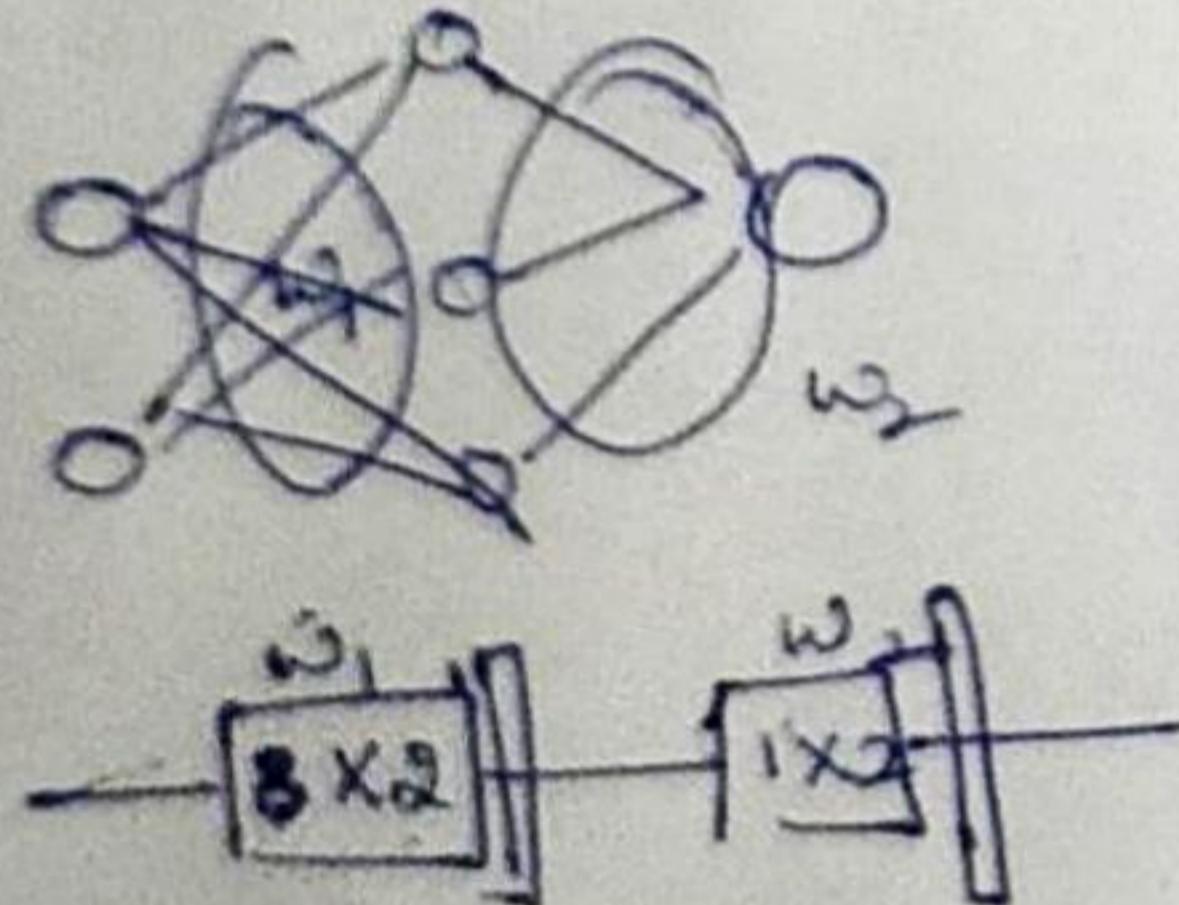


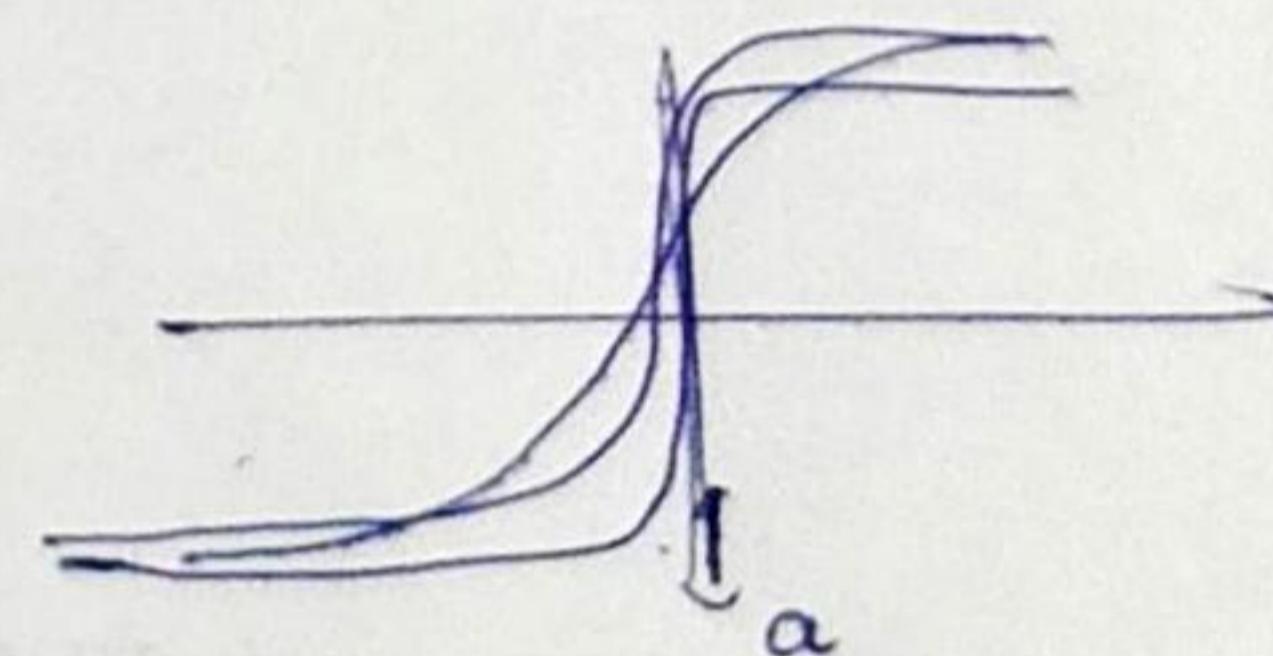
22/10/18

Lec-19

Backpropagation:-



In NLP TB, they generally draw BP's with index



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

~~Sig~~ We are more interested in  $w_1$  and  $w_2$  and what if the  $\phi$  fn has a learnable parameter

Bias  $\leftrightarrow$  Learnable Parameter  
can be incorporated

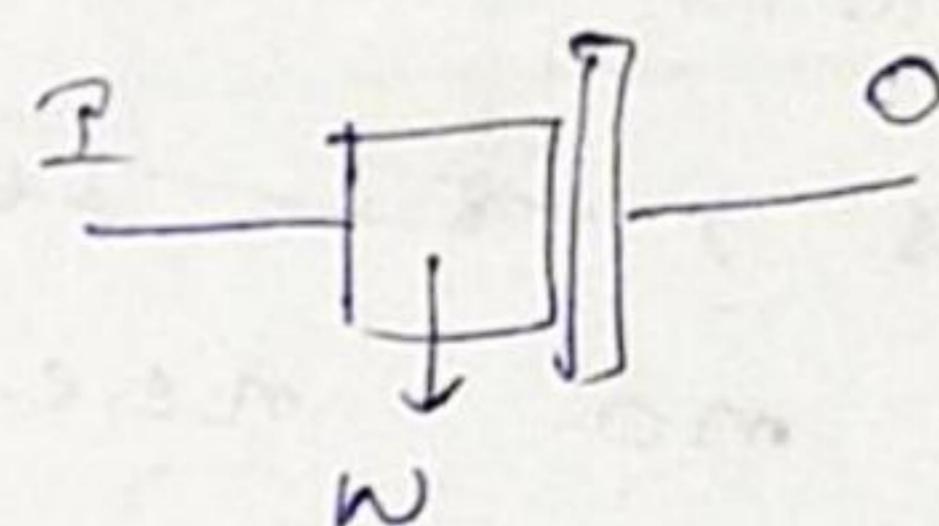
Each layer can be represented as one matrix and a fn, with no params

If  $\frac{\partial L}{\partial O}$  is available and

$O \rightarrow \text{output}$

$$\frac{\partial L}{\partial I} = \frac{\partial L}{\partial O} \times \boxed{\frac{\partial O}{\partial I}}$$

↓  
Needed

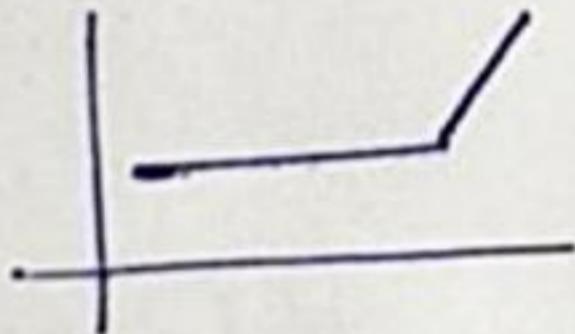


$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial O} \times \boxed{\frac{\partial O}{\partial w}} \rightarrow \text{for every block}$$

$$w^{k+1} \leftarrow w^k - \eta \frac{\partial L}{\partial w}$$

When skip connections are there

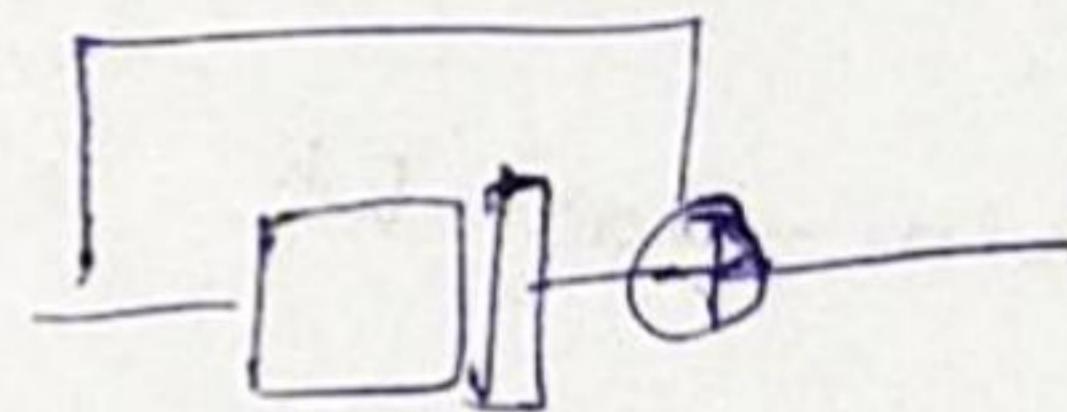
ReLU



Sigmoid

Tanh

were stopped using

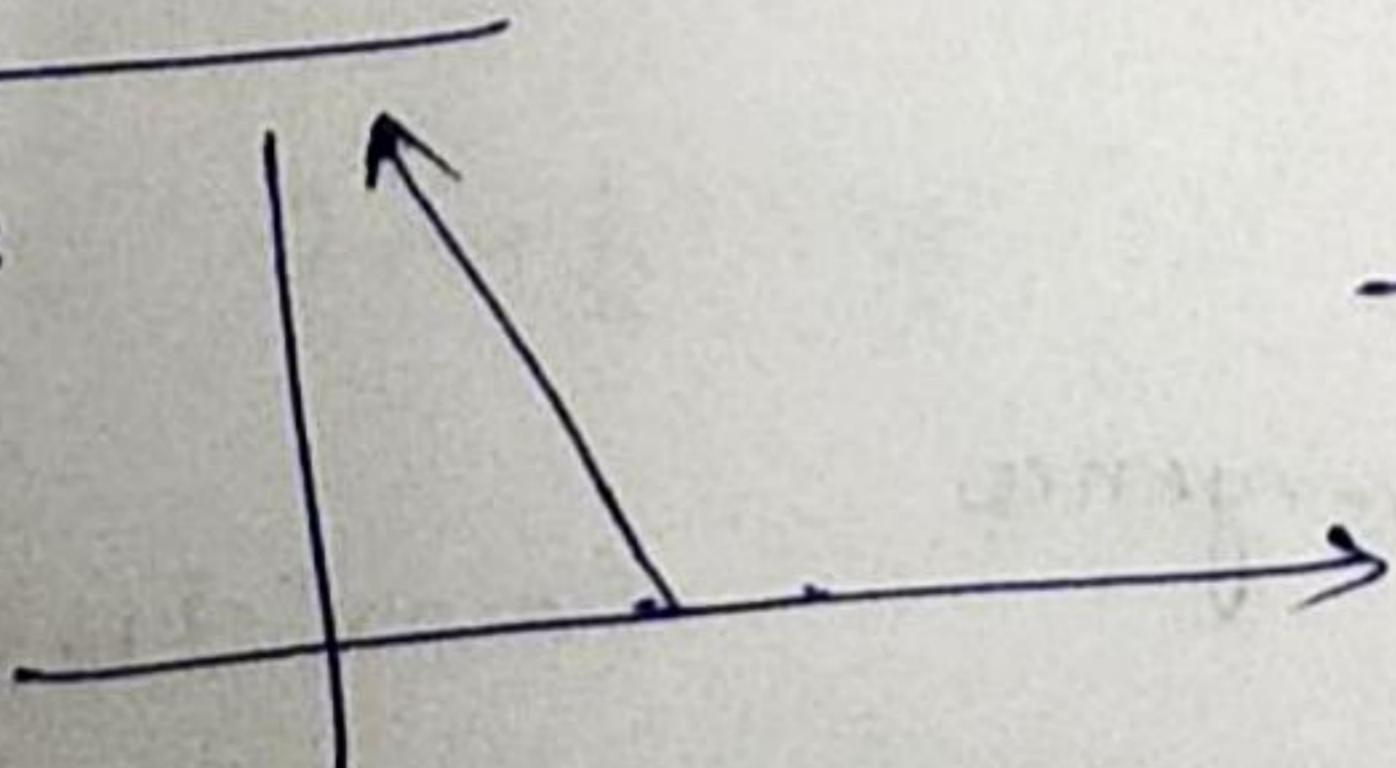


- Easy  $\frac{\partial}{\partial x}$  and non-linearity
- Differentiability at point
- Subgradients are used to solve.

• We want slopes at a point, ~~but~~ and not the equation  $(\frac{\partial L}{\partial w}|_{w=w_0})$

so

Hinge Loss



$$\max(0, (1 - t_B * b))$$

→ A loss that enforces margin?

## Stochastic Versions:-

- If loss is same for any random 100 samples, it ~~is same~~ as total 10000 samples, then, no need of using all samples.

So, That's why stochastic

- "batch" → Send all → find gradient  
find loss,
- "single sample" → 1 forward pass  
Get Loss  
Update
- "mini-batch"

$$\frac{\partial L}{\partial w} = \frac{\partial \sum l_i}{\partial w}$$

Individual samples

$$= \sum \frac{\partial l_i}{\partial w}$$

so every samples

$\frac{\partial}{\partial x}$  is ~~added~~ calculated

If 100 give same  $\frac{\partial}{\partial x}$  and as  
1000 ~~give~~  $\frac{\partial}{\partial x}$

Then 100 is enough

- ① Initialize
- ② Forward pass "find loss"
- ③ Back prop "find w"
- ④ Repeat ②-④ until convergence

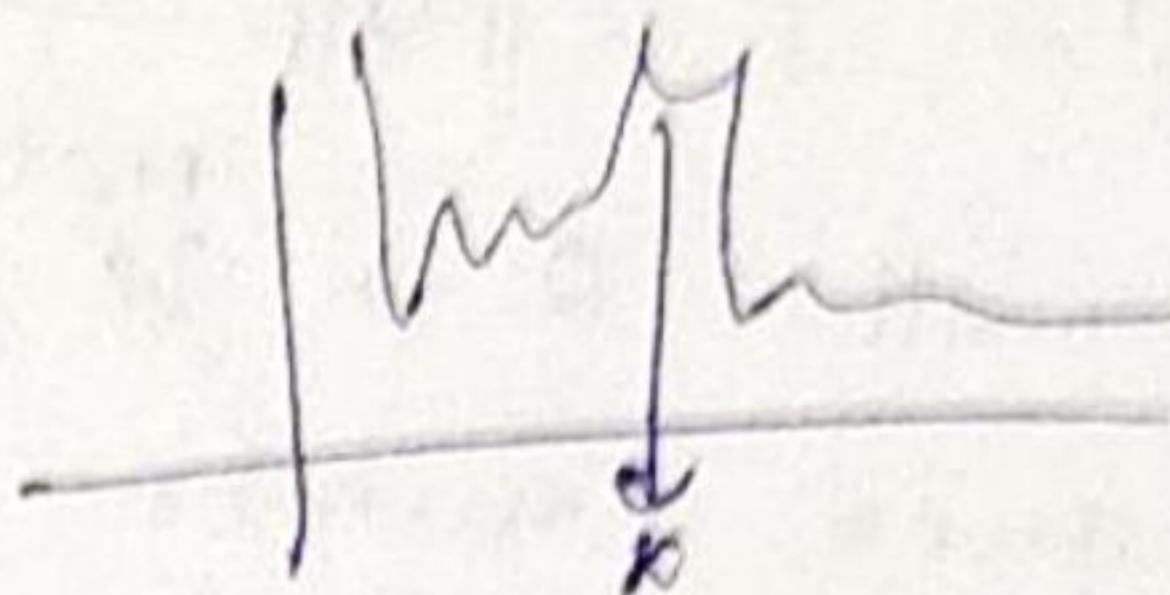
Converges:

$$L \geq 0 \\ L^0 \rightarrow L^1 \rightarrow L^2 \dots L^{n-1} - L^n$$

$$L^0 < L^1 < L^2 \dots$$

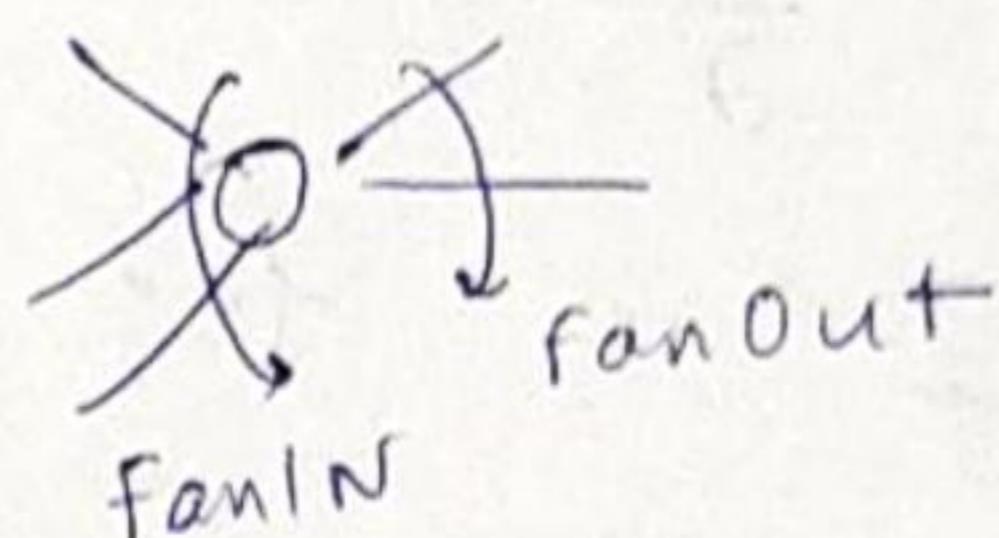
Unless } If Loss ~~should never~~ would never converge.  
u are } a shit head }

. will we reach correct minima?  
 ↗ Don't be worried about getting stuck  
 at  $\pm$ , coz u are seeing 1D  
 but other directions w, be  
 may not be minima  $\rightarrow$  (saddle  
 point)



Initialization  
 ↗ Random (why not zero)? [0 to 1]  
 or [-1 to 1]

- Xavier, He's init ||



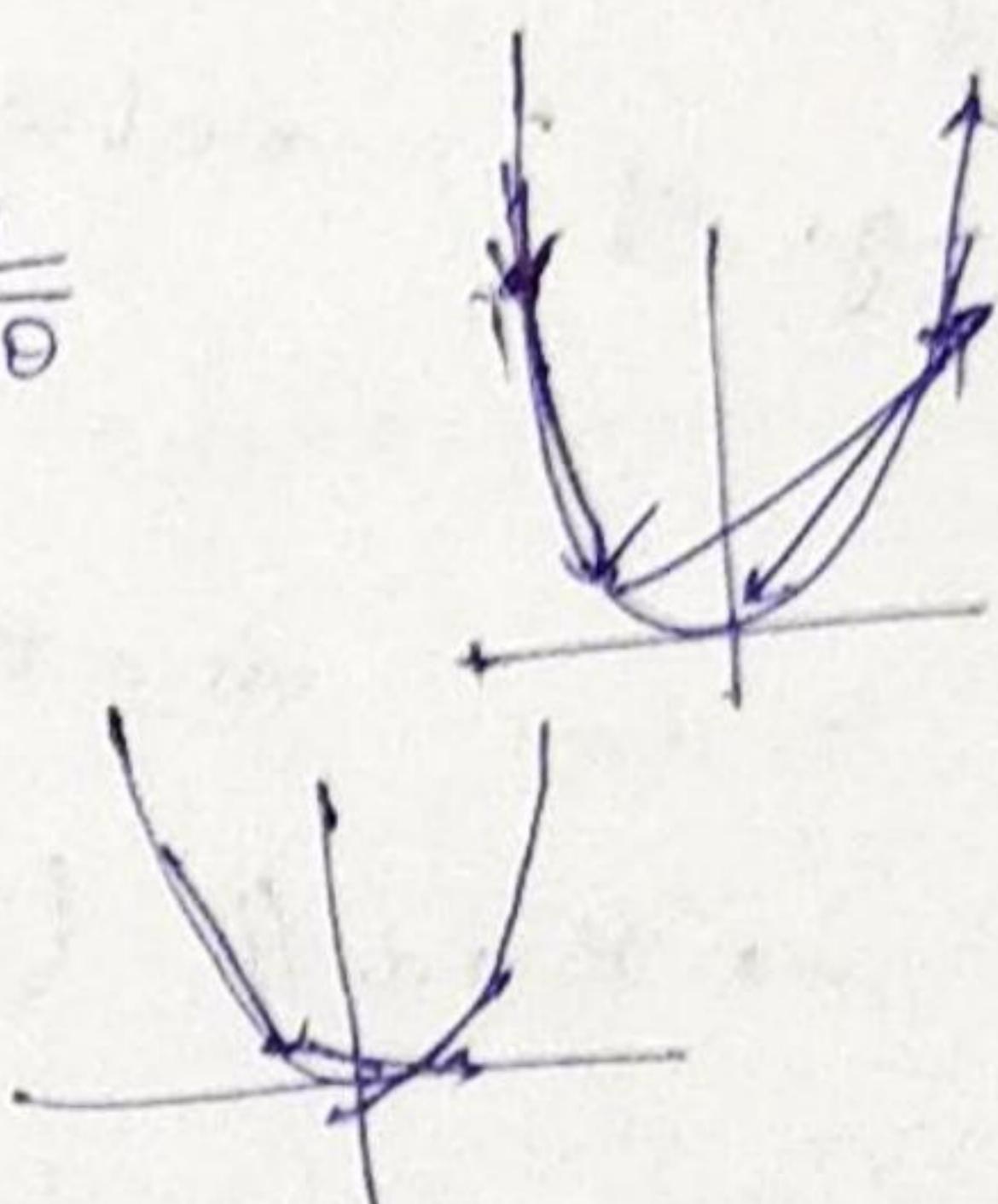
Update rule

- ↗ Momentum

$$\theta^{k+1} \leftarrow \theta^k - \eta \frac{\partial L}{\partial \theta}$$

- Previous direction

- Nesterov momentum



Adaptive Gradient

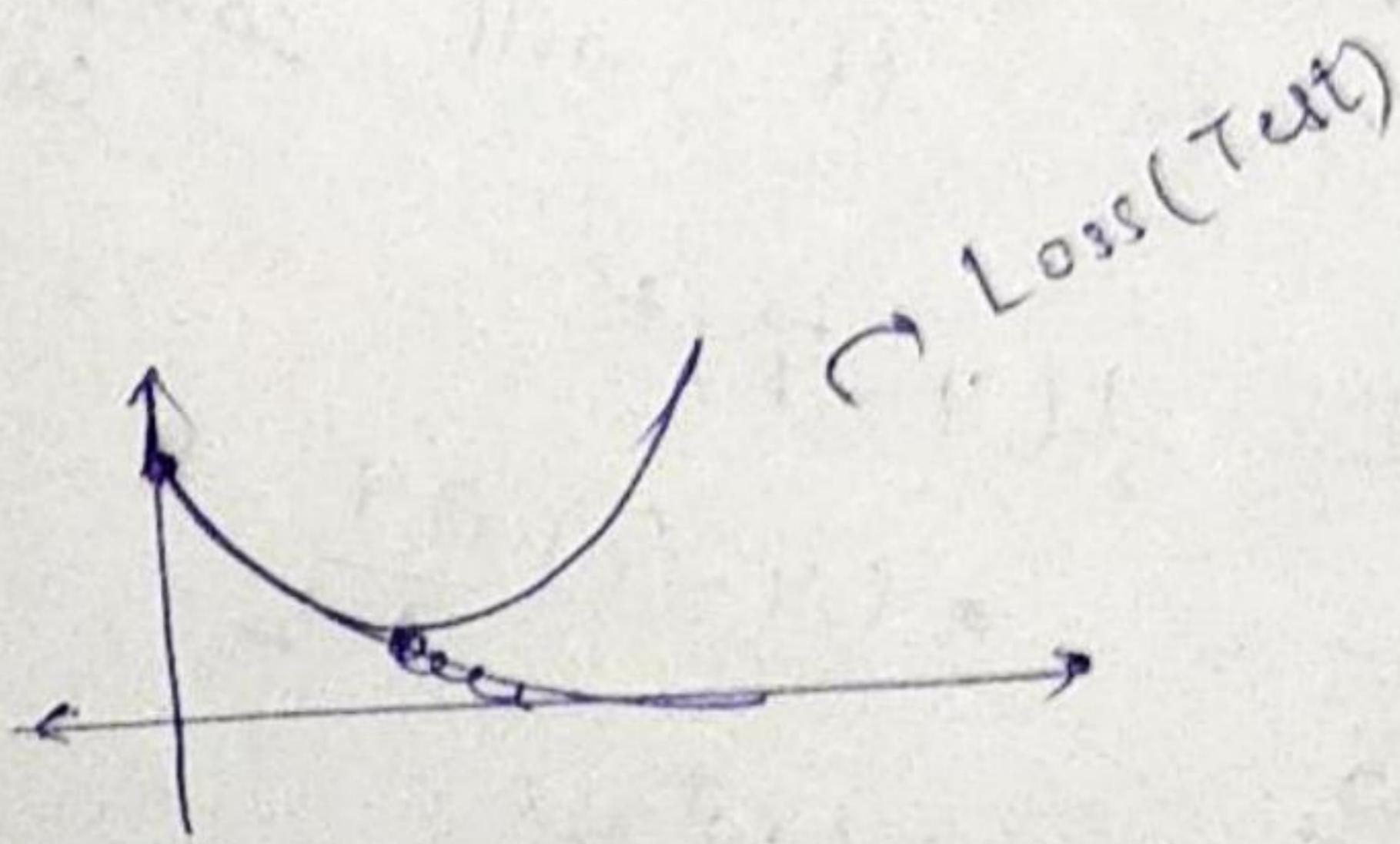
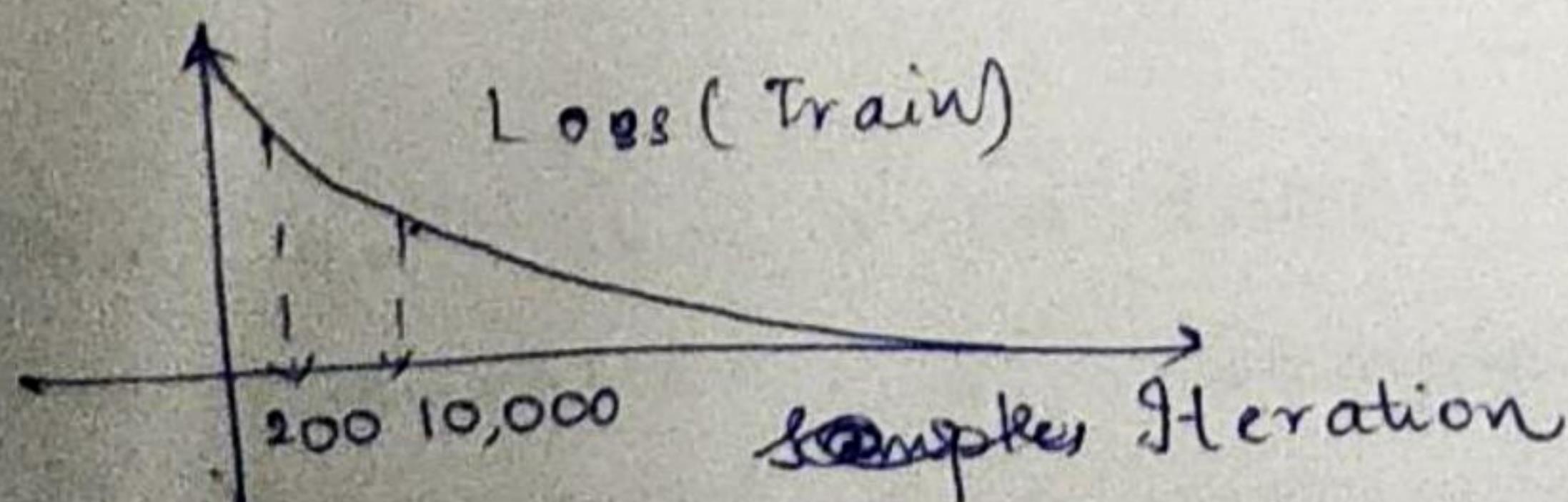
Ada delta

Ada deb

RMS Prop

Adam

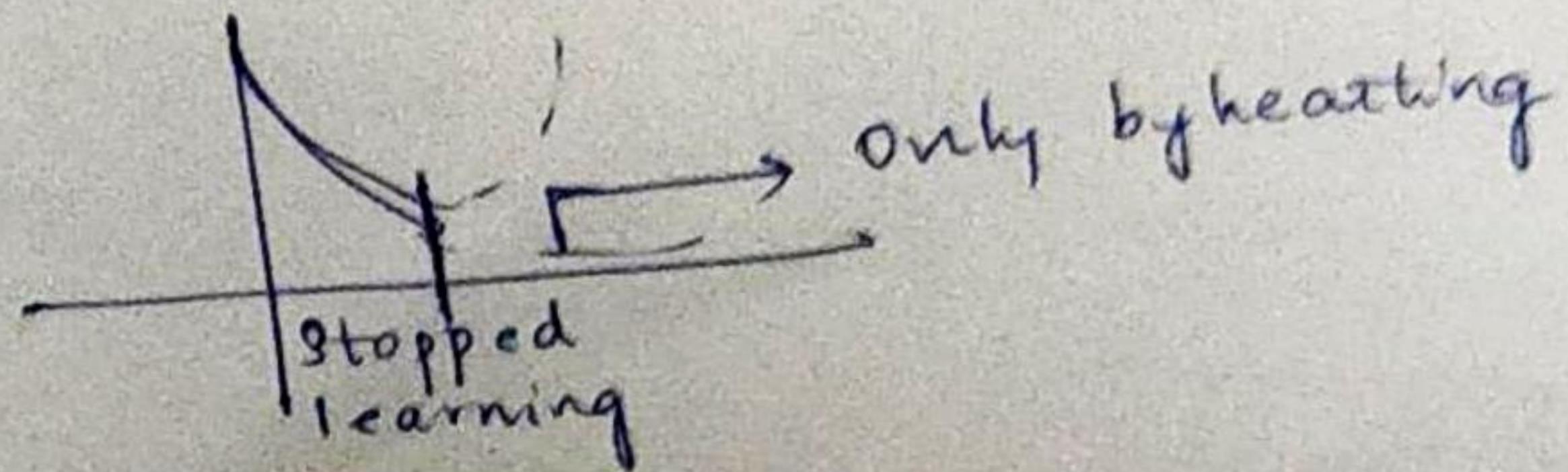
Termination: Overfitting



It should start at the same point

Many samples  $\rightarrow$  that it is super specializing  
 for training data.

so termination



(ii)  $J' = J + \|\theta\|_1 \rightarrow L_0, L_1, L_2$   
 If small  $\theta \rightarrow$  Not enough weights: Sopke No by-heart  
 So, reduce number of params.

Weight decay

If small weights, remove those edges  
 force to recover by back propagation

$$\frac{\partial J'}{\partial \theta} = \frac{\partial J}{\partial \theta} + \frac{\partial \|\theta\|_1}{\partial \theta}$$

$$\frac{\sqrt{x_1^2 + x_2^2 + \dots}}{(x_1 + x_2 + \dots)^{1/2}}$$

In practice  $L_2$  is for  $\|\theta\|_2$

[Pruning of networks]

Second order methods:

$H^{-1} \rightarrow \text{Bad} // (\theta \times \theta)$  → fast convergence (Good)

Neural networks:

$$\phi = 2x$$

↳ Linear fn

Then it will combine collapse into 1

$$E = \frac{1}{2}(y - \tilde{y})^2$$

$$\frac{\partial E}{\partial w_2} = \dots (y - \tilde{y}) \frac{\partial y}{\partial w_2}$$

$$y = \phi(w^T x)$$