# Understanding Deep Neural Networks

## Chapter Four

# BP: An Illustrating Example

Zhang Yi, *IEEE Fellow*
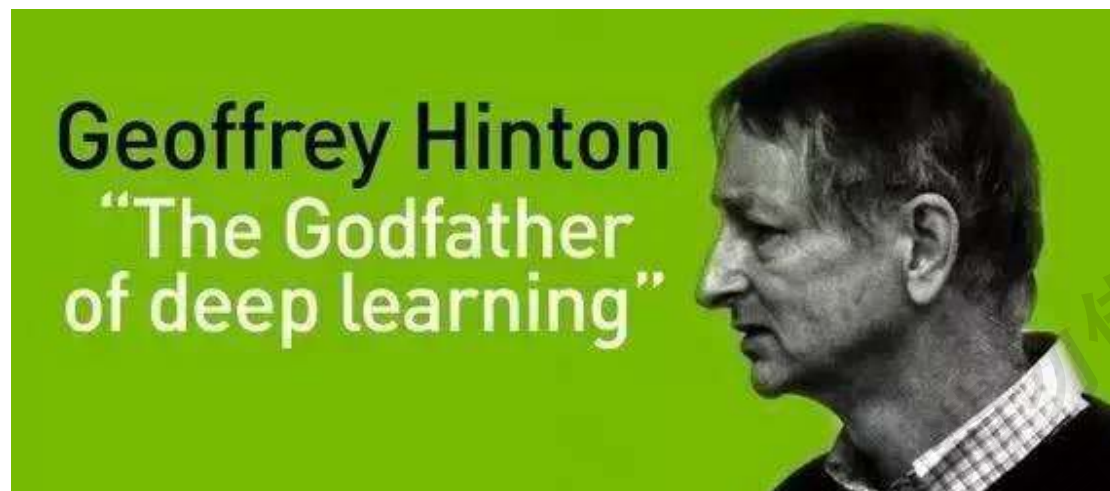
Autumn 2018

# Outline

■ <span style="color:red">Brief Review of Backpropagation Algorithm</span>

■ An Illustrating Example

■ Experiments

■ Assignment

# Brief History of BP



Geoffrey Hinton
"The Godfather of deep learning"
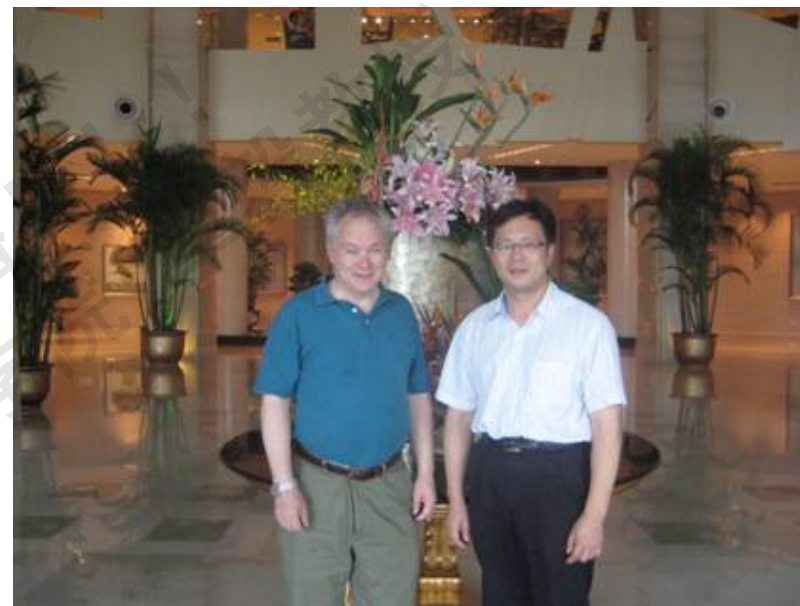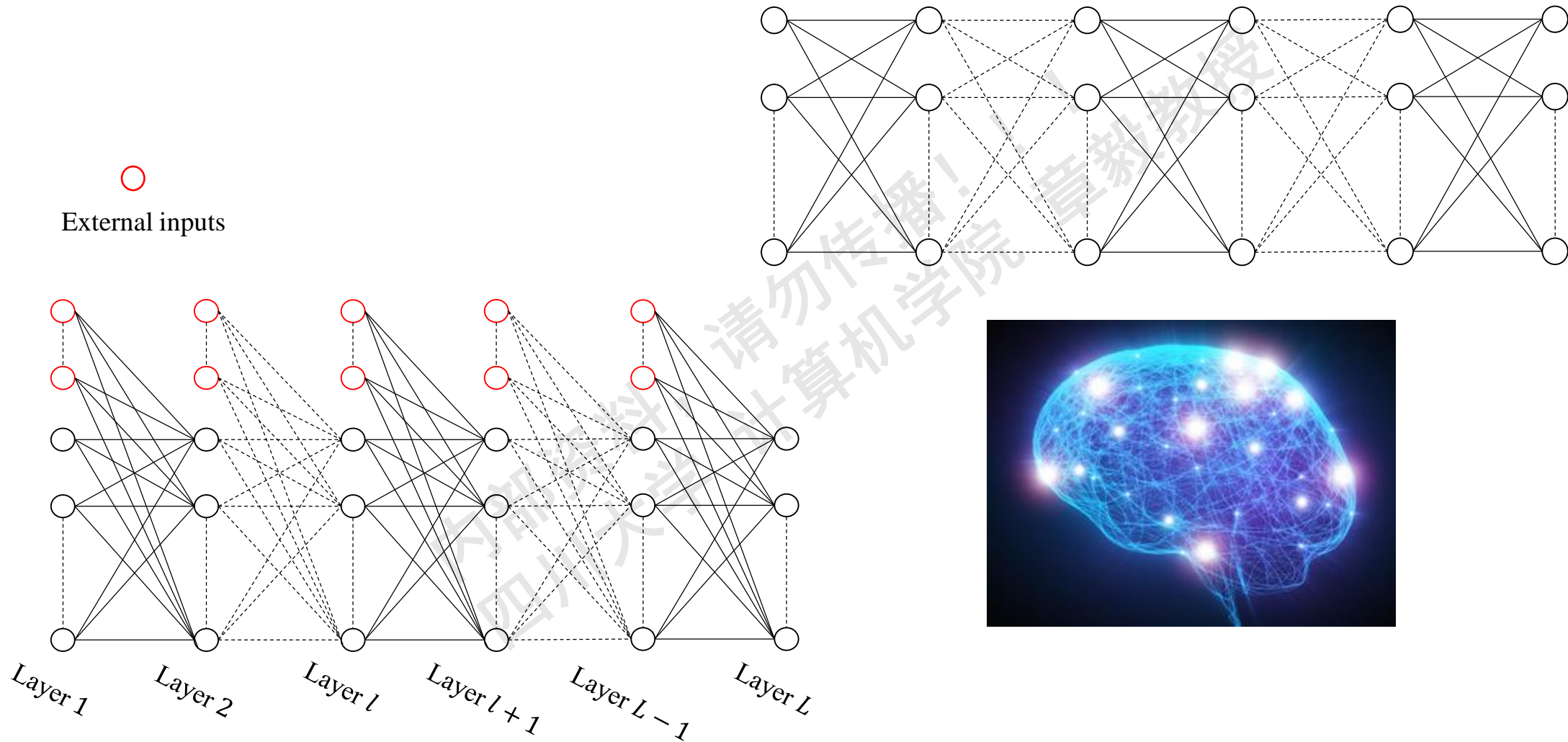
CHAPTER **8**

**Learning Internal Representations by Error Propagation**

D. E. RUMELHART, G. E. HINTON, and R. J. WILLIAMS



Professor P. Werbos

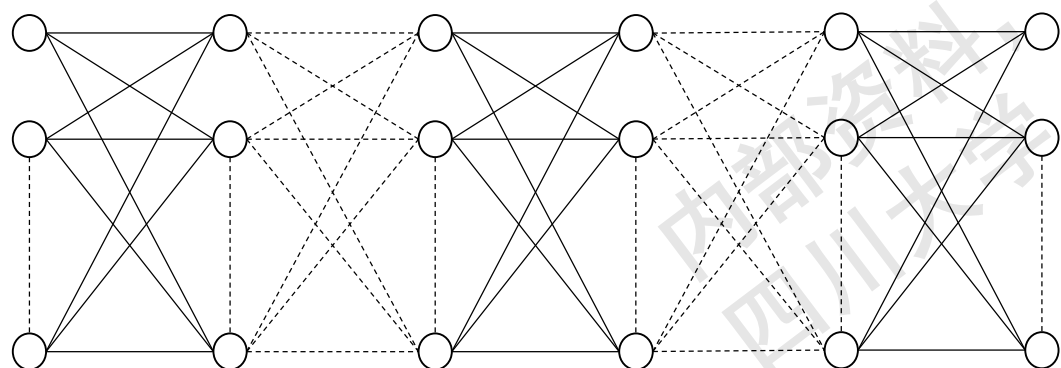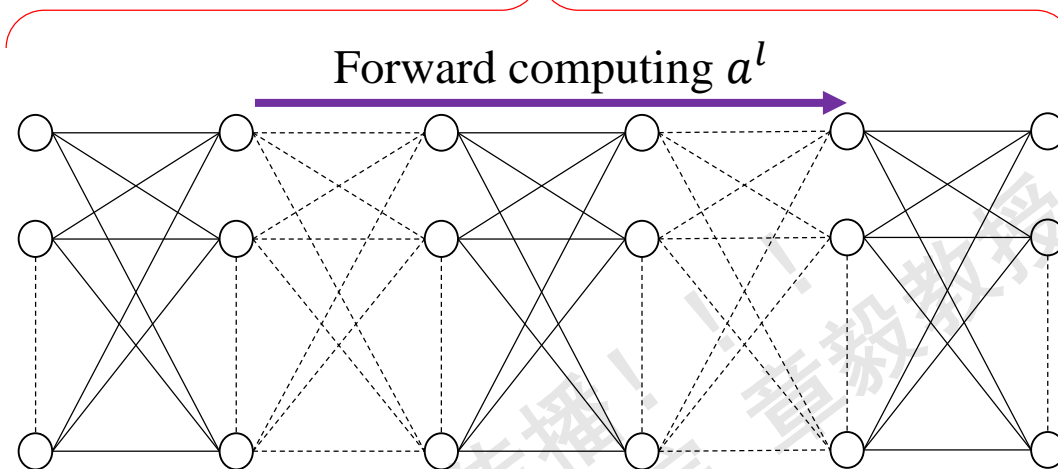# Computational Model of Neural Networks

External inputs

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

Local function defined on neuron

Local activation function $f$



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

$a^l = \begin{bmatrix} a_1^l \\ a_2^l \\ \vdots \\ a_{n_l}^l \end{bmatrix}$
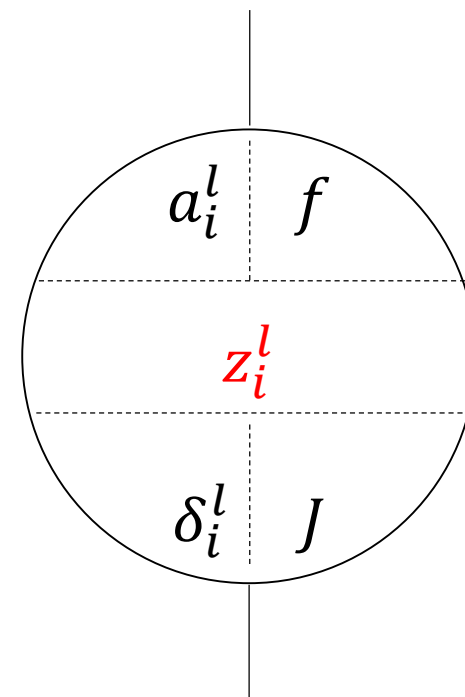
Forward computing $a^l$

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

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

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

$\delta^l = \begin{bmatrix} \delta_1^l \\ \delta_2^l \\ \vdots \\ \delta_{n_l}^l \end{bmatrix}$

Backpropagation $\delta^l$

Global cost function $J$
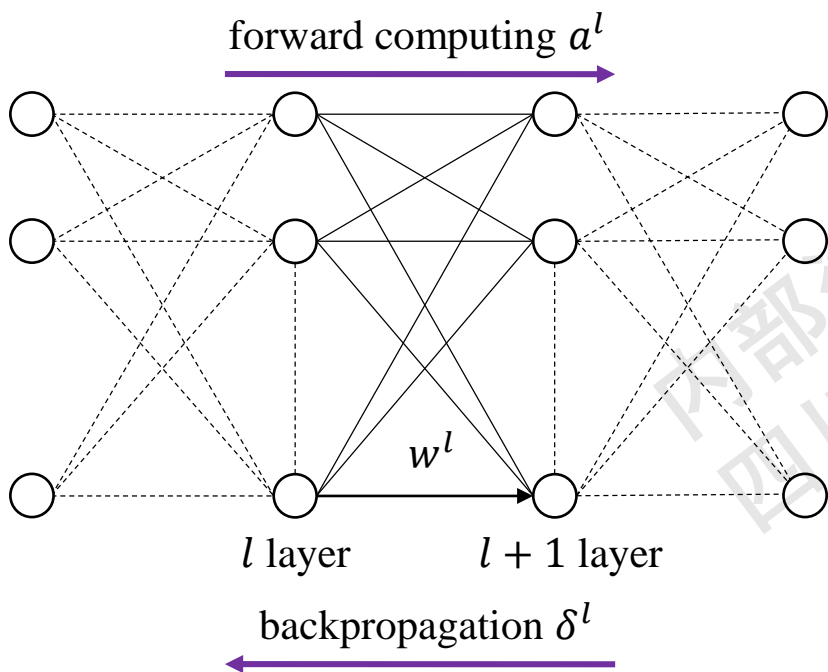
Global function defined on network

# One Page to Understand BP

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

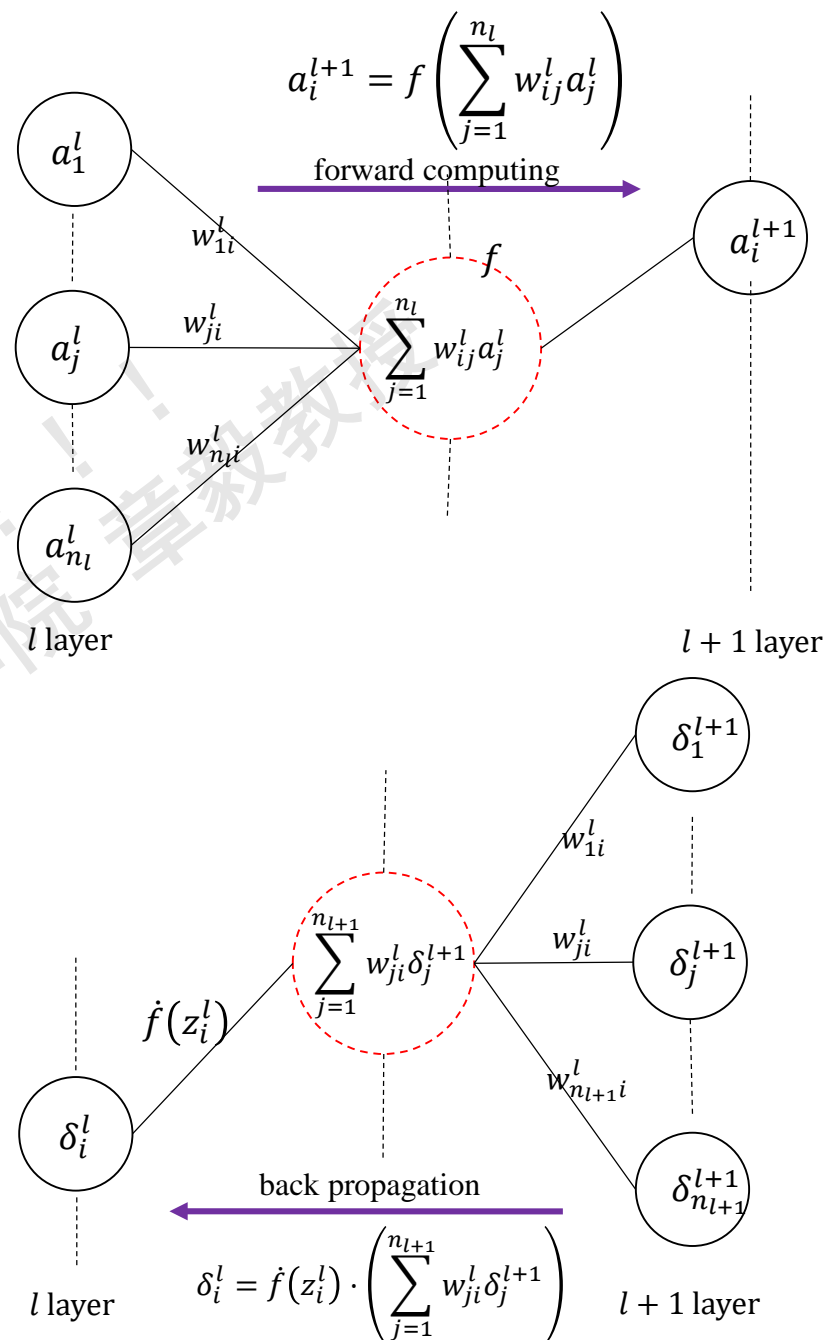Updating rule: $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$

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

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

forward computing

$a_1^l$

$a_j^l$

$a_{n_l}^l$

$w_{1i}^l$

$w_{ji}^l$

$w_{n_l i}^l$

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

$f$

$a_i^{l+1}$

$l$ layer

$l+1$ layer

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

forward computing $a^l$

$w^l$

$l$ layer     $l+1$ layer

backpropagation $\delta^l$

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

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

$\delta_1^{l+1}$

$\delta_j^{l+1}$

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

$w_{1i}^l$

$w_{ji}^l$

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

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

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

$\delta_i^l$

back propagation

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

$l$ layer

$l+1$ layer

# BP Functions

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

forward computing $a^l$
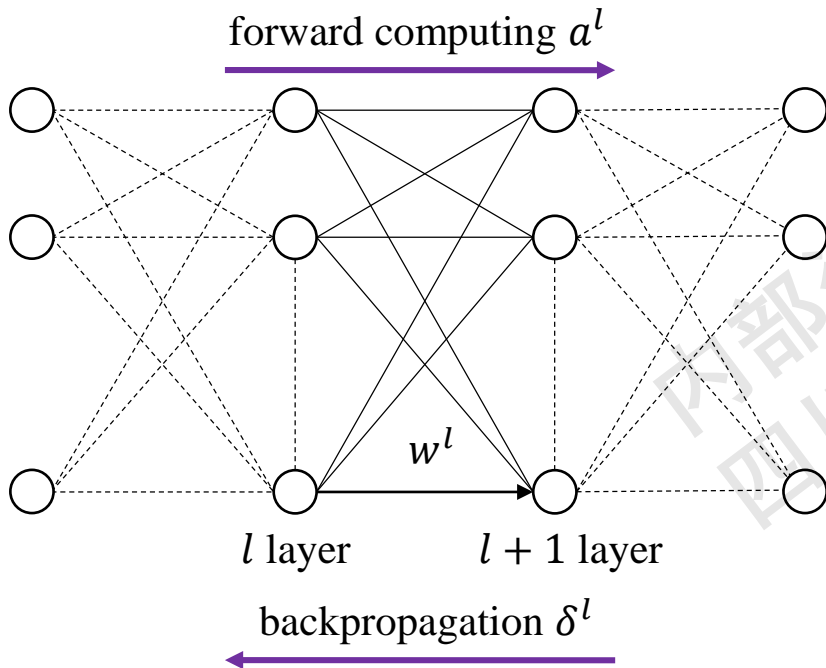
$l$ layer     $l+1$ layer

$w^l$

backpropagation $\delta^l$

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

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

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

%forward computing

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

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

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

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

$end$

%backpropagation

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

$for\ i = 1{:}n_l$

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

$end$

# The BP Algorithm

*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 $D_m \subseteq D$

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

$\qquad$ for $l = 2:L$

$\qquad\qquad a^l \leftarrow fc(w^l, a^l);$

$\qquad$ end

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

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

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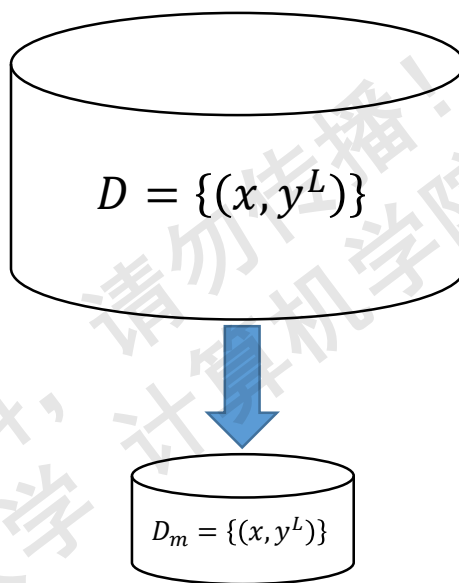
$\qquad$ end

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

Step 4. Updating

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

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

$D = \{(x, y^L)\}$

$D_m = \{(x, y^L)\}$

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

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

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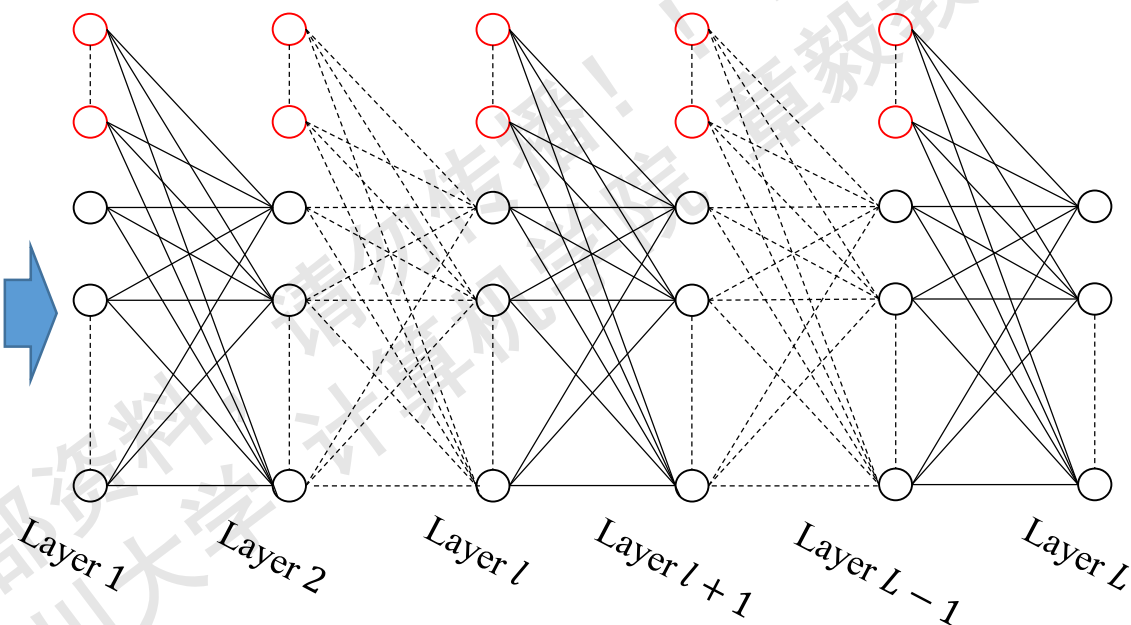$$\delta_i^l = \dot{f}(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} w_{ji}^l \delta_j^{l+1} \right)$$
$$end$$

# Outline

■ Brief Review of Backpropagation Algorithm

■ <span style="color:red">An Illustrating Example</span>

■ Experiments

■ Assignment

# Handwritten digits recognition problem



Task:
Use Backpropagation algorithm to train a neural network to recognize handwritten digits.

# Step 1: Data Preparation

## Data

### Dataset: **MNIST_small**

**MNIST** is a database of handwritten digits created by "re-mixing" the samples from MNIST's original datasets. It contains digits written by high school students and employees of the United States Census Bureau. The digits have been size-normalized and centered in 28 × 28 images.
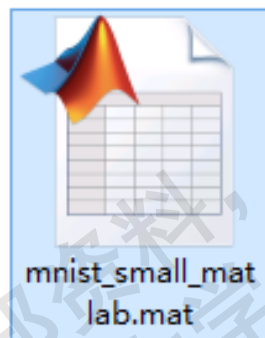
**MNIST_small** dataset is a subset of MNIST containing 10000 training samples and 2000 testing samples.
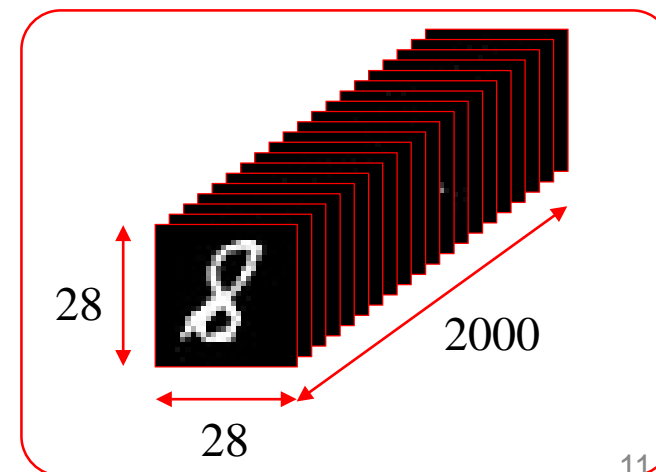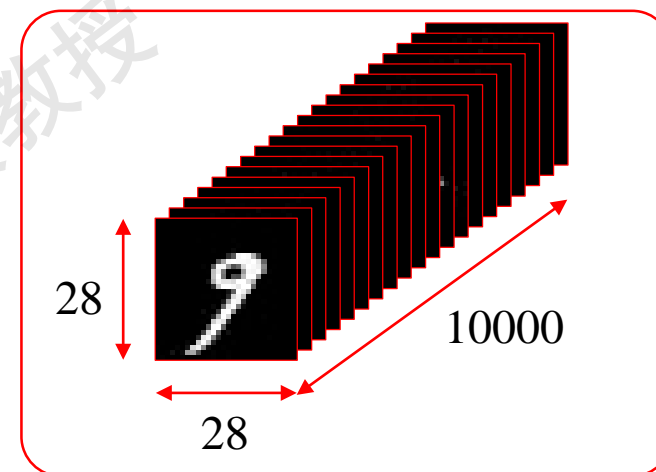
### Training set
- ❑ Used for training network
- ❑ 10000 samples

### Testing set
- ❑ Used for evaluating network performance
- ❑ 2000 samples

Download link:
MNIST http://yann.lecun.com/exdb/mnist/
MNIST_small: https://github.com/kswersky/nnet/blob/master/mnist_small.mat

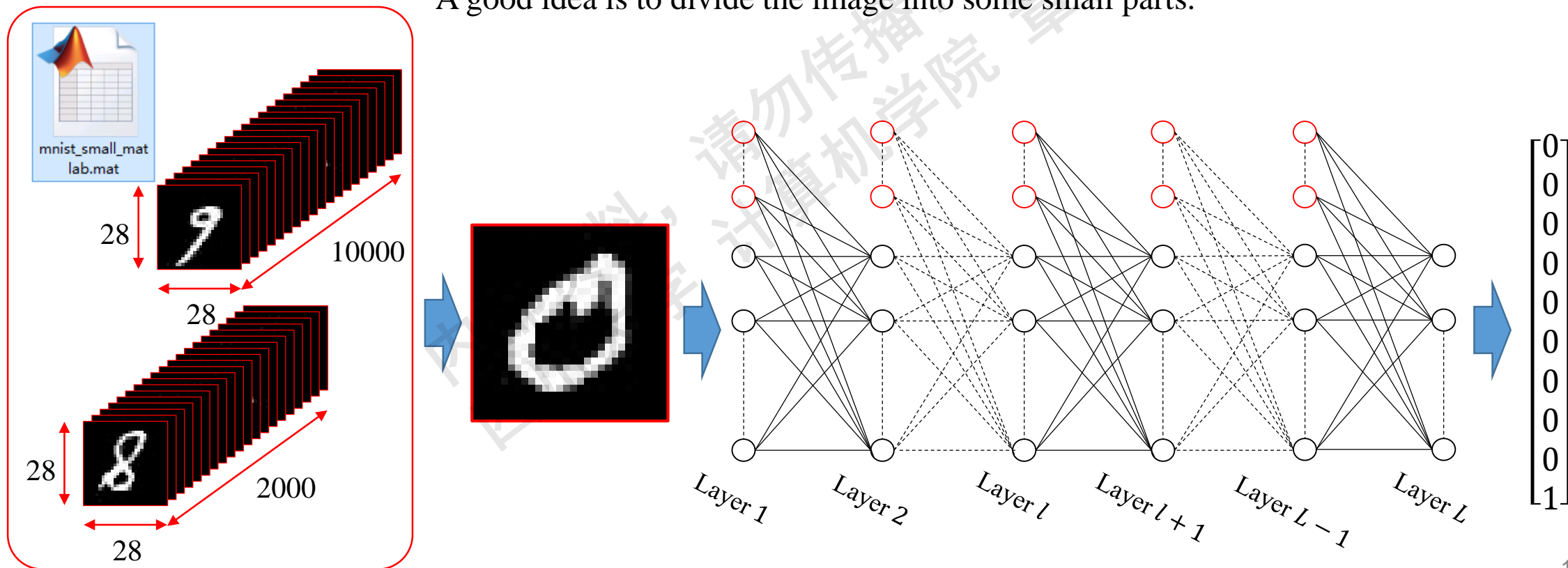mnist_small_mat lab.mat

28
28
10000

28
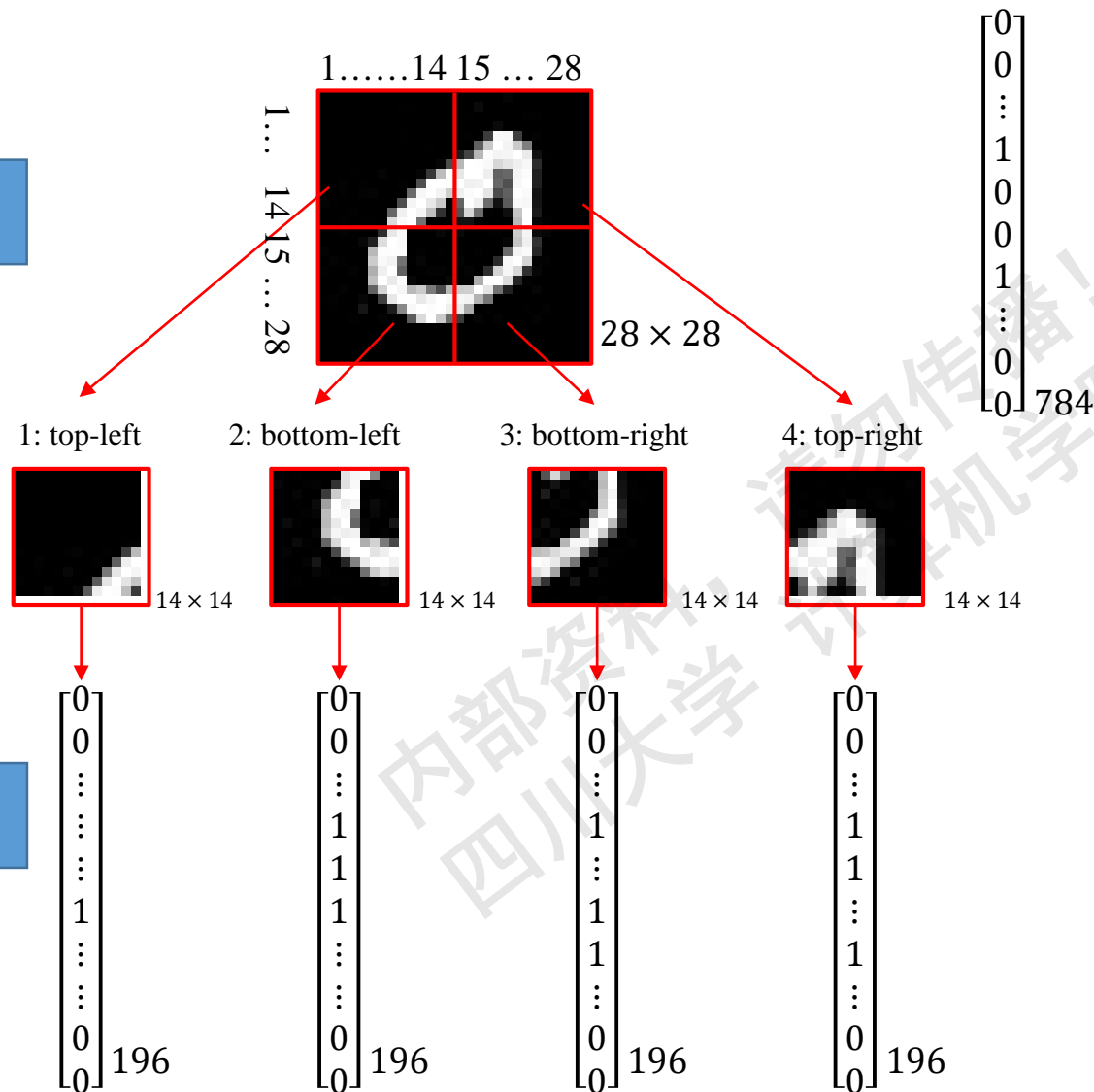28
2000

# Step 1: Data Preparation

# Step 1: Data Preparation

The input image is a $28 \times 28 = 784$ dimensional vector, relatively large in some situation.
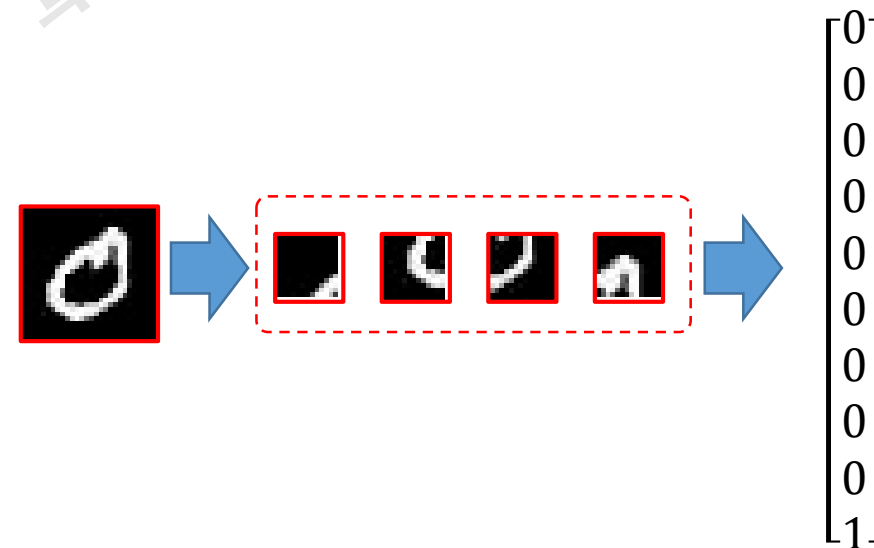A good idea is to divide the image into some small parts.
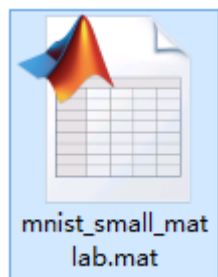
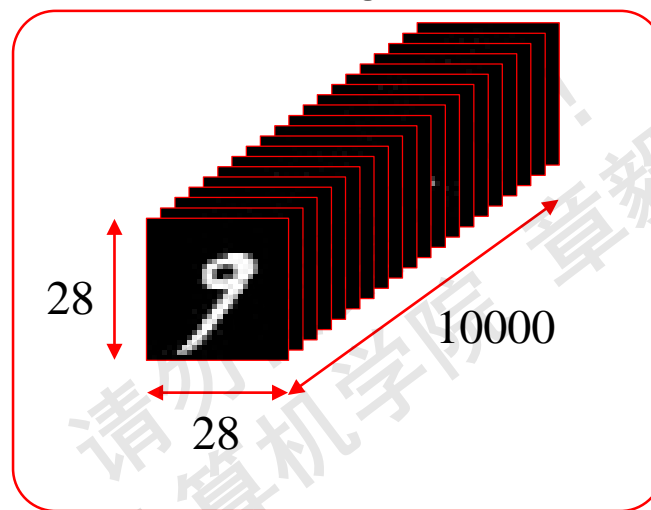# Step 1: Prepare Data Preparation

# Step 1: Data Preparation

**Training Data**

**Training Data**

**Training set**
- ☐ Used for training network
- ☐ 10000 samples
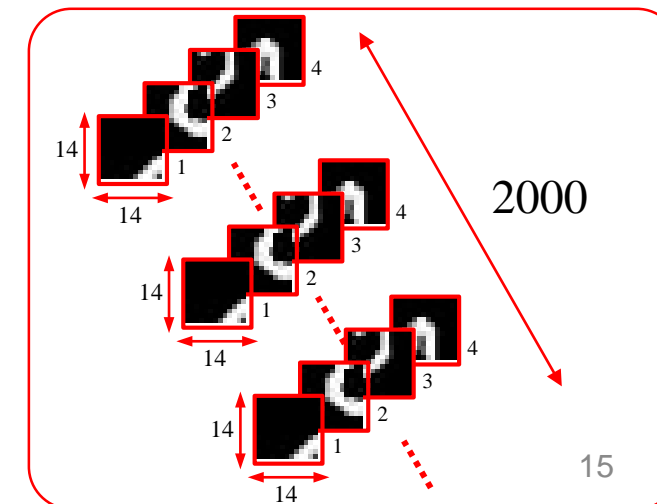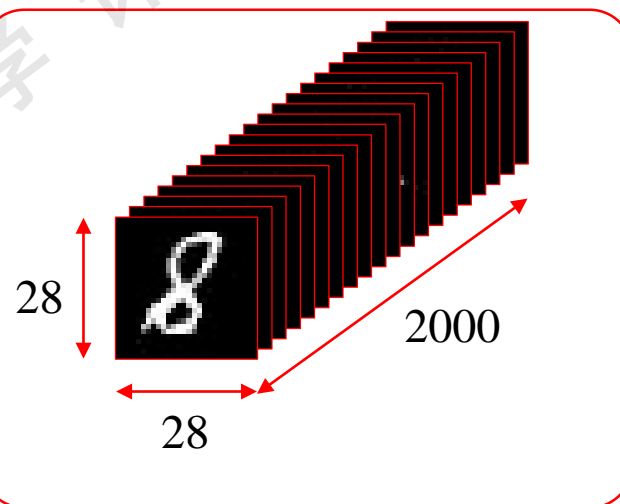- ☐ each sample contains four elements

28

28

10000

10000

**Testing set**
- ☐ Used for evaluating network performance
- ☐ 2000 samples
- ☐ each sample contains four elements

28

28

2000

14

14

14

14

14

14

2000

mnist_small_matlab.mat

# Step 2: Design Network Architecture

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

Activation function

Network architecture design:
1. Number of layers
2. Number of neurons in each layer (external neurons and internal neurons)
3. Activation function

Number of external inputs

Number of internal inputs

Layer $l$

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

# Step 2: Design Network Architecture

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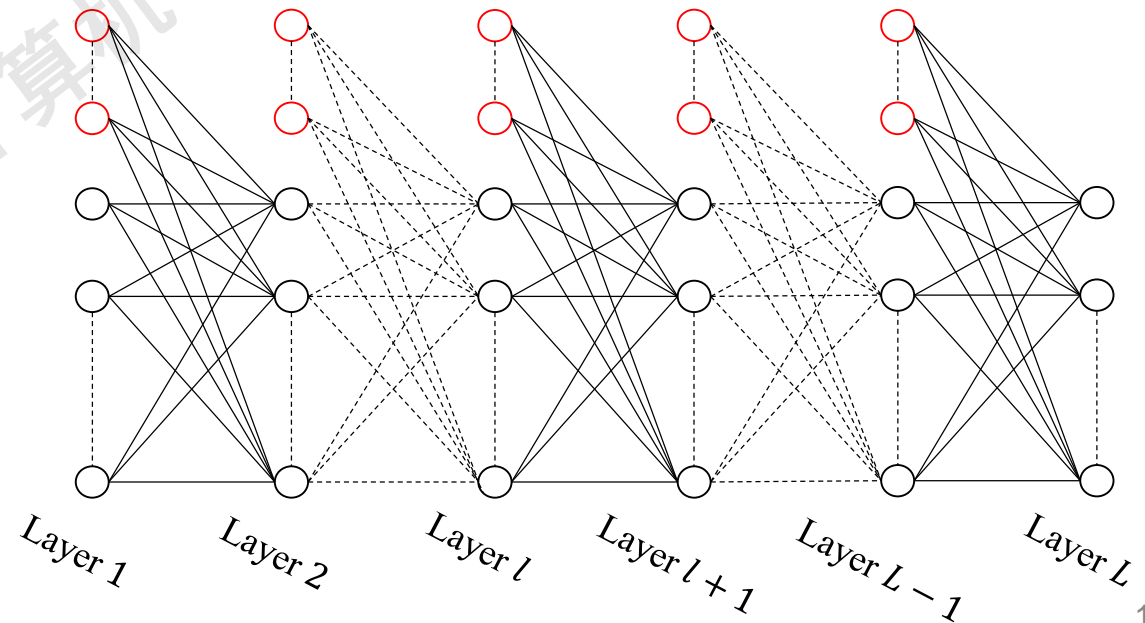$$f(z) = \frac{1}{1 + e^{-z}}$$

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

196-dimension

Layer 1    Layer 2    Layer 3    Layer 4    Layer 5    Layer 6    Layer 7    Layer 8

(196, 32)   (196, 32)   (196, 32)   (196, 32)   (0, 64)    (0, 64)    (0, 64)    (0, 10)

# Step 3: Initial Weights and Learning Rate

**Initialize Weight Connections**

Random initialization:

*Method 1*: Gaussian distribution: $w_{ij}^l \sim N(0,1)$
*Method 2*: Uniform distribution: $w_{ij}^l \sim U(-r^l, r^l)$

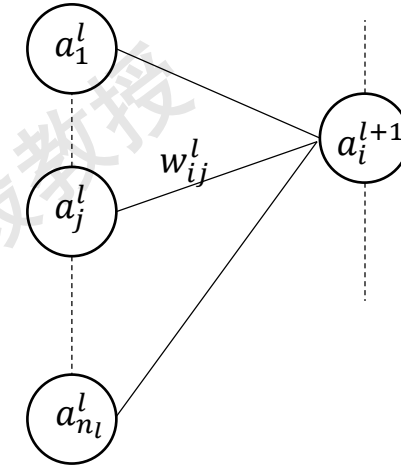$$r^l = \sqrt{\frac{6}{p^l + q^{l+1}}}$$

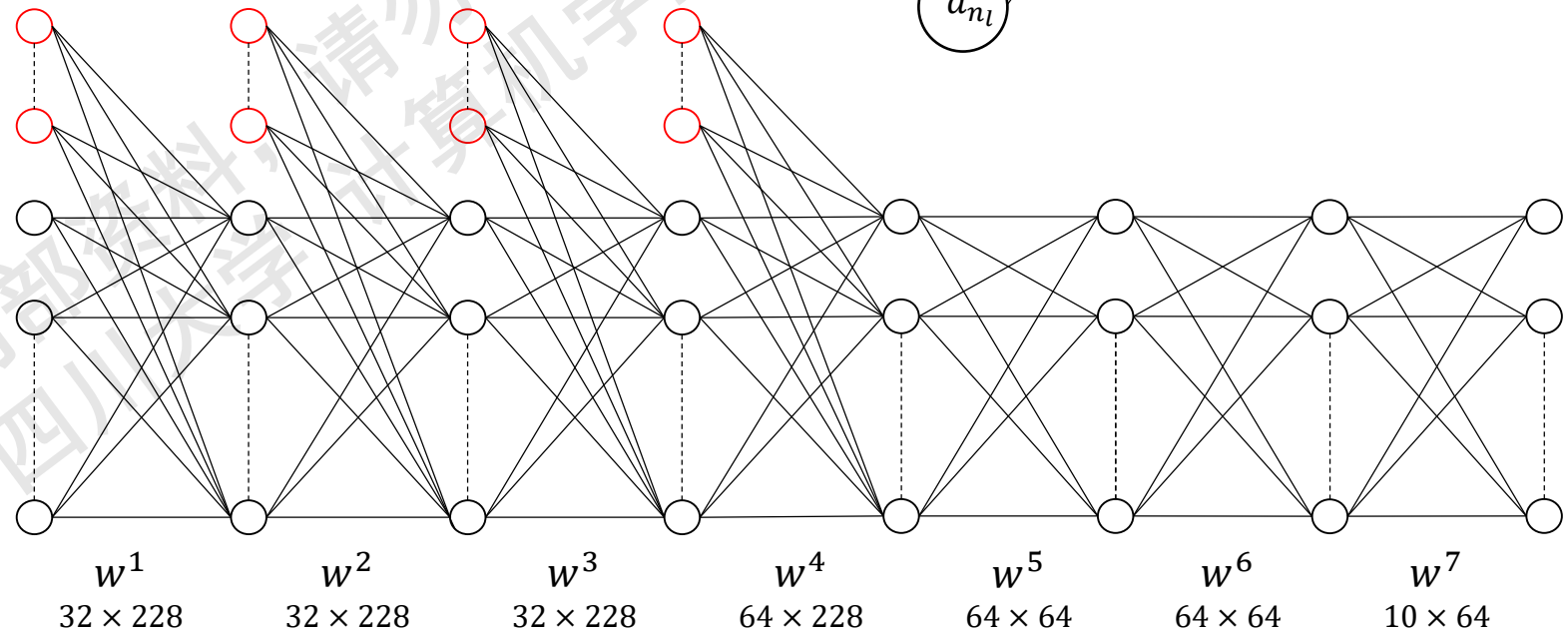$p^l$: number of neurons in $l$ layer
$q^{l+1}$: number of internal neurons in $l + 1$ layer

Initialize Internal Representation of Layer 1:

$$a_i^1 = 0,$$
$$\text{or}$$
$$a_i^1 = 1,$$
$$\dots$$



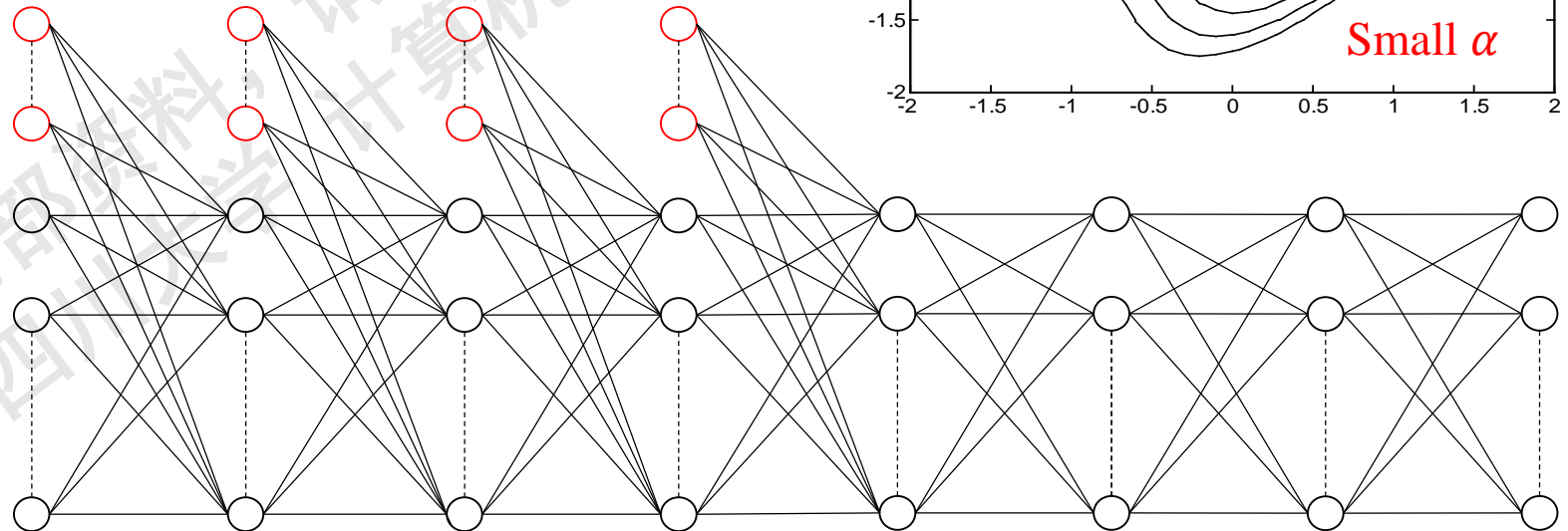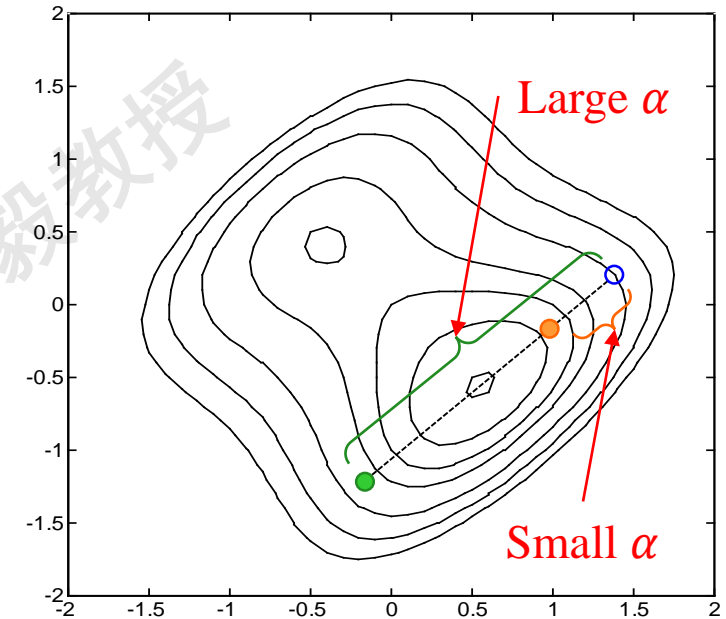| $w^1$ | $w^2$ | $w^3$ | $w^4$ | $w^5$ | $w^6$ | $w^7$ |
|---|---|---|---|---|---|---|
| $32 \times 228$ | $32 \times 228$ | $32 \times 228$ | $64 \times 228$ | $64 \times 64$ | $64 \times 64$ | $10 \times 64$ |

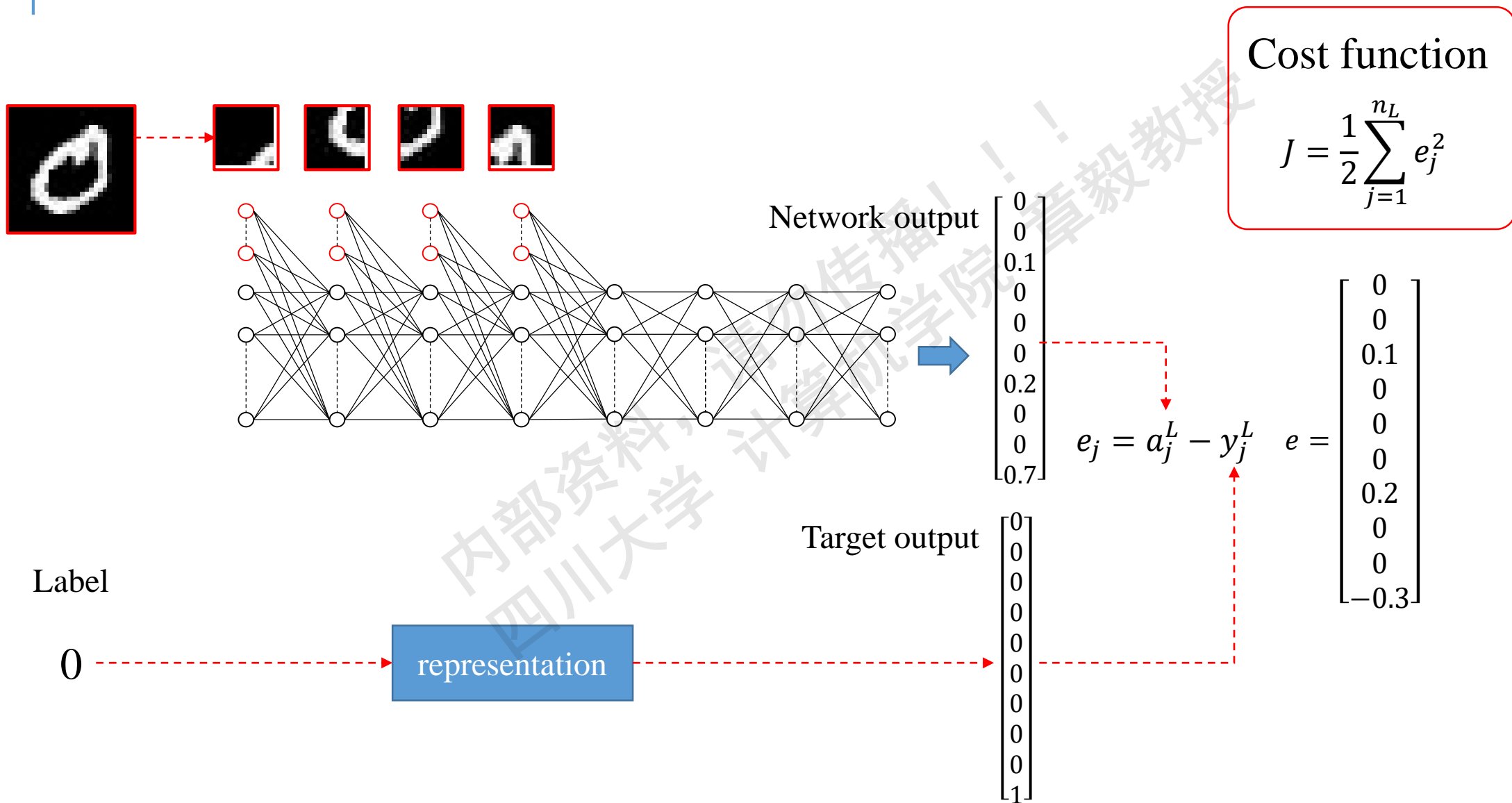# Step 3: Initialization and Learning Rate

**Learning rate:**
- Small: slow learning, long learning time.
- Large: fast learning, possibly not converge to minima.

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

$$\alpha = \cdots, 0.5, 1, 2, 4, \cdots$$

# Step 4: Define Cost Function



Network output

$$\begin{bmatrix} 0 \\ 0 \\ 0.1 \\ 0 \\ 0 \\ 0 \\ 0.2 \\ 0 \\ 0 \\ 0.7 \end{bmatrix}$$

Cost function

$$J = \frac{1}{2}\sum_{j=1}^{n_L} e_j^2$$

$$e_j = a_j^L - y_j^L \qquad e = \begin{bmatrix} 0 \\ 0 \\ 0.1 \\ 0 \\ 0 \\ 0 \\ 0.2 \\ 0 \\ 0 \\ -0.3 \end{bmatrix}$$

Target output

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

Label

0

representation

# Step 5: Define Evaluation Index

$$\text{Accuracy} = \frac{\text{number of correct prediction}}{\text{number of samples}}$$

An example

Tested data 9



Prediction    7    9    0    4    8    8    8    3    1

Correct prediction          Incorrect prediction
7                           2

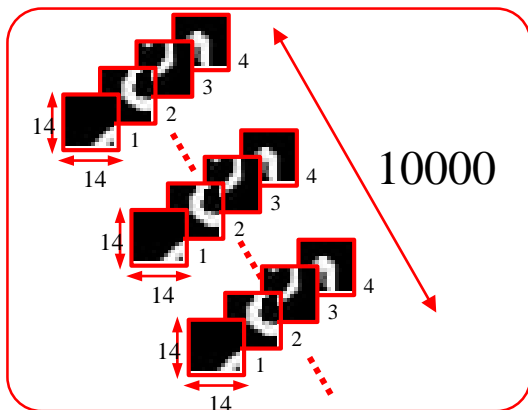$$\text{Accuracy} = \frac{7}{9} = 77.78\%$$

Test on training set:

- Reflect the progress of training.
- Evaluate the ability of the model to fit given data.

Test on testing set:

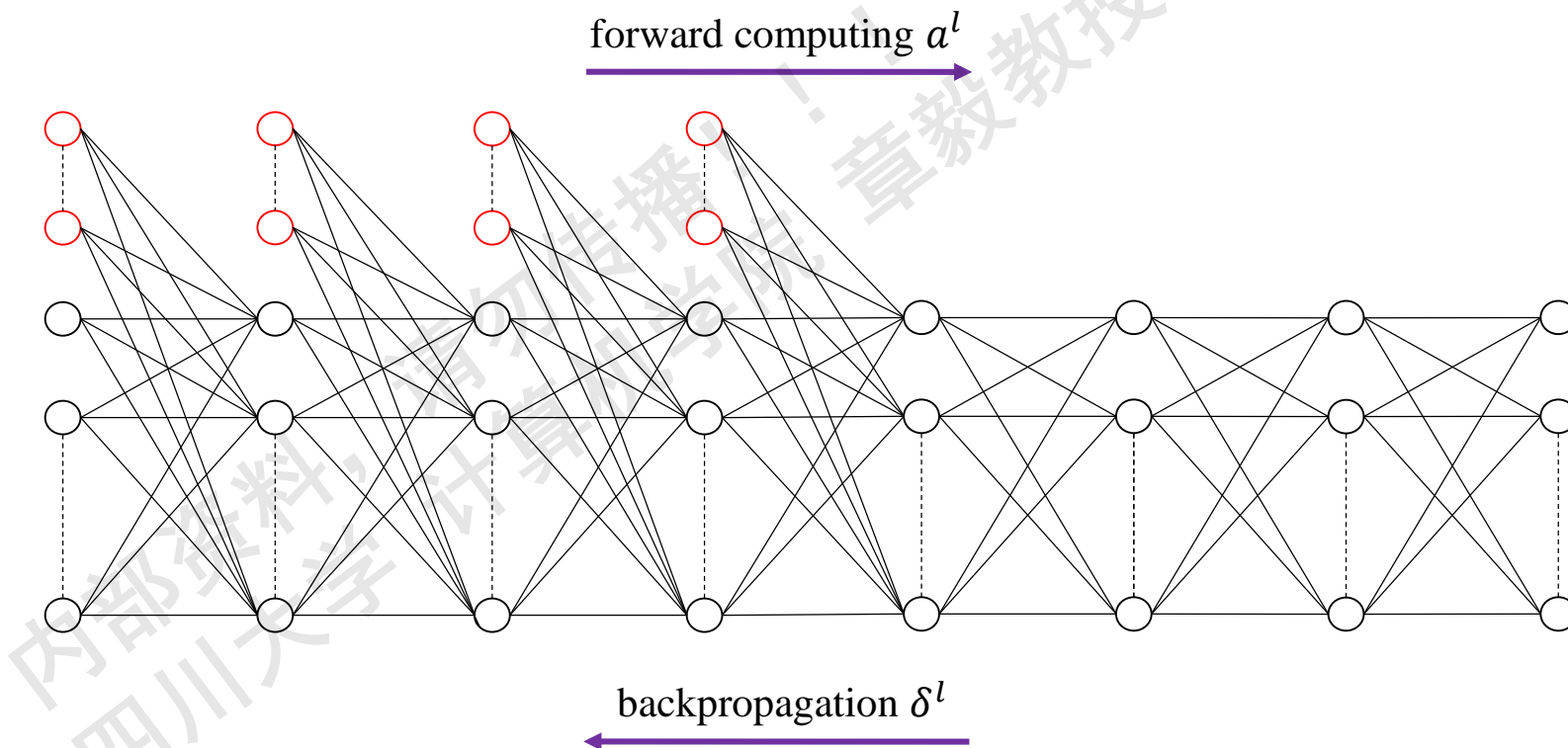- Evaluate the ability of the model to generalize the knowledge.

# Step 6: Train the Network

Training Data



10000
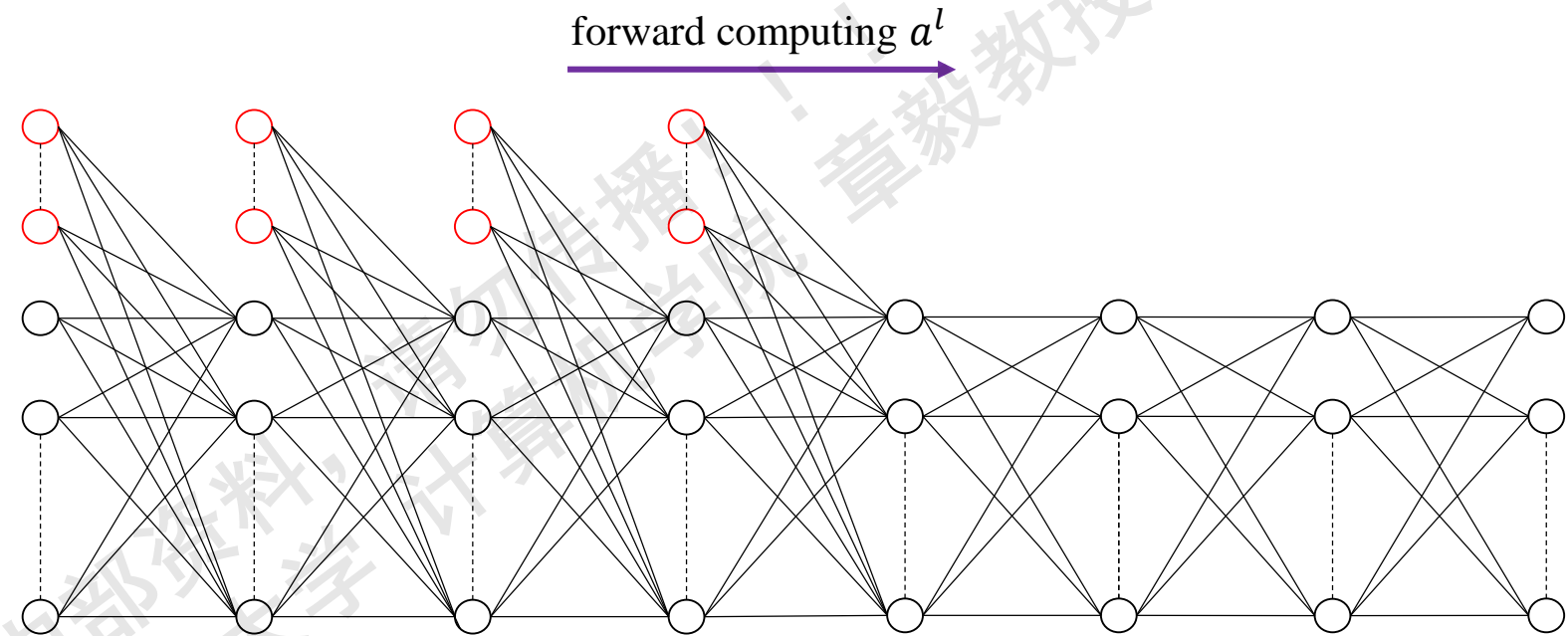
forward computing $a^l$
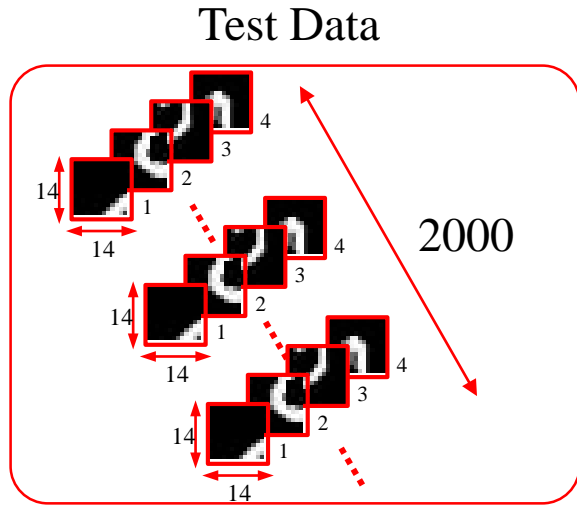
backpropagation $\delta^l$

Updating weights

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

# Step 7: Test the Network

Test Data



2000

forward computing $a^l$

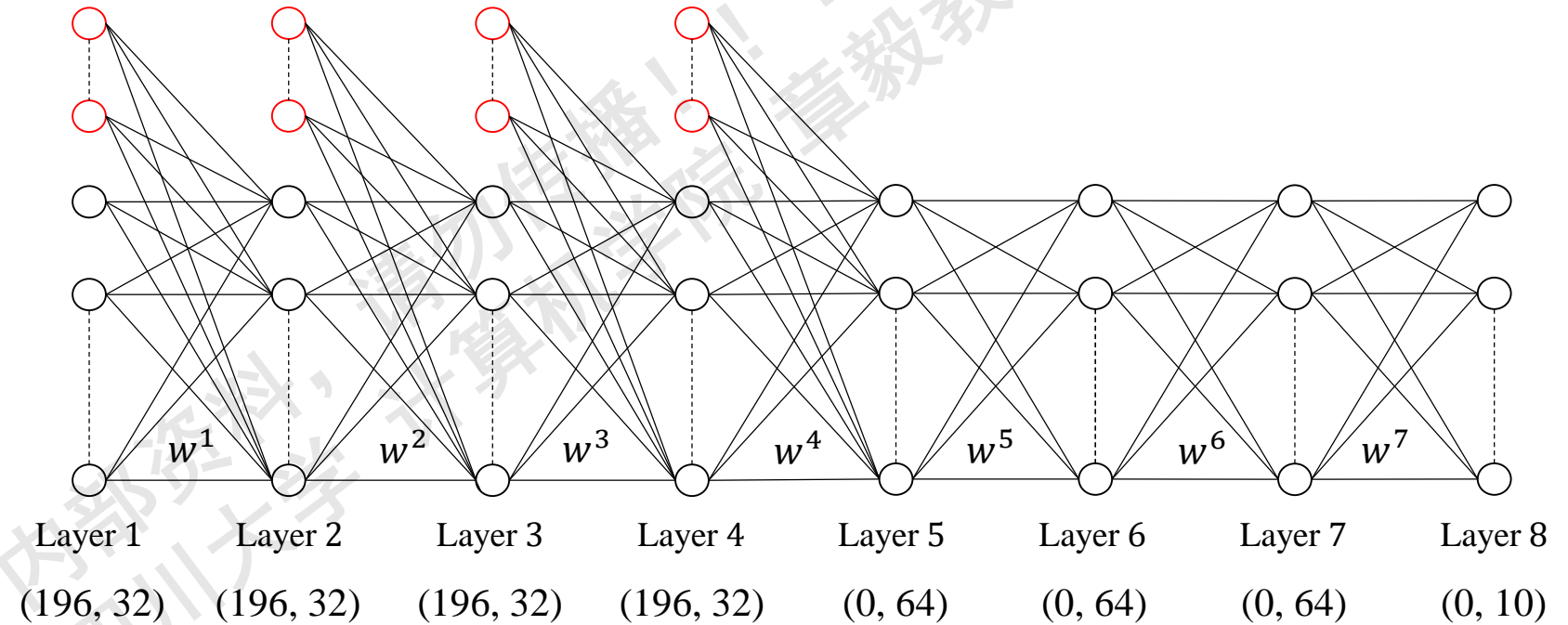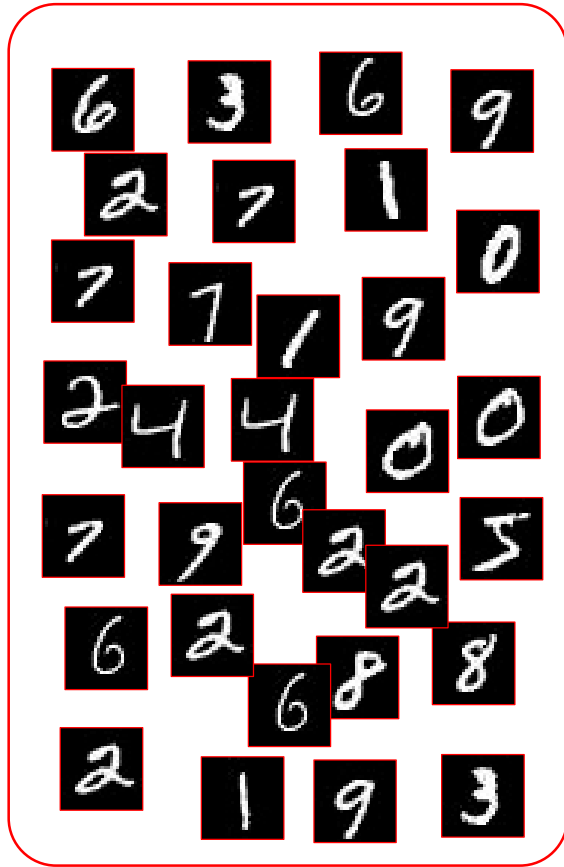Calculating the Evaluation index

$$Accuracy = \frac{number\ of\ correct\ prediction}{2000}$$

# Step 8: Store the Network Parameters

Learning rate $\alpha$

Initialize value of Layer 1



$w^1$  $w^2$  $w^3$  $w^4$  $w^5$  $w^6$  $w^7$

Layer 1    Layer 2    Layer 3    Layer 4    Layer 5    Layer 6    Layer 7    Layer 8

(196, 32)  (196, 32)  (196, 32)  (196, 32)  (0, 64)    (0, 64)    (0, 64)    (0, 10)

# Step 9: Using Trained Network for Applications



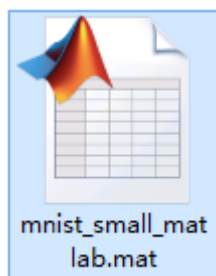| Layer 1 | Layer 2 | Layer 3 | Layer 4 | Layer 5 | Layer 6 | Layer 7 | Layer 8 |
|---------|---------|---------|---------|---------|---------|---------|---------|
| (196, 32) | (196, 32) | (196, 32) | (196, 32) | (0, 64) | (0, 64) | (0, 64) | (0, 10) |

$w^1$ $w^2$ $w^3$ $w^4$ $w^5$ $w^6$ $w^7$

# Outline

■Brief Review of Backpropagation Algorithm

■An Illustrating Example

■<span style="color:red">Experiments</span>

■Assignment

# Experiments: Data Preparation
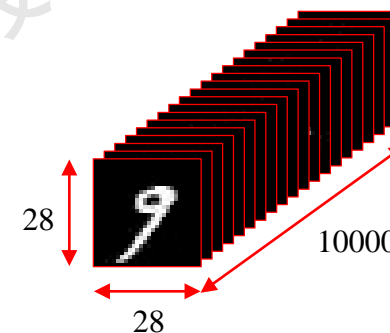


| | 名称 | 值 |
|---|---|---|
| | trainData | 28x28x10000 double |
| | testData | 28x28x2000 double |
| | trainLabels | 10x10000 double |
| | testLabels | 10x2000 double |

mnist_small_mat lab.mat

```
% prepare the data set
load mnist_small_matlab.mat
```

Training Data



```
trainData % 28 * 28 * 10000
```

28

28

10000

Label



```
trainLabels % 10 * 10000
```

10

$$\begin{bmatrix}0\\0\\0\\0\\0\\0\\0\\0\\1\\0\end{bmatrix} \begin{bmatrix}0\\1\\0\\0\\0\\0\\0\\0\\0\\0\end{bmatrix} \cdots \begin{bmatrix}0\\0\\0\\0\\0\\0\\0\\0\\0\\1\end{bmatrix}$$
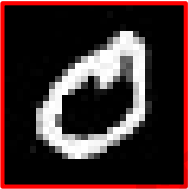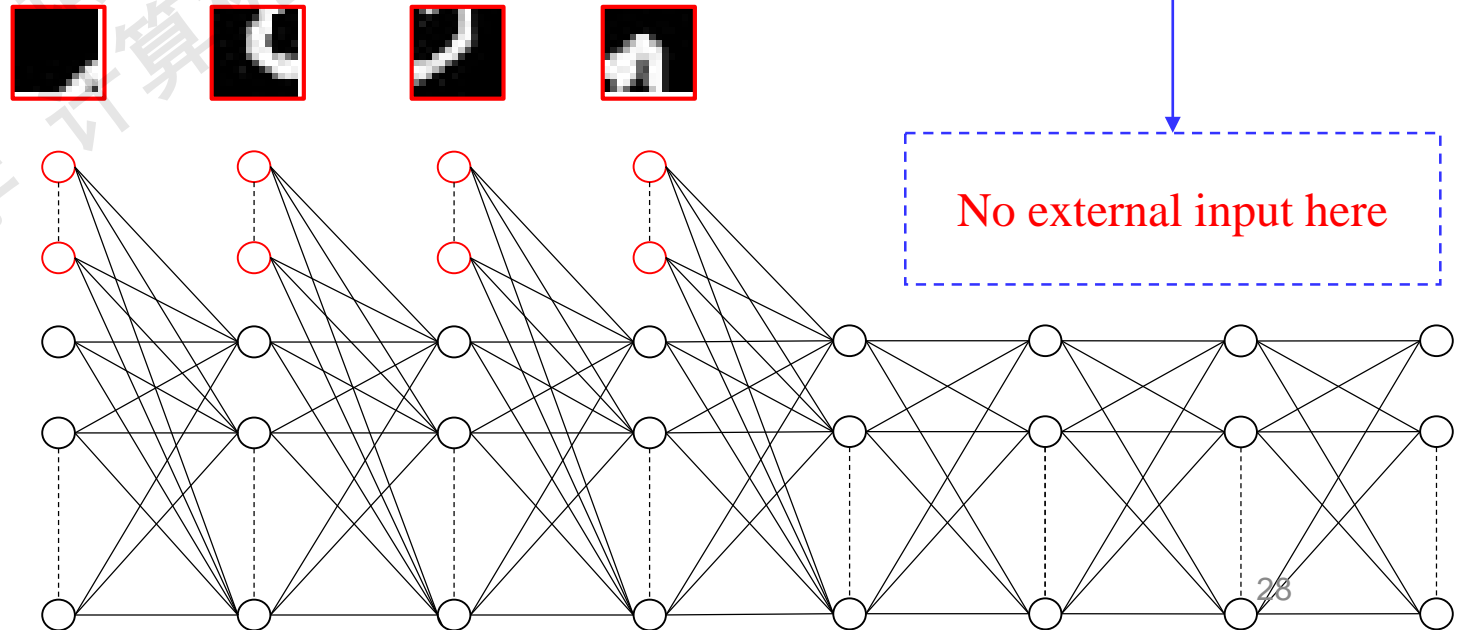
10000

```
% prepare training data
train_size = 10000;
X_train{1} = reshape(trainData(1:14,1:14,:),[],train_size);% top-left
X_train{2} = reshape(trainData(15:28,1:14,:),[],train_size); % bottom-left
X_train{3} = reshape(trainData(15:28,15:28,:),[],train_size); % bottom-right
X_train{4} = reshape(trainData(1:14,15:28,:),[],train_size); % top-right
X_train{5} = zeros(0, train_size);
X_train{6} = zeros(0, train_size);
X_train{7} = zeros(0, train_size);
X_train{8} = zeros(0, train_size);
% prepare testing data
% ...
```
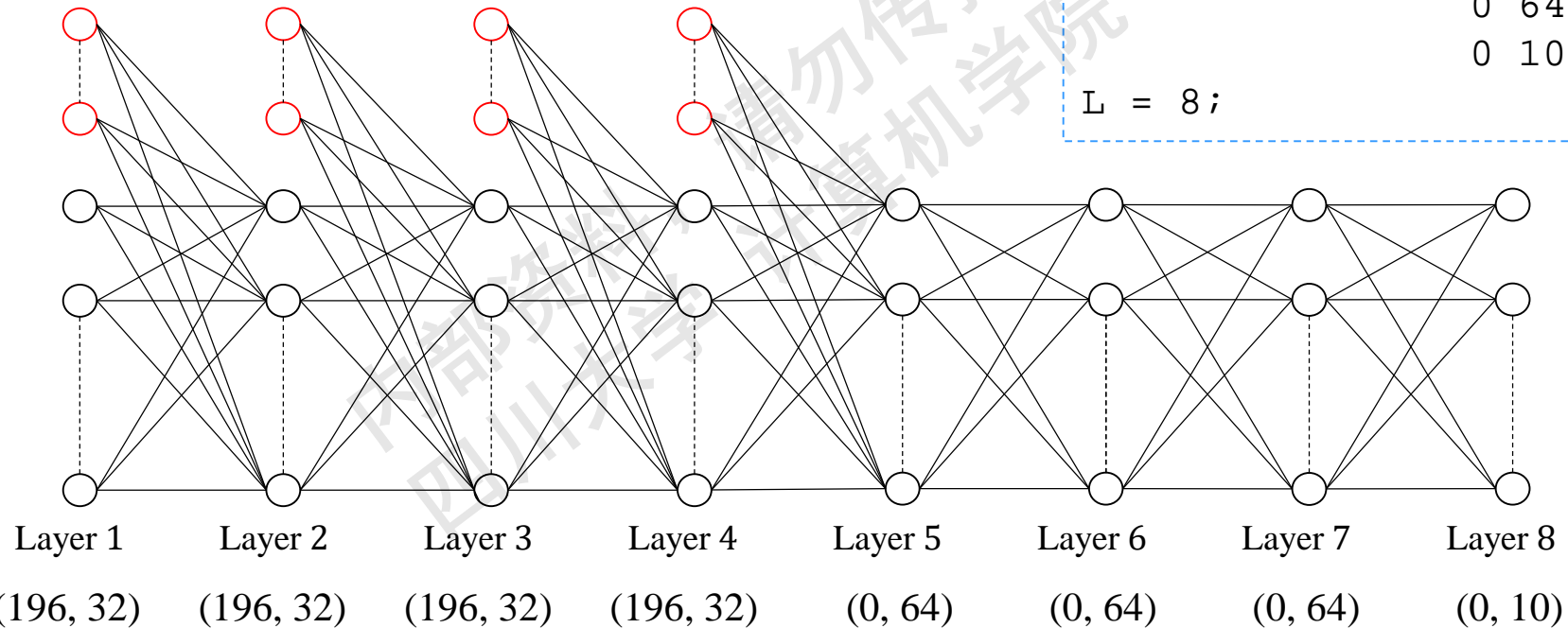
1…14 15 … 28

1…14 15 … 28

1: top-left    2: bottom-left   3: bottom-right   4: top-right
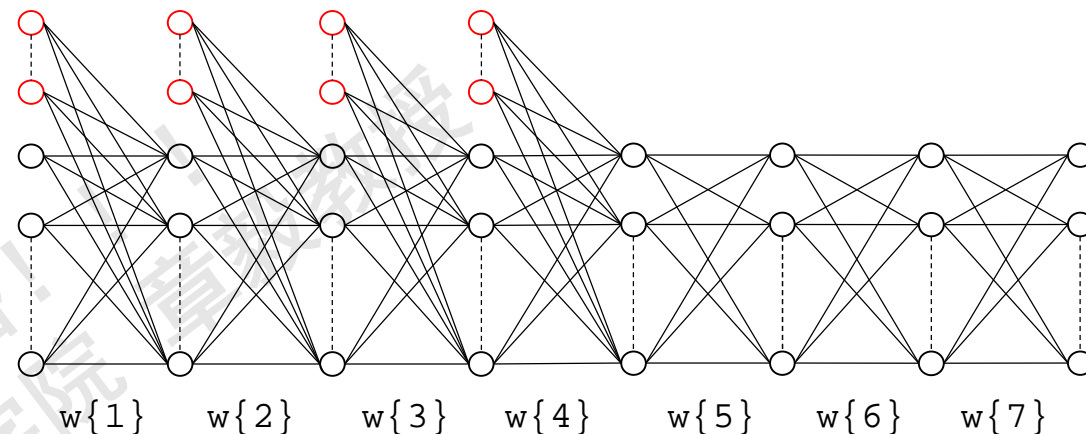
No external input here

# Experiments: Architecture

```
% define network architecture
% - 1st column: external neurons
% - 2nd column: internal neurons
layer_size = [196 32      % layer 1
              196 32      % layer 2
              196 32      % layer 3
              196 32      % layer 4
                0 64      % layer 5
                0 64      % layer 6
                0 64      % layer 7
                0 10];    % layer 8
L = 8;
```



|  | Layer 1 | Layer 2 | Layer 3 | Layer 4 | Layer 5 | Layer 6 | Layer 7 | Layer 8 |
|---|---|---|---|---|---|---|---|---|
| Dimensions | (196, 32) | (196, 32) | (196, 32) | (196, 32) | (0, 64) | (0, 64) | (0, 64) | (0, 10) |

# Experiments: Initialize Weights

Gaussian distribution: $w_{ij}^l \sim N(0,1)$

```
% initialize weights
for l = 1:L-1
    w{l} = randn(layer_size(l+1,2), sum(layer_size(l,:)));
end
```
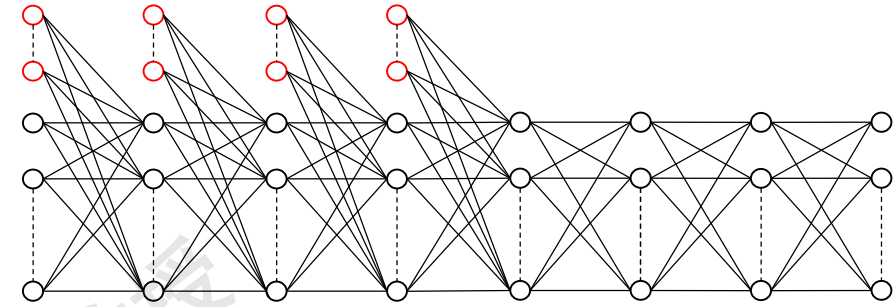
Uniform distribution: $w_{ij}^l \sim U\left(-r^l, r^l\right)$

```
% initialize weights
for l = 1:L-1
    % a tricky, but effective, initialization
    w{l} = (rand(layer_size(l+1,2), sum(layer_size(l,:))) * 2 -1)
            * sqrt(6/(layer_size(l+1,2)+sum(layer_size(l,:)))));
end
```

$$r^l = \sqrt{\frac{6}{p^l + q^{l+1}}}$$

$p^l$: number of neurons in $l$ layer
$q^{l+1}$: number of internal neurons in $l+1$ layer

w{1}  w{2}  w{3}  w{4}  w{5}  w{6}  w{7}

# Experiments: Run the Network



Learning rate

Number of iteration

Number of samples in a batch

```
% choose parameters
alpha = 1;
max_iter = 300;
mini_batch = 100;
```
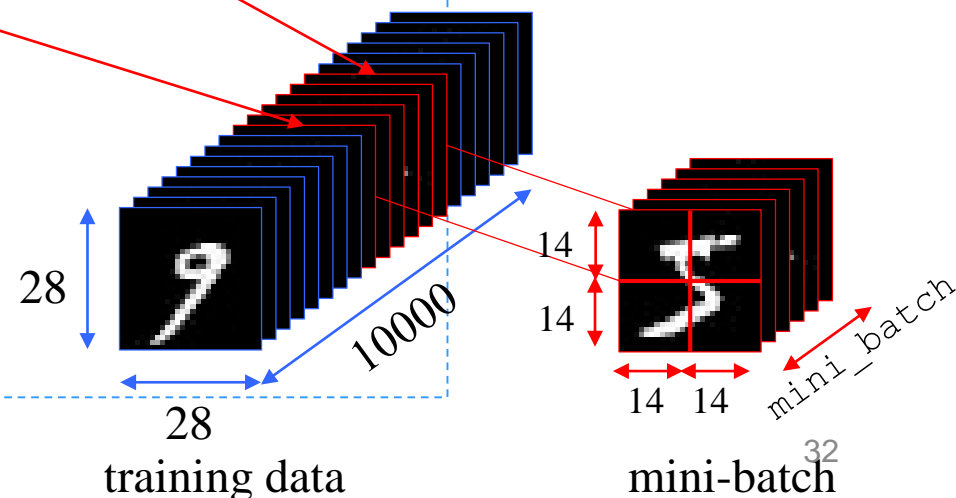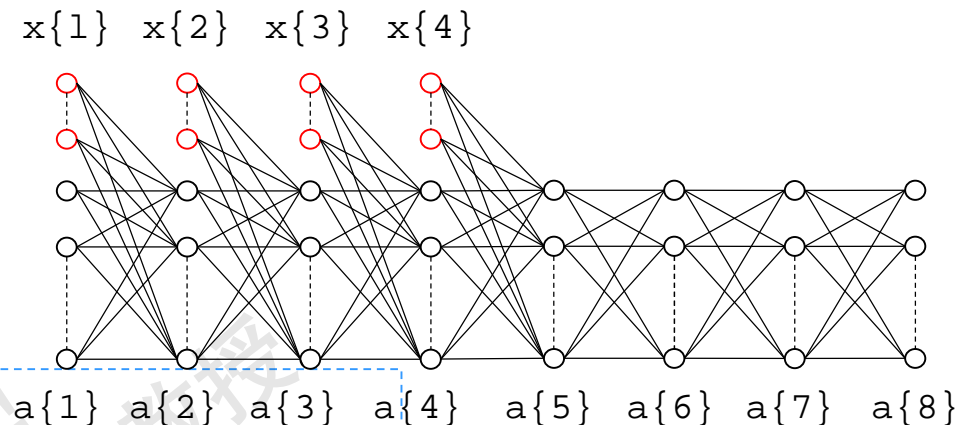
```
% loop until converge
for iter = 1:max_iter
    % for each mini-batch
        % batch forward computation
        % batch backward computation
        % cumulate and update weight
end
```
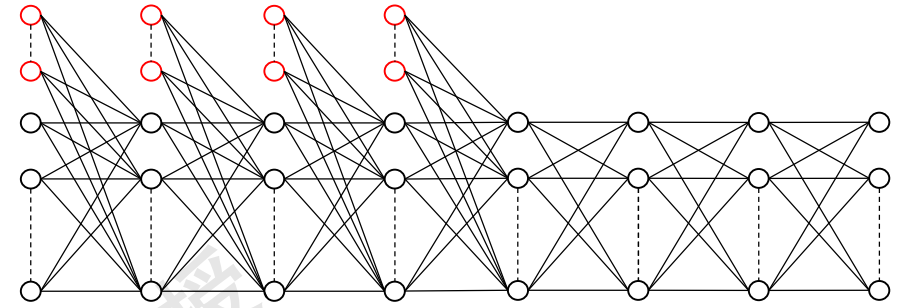
Mini-batch BP implement

_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 $D_m \subset D$

$\quad\quad a^1 \leftarrow samples\ in\ D_m$

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

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

$\quad\quad$ end

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

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

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

$\quad\quad$ end

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

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.

# Experiments: mini-batch BP

$$x\{1\} \quad x\{2\} \quad x\{3\} \quad x\{4\}$$



$$a\{1\} \quad a\{2\} \quad a\{3\} \quad a\{4\} \quad a\{5\} \quad a\{6\} \quad a\{7\} \quad a\{8\}$$

shuffle index

```
ind = randperm(train_size);
% for each mini-batch
for k = 1:ceil(train_size/mini_batch)
    % prepare internal inputs
    a{1} = zeros(layer_size(1,2),mini_batch);
    % prepare external inputs
    for l=1:L
        x{l} = X_train{l}(:,ind((k-1)*mini_batch+1:min(k*mini_batch, train_size)));
    end
    % prepare labels
    y = double(trainLabels(:,ind((k-1)*mini_batch+1:min(k*mini_batch, train_size))));
    % forward computation
    % ...
    % cost function and error
    % ...
    % backward computation
    % ...
    % update weight
    % ...
end
```

28

10000

28

training data

14

14

14 14

mini_batch

mini-batch

# Experiments: BP



```matlab
% forward computation
for l=1:L-1
    [a{l+1}, z{l+1}] = fc(w{l}, a{l}, x{l});
end

% Compute delta of last layer
delta{L} = (a{L} - y).* a{L} .*(1-a{L});

% backward computation
for l=L-1:-1:2
    delta{l} = bc(w{l}, z{l}, delta{l+1});
end

% update weight
for l=1:L-1
    gw = delta{l+1} * [x{l};a{l}]' / mini_batch;
    w{l} = w{l} - alpha * gw;
end
```

*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 $D_m \subset D$

$\qquad a^1 \leftarrow samples\ in\ D_m$

$\qquad$ for $l = 2:L$

$\qquad\qquad a^l \leftarrow fc(w^l, a^l);$

$\qquad$ end

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

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

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

$\qquad$ end

$\qquad \dfrac{\partial J}{\partial w_{ji}^l} \leftarrow \dfrac{\partial J}{\partial w_{ji}^l} + \delta_j^{l+1} \cdot a_i^l;$
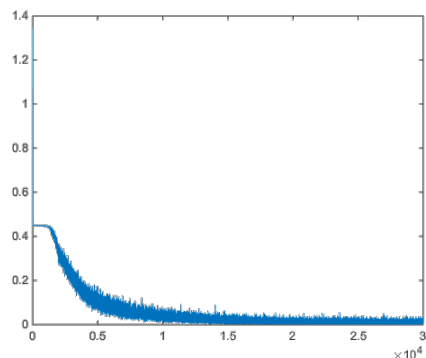
Step 4. Updating

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

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

33

# Experiments: Plotting

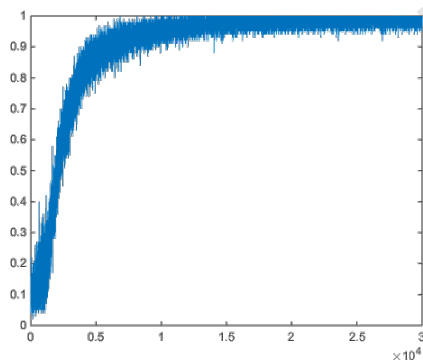**Cost function**

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

```matlab
% cost function
J = [J 1/2/mini_batch*sum((a{L}(:)-y(:)).^2)];
figure
plot(J);
```

**Accuracy**

$$Acc = \frac{number\ of\ correct\ prediction}{number\ of\ samples}$$
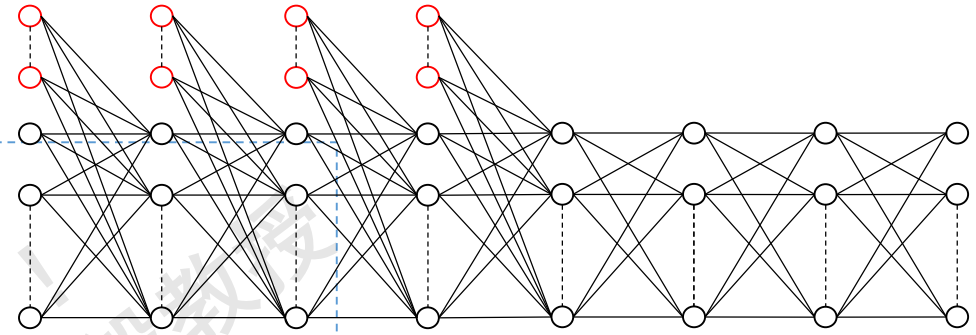
Use max output as prediction

```matlab
% accuary on training batch
[~,ind_train] = max(y);
[~,ind_pred] = max(a{L});
Acc= [Acc sum(ind_train == ind_pred) / mini_batch];
figure
plot(Acc);
```
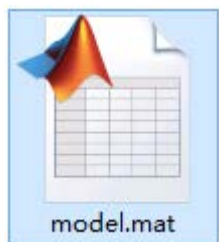
# Experiments: Testing

```
% test on training set
a{1} = zeros(layer_size(1,2),train_size);
for l=1:L-1
    a{l+1} = fc(w{l}, a{l}, X_train{l});
end
[~,ind_test] = max(trainLabels);
[~,ind_pred] = max(a{L});
train_acc = sum(ind_test == ind_pred)/train_size;
fprintf('Accuracy on training dataset is %f%%\n', train_acc*100);
```
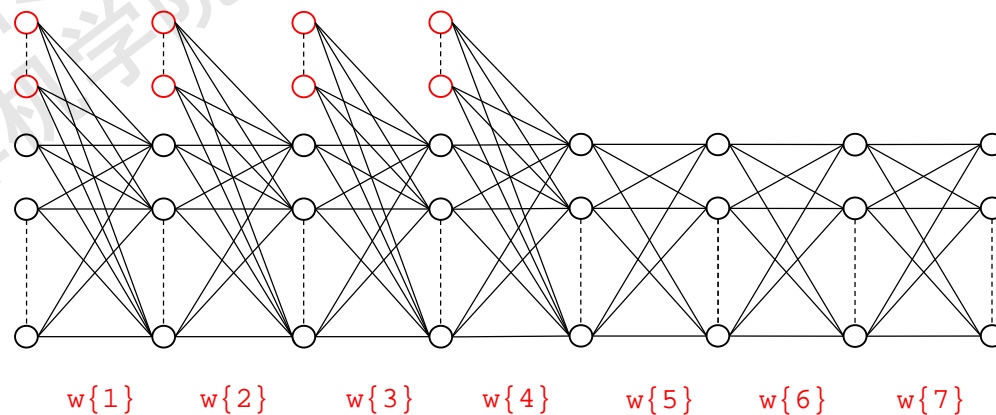
```
% test on testing set
a{1} = zeros(layer_size(1,2),test_size);
for l=1:L-1
    a{l+1} = fc(w{l}, a{l}, X_test{l});
end
[~,ind_test] = max(testLabels);
[~,ind_pred] = max(a{L});
test_acc = sum(ind_test == ind_pred)/test_size;
fprintf('Accuracy on testing dataset is %f%%\n', test_acc*100);
```

# Experiments: Store the Network Parameters

```
% save model
save model.mat w layer_size
```
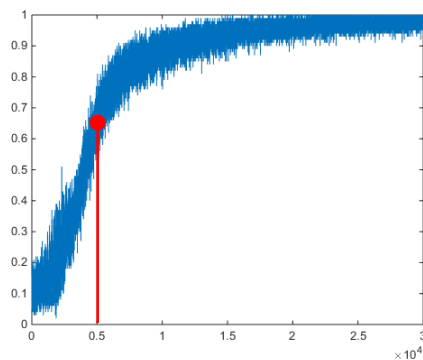
| | 名称 | 值 |
|---|---|---|
| | layer_size | 8x2 double |
| {} | w | 1x7 cell |

model.mat

**This is very important!**

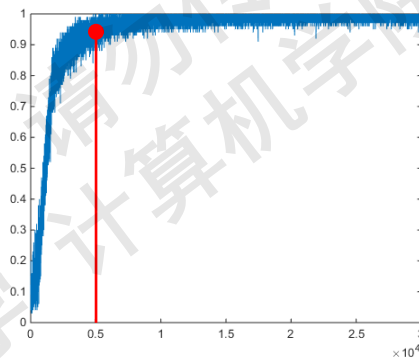w{1}  w{2}  w{3}  w{4}  w{5}  w{6}  w{7}

# Results: Learning Rate

```
% learning rate
alpha = 0.5;
```



Accuracy
- Training=98.05%
- Testing=94.40%

```
% learning rate
alpha = 2;
```



Accuracy
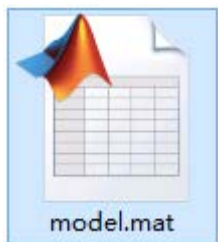- Training=99.14%
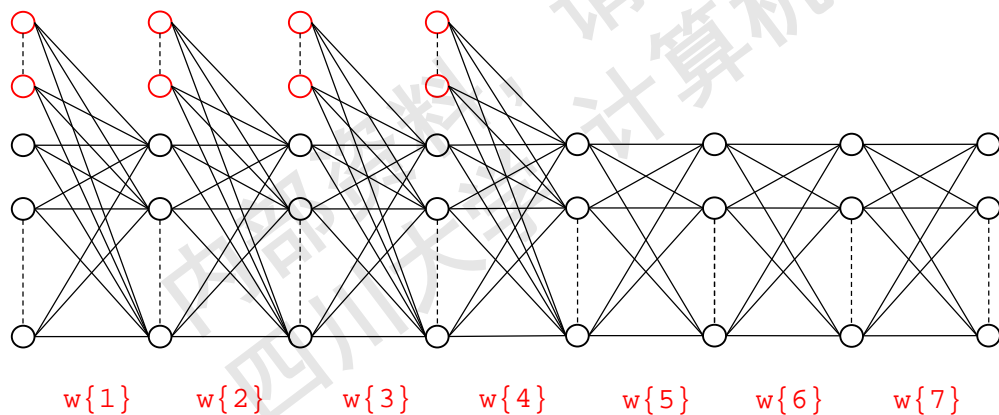- Testing=95.40%

```
% learning rate
alpha = 8;
```


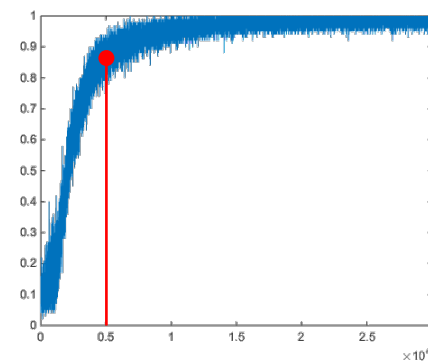
Accuracy
- Training=71.02%
- Testing=69.25%

# Results: Number of Layers

| 名称 | 值 |
|---|---|
| layer_size | 8x2 double |
| w | 1x7 cell |

model.mat

8 layers



Accuracy
Training=98.65%
Testing=95.10%

w{1}  w{2}  w{3}  w{4}  w{5}  w{6}  w{7}

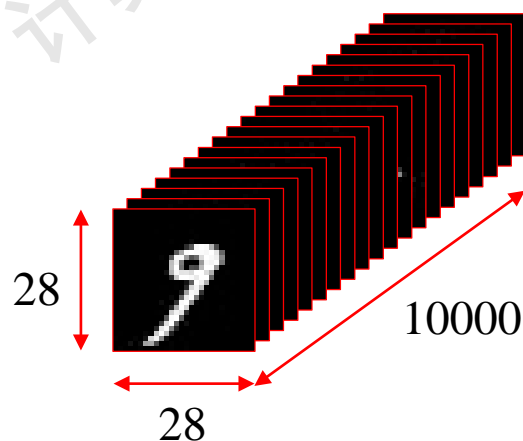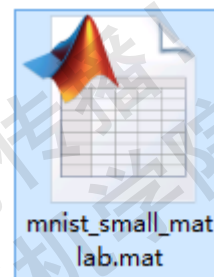# Outline

■ Brief Review of Backpropagation Algorithm

■ An Illustrating Example

■ Experiments

■ <span style="color:red">Assignment</span>

# Assignment

Implement the handwritten digits recognition by MATLAB using only one layer of external input.

*Thanks*