# Lecture 14    ANNs : Part II .

## 1. Machine learning Components.

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y} \longleftrightarrow y$$

1) Data :  $(x, y) \rightarrow$ input data pair.
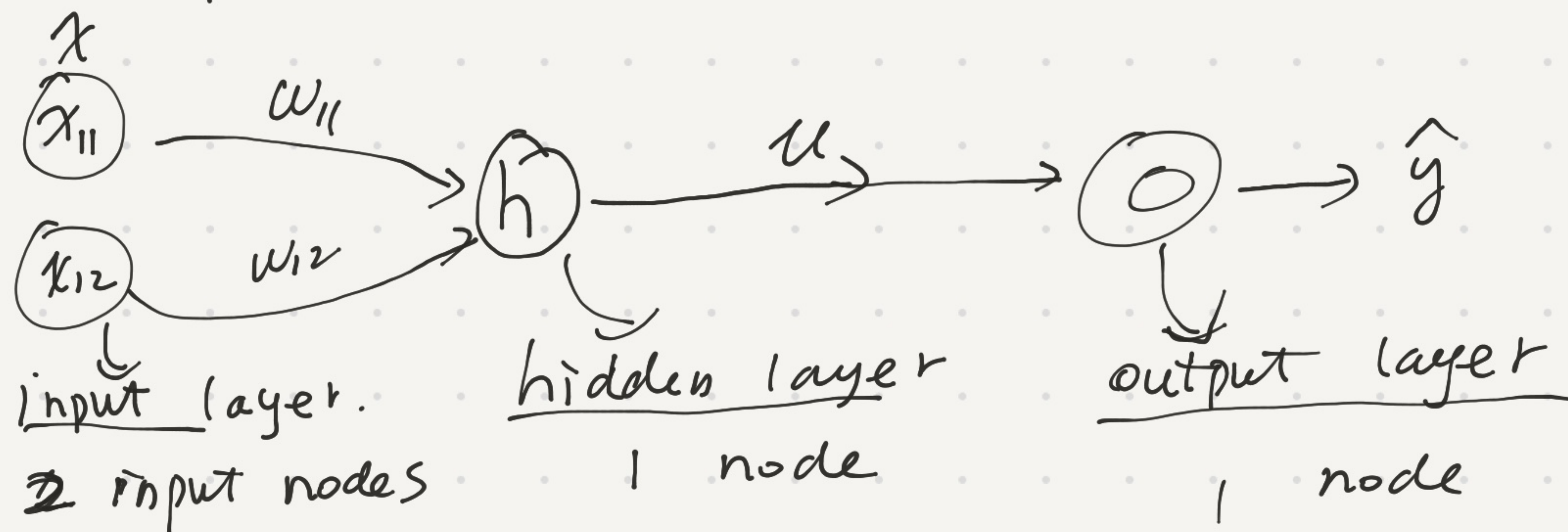
2) Model:  $f$ : maps $x$ to $y$.
   logistic regression :  $f(x) = \dfrac{1}{1 + e^{-(w^T x)}}$

   $$= \sigma(w^T \cdot x)$$

3) Loss function :  $\mathcal{L} = \dfrac{1}{N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)^2$   MSE

4) Optimizer / optimization Algorithms., find a mode
   $f$ than minimizes the loss function $\mathcal{L}$

## 2. A simple NN.



input layer.
2 input nodes

hidden layer
1 node

output layer
1 node

### 2.1. input nodes: No specific operations

the # of input nodes = the # of features

### 2.2 Hidden layer: extracts a set of features from the input data.



We could have multiple hidden layers in a NN; Multiple hidden nodes in each layer.

Net input: $Net_H = x_{11} \cdot w_{11} + x_{12} \cdot w_{12} = w^T \cdot x$

Activation function: $g_H(Net_H) = \dfrac{1}{1 + e^{-net_H}}$

non-linear function

logistic / sigmiod

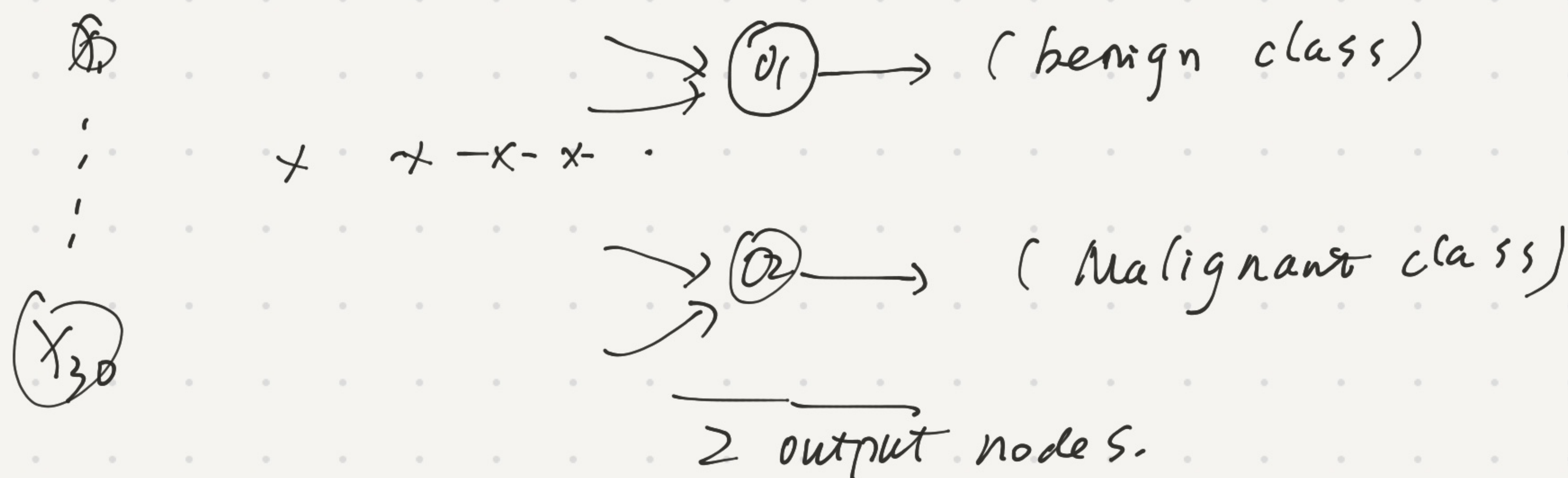## 2.2. Output Layer: Converts hidden output to prediction $\hat{y}$

$h \xrightarrow{U} \boxed{Net_0 | g_0} \xrightarrow{O}$    $Net_0 = h \cdot U$.    $g_0(Net_0)$: activation function

We can have only one output layer in a NN, and we can have multiple output nodes.

How many output nodes do we need in a NN?

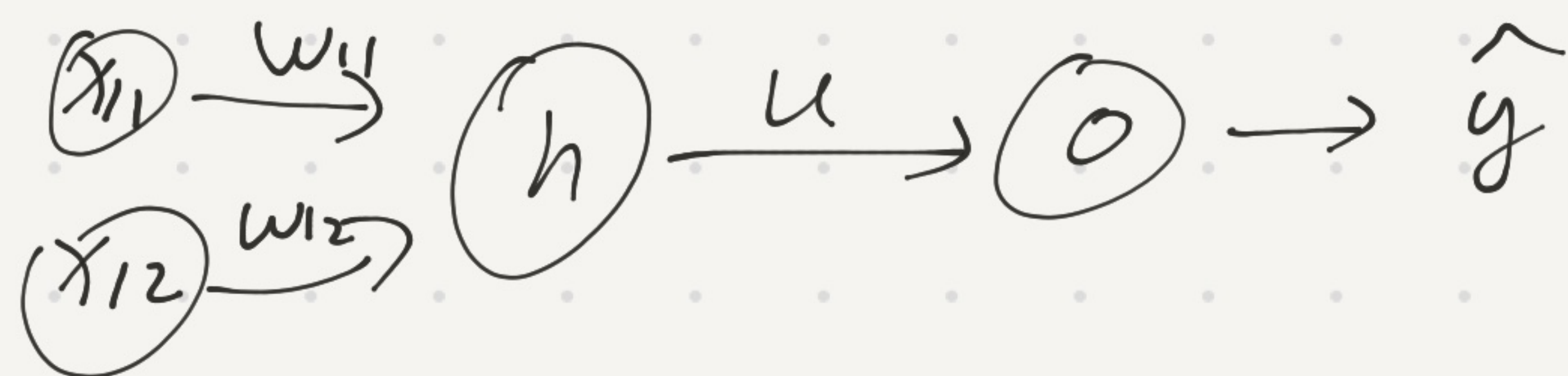Classification: the number of output nodes = the number of classes.

For example: Breast Cancer Classification

$X_1$

$\vdots$

$+$ $+$ $-x- x-$ $\cdot$

$X_{30}$

$\xrightarrow{\quad} \boxed{O_1} \xrightarrow{\quad}$ (benign class)

$\xrightarrow{\quad} \boxed{O_2} \xrightarrow{\quad}$ (Malignant class)

$\underbrace{\qquad\qquad}$
2 output nodes.

## 2.3 Final mode of the simple NN

$$(x_{11}) \xrightarrow{w_{11}} (h) \xrightarrow{u} (0) \to \hat{y}$$
$$(x_{12}) \xrightarrow{w_{12}}$$
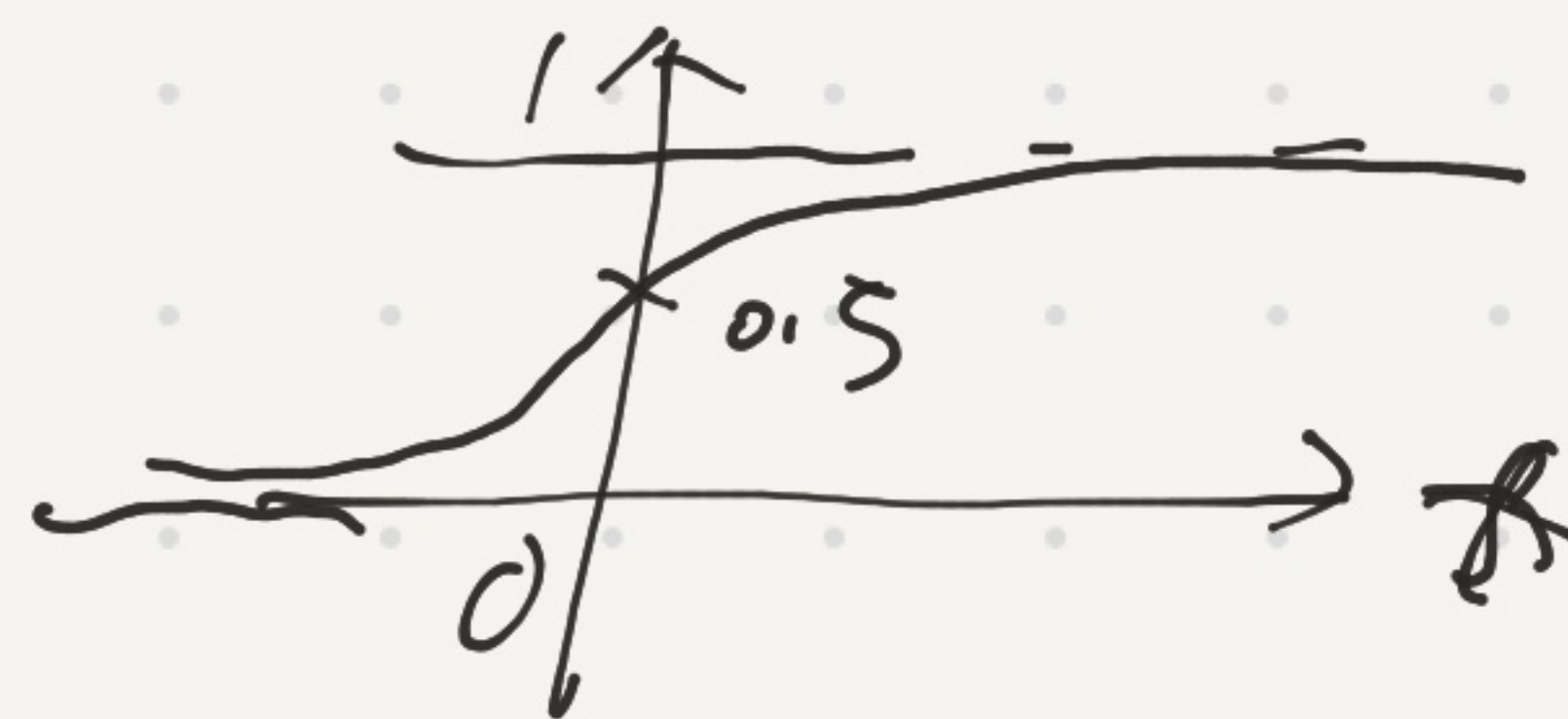
$$f(x) = g_o(Net_o) = g_o(u \cdot \underline{h})$$

$$= g_o(u \cdot g_H(\underline{Net_H}))$$

$$= g_o(u \cdot g_H(x_{11} \cdot w_{11} + x_{12} \cdot w_{12})) = \hat{y}$$

## 2.4 popular activation functions
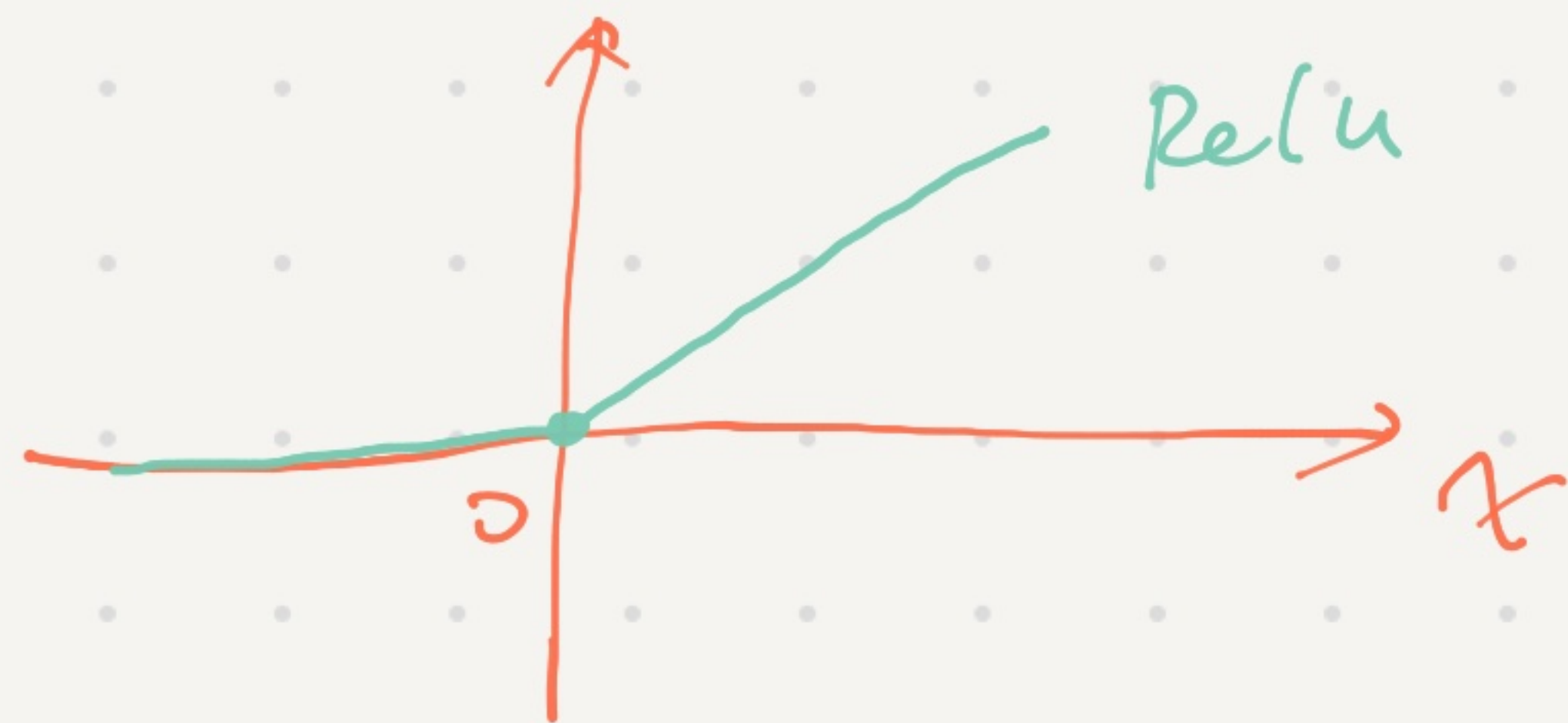
1) Sigmoid function

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



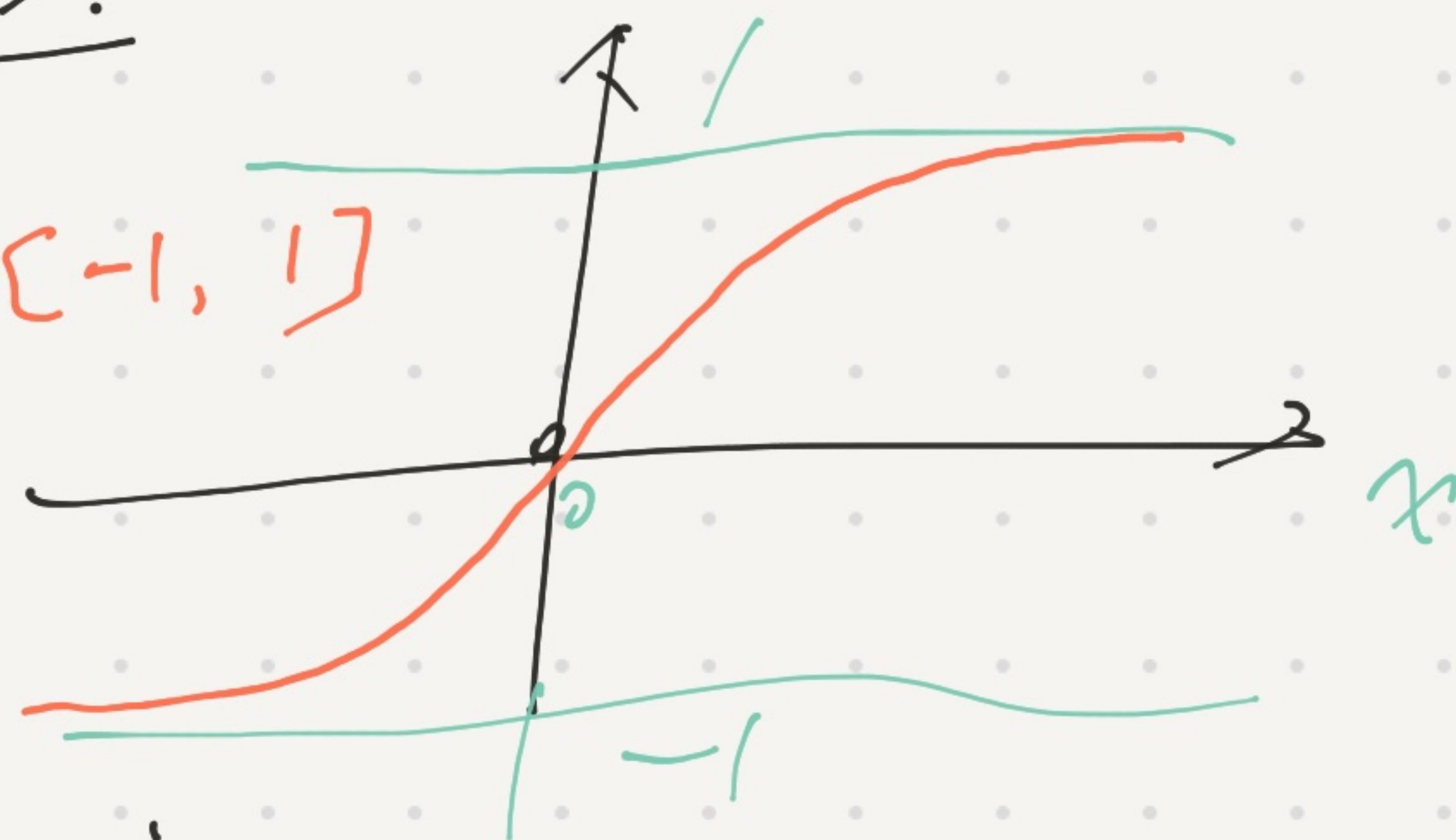we can use it for both hidden and output nodes

2) Rectified linear units (ReLu)

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

ReLu

was proposed for deep NNs. and used **for**

**hidden layers.**

3) $tanh(x) \in [-1, 1]$

both hidden and output layer.

4) **Softmax (x).**

↓

output layer.

input

$$x = \begin{pmatrix} x_1 \\ \vdots \\ \vdots \\ x_n \end{pmatrix}$$

$$\text{Softmax}(x) = \begin{pmatrix} S(x_1) \\ S(x_2) \\ \vdots \\ S(x_n) \end{pmatrix} \quad \forall \, ' \quad S(x_i) = \frac{e^{x_i}}{\sum\limits_{j=1}^{n} e^{x_j}} \in (0,1)$$

$$\begin{cases} S(x_1) + S(x_2) + \cdots S(x_n) \\ = 1 \end{cases}$$

## 3. Loss function

Binary Cross-entropy

Classification: $\mathcal{L} = - \sum\limits_{i=1}^{N} \left[ \underbrace{y_i \log \widehat{y}_i}_{①} + \underbrace{(1-y_i) \log(1-\widehat{y}_i)}_{②} \right]$
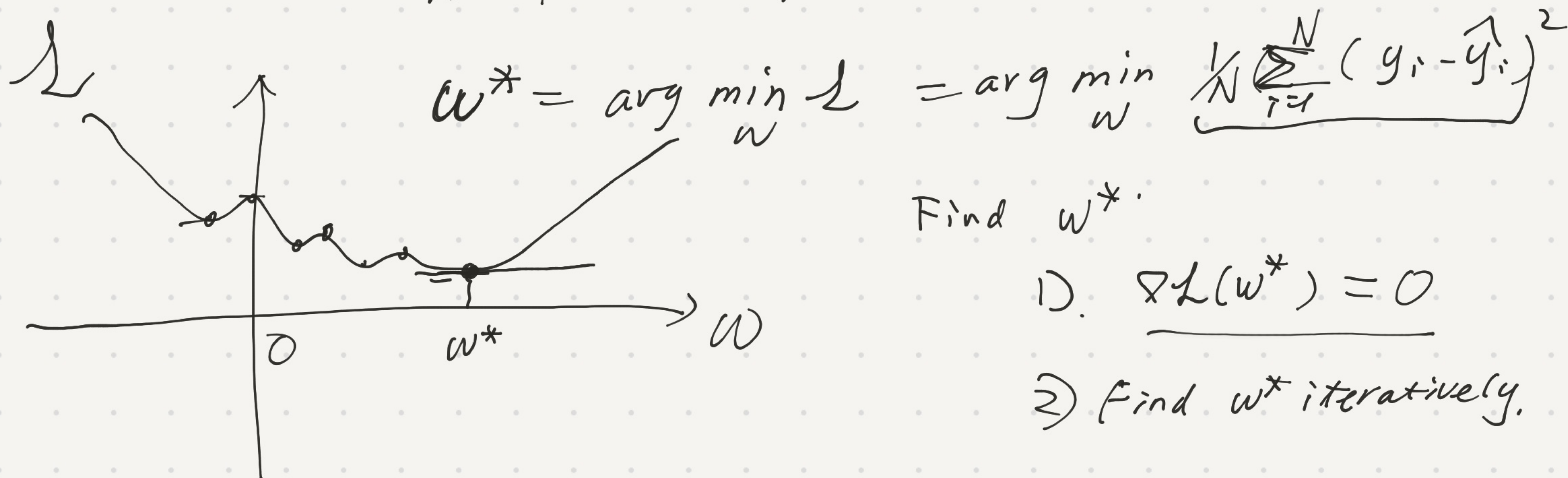
Regression: MSE

$$\mathcal{L} = \frac{1}{N} \sum\limits_{i=1}^{N} (y_i - \widehat{y}_i)^2$$

## 4. Optimizers: learn model paramers than minimize $\mathcal{L}$

$f(x) \longrightarrow (W_{11}, W_{12} \text{ and } u)^T = \mathbf{w}$
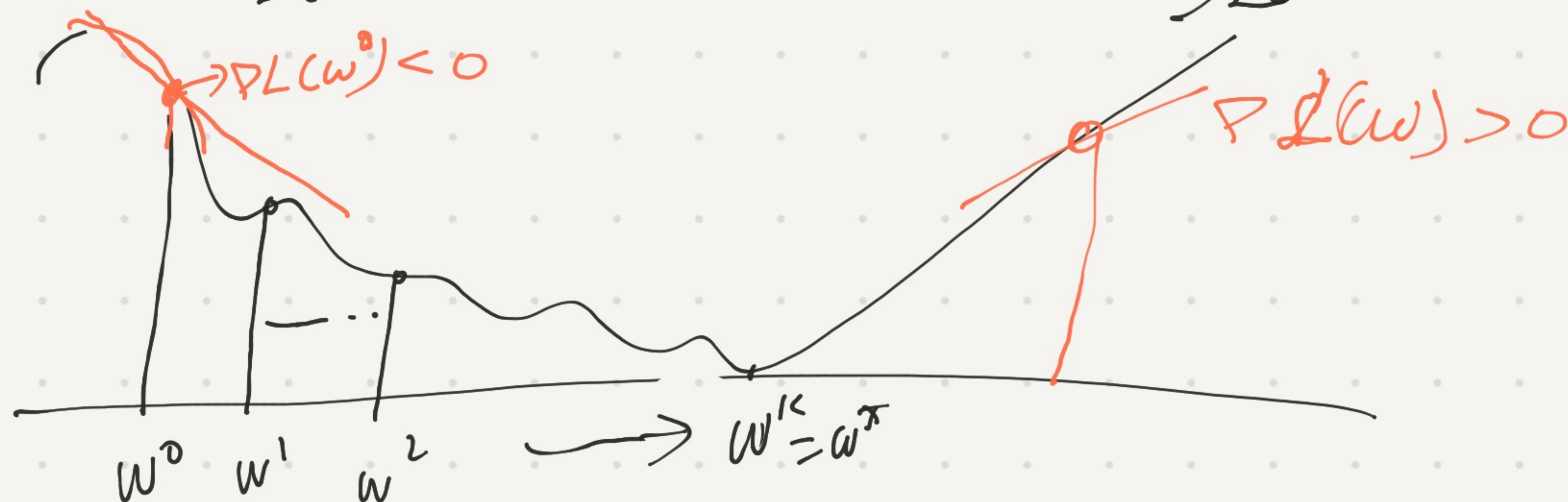model          parameters:

$$w^* = \arg\min_{w} \mathcal{L} = \arg\min_{w} \frac{1}{N} \sum\limits_{i=1}^{N} (y_i - \widehat{y}_i)^2$$



Find $w^*$.

1) $\nabla \mathcal{L}(w^*) = 0$

2) Find $w^*$ iteratively.

Find a sequence.  $\longrightarrow$ initial parameter.

$\omega^0, \omega^1, \omega^2, \omega^3, \cdots, \omega^k = \omega^*$

$L(\omega^0) > L(\omega^1) > L(\omega^2) > L(\omega^3) \cdots > L(\omega^*)$

$\nabla L(\omega^0) < 0$

$\nabla L(\omega) > 0$

$\omega^0 \quad \omega^1 \quad \omega^2 \quad \longrightarrow \quad \omega^k = \omega^*$

produce the sequence:

1) moving direction $\Bigg\}$ $\longrightarrow$ Gradient descent: (GD)

1) step size

$\nabla L(\omega^0) > 0 \longrightarrow$ move to negative direct.

$\nabla L(\omega^0) < 0, \longrightarrow$ move to the positive direction