

Lecture started with a brief information about project submission, paper presentation dates and submission and Homework 4 brief information and read map to follow to solve it.

* A problem:- The vanishing gradient:

$$\text{With large no. of layers } (\delta^{(l)}) \rightarrow 0$$

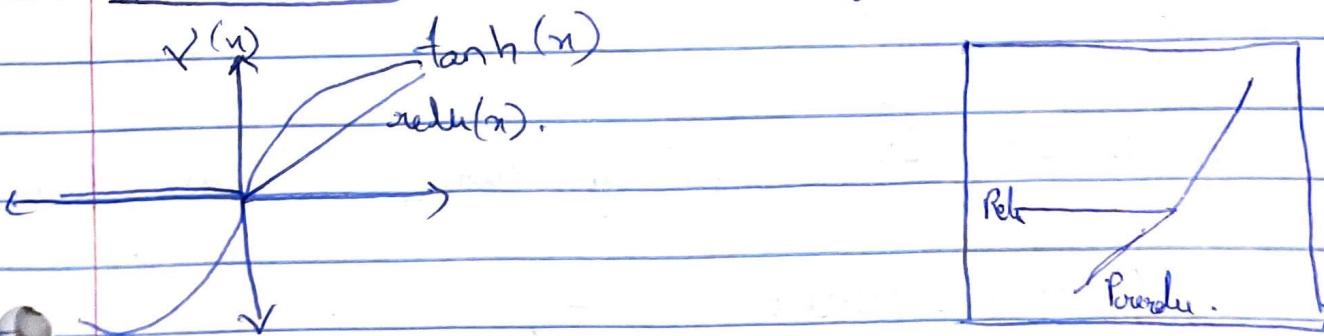
$$G^{(l)}(z_n) = [n^{(l)} (\delta^{(l)})^T]$$

$$w^{(l)} \leftarrow w^{(l)} - \eta G^{(l)}$$

As $\delta^{(l)} \rightarrow 0$, then $G^{(l)} \rightarrow 0$.

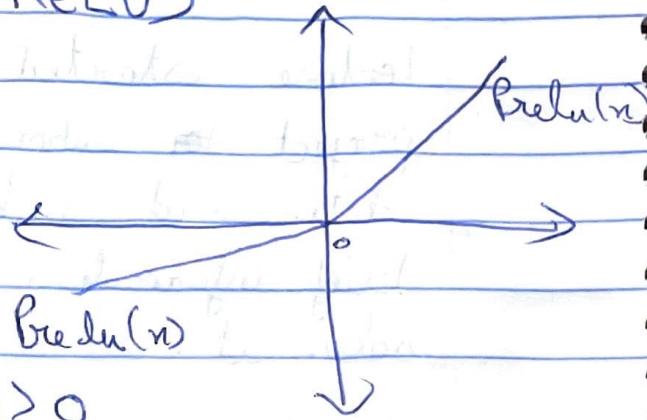
So when we try to update weight using above formul. There is no change noticed.

* A Solution:- The ReLU (Rectified Linear Unit)



Benefit:- Helps in creating sparse values as
derivative is 0 if values < 0

* Parameterised ReLU (PReLU)



$$\frac{d}{dn} \text{relu}(n) = \begin{cases} 1 & \text{if } n \geq 0 \\ 0 & \text{else.} \end{cases}$$

$$\text{Prelu}(n) = \begin{cases} n & \text{if } n \geq 0 \\ \text{learnable parameter} & \text{else.} \end{cases}$$

learnable parameter.

* A similar Problem:- Exploding Gradient.

Artificially scaling down to a ~~so~~ so a norm is under budget.

* Breaking the symmetry during initialization.

Ways to handle are :-

- 1] Random initialization (Uniform Distribution)
- 2] Xavier initialization (Uniform distribution from network statistics)

* Dropout Regularization

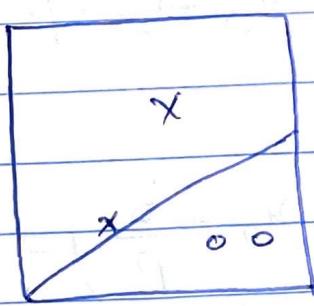
Masking few nodes at random

$$x_i^{(l)} = \Theta(p_i^{(l)} s_i^{(l)})$$

randomly masking \rightarrow the nodes by making the output 0.

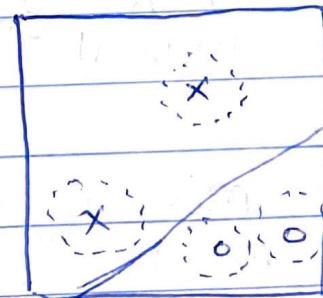
Dropout rate needs ^{can} to be specified at each layer.

* Robustness to Noisy Data.



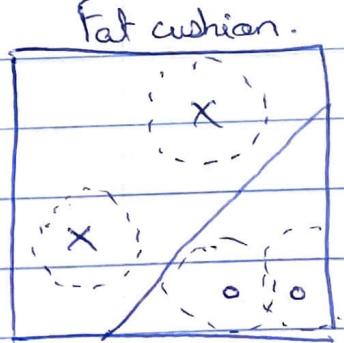
$$E_{in} = 0$$

But less robust
(Can handle lesser error)



$$E_{in} = 0$$

more robust



$$E_{in} = 0$$

has the most robust
(Can handle
noisy data
more efficiently)

Fat cushion.

Increasing the thickness effectively reduces the hypothesis set.

Thick cushion makes our model more efficient it is equivalent to picking ~~to~~ the hypothesis which is regularised.

Separating Data.

$$\text{sign}(\omega^T x_n + b) = y_n$$

$$y_n (\omega^T x_n + b) \geq 0 \quad (\text{always } +\text{ve})$$

\Rightarrow there is a p such that.

$$\min_{n=1, \dots, N} y_n (\omega^T x_n + b) = p$$

$$\min_{n=1, \dots, N} y_n \left(\frac{\omega^T x_n + b}{p} \right) = 1.$$

rescaling the weights, is an identical hyperplane

Remember this fact.