

$\nwarrow$   $\nwarrow$   $\nwarrow$   
 $j$ th hidden unit     $j$ th input unit    bias term with  $j$ th unit



Output of the  $j$ th hidden unit:  $i_j = f_j^h(z_j^h)$

Net input to the  $k$ th output unit:

transfer function

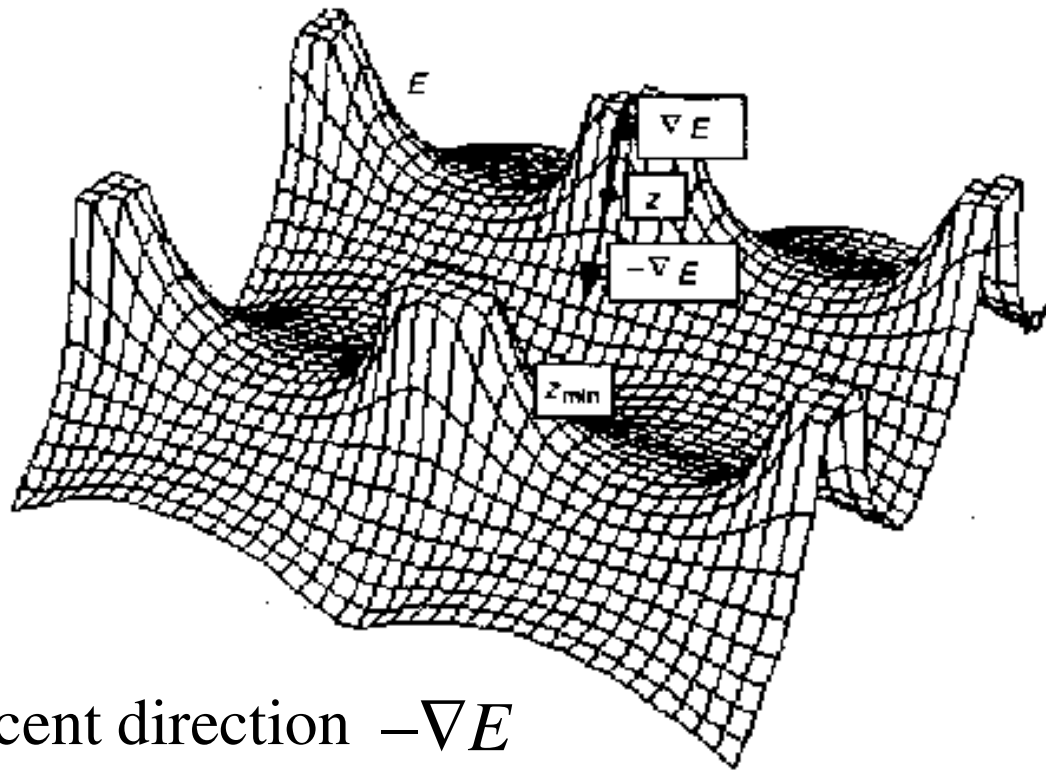
$$z_k^o = \sum_{j=1}^L w_{kj}^o i_j + b_k^o$$

Output of the  $k$ th output unit:  $o_k = f_k^o(z_k^o)$

- Update of output layer weights  $W^o$

The error to be minimized:

$$E = \frac{1}{2} \sum_{k=1}^M (y_k - o_k)^2, \text{ where } M: \# \text{ output units}$$



The descent direction  $-\nabla E$

The learning rule:

$$W^o(t+1) = W^o(t) + \Delta W^o = W^o(t) - \eta \nabla E$$

where:  $\eta$  learning rate

- Chain Rule:  $\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g)}{\partial g} \frac{\partial g(x)}{\partial x} = f'g'$

$$\frac{\partial f(g(h(x)))}{\partial x} = \frac{\partial f(g)}{\partial g} \frac{\partial g(h)}{\partial h} \frac{\partial h(x)}{\partial x} = f'g'h'$$

- Determine  $\nabla E$

$$E = \frac{1}{2} \sum_{k=1}^M (y_k - o_k)^2, \quad \frac{\partial E}{\partial w_{kj}^o} = -(y_k - o_k) \frac{\partial o_k}{\partial w_{kj}^o}, \quad o_k = f_k^o(z_k^o)$$

$$\frac{\partial o_k}{\partial w_{kj}^o} = \frac{\partial f_k^o(z_k^o)}{\partial w_{kj}^o} = \frac{\partial f_k^o(z_k^o)}{\partial z_k^o} \frac{\partial z_k^o}{\partial w_{kj}^o} = f_k^{o'}(z_k^o) \frac{\partial z_k^o}{\partial w_{kj}^o}$$

$$\frac{\partial z_k^o}{\partial w_{kj}^o} = \frac{\partial}{\partial w_{kj}^o} \left( \sum_{j=1}^L w_{kj}^o i_j + b_j^o \right) = i_j \quad \therefore \quad -\frac{\partial E}{\partial w_{kj}^o} = (y_k - o_k) f_k^{o'} i_j$$



The weights on the output layer are updated as

$$\begin{aligned}w_{kj}^o(t+1) &= w_{kj}^o(t) + \Delta w_{kj}^o(t) = w_{kj}^o(t) - \eta \frac{\partial E}{\partial w_{kj}^o} \\ &= w_{kj}^o(t) + \eta(y_k - o_k)f_k^{o'}i_j \quad \text{--- (A)}\end{aligned}$$

- Consider  $f_k^{o'}$

Two forms for the output functions  $f_k^o$

i) Linear  $f_k^o(x) = x$

ii) Sigmoid  $f_k^o(x) = (1 + e^{-\lambda x})^{-1}$  or  
 $f_k^o(x) = \frac{1}{2}[1 - \tanh(\lambda x)]$

- Sigmoid:

$$f(x) = \frac{1}{1 + e^{-x}} \Rightarrow \frac{df}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

- Tanh:

$$f(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}} \Rightarrow \frac{df}{dx} = 1 - \left( \frac{1 - e^{-2x}}{1 + e^{-2x}} \right)^2$$

- Softsign:

$$f(x) = \frac{x}{1 + |x|} \Rightarrow \frac{df}{dx} = \frac{1}{(1 + |x|)^2}$$

i) For linear function  $f_k^o(x) = x$ ,  $f_k^{o'} = 1$

$$(A) \Rightarrow w_{kj}^o(t+1) = w_{kj}^o(t) + \eta(y_k - o_k)i_j$$

ii) For sigmoid function  $f_k^o(x) = (1 + e^{-\lambda x})^{-1}$ ,

$$\begin{aligned} f_k^{o'}(z_k^o) &= -(1 + e^{-\lambda z_k^o})^{-2} \frac{\partial e^{-\lambda z_k^o}}{\partial z_k^o} \\ &= -(1 + e^{-\lambda z_k^o})^{-2} e^{-\lambda z_k^o} (-\lambda) = \lambda(1 + e^{-\lambda z_k^o})^{-2} e^{-\lambda z_k^o} \\ &= \lambda f_k^{o2} (f_k^{o-1} - 1) = \lambda f_k^o (1 - f_k^o) = \lambda o_k (1 - o_k) \end{aligned}$$

$$(A) \Rightarrow w_{kj}^o(t+1) = w_{kj}^o(t) + \eta \lambda (y_k - o_k) o_k (1 - o_k) i_j$$

© Updates of hidden-layer weights

Difficulty: Unknown outputs of the hidden-layer units

Idea: Relate error  $E$  to the output of the hidden layer

$$\begin{aligned} E(w_{ji}^h) &= \frac{1}{2} \sum_{k=1}^M (y_k - o_k)^2 = \frac{1}{2} \sum_k [y_k - f_k^o(z_k^o)]^2 \\ &= \frac{1}{2} \sum_k [y_k - f_k^o(\sum_{j=1}^L w_{kj}^o i_j + b_k^o)]^2 \\ &= \frac{1}{2} \sum_k [y_k - f_k^o(\sum_{j=1}^L w_{kj}^o f_j^h(z_j^h) + b_k^o)]^2 \\ &= \frac{1}{2} \sum_k [y_k - f_k^o(\sum_{j=1}^L w_{kj}^o f_j^h(\sum_{i=1}^N w_{ji}^h x_i + b_j^h) + b_k^o)]^2 \end{aligned}$$

$$\begin{aligned}
\frac{\partial E}{\partial w_{ji}^h} &= \frac{1}{2} \frac{\partial}{\partial w_{ji}^h} \left[ \sum_k (y_k - o_k)^2 \right] = \frac{1}{2} \sum_k \left[ \frac{\partial}{\partial w_{ji}^h} (y_k - o_k)^2 \right] \\
&= - \sum_k (y_k - o_k) \frac{\partial o_k}{\partial w_{ji}^h} \\
\frac{\partial o_k}{\partial w_{ji}^h} &= \frac{\partial f_k^o(z_k^o)}{\partial w_{ji}^h} = \frac{\partial f_k^o(z_k^o)}{\partial z_k^o} \frac{\partial z_k^o}{\partial i_j} \frac{\partial i_j}{\partial w_{ji}^h} \\
&= f_k^{o'}(z_k^o) \frac{\partial z_k^o}{\partial i_j} \frac{\partial i_j}{\partial w_{ji}^h} \\
\frac{\partial i_j}{\partial w_{ji}^h} &= \frac{\partial f_j^h(z_j^h)}{\partial w_{ji}^h} = \frac{\partial f_j^h(z_j^h)}{\partial z_j^h} \frac{\partial z_j^h}{\partial w_{ji}^h} = f_j^{h'}(z_j^h) \frac{\partial z_j^h}{\partial w_{ji}^h}
\end{aligned}$$

$$\frac{\partial o_k}{\partial w_{ji}^h} = f_k^{o'}(z_k^o) \frac{\partial z_k^o}{\partial i_j} f_j^{h'}(z_j^h) \frac{\partial z_j^h}{\partial w_{ji}^h}$$

$$z_k^o = \sum_j w_{kj}^o i_j + b_k^o, \quad \frac{\partial z_j^o}{\partial i_j} = w_{kj}^o$$

$$z_j^h = \sum_i w_{ji}^h x_i + b_j^h, \quad \frac{\partial z_j^h}{\partial w_{ji}^h} = x_i$$

$$\frac{\partial o_k}{\partial w_{ji}^h} = f_k^{o'}(z_k^o) w_{kj}^o f_j^{h'}(z_j^h) x_i$$

Consider sigmoid output function  $f(x) = (1 + e^{-\lambda x})^{-1}$

$$f_k^{o'}(z_k^o) = \lambda o_k (1 - o_k), \quad f_j^{h'}(z_j^h) = \lambda i_j (1 - i_j)$$

$$\frac{\partial o_k}{\partial w_{ji}^h} = \lambda^2 o_k (1 - o_k) w_{kj}^o i_j (1 - i_j) x_i$$

$$\begin{aligned} \frac{\partial E}{\partial w_{ji}^h} &= - \sum_k (y_k - o_k) \frac{\partial o_k}{\partial w_{ji}^h} \\ &= -\lambda^2 \sum_k (y_k - o_k) o_k (1 - o_k) w_{kj}^o i_j (1 - i_j) x_i \\ &= -\lambda^2 i_j (1 - i_j) x_i \sum_k (y_k - o_k) o_k (1 - o_k) w_{kj}^o \end{aligned}$$

$$\Delta w_{ji}^h = -\eta \frac{\partial E}{\partial w_{ji}^h} = -\eta' i_j (1 - i_j) x_i \sum_k (y_k - o_k) o_k (1 - o_k) w_{kj}^o$$

$$\text{where } \eta' = \eta \lambda^2.$$

The update rule:  $w_{ji}^h(t+1) = w_{ji}^h(t) + \Delta w_{ji}^h$

- Summary

1. Apply the input vector,  $\mathbf{x}_p = (x_{p1}, x_{p2}, \dots, x_{pN})^t$  to the input units.

2. Calculate the net-input values to the hidden layer units:

$$z_{pj}^h = \sum_{i=1}^N w_{ji}^h x_{pi} + b_j^h$$

3. Calculate the outputs from the hidden layer:  $i_{pj} = f_j^h(z_{pj}^h)$

4. Calculate the net-input values to each output unit:

$$z_{pk}^o = \sum_{j=1}^L w_{kj}^o i_{pj} + b_k^o$$

5. Calculate the outputs:  $o_{pk} = f_k^o(z_{pk}^o)$



6. Calculate the error terms for the output units:

$$\delta_{pk}^o = (y_{pk} - o_{pk}) f_k^{o'}$$

7. Calculate the error terms for the hidden units:

$$\delta_{pj}^h = f_j^{h'} \sum_i \delta_{pk}^o w_{kj}^o$$

8. Update weights on the output layer:

$$w_{kj}^o(t+1) = w_{kj}^o(t) + \eta \delta_{pk}^o i_{pj}$$

9. Update weights on the hidden layer:

$$w_{ji}^h(t+1) = w_{ji}^h(t) + \eta \delta_{pj}^h x_i$$

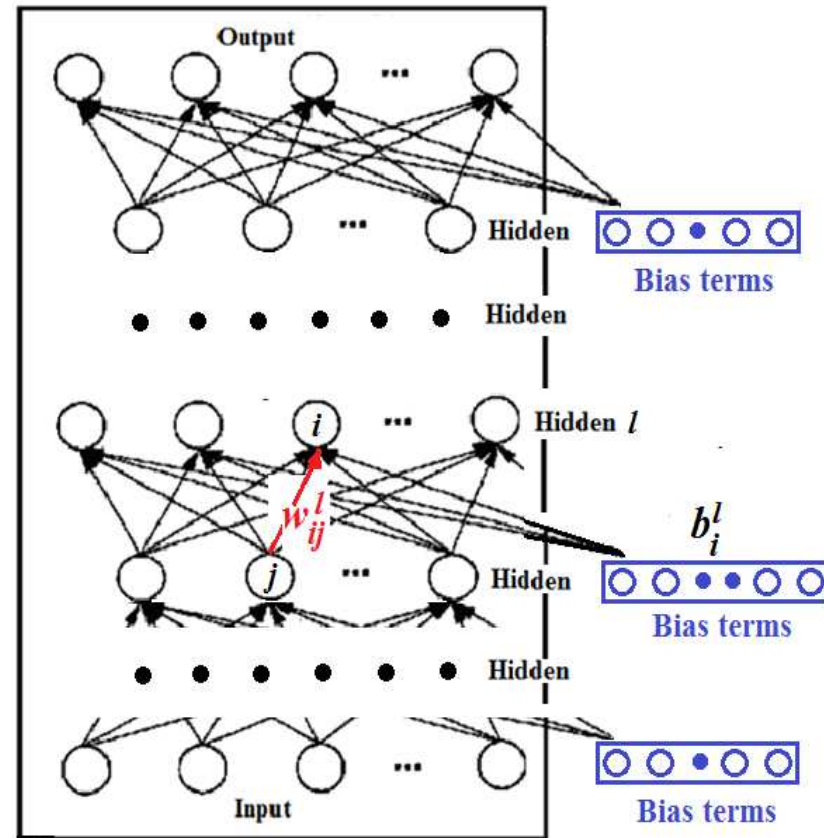
◎ Updates of multiple hidden-layers

$$E = \frac{1}{2} \sum_{k=1}^M (y_k - o_k)^2, \quad o_k = f_k^o(z_k^o),$$

$$z_k^o = \sum_{j=1}^L w_{kj}^o i_j + b_k^o, \quad i_j = f_j^h(z_j^h)$$

$$z_k^o = \sum_{j=1}^L w_{kj}^o f_j^h(z_j^h) + b_k^o,$$

$$o_k = f_k^o \left( \sum_{j=1}^L w_{kj}^o f_j^h(z_j^h) + b_k^o \right)$$



$$\begin{aligned}
E(w_{kj}^o) &= \frac{1}{2} \sum_{k=1}^M \left( y_k - f_k^o \left( \sum_{j=1}^L w_{kj}^o f_j^h(z_j^h) + b_k^o \right) \right)^2 \\
&\quad \left( z_j^h = \sum_{i=1}^N w_{ji}^h x_i + b_j^h \right) \\
&= \frac{1}{2} \sum_k [y_k - f_k^o (\sum_{j=1}^L w_{kj}^o f_j^h (\sum_{i=1}^N w_{ji}^h x_i + b_j^h) + b_k^o)]^2 \\
E(w_{ji}^h) &= \frac{1}{2} \sum_k [y_k - f_k^o (\sum_{j=1}^L w_{kj}^o f_j^h (\sum_{i=1}^N w_{ji}^h x_i + b_j^h) + b_k^o)]^2 \\
E(w_{j_1 i}^{h_1}) &= \frac{1}{2} \sum_k [y_k - f_k^o (\sum_{j_1=1}^{L_1} w_{kj_1}^o f_{j_1}^{h_1} (\sum_{i=1}^N w_{j_1 i}^{h_1} x_i + b_{j_1}^{h_1}) + b_k^o)]^2
\end{aligned}$$

$$E(w_{j_2 j_1}^{h_2}) = \frac{1}{2} \sum_k [y_k - f_k^o(\sum_{j_1=1}^{L_1} w_{kj_1}^o f_{j_1}^{h_1}(\sum_{j_2=1}^{L_2} w_{j_1 j_2}^{h_1} f_{j_2}^{h_2}(\sum_{i=1}^N w_{j_2 i}^{h_2} x_i + b_{j_2}^{h_2}) + b_{j_1}^{h_1}) + b_k^o)]^2$$

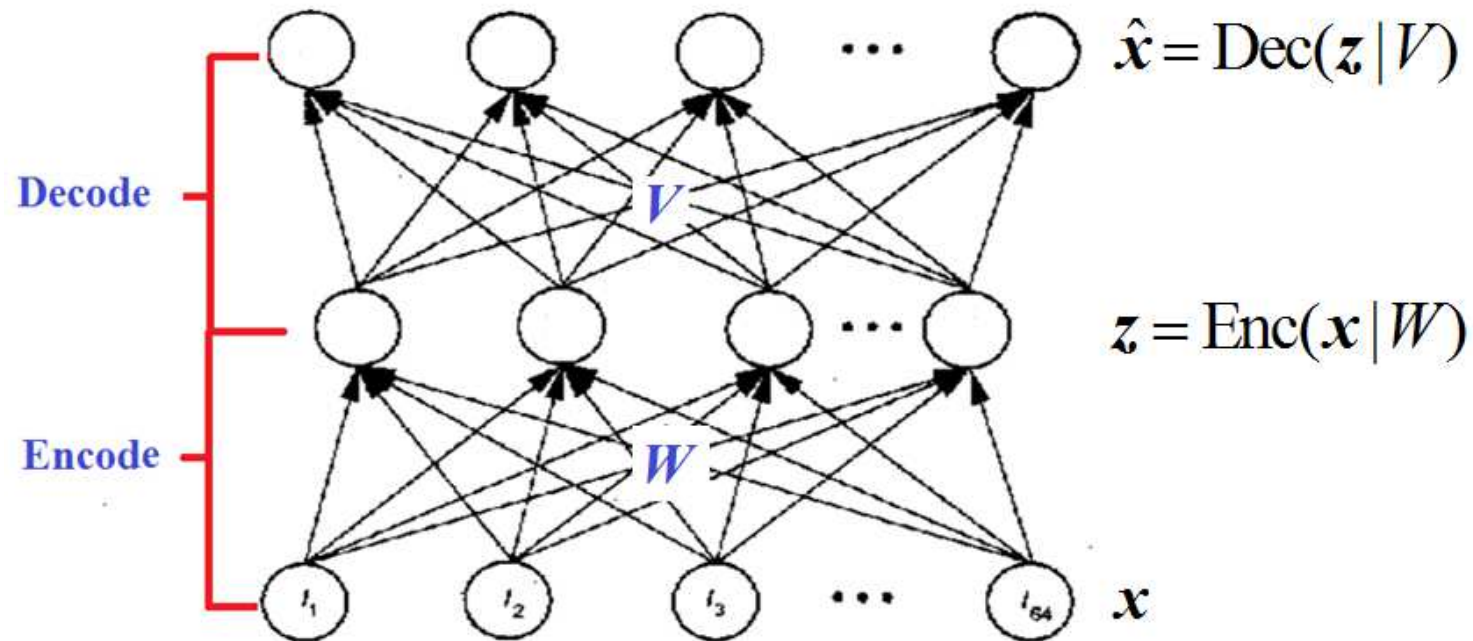
.....

.....

$$E(w_{j_n j_{n-1}}^{h_n}) = \frac{1}{2} \sum_k [y_k - \dots]$$

**Loss functions:** Error function, cross-entropy, differential entropy, negative log-likelihood, information gain

## 11.7 Autoencoders



Variants: Denoising autoencoder

Sparse autoencoder

## 11.8 ANNs

Hopfield neural model

Associative memory

Counterpropagation neural networks

Self-organization feature map

Adaptive resonance theory

Neocognitron

Recurrent neural networks

Boltzmann machine

Support vector machine

Kernel machines