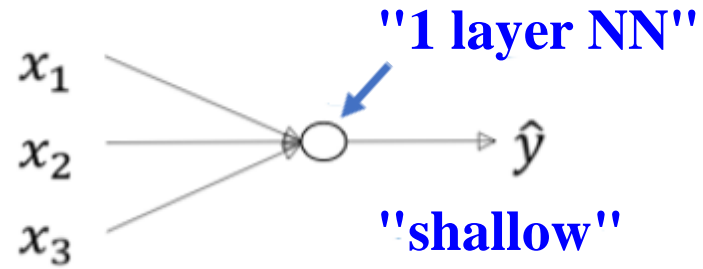


# **Deep Neural Network**

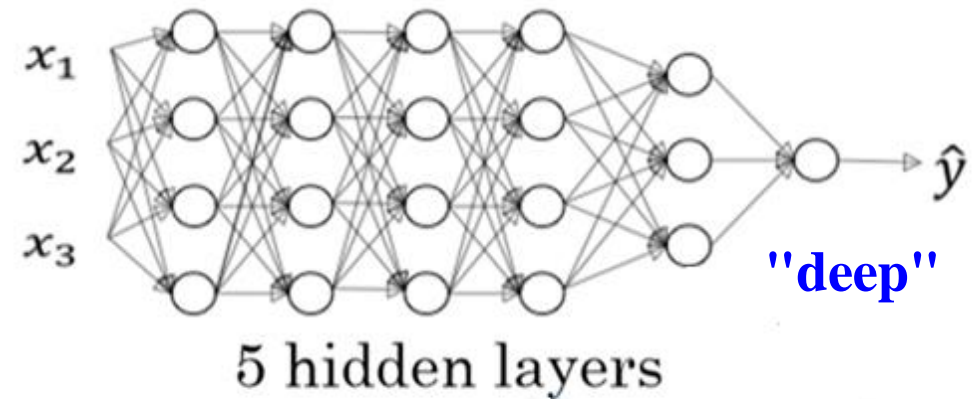
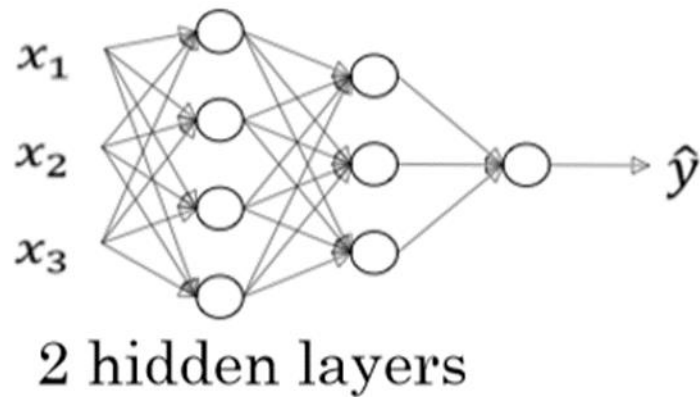
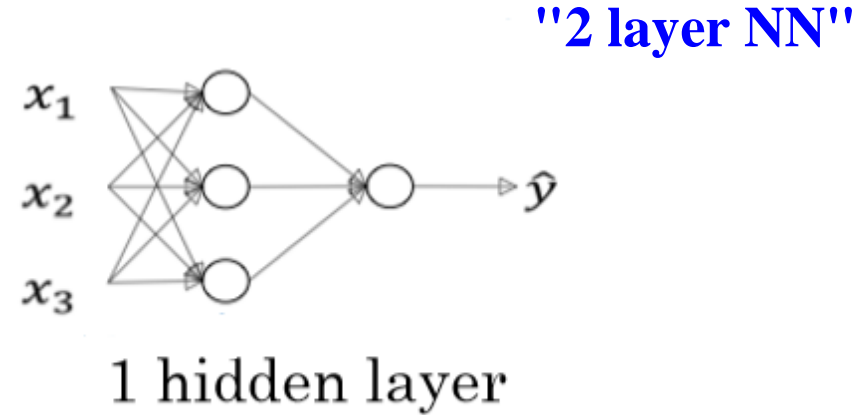
---

**Deep L-layer  
Neural Network**

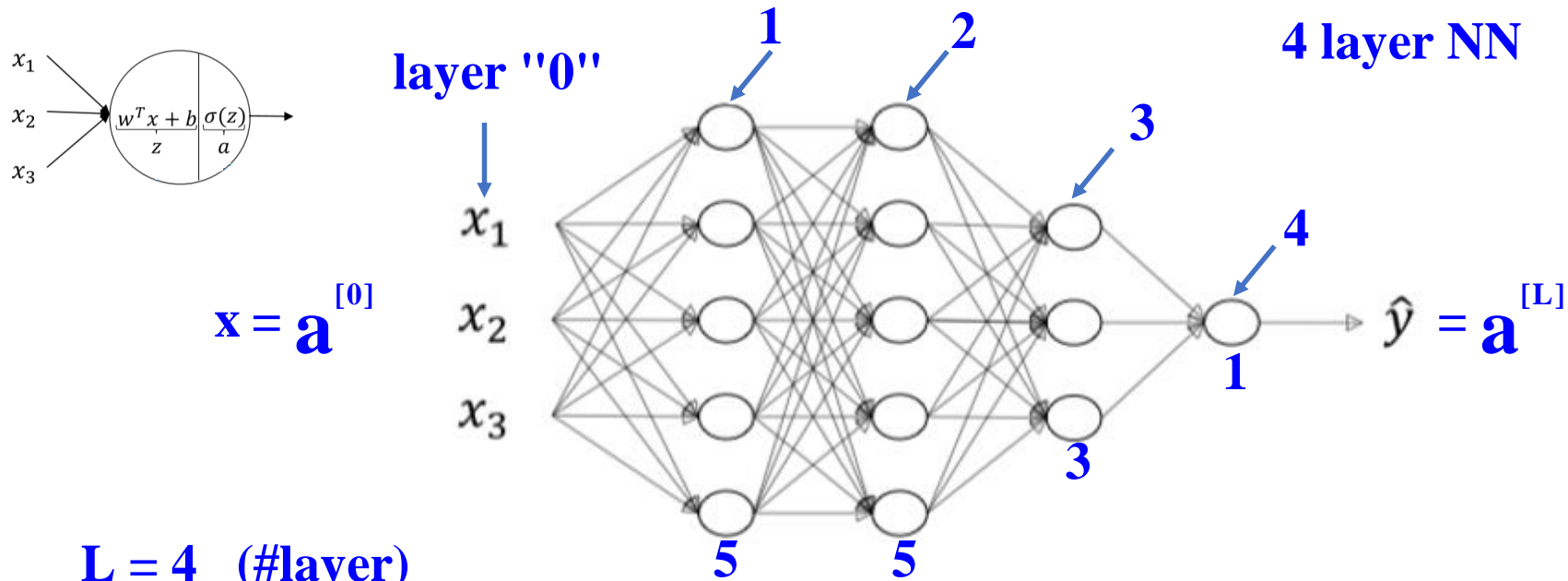
# What is a deep neural network



logistic regression



# Deep neural network notation



$L = 4$  (#layer)

$\mathbf{n}^{[l]} = \text{\#units in layer } l$

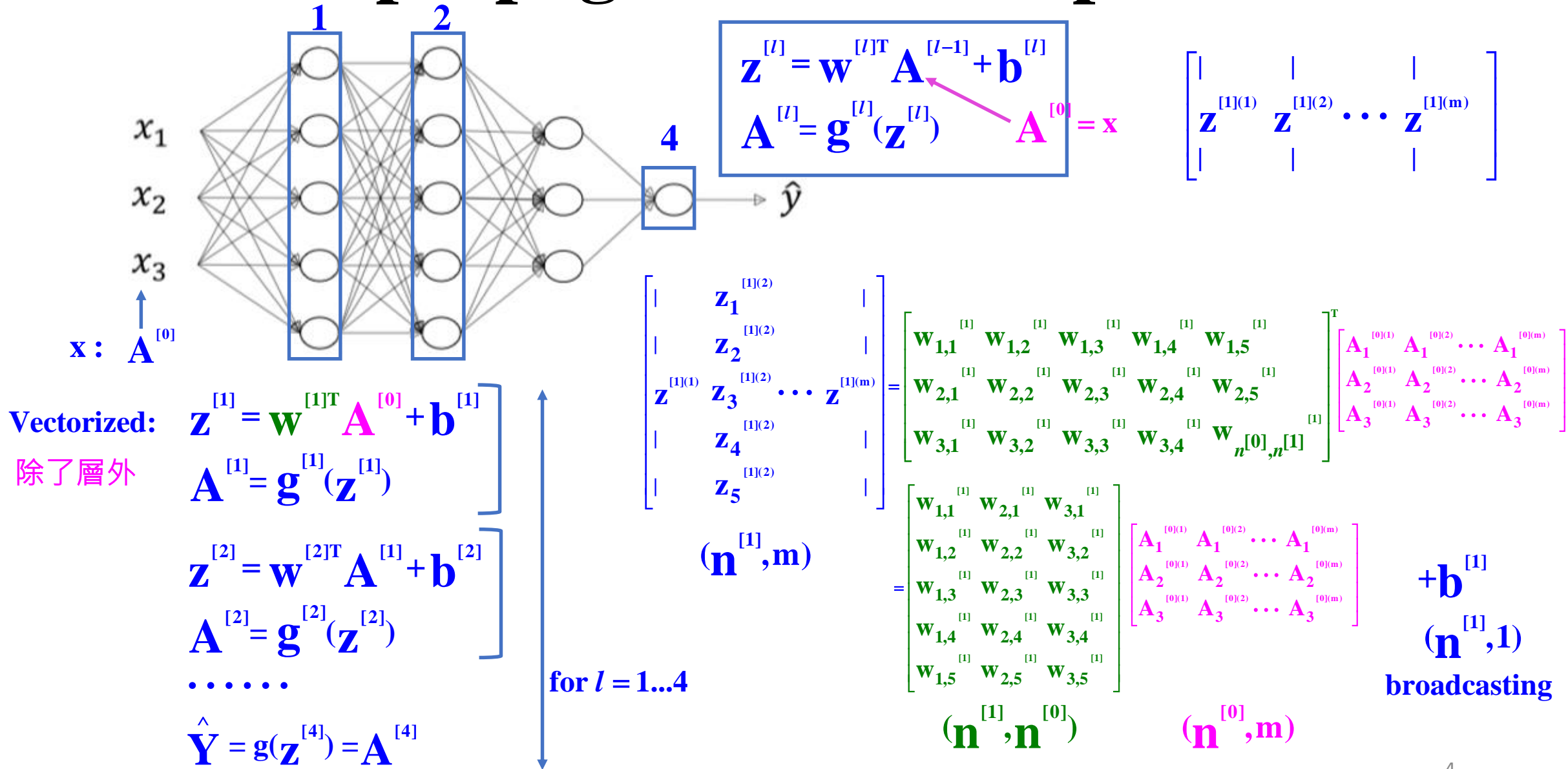
$\mathbf{a}^{[l]} = \text{activation in layer } l$

$\mathbf{a}^{[l]} = \mathbf{g}^{[l]}(\mathbf{z}^{[l]}), \quad \mathbf{w}^{[l]} = \text{weights for } \mathbf{z}^{[l]}$   
 $\mathbf{b}^{[l]}$

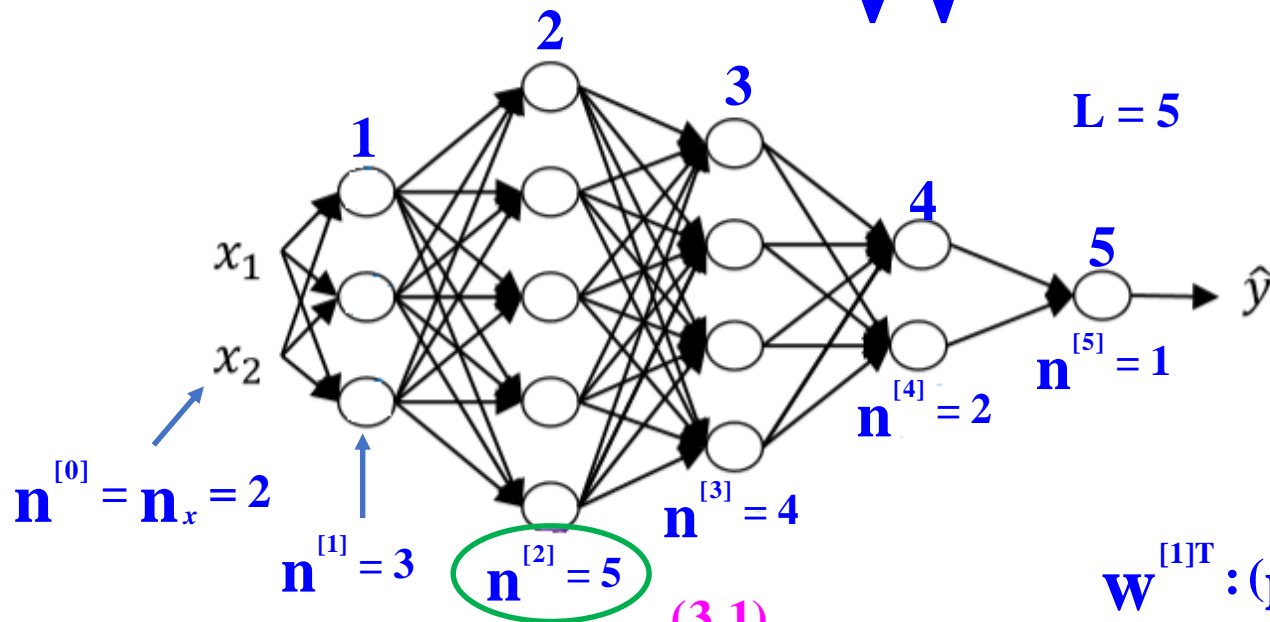
$\mathbf{n}^{[1]} = 5, \mathbf{n}^{[2]} = 5, \mathbf{n}^{[3]} = 3, \mathbf{n}^{[4]} = \mathbf{n}^{[L]} = 1$

$\mathbf{n}^{[0]} = \mathbf{n}_x = 3$

# Forward propagation in a deep network



# Parameters: $\mathbf{W}^l$ and $\mathbf{b}^l$



$$\left[ \begin{array}{l} \mathbf{w}^{[L]T} : (n^{[L]}, n^{[L-1]}) \\ \mathbf{b}^{[L]} : (n^{[L]}, 1) \\ d\mathbf{w}^{[L]T} : (n^{[L]}, n^{[L-1]}) \\ d\mathbf{b}^{[L]} : (n^{[L]}, 1) \end{array} \right]$$

有d無d,  
同矩陣大小(以此記)

$$\mathbf{z}^{[1]} = \mathbf{w}^{[1]T} \mathbf{x} + \mathbf{b}^{[1]}$$

$(3,1) = (3,2) (2,1)$

$(\mathbf{n}^{[1]}, 1) = (\mathbf{n}^{[1]}, \mathbf{n}^{[0]})(\mathbf{n}^{[0]}, 1)$

Assume 1 sample

$$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$$\mathbf{w}^{[1]T} : (\mathbf{n}^{[1]}, \mathbf{n}^{[0]})$$

$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}, \quad \mathbf{w}^{[2]T} : (5,3), (\mathbf{n}^{[2]}, \mathbf{n}^{[1]})$$

$\uparrow \quad \uparrow \quad \uparrow$

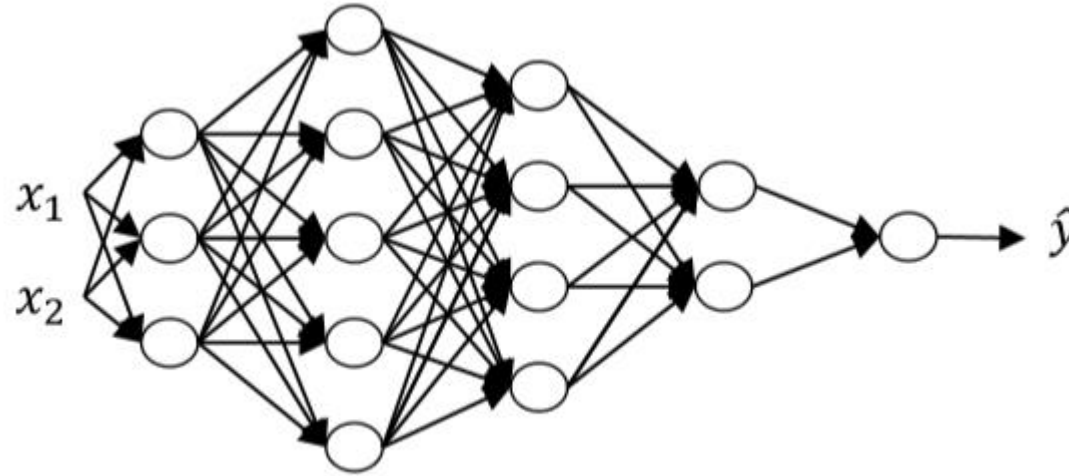
$(5,1) \quad (5,3) \quad (3,1) \quad (5,1)$

$(\mathbf{n}^{[2]}, 1)$

$$\mathbf{w}^{[3]T} : (4,5)$$

$$\mathbf{w}^{[4]T} : (2,4), \quad \mathbf{w}^{[5]T} : (1,2)$$

# Vectorized implementation



$$\underset{(\mathbf{n}^{[1]},1)}{\mathbf{z}^{[1]}} = \underset{(\mathbf{n}^{[1]},\mathbf{n}^{[0]})}{\mathbf{W}^{[1]}} \underset{(\mathbf{n}^{[0]},1)}{\mathbf{x}} + \underset{(\mathbf{n}^{[1]},1)}{\mathbf{b}^{[1]}}$$

$$\begin{bmatrix} | & | & & | \\ \mathbf{z}^{[1](1)} & \mathbf{z}^{[1](2)} & \cdots & \mathbf{z}^{[1](m)} \\ | & | & & | \end{bmatrix}$$

$$\underset{(\mathbf{n}^{[1]},\mathbf{m})}{\mathbf{Z}^{[1]}} = \underset{(\mathbf{n}^{[1]},\mathbf{n}^{[0]})}{\mathbf{W}^{[1]}} \underset{(\mathbf{n}^{[0]},\mathbf{m})}{\mathbf{X}} + \underset{(\mathbf{n}^{[1]},\mathbf{m})}{\mathbf{b}^{[1]}}$$

$$\begin{array}{c} \frac{(\mathbf{n}^{[1]},1)}{\downarrow} \\ (\mathbf{n}^{[1]},\mathbf{m}) \end{array}$$

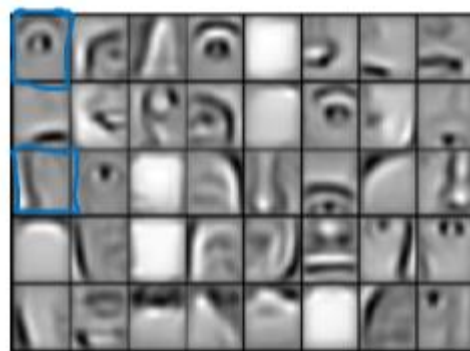
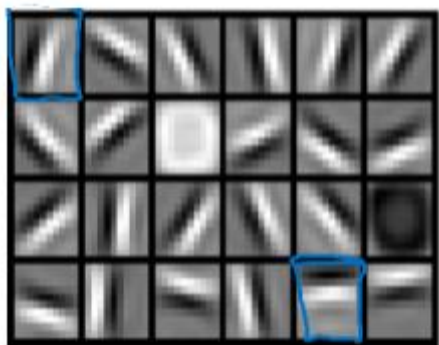
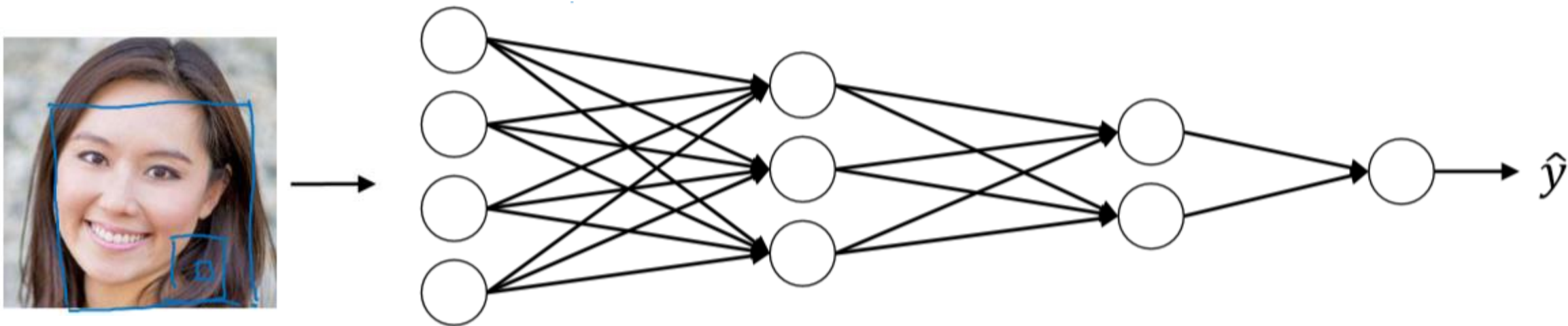
$$\mathbf{z}^{[l]}, \mathbf{a}^{[l]} : (\mathbf{n}^{[l]}, 1)$$

$$\mathbf{Z}^{[l]}, \mathbf{A}^{[l]} : (\mathbf{n}^{[l]}, \mathbf{m})$$

$$l = 0 \quad \mathbf{A}^{[0]} = \mathbf{X} = (\mathbf{n}^{[0]}, \mathbf{m})$$

$$d\mathbf{Z}^{[l]}, d\mathbf{A}^{[l]} : (\mathbf{n}^{[l]}, \mathbf{m})$$

# Intuition about deep representation

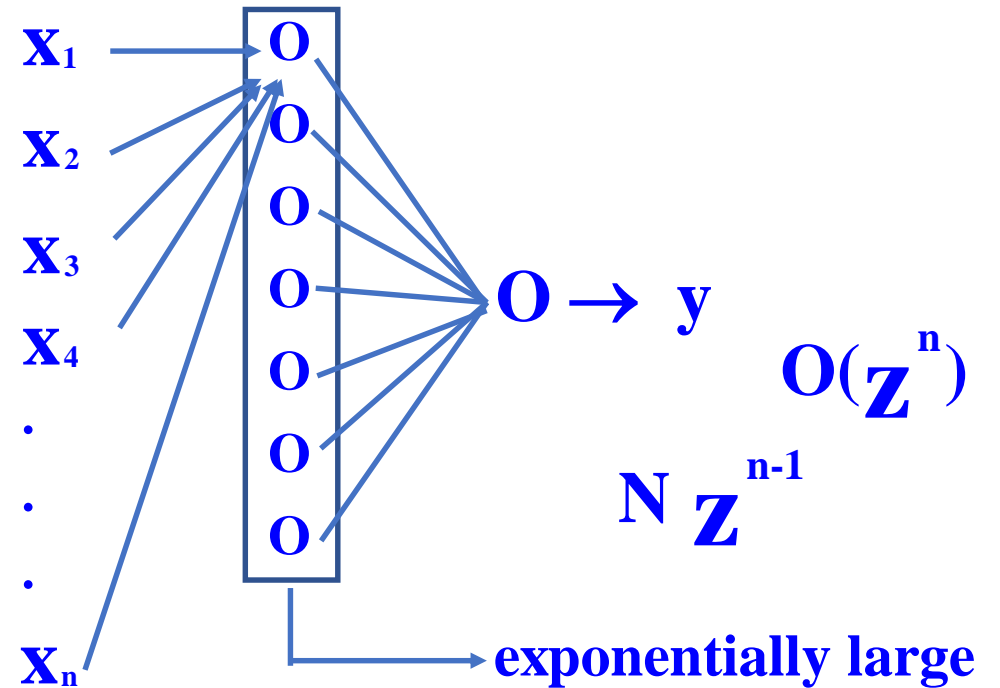
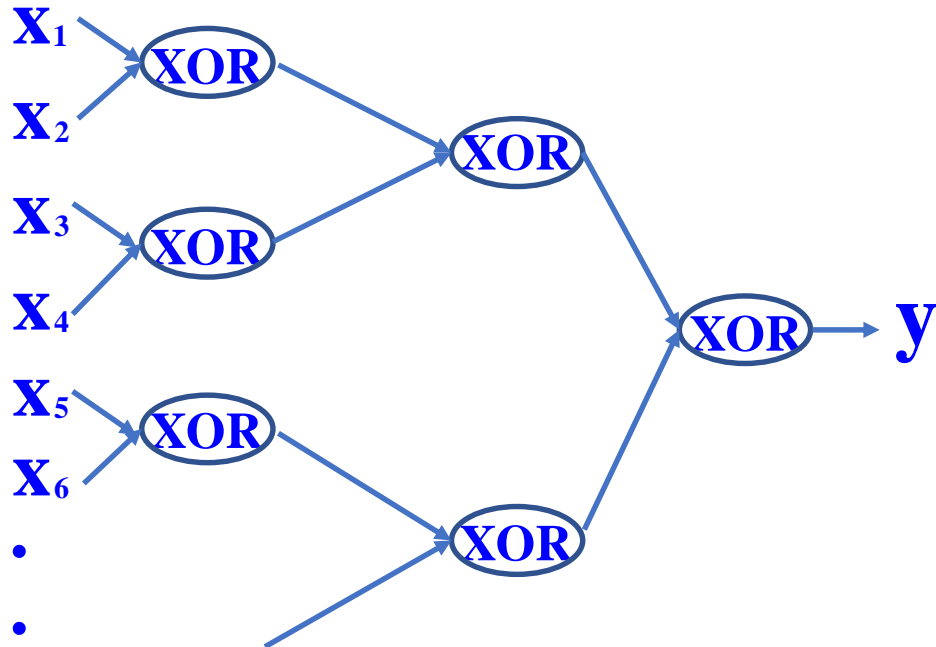


# Circuit theory and deep learning

Informally: There are functions you can compute with a "*small*" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

$$y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$$

$\longleftarrow O(\log n) \longrightarrow$



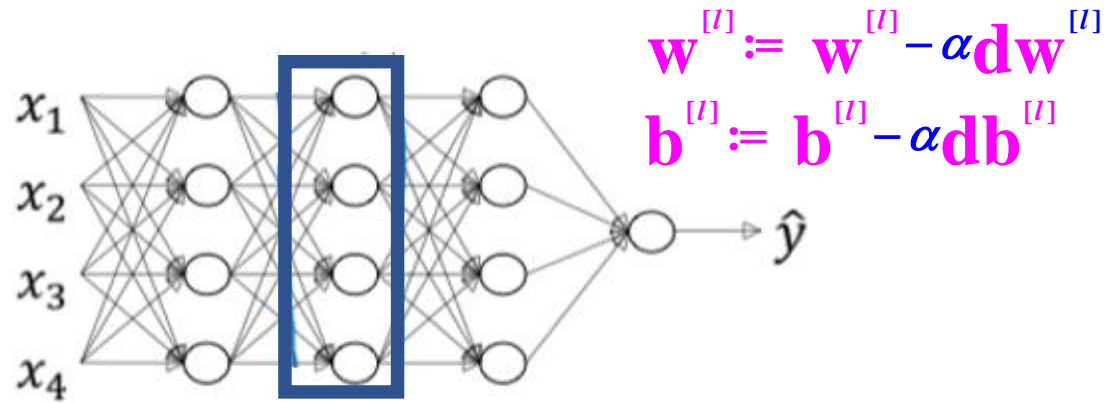


# **Deep Neural Network**

---

**Building blocks of  
deep neural networks**

# Forward and backward functions



layer  $l$ :  $\mathbf{w}^{[l]}, \mathbf{b}^{[l]}$

Forward: Input  $\mathbf{a}^{[l-1]}$ , Output  $\mathbf{a}^{[l]}$

$$\mathbf{z}^{[l]} = \mathbf{w}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]} \quad \text{cache } \mathbf{z}^{[l]} \dots$$

$$\mathbf{a}^{[l]} = \mathbf{g}^{[l]}(\mathbf{z}^{[l]})$$

Backward: Input  $\mathbf{da}^{[l]}$ , Output  $\mathbf{da}^{[l-1]}$

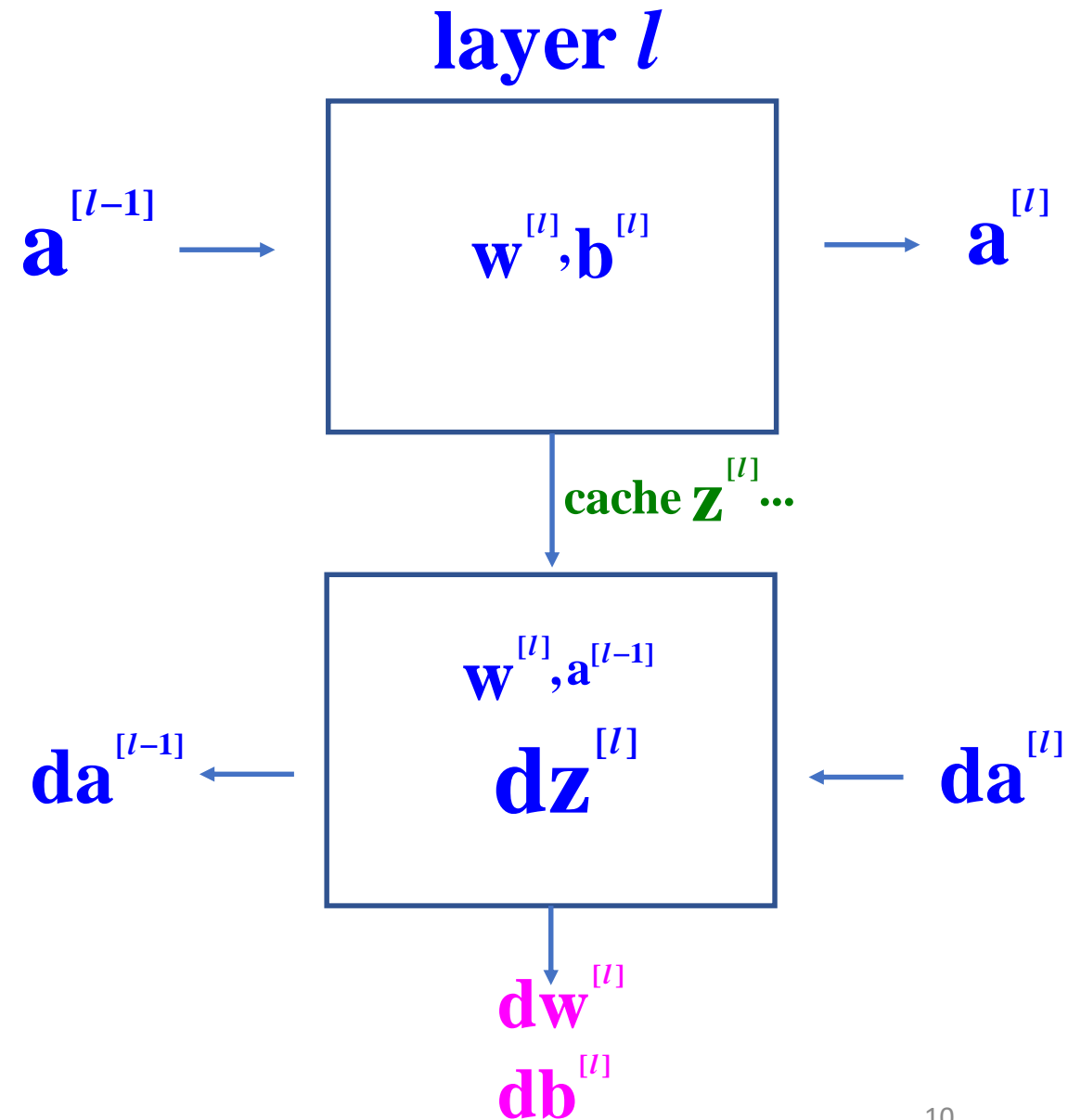
cache( $\mathbf{z}^{[l]}$ )...

$$\mathbf{dz}^{[2]} = \mathbf{dA}^{[2]} * \mathbf{g}^{[2]'}(\mathbf{z}^{(2)}) \quad \mathbf{dw}^{[1]} \quad \mathbf{db}^{[1]}$$

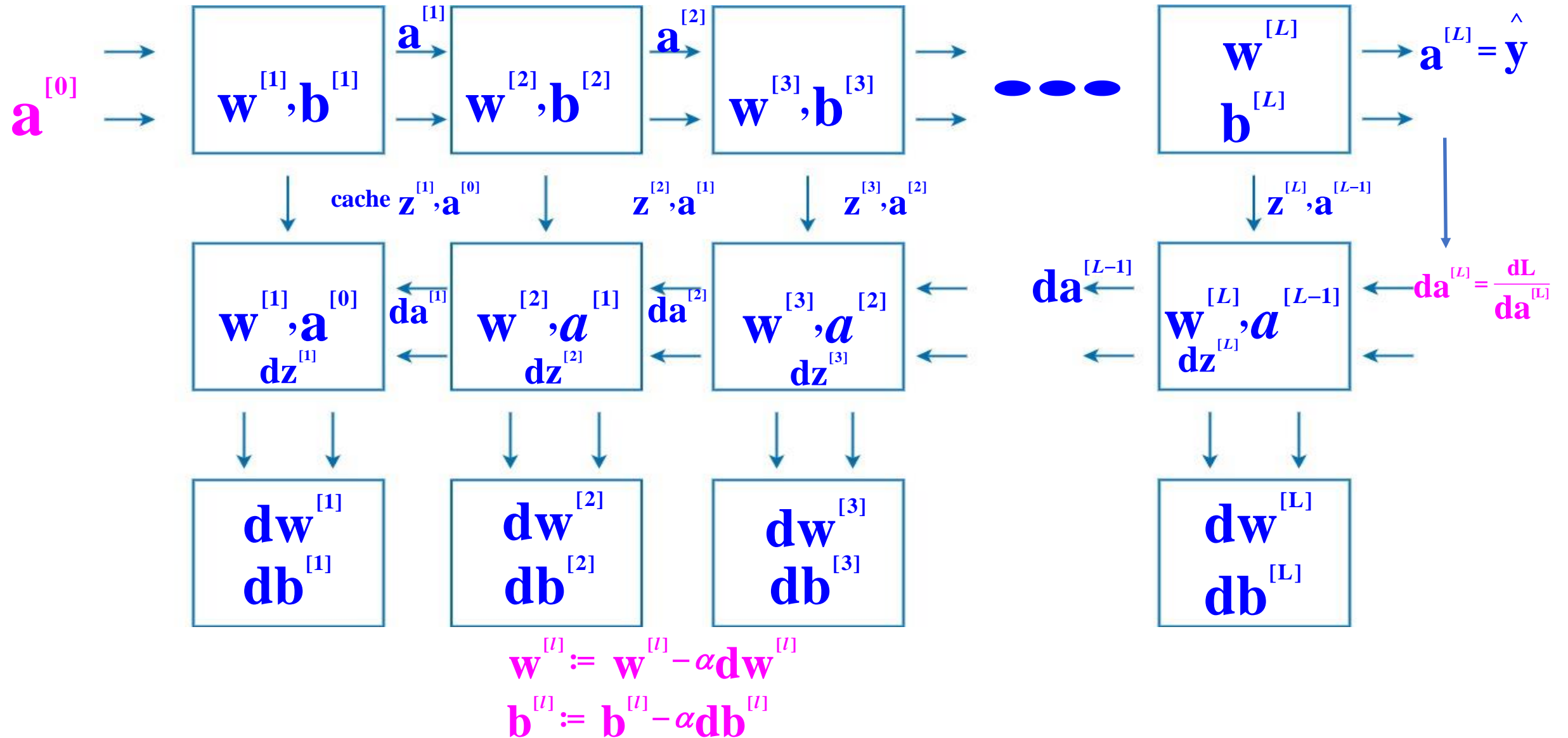
$$\mathbf{dw}^{[2]} = \frac{1}{m} \mathbf{dz}^{[2]} \mathbf{x}^T = \frac{1}{m} \mathbf{dz}^{[2]} \mathbf{A}^{[1]T}$$

$$\mathbf{db}^{[2]} = \frac{1}{m} \text{np.sum}(\mathbf{dz}^{[2]}, \dots)$$

$$\mathbf{dA}^{[1]} = \mathbf{w}^{[2]} \mathbf{dz}^{[2]}$$



# Forward and backward functions



# **Deep Neural Network**

---

**Forward and backward  
propagation**

# Backward propagation for layer $l$

Input  $\mathbf{da}^{[l]}$  **cache:**  $\mathbf{z}^{[l]}, \mathbf{a}^{[l-1]}$

Output  $\mathbf{da}^{[l-1]}, \mathbf{dw}^{[l]}, \mathbf{db}^{[l]}$

1 sample

$$\mathbf{dz}^{[l]} = \mathbf{da}^{[l]} * \mathbf{g}^{[l]'}(\mathbf{z}^{[l]})$$

$$\mathbf{dw}^{[l]} = \mathbf{a}^{[l-1]} \mathbf{dz}^{[l]T}$$

$$\mathbf{db}^{[l]} = \mathbf{dz}^{[l]}$$

$$\mathbf{da}^{[l-1]} = \mathbf{w}^{[l]} \mathbf{dz}^{[l]} \\ (\mathbf{n}^{[l-1]}, \mathbf{n}^{[l]}) (\mathbf{n}^{[l]}, 1)$$

m samples

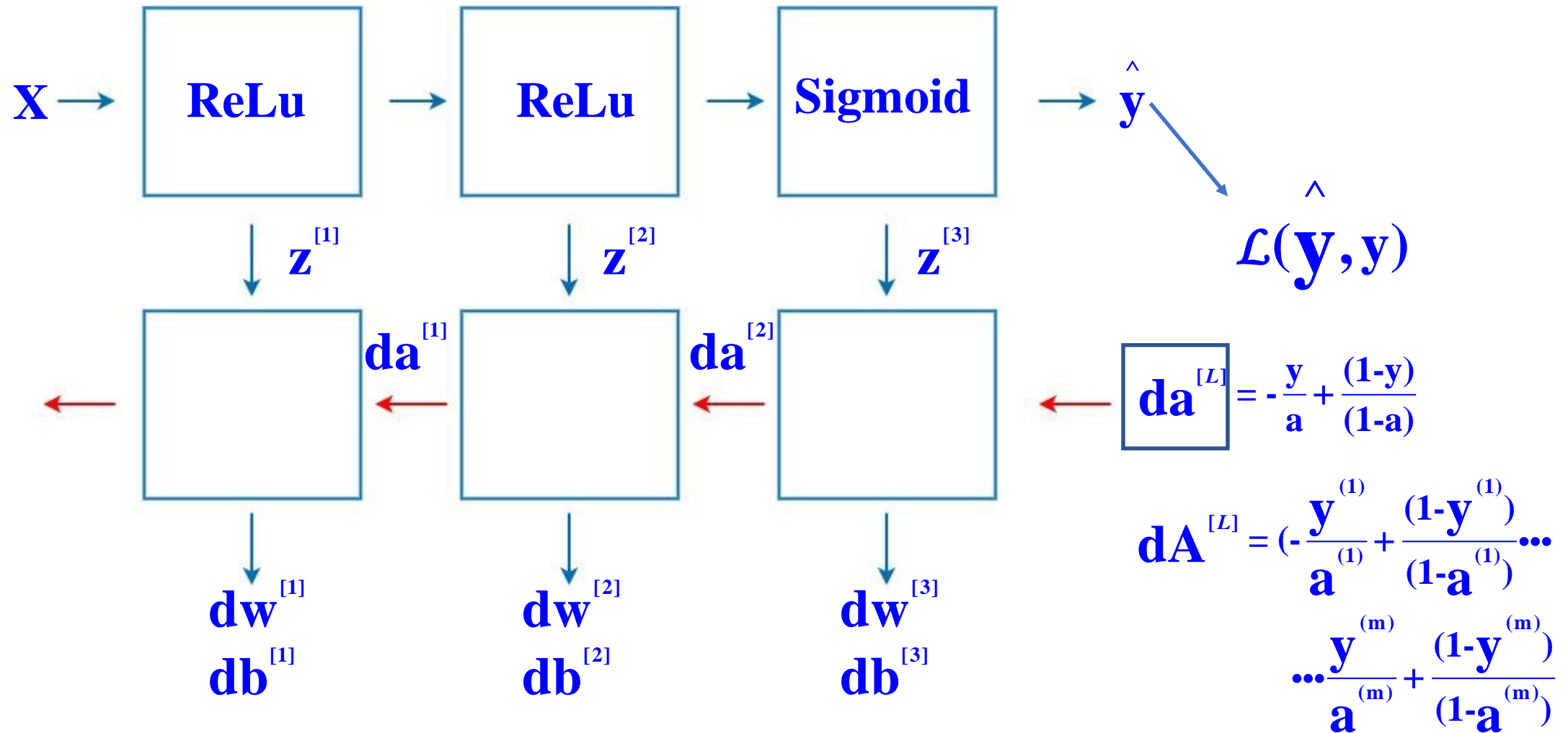
$$\mathbf{dz}^{[l]} = \mathbf{dA}^{[l]} * \mathbf{g}^{[l]'}(\mathbf{z}^{[l]})$$

$$\mathbf{dw}^{[l]} = \frac{1}{m} \mathbf{A}^{[l-1]} \mathbf{dz}^{[l]T} \\ (\mathbf{n}^{[l-1]}, \mathbf{n}^{[l]}) (\mathbf{n}^{[l-1]}, m) (m, \mathbf{n}^{[l]})$$

$$\mathbf{db}^{[l]} = \frac{1}{m} \text{np.sum}(\mathbf{dz}^{[l]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$\mathbf{dA}^{[l-1]} = \mathbf{w}^{[l]} \mathbf{dz}^{[l]} \\ (\mathbf{n}^{[l-1]}, m) (\mathbf{n}^{[l-1]}, \mathbf{n}^{[l]}) (\mathbf{n}^{[l]}, m)$$

# Summary



# **Deep Neural Network**

---

## **Parameters vs Hyperparameters**

# What are hyperparameters?

Parameters :  $\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \mathbf{W}^{[2]}, \mathbf{b}^{[2]}, \mathbf{W}^{[3]}, \mathbf{b}^{[3]}$

Hyperparameters : learning rate  $\alpha$

#iterations

#hidden layers  $L$

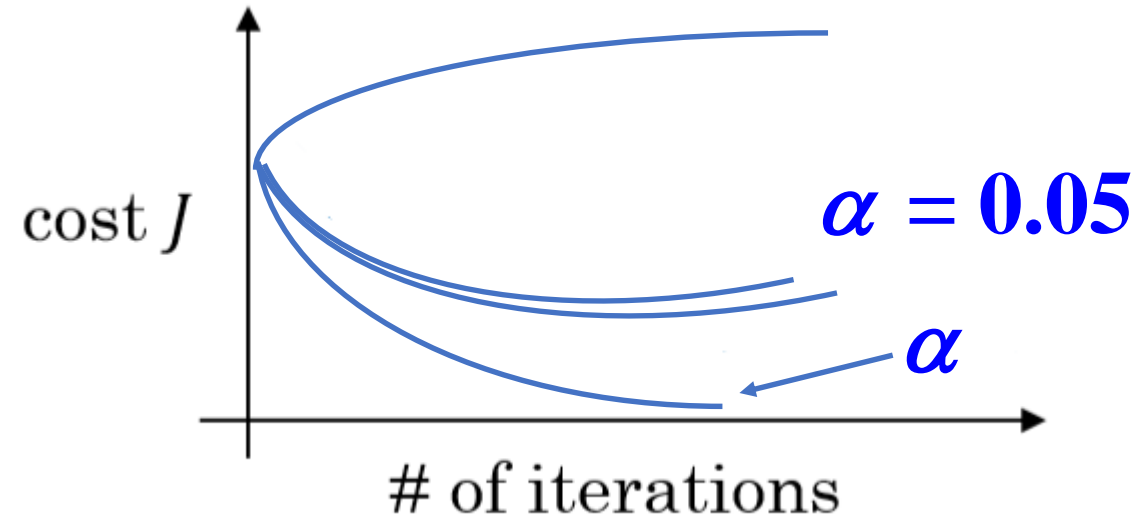
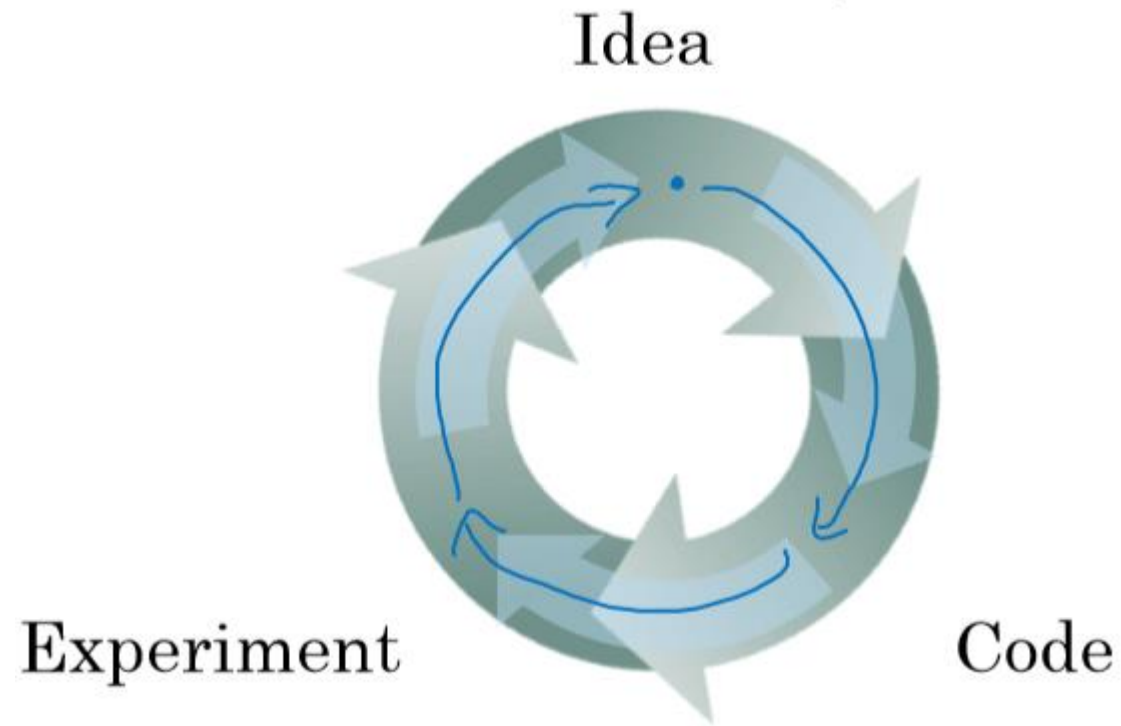
#hidden units  $\mathbf{n}^{[1]}, \mathbf{n}^{[2]}, \dots$

choice of activation function

Later: momentum, minibatch size, regularization, ...



# Applied deep learning is a very empirical process



**Vision, Speech, NLP, Ad, Search, Recommendations**

# **Deep Neural Network**

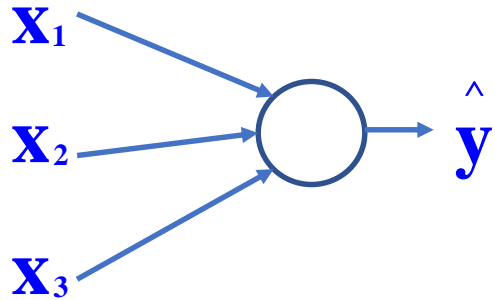
---

**What does this  
have to do with  
the brain?**

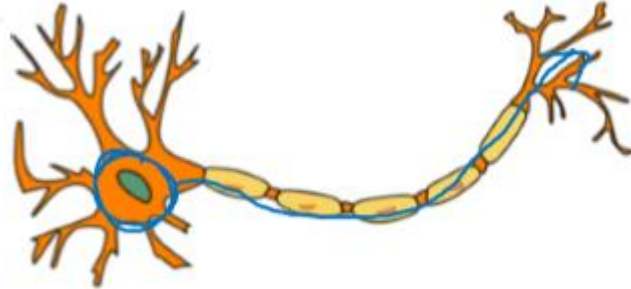
# Forward and backward propagation

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{Y} \end{aligned}$$

"It's like the brain"



$X_1$   
 $X_2$   
 $X_3$



$x \rightarrow y$

$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \\ db^{[L]} &= \frac{1}{m} np.sum(dZ^{[L]}, axis = 1, keepdims = True) \\ dZ^{[L-1]} &= dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]}) \\ &\vdots \\ dZ^{[1]} &= dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[1]T} \\ db^{[1]} &= \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True) \end{aligned}$$

$dw^{[l]} = \frac{1}{m} A^{[l-1]} dz^{[l]T}$   
看W怎樣定義

# 3 – layer neural network programming

1. Change alpha
2. Change initial condition to others, ex. All zero, can you make it faster
3. Reduce layer numbers to two, one layer and redraw the distribution plot
4. Change A0 by:

```
import sklearn; import sklearn.datasets
# load image dataset: blue/red dots in circles
train_X, train_Y = sklearn.datasets.make_circles(n_samples=35, noise=.05)
A0=train_X.T; Y=train_Y.reshape(1,35)
print("np.shape(Y):",np.shape(Y),"np.shape(A0):",np.shape(A0))
for i in range(35):
    if(Y[0][i]<0.5):
        plt.scatter(A0[0][i], A0[1][i],c="blue",marker="o")
    else:
        plt.scatter(A0[0][i], A0[1][i],c="purple",marker="x")
plt.title('x1, x2, y')
plt.show()
```

# 3-layer neural network programming

a standard Normal distribution (mean=0, stdev=1).

```
1 import matplotlib
2 import time
3 import numpy as np
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 n0=2;n1=3;n2=4;n3=1;m=35
7 alpha=0.01
```

```
1 # Prepare input A0 and labelling output Y
2 x1=[0.4,1.2,1.1,1.9,1.8,1.5,2.6,2.2,3.3,3.0,4.1,4.2,5.6,0.7,0.7
3 x2=[4.1,1.5,4.8,0.7,2.2,3.8,1.8,4.8,1.0,4.3,1.8,3.1,6.0,5.9,7.0
4 y0=[0,0,0,0,0,0,0,0,0,0,0,0,0,0,0];y1=[1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,
5 A0=np.array([x1,x2])
6 Y= np.concatenate((y0,y1),axis=None).reshape(1,35)
7 print("np.shape(Y):",np.shape(Y),"np.shape(A0):",np.shape(A0),
8       "A0[0][0]:",A0[0][0],"A0[0][1]:",A0[0][1],"Y",Y[0][0:3])
9 plt.scatter(A0[0][0:13], A0[1][0:13],marker="o")
10 plt.scatter(A0[0][13:35], A0[1][13:35],marker="x")
11 plt.title('x1, x2, y')
12 plt.show()
```

```
1 # Initialization
2 np.random.seed(1)
3 WT1=np.random.randn(n1,n0)*0.01;b1=np.random.randn(n1,1)
4 WT2=np.random.randn(n2,n1)*0.01;b2=np.random.randn(n2,1)
5 WT3=np.random.randn(n3,n2)*0.01;b3=np.random.randn(n3,1)
```

# 3 – layer neural network programming

```
1 # Initialization
2 np.random.seed(1)
3 WT1=np.random.randn(n1,n0)*0.01;b1=np.random.randn(n1,1)
4 WT2=np.random.randn(n2,n1)*0.01;b2=np.random.randn(n2,1)
5 WT3=np.random.randn(n3,n2)*0.01;b3=np.random.randn(n3,1)
```

```
1 def forwardfunc():
2     global A0,A1,A2,A3,WT1,WT2,WT3,b1,b2,b3
3     Z1=np.dot(WT1,A0)+b1;
4     A1=(np.exp(Z1)-np.exp(-Z1))/(np.exp(Z1)+np.exp(-Z1))
5     Z2=np.dot(WT2,A1)+b2;
6     A2=(np.exp(Z2)-np.exp(-Z2))/(np.exp(Z2)+np.exp(-Z2))
7     Z3=np.dot(WT3,A2)+b3;
8     A3=1/(1+np.exp(-Z3))
```

# 3 – layer neural network programming

```
1  def backwardprop():
2      global Y,A3,A2,A1,A0,dA2,dA1,dZ3,dZ2,dZ1,
3          dWT3,dWT2,dWT1,db3,db2,db1
4      #   dZ3=A3-Y   #dZ3 comes from dA3 and dA3/dZ3
5      dA3=-Y/A3+(1-Y)/(1-A3)
6      dZ3=dA3*(A3*(1-A3))
7      dWT3=1/m*np.dot(dZ3,A2.T)
8      db3=1/m*np.sum(dZ3,axis=1,keepdims=True)
9      dA2=np.dot(WT3.T,dZ3)
10     dZ2=dA2*(1-A2**2)
11     dWT2=1/m*np.dot(dZ2,A1.T)
12     db2=1/m*np.sum(dZ2,axis=1,keepdims=True)
13     dA1=np.dot(WT2.T,dZ2)
14     dZ1=dA1*(1-A1**2)
15     dWT1=1/m*np.dot(dZ1,A0.T)
16     db1=1/m*np.sum(dZ1,axis=1,keepdims=True)
```

# 3-layer neural network programming

```
1 itera=0
2 cost = 100
3 while (cost > 0.05): #0.05, 0.2
4     forwardfunc()
5     cost=-1/m*np.sum((Y*np.log(A3)+(1-Y)*np.log(1-A3)),axis=1,keepdims=True)
6     if(itera%10000==0):
7         print("itera ", itera, "cost ", cost)
8     backwardprop()
9     WT3=WT3-alpha*dWT3; b3=b3-alpha*db3; WT2=WT2-alpha*dWT2; b2=b2-alpha*db2
10    WT1=WT1-alpha*dWT1; b1=b1-alpha*db1
11    itera=itera+1
12 print("np.shape(Y)",np.shape(Y), "np.shape(A3)",np.shape(A3))
```

```
itera 0 cost [[0.65971979]]
itera 10000 cost [[0.65971121]]
itera 20000 cost [[0.65970933]]
itera 30000 cost [[0.65967028]]
```



# 3-layer neural network programming

```
1  # Inference
2  A0 = np.zeros((n0,1600))#[xx1, xx2][xy1,xy2]
3  for i in range(40):
4      for j in range(40):
5          A0[0][i*39+j]=i*0.2
6          A0[1][i*39+j]=j*0.2
7  forwardfunc()
8  print("A3 ",A3)
9  plt.scatter(A0[0][0:2], A0[1][0:2],c="red",marker=".")
10 plt.scatter(A0[0][2], A0[1][2],c="green",marker=".")
11 for i in range(40):
12     for j in range(40):
13         if(A3[0][i*39+j]>0.5):
14             plt.scatter(A0[0][i*39+j],
15                         A0[1][i*39+j],c="red",marker=".")
16         else:
17             plt.scatter(A0[0][i*39+j],
18                         A0[1][i*39+j],c="green",marker=".")
```

# 3-layer neural network programming

```
19 A0=np.array([x1,x2])
20 forwardfunc()
21 for i in range(m):
22     if(A3[0][i]<0.5):
23         plt.scatter(A0[0][i], A0[1][i],c="blue",marker="o")
24     else:
25         plt.scatter(A0[0][i], A0[1][i],c="purple",marker="x")
26 plt.title('xx, xy, y')
27 plt.show()
```

