# Outline
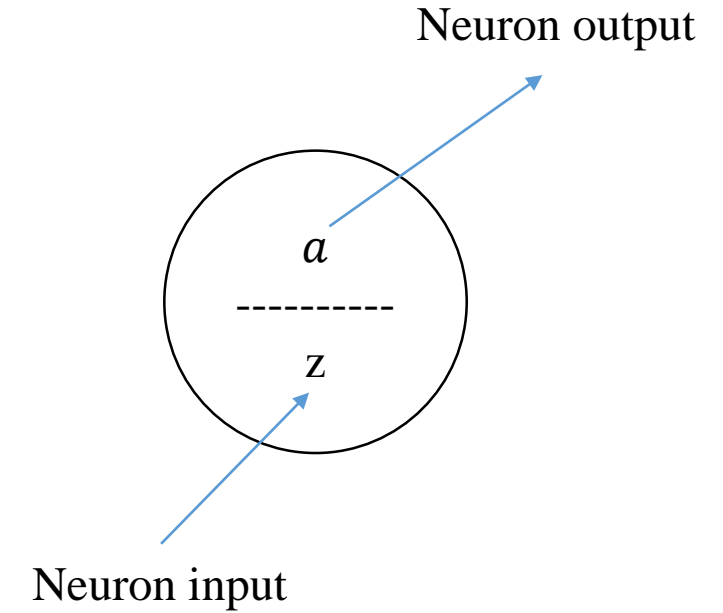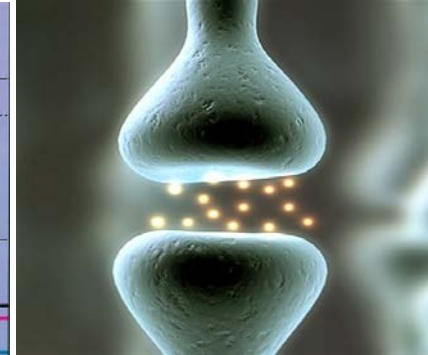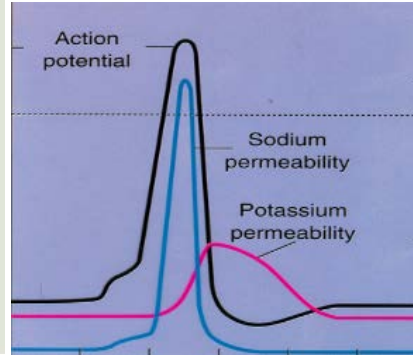
■Brief Review of Computational Model of Neural Networks
■Network Performance: Cost Function
■Steepest Gradient Method
■Backpropagation
■Three Pages to Understand BP
■Only One Page to Understand BP
■The BP Algorithm
■Assignment

# Computational Model of Neurons

Neuron output

$$a$$
$$----------$$
$$z$$

Neuron input

$x_1$

Connection weights

$w_1$

$x_i$

$w_i$

$x_n$

$w_n$

inputs

Activation function

$$z = \sum_{i=1}^{n} w_i x_i \quad f$$

Total input

$$a = f(z)$$

Neuron output

$$y = f\left(\sum_{i=1}^{n} w_i x_i\right)$$

# Computational Model of Neural Networks

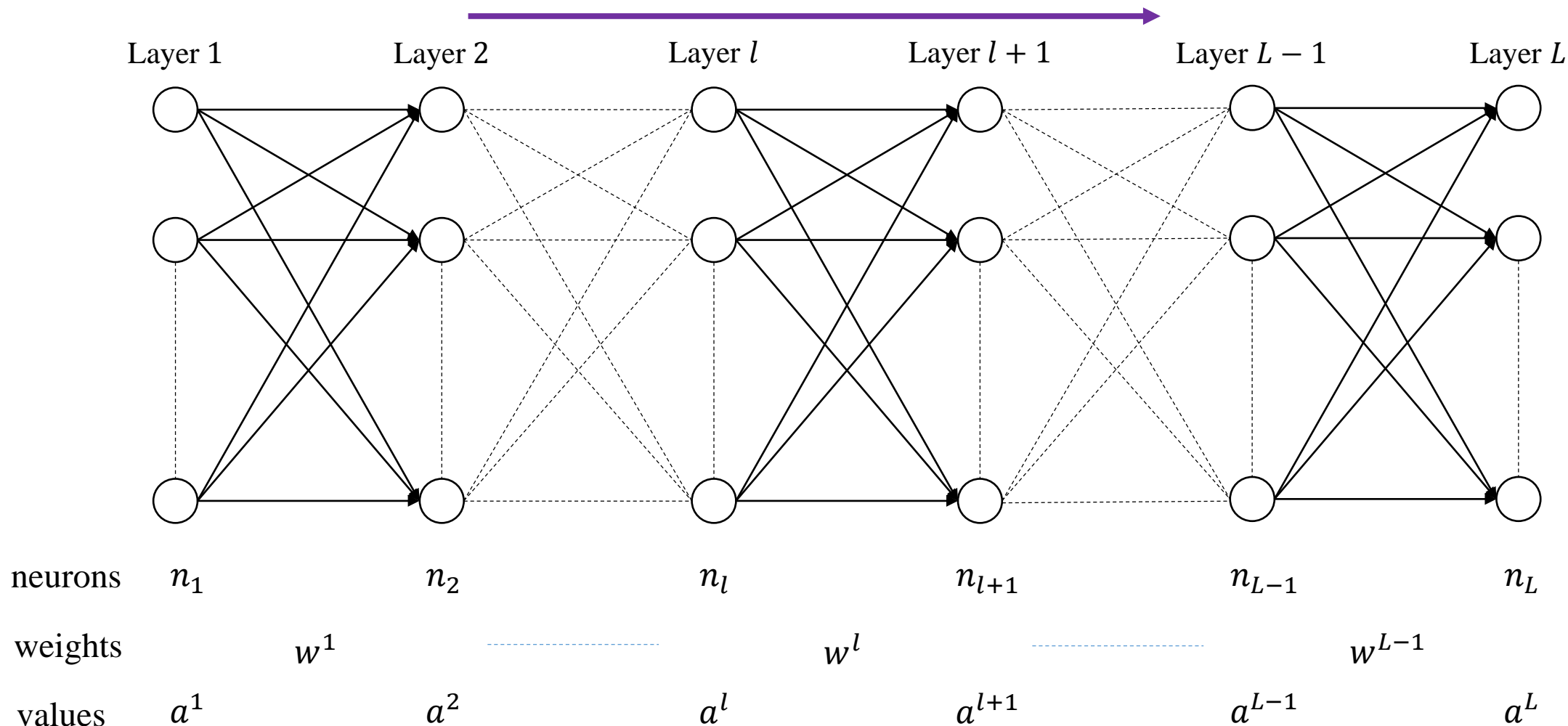Forward computing

| Layer 1 | Layer 2 | Layer $l$ | Layer $l+1$ | Layer $L-1$ | Layer $L$ |
|---------|---------|-----------|-------------|-------------|-----------|



| | | | | | | |
|---------|---------|---------|-----------|-------------|-------------|-----------|
| neurons | $n_1$ | $n_2$ | $n_l$ | $n_{l+1}$ | $n_{L-1}$ | $n_L$ |
| weights | | $w^1$ | $w^l$ | | $w^{L-1}$ | |
| values | $a^1$ | $a^2$ | $a^l$ | $a^{l+1}$ | $a^{L-1}$ | $a^L$ |

# Computational Model of Neural Networks



layer $l$ contains $n_l$ neurons.

The neuron located in $l$ layer $j^{th}$ place, $a_j^l$ denotes the output value of the neuron.

$$a_j^l = f(z_j^l)$$

Vector form

$$a^l = \begin{bmatrix} a_1^l \\ \vdots \\ a_j^l \\ \vdots \\ a_{n_l}^l \end{bmatrix}$$

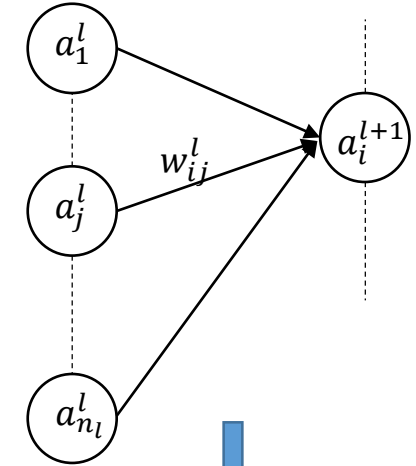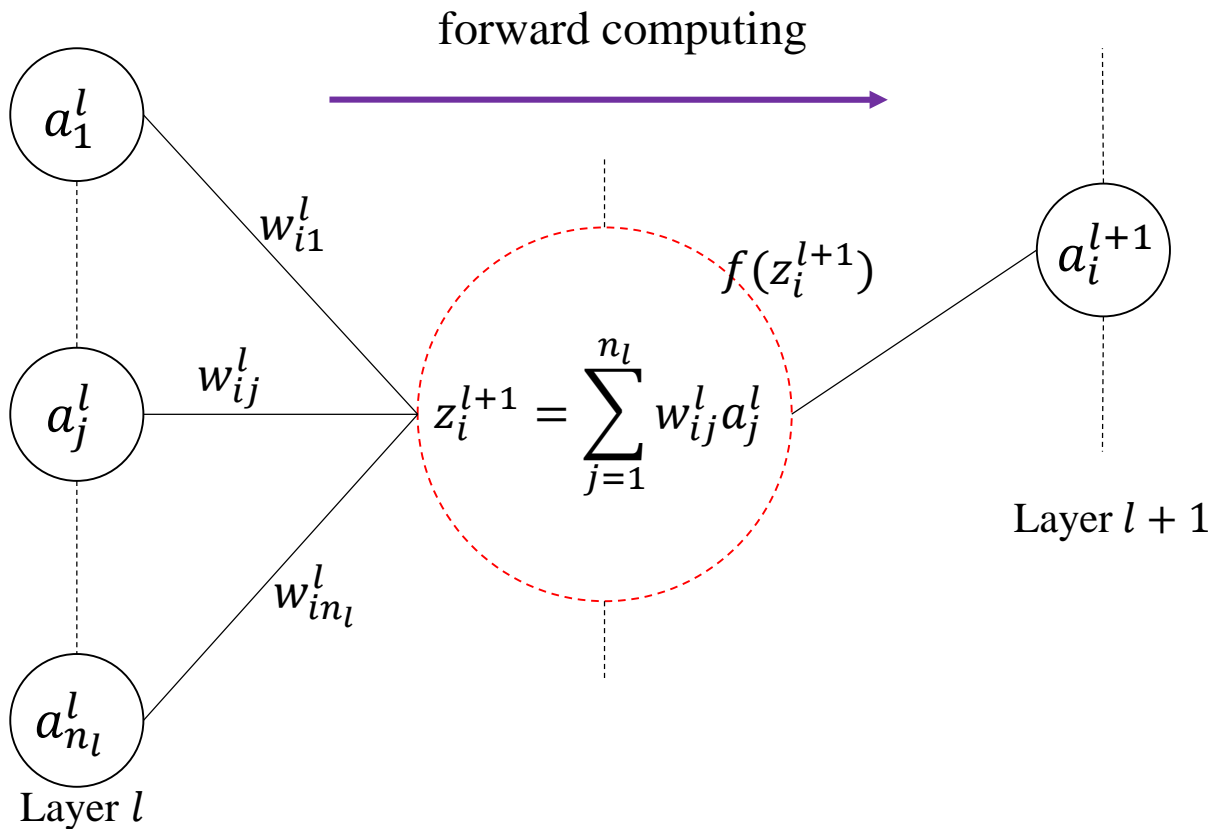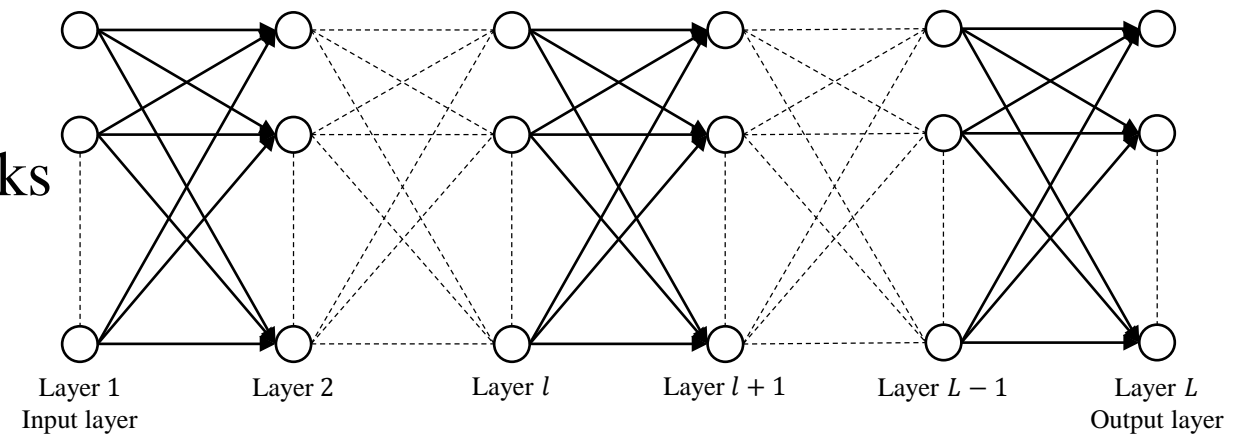Layer $l$　　　Layer $l+1$

$a^l$　　$w^l$　　$a^{l+1}$

$a^l$: input of $l+1$ layer
$a^{l+1}$: representation of $a^l$

$$W^l = \begin{bmatrix} w_{11}^l & \cdots & w_{1n_l}^l \\ \vdots & w_{ij}^l & \vdots \\ w_{n_{l+1}1}^l & \cdots & w_{n_{l+1}n_l}^l \end{bmatrix}_{n_{l+1} \times n_l}$$

# Computational Model of Neural Networks



forward computing

$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l$$

$f(z_i^{l+1})$

$a_i^{l+1}$

Layer $l+1$

$a_1^l$, $w_{i1}^l$

$a_j^l$, $w_{ij}^l$

$a_{n_l}^l$, $w_{in_l}^l$

Layer $l$

Component form
$$\begin{cases} a_i^{l+1} = f(z_i^{l+1}) \\ z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l \end{cases}$$

Vector form
$$\begin{cases} a^{l+1} = f(z^{l+1}) \\ z^{l+1} = w^l a^l \end{cases}$$

$$a_i^{l+1} = f\left(\sum_{j=1}^{n_l} w_{ij}^l a_j^l\right)$$

# Computational Model of Neural Networks

$$a^L = f(w^{L-1}a^{L-1}) = f\left(w^{L-1}f\left(w^{L-2}f(w^{L-3}\cdots f(w^1a^1))\right)\right)$$



$R^{n_1}$

$$f\left(w^{L-1}f\left(w^{L-2}f(w^{L-3}\cdots f(w^1a^1))\right)\right)$$

$R^{n_L}$

$a^1 \quad w^1 \quad a^2 \qquad a^l \quad w^l \quad a^{l+1} \qquad a^{L-1} \quad w^{L-1} \quad a^L$

Input

Output

7

# Computational Model of Neural Networks

**External inputs:**

If neurons in $l$ layer are not connected to any neurons in the $l-1$ layer, these neurons are called external inputs of $l+1$ layer.

External inputs can exist in any layer except the last one.

# Example: XOR problem

Problem: How to design the NN? Are there any methods to design the connection weights?

# Example: Handwritten Digits Recognition



$a^1$  $w^1$  $a^2$   $a^l$  $w^l$  $a^{l+1}$   $a^{L-1}$ $w^{L-1}$  $a^L$

Input   Output

Problem: How to design the NN? Are there any methods to design the connection weights?

# Outline

■Brief Review of Computational Model of Neural Networks

■<span style="color:red">Network Performance: Cost Function</span>

■Steepest Gradient Method

■Backpropagation

■Three Pages to Understand BP

■Only One Page to Understand BP

■The BP Algorithm

■Assignment

# Network Performance: Cost Function



Training

Good Performance!
*The father knows the correct answer.*

Two important factors:

1. There must be a measure to measure the correctness between correct answer and the boy's real output. ----- <span style="color:red">Performance function</span>.

2. There must be a mechanism to change the knowledge system of the boy. ---- <span style="color:red">Learning algorithm</span>.

# Network Performance: Cost Function



Network Output

$$a^L = \begin{bmatrix} a_1^L \\ \vdots \\ a_{n_L}^L \end{bmatrix}$$

Target Output

$$y^L = \begin{bmatrix} y_1^L \\ \vdots \\ y_{n_L}^L \end{bmatrix}$$

$a^1$   $w^1$   $a^2$     $a^l$   $w^l$   $a^{l+1}$     $a^{L-1}$ $w^{L-1}$   $a^L$

Input                                       Output

Changing the weights: Learning algorithm

$J(a^L, y^L)$
Performance function $J(a^L, y^L)$, or cost function, is used to describe the distance between $a^L$ and $y^L$, $J(a^L, y^L)$ is indeed a function of $(w^1, \cdots, w^L)$, i.e.,
$$J = J(w^1, \cdots, w^L).$$

Problem: How to construct a cost function?

13

# Network Performance: Cost Function



$$a^1 \quad w^1 \quad a^2 \qquad a^l \quad w^l \quad a^{l+1} \qquad a^{L-1} \, w^{L-1} \quad a^L$$

Input

Output

Target Output

$$y^L = \begin{bmatrix} y_1^L \\ \vdots \\ y_{n_L}^L \end{bmatrix}$$

Network Output

$$a^L = \begin{bmatrix} a_1^L \\ \vdots \\ a_{n_L}^L \end{bmatrix}$$

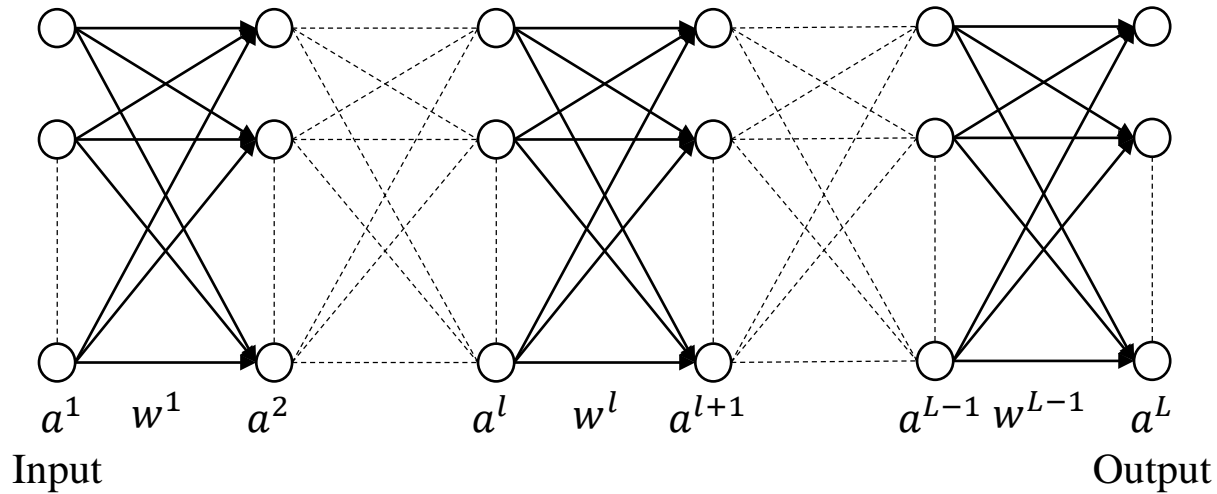A cost function $J$ describes the performance of the network. If the $J$ is small, it implies that the network output $a^L$ close to the target output $y^L$, the network is called in good performance. Since $J$ is a function with variables $(w^1, \cdots, w^L)$, good performance means to find suitable $(w^1, \cdots, w^L)$ such that $J$ is small. The process of looking for suitable $(w^1, \cdots, w^L)$ is called network learning.
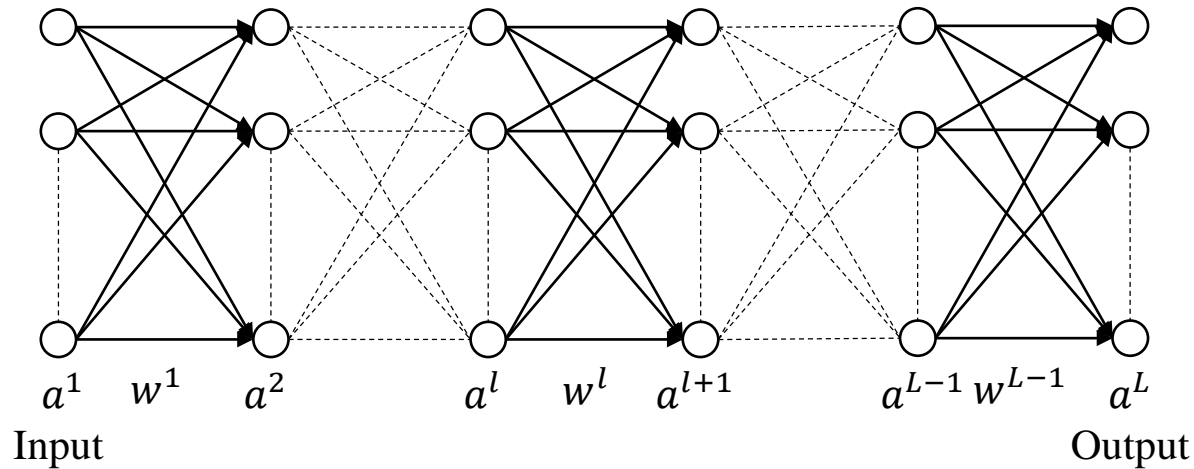
Problem: How to learn?

There are many ways to construct a cost function. A frequently used cost is as follows:
$$e_j = a_j^L - y_j^L$$

$$J = \frac{1}{2} \sum_{j=1}^{n_L} e_j^2 = J(w^1, \cdots, w^L)$$

Clearly, $J$ is a function of $w^1, \cdots, w^L$.

# Network Performance: Cost Function



$a^1$ $w^1$ $a^2$ $\quad$ $a^l$ $w^l$ $a^{l+1}$ $\quad$ $a^{L-1}$ $w^{L-1}$ $a^L$

Input $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Output

Target Output

$$y^L = \begin{bmatrix} y_1^L \\ \vdots \\ y_{n_L}^L \end{bmatrix}$$

Network Output

$$a^L = \begin{bmatrix} a_1^L \\ \vdots \\ a_{n_L}^L \end{bmatrix}$$

Learning is a process such that $a^L$ is close to $y^L$, i.e., the cost function $J$ reaches minimum. A cost function $J = J(w^1, \cdots, w^{L-1})$ is a function with variables $w^l (l = 1, \cdots, L)$, thus the network learning is to looking for some $w^l (l = 1, \cdots, L)$ such that $w^l (l = 1, \cdots, L)$ is a minimum point of $J$.

A frequently used cost function:

$$J = \frac{1}{2} \sum_{j=1}^{n_L} e_j^2 = J(w^1, \cdots, w^L)$$

$J$ is a function of $w^1, \cdots, w^L$.

Problem: How to find out the minimum points of $J$ ?

# Outline

■Brief Review of Neural Networks Structure

■Network Performance: Cost Function

■<span style="color:red">Steepest Gradient Method</span>

■Backpropagation

■Three Pages to Understand BP

■Only One Page to Understand BP

■The BP Algorithm

■Assignment

# Minimum Points

$$F(x) = (x_2 - x_1)^4 + 8x_1x_2 - x_1 + x_2 + 3$$

$$F(x) = (x_1^2 - 1.5x_1x_2 + 2x_2^2)x_1^2$$



Minima

Contour

Minimum points

Problem:
How to find the minimum points?

# Iteration Method

Finding a minimum point step by step

$$x^{k+1} = x^k + \alpha_k \cdot p_k$$

*To begin the iteration, you must need a given starting point $x_0$.*

$p_k$ is called searching direction

$\alpha_k$ is learning rate at step $k$.

Problem: How to get the searching direction $p_k$?

# Steepest Descent Method

Slowest changing direction

Fastest increasing direction

Steepest descent direction

Gradient:

$$g_k = \nabla F(x)\Big|_{x^k} = \frac{\partial F}{\partial x}\Big|_{x^k} = \begin{pmatrix} \dfrac{\partial F}{\partial x_1} \\ \vdots \\ \dfrac{\partial F}{\partial x_n} \end{pmatrix}\Bigg|_{x^k}$$

Steepest Descent Algorithm:

$$p_k = -g_k$$
$$x^{k+1} = x^k - \alpha_k \cdot g_k$$

or

$$x^{k+1} = x^k - \alpha_k \cdot \frac{\partial F}{\partial x}\Big|_{x^k}$$

$\alpha_k \cdot p_k$

$x^k$

$x^{k+1}$

# Steepest Descent Method



Updating weights
$$w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$$

Computing gradient
$$\frac{\partial J}{\partial w_{ji}^l}$$

Construct cost function
$$J = \frac{1}{2} \sum_{i=1}^{n_L} (y_i - a_j^L)^2$$

Net output

Target output

# Steepest Descent Method



$a^1$  $w^1$  $a^2$  $a^l$  $w^l$  $a^{l+1}$  $a^{L-1}$  $w^{L-1}$  $a^L$

Input

Output

Target Output

$$y^L = \begin{bmatrix} y_1^L \\ \vdots \\ y_{n_L}^L \end{bmatrix}$$

Network Output

$$a^L = \begin{bmatrix} a_1^L \\ \vdots \\ a_{n_L}^L \end{bmatrix}$$

Steepest Descent Method

$$J = \frac{1}{2}\sum_{j=1}^{n_L} e_j^2 = \frac{1}{2}\sum_{j=1}^{n_L}\left(a_j^L - y_j^L\right)^2$$

1. Computing

$$\frac{\partial J}{\partial w_{ji}^l}$$

2. Iterating

$$w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$$

$$a^L = f(w^{L-1}a^{L-1}) = f\left(w^{L-1}f\left(w^{L-2}f(w^{L-3}\cdots f(w^1 a^1))\right)\right)$$

<span style="color:red">Problem: How to compute $\dfrac{\partial J}{\partial w_{ji}^l}$?</span>

Answer:
Using the well-known BP method.

# Outline

- Brief Review of Neural Networks Structure
- Network Performance: Cost Function
- Steepest Gradient Method
- <span style="color:red">Backpropagation</span>
- Three Pages to Understand BP
- Only One Page to Understand BP
- The BP Algorithm
- Assignment

# Backpropagation

Forward computing $a^l$



Back propagation $\delta^l$

Layer 1

Layer L

Backpropagation is a efficient way to calculate
$$\frac{\partial J}{\partial w_{ji}^l}$$

Cost function:

$$J = \frac{1}{2}\sum_{j=1}^{n_L} e_j^2 = \frac{1}{2}\sum_{j=1}^{n_L}\left(a_j^L - y_j^L\right)^2$$

$l$ layer

Problem: What's the relation between $\delta_i^l$ and $\frac{\partial J}{\partial w_{ji}^l}$ ?

$l + 1$ layer

$l$ layer

$a_i^l = f(z_i^l)$

define $\delta_i^l = \frac{\partial J}{\partial z_i^l}$

$a_i^l = f(z_i^l)$
------------------
$\delta_i^l = \frac{\partial J}{\partial z_i^l}$

$w_{ji}^l$

$a_j^{l+1} = f(z_j^{l+1})$
--------------
$\delta_j^{l+1} = \frac{\partial J}{\partial z_j^{l+1}}$

$J(w^1, \cdots, w^{L-1})$

Relation between $\delta_i^l$ and $\frac{\partial J}{\partial w_{ji}^l}$

$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

Why?

$$\frac{\partial J}{\partial w_{ji}^l} = \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

$a^1 \quad w^1 \quad a^2 \qquad a^l \quad w^l \quad a^{l+1} \qquad a^{L-1} w^{L-1} \ a^L$

Input

Output

Next problem:

What's the relation between $\delta_i^l$ and $\delta_j^{l+1}$ ?

24

Relation between $\delta_i^l$ and $\delta_j^{l+1}$

$\delta_1^{l+1}$

$w_{1i}^l$

$\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}$

$w_{ji}^l$

$\delta_j^{l+1}$

$\dot{f}(z_i^l)$

$w_{n_{l+1}i}^l$

$\delta_i^l$

$\delta_{n_{l+1}}^{l+1}$

$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$

$l$ layer    &larr;    $l+1$ layer

$a_1^{l+1} = f(z_1^{l+1})$
$\cdots\cdots\cdots$
$\delta_1^{l+1} = \dfrac{\partial J}{\partial z_1^{l+1}}$

$w_{1i}^l$

$a_i^l = f(z_i^l)$
$\cdots\cdots\cdots$
$\delta_i^l = \dfrac{\partial J}{\partial z_i^l}$

$w_{ji}^l$

$a_j^{l+1} = f(z_j^{l+1})$
$\cdots\cdots\cdots$
$\delta_j^{l+1} = \dfrac{\partial J}{\partial z_j^{l+1}}$

$w_{n_{l+1}i}^l$

$a_{n_{l+1}}^{l+1} = f(z_{n_{l+1}}^{l+1})$
$\cdots\cdots$
$\delta_{n_{l+1}}^{l+1} = \dfrac{\partial J}{\partial z_{n_{l+1}}^{l+1}}$

$$z_j^{l+1} = \sum_{i=1}^{n_l} w_{ji}^l a_i^l = \sum_{i=1}^{n_l} w_{ji}^l f(z_i^l)$$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \, \dot{f}(z_i^l) = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right)$$

Relation between $\delta_i^l$ and $\delta_j^{l+1}$

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right)$$

$$\delta_i^L = \frac{\partial J}{\partial z_i^L}$$

If

$$J = \frac{1}{2} \sum_{j=1}^{n_L} \left( a_j^L - y_j^L \right)^2 \qquad \textbf{\color{red}3}$$

then,

$$\delta_i^L = \frac{\partial J}{\partial z_i^L} = \left( a_i^L - y_i^L \right) \cdot \frac{\partial a_j^L}{\partial z_i^L} = \left( a_i^L - y_i^L \right) \cdot \dot{f}(z_i^L)$$

$$\textbf{\color{red}4}$$

Backpropagation $\delta$

# Outline

- Brief Review of Neural Networks Structure
- Network Performance: Cost Function
- Steepest Gradient Method
- Backpropagation
- <span style="color:red">Three Pages to Understand BP</span>
- Only One Page to Understand BP
- The BP Algorithm
- Assignment

# Three Pages to Understand BP: *The first page*

Cost function: $J(w^1, \cdots, w^L)$

Updating rule: $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$

Relationship: $\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$

**6**

**7**



$l$ layer

$$a_i^l = f(z_i^l)$$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l}$$

Forward computing

$a_1^l$

$a_j^l$

$a_{n_l}^l$

$w_{ij}^l$

$a_i^{l+1}$

$\delta_1^{l+1}$

$\delta_j^{l+1}$

$\delta_i^l$

$w_{ji}$

$\delta_{n_{l+1}}^{l+1}$

Back propagation

$a^l \quad w^l \quad a^{l+1}$

# Three Pages to Understand BP: *The second page*



forward computing

$l$ layer

$l + 1$ layer

$$a_i^{l+1} = f\left(\sum_{j=1}^{n_l} w_{ij}^l a_j^l\right)$$

or ------------------

$$a_i^{l+1} = f\left(z_i^{l+1}\right) \qquad \textbf{1}$$

$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l \qquad \textbf{2}$$

# Three Pages to Understand BP: *The third page*

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$$

**5**



$\dot{f}(z_i^l)$

$\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}$

$\delta_i^l$

$\delta_1^{l+1}$

$w_{1i}^l$

$w_{ji}^l$

$\delta_j^{l+1}$

$w_{n_{l+1}i}^l$

$\delta_{n_{l+1}}^{l+1}$
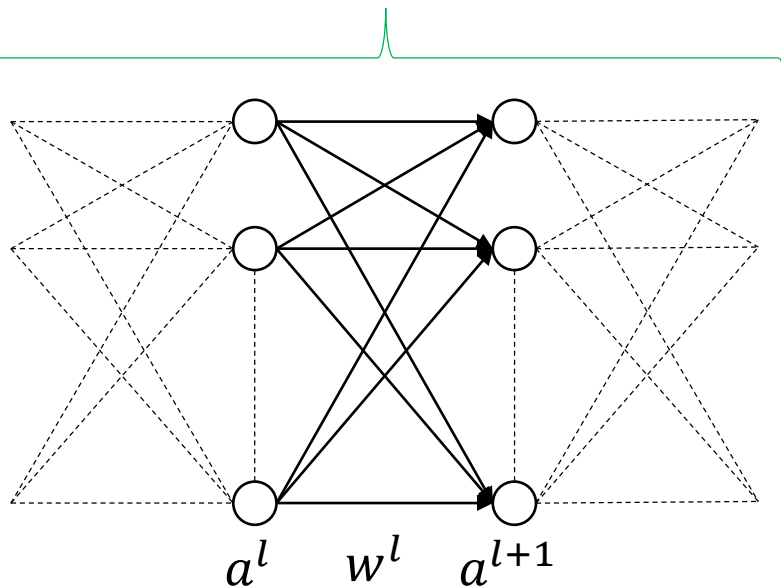
back propagation

$l$ layer

$l + 1$ layer

# Outline

■Brief Review of Neural Networks Structure

■Network Performance: Cost Function

■Steepest Gradient Method

■Backpropagation

■Three Pages to Understand BP

■<span style="color:red">Only One Page to Understand BP</span>
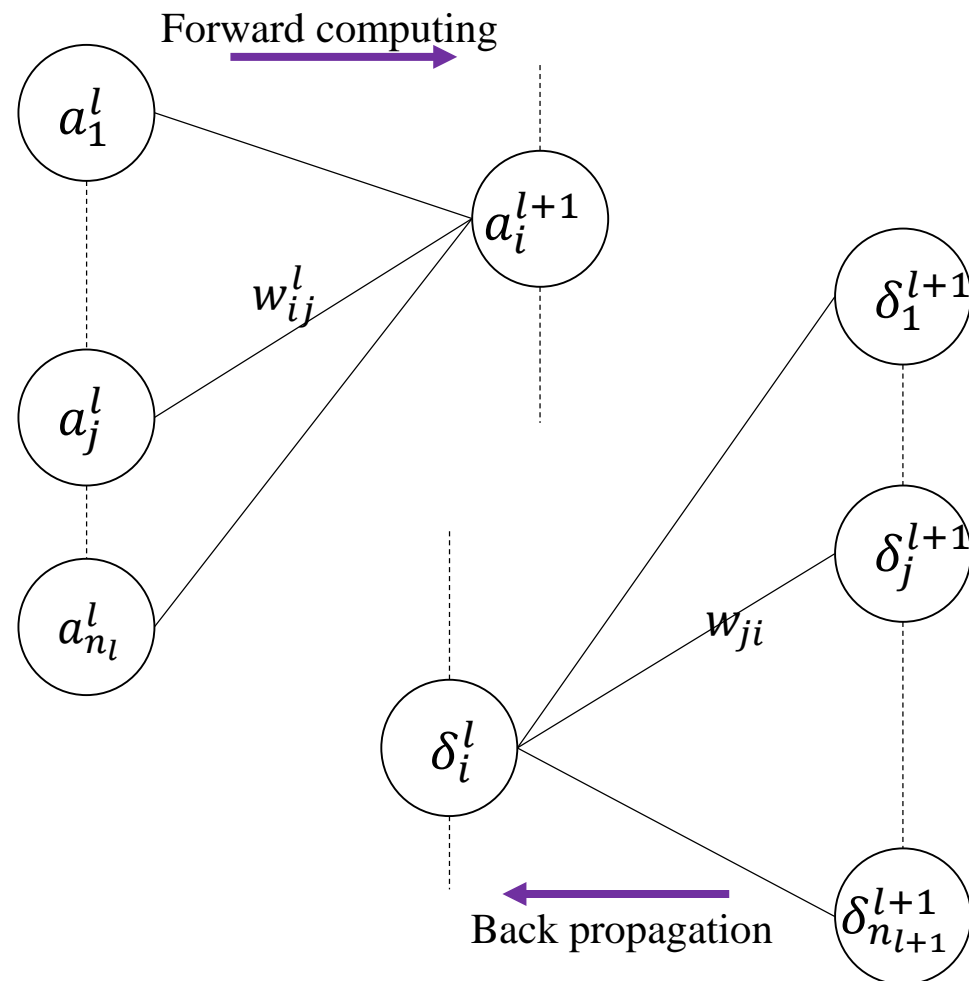
■The BP Algorithm

■Assignment

# Only One Pages to Understand BP

Cost function: $J(w^1, \cdots, w^L)$

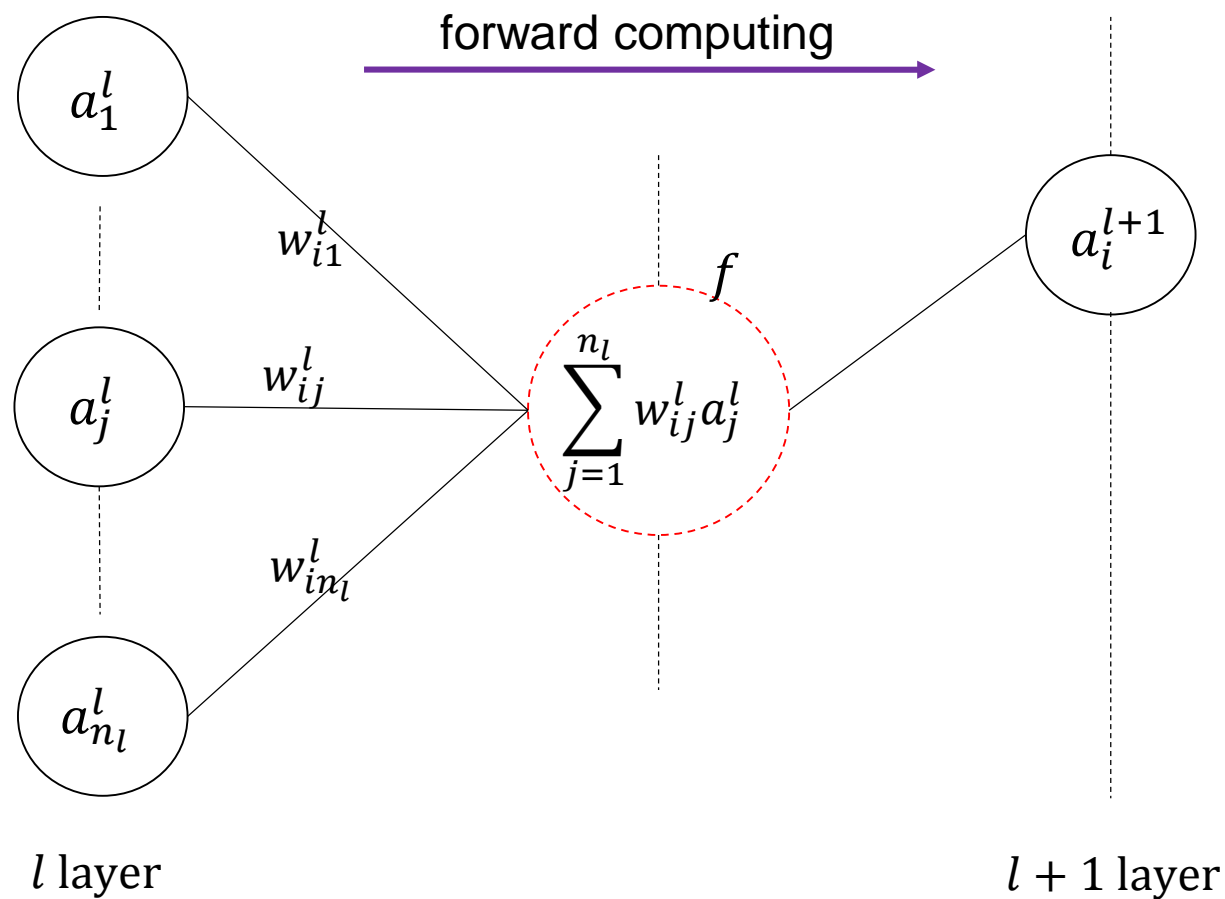Updating rule: $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \dfrac{\partial J}{\partial w_{ji}^l}$

Relationship: $\dfrac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$

forward computing $a^l$

$a^l \qquad w^l \qquad a^{l+1}$

backpropagation $\delta$

$l$ layer $i^{th}$ neuron

$$a_i^l = f(z_i^l)$$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l}$$

$$a_i^{l+1} = f\left( \sum_{j=1}^{n_l} w_{ij}^l a_j^l \right)$$

forward computing

$a_1^l$

$w_{i1}^l \qquad f$

$a_j^l \qquad w_{ij}^l \qquad \sum_{j=1}^{n_l} w_{ij}^l a_j^l \qquad a_i^{l+1}$

$w_{in_l}^l$

$a_{n_l}^l$

$l$ layer $\qquad\qquad l + 1$ layer

$\delta_1^{l+1}$

$w_{1i}^l$

$\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \qquad w_{ji}^l \qquad \delta_j^{l+1}$

$\dot{f}(z_i^l)$

$w_{n_{l+1}i}^l$

$\delta_i^l$

back propagation

$\delta_{n_{l+1}}^{l+1}$

$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$

$l$ layer $\qquad\qquad l + 1$ layer

32

# Outline

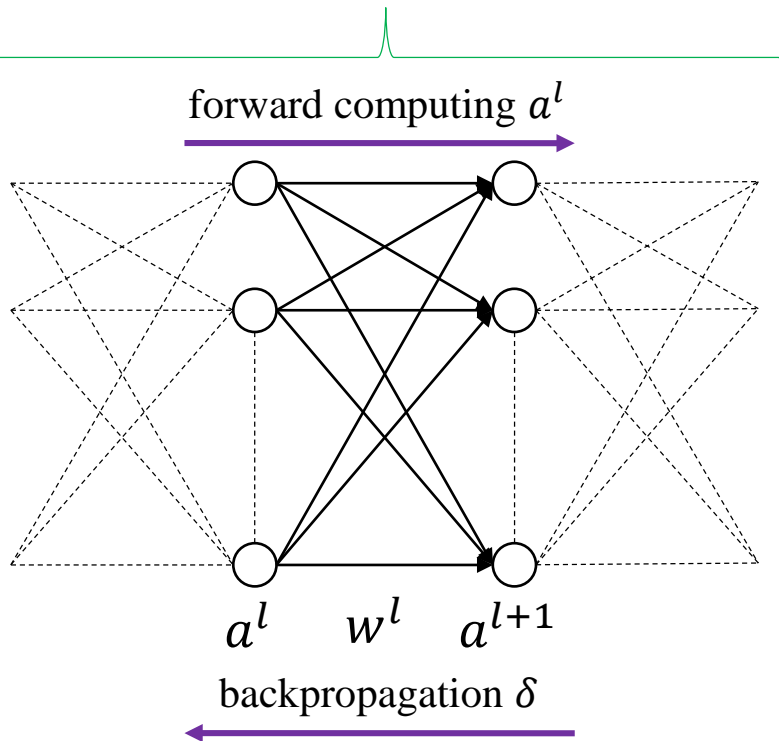■Brief Review of Neural Networks Structure

■Network Performance: Cost Function

■Steepest Gradient Method

■Backpropagation

■Three Pages to Understand BP

■Only One Page to Understand BP

■The BP Algorithm

■Assignment

# The BP Algorithm



Forward computing $a^l$

Back propagation $\delta^l$

Layer 1

Layer L

Forward computing $a^l$

layer $l$

layer $l + 1$

Backpropagation $\delta^l$

34

# The BP Algorithm

$$a_i^{l+1} = f\left(\sum_{j=1}^{n_l} w_{ij}^l a_j^l\right)$$

forward computing

$$a_i^{l+1} = f\left(\sum_{j=1}^{n_l} w_{ij}^l a_j^l\right)$$

$w_{i1}^l$

$f$

$\sum_{j=1}^{n_l} w_{ij}^l a_j^l$

$w_{ij}^l$

$a_j^l$

$w_{in_l}^l$

$a_{n_l}^l$

$a_1^l$

$a_i^{l+1}$

$l$ layer

$l+1$ layer

function $fc(W^l, a^l)$
$for\ i = 1: n_{l+1}$
$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l$$
$$a_i^{l+1} = f(z_i^{l+1})$$
end

function $bc(w^l, \delta^{l+1})$
$for\ i = 1: n_l$
$$\delta_i^l = \dot{f}(z_i^l) \cdot \left(\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}\right)$$
end

$\delta_1^{l+1}$

$w_{1i}^l$

$\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}$

$w_{ji}^l$

$\delta_j^{l+1}$

$\dot{f}(z_i^l)$

$w_{n_{l+1}i}^l$

$\delta_i^l$

$\delta_{n_{l+1}}^{l+1}$

back propagation

$l$ layer

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left(\sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}\right)$$

$l+1$ layer

35

# The BP Algorithm

The training data set
$$D = \{(x, y) | m \; samples\}$$

$x$: $input \; sample$
$y$: $target \; output$

There are two ways to train the network.

1. Online training: For each sample $(x, y) \in D$, define a cost function, for example, as
$$J(x, y) = \frac{1}{2} \sum_{j=1}^{n_L} (a_j^L - y_j^L)^2$$

2. Batch training: Define cost function as
$$J = \frac{1}{m} \sum_{(x,y) \in D} J(x, y)$$

Forward computing $a^l$

Layer 1      Back propagation $\delta^l$      Layer L

# The BP Algorithm

*Online BP Algorithm:*

Step 1. Input the training data set $D = \{(x, y)\}$

Step 2. Initial each $w_{ij}^l$, and choose a learning rate $\alpha$.

Step 3. Choose a sample $(x, y) \in D$, define $J(x, y)$, set $a^1 = x$

$\quad$ for $l = 1:L$

$\quad\quad fc(w^l, a^l);$

$\quad$ end

$\quad \delta^L = \dfrac{\partial J(x, y)}{\partial z^L};$

$\quad$ for $l = L - 1:1$

$\quad\quad bc(w^l, \delta^{l+1});$

$\quad$ end

Step 4. Updating

$\quad \dfrac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l;$

$\quad w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \dfrac{\partial J(x, y)}{\partial w_{ji}^l};$

Step 5. Return to Step 3 until each $w^l$ converge.

---

function $fc(w^l, a^l)$

$for\ i = 1:n_{l+1}$

$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l$$

$$a_i^{l+1} = f(z_i^{l+1})$$

$end$

---

Relationship:

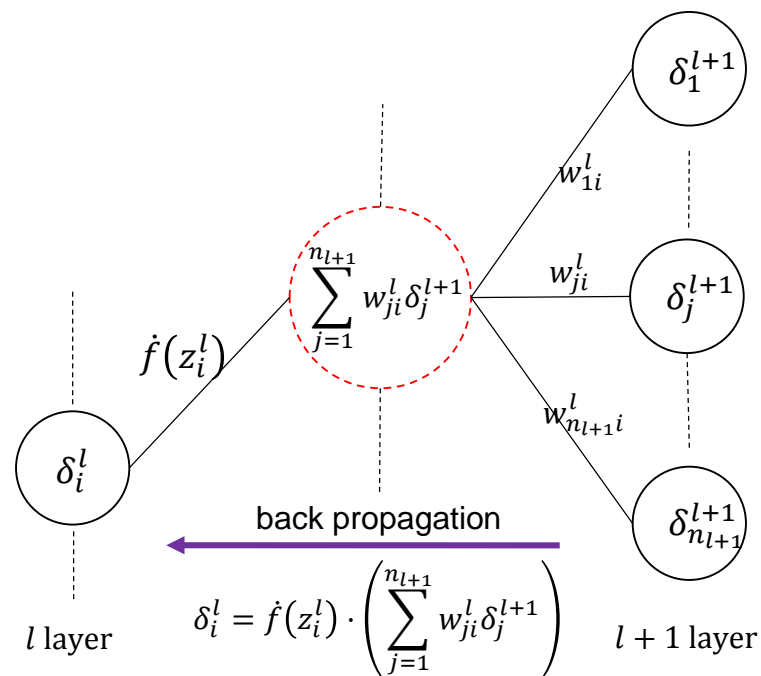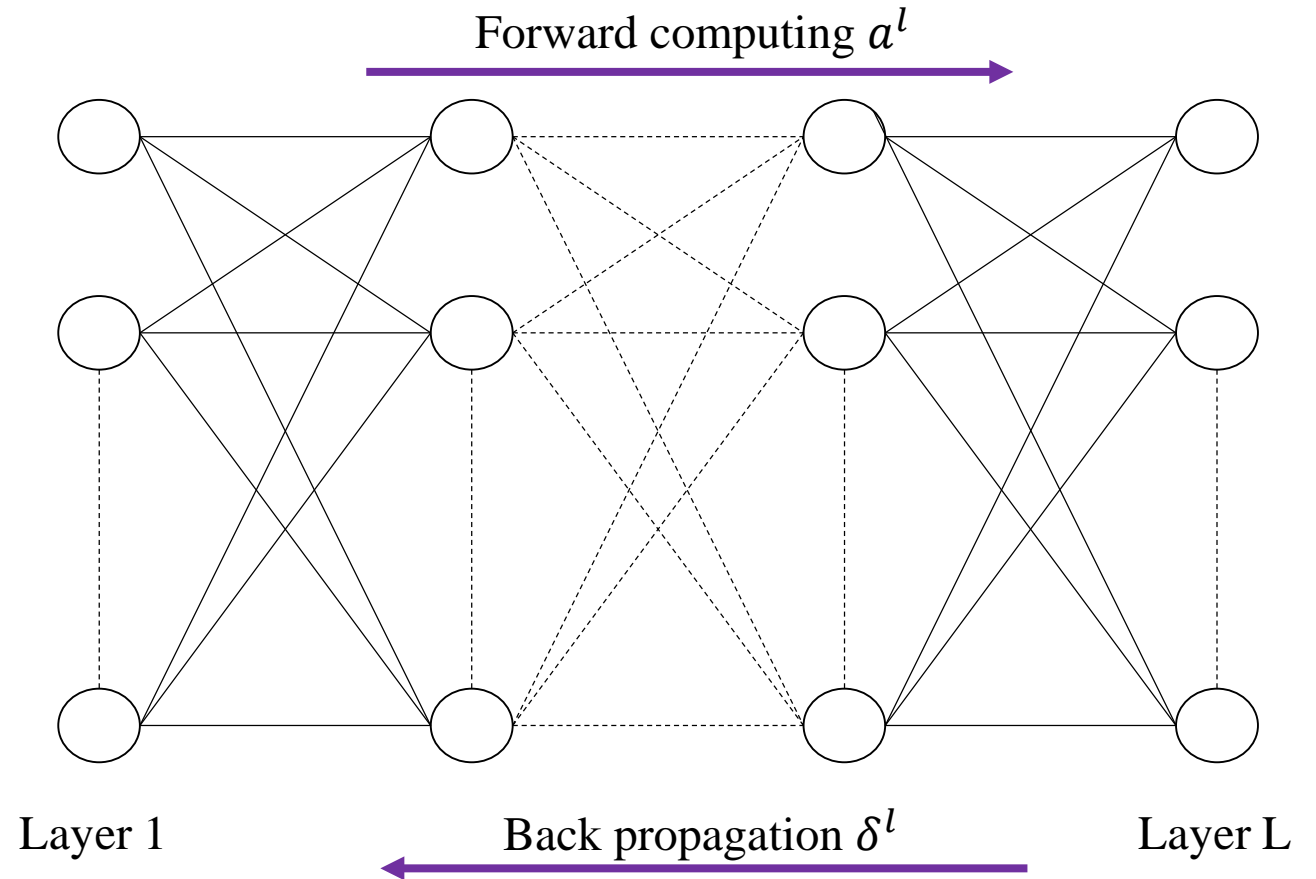$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

---

function $bc(w^l, \delta^{l+1})$

$for\ i = 1:n_l$

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$$

$end$

37

# The BP Algorithm

Step 1. Input the training data set $D = \{(x, y)\}$

Step 2. Initial each $w_{ij}^l$, and choose a learning rate $\alpha$.

Step 3. For each sample $(x, y) \in D$, set $a^1 = x$

   for $l = 1:L$

    $fc(w^l, a^l)$;

   end

   $\delta^L = \dfrac{\partial J}{\partial z^L}$;

   for $l = L - 1:1$

    $bc(w^l, \delta^{l+1})$;

   end

   $\dfrac{\partial J}{\partial w_{ji}^l} \leftarrow \dfrac{\partial J}{\partial w_{ji}^l} + \delta_j^{l+1} \cdot a_i^l$;

Step 4. Updating

   $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \dfrac{\partial J}{\partial w_{ji}^l}$;

Step 5. Return to Step 3 until each $w^l$ converge.

function $fc(w^l, a^l)$

$for \ i = 1:n_{l+1}$

$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l$$

$$a_i^{l+1} = f(z_i^{l+1})$$

$end$

Relationship:
$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

function $bc(w^l, \delta^{l+1})$

$for \ i = 1:n_l$

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$$

$end$

# Outline

■Brief Review of Neural Networks Structure

■Network Performance: Cost Function

■Steepest Gradient Method

■Backpropagation

■Three Pages to Understand BP

■Only One Page to Understand BP

■The BP Algorithm

■Assignment

# Assignment

$$
\text{function } fc(w^l, a^l) \\
for\ i = 1:n_{l+1} \\
z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l \\
a_i^{l+1} = f(z_i^{l+1}) \\
end
$$

$$
\text{function } bc(w^l, \delta^{l+1}) \\
for\ i = 1:n_l \\
\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right) \\
end
$$

Assignment: BP algorithms by MATLAB.

*Batch BP Algorithm:*

Step 1. Input the training data set $D = \{(x, y)\}$

Step 2. Initial each $w_{ij}^l$, and choose a learning rate $\alpha$.

Step 3. For each sample $(x, y) \in D$, set $a^1 = x$

        for $l = 1:L$

          $fc(w^l, a^l)$;

        end

        $\delta^L = \dfrac{\partial J}{\partial z^L}$;

        for $l = L - 1:1$

          $bc(w^l, \delta^{l+1})$;

        end

        $\dfrac{\partial J}{\partial w_{ji}^l} \leftarrow \dfrac{\partial J}{\partial w_{ji}^l} + \delta_j^{l+1} \cdot a_i^l$;

Step 4. Updating

        $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \dfrac{\partial J}{\partial w_{ji}^l}$;

Step 5. Return to Step 3 until each $w^l$ converge.

*Online BP Algorithm:*

Step 1. Input the training data set $D = \{(x, y)\}$

Step 2. Initial each $w_{ij}^l$, and choose a learning rate $\alpha$.

Step 3. Choose a sample $(x, y) \in D$, define $J(x, y)$, set $a^1 = x$

        for $l = 1:L$

          $fc(w^l, a^l)$;

        end

        $\delta^L = \dfrac{\partial J(x, y)}{\partial z^L}$;

        for $l = L - 1:1$

          $bc(w^l, \delta^{l+1})$;

        end

Step 4. Updating

        $\dfrac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$;

        $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \dfrac{\partial J(x, y)}{\partial w_{ji}^l}$;

Step 5. Return to Step 3 until each $w^l$ converge.

*Thanks*