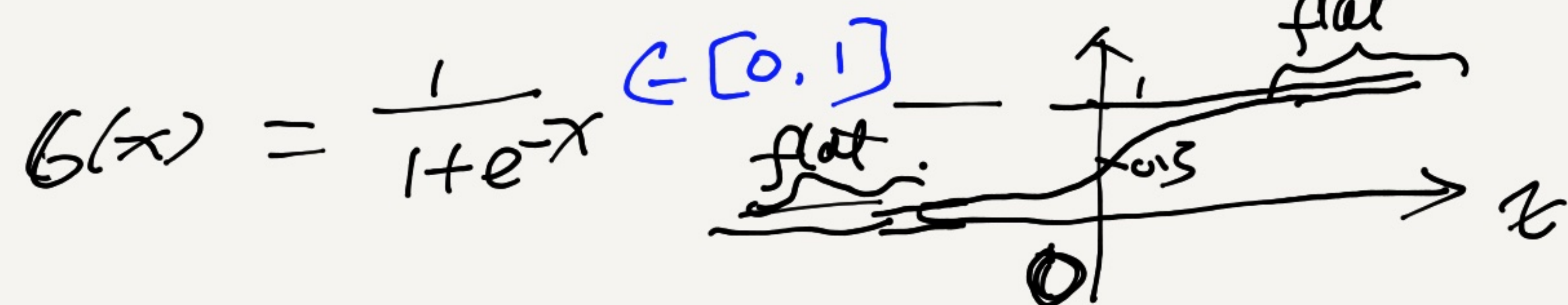


Lecture 13.

1. Vanishing gradient problem

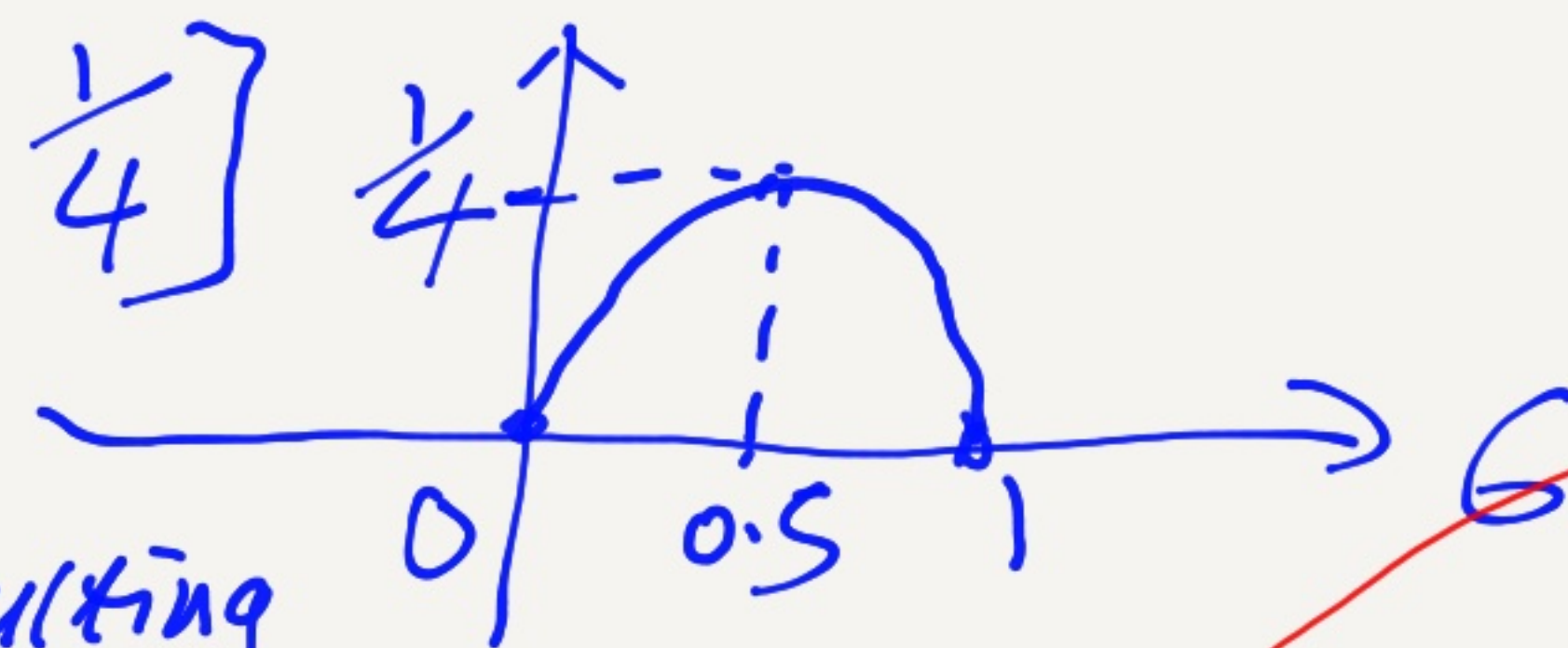
1) Assume a multi-layer NN, with sigmoid activation functions. (all nodes)



- ① If x is too large or too small, $\sigma(x)$ will saturate (0 or 1) and $\sigma'(x)$ will approach 0.

$$\nabla_{w_{kj}} L = \delta_k \cdot h_j = (o_k - y_k) \cdot \sigma' \cdot h_j \rightarrow 0 \text{ if } \sigma' \rightarrow 0$$

- ② If we use the BP through σ reduce the gradient by a factor of at least 4. $\sigma' = \sigma(1-\sigma) \in [0, \frac{1}{4}]$



- ③ Propagating through several layers, the resulting gradients become very small.

Weights between input and hidden layer

$$\nabla_{w_{ij}} L = \delta_j \cdot x_i = \sum_{k=1}^c (\delta_k \cdot w_{kj}) \cdot \sigma' \cdot x_i$$

activation of hidden layer

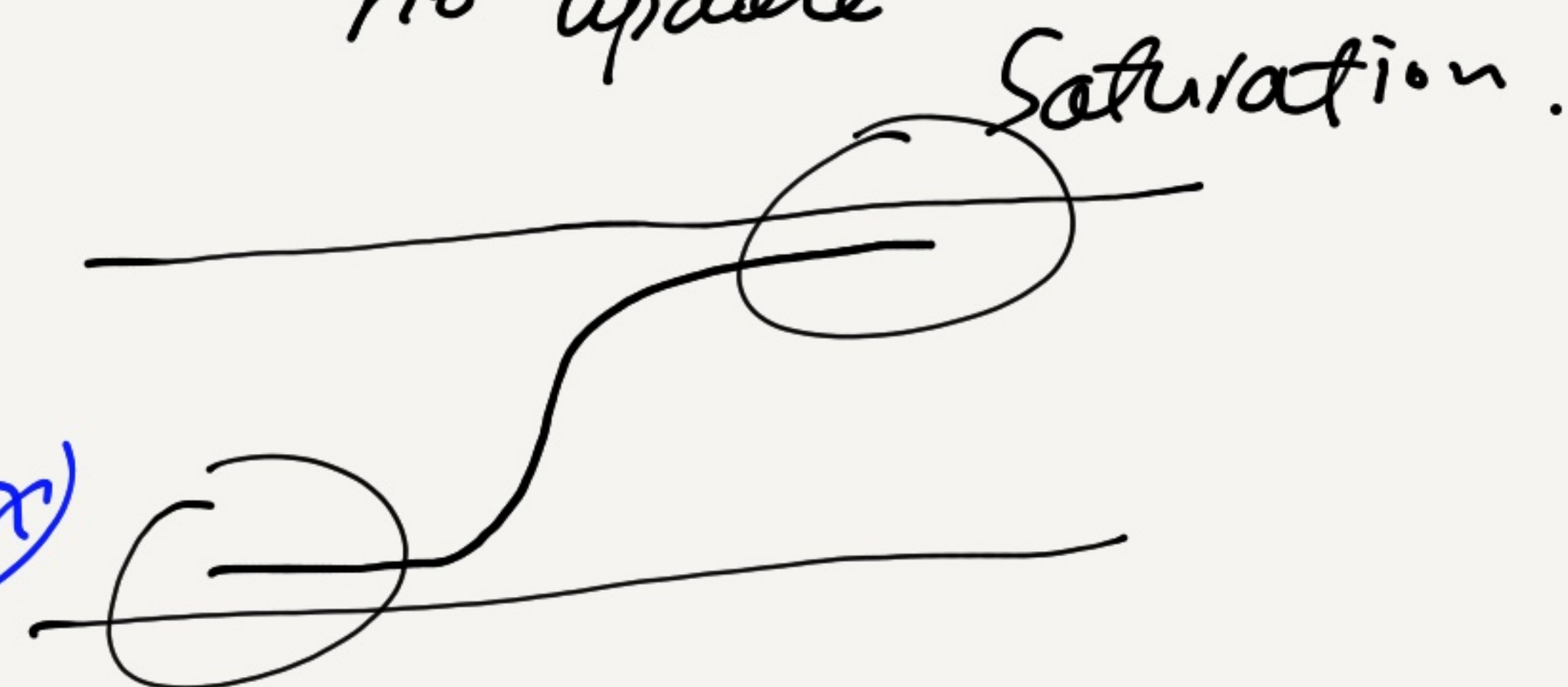
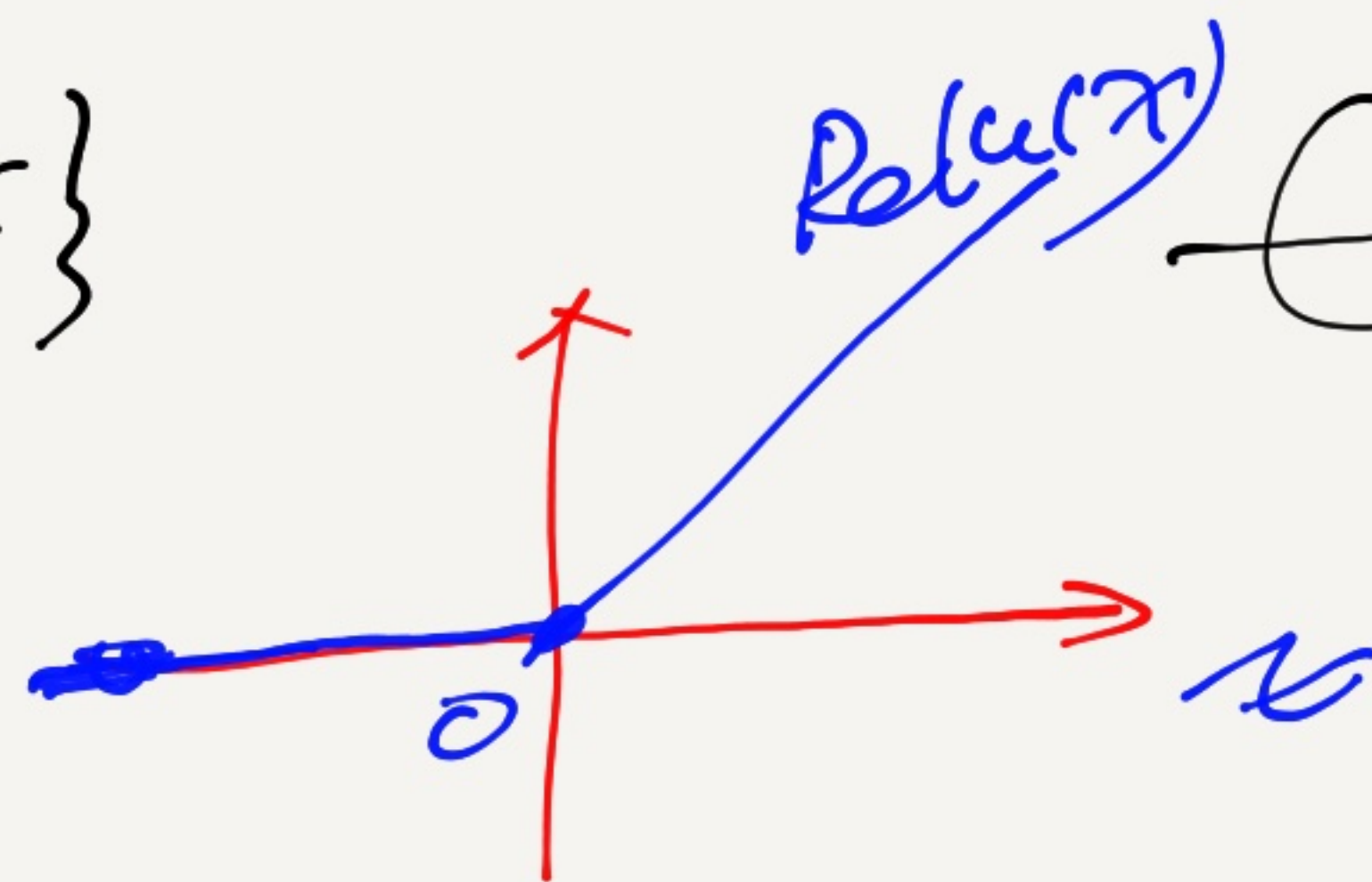
GD/SGD: $\omega^{i+1} = \omega^i - \underbrace{\epsilon \cdot \nabla_{\omega^i} L}_{\rightarrow 0} \rightarrow \omega^{i+1} = \omega^i$
no update

2) solve the issue.

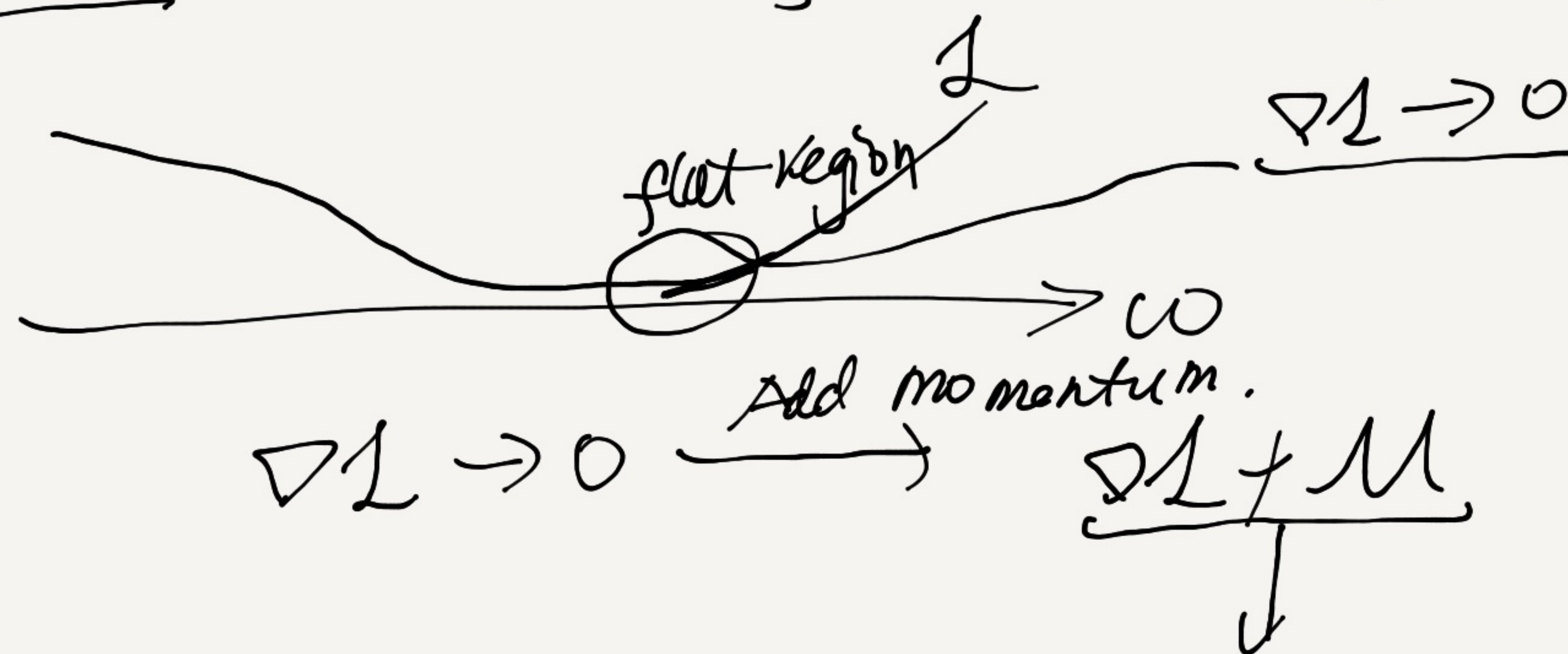
① use different activation function

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

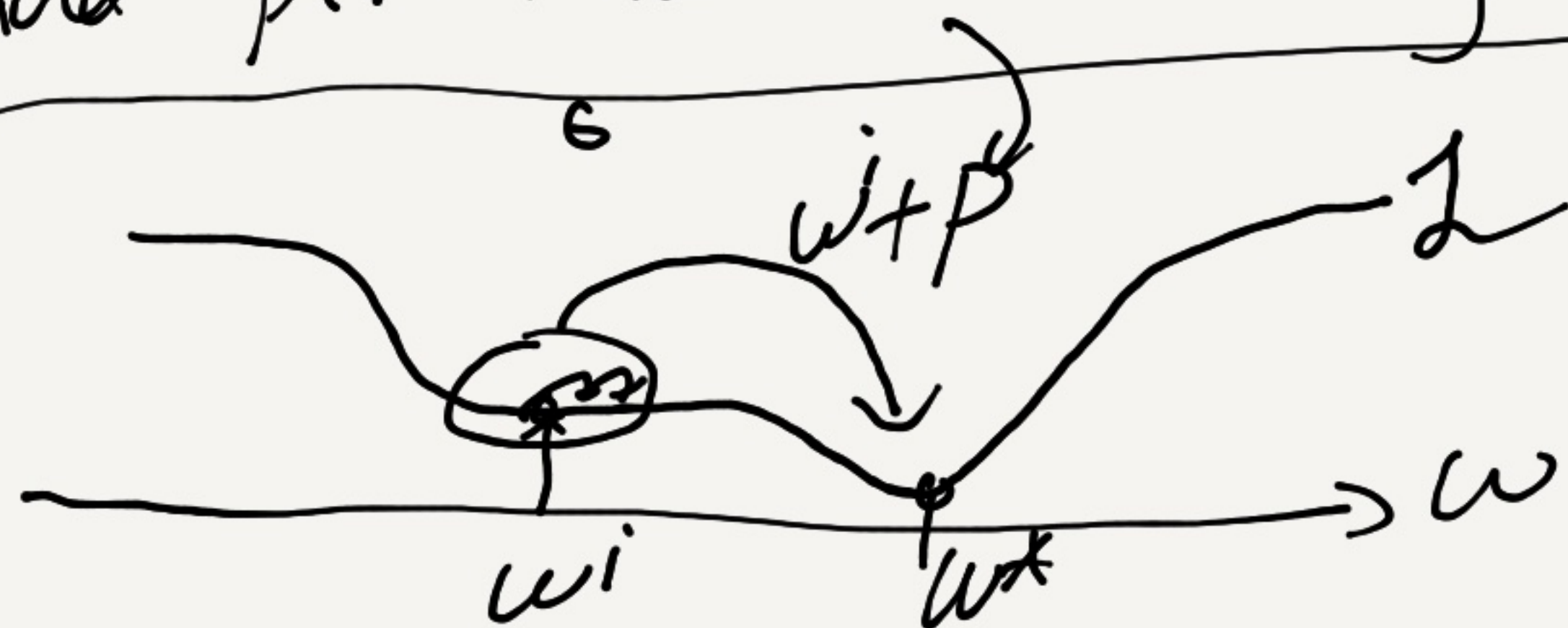
$$\nabla \text{ReLU}(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$



② Add momentum term during the training



③ Add perturbations during training



$$\omega^{i+1} = \omega^i - (\underbrace{\epsilon \cdot \nabla_{\omega^i} L}_0 + \mu)$$

2. Multiple factors control the learning process

$$w^{i+1} = w^i - \underbrace{\epsilon \cdot \underbrace{\nabla_{w^i} L}}_{\text{gradient}}$$

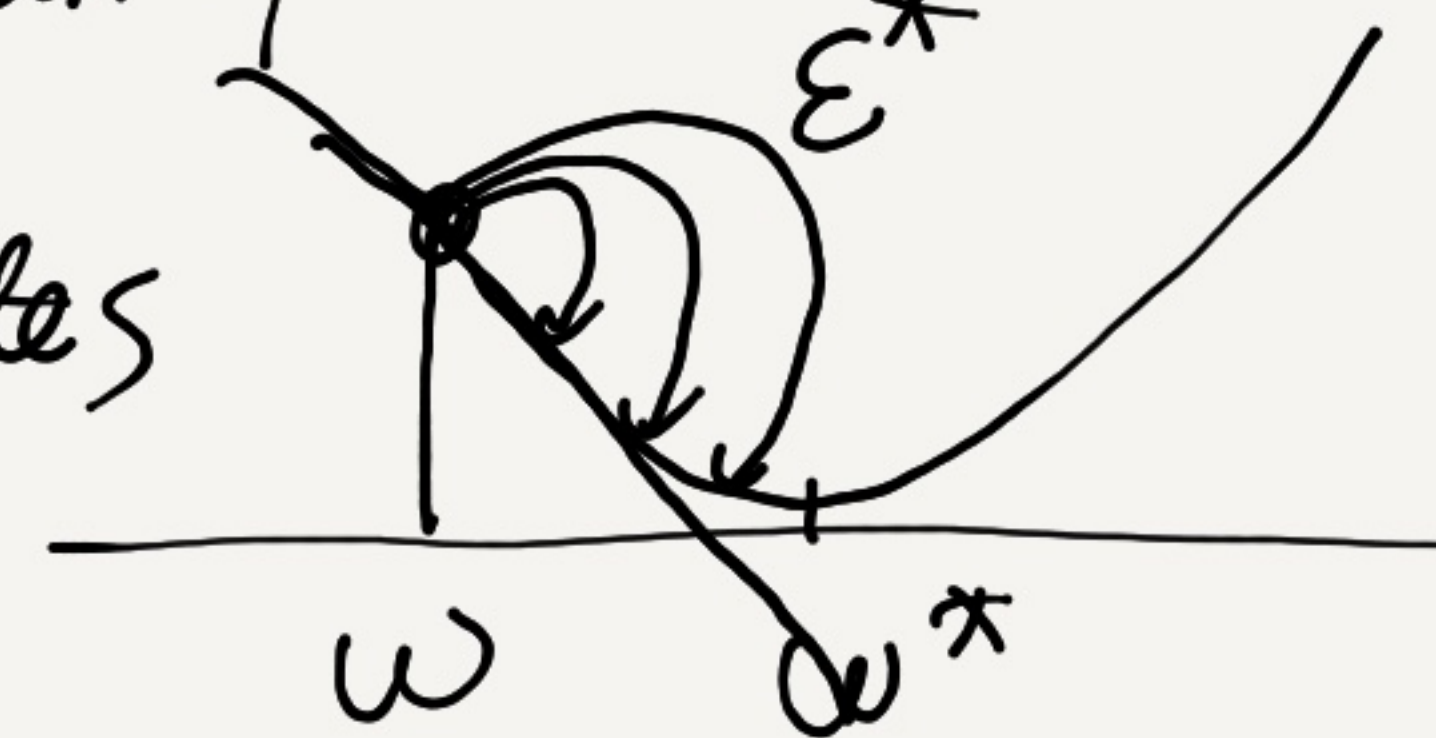
1) learning rate : $\begin{cases} \text{too small: training converges slowly.} \\ \text{too large: not converge.} \end{cases}$

① Use another algorithm to identify a 'good' learning rate.

• line search. \rightarrow search an optimal ϵ^*

• Try different learning rates

and choose the best,



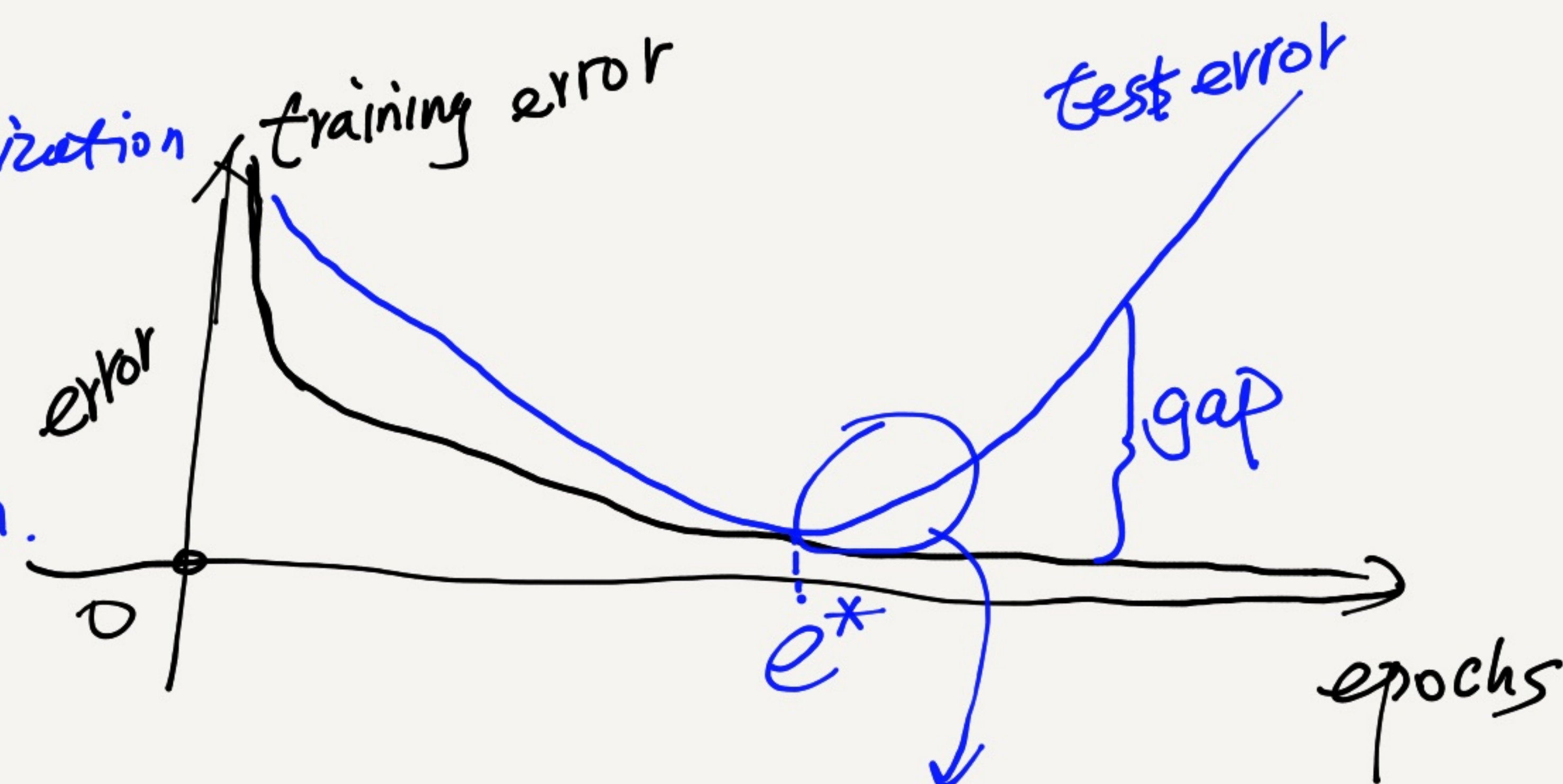
② Determine the learning rate adaptively. \rightarrow Adam

2) ∇L (avoid the vanishing gradient issue)

3. Early stopping and evaluation.

Excessive training leads to poor generalization performance (test perform.)

① stop training at the right epoch.



If the test/validation performance start to decrease, we get the signal to stop the training.

→ If we involve test set into the training (test curve), we will leak our test data (model observed the test data during training)

↓ Solve this issue by splitting the whole data sets into

3 subsets: training, validation, test sets

