# The Introduction To Artificial Intelligence

Yuni Zeng yunizeng@zstu.edu.cn
2022-2023-1

# The Introduction to Artificial Intelligence

- Part I Brief Introduction to AI & Different AI tribes
- Part II Knowledge Representation & Reasoning
- Part III AI GAMES and Searching
- Part IV Model Evaluation and Selection
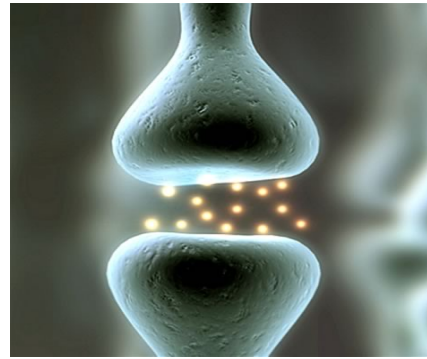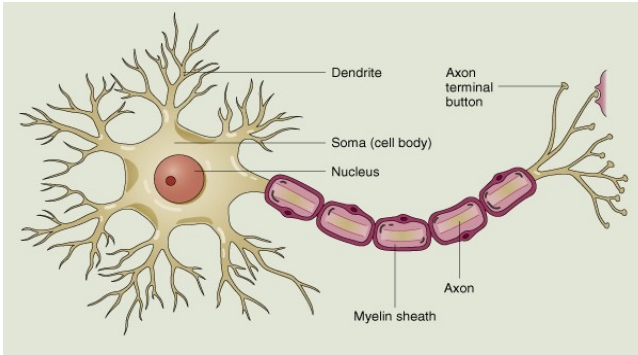- Part V Machine Learning
- Part VI Neural Networks

# Neural Networks

- *Brief review*

- Sequence Learning

# Brief review

☐ Artificial Neuron



Neuron output

$$a = f(z)$$

Neuron input

Connection weights

$w_1$

$x_1$

$w_i$

$x_i$

$w_n$

$x_n$

inputs

Total input

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

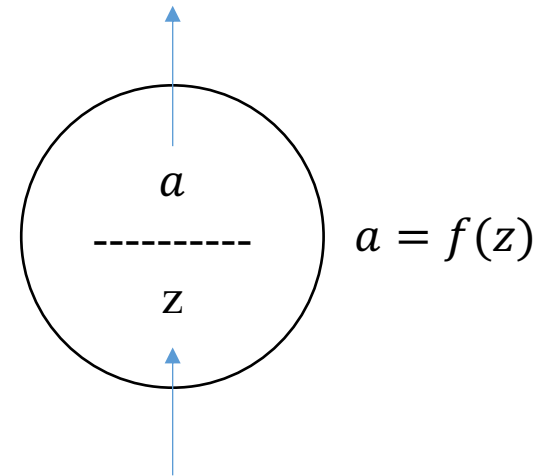Activation function

$$a = f(z)$$

Neuron output

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

4

# Computational Model of Neural Network

☐ Neural Networks

| Feedforward neural network |
| --- |

**=**

| neurons + feedforward connections |
| --- |



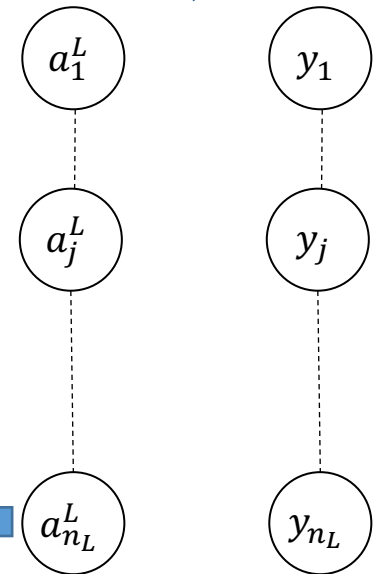| Recurrent neural network |
| --- |

**=**

| neurons + recurrent connections |
| --- |

# Steepest Descent Method

☐ Deep learning

Steepest Descent Algorithm:

$$w^{k+1} = w^k - \alpha_k \cdot \left.\frac{\partial F}{\partial w}\right|_{w^k}$$



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

Input

Output

$a_1^L$   $y_1$

$a_j^L$   $y_j$

| Updating weights | Computing gradient | Construct cost function | $a_{n_L}^L$ | $y_{n_L}$ |

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

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

Net output    Target output

# Backpropagation

□ Conclusion: BP for FNN

Forward computing: $y = f(\sum_{i=1}^{n} w_i x_i)$

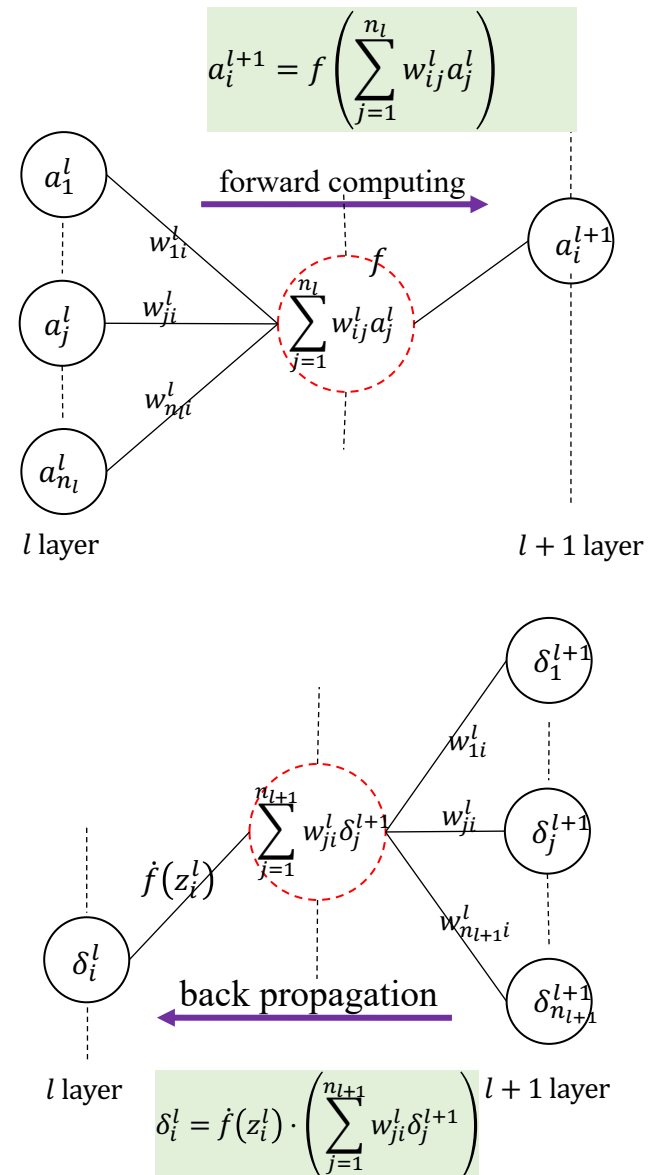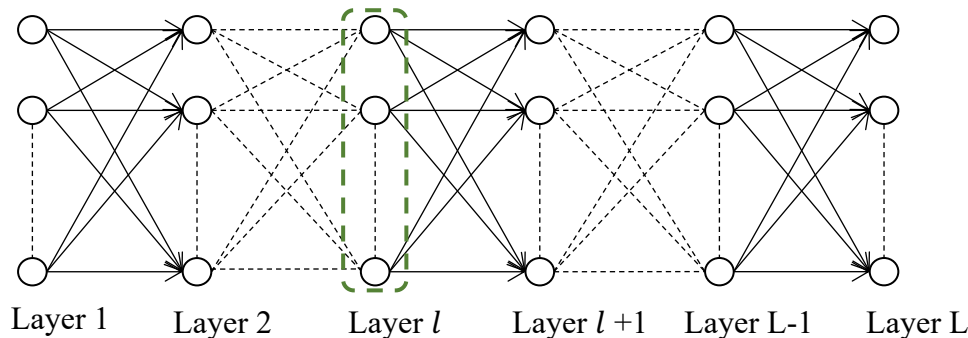Define cost function: $J = J(w^1, \cdots, w^{L-1})$

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

Define $\delta$: $\delta_i^l = \dfrac{\partial J}{\partial z_i^l}$

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

Back propagation: $\delta_i^L = \dfrac{\partial J}{\partial z_i^L} = \left(a_i^L - y_i^L\right) \cdot \dot{f}(z_i^L)$

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



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

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



forward computing

$l$ layer                    $l + 1$ layer



back propagation

$l$ layer                    $l + 1$ layer

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

7

# Neural Networks

- Brief review

- Sequence Learning

# Sequence Learning

□ A Sequence Recognizing Example

*Recognize A followed by B Problem*

The task is to recognize A followed by B.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
…………......

# Sequence Learning

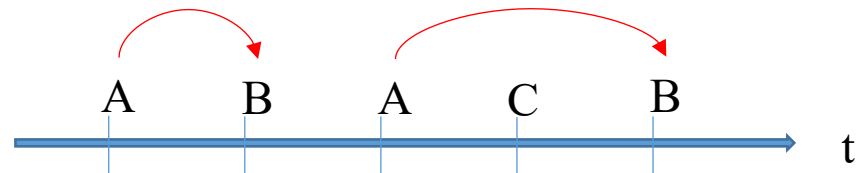□ **A Sequence Recognizing Example**

*Recognize A followed by B Problem*
The task is to recognize A followed by B.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
…………......



Problem: Can we use neural network to solve this problem?

# Sequence Learning

☐ A Sequence Recognizing Example

*Recognize A followed by B Problem*
The task is to recognize A followed by B.

The brain can solve this problem simply by memorizing the last A.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
…………......



0    1    0    0    1

A    B    A    C    B        t

Problem: Can we use neural network to solve this problem?
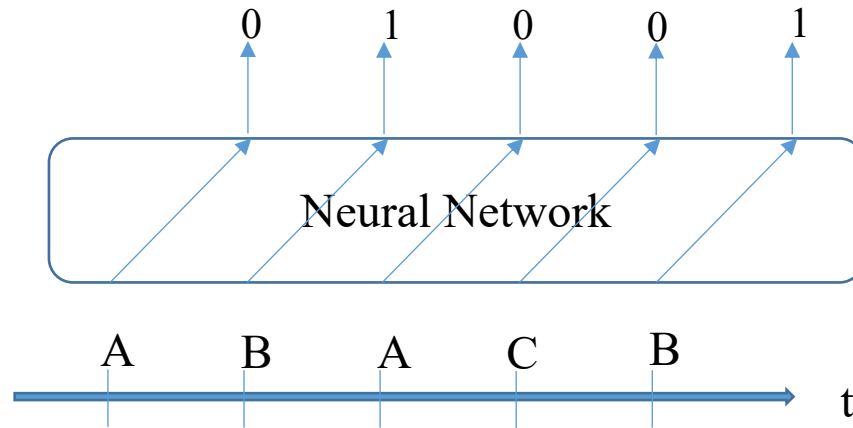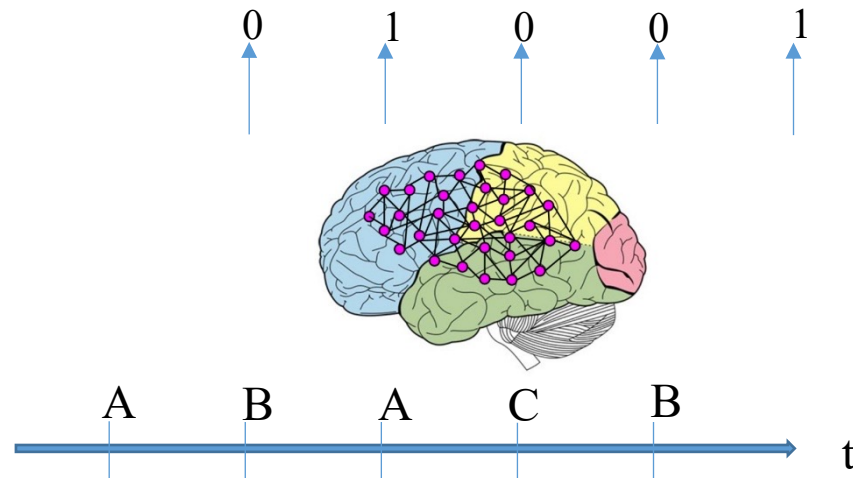
# Sequence Learning

☐ A Sequence Recognizing Example

*Recognize A followed by B Problem*

The task is to recognize A followed by B.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

……….......

# Sequence Learning

☐ A Sequence Recognizing Example



Mono-target output

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

In the last layer only

Mono-target output network cannot solve the sequence recognizing problem.
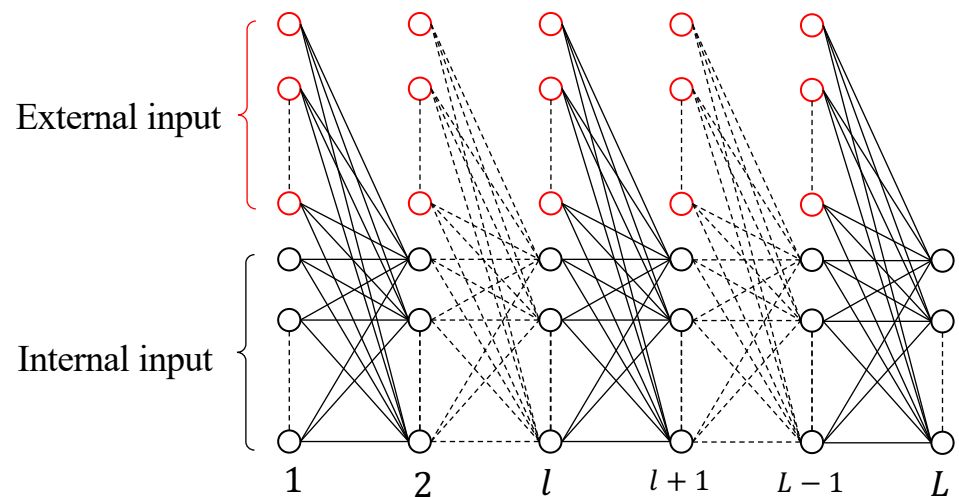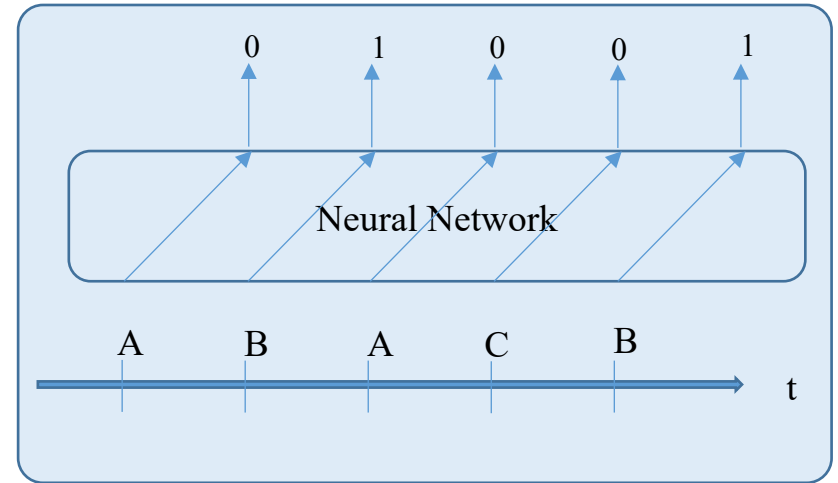
# Sequence Learning

☐ A Sequence Recognizing Example

*Recognize A followed by B Problem*
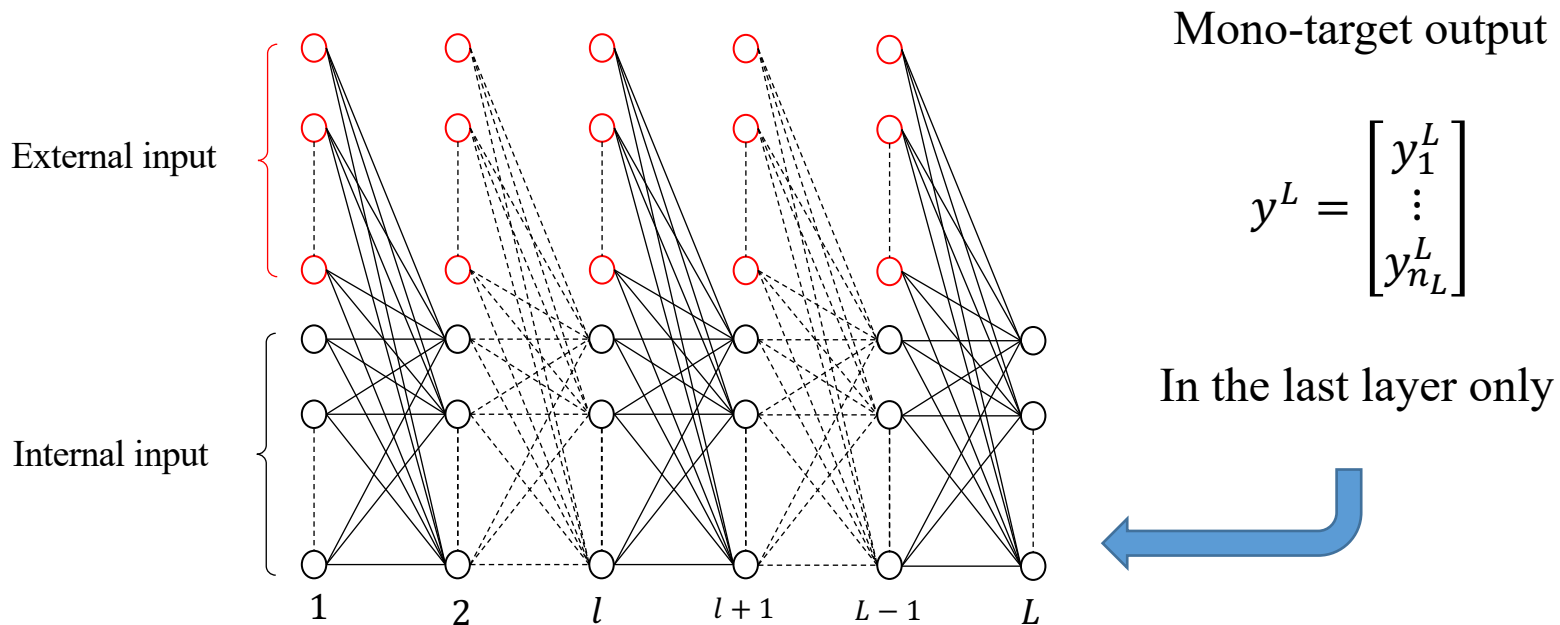The task is to recognize A followed by B.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
…………......

Multi-target outputs



Problem: Can we use neural network to solve this problem?

# Sequence Learning

□ A Sequence Recognizing Example



External input

Internal input

1   2   $l$   $l+1$   $L-1$   $L$

Multi-target outputs

Multi-target outputs

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

$(l = 2, \cdots, L)$

# Sequence Learning

□ A Sequence Recognizing Example



Multi-target outputs

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

$$(l = 2, \cdots, L)$$

Problem:
How to develop algorithm to train the network?

# Sequence Learning

- A Sequence Recognizing Example



External input

Internal input

Mono-target output

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

In the last layer only

We already developed BP Algorithm for mono-target output.

# Sequence Learning

□ A Sequence Recognizing Example



External input

Internal input

$y^2$  $y^l$  $y^{l+1}$  $y^{L-1}$  $y^L$

Multi-target outputs

*Problem:*
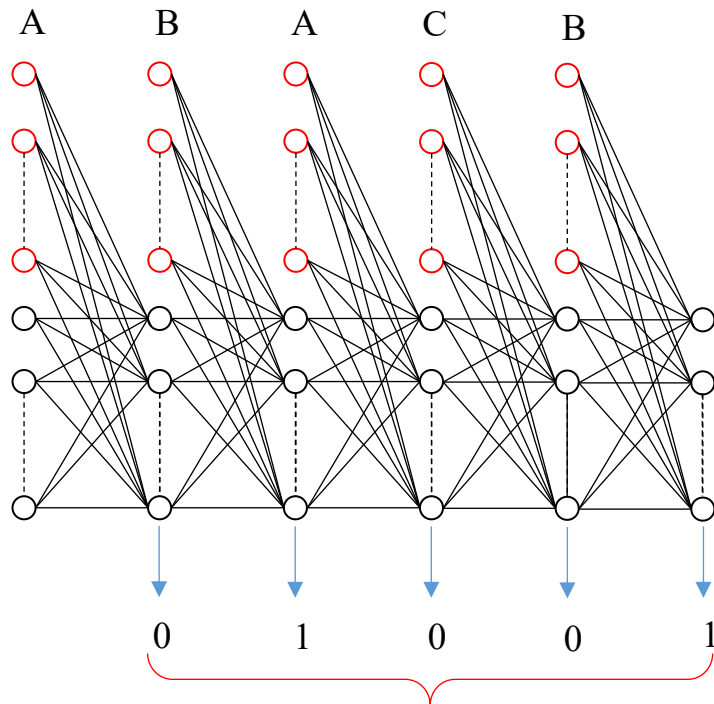Can we develop learning algorithms similar to BP for multi-target output?

Multi-target outputs

$$y^l = \begin{bmatrix} y^l_1 \\ \vdots \\ y^l_{n_L} \end{bmatrix}$$

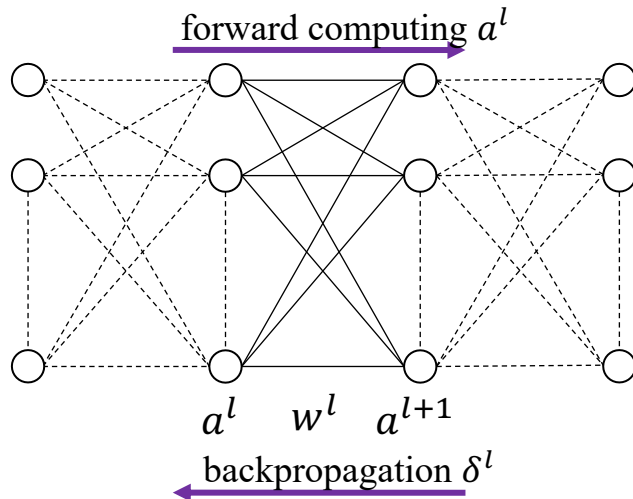$(l = 2, \cdots, L)$

# Sequence Learning

□ Review of BP algorithm for mono-output NNs
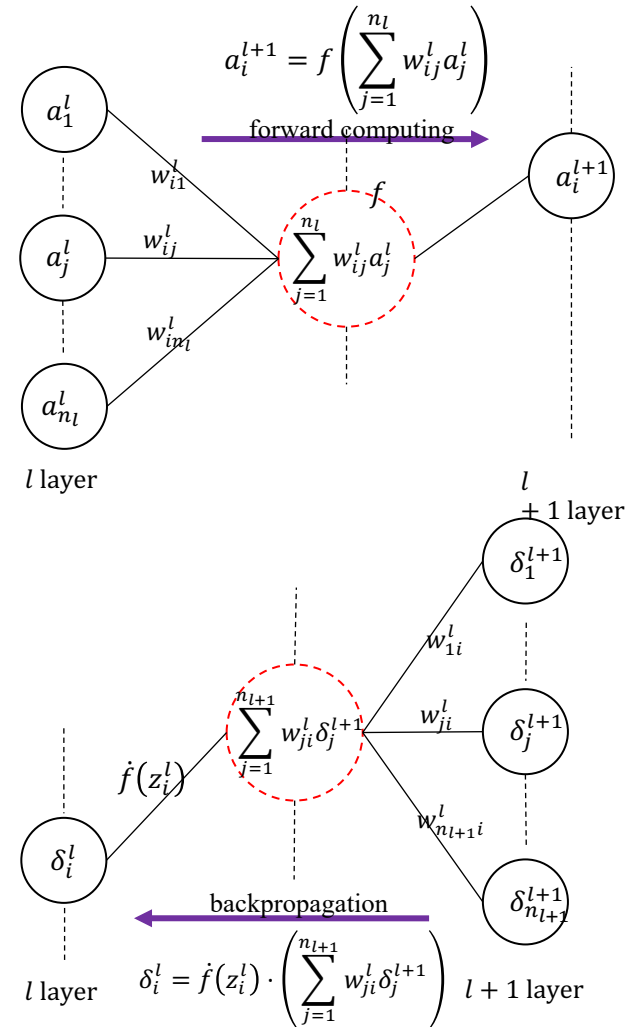
Cost function: $J(w^1, \cdots, w^{L-1})$

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

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

forward computing $a^l$

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

backpropagation $\delta^l$

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

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

Global     Local

$J$       $f$

$\delta_i^l \longleftrightarrow z_i^l \longleftrightarrow a_i^l$

Bridge

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

forward computing

$a_1^l$

$w_{i1}^l$

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

$a_i^{l+1}$

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

$w_{in_l}^l$

$a_{n_l}^l$

$l$ layer

$l + 1$ layer

$\delta_1^{l+1}$

$w_{1i}^l$

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

$w_{ji}^l$   $\delta_j^{l+1}$

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

$\delta_i^l$

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

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

backpropagation

$l$ 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

# Sequence Learning

□ Review of BP algorithm for mono-output NNs



Cost function

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L} J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

Network Outputs    Multi-target outputs

$$a^l = \begin{bmatrix} a_1^l \\ \vdots \\ a_{n_l}^l \end{bmatrix} \qquad y^l = \begin{bmatrix} y_1^l \\ \vdots \\ y_{n_l}^l \end{bmatrix}$$

$(l = 2, \cdots, L)$        $(l = 2, \cdots, L)$

# Sequence Learning

☐ Review of BP algorithm for mono-output NNs

Steepest Descent Method

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L} J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

1. Computing

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

2. Iterating

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

How to calculate me?

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

$J$

$J^2 \qquad J^l \qquad J^{l+1} \qquad J^{L-1} \qquad J^L$

$y^2 \qquad y^l \qquad y^{l+1} \qquad y^{L-1} \qquad y^L$

# Sequence Learning

☐ Review of BP algorithm for mono-output NNs

$l$ layer

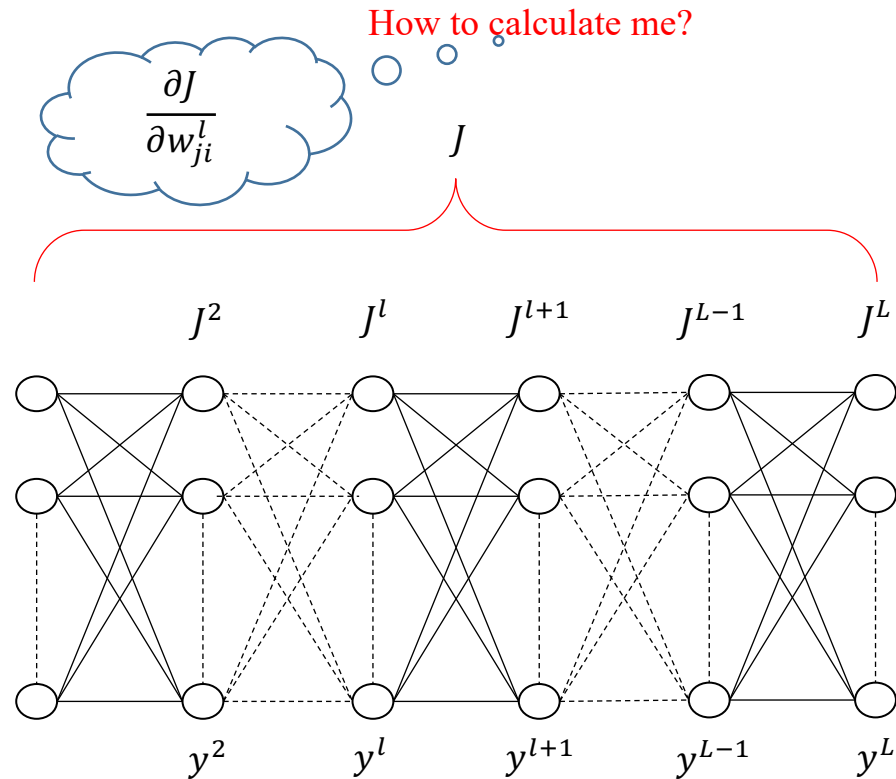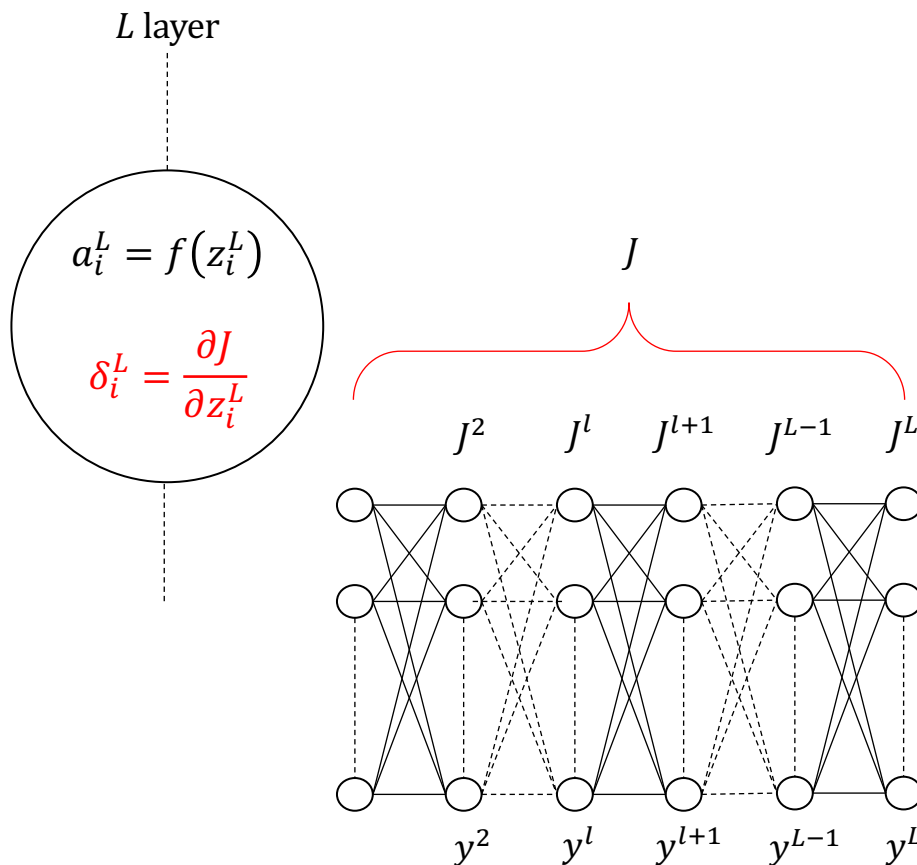$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L} J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

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

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

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

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

Why?

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

$l + 1$ layer

$J$

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

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

Problem: Can we back propagate $\delta^l$?

$l$ layer

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

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

$w_{ji}^l$

$$\delta_j^{l+1} = \frac{\partial J}{\partial z_j^{l+1}}$$

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

# Sequence Learning

☐ Step 1: Calculating $\delta^{\wedge}L$ in Last Layer

$L$ layer

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

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

$J$

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



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

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L} J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

It holds that,

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

☐ Step 2: Relation Between $\delta^l$ and $\delta^{(l+1)}$

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L}J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

$J$ may have an explicit dependence on , it may also have an implicit dependence on through later output values. To avoid ambiguity in interpreting partial derivatives, define $z_i^l(*) = z_i^l$.

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}}\frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

☐ Step 2: Relation Between $\delta^{\wedge}l$ and $\delta^{\wedge}(l+1)$

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}(a_i^l - y_i^l)^2 , (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L}J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}(a_i^l - y_i^l)^2$$

An Illustrate Example

$J$ may have an explicit dependence on , it may also have an implicit dependence on through later output values. To avoid ambiguity in interpreting partial derivatives, define $z_i^l(*) = z_i^l$.

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)}\cdot\frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}}\frac{\partial J}{\partial z_j^{l+1}}\cdot\frac{\partial z_j^{l+1}}{\partial z_i^l}$$

$$J = x + y, \; y = \exp(x)$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial x} + \frac{\partial J}{\partial y}\cdot\frac{\partial y}{\partial x}$$

$$x^* = x$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial x^*}\cdot\frac{\partial x^*}{\partial x} + \frac{\partial J}{\partial y}\cdot\frac{\partial y}{\partial x}$$
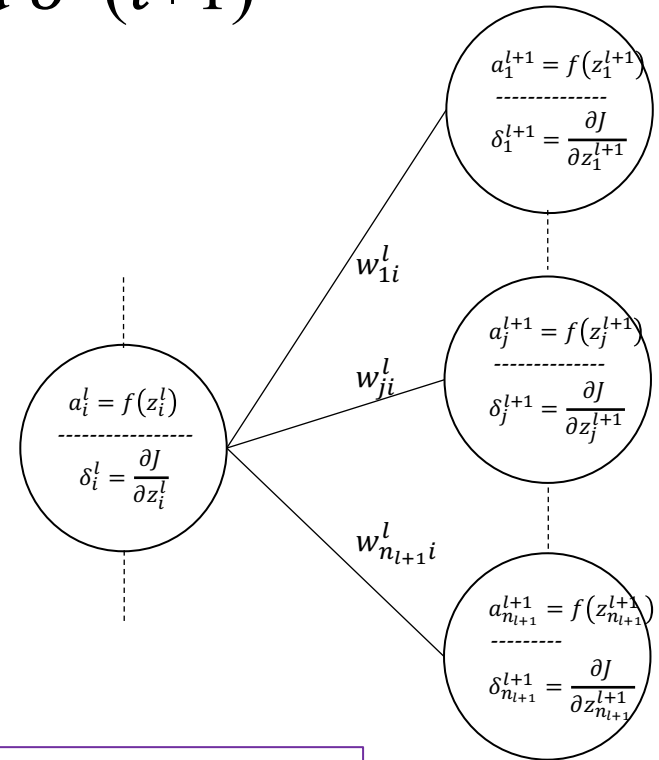
# Sequence Learning

□ Step 2: Relation Between $\delta^{\wedge}l$ and $\delta^{\wedge}(l+1)$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

$$\frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} = \frac{\partial J^l}{\partial z_i^l} = (a_i^l - y_i^l) \cdot \dot{f}(z_i^l)$$

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

$$\delta_i^l = \dot{f}(z_i^l) \cdot \left[ (a_i^l - y_i^l) + \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right) \right]$$
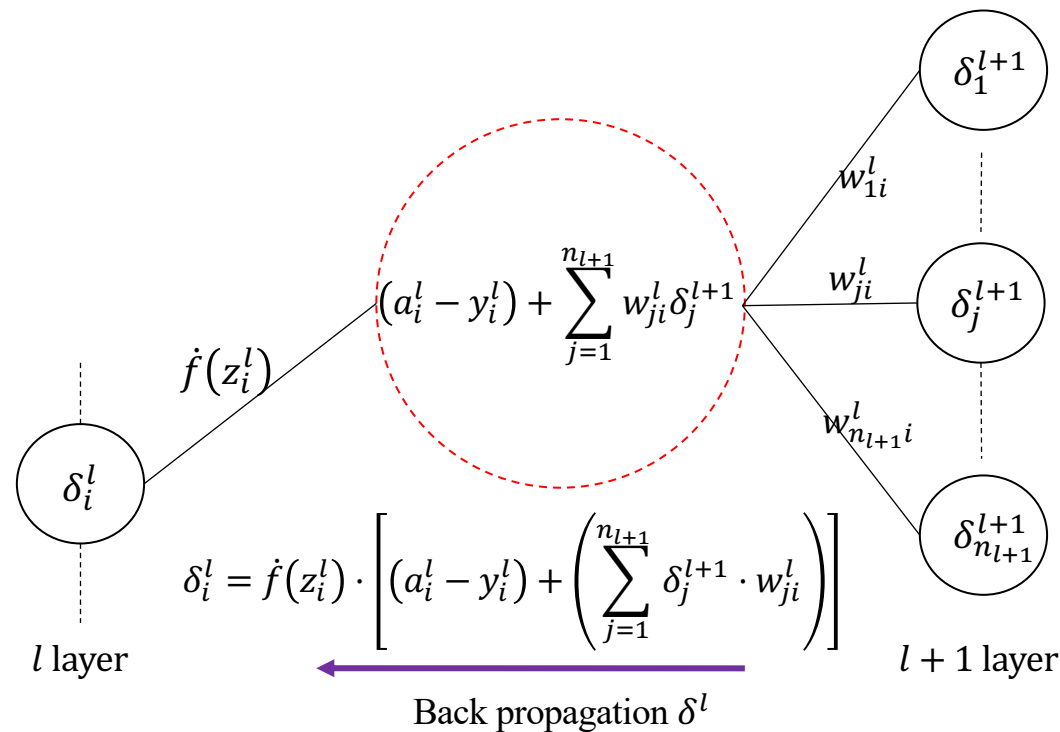
$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \cdots, L)$$

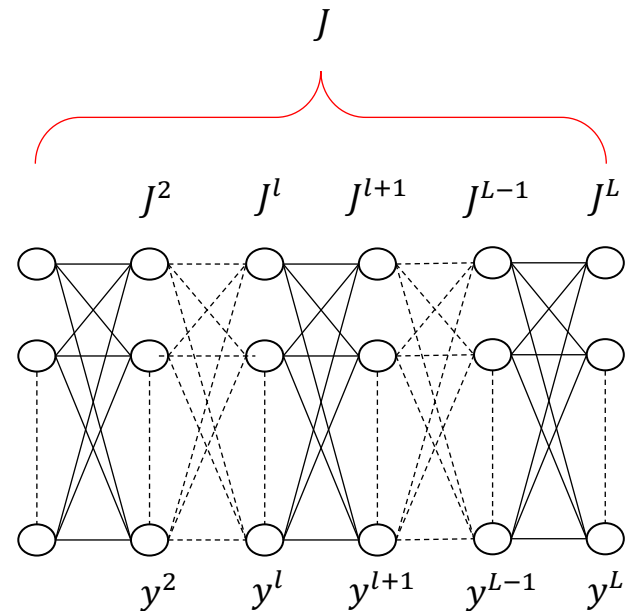$$J = \sum_{l=2}^{L} J^l = \frac{1}{2} \sum_{l=2}^{L} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

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

$$w_{1i}^l$$

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

$$w_{ji}^l$$

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

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

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

# Sequence Learning

☐ Step 3: Backpropagation $\delta$

$$J^l = \frac{1}{2}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2, (l = 2, \cdots, L)$$

$$J = \sum_{l=2}^{L} J^l = \frac{1}{2}\sum_{l=2}^{L}\sum_{i=1}^{n_l}\left(a_i^l - y_i^l\right)^2$$

$\delta_1^{l+1}$

$w_{1i}^l$

$(a_i^l - y_i^l) + \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1}$

$w_{ji}^l$

$\delta_j^{l+1}$

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

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

$\delta_i^l$

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

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

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

Back propagation $\delta^l$

$J$

$J^2 \qquad J^l \qquad J^{l+1} \qquad J^{L-1} \qquad J^L$

$y^2 \qquad y^l \qquad y^{l+1} \qquad y^{L-1} \qquad y^L$

# Sequence Learning

□ **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 mini-batch sample $D_m \subseteq D$

$\quad a^1 \leftarrow x \in D_m$;

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

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

$\quad\quad$ end

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

$\quad\quad$ for $l = L - 1: 2$

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

$\quad\quad$ end

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

Step 4. Updating

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

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



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

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

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

# Sequence Learning

☐ Illustrative Example

*Recognize A followed by B Problem*
The task is to recognize A followed by B.

Generated Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

…………........

# Sequence Learning

□ Illustrative Example

$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad C = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$
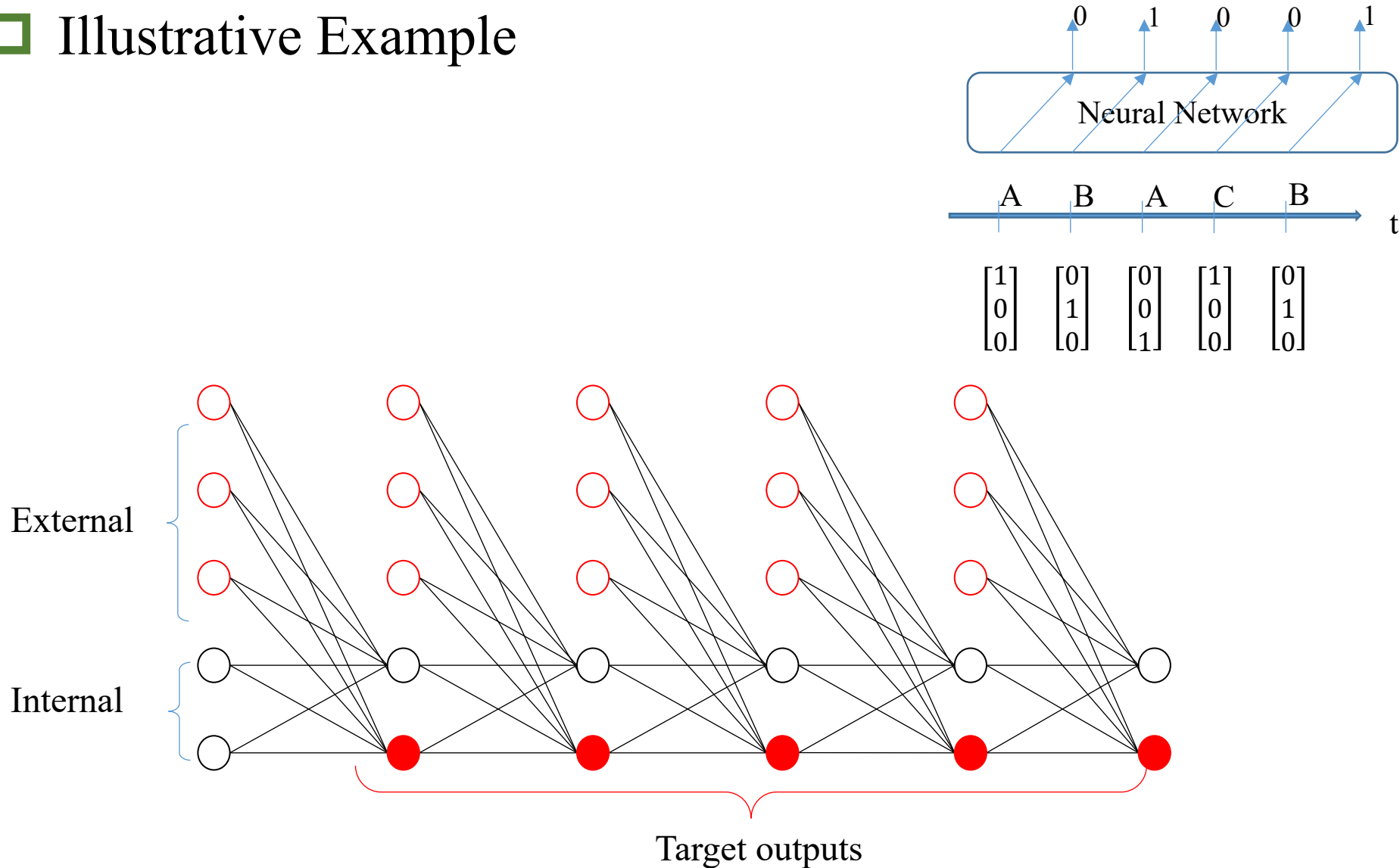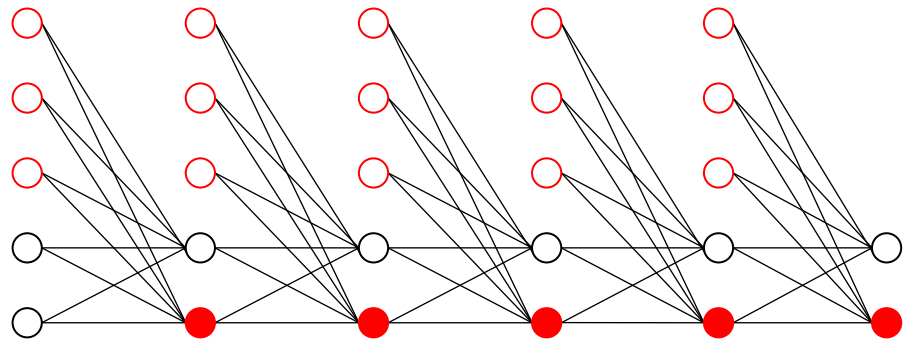
Generated Sequences

1.     A   B   C   A   B

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

2.     C   A   C   C   B

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

     0      1      0      0      1

Neural Network

   A      B      A      C      B

t

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

# Sequence Learning

☐ Illustrative Example



Target outputs

# Sequence Learning

☐ Illustrative Example

Generated Training Sequences
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
…………......



Generated Testing Sequences
1. CBCAC
2. ACBBA
3. BACCB
4. ACBCB
5. AACBC
6. BAACB
7. AAACB
8. CCBAB
9. BBCAB
10. AABAC
…………......

# Sequence Learning

□ Illustrative Example

$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad C = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Example outputs:

| | | | | | | |
|---|---|---|---|---|---|---|
| CAAC**B** | → | 0.0184 | 0.0001 | 0.0211 | 0.0801 | 0.9928 |
| A**B**BCA | → | 0.0179 | 0.9375 | 0.0267 | 0.0012 | 0.0000 |
| AAC**B**A | → | 0.0179 | 0.0336 | 0.0286 | 0.8722 | 0.0000 |
| CAC**B**B | → | 0.0184 | 0.0001 | 0.0170 | 0.8494 | 0.0013 |
| BCAAA | → | 0.0182 | 0.0001 | 0.0001 | 0.0622 | 0.0018 |

# Sequence Learning

□ Another Example: Image Caption

*Image Caption:*
The task is to describe the content of an image using properly formed English sentence.



image          Neural Network

English sentence    a    cute  panda  lies    on    a    tree

# Sequence Learning

□ Another Example: Image Caption

Dataset：COCO

COCO is a new image recognition, segmentation, and captioning dataset sponsored by Microsoft.
http://mscoco.org/dataset/#download
There are：
    80,000 training samples
    40,000 validation samples
    40,000 test samples



Dataset

$$D = \{(x, s_1, s_2, s_3, s_4, s_5)\}$$

1. Two person drive a small race car .
2. Two racer drive a white bike down a road .
3. Two motorist be ride along on their vehicle that be oddly design and color .
4. Two person be in a small race car drive by a green hill .
5. Two person in race uniform in a street car .

# Sequence Learning

□ Another Example: Image Caption

# Sequence Learning

☐ Another Example: Image Caption

Coding the Inputs：

1. building the vocabulary

End of sentence

| a | an | two | car | EOS | ⋯ |

One-hot word vector

2. coding the sentence

$s =$ two  person  drive  a  small  race  car  EOS

3. digitizing the image

# Sequence Learning

☐ Another Example: Image Caption

# Sequence Learning

☐ Other RNN

# Sequence Learning

☐ Other RNN



长期依赖(Long Term Dependencies)

eg1: The cat, which already ate a bunch of food, was full.
eg2: The cats, which already ate a bunch of food, were full.

# Sequence Learning

□ Other RNN - LSTM

LSTM: Long Short Term Memory，顾名思义，它具有记忆长短期信息的能力的神经网络

# Sequence Learning

□ Other RNN - LSTM



LSTM解决RNN的长期依赖问题:LSTM引入了门（gate）机制
用于控制特征的流通和损失

# Sequence Learning

☐ Other RNN - LSTM



$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \; + \; b_f \right)$$

$f_t$ 遗忘门

# Sequence Learning

☐ Other RNN - LSTM



$$i_t = \sigma \left( W_i \cdot [h_{t-1}, x_t] \ + \ b_i \right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \ + \ b_C)$$

# Sequence Learning

□ Other RNN - LSTM



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
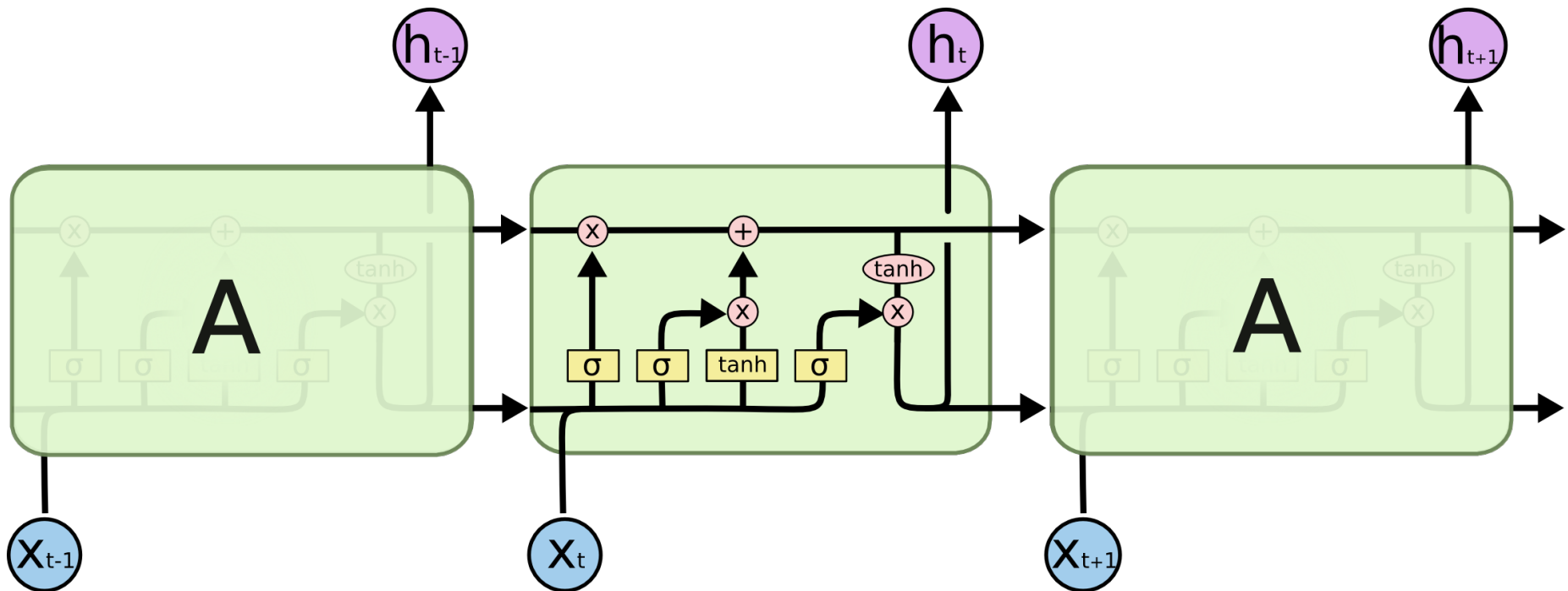
# Sequence Learning

☐ Other RNN - LSTM



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

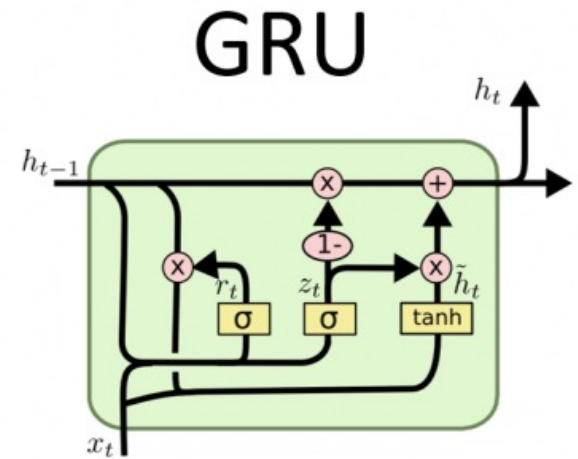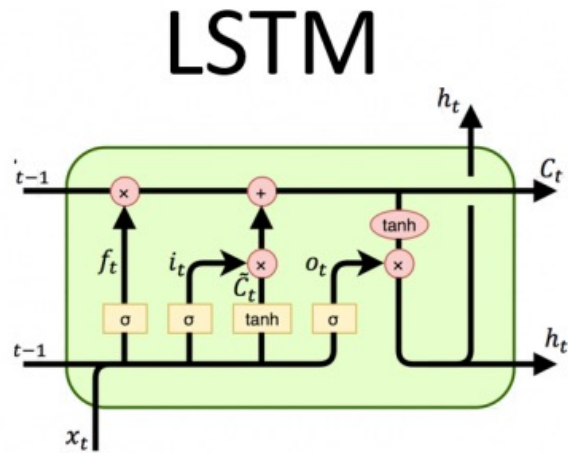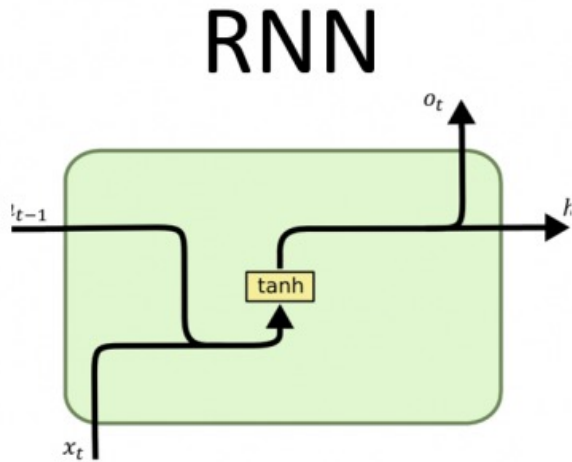$$h_t = o_t * \tanh \left( C_t \right)$$

# Sequence Learning

☐ Other RNN - LSTM
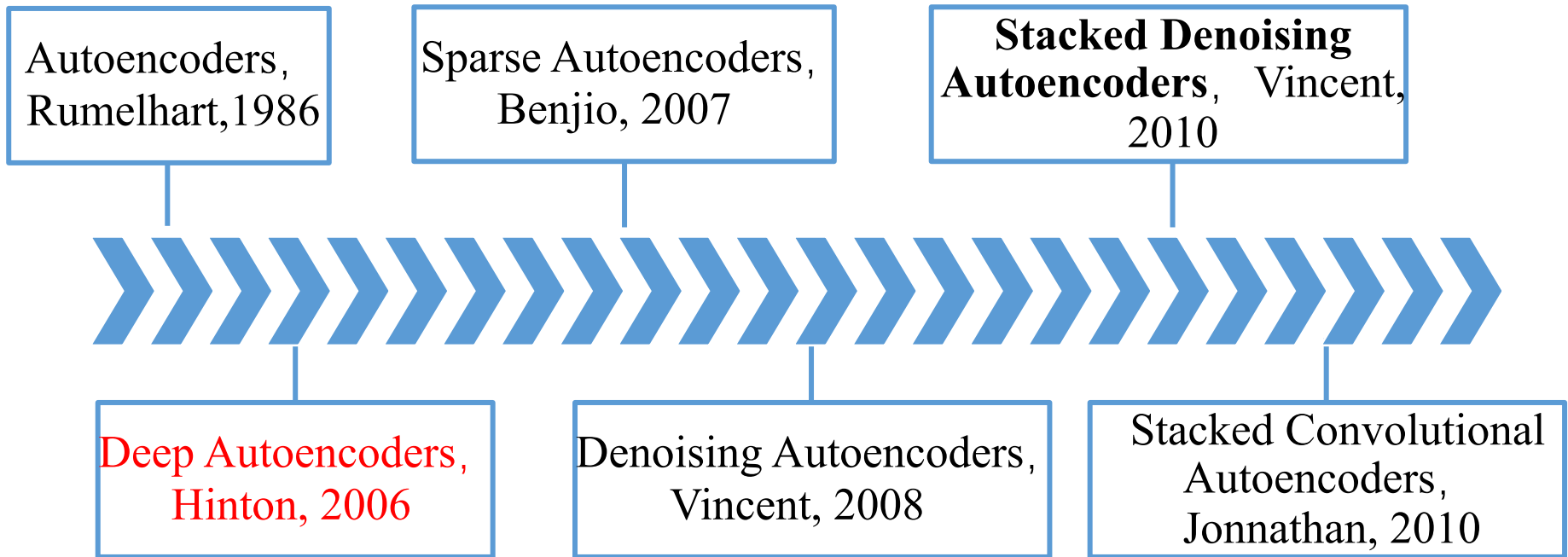
# Sequence Learning

☐ Other RNNs

# Neural Networks

- Brief review

- Sequence Learning

- *Representation learning*

# Representation learning

□ Autoencoders

# Representation learning

- Autoencoders



**Reducing the Dimensionality of Data with Neural Networks**

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such "autoencoder" networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer "encoder" network
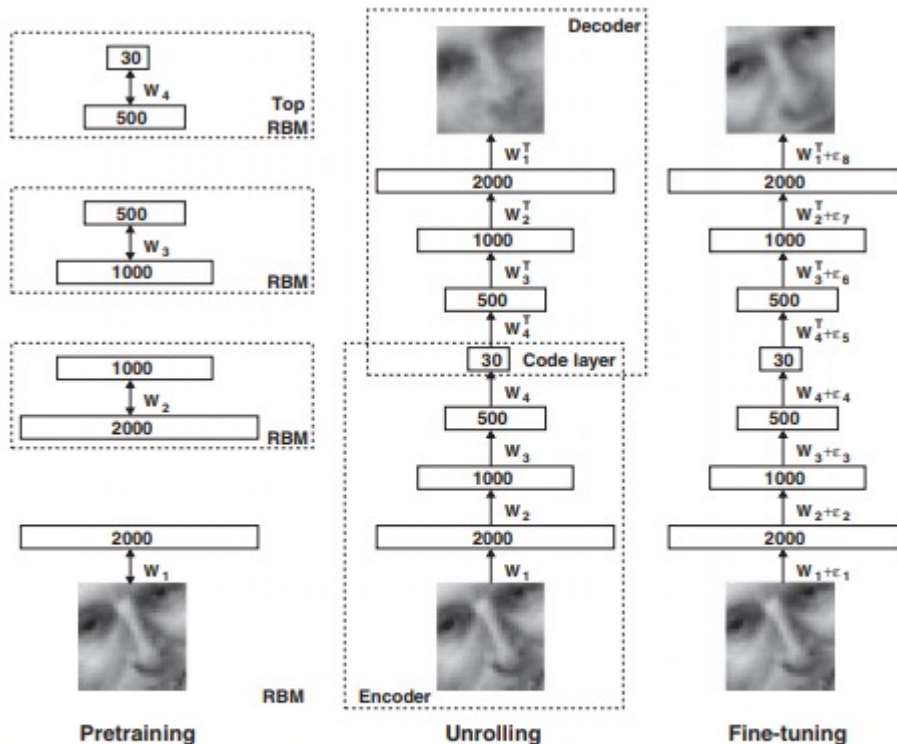
**Fig. 1.** Pretraining consists of learning a stack of restricted Boltzmann machines (RBMs), each having only one layer of feature detectors. The learned feature activations of one RBM are used as the "data" for training the next RBM in the stack. After the pretraining, the RBMs are "unrolled" to create a deep autoencoder, which is then fine-tuned using backpropagation of error derivatives.

# Representation learning

☐ Autoencoders

**Encoder**

$d' < d \longrightarrow$ **dimensionality reduction**



Input

Hidden

Output

$\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$

**Code**

$y = f_\theta(x)$
$\in [0,1]^{d'}$

$x \in [0,1]^d$

$z = s_{\theta'}(f_\theta(x)) \in [0,1]^d$

# Representation learning

☐ Autoencoders

$d' < d \longrightarrow$ **dimensionality reduction**

**Encoder**

**Decoder**

Input

Hidden

Output

$\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$

$\boldsymbol{\theta}' = \{\boldsymbol{W}', \boldsymbol{b}'\}$
$\boldsymbol{W}' = \boldsymbol{W}^T$

**Code**

$y = f_\theta(x)$

$z = s_{\theta'}(f_\theta(x)) \in [0,1]^d$

# Representation learning

☐ Autoencoders

**Encoder**                                    **Decoder**



Input

Hidden

Output

$\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$

$\boldsymbol{\theta}' = \{\boldsymbol{W}', \boldsymbol{b}'\}$
$\boldsymbol{W}' = \boldsymbol{W}^T$

**Code**

$y = f_\theta(x)$

$\theta *, \theta' *$
$= \underset{\theta, \theta'}{\mathrm{argmin}} L_H(x, z)$

$z = s_{\theta'}(f_\theta(x)) \in [0,1]^d$

# Representation learning

☐ Denoising Autoencoders



Input        Hidden        Output

$\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$

$\boldsymbol{\theta}' = \{\boldsymbol{W}', \boldsymbol{b}'\}$
$\boldsymbol{W}' = \boldsymbol{W^T}$

**Code**

$y = f_\theta(x)$

$$\theta, \theta'^* = \underset{\theta, \theta'}{\operatorname{argmin}} L_H(x, z)$$

$x \longrightarrow \tilde{x}$
corrupted

$z \in [0,1]^d$

# Representation learning

□ Denoising Autoencoders



**Layer 1**

Input    Hidden    Output

$\theta = \{W, b\}$    $\theta' = \{W', b'\}$
$W' = W^T$

**Code**

$y = f_\theta(x)$

$x \xrightarrow{\text{corrupted}} \tilde{x}$    $z \in [0,1]^d$

# Representation learning

☐ Denoising Autoencoders

Representation

$$\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$$

**Code**

$y$

# Representation learning

☐ Denoising Autoencoders

Representation

$\theta_1$
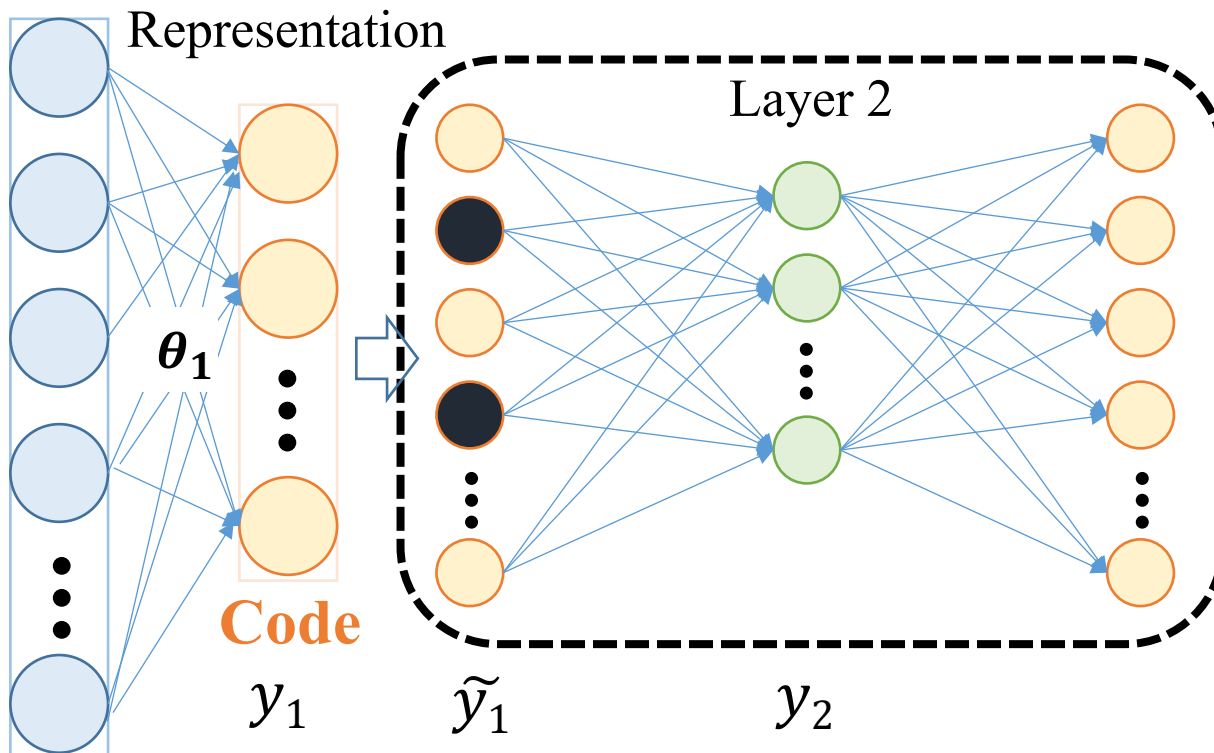
**Code**

$y_1$

$x$

Advantages of deep architectures:
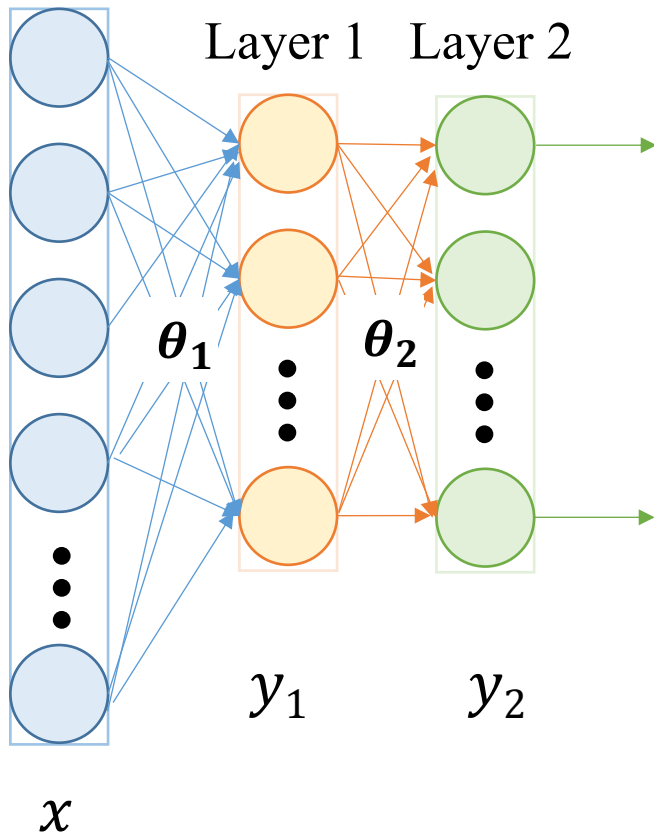  ▫ Re-use of features
  ▫ More abstract features

# Representation learning
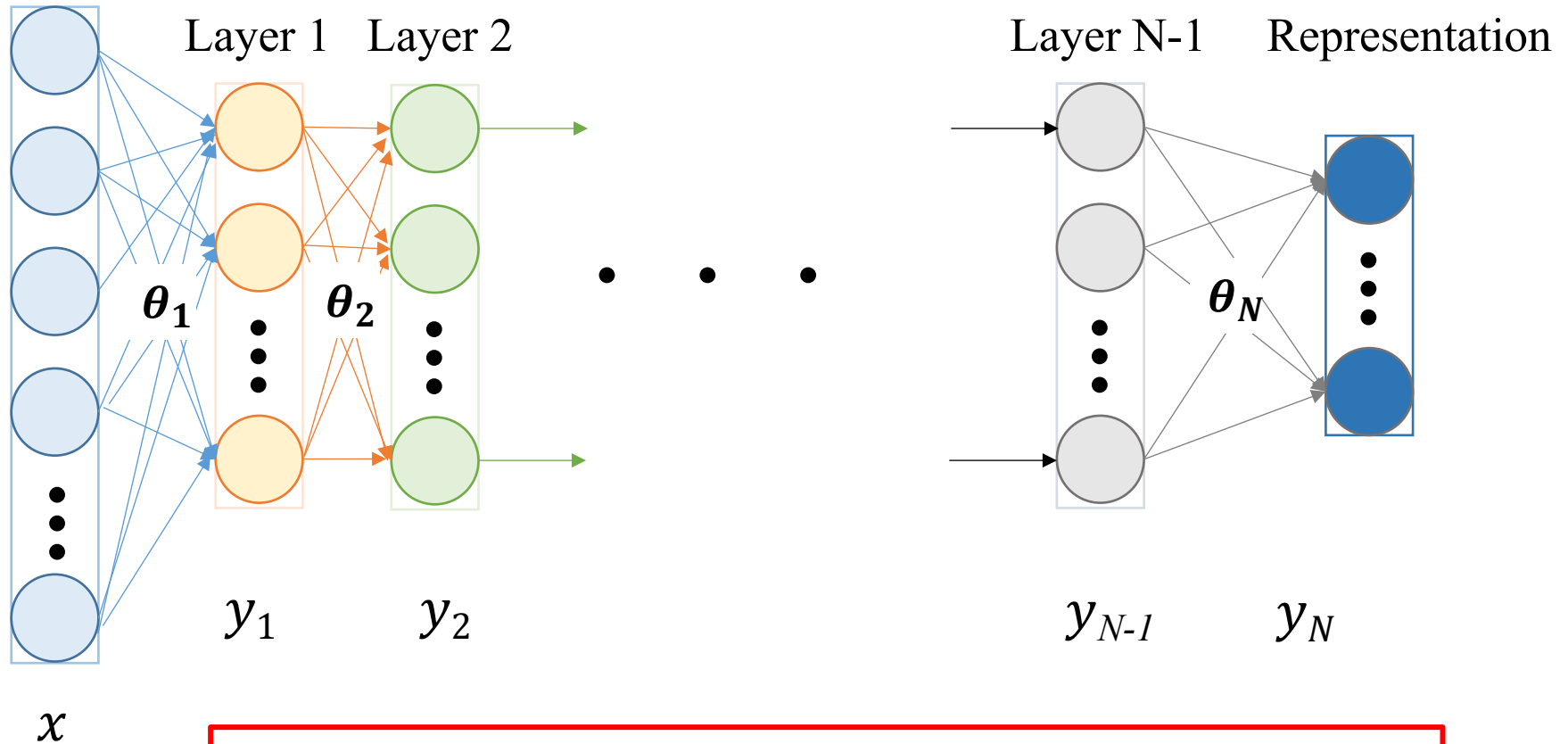
☐ Stacked Denoising Autoencoders

# Representation learning

□ Stacked Denoising Autoencoders

# Representation learning

□ Stacked Denoising Autoencoders



Each layer of the network is trained to produce a higher- level representation of the observed patterns.