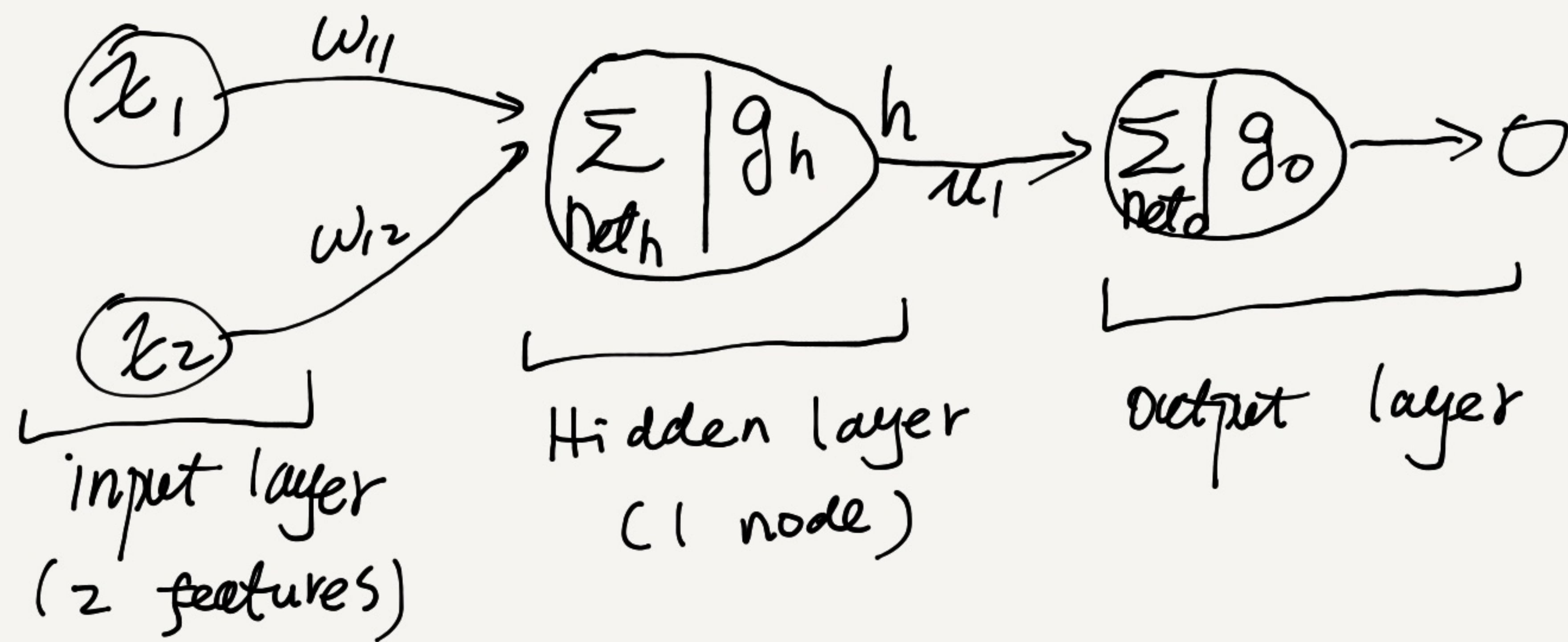


Lecture 9. ANNs (Chapter 2. The science of DL)



3-layer NN:

1. NN model

1) Input layer. no calculations.

input nodes = # features.

2) Hidden layer(s)

Net_h (pre-activation / net input): $\text{Net}_h = w_{11} \cdot x_1 + w_{12} \cdot x_2$

g_h (activation function): $h = g_h(\text{Net}_h) = g_h(w_{11}x_1 + w_{12}x_2)$

3) Output layer. $\text{Net}_o = u_1 \cdot h$, $O = g_o(\text{Net}_o) = g_o(u_1 h)$

4) Final model function

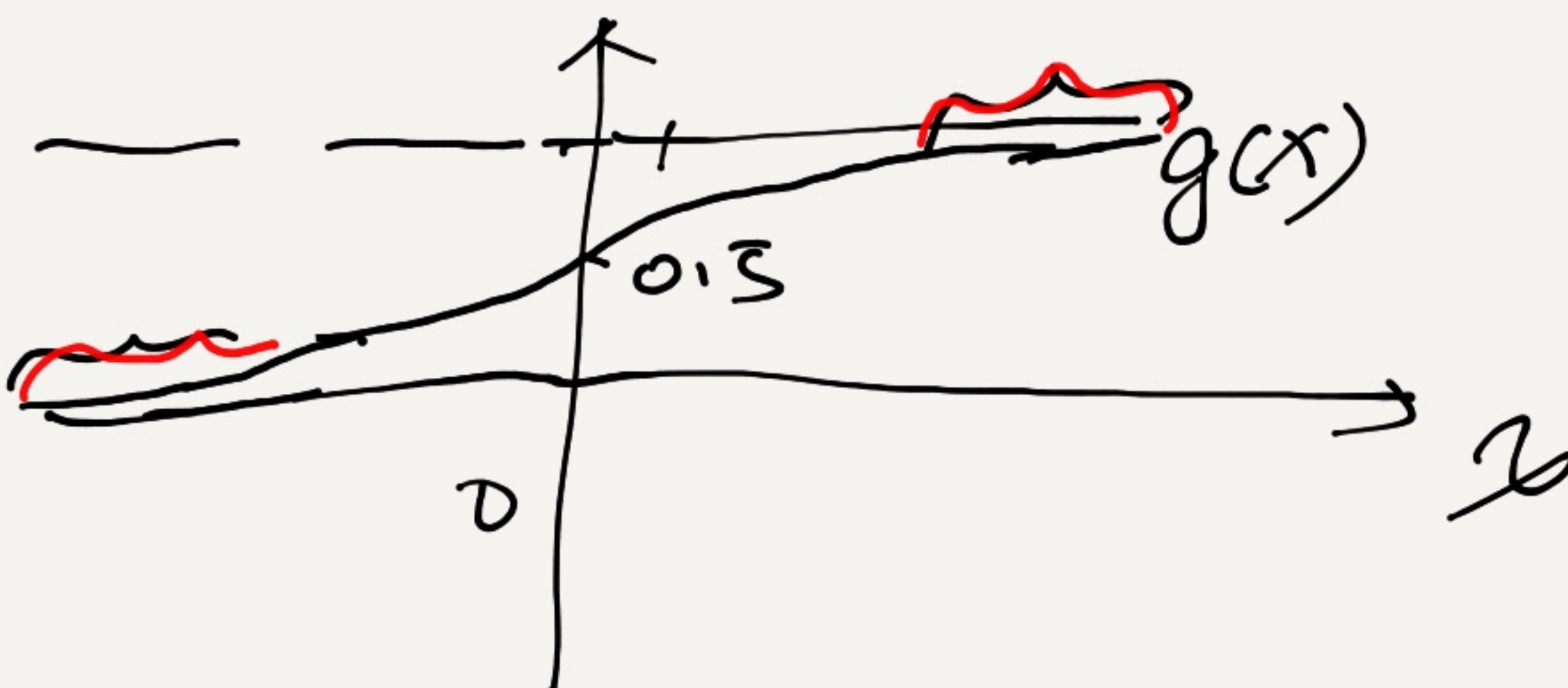
$$F(x) = g_o(u_1 \cdot h) = g_o(u_1 \cdot g_h(\text{Net}_h))$$

$$= g_o(u_1 \cdot g_h(w_{11}x_1 + w_{12}x_2))$$

2. Activation functions.

1) Sigmoid (logistic) function.

$$g(x) = \frac{1}{1+e^{-x}} \in (0, 1)$$



prior to 2011, it was the most popular activation function

It can be used for both the hidden layers and output layer
 g_h g_o

If $x \rightarrow +\infty$, or $x \rightarrow -\infty$, $g(x)$ becomes very flat.

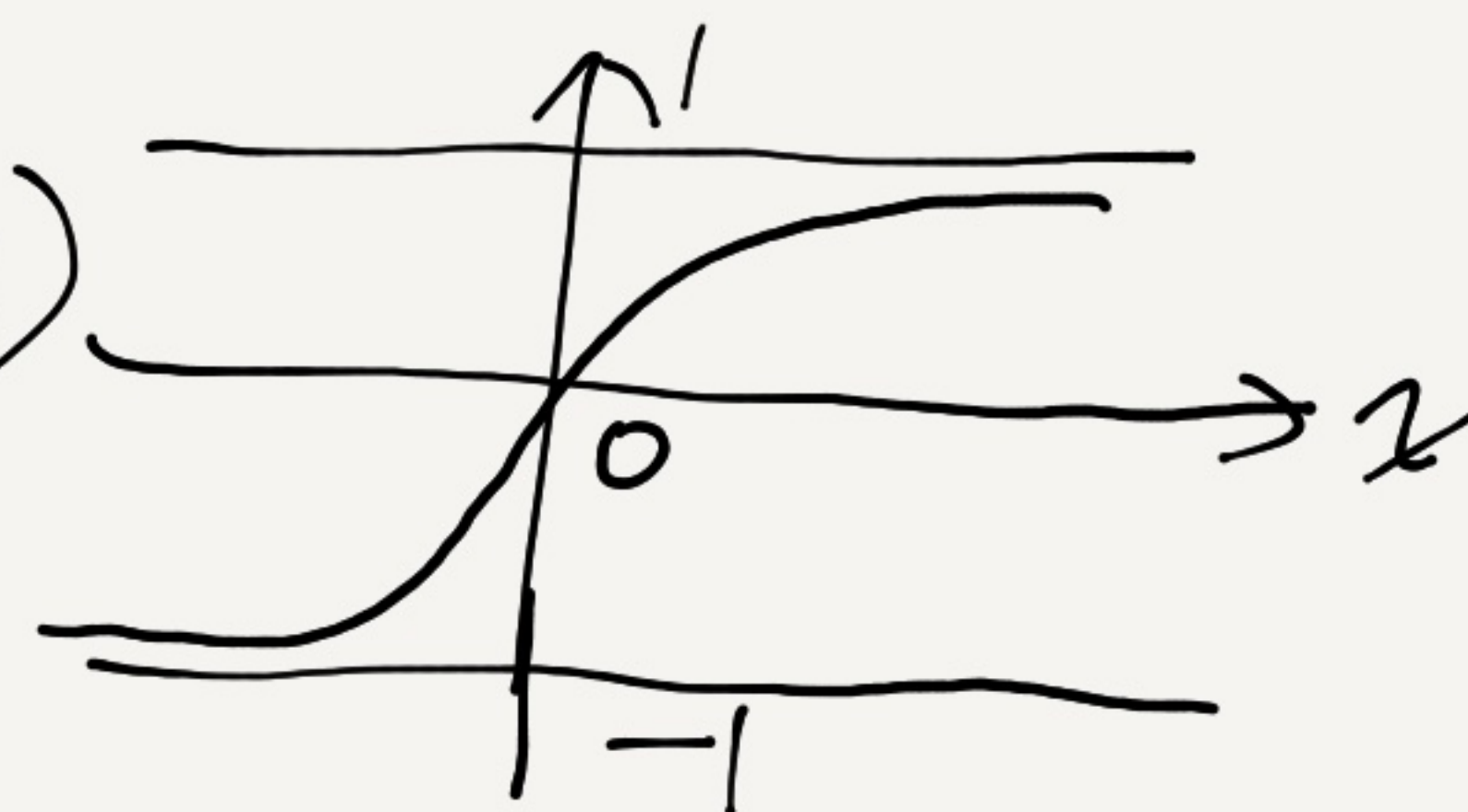
$$\rightarrow \nabla g(x) \rightarrow 0$$

$$w^{i+1} = w^i - \epsilon \cdot \overbrace{\nabla g(x)}^{\rightarrow 0} \rightarrow w^{i+1} = w^i$$

gradient vanishing

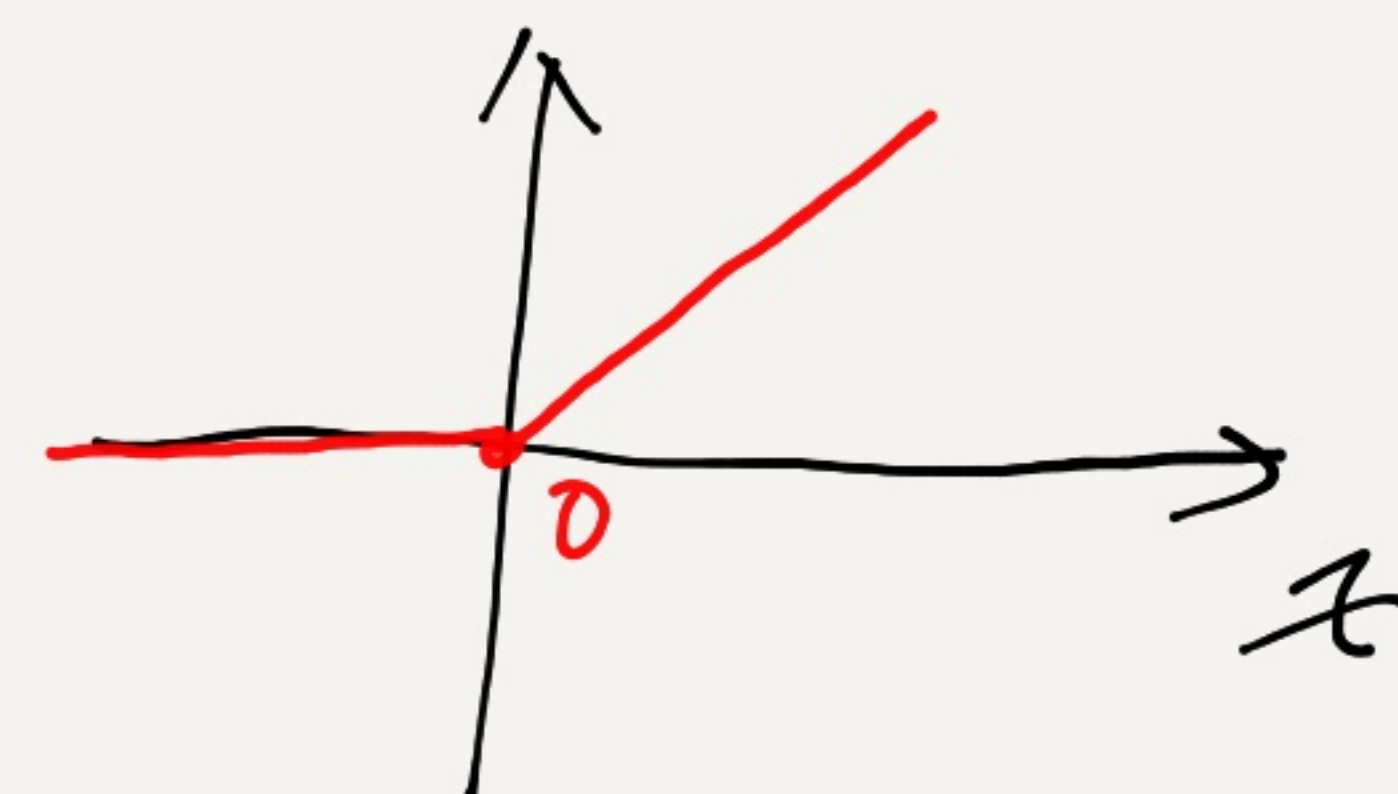
2) hyperbolic tangent.

$$g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \in (-1, 1)$$



3) Rectified Linear Units (ReLU)

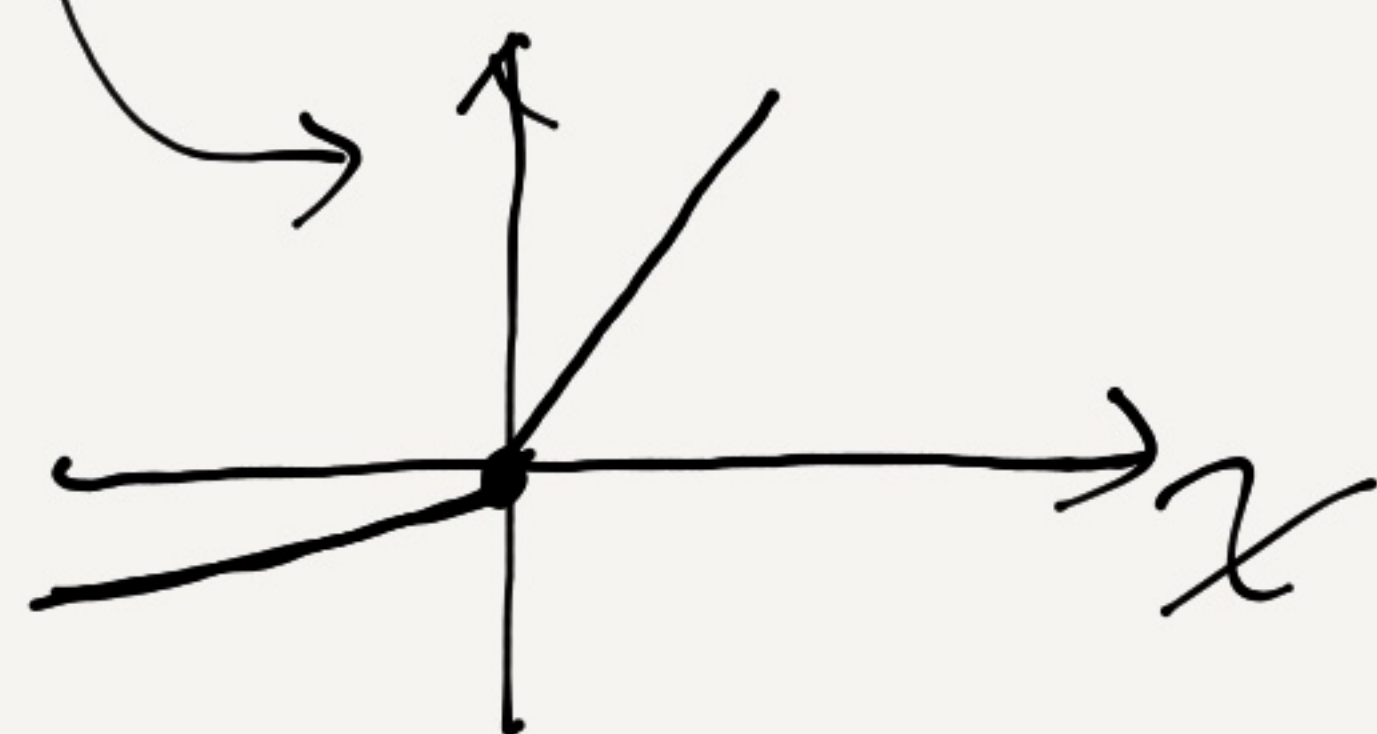
$$g(x) = \max\{0, x\} = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise.} \end{cases}$$



If $x > 0$, $\Rightarrow g(x) = 1$

ReLU becomes the default activation function for hidden layers

leaky ReLU: $g(x) = \underbrace{\max\{0, x\}}_{x > 0} - 2 \cdot \underbrace{\max\{0, -x\}}_{x < 0} = \begin{cases} x, & x > 0 \\ -2x, & x \leq 0 \end{cases}$



Swiss (2017): $f(x) = x \cdot \text{sigmoid}(\beta x)$

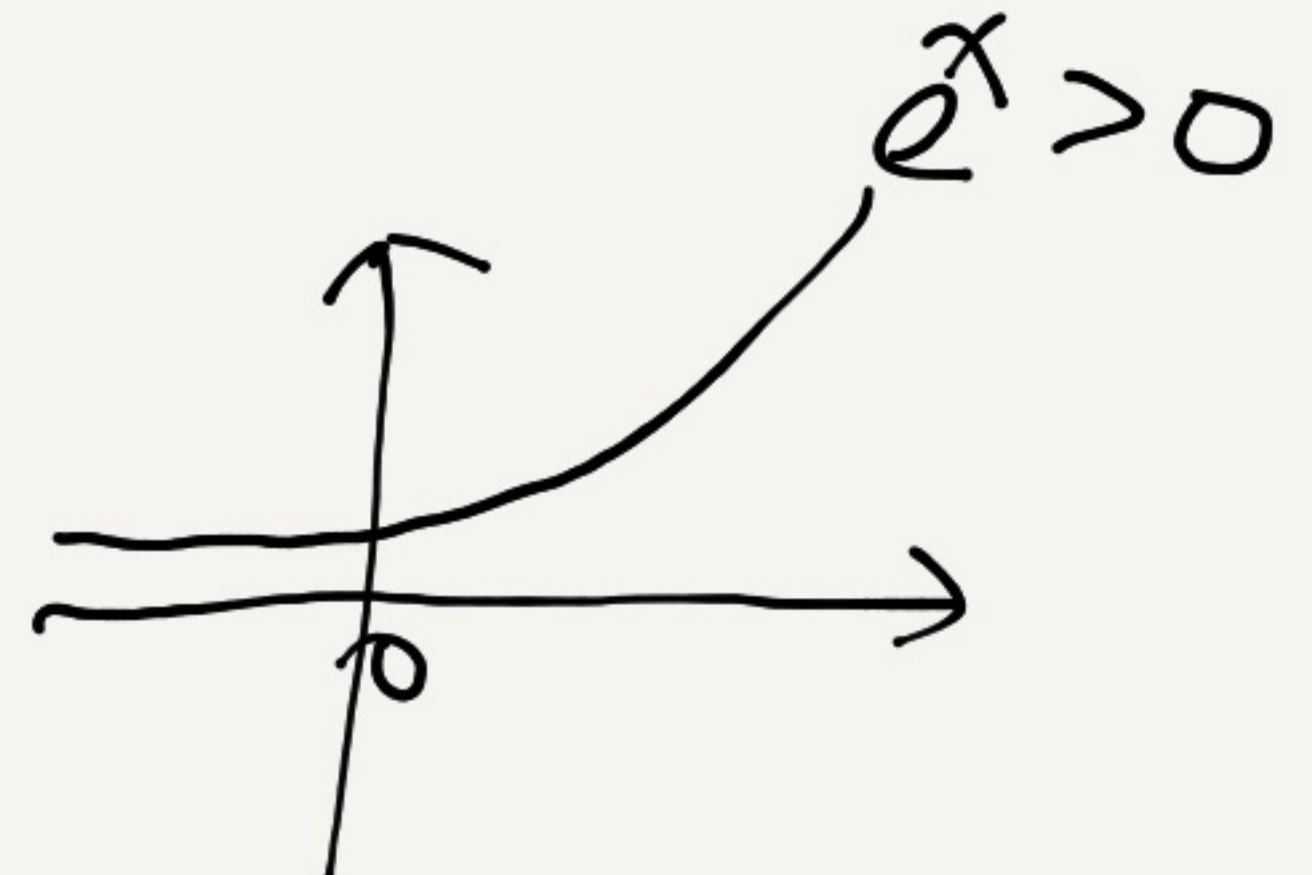
Smooth transition at point 0.

4) Softmax function. (output layer)

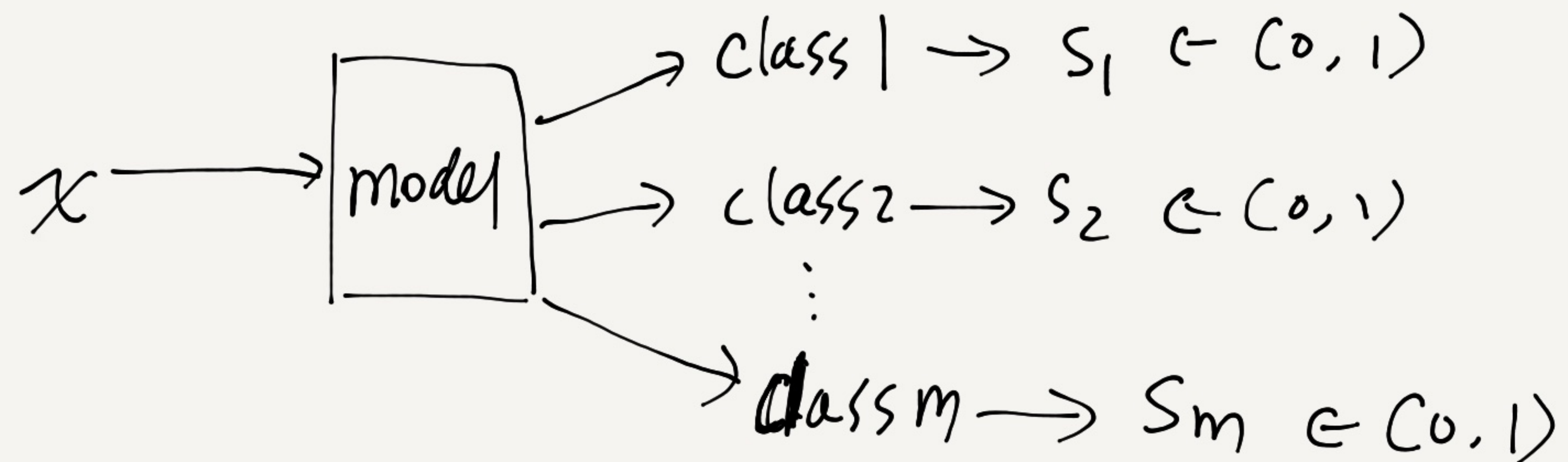
$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \quad \text{softmax}(x) = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_m \end{pmatrix}$$

$$s_i = \frac{e^{x_i}}{\sum_{i=1}^m e^{x_i}}$$

$s_i > 0, \quad \sum_{i=1}^m s_i = 1$

A hand-drawn graph of the exponential function y = e^x. The horizontal axis is labeled with '0' at the origin. The curve starts near the origin and rises steeply into the first quadrant. An arrow points to the curve with the label 'e^x > 0'.

Softmax is good for multi-class classification problem.



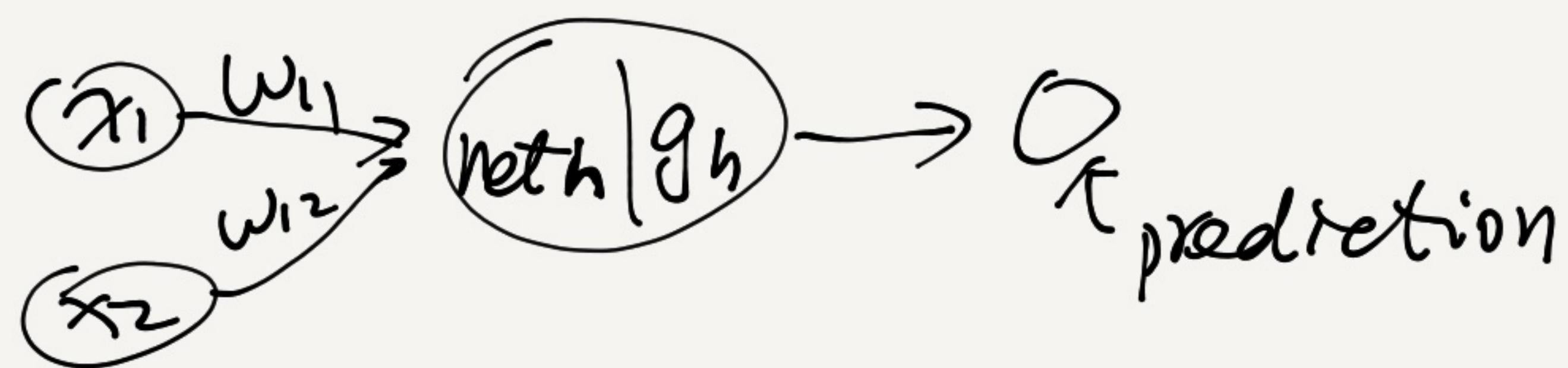
3, Loss function

$$L_0 = \frac{1}{2} (y - \underset{\substack{\uparrow \\ \text{prediction}}}{o})^2 \quad (\text{single data sample})$$

$$L = \frac{1}{n} \sum_{i=1}^n (y_i - o_i)^2 \quad (\text{multiple samples})$$

4. Optimization.

▷ 2-layer NN



$$\text{model: } O = g_h(w_{11} \cdot x_1 + w_{12} \cdot x_2)$$

$$\text{Loss: } L(w_{11}, w_{12}) = \frac{1}{2} (y - O)^2$$

$$= \frac{1}{2} (y - \underbrace{g_h(w_{11}x_1 + w_{12}x_2)}_O)^2$$

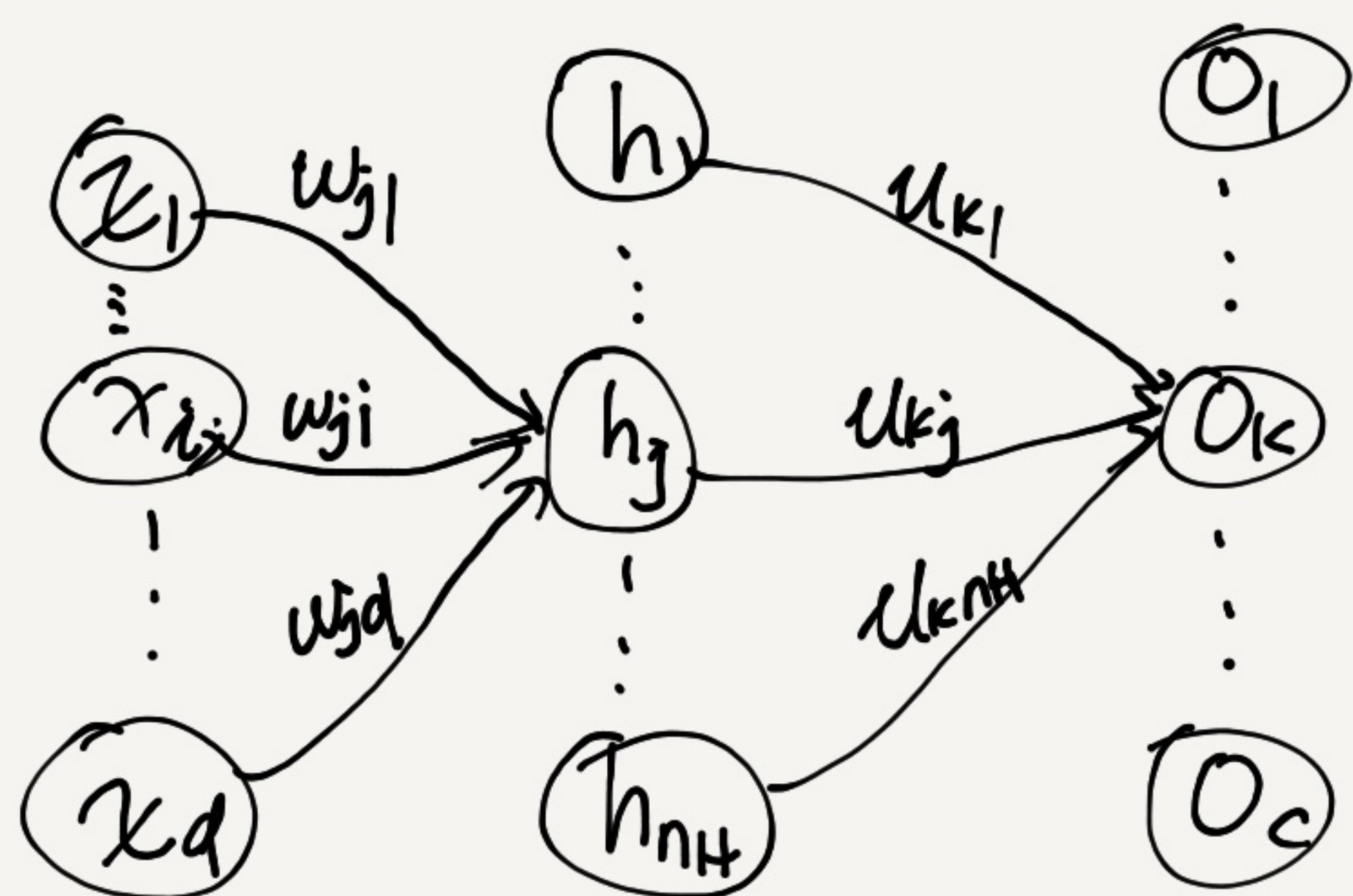
Optimization using GD:

$$\text{Weight update} \begin{cases} w_{11}^{i+1} = w_{11}^i - \epsilon \cdot \nabla_{w_{11}} L \\ w_{12}^{i+1} = w_{12}^i - \epsilon \cdot \nabla_{w_{12}} L \end{cases}$$

$$\begin{aligned} \nabla_{w_{11}} L &= \frac{\partial L}{\partial O} \cdot \frac{\partial O}{\partial w_{11}} = \overbrace{2 \cdot \frac{1}{2} (y - O) \cdot (-1)}^{\frac{\partial L}{\partial O}} \cdot g'_h \cdot x_1 \\ &= \underbrace{(O - y) \cdot g'_h}_{\text{red underline}} \cdot x_1 \end{aligned}$$

$$\nabla_{w_{12}} L = \underbrace{(O - y) \cdot g'_h}_{\text{red underline}} \cdot x_2$$

2) 3-layer NN



d : # of features

n_H : # of hidden nodes.

C : # of output nodes.

model:

$$h_j = g_h(\text{Net}_j) \quad \text{Net}_j = \sum_{i=1}^d w_{ji} \cdot x_i, \quad j=1, \dots, n_H$$

$$o_k = g_o(\text{Net}_k) \quad \text{Net}_k = \sum_{j=1}^{n_H} u_{kj} \cdot h_j, \quad k=1, \dots, C$$

$$\text{Loss} : L(W, U) = \frac{1}{2} \sum_{k=1}^C (y_k - \underbrace{o_k}_{\text{prediction}})^2$$

model parameters.

$$n_H \times d + n_H \times C$$

$$= n_H \times (d + C)$$

GD / SGD: $w_{ji} = w_{ji} - \epsilon \cdot \nabla_{w_{ji}} L$?

$u_{kj} = u_{kj} - \epsilon \cdot \nabla_{u_{kj}} L$?