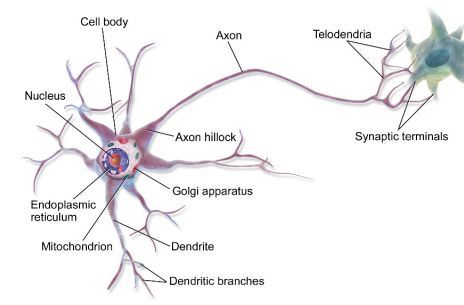
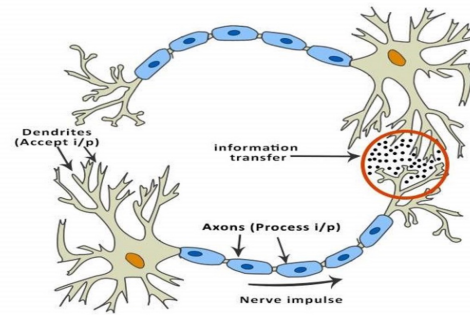


ANN and Handwritten Character Recognition

YAO ZHAO

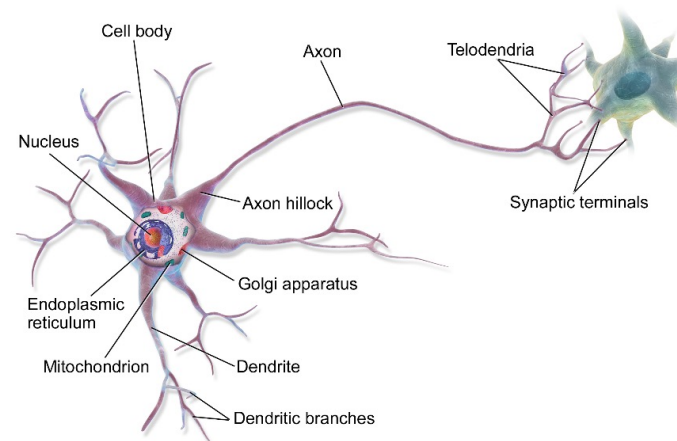
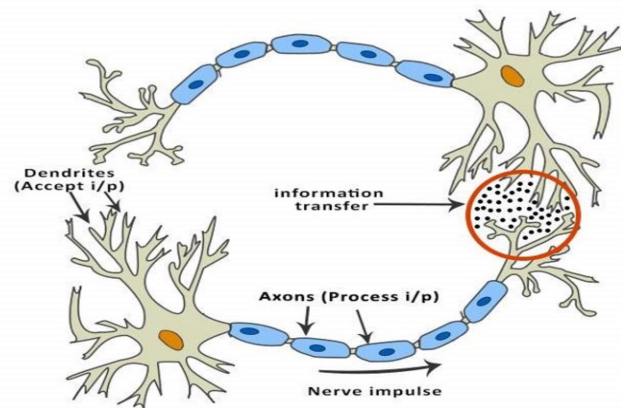
ANN

- ▶ The study of artificial neural networks was inspired by attempts to simulate biological neural systems (e. g. human brain).
- ▶ Basic structural and functional unit: **nerve cells called neurons**
- ▶ Work Mechanism
 - ▶ Different neurons are linked together via axons (轴突) and dendrite (树突)
 - ▶ When one neuron is excited after stimulation, it sends chemicals to the connected neurons, thereby changing the potential (电位) within these neurons.
 - ▶ If the potential (电位) of a neuron exceeds a “threshold”, then it is activated and send chemicals to other neurons.



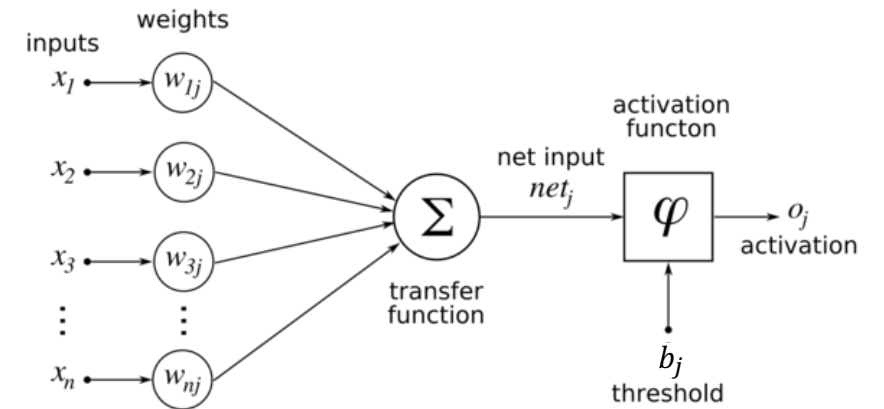
ANN

- ▶ A neuron is connected to the axons (轴突) of other neurons via dendrites (树突), which are extensions from the cell body of the neuron.
- ▶ The contact point between a dendrite (树突) and an axon (轴突) is called a synapse (突触).
- ▶ The human brain learns by changing the strength of the synaptic connection between neurons upon repeated stimulation by the same impulse.



Artificial Neuron Mathematical Model

- ▶ Input: x_i from the i-th neuron
- ▶ Weights: connection weights (synapse)
- ▶ Output: $o_j = \varphi(\sum_{i=1}^n w_{ij}x_i + b_j)$



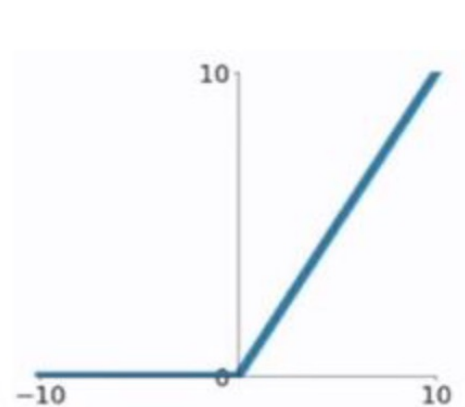
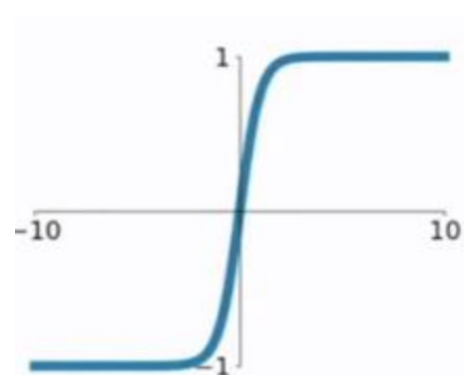
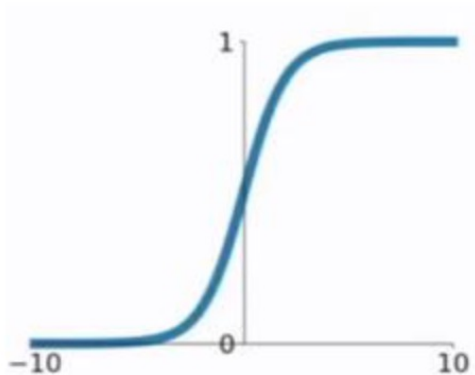
Artificial Neuron Model

- ▶ Output: $o_j = \varphi(\sum_{i=1}^n w_{ij}x_i + b_j)$
- ▶ Ideal activation function: step function but inapplicable
- ▶ Common activation function: sigmoid, tanh, ReLU

$$g(z) = \frac{1}{1 + e^{-z}}$$

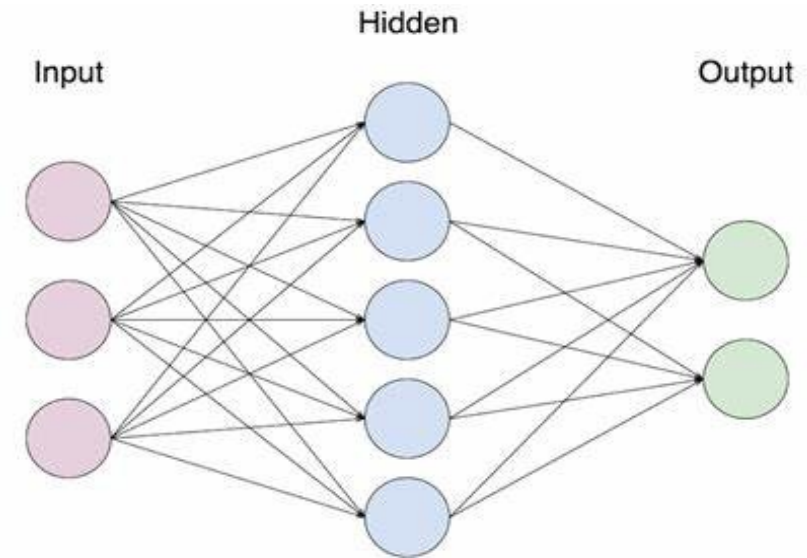
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g(z) = \begin{cases} z, & \text{if } z > 0 \\ 0, & \text{if } z < 0 \end{cases}$$



Artificial Neural Networks

- ▶ Consist of multiple artificial neurons
- ▶ Usually have the structure of an input layer, multiple hidden layers, an output layer
- ▶ The design of an NN or AutoML aims to design appropriate hidden layers and connection weights.



3-layer Feedforward neural networks

One Inference Process

$$x = x^{(0)} \rightarrow a^{(1)} \rightarrow x^{(1)} \rightarrow a^{(2)} \rightarrow x^{(2)} \rightarrow a^{(3)} \rightarrow x^{(3)} = y$$

Annotations:

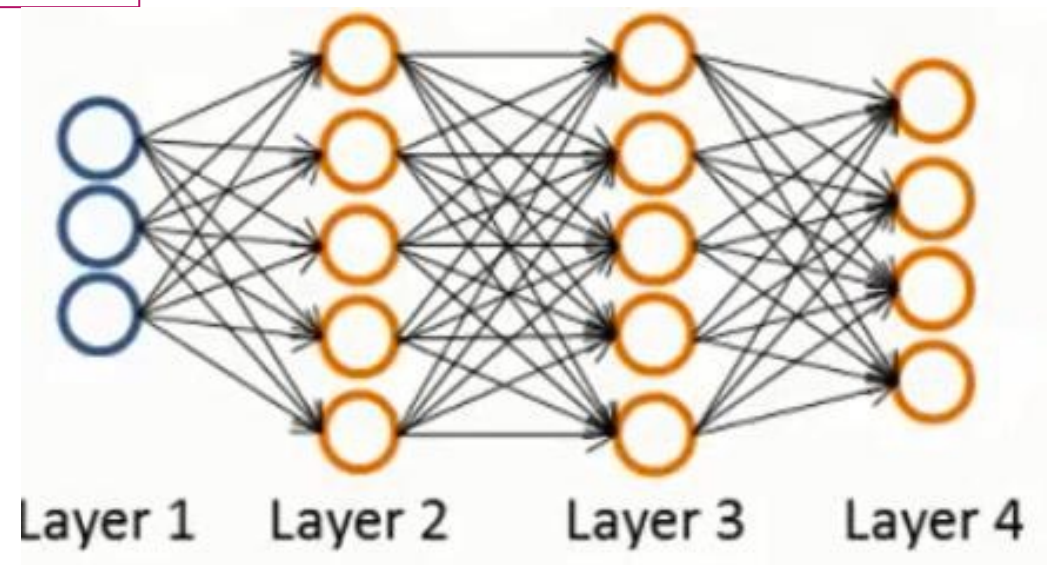
- $x_j^{(1)} = f(a_j^{(1)})$
- $x_k^{(2)} = f(a_k^{(2)})$
- $y_l = x_l^{(3)} = f(a_l^{(3)})$

$$a_j^{(1)} = \sum_{i=1}^D w_{ji}^{(1)} x_i^{(0)} + b_j^{(1)}$$

$$a_k^{(2)} = \sum_{j=1}^M w_{kj}^{(2)} x_j^{(1)} + b_k^{(2)}$$

Superscript of w: layer index

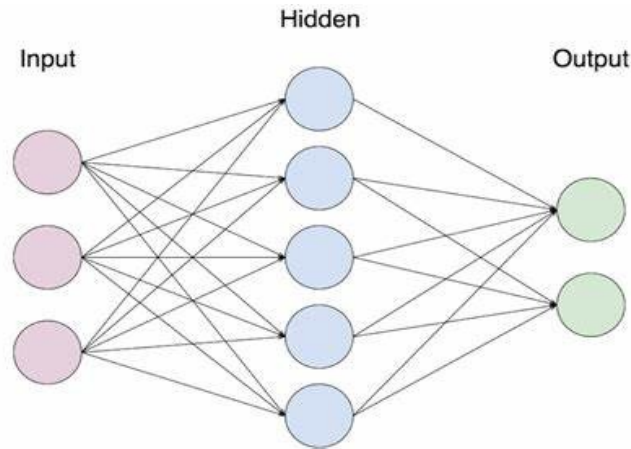
$$a_l^{(3)} = \sum_{k=1}^N w_{lk}^{(3)} x_k^{(2)} + b_l^{(3)}$$



Training of NN

- ▶ W and Threshold values decide the output of NN
- ▶ The training is to find appropriate values for W and Threshold
- ▶ The learning process is to tune weight matrix

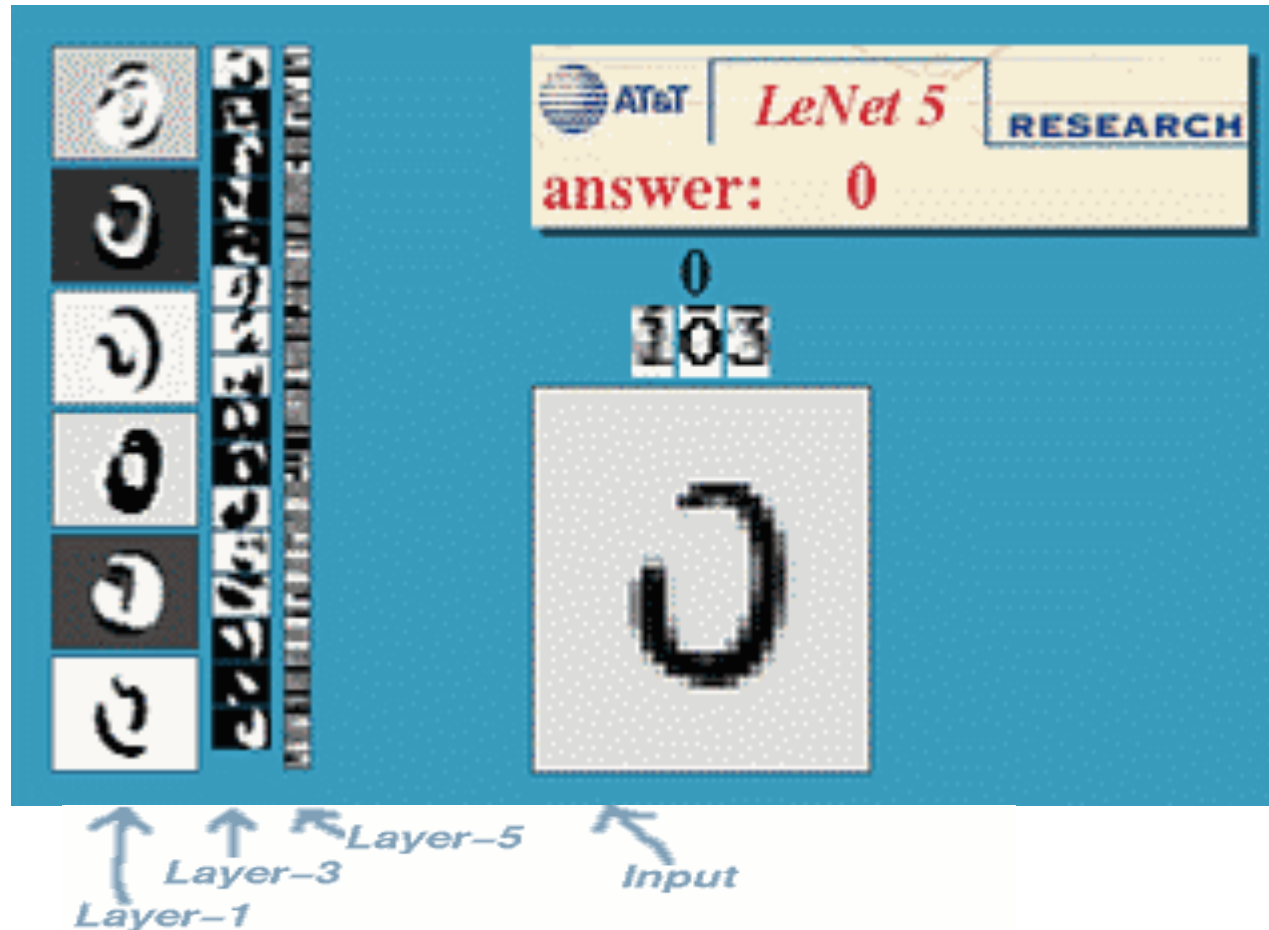
x1	x2	x3	l1	l2
1.0	0.1	0.3	1	0
0.1	1.5	1.2	1	0
1.1	1.1	2.0	0	1
0.2	0.2	0.3	0	1



$$\begin{aligned} \text{Error Function } (W, b) \\ &= \frac{1}{2} [(o1(W, b) - l1)^2 \\ &\quad + (o2(W, b) - l2)^2] \end{aligned}$$

Calculate the gradient for (W, b) , then tune them

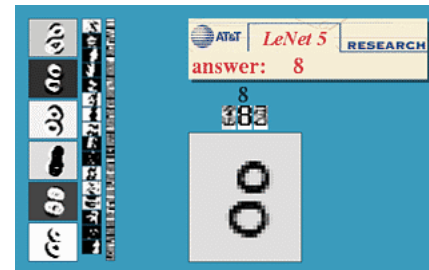
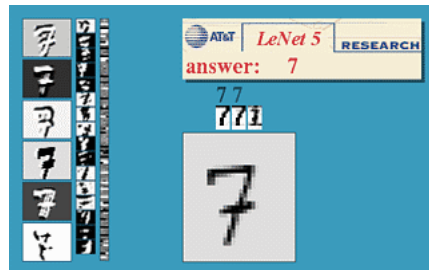
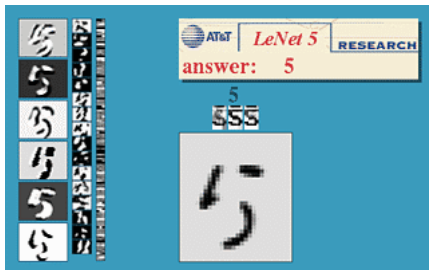
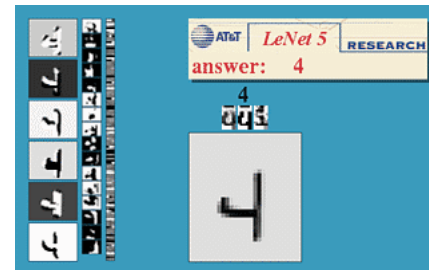
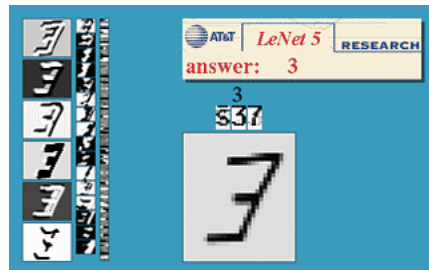
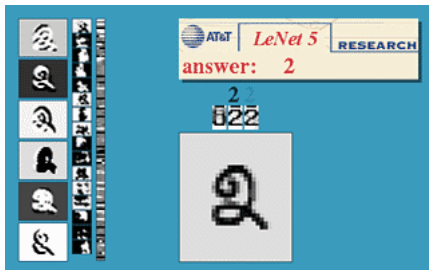
ANN Application: Handwritten character recognition



<http://yann.lecun.com/exdb/lenet/>

Performance of ANN

- ▶ Anti-interference ability, such as different sizes, digital distortion



The practice: Handwritten character recognition

- ▶ The dataset of this practice is a subset of MNIST

A 28*28 grayscale image:



Label:

7

One-hot vector t :

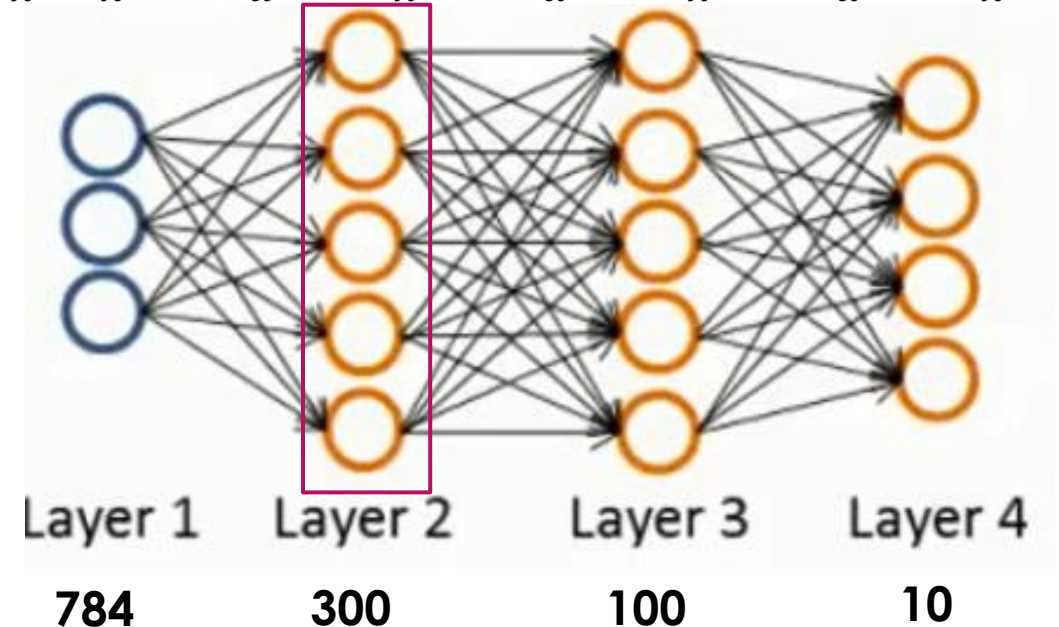
0	→	0
0	→	1
0	→	2
0	→	3
0	→	4
0	→	5
0	→	6
1	→	7
0	→	8
0	→	9

Handwritten character recognition

ANN network structure

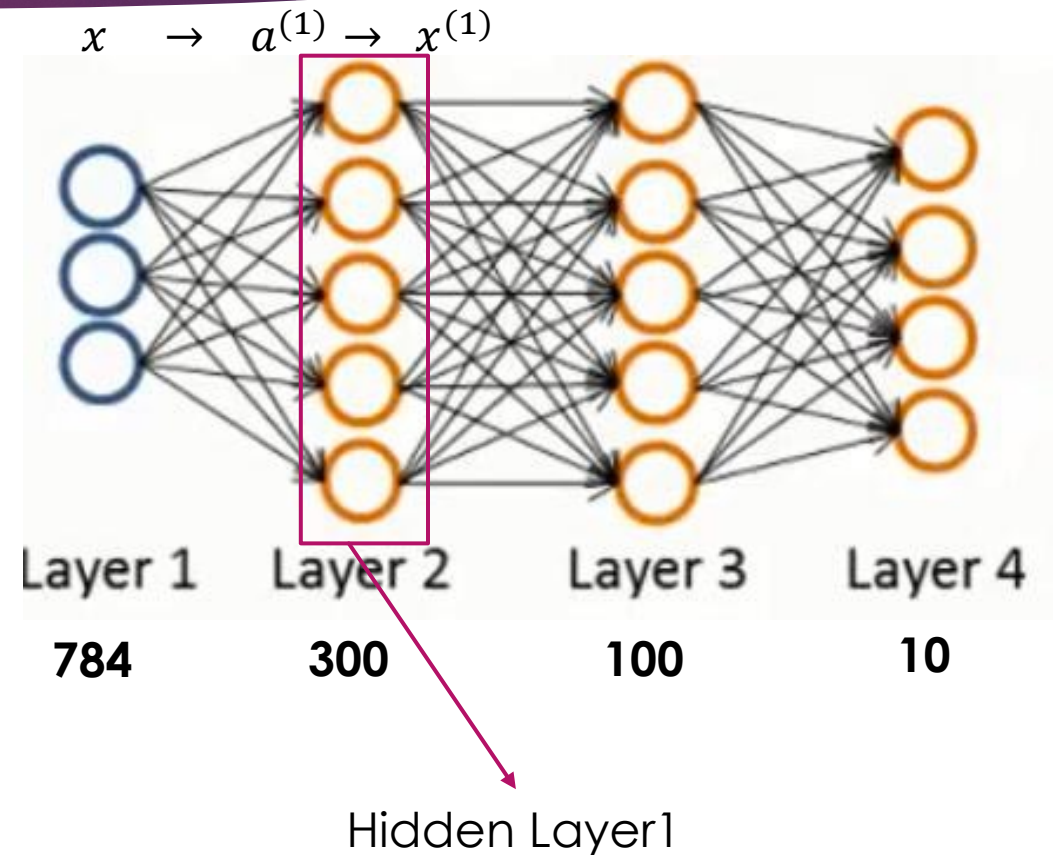
- ▶ Suppose 100 training samples, each sample: 28*28, feature: 784
- ▶ Layer 1: Input Layer, size D = 784
- ▶ Layer 2: Hidden Layer 1, size M = 300
- ▶ Layer 3: Hidden Layer 2, size N = 100
- ▶ Layer 4: Output Layer, size K = 10

$$x = x^{(0)} \rightarrow a^{(1)} \rightarrow x^{(1)} \rightarrow a^{(2)} \rightarrow x^{(2)} \rightarrow a^{(3)} \rightarrow x^{(3)} = y$$



ANN-Input Layer to Hidden Layer 1

- ▶ Input: 100(samples number) * 784(feature)
- ▶ Input to Hidden Layer 1: Weight Matrix: 784*300, b : 1*300
- ▶ Weighted Sum for j - th neurons of hidden layer 1:
$$a_j^{(1)} = \sum_{i=1}^D w_{ji}^{(1)} x_i^{(0)} + b_j^{(1)}$$
- ▶ After activation function, output of j -th neuron: $x_j^{(1)} = f(a_j^{(1)})$
- ▶ Output of Hidden Layer 1, $x^{(1)}$: 100(samples number)* 300 (feature) matrix
- ▶ f : Relu function

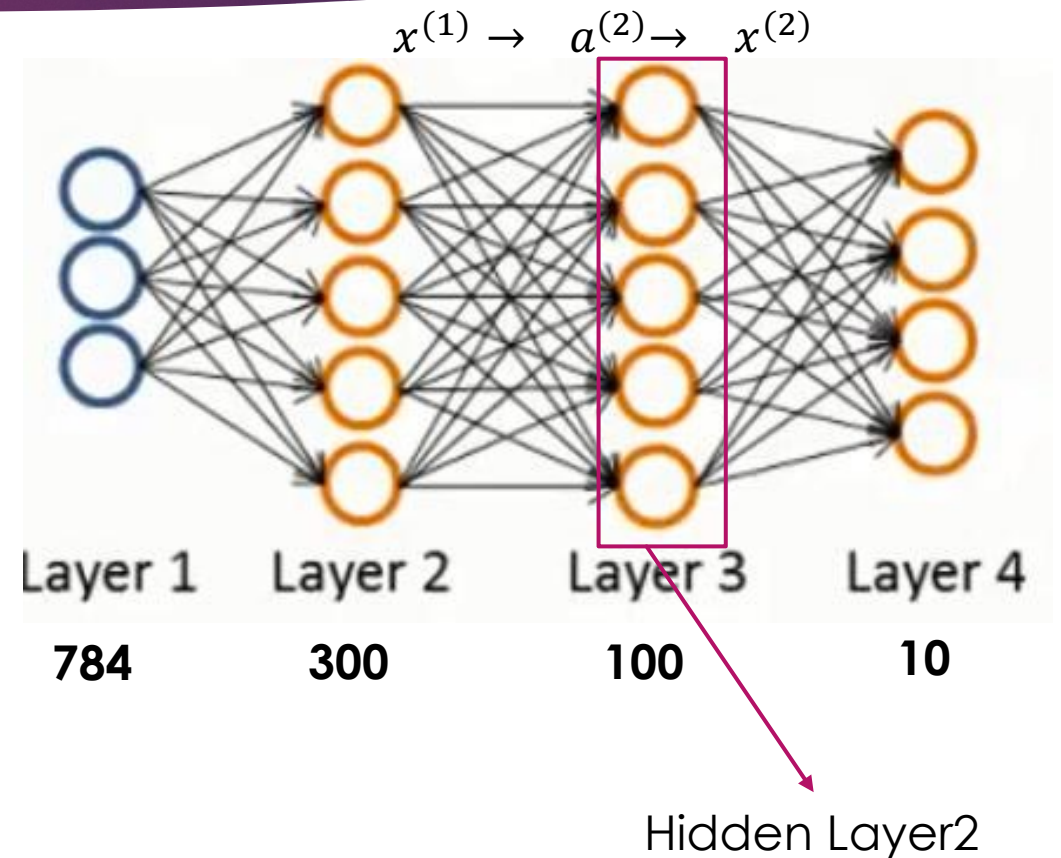


ANN-Hidden Layer 1 to Hidden Layer 2

- ▶ For Hidden Layer 2: Input data $x^{(1)}$: 100(samples number)* 300 (feature)
- ▶ Hidden Layer 1 to Hidden Layer 2: Weight Matrix: 300*100, b : 1*100
- ▶ Weighted Sum for k - th neurons of hidden layer 2:

$$a_k^{(2)} = \sum_{j=1}^M w_{kj}^{(2)} x_j^{(1)} + b_k^{(2)}$$

- ▶ After activation function, output of k-th neuron: $x_k^{(2)} = f(a_k^{(2)})$
- ▶ Output of Hidden Layer 2, $x^{(2)}$: 100(samples number) *100 (feature) matrix
- ▶ f : Relu function

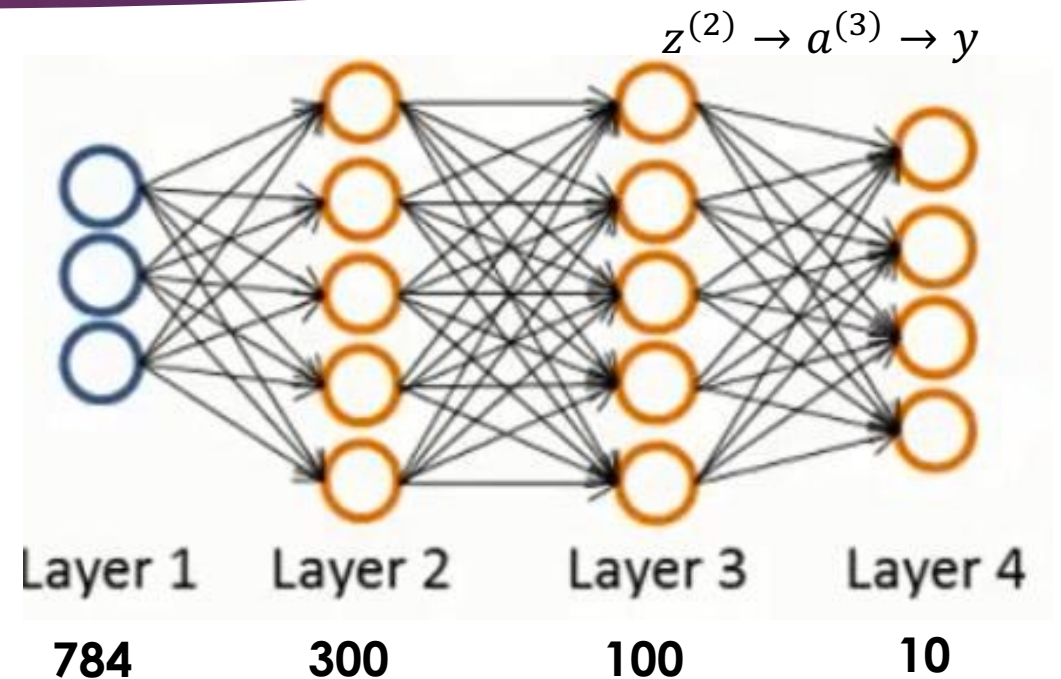


ANN-Hidden Layer 2 to Output Layer

- ▶ For output layer: Input data $x^{(2)}$: 100 (samples number) * 100 (feature)
- ▶ Hidden Layer 2 to Output Layer: Weight Matrix: 100*10, b : 10*1
- ▶ Weighted Sum for l - th neurons of output layer:

$$a_l^{(3)} = \sum_{k=1}^N w_{lk}^{(3)} x_k^{(2)} + b_l^{(3)}$$

- ▶ After activation function, output of l -th neuron: $y_l = x_l^{(3)} = \sigma(a_l^{(3)})$
- ▶ Output of Output Layer, y : 100(samples number) * 10 (feature) matrix
- ▶ σ : softmax function



Training of NN

- ▶ W and Threshold values decide the output of NN
- ▶ The training is to find appropriate values for W and Threshold so that the output is close to the true value
- ▶ Given the structure of NN, the learning process is to tune weight matrix in order to in order to minimize the difference between true value and prediction value.

Loss Function and Gradient Descent

- ▶ How to evaluate the difference? Loss Function
- ▶ How to tune weight matrix? optimization method (e. g. Gradient Descent)

As an Example

- ▶ W and b initialization:

- ▶ Xavier: is an initialization scheme for neural networks, make w uniform distributed over $(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}})$
(m: number of rows of the matrix, n: number of columns of the matrix)

- ▶ Loss Function: $E(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^{10} -t_{n,k} \ln y_{n,k}$

N is the samples number

- ▶ To avoid overfitting, add L2 Regularization term:

$$E_{total}(\mathbf{y}) = E(\mathbf{y}) + \frac{\lambda}{2} ||W||_2^2$$

W contains all the weights.

- ▶ Gradient Descent: $\min f(x) \rightarrow x(t+1) = x(t) - \eta * f'(x(t))$ (η : learning rate)

Loss

A 28*28 grayscale image:



Label:

7

t :

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

y :

$\begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$

$$E(y) = -\ln(0.1)$$

Gradient Descent

$$\begin{aligned} &\blacktriangleright \frac{\partial E_{total}(W^{(t)}, b^{(t)})}{\partial W^{(t)}} \\ &= \frac{\partial (\frac{1}{N} \sum_{n=1}^N E_n(W^{(t)}, b^{(t)}) + \frac{\lambda}{2} \|W^{(t)}\|^2)}{\partial W^{(t)}} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial E_n(W^{(t)}, b^{(t)})}{\partial W^{(t)}} + \lambda W^{(t)} \\ &\blacktriangleright \frac{\partial E_{total}(W^{(t)}, b^{(t)})}{\partial b^{(t)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial E_n(W^{(t)}, b^{(t)})}{\partial b^{(t)}} \end{aligned}$$

Chain Rule Review

- ▶ Chain rule: to find the derivative of a composite function
- ▶ If $h(x)=f(g(x))$, then $h'(x)=f'(g(x))g'(x)$
- ▶ E.g.: $f(x)=2x+2$, $g(x)=3x+3$, $g(f(x))$ is a composite function, $g'(f(x))=6$

One Inference Process

$$x = x^{(0)} \rightarrow a^{(1)} \rightarrow x^{(1)} \rightarrow a^{(2)} \rightarrow x^{(2)} \rightarrow a^{(3)} \rightarrow x^{(3)} = y$$

Annotations:

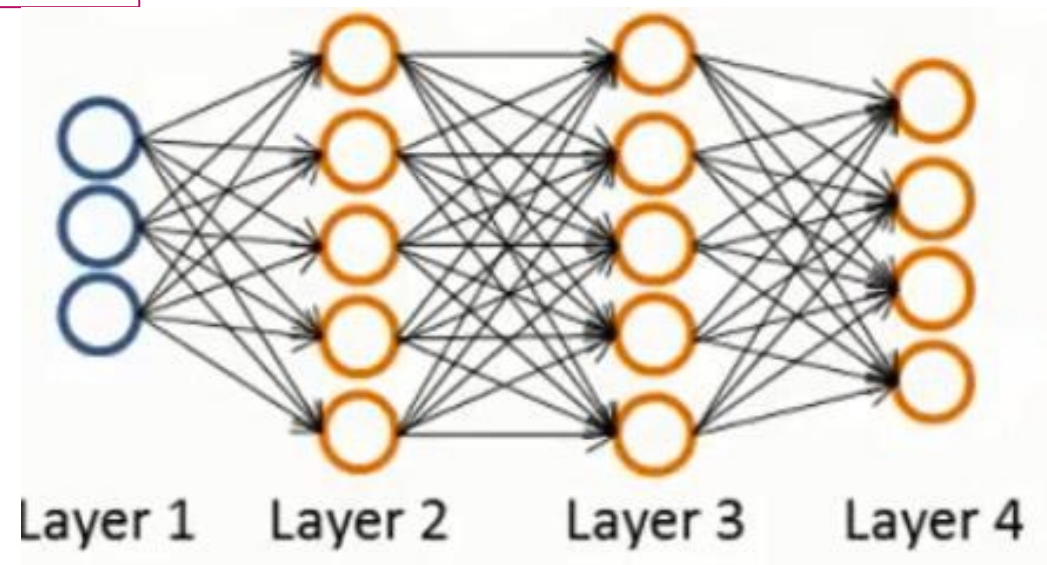
- $x_j^{(1)} = f(a_j^{(1)})$
- $x_k^{(2)} = f(a_k^{(2)})$
- $y_l = x_l^{(3)} = f(a_l^{(3)})$

$$a_j^{(1)} = \sum_{i=1}^D w_{ji}^{(1)} x_i^{(0)} + b_j^{(1)}$$

$$a_k^{(2)} = \sum_{j=1}^M w_{kj}^{(2)} x_j^{(1)} + b_k^{(2)}$$

Superscript of w: layer index

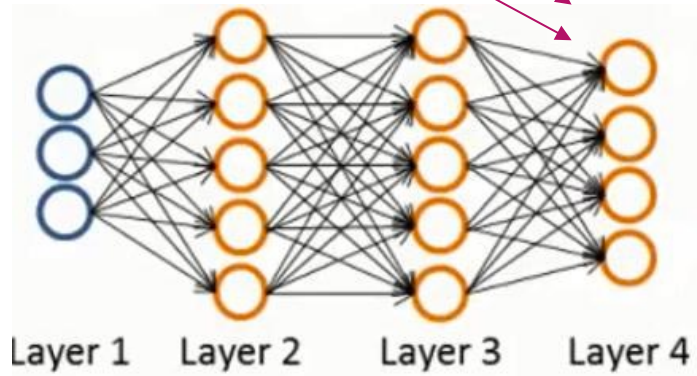
$$a_l^{(3)} = \sum_{k=1}^N w_{lk}^{(3)} x_k^{(2)} + b_l^{(3)}$$



Derivative of E on $w^{(3)}$, $b^{(3)}$

► $\frac{\partial E}{\partial W^{(3)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial W^{(3)}}$

► $\frac{\partial E}{\partial b^{(3)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial b^{(3)}}$



$$x_j^{(1)} = f(a_j^{(1)})$$

$$x_k^{(2)} = f(a_k^{(2)})$$

$$y_l = x_l^{(3)} = f(a_l^{(3)})$$

$$x = x^{(0)} \rightarrow a^{(1)} \rightarrow x^{(1)} \rightarrow a^{(2)} \rightarrow x^{(2)} \rightarrow a^{(3)} \rightarrow x^{(3)} = y$$

$$a_j^{(1)} = \sum_{i=1}^D w_{ji}^{(1)} x_i^{(0)} + b_j^{(1)}$$

$$a_k^{(2)} = \sum_{j=1}^M w_{kj}^{(2)} x_j^{(1)} + b_k^{(2)}$$

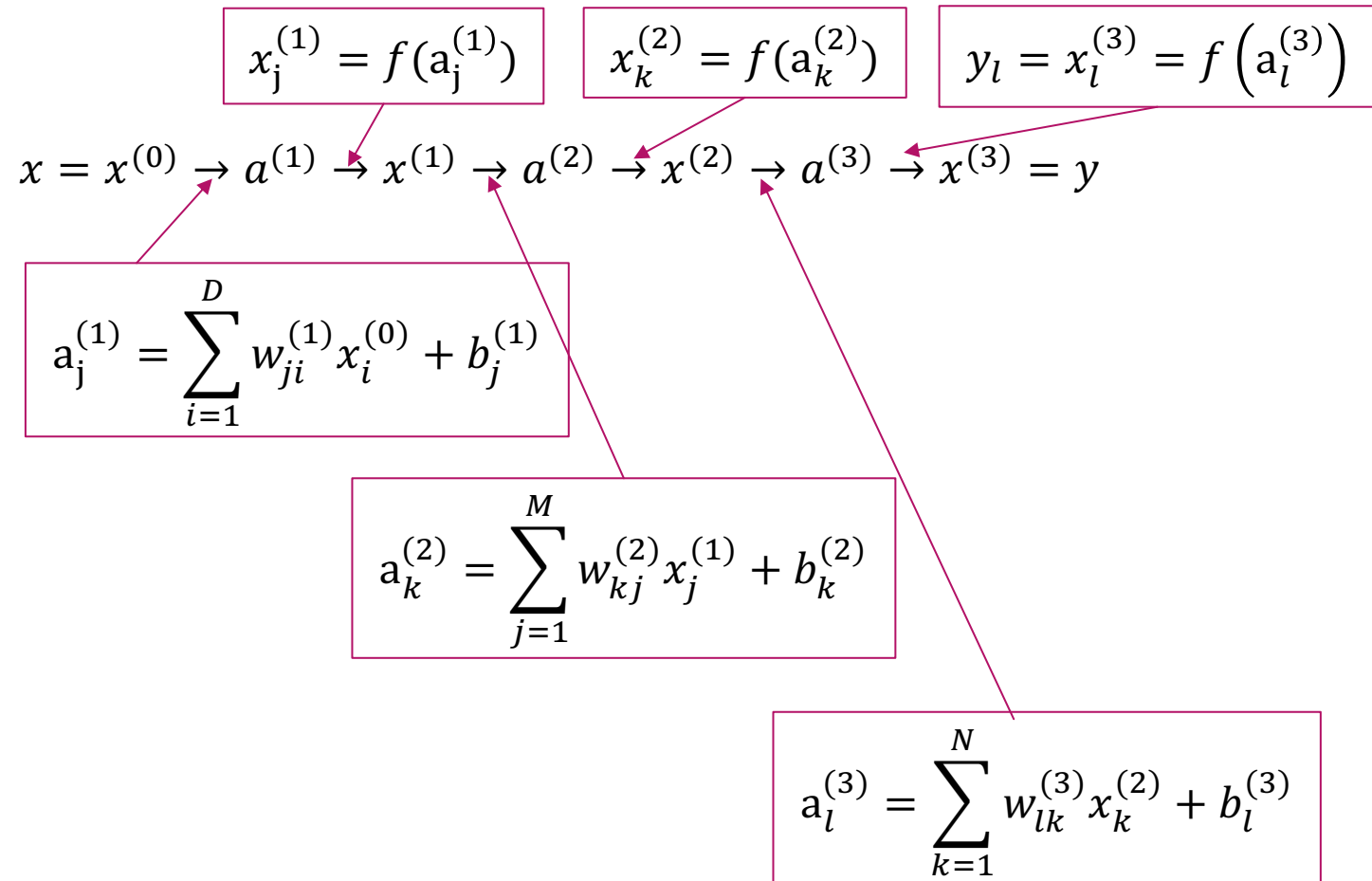
$$a_l^{(3)} = \sum_{k=1}^N w_{lk}^{(3)} x_k^{(2)} + b_l^{(3)}$$

Derivative of E on $w^{(3)}, b^{(3)}$

$$\triangleright \frac{\partial E}{\partial w^{(3)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial w^{(3)}}$$

$$\triangleright \frac{\partial E}{\partial b^{(3)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial b^{(3)}}$$

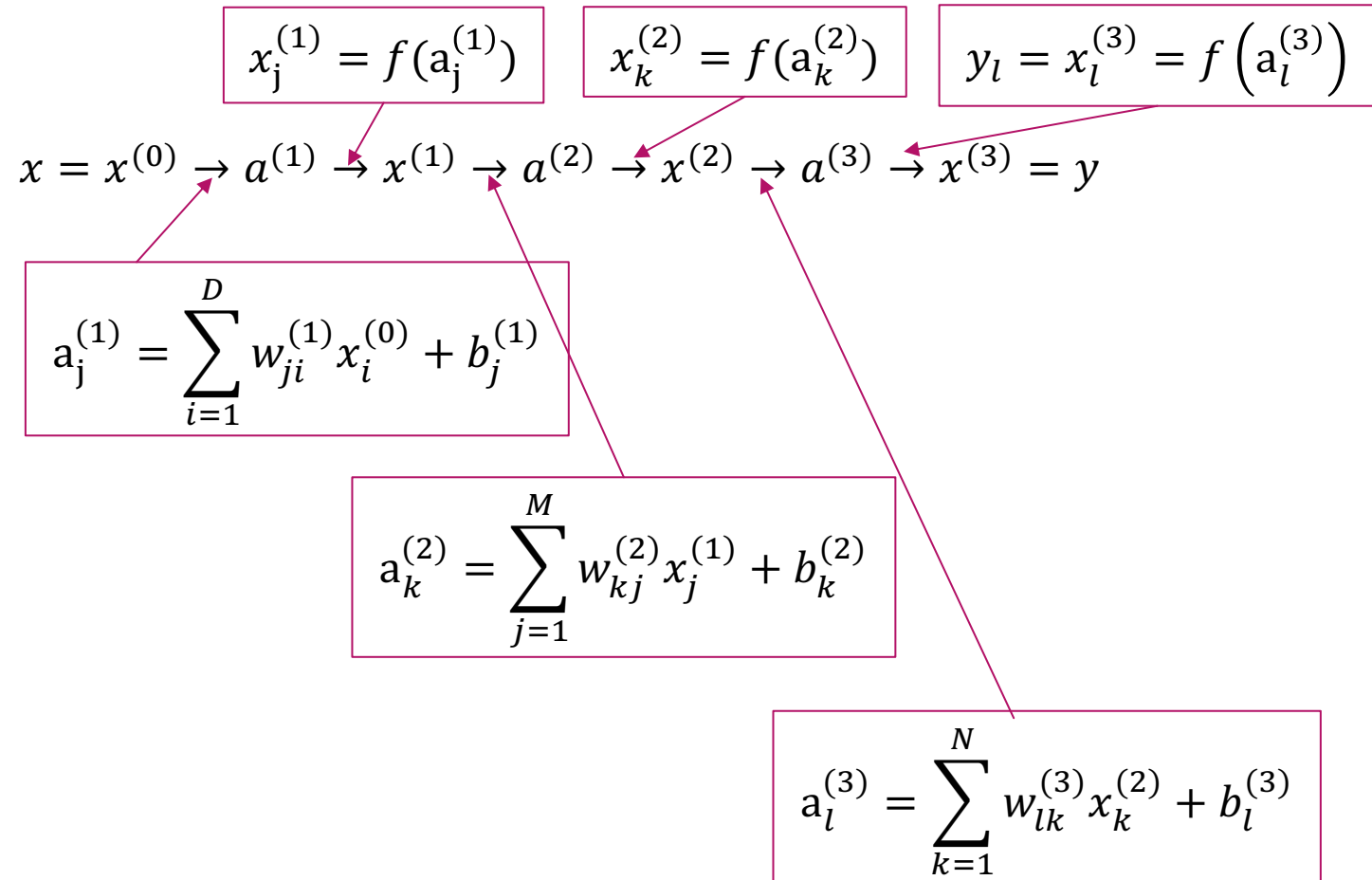
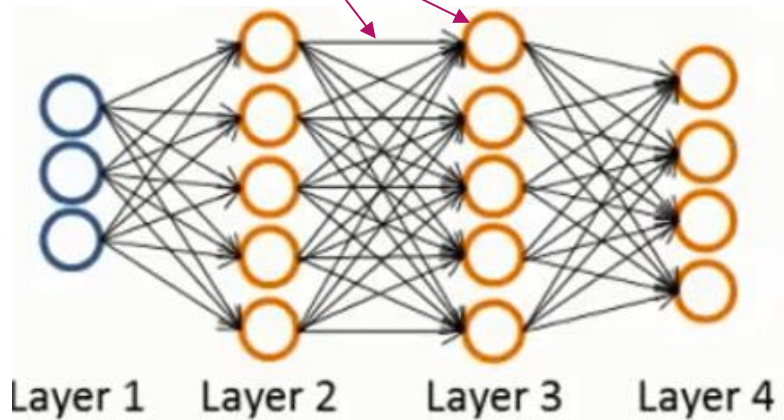
$$\text{Let } \delta^{(3)} = \frac{\partial E}{\partial a^{(3)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}}$$



Derivative of E on $w^{(2)}$, $b^{(2)}$

$$\frac{\partial E}{\partial W^{(2)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial W^{(2)}}$$

$$\frac{\partial E}{\partial b^{(2)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial b^{(2)}}$$

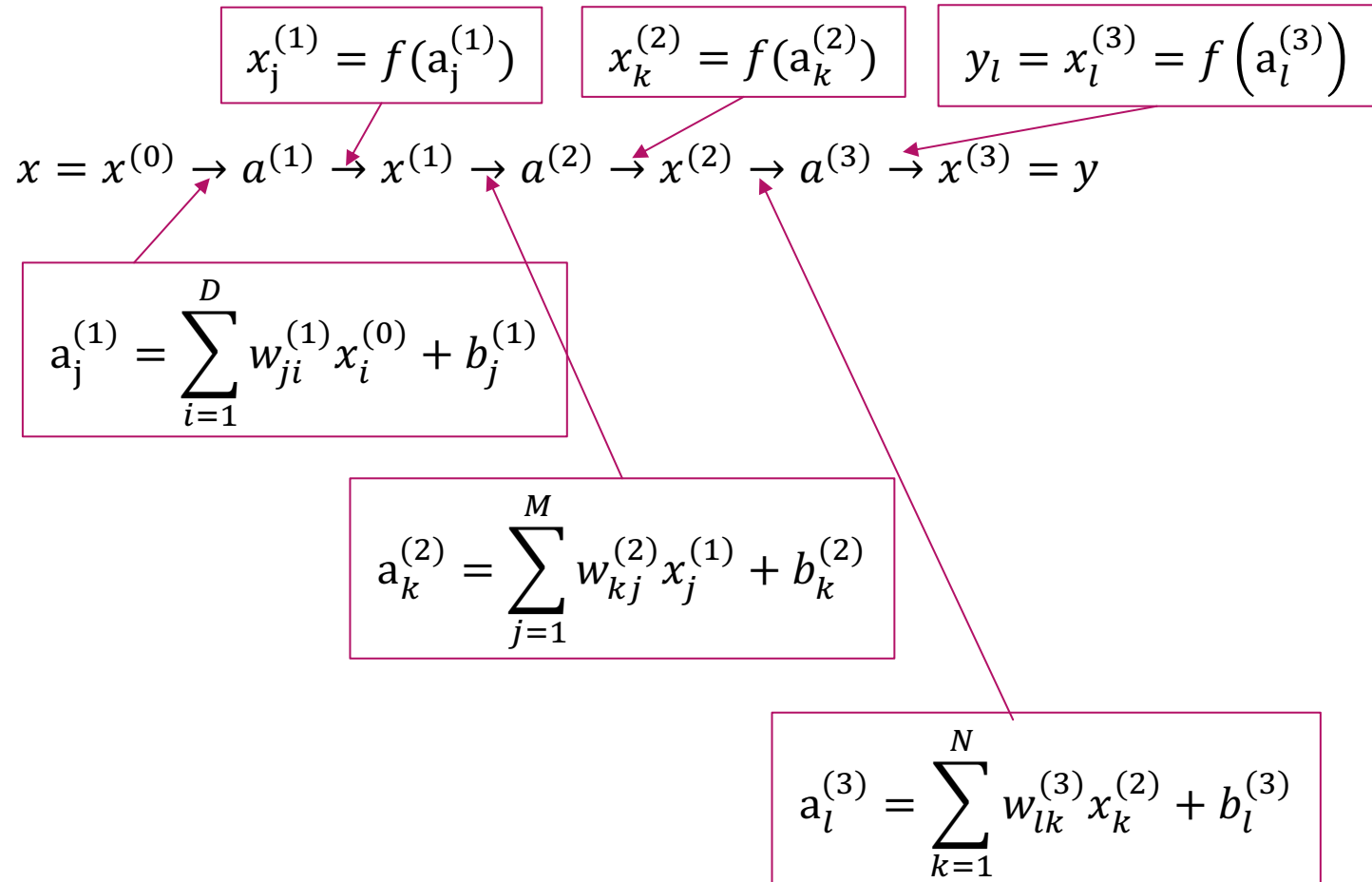


Derivative of E on $w^{(2)}, b^{(2)}$

$$\triangleright \frac{\partial E}{\partial W^{(2)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial W^{(2)}}$$

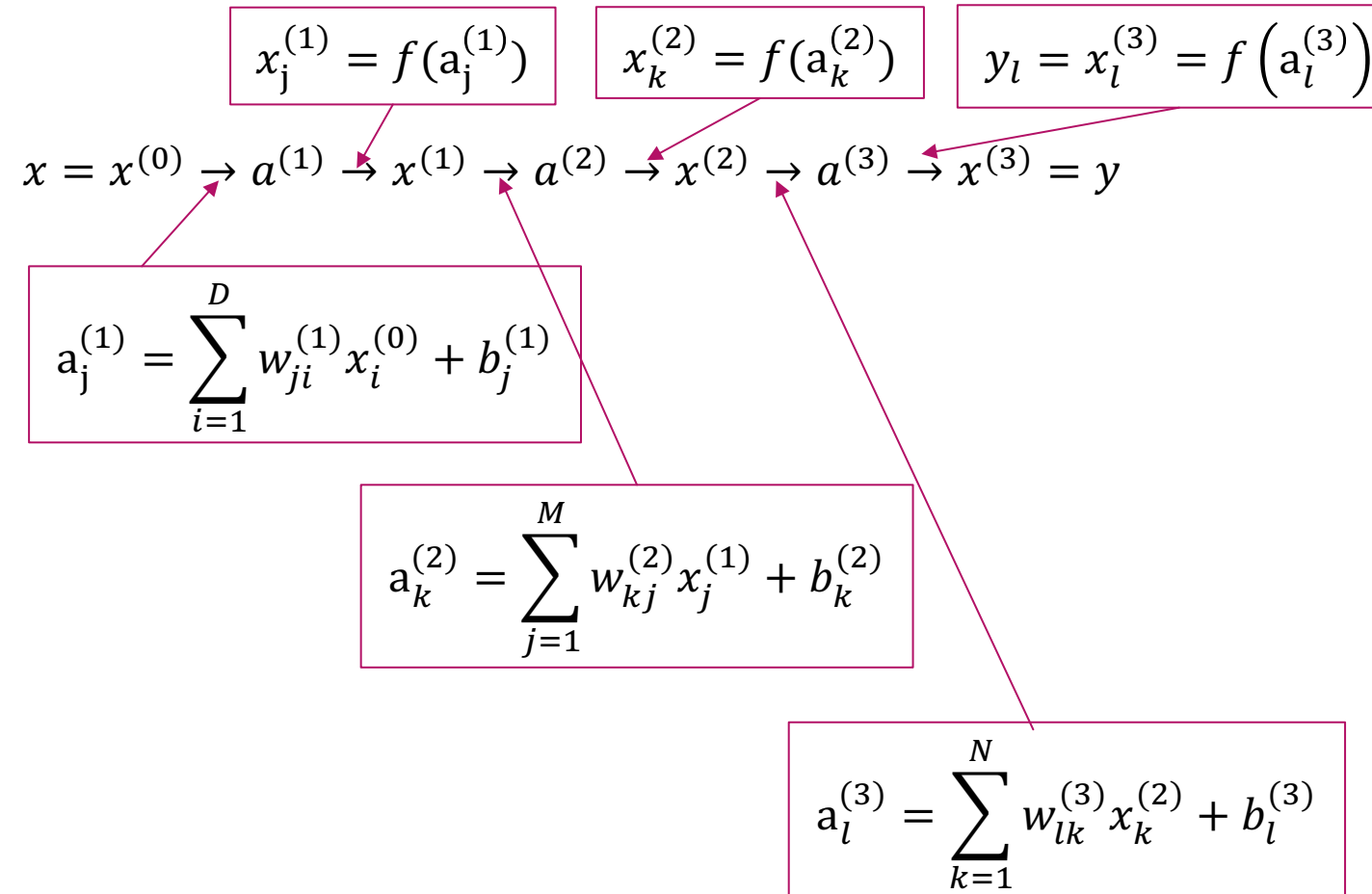
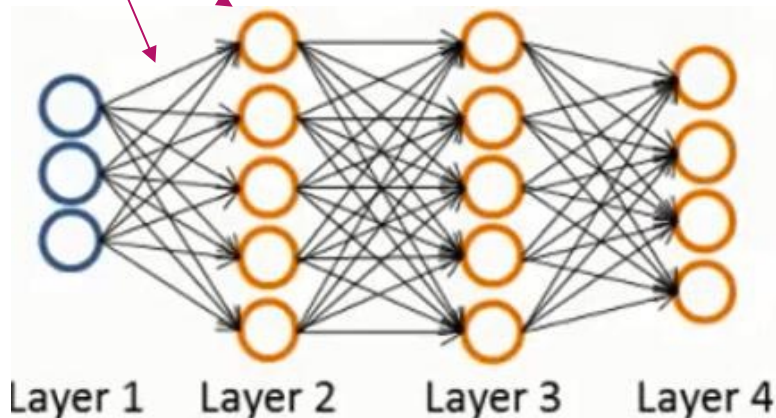
$$\triangleright \frac{\partial E}{\partial b^{(2)}} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial b^{(2)}}$$

$$\text{Let } \delta^{(2)} = \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} = \frac{\partial E}{\partial a^{(2)}}$$



Derivative of E on $w^{(1)}, b^{(1)}$

$$\begin{aligned} \triangleright \frac{\partial E}{\partial W^{(1)}} &= \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(1)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial W^{(1)}} \\ \triangleright \frac{\partial E}{\partial b^{(1)}} &= \frac{\partial E(y)}{\partial y} \frac{\partial y(a^{(3)})}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial x^{(2)}} \frac{\partial x^{(2)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(1)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial b^{(1)}} \end{aligned}$$



Gradient Descent Process (BP Algorithm)

1. Set values to max_iterations and η
2. Randomly generate $W^{(1)}$, $W^{(2)}$, $W^{(3)}$ and $b^{(1)}$, $b^{(2)}$, $b^{(3)}$
3. For $t = 1$ to max_iterations
4. Calculate the derivatives: $\frac{\partial E_{\text{total}}}{\partial W^{(1\sim3)}}$ and $\frac{\partial E_{\text{total}}}{\partial b^{(1\sim3)}}$
5. Update $W^{(1\sim3)} = W^{(1\sim3)} - \eta \frac{\partial E_{\text{total}}}{\partial W^{(1\sim3)}}$, $B^{(1\sim3)} = b^{(1\sim3)} - \eta \frac{\partial E_{\text{total}}}{\partial b^{(1\sim3)}}$
6. EndFor

Note: **In this practice, E is a batch of randomly selected samples.**

(E can be loss for a single sample, a batch of randomly selected samples, or all training samples.)