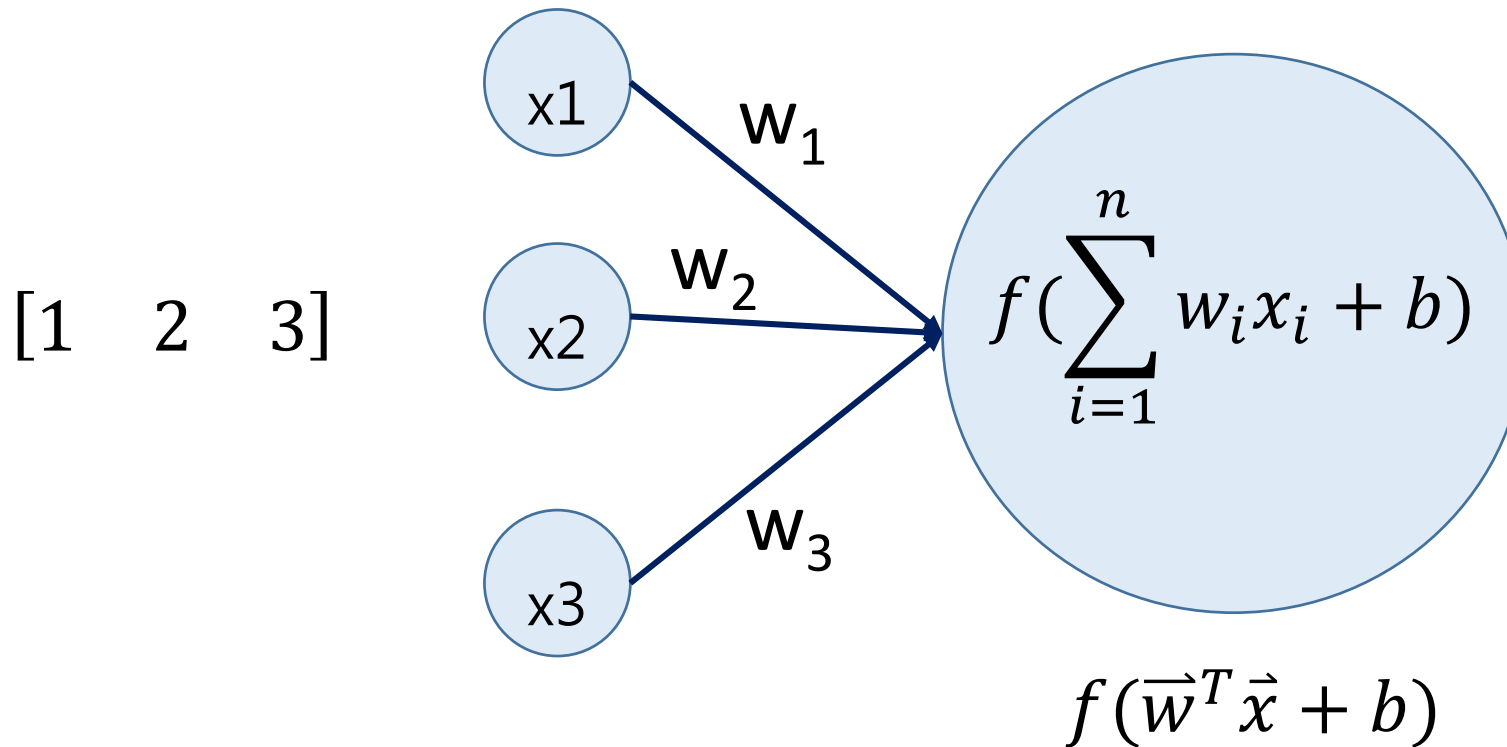# Introduction to artificial neural network (deep NN)

MIT Introduction to Deep Learning | 6.S191   https://youtu.be/5tvmMX8r_OM

# Neuron (perceptron)

- Neuron performs weighted linear combination with bias and activation function

1.1. Perceptron.ipynb

$$[1 \quad 2 \quad 3]$$

x1

$w_1$

x2

$w_2$

x3

$w_3$

$$f(\sum_{i=1}^{n} w_i x_i + b)$$

$$f(\vec{w}^T \vec{x} + b)$$

# Neuron (perceptron)

```
lstX   =   [[1, 2, 3],   [4, 5, 6], [7, 8, 9]]
tensorX   =   torch.FloatTensor(lstX)
print(tensorX,   "\n",   tensorX.shape)
```
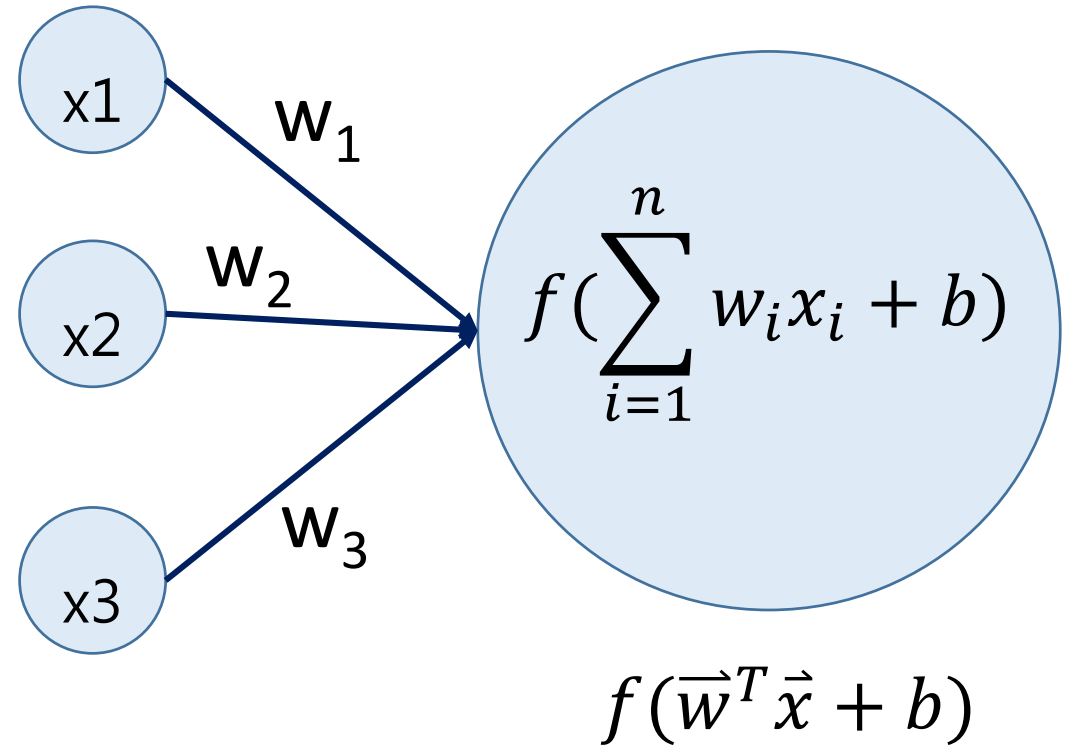
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

```
tensorX.mm(W0t)
```

```
tensor([[0.6568],
        [1.8483],
        [3.0397]],
```

$$\begin{bmatrix} x_1^1 & x_2^1 & x_3^1 \\ x_1^2 & x_2^2 & x_3^2 \\ x_1^3 & x_2^3 & x_3^3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$
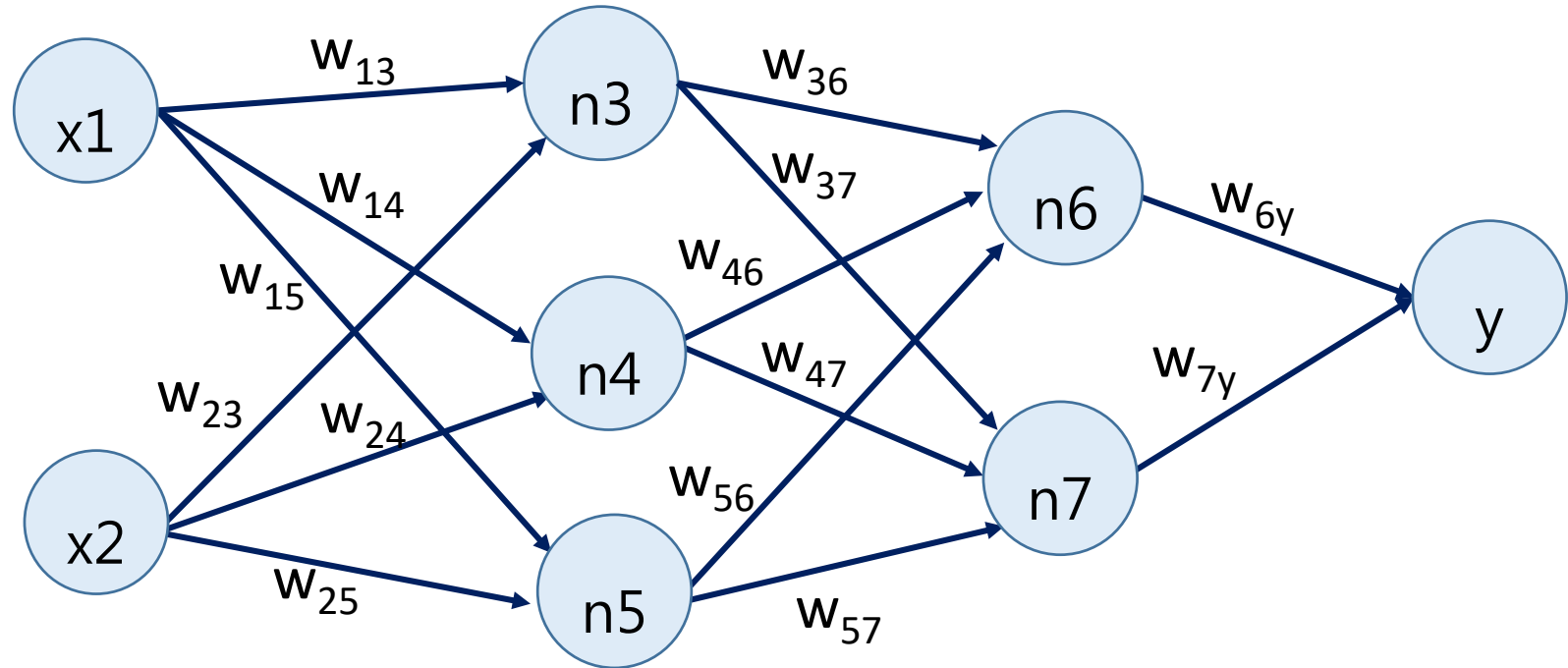
x1 — $w_1$

x2 — $w_2$

x3 — $w_3$

$$f(\sum_{i=1}^{n} w_i x_i + b)$$

$$f(\vec{w}^T \vec{x} + b)$$

# Multiple-layer perception (MLP)

1.2 MLP forward propagation.ipynb
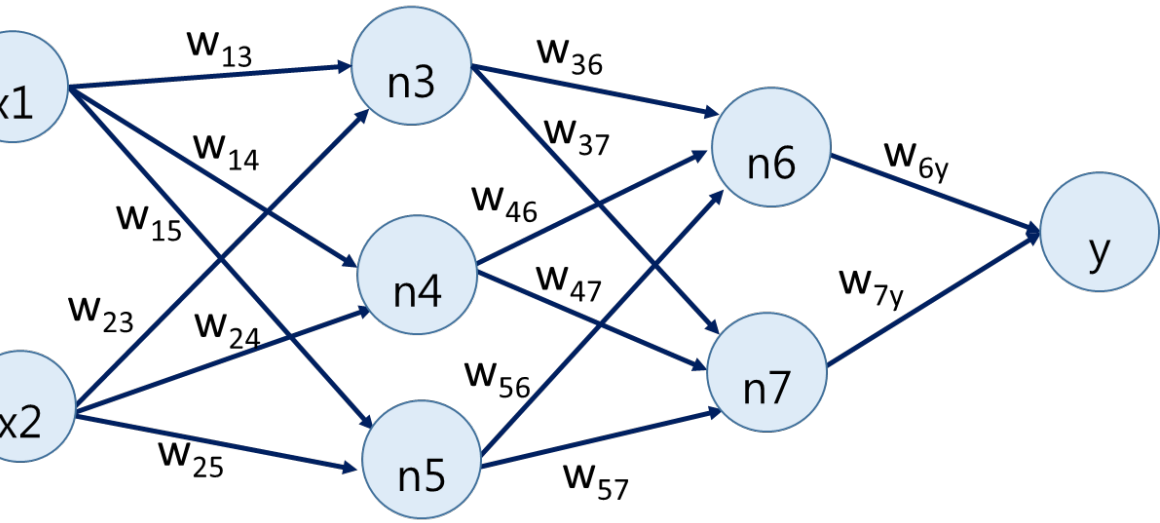
```
MyNet  =  nn.Sequential(
    nn.Linear(2,   3),
    nn.Linear(3,   2),
    nn.Linear(2,   1)
)
print(MyNet)
```

# Weights and bias



```python
for param in MyNet.parameters():
    if param.requires_grad:
        print(param.data)
```

$$\text{tensor}\left(\begin{bmatrix} 0.4727 & -0.5188 \\ -0.5681 & -0.6032 \\ -0.0252 & -0.3011 \end{bmatrix}\right) \quad \begin{bmatrix} w_{13} & w_{23} \\ w_{14} & w_{24} \\ w_{15} & w_{25} \end{bmatrix}$$

$$\text{tensor}\left(\begin{bmatrix} -0.6986 & -0.6602 & -0.4860 \end{bmatrix}\right) \quad \begin{bmatrix} b_3 & b_4 & b_5 \end{bmatrix}$$

$$\text{tensor}\left(\begin{bmatrix} -0.5549 & 0.2550 & 0.4584 \\ 0.2930 & 0.0849 & -0.3146 \end{bmatrix}\right) \quad \begin{bmatrix} w_{36} & w_{46} & w_{56} \\ w_{37} & w_{47} & w_{57} \end{bmatrix}$$

$$\text{tensor}\left(\begin{bmatrix} 0.1677 & 0.0736 \end{bmatrix}\right) \quad \begin{bmatrix} b_6 & b_7 \end{bmatrix}$$

$$\text{tensor}\left(\begin{bmatrix} 0.4106 & -0.3618 \end{bmatrix}\right) \quad \begin{bmatrix} w_{6y} & w_{7y} \end{bmatrix}$$

$$\text{tensor}\left(\begin{bmatrix} -0.2270 \end{bmatrix}\right) \quad \begin{bmatrix} b_y \end{bmatrix}$$

# Prepare input

```
1stX  =  [[1,   2],   [2,   3],   [10,   5]]
tensorX  =  torch.FloatTensor(1stX)
print(tensorX,  "\n",   tensorX.shape)
```

```
tensor([[ 1.,   2.],
        [ 2.,   3.],
        [10.,   5.]])
 torch.Size([3,  2])
```



$$\vec{X} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 10 & 5 \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 \\ x_1^2 & x_2^2 \\ x_1^3 & x_2^3 \end{bmatrix}$$
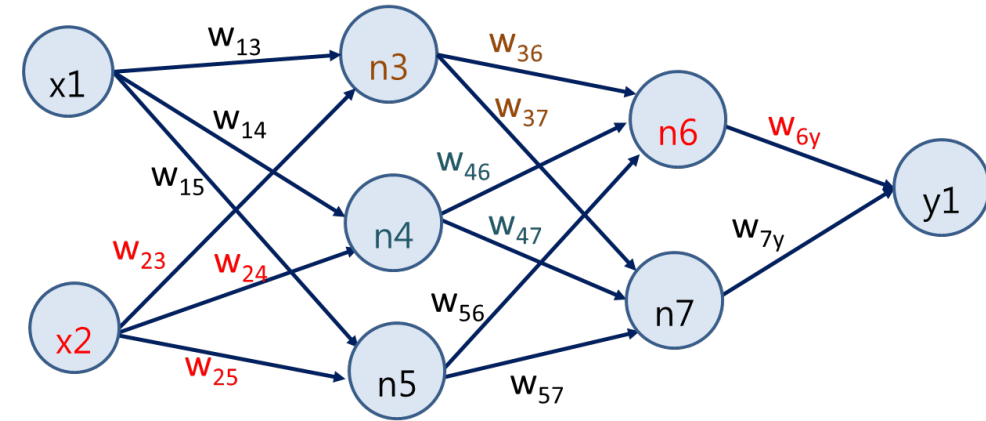
# Forward propagation (layer 1)

```
#Calculate n3, n4, n5 using Pytorch matrix operation
Layer1_1 = tensorX.mm(torch.transpose(W0, 1, 0)) + b0
print(Layer1_1)
```
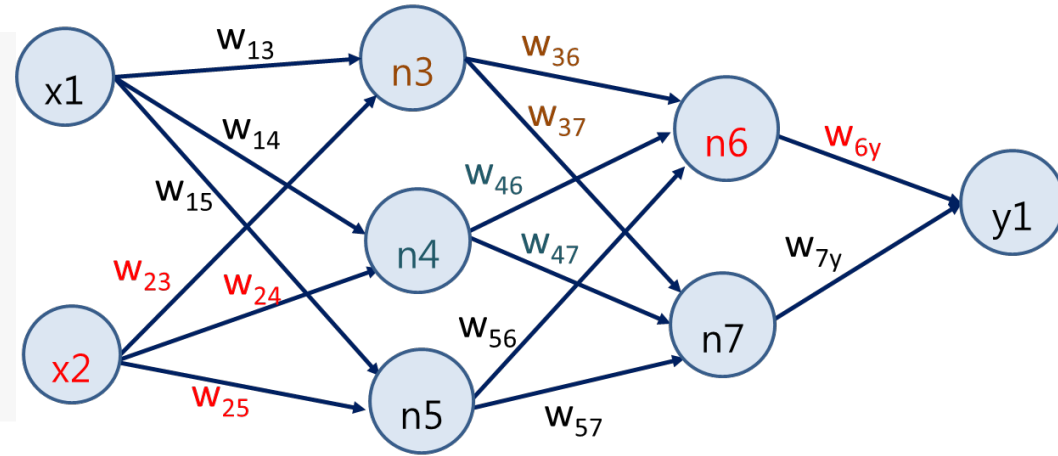


$$\begin{bmatrix} w_{13} & w_{23} \\ w_{14} & w_{24} \\ w_{15} & w_{25} \end{bmatrix}^T = \begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix}$$

$$\begin{bmatrix} x_1^1 & x_2^1 \\ x_1^2 & x_2^2 \\ x_1^3 & x_2^3 \end{bmatrix} \begin{bmatrix} w_{13} & w_{14} & w_{15} \\ w_{23} & w_{24} & w_{25} \end{bmatrix} + \begin{bmatrix} b_3 & b_4 & b_5 \end{bmatrix} = \begin{bmatrix} k_3^1 & k_4^1 & k_5^1 \\ k_3^2 & k_4^2 & k_5^2 \\ k_3^3 & k_4^3 & k_5^3 \end{bmatrix} + \begin{bmatrix} b_3 & b_4 & b_5 \\ b_3 & b_4 & b_5 \\ b_3 & b_4 & b_5 \end{bmatrix} = \begin{bmatrix} n_3^1 & n_4^1 & n_5^1 \\ n_3^2 & n_4^2 & n_5^2 \\ n_3^3 & n_4^3 & n_5^3 \end{bmatrix}$$

# Forward propagation (layer 2)

```
#Calculate  n6,  n7  using  PyTorch  matrix  operation
W1  =  MyNet[1].weight
b1  =  MyNet[1].bias
Layer2_1  =  Layer1.mm(torch.transpose(W1,  1,  0))  +b1
print(Layer2_1)
```



$$\begin{bmatrix} w_{36} & w_{46} & w_{56} \\ w_{37} & w_{47} & w_{57} \end{bmatrix}^T = \begin{bmatrix} w_{36} & w_{37} \\ w_{46} & w_{47} \\ w_{56} & w_{57} \end{bmatrix}$$

$$\begin{bmatrix} n_3^1 & n_4^1 & n_5^1 \\ n_3^2 & n_4^2 & n_5^2 \\ n_3^3 & n_4^3 & n_5^3 \end{bmatrix} \begin{bmatrix} w_{36} & w_{37} \\ w_{46} & w_{47} \\ w_{56} & w_{57} \end{bmatrix} + \begin{bmatrix} b_6 & b_7 \end{bmatrix} = \begin{bmatrix} k_6^1 & k_7^1 \\ k_6^2 & k_7^2 \\ k_6^3 & k_7^3 \end{bmatrix} + \begin{bmatrix} b_6 & b_7 \\ b_6 & b_7 \\ b_6 & b_7 \end{bmatrix} = \begin{bmatrix} n_6^1 & n_7^1 \\ n_6^2 & n_7^2 \\ n_6^3 & n_7^3 \end{bmatrix}$$
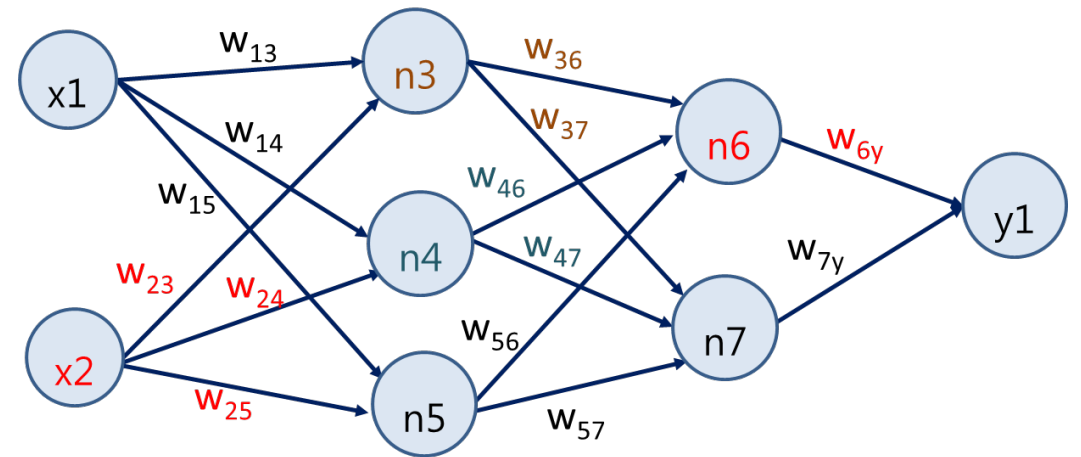
# Forward propagation (output)

```
#Calculate  y  by  matrix  operation
W2  =  MyNet[2].weight
b2  =  MyNet[2].bias
tensorY_1  =  Layer2.mm(torch.transpose(W2,  1,  0))  +b2
print(tensorY_1)
```

$$[w_{6y} \quad w_{7y}]^T = \begin{bmatrix} w_{6y} \\ w_{7y} \end{bmatrix}$$

$$\begin{bmatrix} n_6^1 & n_7^1 \\ n_6^2 & n_7^2 \\ n_6^3 & n_7^3 \end{bmatrix} \begin{bmatrix} w_{6y} \\ w_{7y} \end{bmatrix} + [b_y] = \begin{bmatrix} k_y^1 \\ k_y^2 \\ k_y^3 \end{bmatrix} + \begin{bmatrix} b_y \\ b_y \\ b_y \end{bmatrix} = \begin{bmatrix} y^1 \\ y^2 \\ y^3 \end{bmatrix}$$
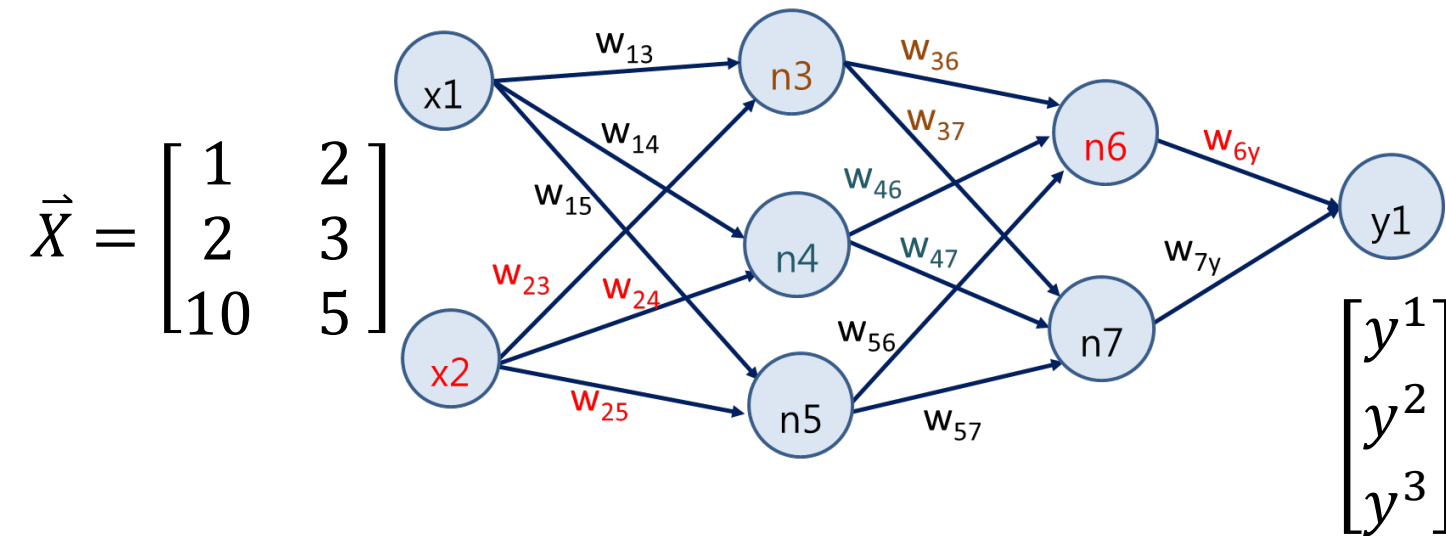
# Calculate prediction error

```
lstX  =  [[1,   2],   [2,   3],   [10,   5]]
lstY  =  [[7],   [12],   [40]]      #   y=3x1+2x2
tensorX  =  torch.FloatTensor(lstX)
tensorY_hat  =  torch.FloatTensor(lstY)
```

$$\begin{bmatrix} \hat{y}^1 \\ \hat{y}^2 \\ \hat{y}^3 \end{bmatrix} = \begin{bmatrix} 7 \\ 12 \\ 40 \end{bmatrix}$$

$$\vec{X} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 10 & 5 \end{bmatrix}$$



$$\begin{bmatrix} y^1 \\ y^2 \\ y^3 \end{bmatrix}$$

```
loss  =  loss_func(tensorY  ,   tensorY_hat)
print(loss)
```

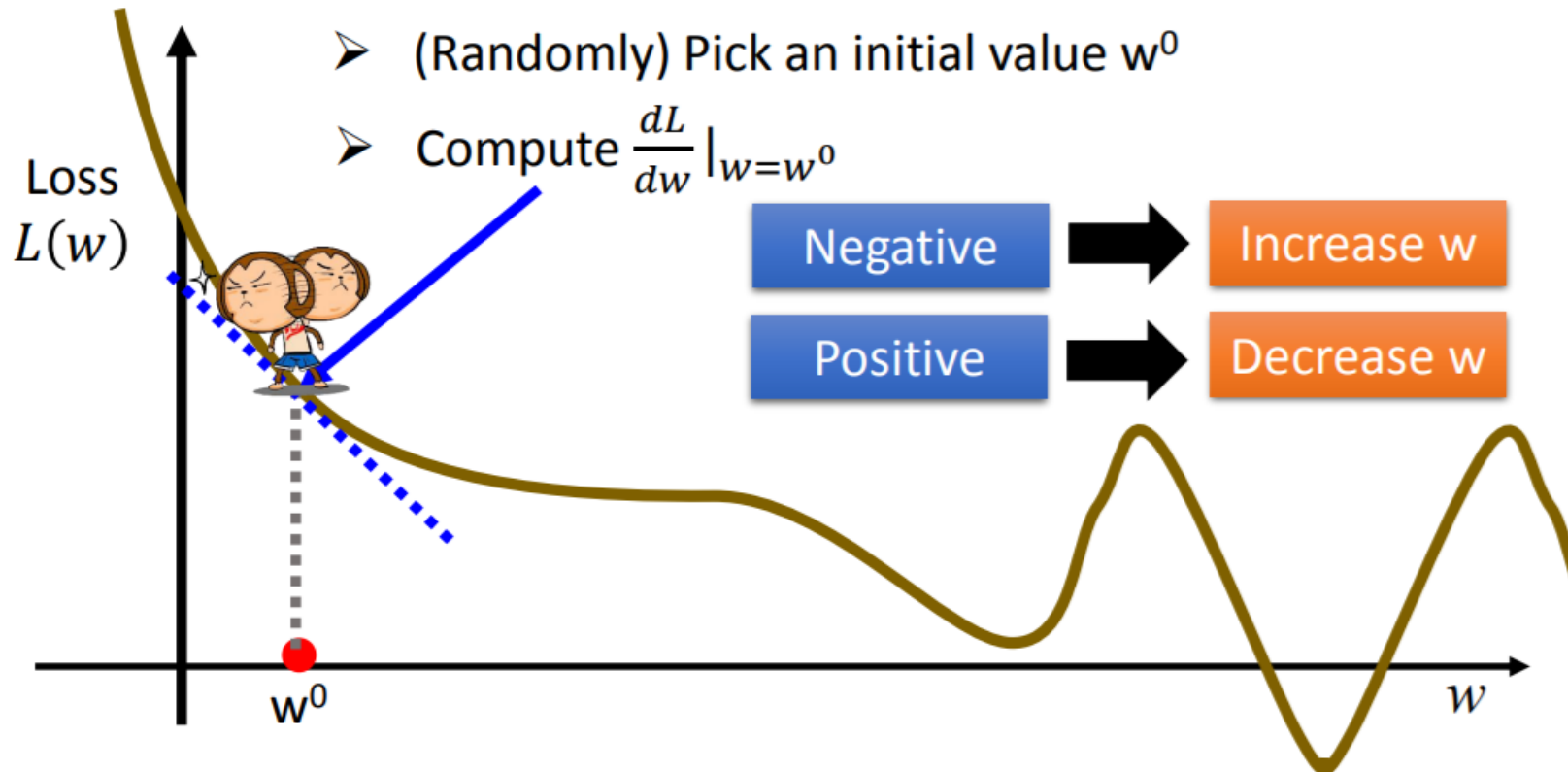$$L = \frac{1}{N} \sum_{i=1}^{N} (y^i - \hat{y}^i)^2$$

```
tensorY=  MyNet(tensorX)
print(tensorY)
```

# Use gradient decent to find optimal parameters

3. Find the optimal parameters that minimize $\mathcal{L}(f)$
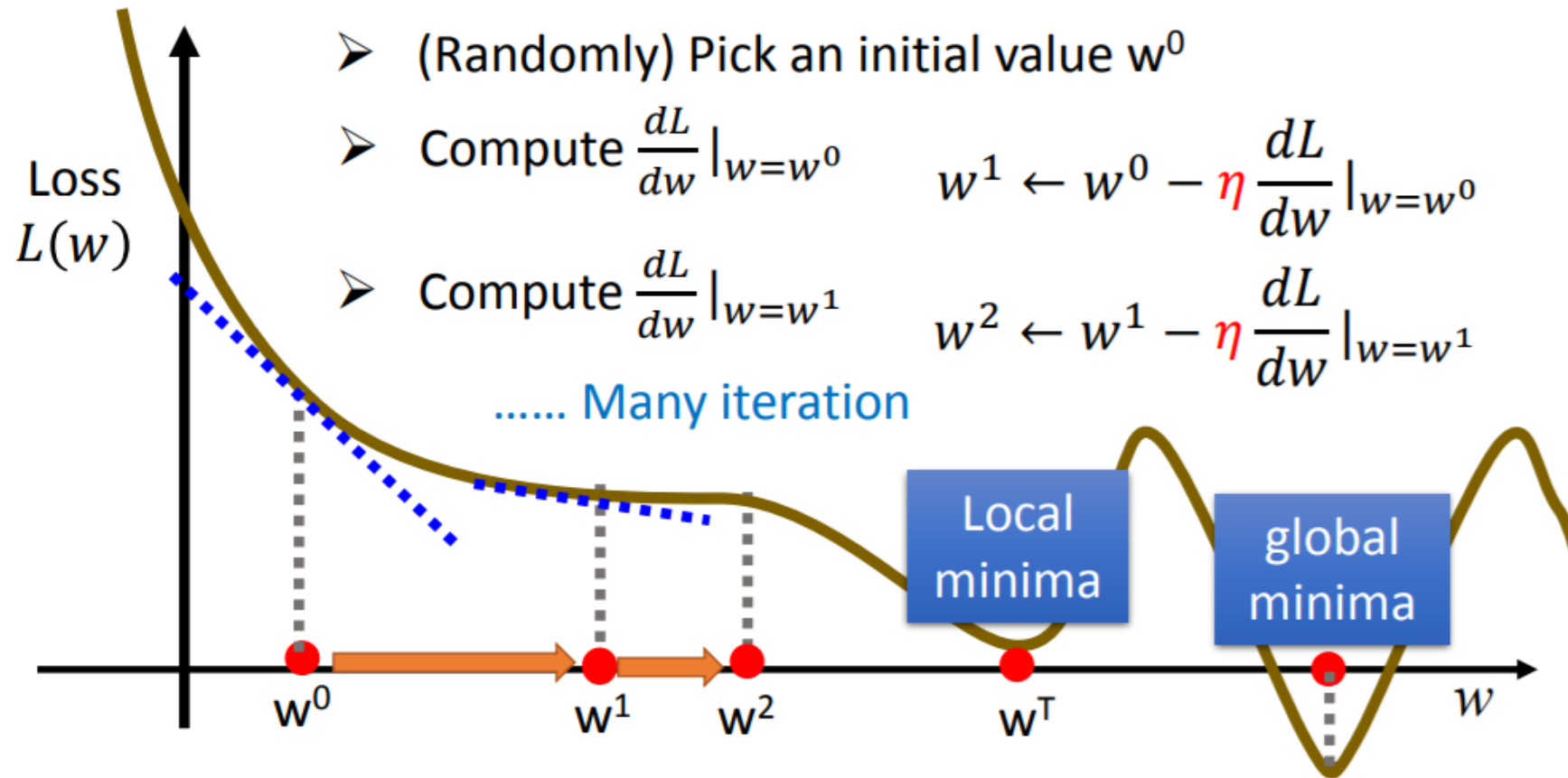
$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:

  ➢ (Randomly) Pick an initial value $w^0$

  ➢ Compute $\frac{dL}{dw}|_{w=w^0}$

| Negative | ➡ | Increase w |
| Positive | ➡ | Decrease w |

Loss $L(w)$

$w^0$

$w$

# Use gradient decent to find optimal parameters

$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:

  ➤ (Randomly) Pick an initial value $w^0$

  ➤ Compute $\frac{dL}{dw}\big|_{w=w^0}$     $w^1 \leftarrow w^0 - \eta \frac{dL}{dw}\big|_{w=w^0}$

  ➤ Compute $\frac{dL}{dw}\big|_{w=w^1}$     $w^2 \leftarrow w^1 - \eta \frac{dL}{dw}\big|_{w=w^1}$

  ...... Many iteration

Loss $L(w)$

Local minima

global minima

$w^0$     $w^1$   $w^2$     $w^T$     $w$

# Gradient decent to find two parameters $w^*$ and $b^*$

- How about two parameters?   $w^*, b^* = arg \min_{w,b} L(w, b)$

$$\begin{bmatrix} \dfrac{\partial L}{\partial w} \\ \dfrac{\partial L}{\partial b} \end{bmatrix} \text{gradient}$$

  ➢ (Randomly) Pick an initial value $w^0, b^0$

  ➢ Compute $\dfrac{\partial L}{\partial w}\big|_{w=w^0, b=b^0}, \dfrac{\partial L}{\partial b}\big|_{w=w^0, b=b^0}$
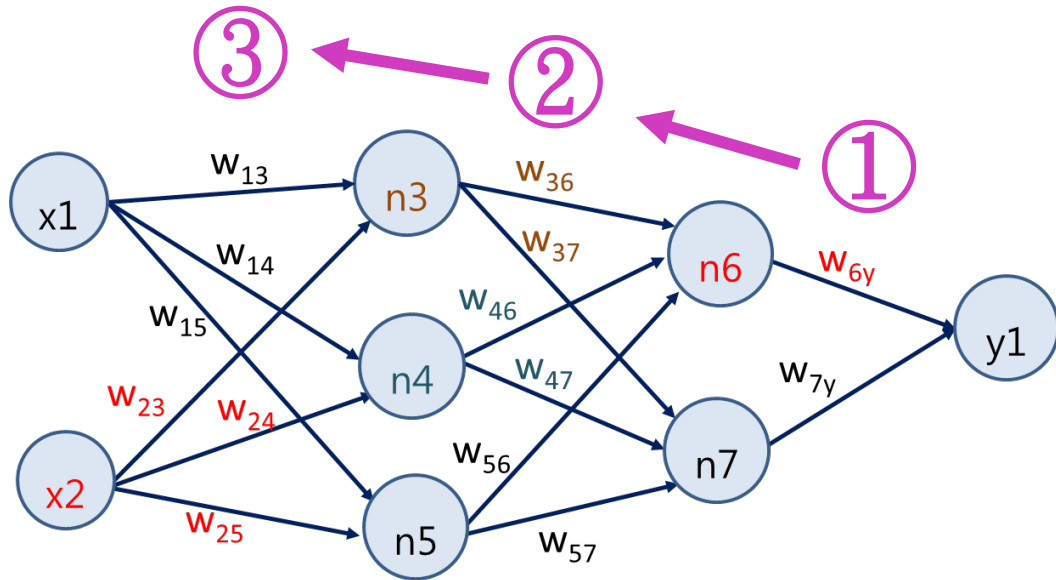
$$w^1 \leftarrow w^0 - \eta \dfrac{\partial L}{\partial w}\big|_{w=w^0, b=b^0} \qquad b^1 \leftarrow b^0 - \eta \dfrac{\partial L}{\partial b}\big|_{w=w^0, b=b^0}$$

  ➢ Compute $\dfrac{\partial L}{\partial w}\big|_{w=w^1, b=b^1}, \dfrac{\partial L}{\partial b}\big|_{w=w^1, b=b^1}$

$$w^2 \leftarrow w^1 - \eta \dfrac{\partial L}{\partial w}\big|_{w=w^1, b=b^1} \qquad b^2 \leftarrow b^1 - \eta \dfrac{\partial L}{\partial b}\big|_{w=w^1, b=b^1}$$

# Use gradient decent to find optimal NN weights



$$L = g(y - \hat{y}) \qquad y = \sigma(n_6 * w_{6y} + n_7 * w_{7y} + b_y)$$

① $w_{6y} \leftarrow w_{6y} - \eta \dfrac{\partial L}{\partial w_{6y}} \qquad \dfrac{\partial L}{\partial w_{6y}} = \dfrac{\partial L}{\partial y} \dfrac{\partial y}{\partial w_{6y}}$

$w_{7y} \leftarrow w_{7y} - \eta \dfrac{\partial L}{\partial w_{7y}} \qquad \dfrac{\partial L}{\partial w_{7y}} = \dfrac{\partial L}{\partial y} \dfrac{\partial y}{\partial w_{7y}}$

$$w_i \leftarrow w_i - \eta \dfrac{\partial e}{\partial w_i}$$

② $w_{57} \leftarrow w_{57} - \eta \dfrac{\partial L}{\partial w_{57}} \qquad \dfrac{\partial L}{\partial w_{57}} = \dfrac{\partial L}{\partial y} \dfrac{\partial y}{\partial n_7} \dfrac{\partial n_7}{\partial w_{57}}$
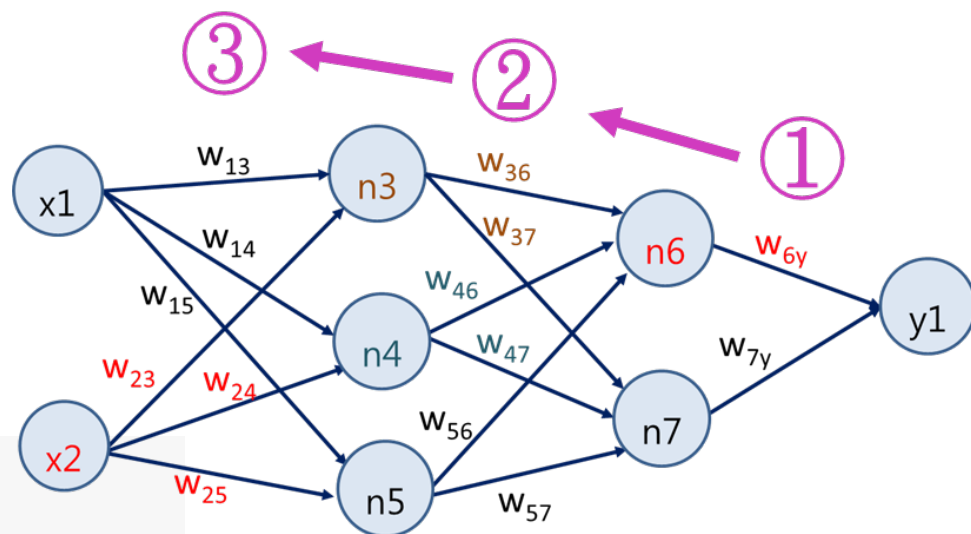
$$n_7 = f(n_3 * w_{37} + n_4 * w_{47} + n_5 * w_{57} + b_7)$$

14

# Back propagation
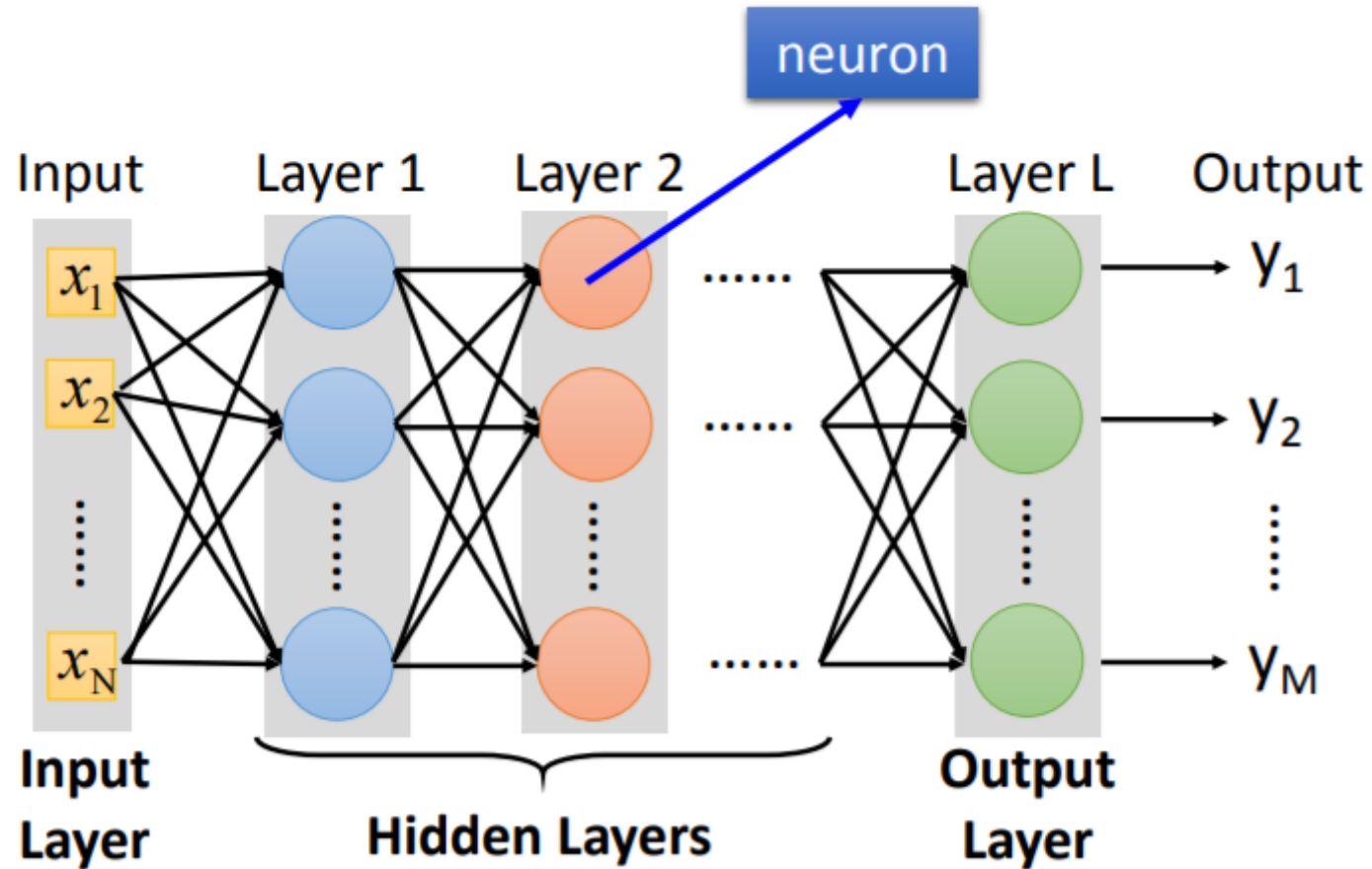
```
loss.backward()
```

```
for name, param in MyNet.named_parameters():
    if param.requires_grad:
        print(name, param.data, param.grad)
```

```
for name, param in MyNet.named_parameters():
    if param.requires_grad:
        param = param - learning_rate*param.grad

        nameLst = name.split(".")    #"0.weight" -> ['0', 'weight']
        layerNo = int(nameLst[0])
        s = nameLst[1]
        if(s=="weight"):
            MyNet[layerNo].weight = torch.nn.parameter.Parameter(param)
        elif(s=="bias"):
            MyNet[layerNo].bias = torch.nn.parameter.Parameter(param)
        else:
            print("wrong label")
```
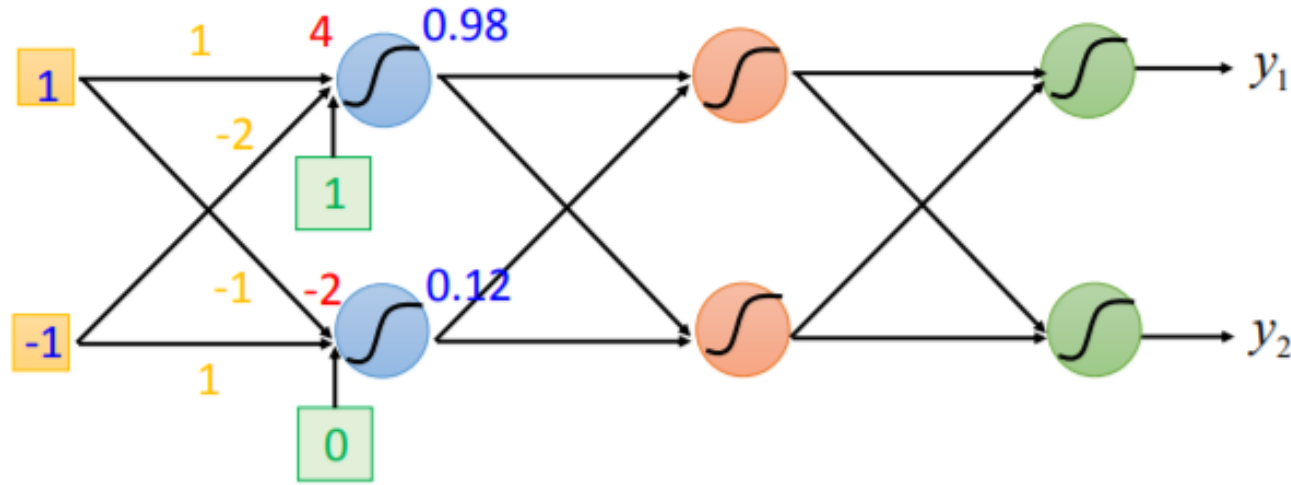


$$w_i \leftarrow w_i - \eta \frac{\partial e}{\partial w_i}$$
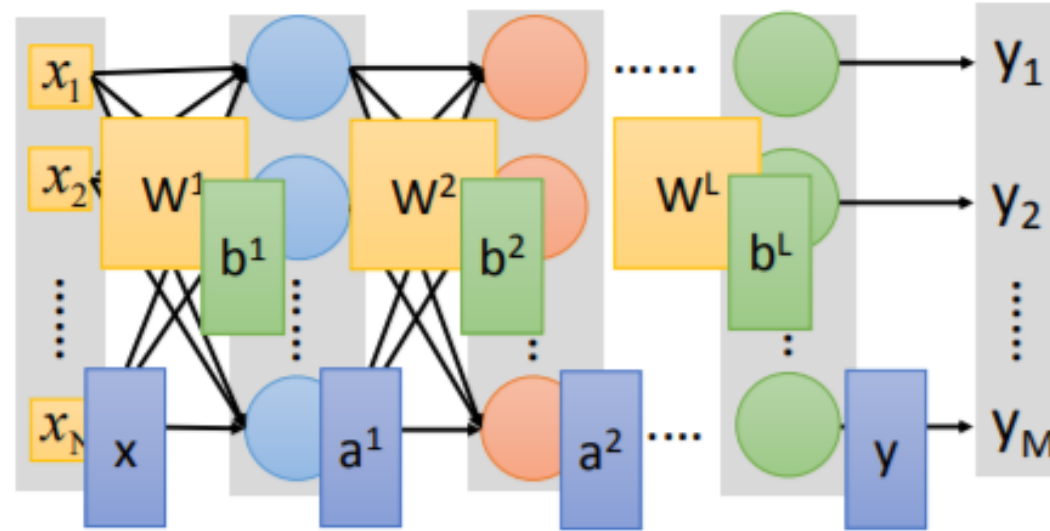
# MLP is a fully connected feedforward network

# Fully connected feed forward network is implemented as matrix operation



$$y = \sigma(w \cdot x + b)$$

$$\sigma\left( \begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Use parallel computing to speed up matrix operation



$$\boxed{y} = f(\boxed{x})$$

Using parallel computing techniques to speed up matrix operation

$$= \sigma(\boxed{W^L} \cdots \sigma(\boxed{W^2} \sigma(\boxed{W^1} \boxed{x} + \boxed{b^1}) + \boxed{b^2}) \cdots + \boxed{b^L})$$