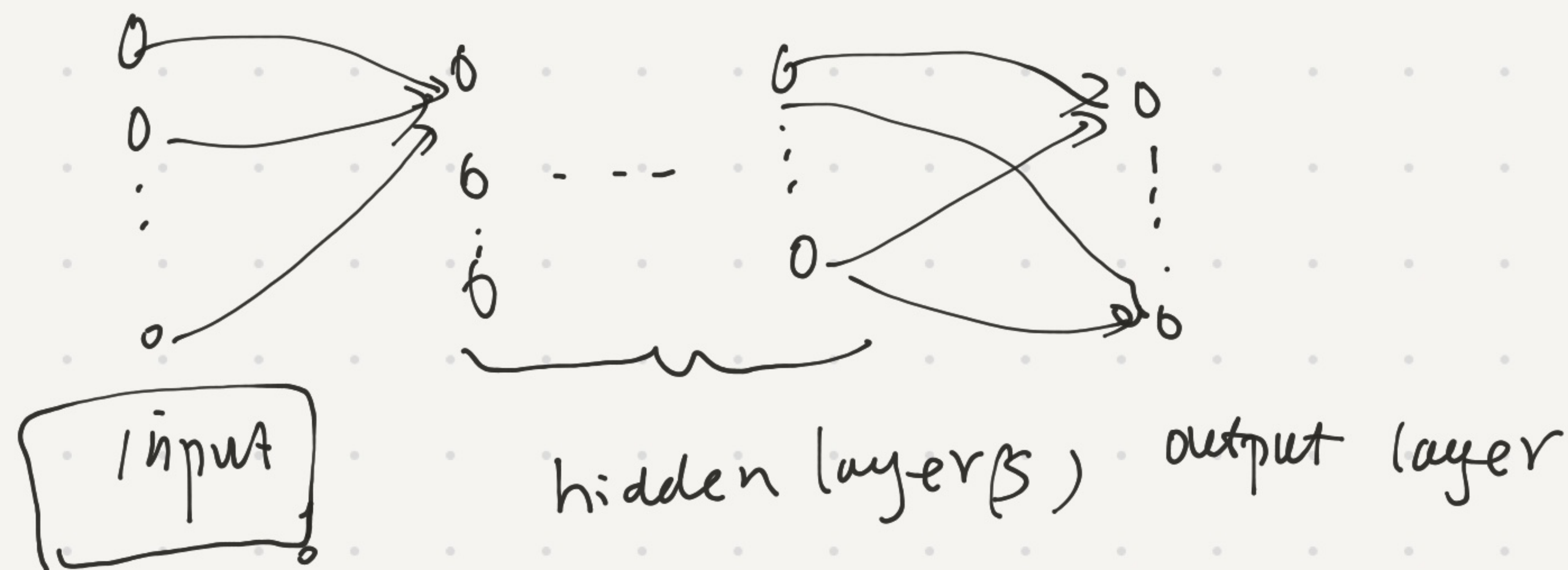


Lecture 16 & 17 Tips for NN design



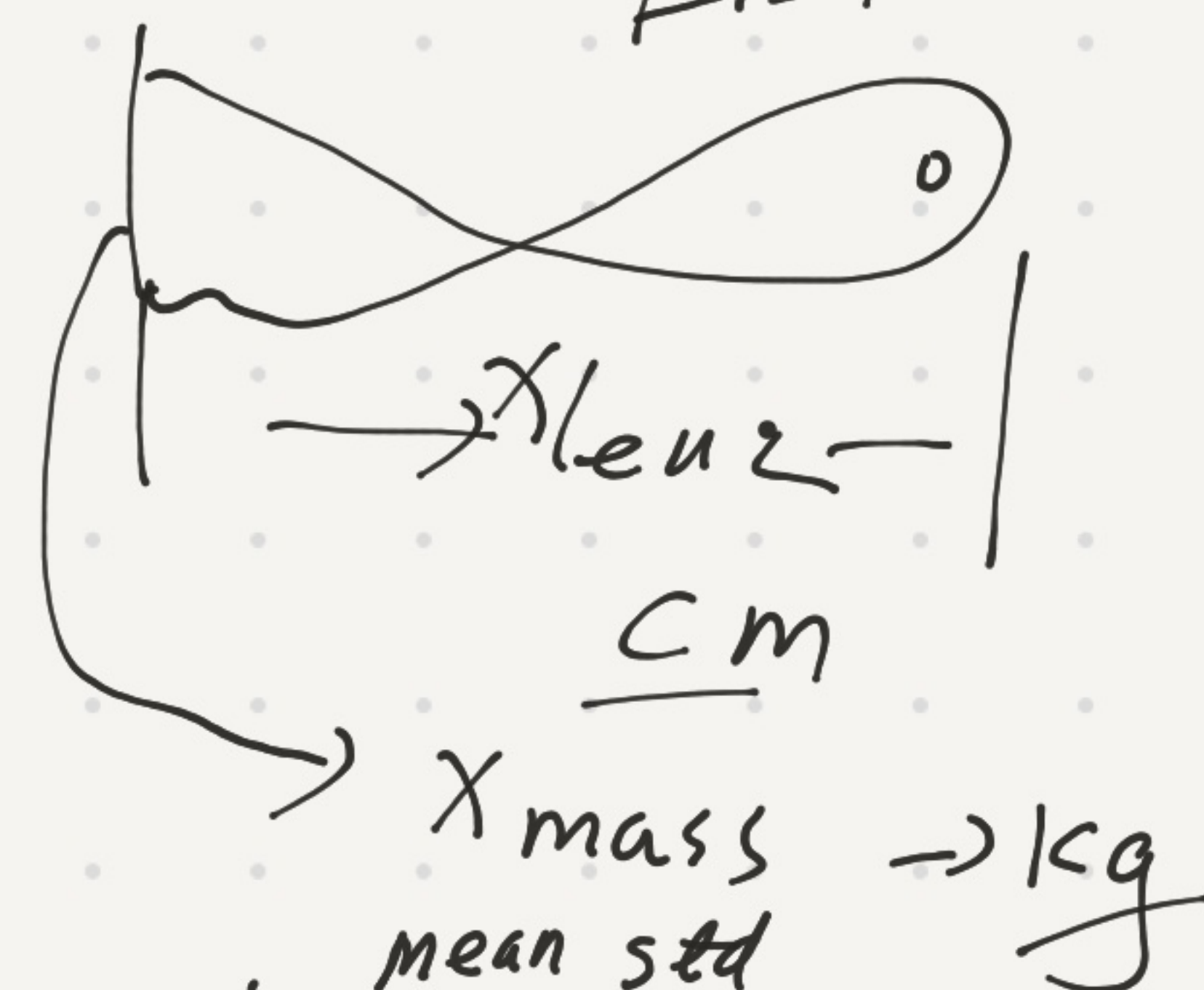
① Scaling input vectors.

$X = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$ → features may have different scales

feature vectors $x_{len} \in [0, 200] \text{ cm}$

$x_{mass} \in [0, 40] \text{ kg}$

Fish



We can use standardization tech.

$$x'_{len} = \frac{x_{len} - \text{mean}(x_{len})}{\text{std}(x_{len})} \sim N(0, 1)$$

normal distribution

$$x'_{mass} = \frac{x_{mass} - \text{mean}(x_{mass})}{\text{std}(x_{mass})} \sim N(0, 1)$$

Normalization \rightarrow convert values to $[0, 1]$

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \in [0, 1]$$

x : input feature.

x' : output feature.

(2) Determine the # hidden nodes.

hidden nodes determines # of parameters.

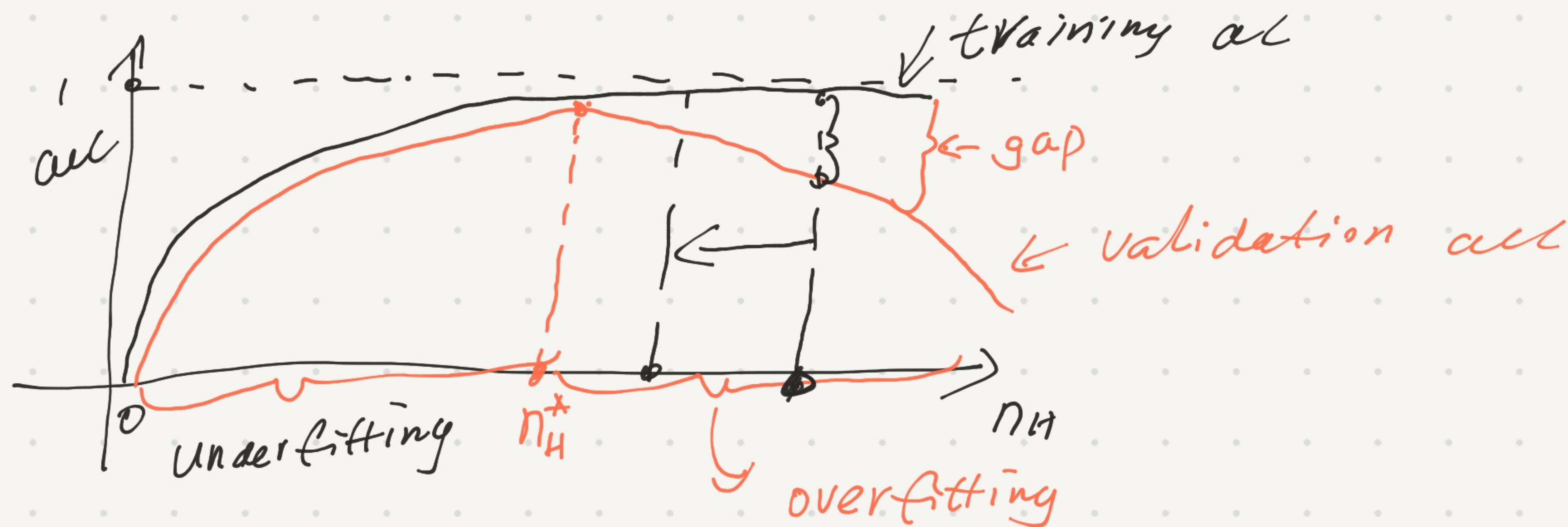
- - - the complexity of the model.

3-layer NN: $\underbrace{(N_{in} \times N_H)}_{\text{input nodes}} + (N_H \times \underbrace{N_{out}}_{\text{output nodes}})$: Parameters.

$$= \underbrace{(N_H)}_{\text{hidden nodes}} \times (\underbrace{N_{in}}_{\text{input nodes}} + \underbrace{N_{out}}_{\text{output nodes}})$$

How to choose N_H :

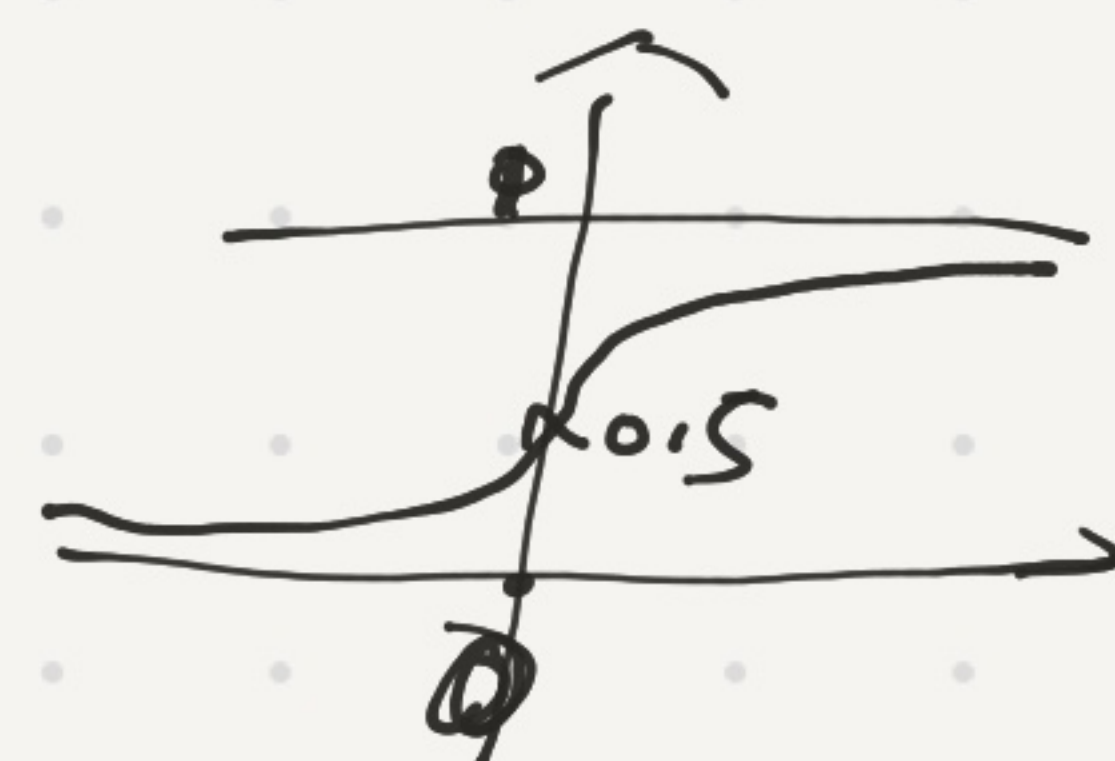
- Small-scale NN: $\begin{cases} N_H < \text{total \# of training samples} \\ \text{choose an } N_H \text{ to make total \# of parameters} \approx \frac{n}{10} \end{cases}$
- Large-scale NN: start from large N_H , then 'decay' N_H \rightarrow training samples



③ Activation functions

Sigmoid / logistic function:
most popular act before 2011.

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



hyperbolic tan. $\tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$



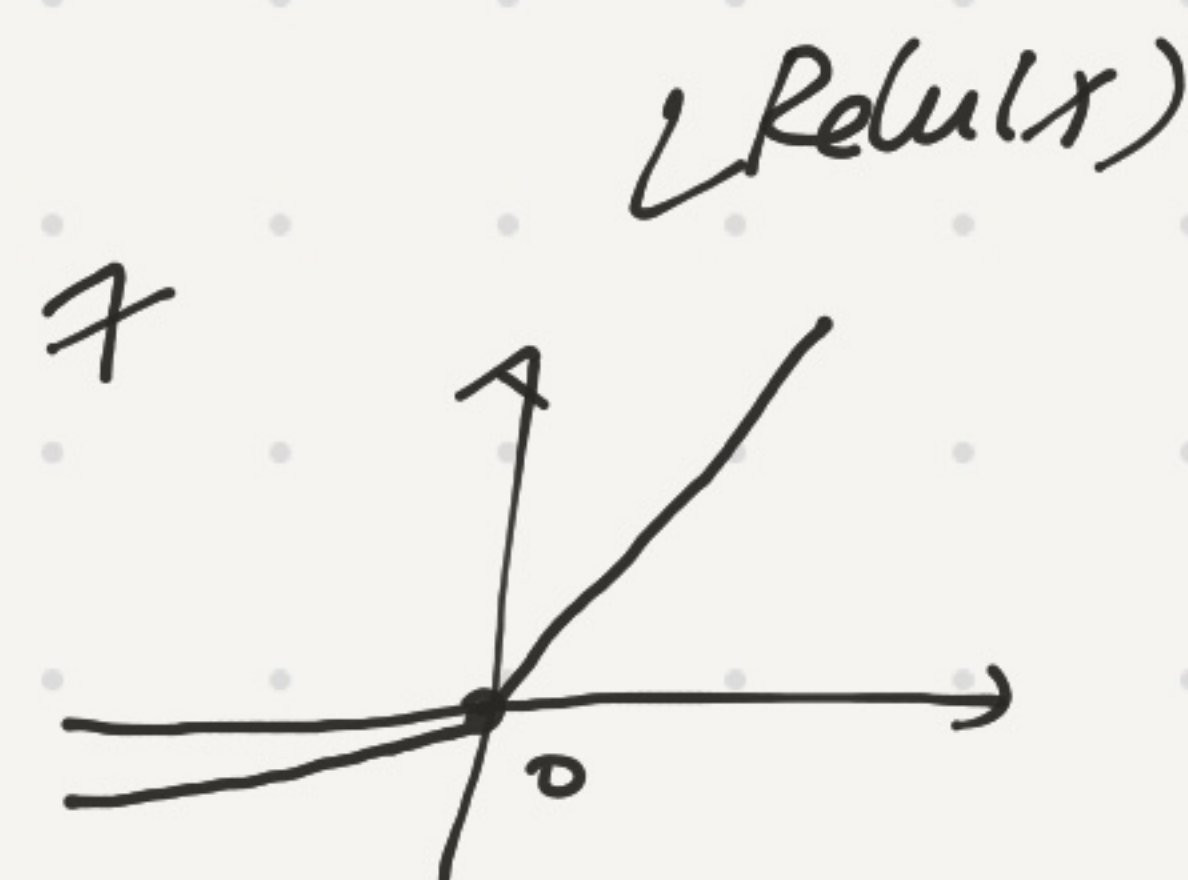
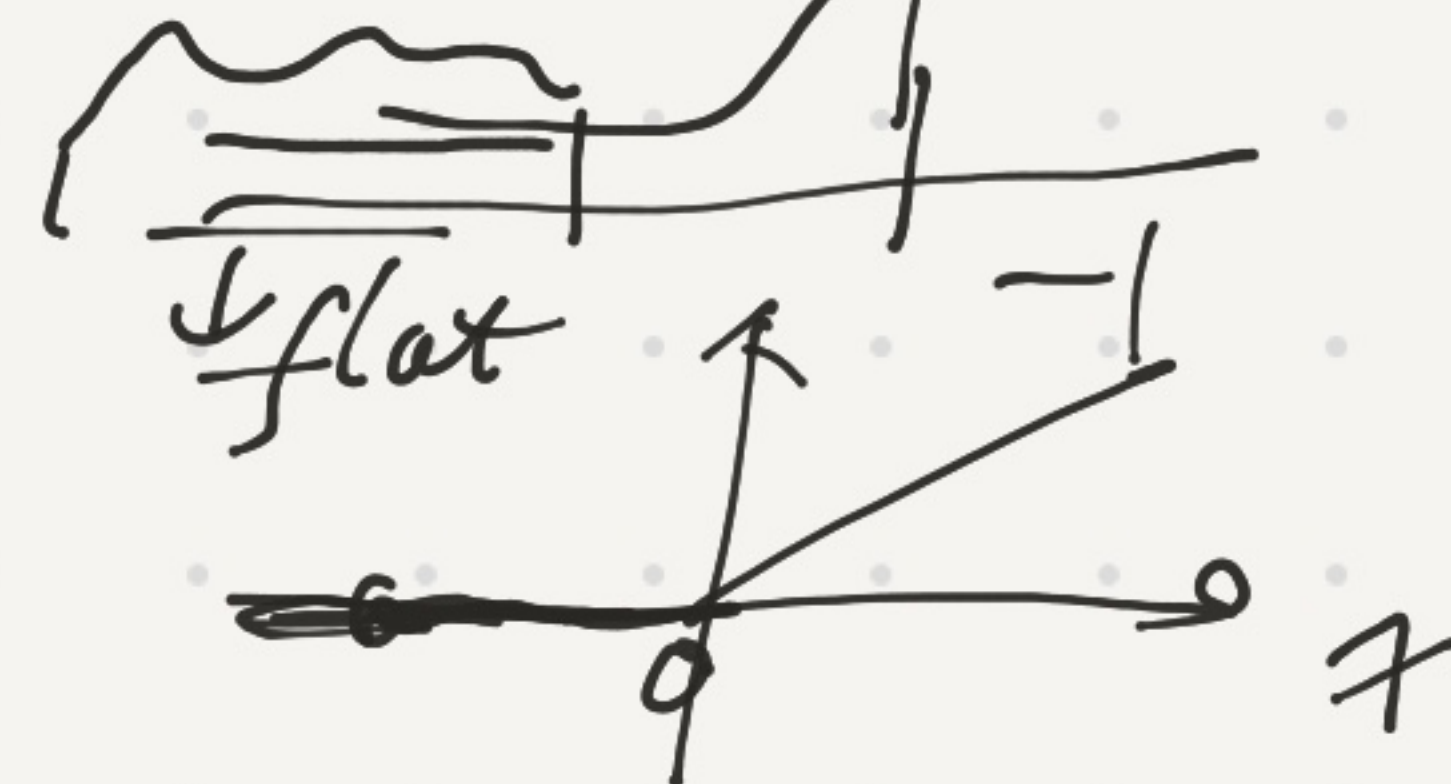
Deep neural Networks.

Relu (2011) $\text{Relu}(x) = \max\{0, x\}$

leaky Relu (2014) $\text{LRelu}(x) = \begin{cases} x & x > 0 \\ 0.01x & \text{otherwise} \end{cases}$

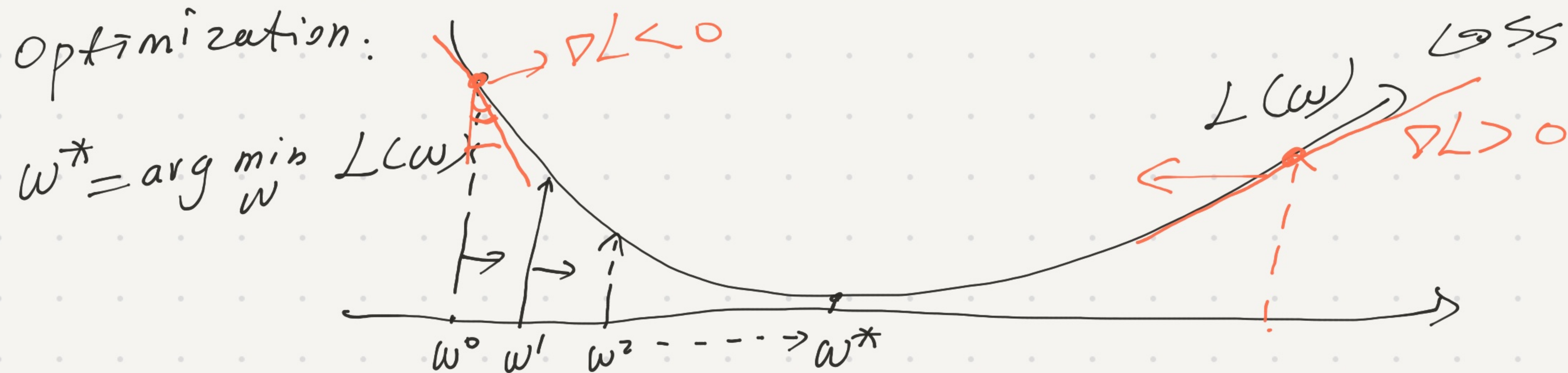
Swish(x) = $x \cdot \text{Sigm}(x)$

Elu(x) = $\begin{cases} x & x > 0 \\ 0(e^x-1) & \text{otherwise} \end{cases}$



④ Batch and epochs

Optimization:



Gradient descent approach: $w^0 \rightarrow w^1 \rightarrow w^2 \rightarrow \dots \rightarrow w^k = w^*$

w^0 : initial parameter(s): generate randomly,

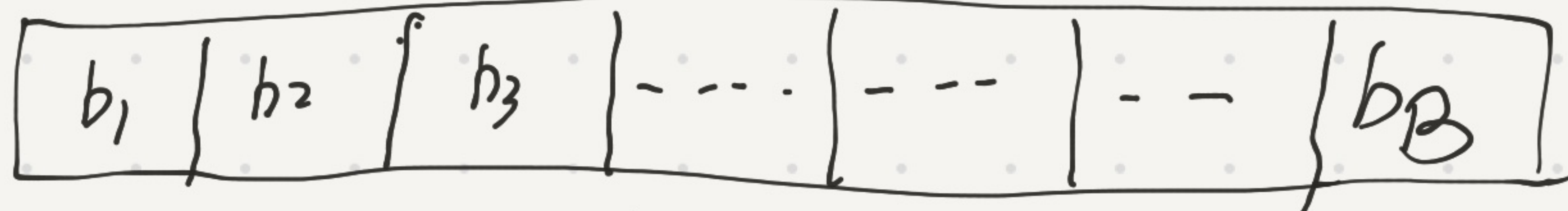
$$w^{i+1} = \underbrace{w^i}_{\text{previous parameter}} + \underbrace{\epsilon \cdot (-\nabla_{w^i} L)}_{\text{learning rate : } 0.0001} \rightarrow \text{gradient of } L \text{ over } w^i$$

$$w^1 = w^0 + \epsilon \cdot (-\nabla_{w^0} L) \rightarrow \text{can make sure that we are approaching } w^*$$

$\nabla_{w^i} L \rightarrow$ will involve all training data \rightarrow expensive and infeasible.

mini-batch optimization approach.

training set



split the training set into B mini-batches; each iteration, only one

Batch is involved in $(\nabla_{w^i} L)$

For epoch in # of max epochs:
shuttle the training set; prepare the minibatches

For l to B :

$$w^{i+1} = w^i + \nabla_{w^i} L$$