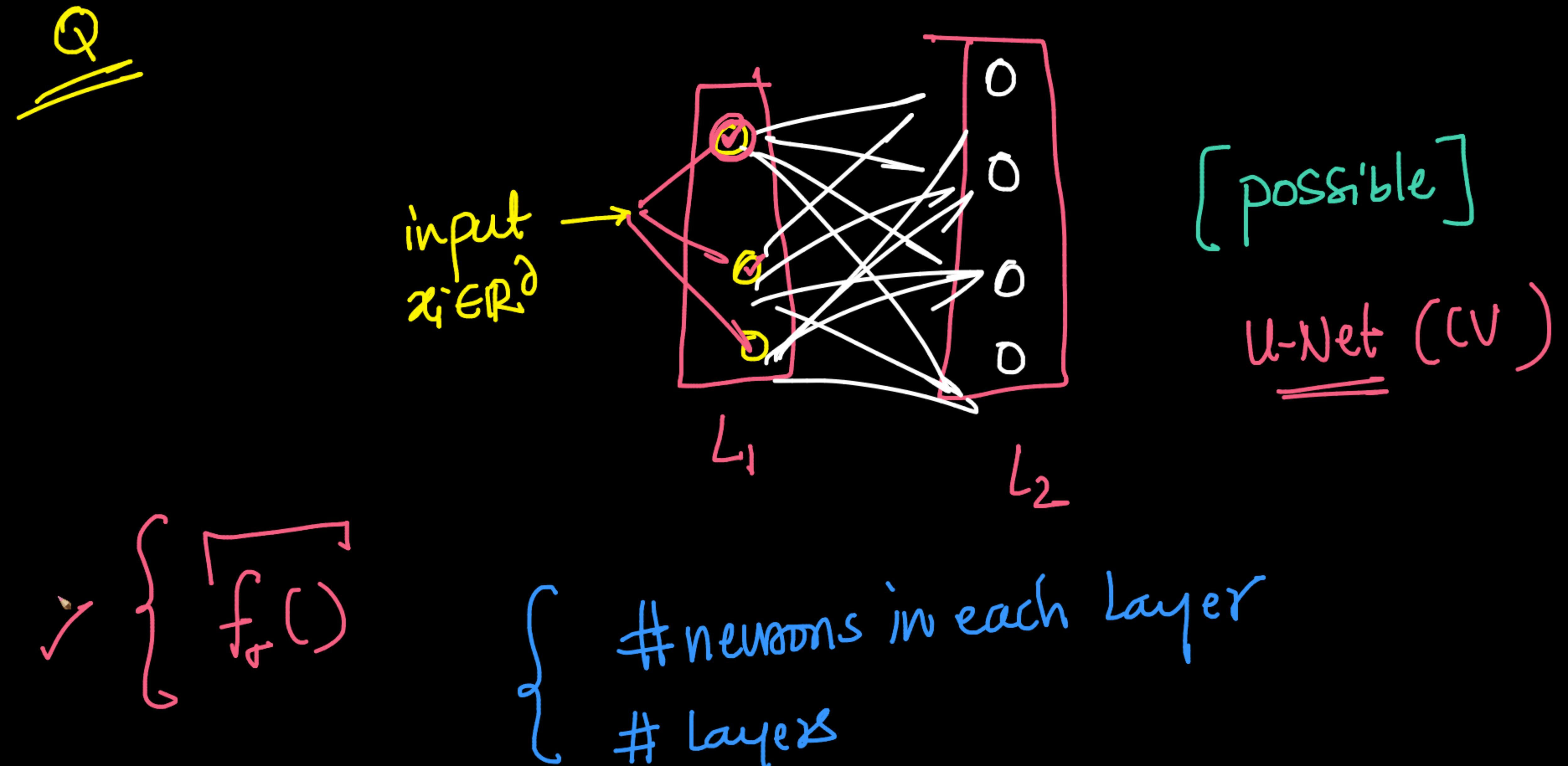


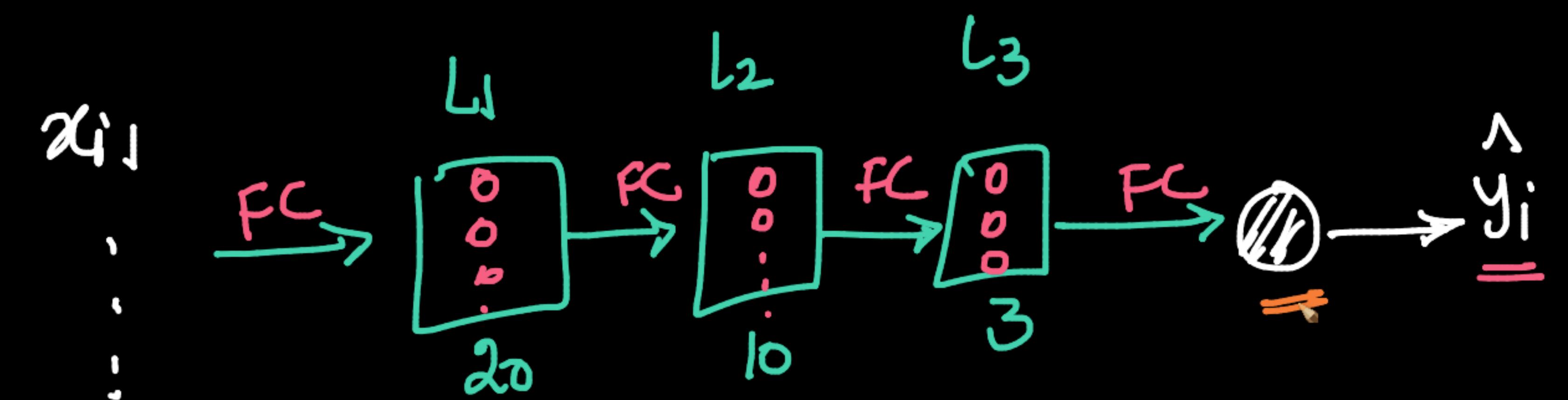
Understand + Code ]  
Back-propagation → 1980's (Hinton & others)  
[ ✓ Back-propagation  
(How MLPs are trained) → Most of DL-NN

{ → CV  
→ NLP / text ; speech ...



$\left\{ \begin{array}{l} f_{11} : f_\sigma \\ f_{12} : \text{anolkin elvin} \end{array} \right.$

$\Rightarrow$  tons of More  
hyper-param



$x_{i,1} \dots x_{i,n}$

# Recap:

① MLP ≈ function composition

② Activation fns have to be non-linear

③ Logistic reg & Lr-SVM → single Neuron model

④ Brute force : Compute  $\frac{\partial \text{Loss}}{\partial w_j}$  + weight<sub>j</sub> + GD

Recap:

✓ ① MLP ≈ function composition  $f \circ g; g \circ f$

✓ ② Activation fns have to be non-linear

③ Logistic reg & Lr-SVM → single Neuron model

{ ④ Brute force : Compute  $\frac{\partial \text{Loss}}{\partial w_{ij}^k}$  + weight<sub>ij</sub>  
=  $w_{ij}^k$

+ GD

Q

why should all activ fun be non-linear?

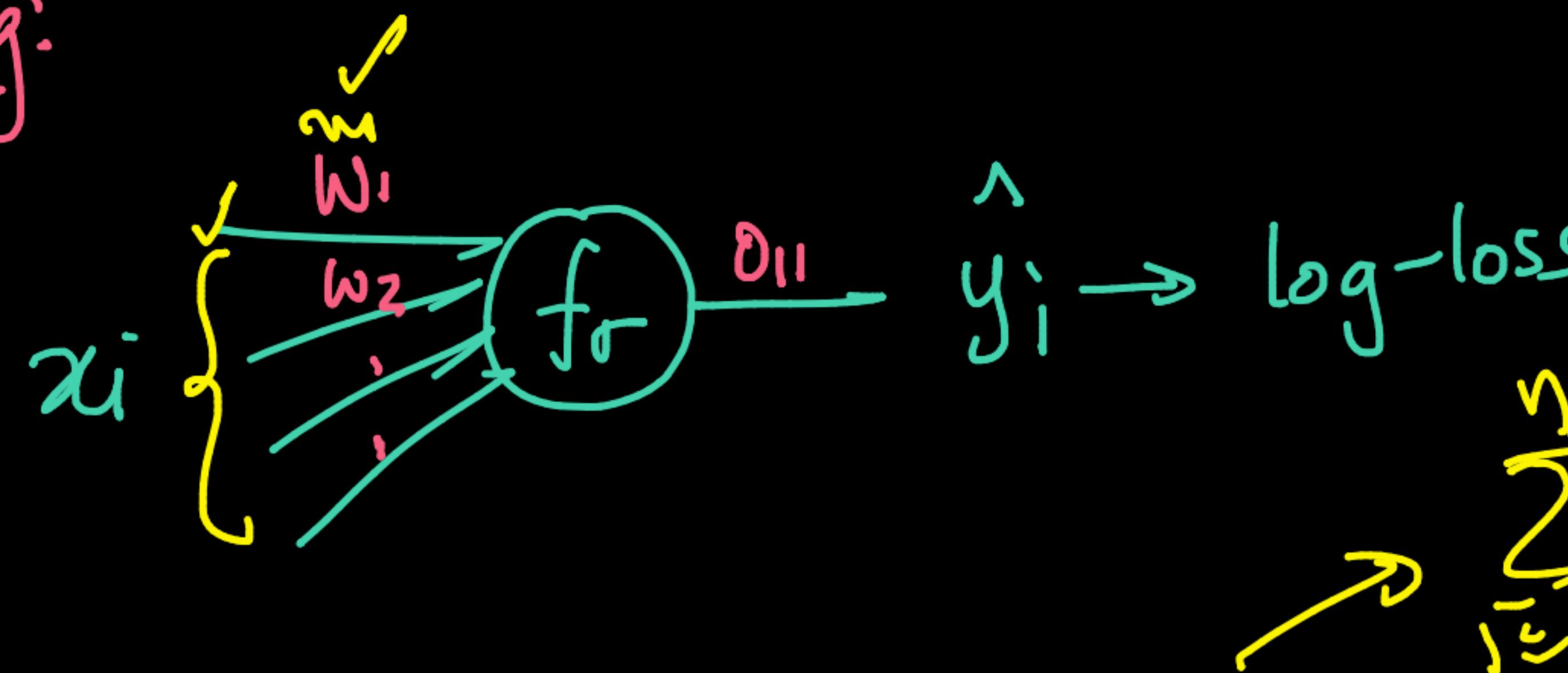
$$d' \left( d \left( \underline{\omega^T x_i + b} \right) + \beta \right) + \beta'$$

$$x_i \xrightarrow{f_c} 0 \xrightarrow{f_c} 0$$

$$\begin{matrix} \vdots & \vdots \\ \vdots & \vdots \end{matrix}$$

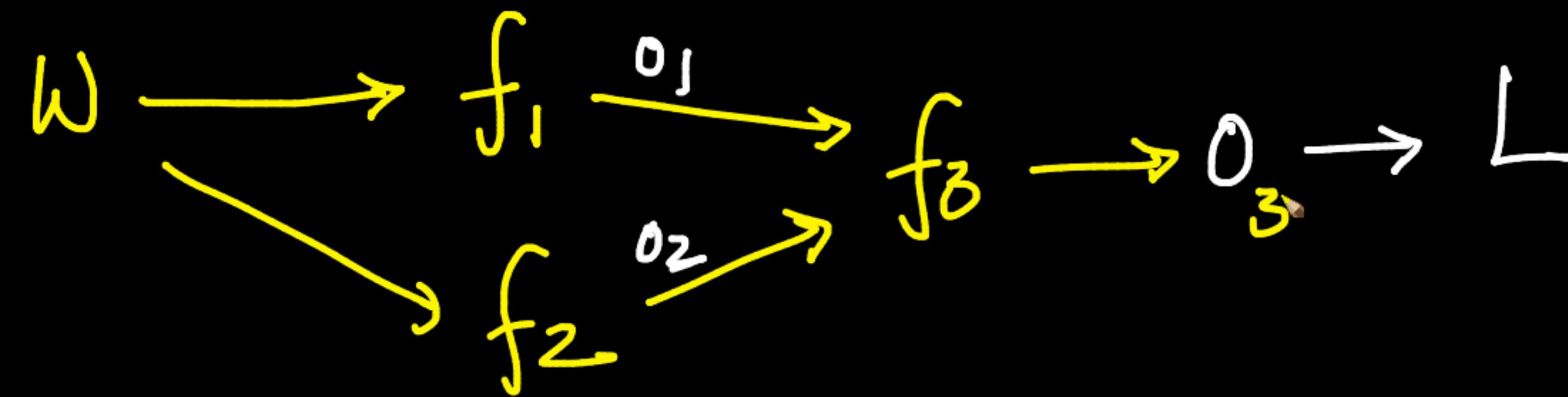
$$f \circ g(x) = f_L(g_L(x)) = h_L(x)$$

logistic reg:



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_1}$$

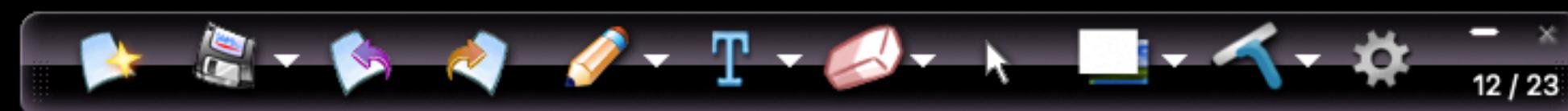
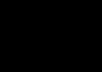
e.g.:



$$\begin{aligned}\text{Loss} &= L(y, f_3(o_1, o_2)) \\ &= L(y, f_3(f_1(w), \underline{f_2(w)})\end{aligned}$$

$$\text{Loss} = L \left( y, \underline{\underline{f_3}} \left( \underline{\underline{f_1(\omega)}}, \underline{\underline{f_2(\omega)}} \right) \right)$$

$\cdot o_3$



obj:

$$\frac{\partial \underline{L}}{\partial \underline{w}_{11}} = ?$$

$$\frac{\partial \underline{w}_{11}^1}{\partial \underline{w}_{11}} =$$

$$\frac{\partial \underline{L}}{\partial \underline{w}_{11}^1} = \frac{\partial \underline{L}}{\partial \underline{o}_{31}} \cdot \frac{\partial \underline{o}_{31}}{\partial \underline{w}_{11}}$$

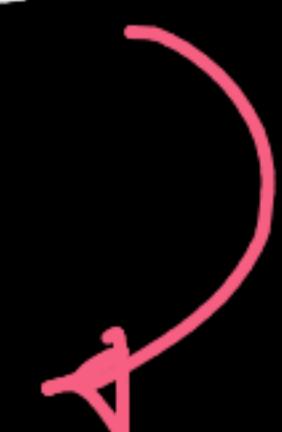
$$\frac{\partial \underline{o}_{21}}{\partial \underline{w}_{11}^2}$$



chain-rule!

$$\left\{ \frac{\partial L}{\partial \dot{O}_{31}} \cdot \frac{\partial \dot{O}_{31}}{\partial \dot{O}_{21}} \cdot \frac{\partial \dot{O}_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}} \right\} \checkmark$$

One-path  
not



④  $\nabla_{ij}^k$

$$w_{ij}^k \text{ new} = w_{ij}^k - \eta \frac{\partial L}{\partial w_{ij}^k}$$

$$b_{ij}^k \text{ new} = b_{ij}^k - \eta \frac{\partial L}{\partial b_{ij}^k}$$

$\hookrightarrow$  back-prop

$n = 100 \text{ Million}$   
 $d = 1000\text{-dim}$

GD:

$$\frac{\partial L}{\partial w_{ij}^k}$$

using all  $\overset{n}{\underset{^}{\text{points}}}$

SGD:

$$\frac{\partial L}{\partial w_{ij}^k}$$

using one random pt

✓ batch-GD :

$$\frac{\partial L}{\partial w_{ij}^k}$$

using a random  $A \in \mathbb{R}^{m \times k}$

Challenge:  $\frac{\partial L}{\partial w_{ij}^k}$

$$\frac{\partial L}{\partial w_{ii}^l} =$$

Challenge:  $\frac{\partial L}{\partial w_{ij}^k}$

$$\frac{\partial L}{\partial w_{ii}^l} =$$

④  $\nabla_{ij}^k$

$$w_{ij}^k \text{ new} = w_{ij}^k - \eta \frac{\partial L}{\partial w_{ij}^k}$$

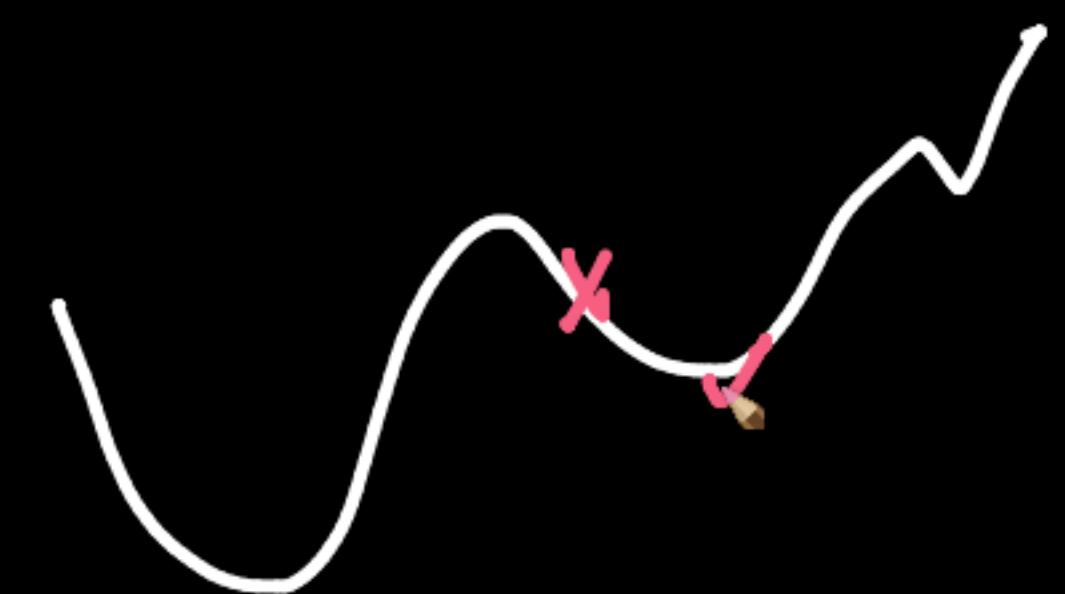
$$b_{ij}^k \text{ new} = b_{ij}^k - \eta \frac{\partial L}{\partial b_{ij}^k}$$

$\hookrightarrow$  back-prop

⑤

Repeat 2,3 & 4 till convergence

$$\left\{ \begin{array}{l} w_{ij}^k \text{ new } \approx w_{ij}^k \\ b_{ij}^k \text{ new } \approx b_{ij}^k \end{array} \right.$$

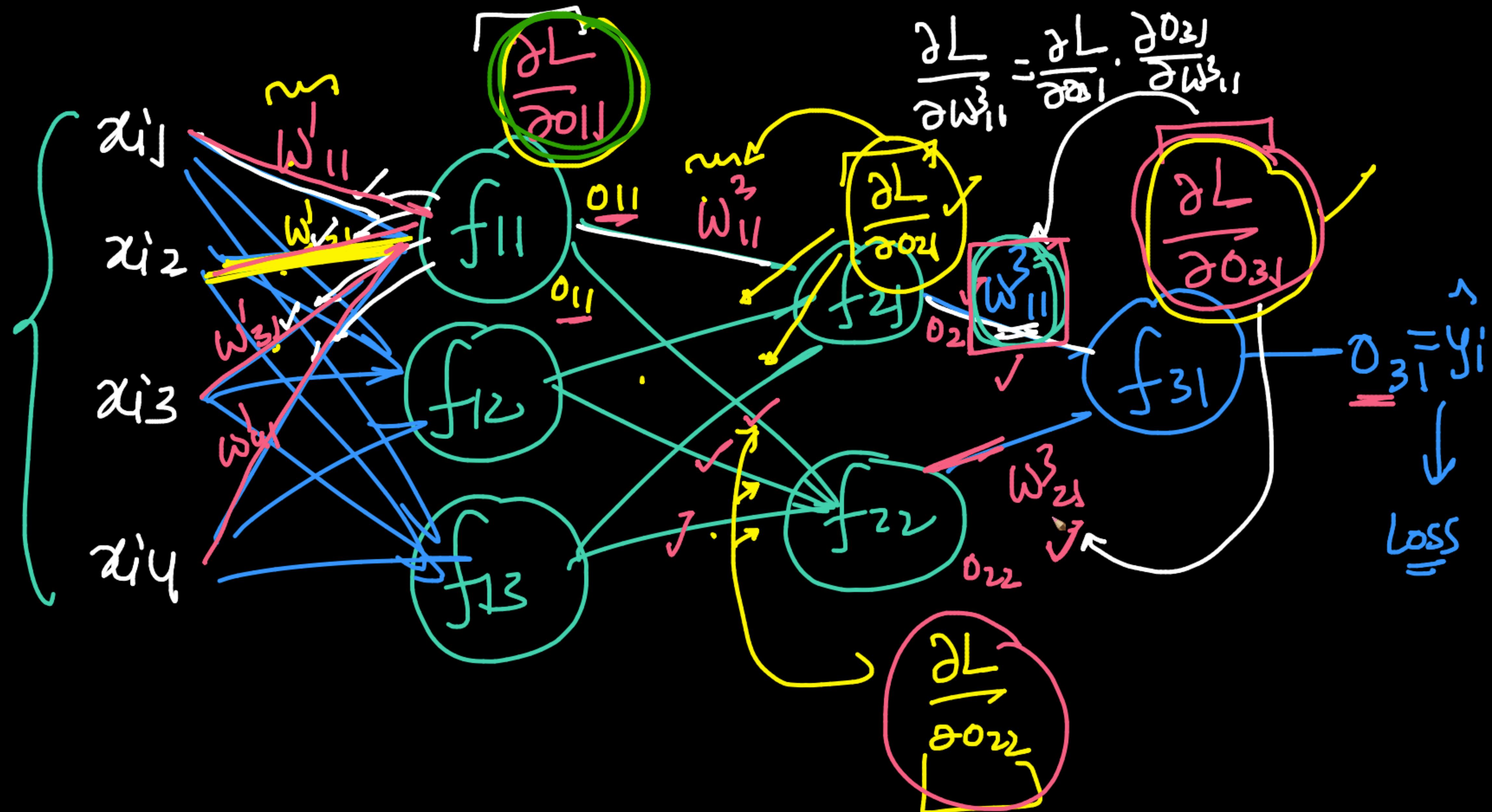


Challenge:  $\frac{\partial L}{\partial w_{ij}^k}$

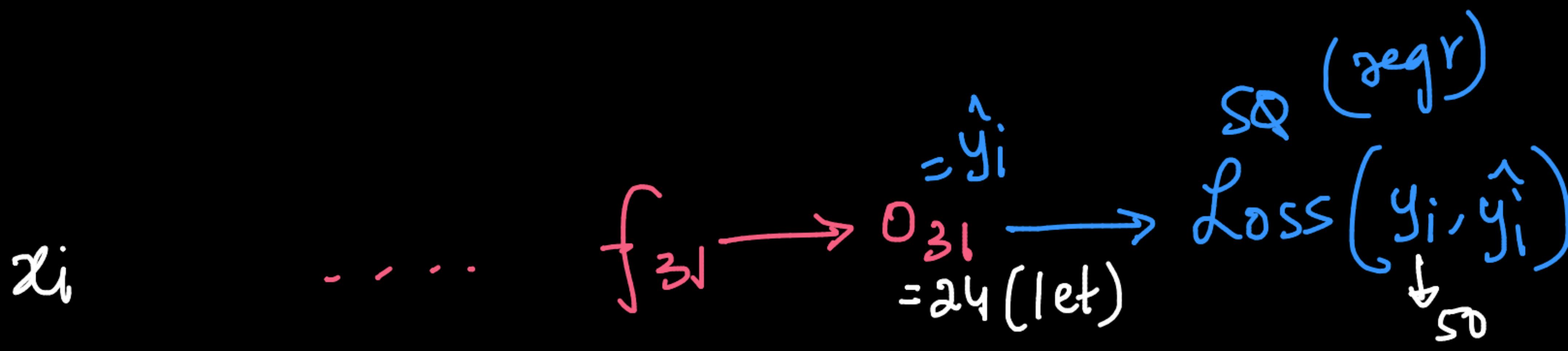
$\frac{\partial L}{\partial w_{ii}^l}$ : backwards (using chain rule)  
from loss till  
first layer

input —  $f$  — output

$$\frac{\partial \theta}{\partial i}$$



$$\left. \begin{array}{l} \omega'_{11} : \\ \omega'_{21} : \end{array} \right\} \quad \begin{aligned} \frac{\partial L}{\partial \omega'_{11}} &= \boxed{\frac{\partial L}{\partial \omega_{11}}} & \frac{\partial \omega_{11}}{\partial \omega'_{11}} \\ \frac{\partial L}{\partial \omega'_{21}} &= \boxed{\frac{\partial L}{\partial \omega_{11}}} & \frac{\partial \omega_{11}}{\partial \omega'_{21}} \end{aligned}$$



$$\frac{\partial L}{\partial o_{31}} = \frac{\partial \left( \underline{y_i - o_{31}} \right)^2}{\partial o_{31}}$$

$$= -2\underline{(y_i - o_{31})}$$

Backpropagation:

①

Chain rule to compute

$$\frac{\partial L}{\partial w_k}_{ij}$$

update

②

memoization:

[Dynamic programming]

Show the computational results if they are to be computed over & over again



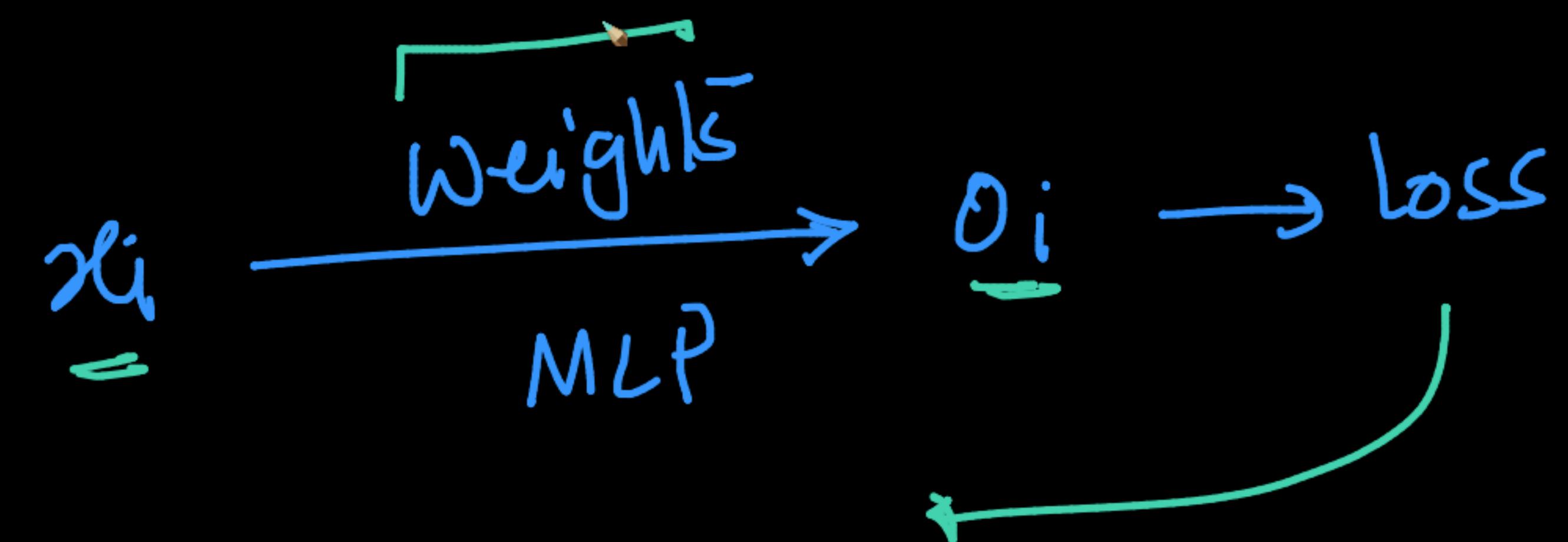


(c.g)

$$\theta = w_1x_1 + \overbrace{w_2x_2}$$

$$\frac{\partial \theta}{\partial w_1} = x_1$$

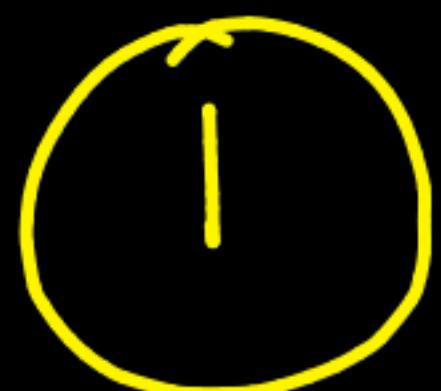
$$\frac{\partial \theta}{\partial w_2} = x_2$$



back-prop

Given:  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Task:  $\omega_{ij}^k$



initialize  $\omega_{ij}^k$  &  $n \rightarrow$  learning rate  
(randomly)

2 Forward Prop (batch-SGD)

3 loss

4 Update weights  $\xrightarrow{\text{backwards}}$  (op input)

↳ Keep them in memory & reuse

$\frac{\partial L}{\partial o_{ij}}$  → chain rule

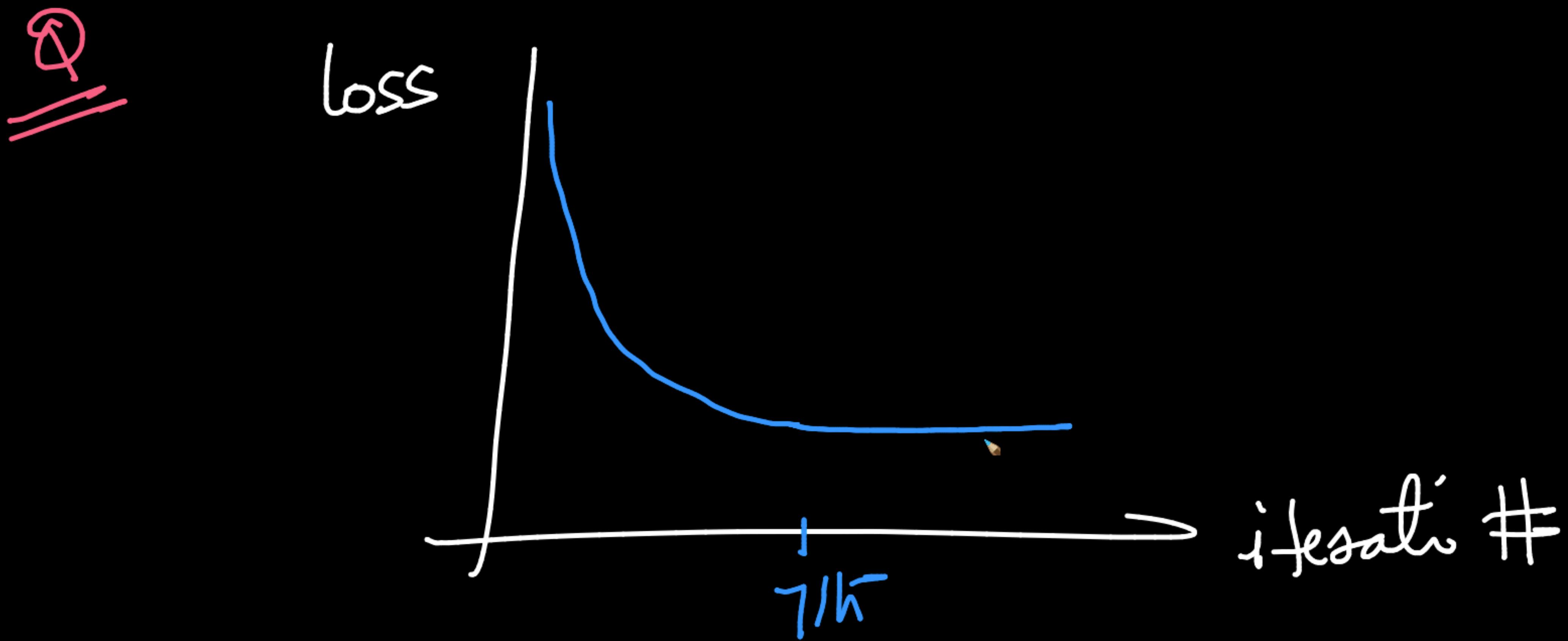
⑤ Repeat 2,3,4 till convergence

1

activation fns & loss fns  
should be differentiable

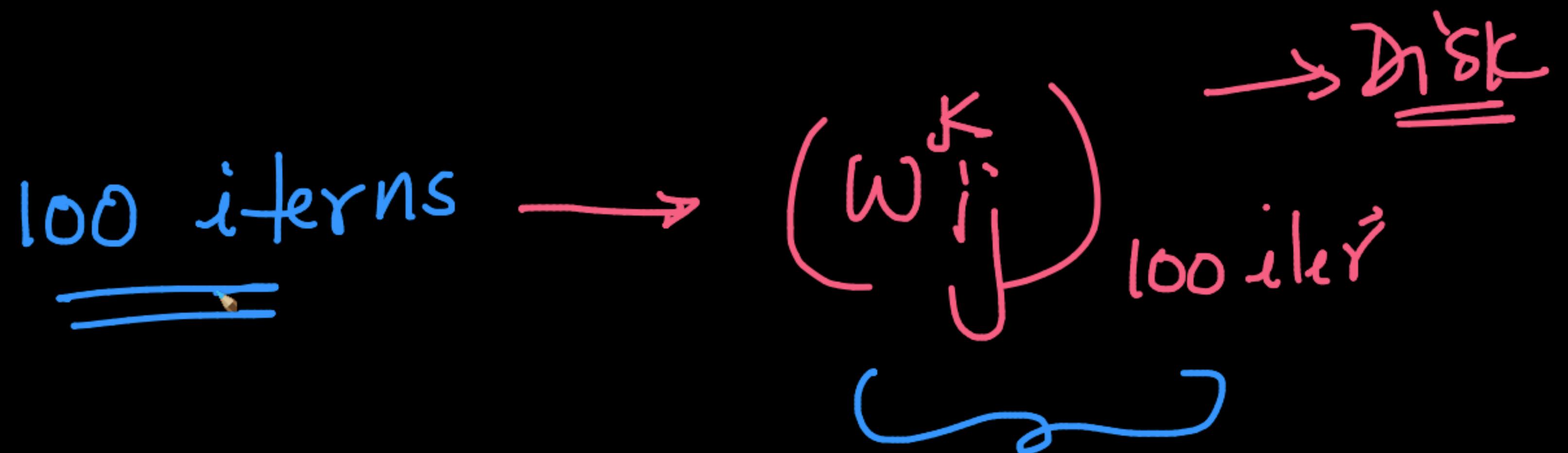
2

Slow ← simple back-prop  
(later)





mini-batch  $\stackrel{\text{SGD}}{=}$   
↑



Actual:req: non-linear fu:

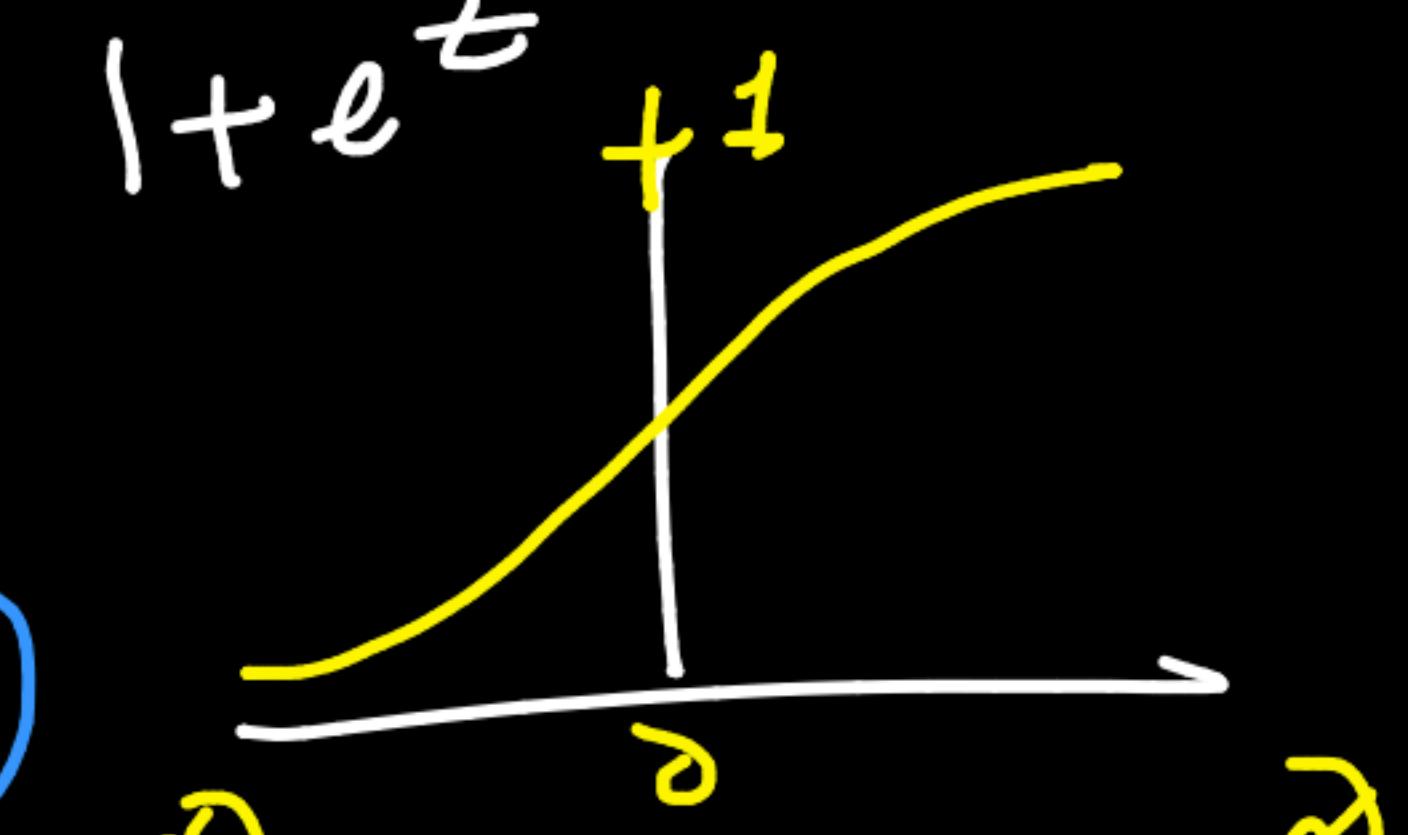
$$\omega^T x + b = z$$

1980's & 90's:

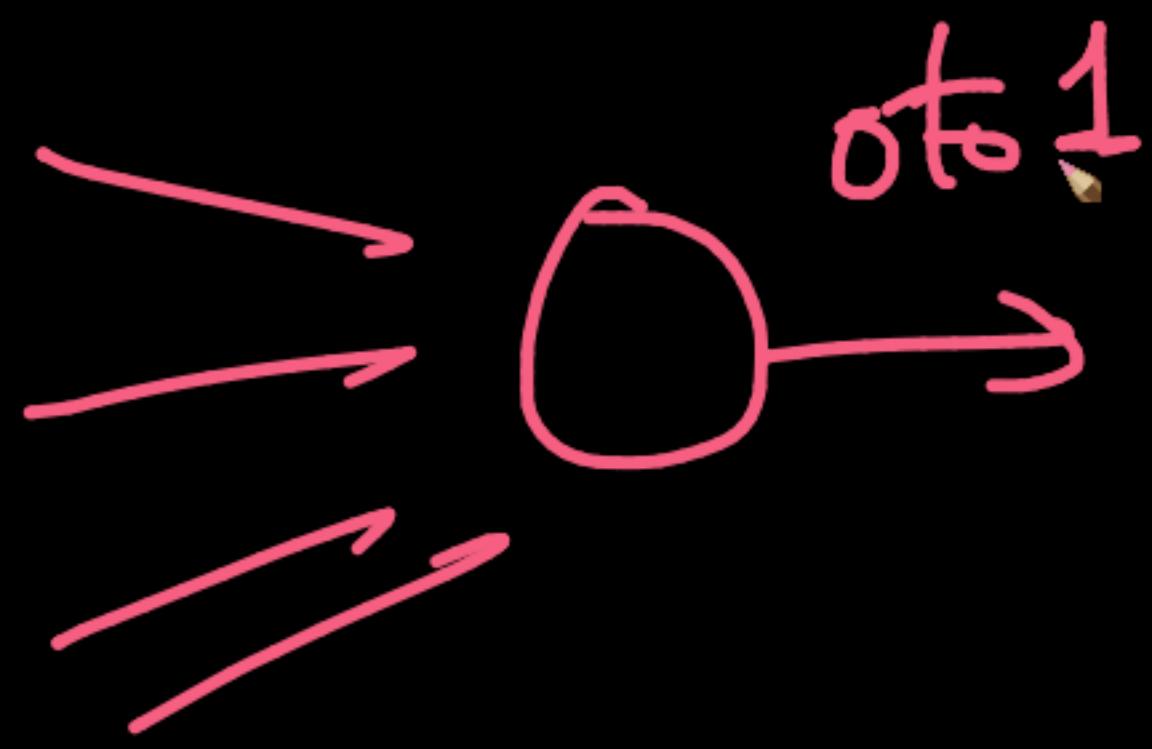
Sigmoid

logistic segn..

$$\sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$$



$$\sigma'(z) = \frac{\partial \sigma}{\partial z} = \underline{\sigma(z)(1-\sigma(z))}$$





Arunava Chakraborty | LinkedIn

View Arunava Chakraborty's profile  
LinkedIn is the world's largest professional network.

41 / 41

Valleywood AI in Nerd For Tech  
What is a Machine Learning Model?

Janko

BackpropFromScratch.ipynb - X | Derivative of the Softmax Func X | derivative of sigmoid - Google X | Derivative of the Sigmoid funct X | tanh and derivative - Google S X | Activation Functions with Deriv X +

medium.com/@omkar.nallagoni/activation-functions-with-derivative-and-python-code-sigmoid-vs-tanh-vs-relu-44d23915c1f4

Get started Sign In

Search

You will learn how to implement activation functions in Python. You will also learn how to calculate their derivatives.

Like the sigmoid function, the tanh function is also a logistic function. The main difference is that the tanh function is centered at zero, while the sigmoid function is centered at one. This means that the tanh function maps values between -1 and 1, while the sigmoid function maps values between 0 and 1.

Y<sub>i</sub> ∈ {0, 1}

Y<sub>i</sub> ∈ {-1, 1}

-∞

Tanh

Derivative of tanh(z):

a = (e<sup>z</sup> - e<sup>(-z)</sup>) / (e<sup>z</sup> + e<sup>(-z)</sup>)

use same up

Help Status Writers Blog Careers Privacy Terms About Knowable



Like the sigmoid function, one of the interesting properties of the tanh function is that the derivative can be expressed in terms of the function itself. Below is the actual formula for the tanh function along with the formula for calculating its derivative.

$$\left\{ \begin{array}{l} a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \\ \frac{da}{dz} = 1 - a^2 \end{array} \right.$$

$$= \frac{1 - e^{-z}}{1 + e^{-z}}$$

### Derivative of $\tanh(z)$ :

$$a = (e^z - e^{-z}) / (e^z + e^{-z})$$

use same u/v rule

$$da = [(e^z + e^{-z}) * d(e^z - e^{-z})] - [(e^z - e^{-z}) * d(e^z + e^{-z})] / [(e^z + e^{-z})]^2$$

$$da = [(e^z + e^{-z}) * (e^z + e^{-z})] - [(e^z - e^{-z}) * (e^z - e^{-z})] / [(e^z + e^{-z})]^2$$

$$da = [(e^z + e^{-z})]^2 - [(e^z - e^{-z})]^2 / [(e^z + e^{-z})]^2$$

$$da = 1 - [(e^z - e^{-z}) / (e^z + e^{-z})]^2$$

[Get started](#)
[Sign In](#)


**Nallagoni Omkar**

17 Followers

[Follow](#)


### More from Medium

Stanghong

[Explain Neural Net Model Using One Factor A Time Design \(OFAT\)](#)


Pratikbais

[ML | Linear Regression](#)


Mukesh Kumar in INSAID

[K-Fold Target Encoding for High Cardinality Features](#)


Moosa Ali in Geek Culture

[Boruta Feature Selection Explained in Python](#)




Like the sigmoid function, one of the interesting properties of the tanh function

is that the derivative can be expressed in terms of the function itself. Below is the actual formula for the tanh function along with the formula for calculating its derivative.

[Get started](#)[Sign In](#)
 Search


Nallagoni Omkar

17 Followers

[Follow](#)

## More from Medium

Stanghong

[Explain Neural Net Model Using One Factor A Time Design \(OFAT\)](#)

Pratikbais

[ML | Linear Regression](#)

Mukesh Kumar in INSAID

[K-Fold Target Encoding for High Cardinality Features](#)

Moosa Ali in Geek Culture

[Boruta Feature Selection Explained in Python](#)

$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{da}{dz} = 1 - a^2$$

Derivative of  $\tanh(z)$ :

$$a = (e^z - e^{-z}) / (e^z + e^{-z})$$

use same u/v rule

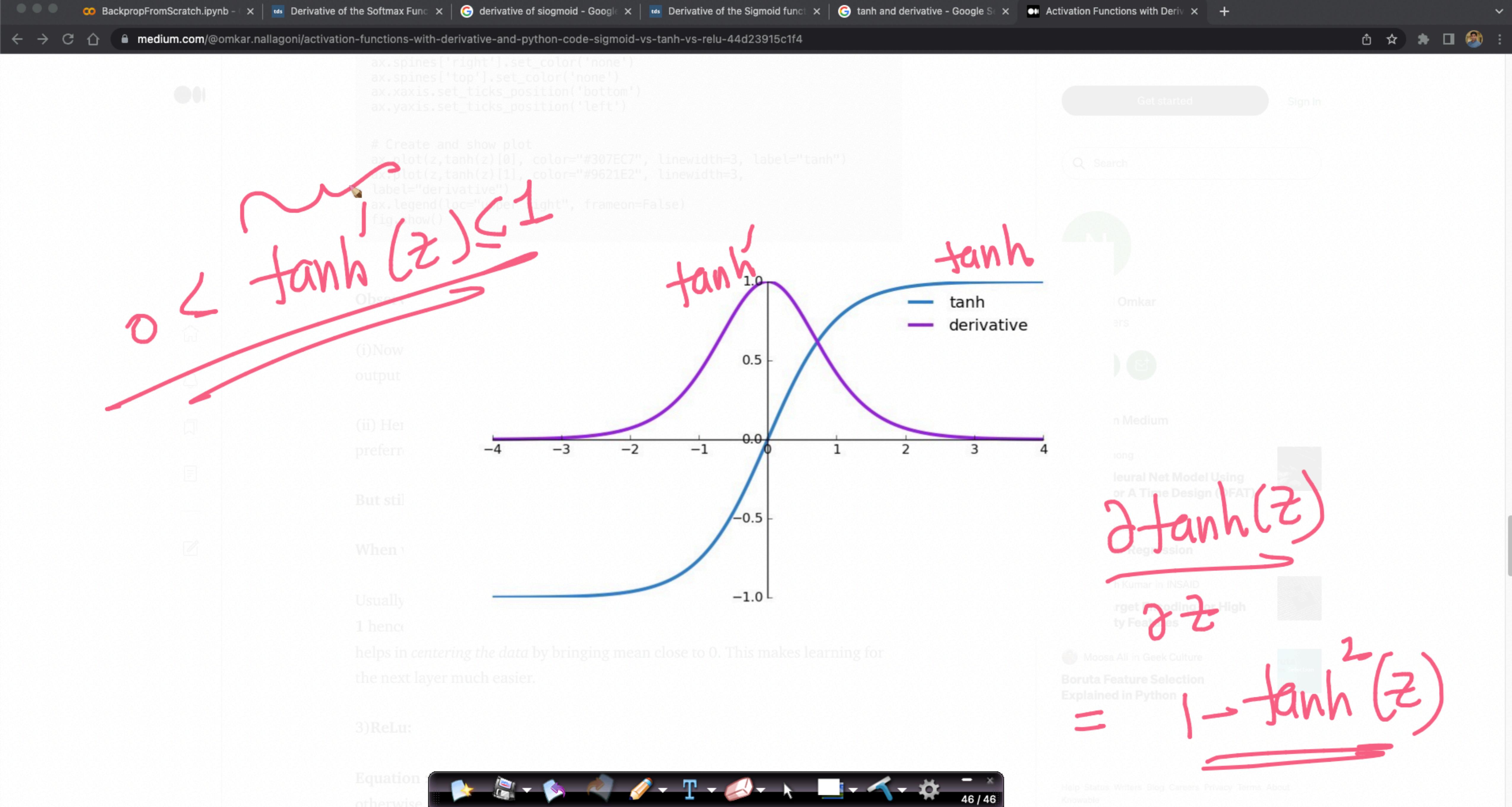
$$da = [(e^z + e^{-z}) * d(e^z - e^{-z})] - [(e^z - e^{-z}) * d(e^z + e^{-z})] / [(e^z + e^{-z})]^2$$

$$da = [(e^z + e^{-z}) * (e^z + e^{-z})] - [(e^z - e^{-z}) * (e^z - e^{-z})] / [(e^z + e^{-z})]^2$$

$$da = [(e^z + e^{-z})]^2 - [(e^z - e^{-z})]^2 / [(e^z + e^{-z})]^2$$

$$da = 1 - [(\dots)]^2$$

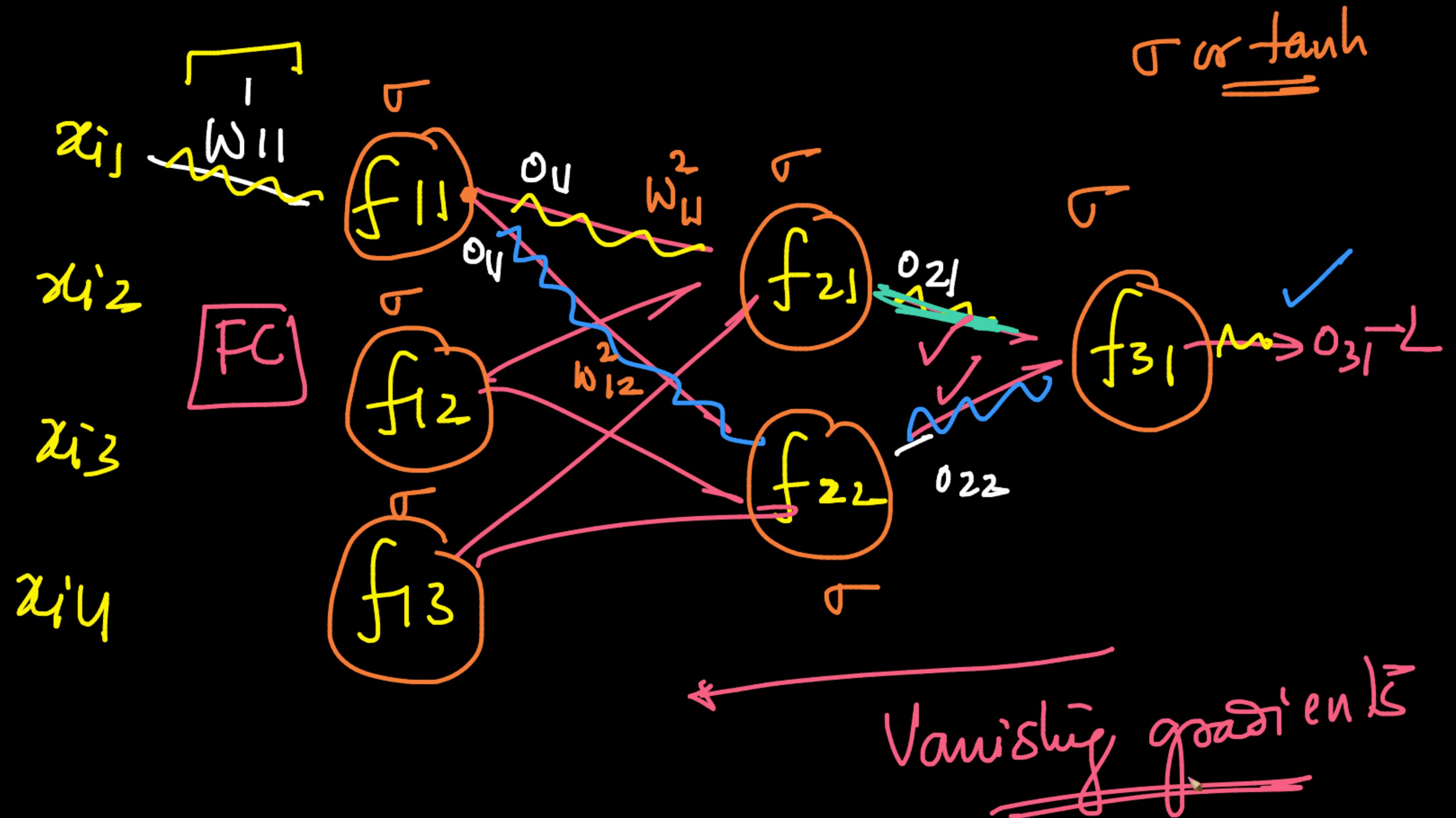




Traditionally:

Sigmoid:  $\sigma(z)$ ;  $0 < \sigma(z) < 1$

Tanh(z):  $-1 < \tanh(z) < 1$



$\sigma$  of tank

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial o_3} \cdot \left[ \frac{\frac{\partial o_3}{\partial o_2}}{\frac{\partial o_2}{\partial o_1}}, \frac{\frac{\partial o_2}{\partial o_1}}{\frac{\partial o_1}{\partial w'_{11}}}, \frac{\frac{\partial o_1}{\partial o_1}}{\frac{\partial o_1}{\partial w'_{11}}} \right] + \underbrace{o_1}_{\text{oto 1}}$$

$$\left[ \frac{\frac{\partial o_3}{\partial o_2}}{\frac{\partial o_2}{\partial o_1}}, \frac{\frac{\partial o_2}{\partial o_1}}{\frac{\partial o_1}{\partial w'_{11}}}, \frac{\frac{\partial o_1}{\partial o_1}}{\frac{\partial o_1}{\partial w'_{11}}} \right]$$

$\underbrace{o_2}_{\text{oto 1}} \quad \underbrace{o_1}_{\text{oto 1}} \quad \underbrace{o_1}_{\text{oto 1}}$

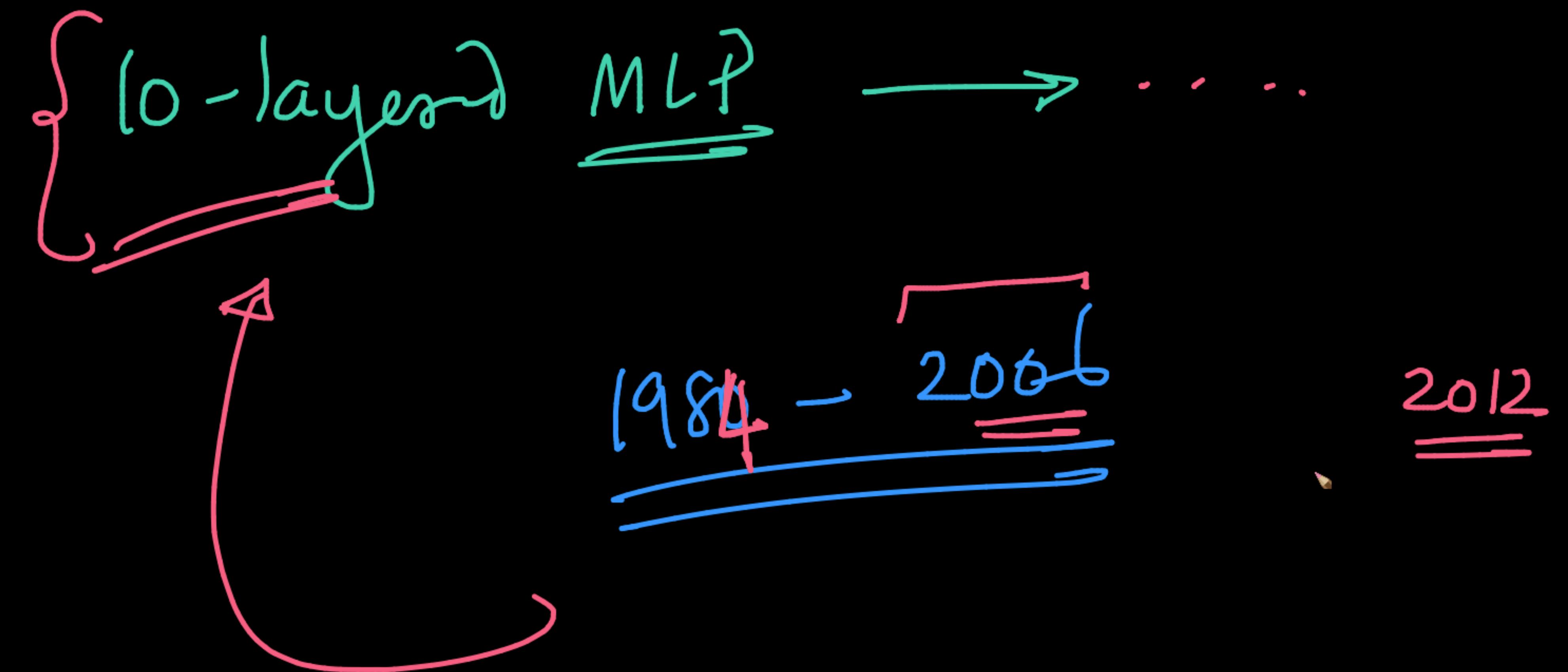
$$\omega_{||, \text{new}}^I = \omega_{||, \text{old}}^I - \eta \frac{\partial L}{\partial \omega_{||}^I}$$

update

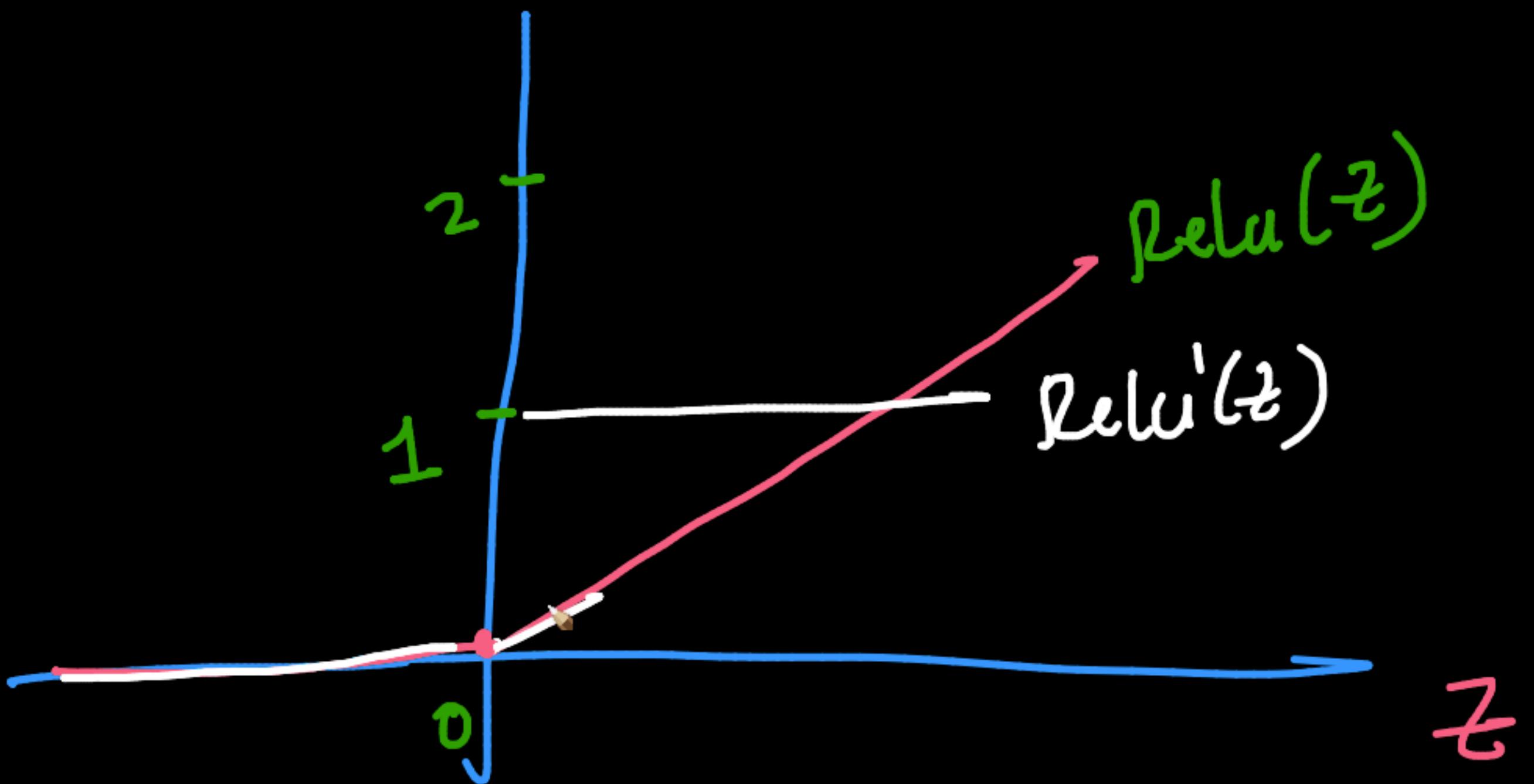
v.v.v. slowly  
because A

miniscale

vanning gradients



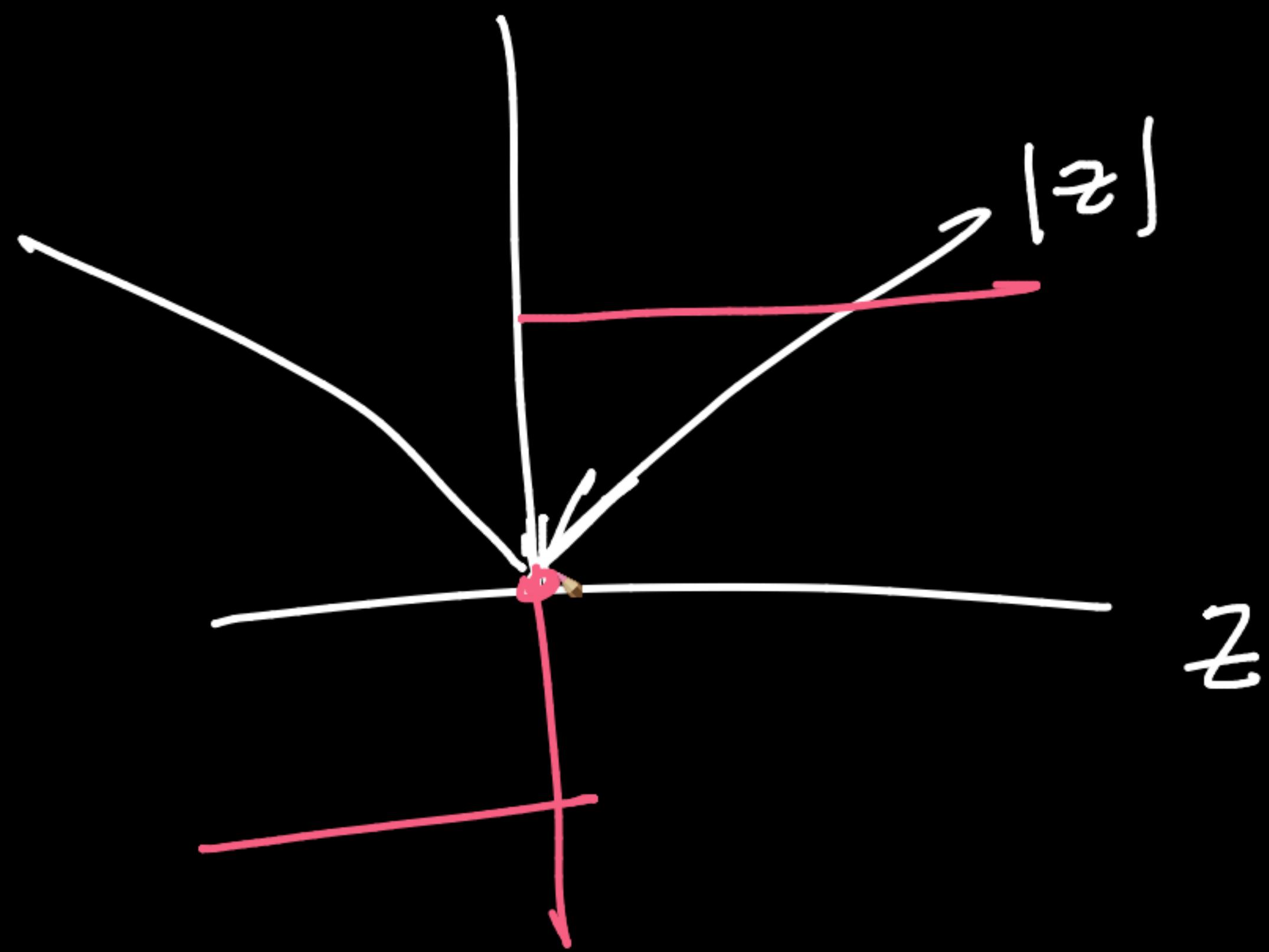
ReLU



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

$$\text{ReLU}'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

~~Simple & fast~~



Most  
widely  
used

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

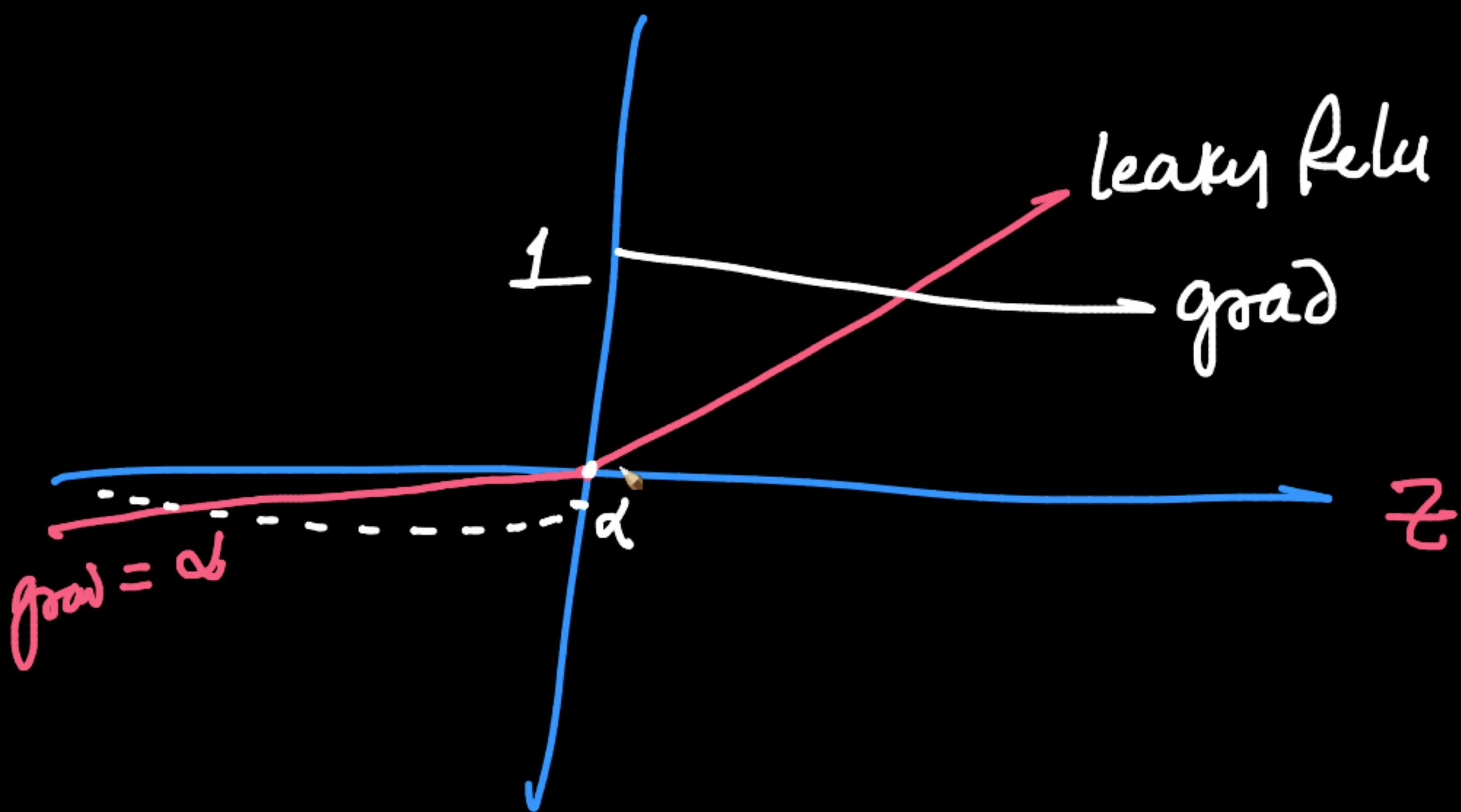
1x1x1 ...

no vanishing grad

1x1x1x0 ...

grad = 0  $\Rightarrow$  no update

leaky - Relu

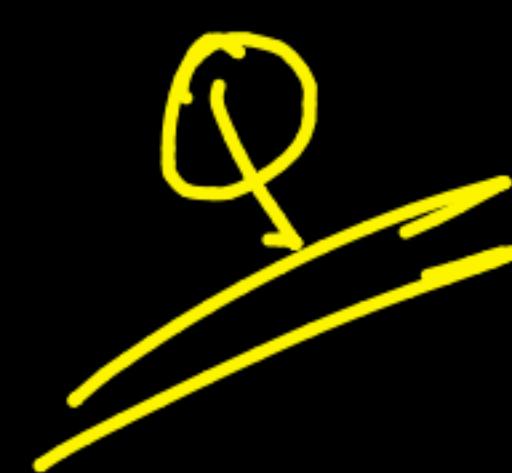


Next-class!

→ Softmax → K-classes

→ Code (Softmax, CE, Relu)

Backprop from scratch ...



Backprop: chain-rule + memoization

Multiplying lots of partial -der

tanh, sig mod, ReLU / leaky  
ReLU

0 to 1 → Vanishing  
gradient

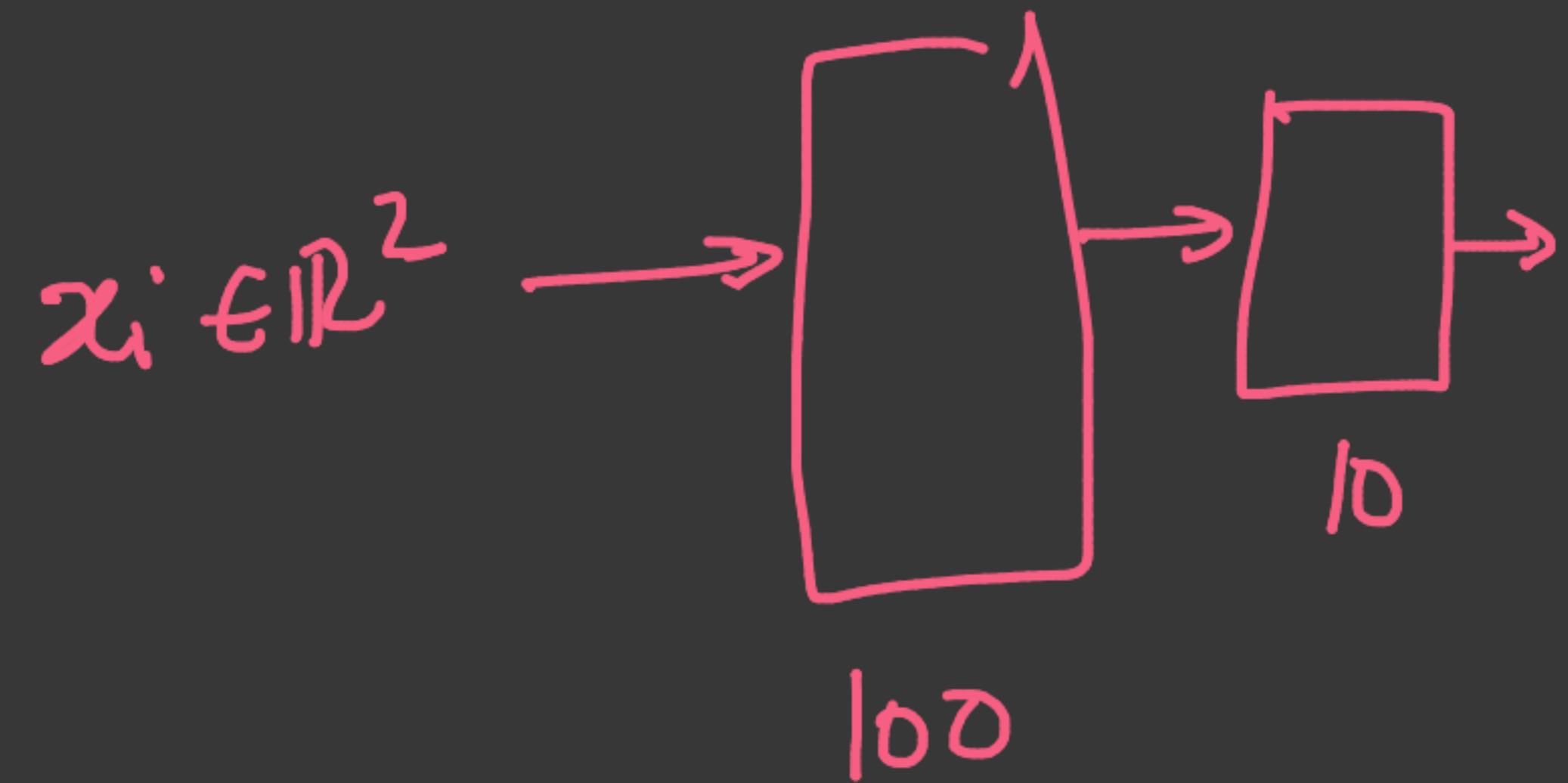
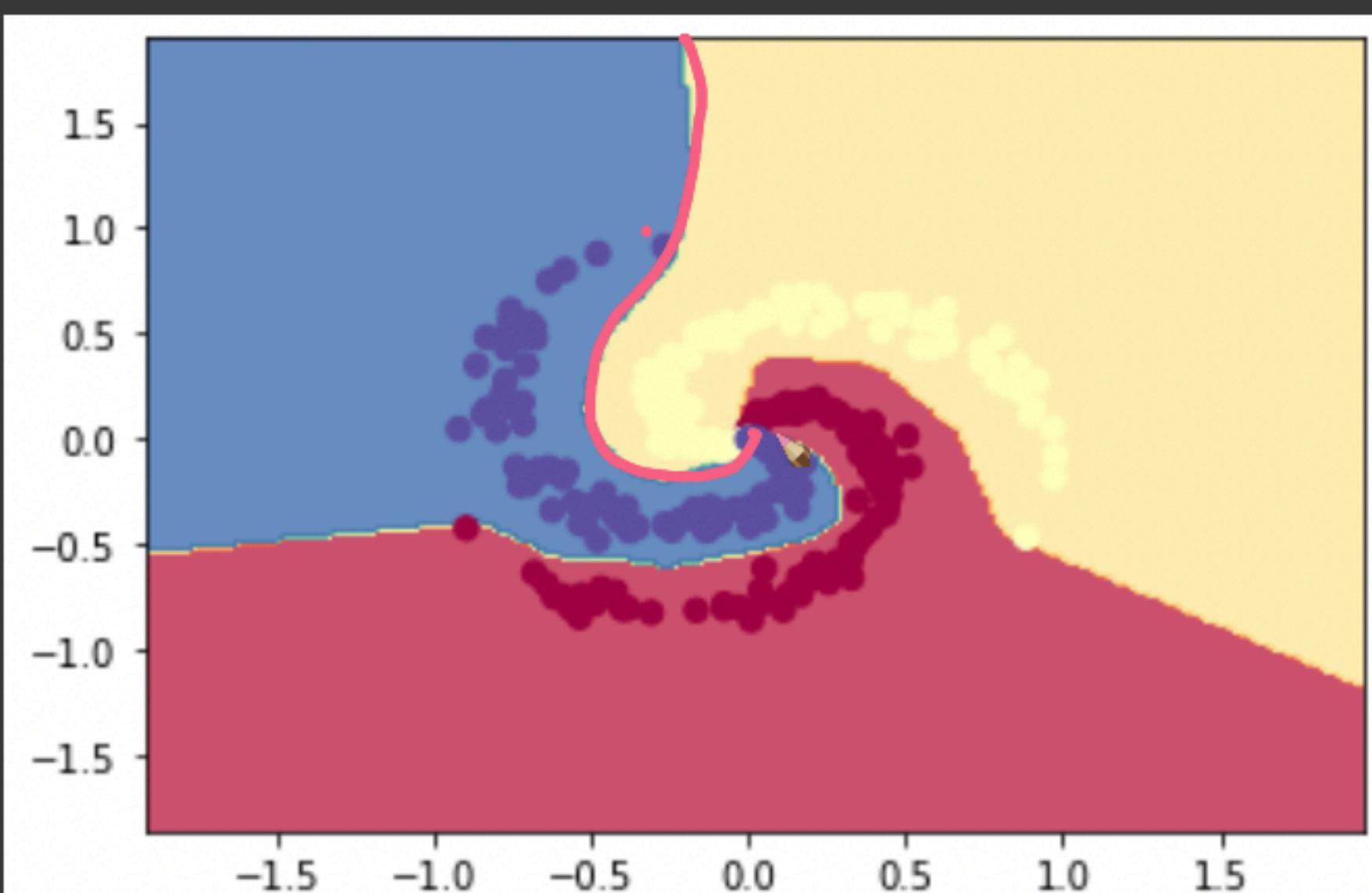
$> 1$  → exploding  
gradients

$$1.2 \times 1.3 \times 1.4 \times \dots$$

↓  
NLP(RNN)

Q  
—

—



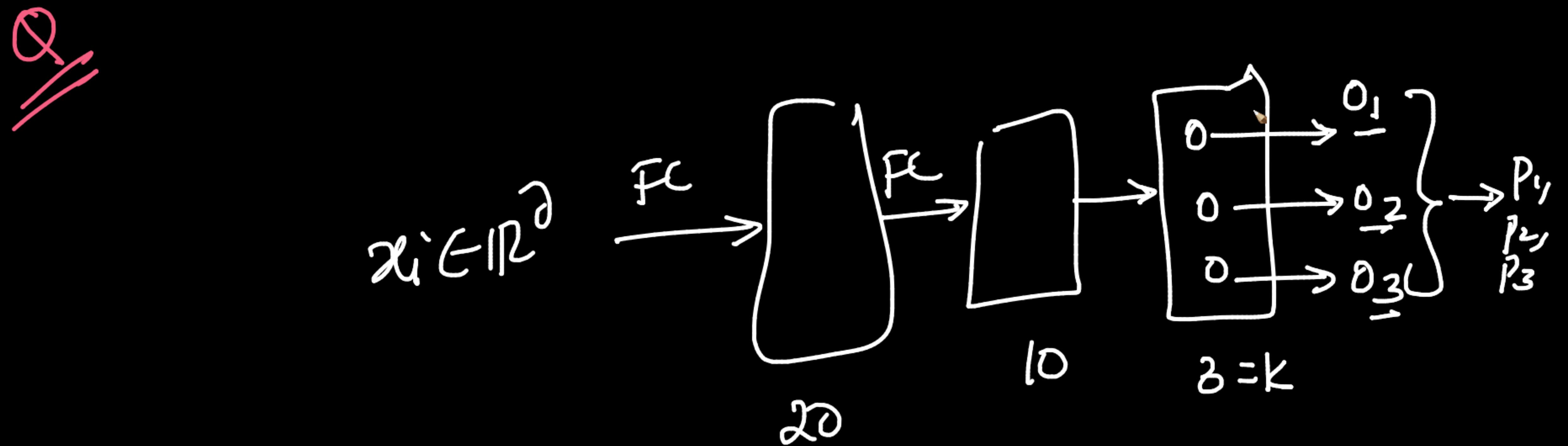
ReLU(....)

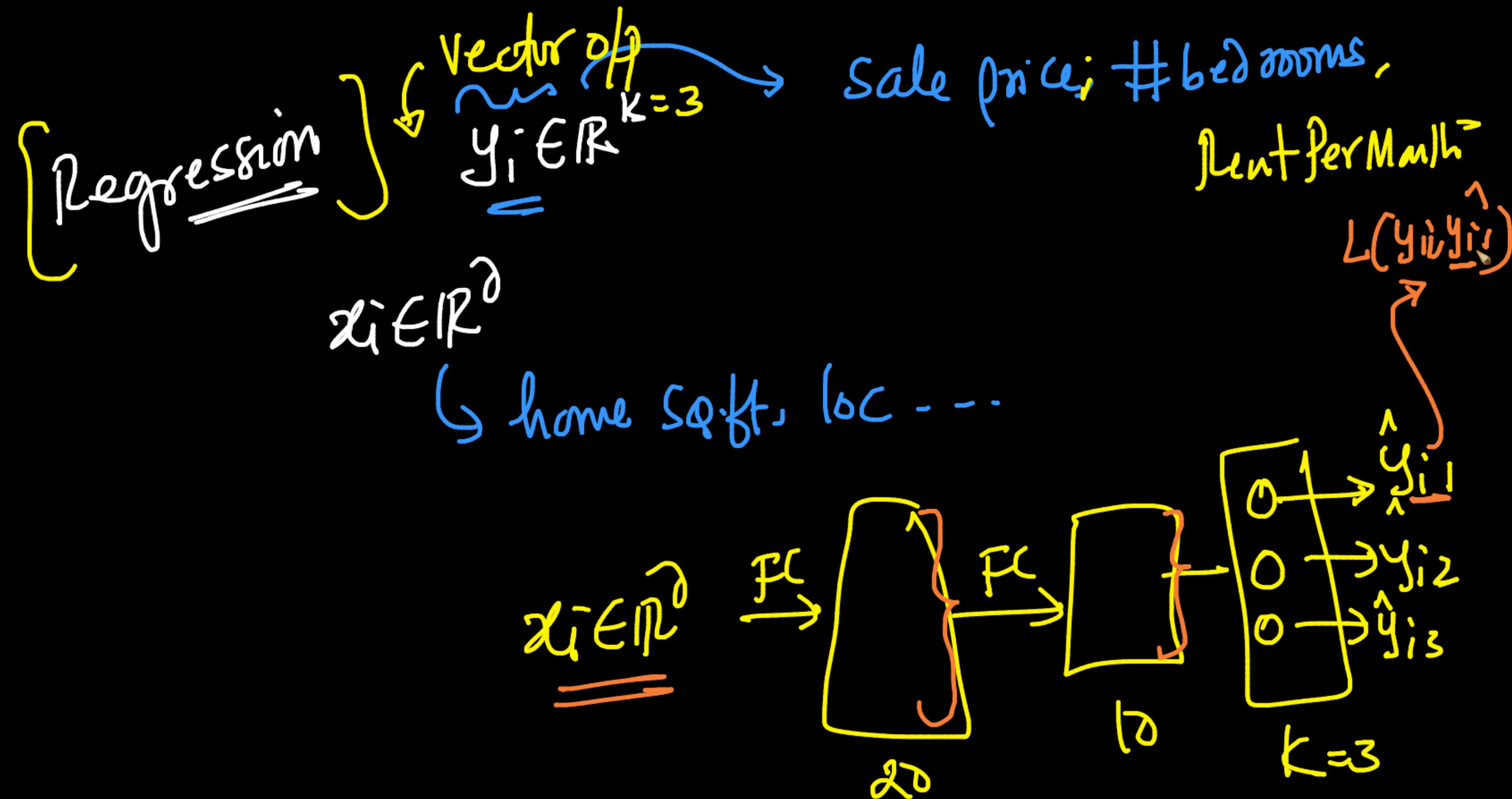
$\hat{=}$  Code  $\rightarrow$  next class

TF

2-layer ReLU

$\rho$   
 $\rho$   
 $\rho$   
 $\rho$   
 $\rho$

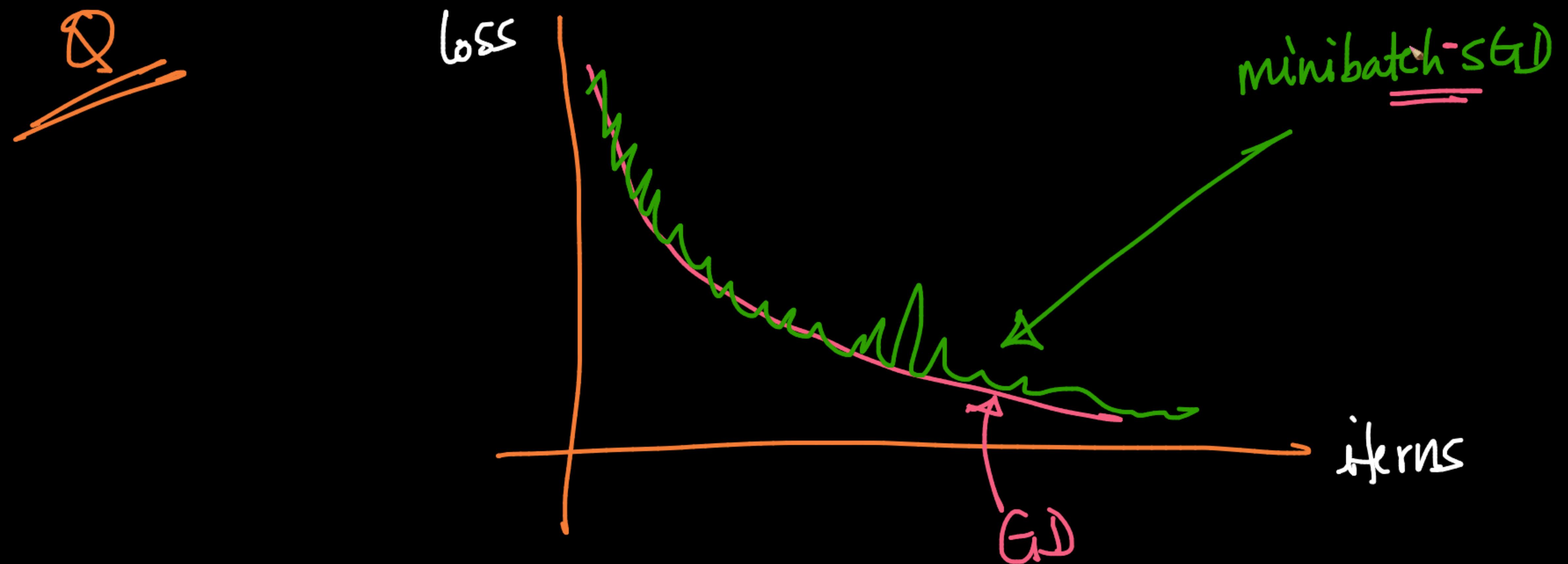


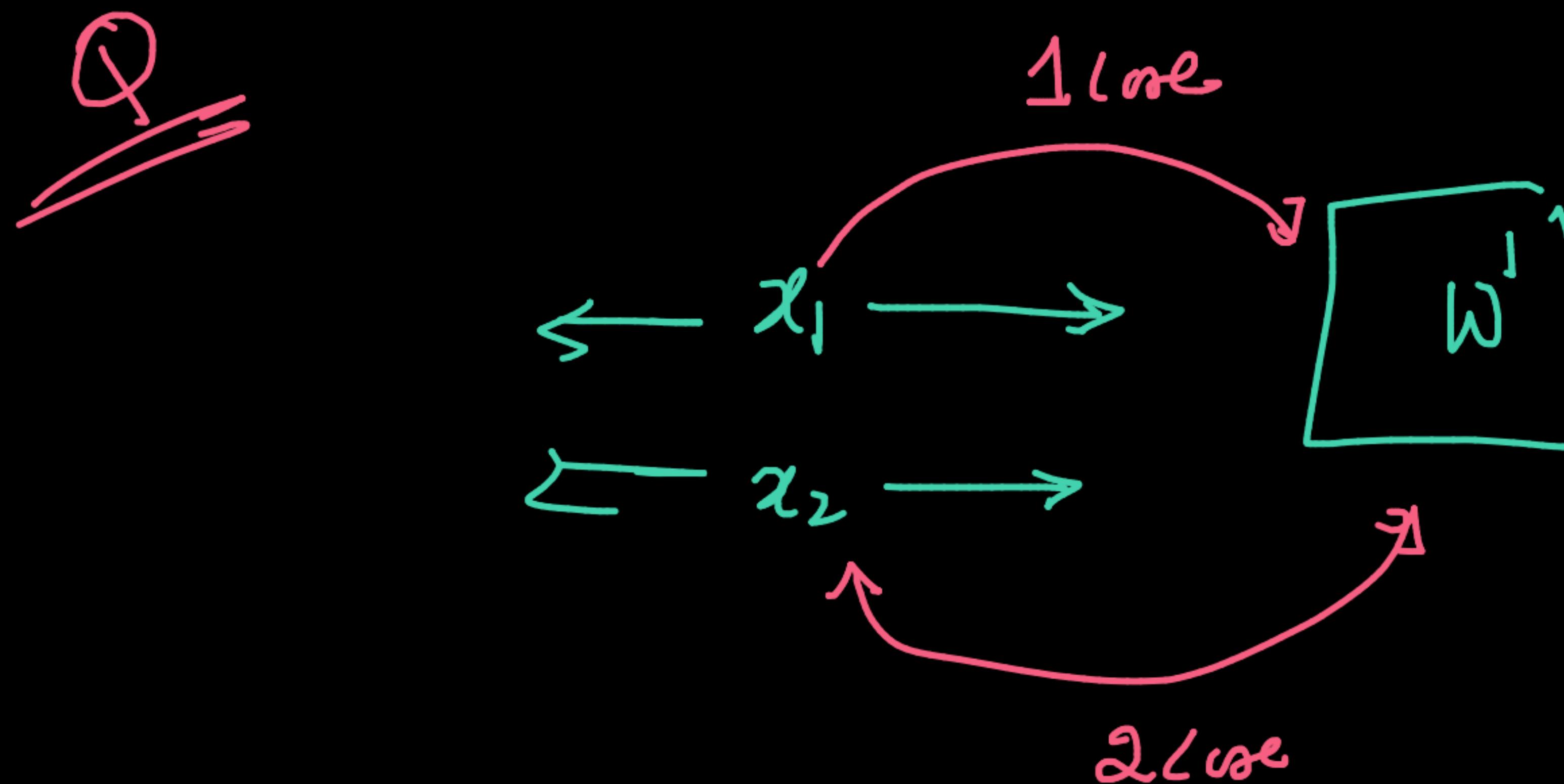


mini-batch =  $\frac{1000}{n}$  random data points for each iter<sup>n</sup>

$n=10^7$

epoch: (later) {  
1. shuffle  $10^7 = n$  pls randomly  
2.  $\frac{1000 \text{ pls}}{10^4 \text{ times /itern}}$  every time per iter<sup>n</sup>  
 $\frac{10^4 \text{ times /itern}}{10^4 \text{ times /itern}}$





0  
0  
 $\rho$   
G

CPU: 8-cores

✓ GPU: thousands of cores (lucky)

