

1 Vanishing grad | Exploding grad. ✓

2 Data Scaling → Batch Norm. ✓

3 Drop out ↵

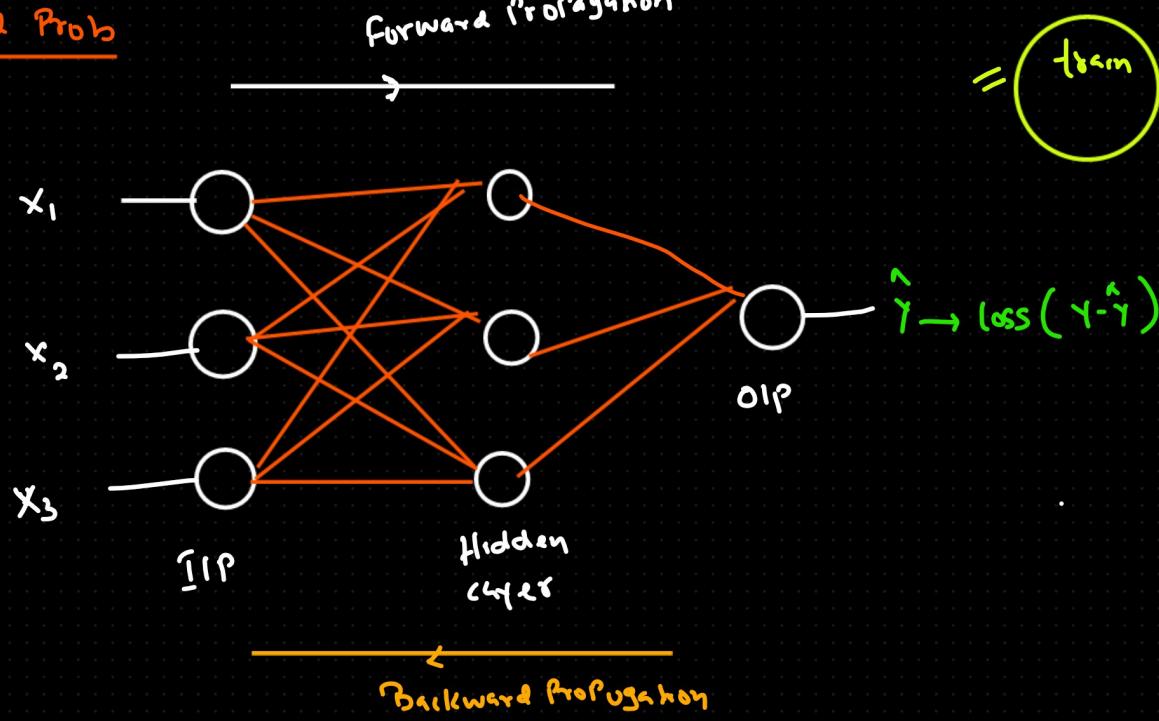
4 Regularization ↵

5 Weight Initialization ✓

** Start from how to improve the performance of the NN

Vanishing grad Prob
"

Forward Propagation

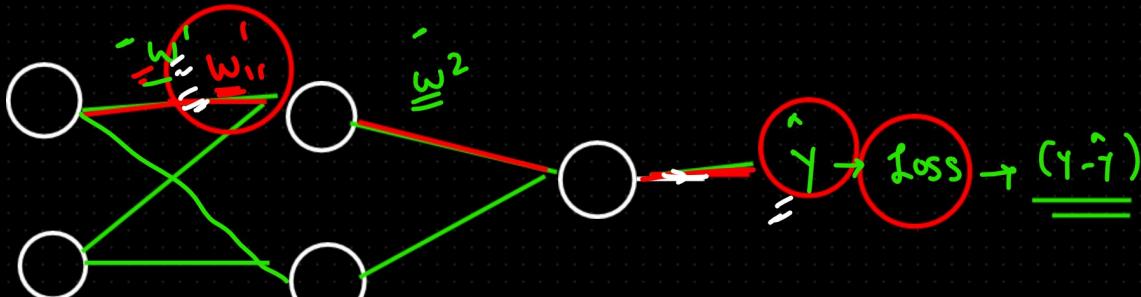


Gradient will be very very small \rightarrow There won't be any sort of a weight update

Vanishing

Vanishing gradient

Neural Network \Rightarrow Layers Nodes \Rightarrow Multiple \Rightarrow Deep neural network



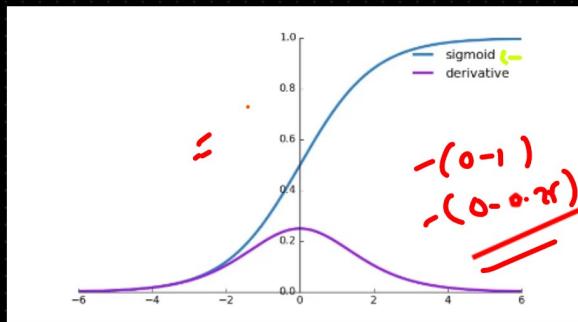
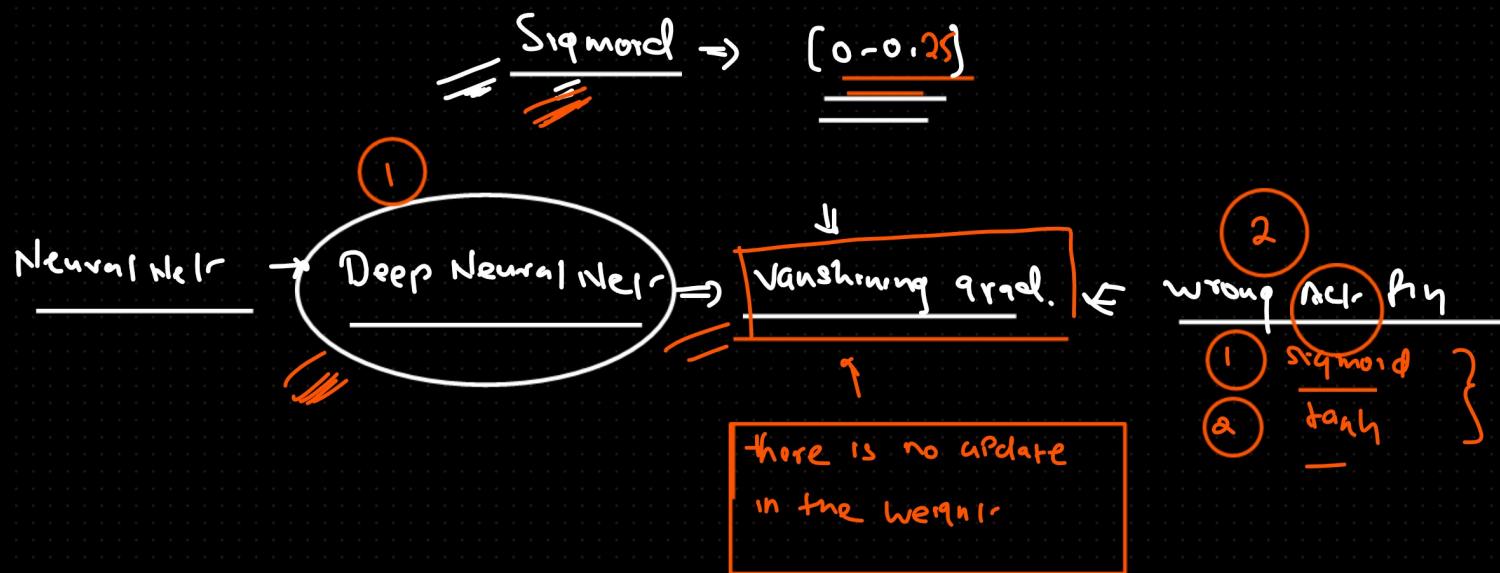
$$w_n = w_0 - h \frac{\partial L}{\partial w}$$

$$\left[\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial I} \times \frac{\partial I}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial w_i} \right]$$

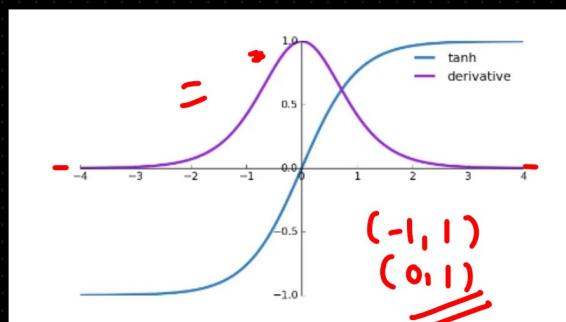
$$= \text{curr}(I/p \times \text{weights}) = 0.12 \Rightarrow \text{Backpropagation}$$

Derivation / Partial Derivation

[0-1]

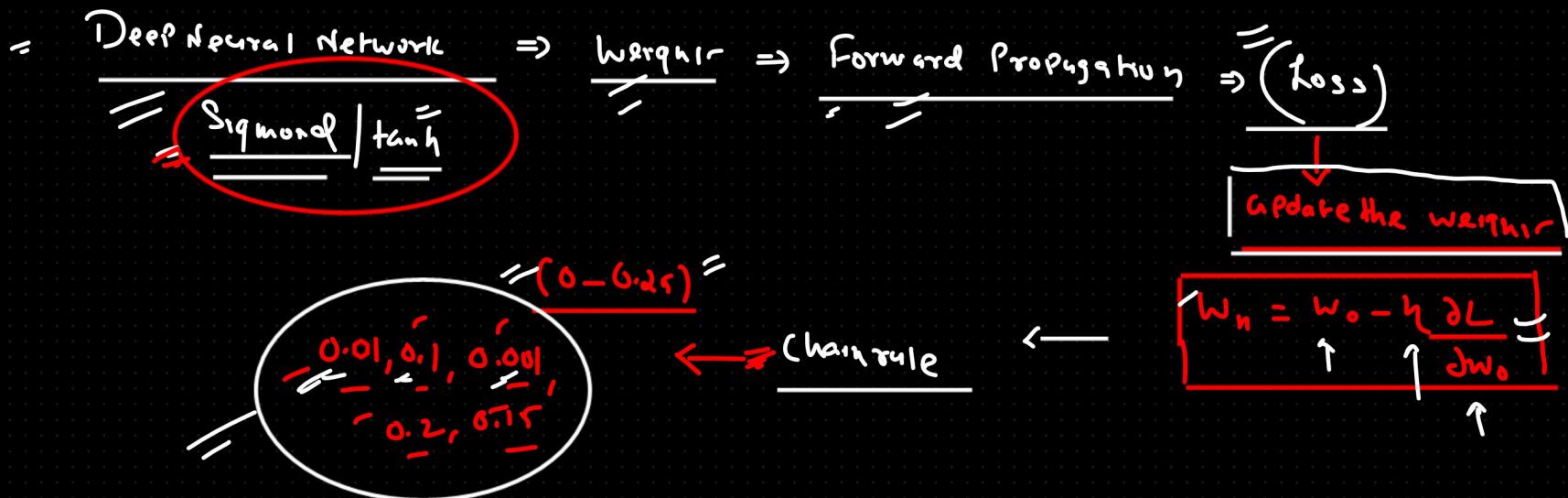


Sigmoid



tanh.

$$= \overbrace{0.1 \times 0.1 \times 0.1 \times 0.1}^{\text{Forward Propagation}} = \underline{\underline{0.0001}}$$



Small \times Small \times Small $\dots =$ very small value

When $= \frac{1}{w_0 + \eta \frac{\partial L}{\partial w}}$

$$= 1 - 0.00001$$

$$= 0.9999$$

negligible

$0.9999 - 0.00001$

gradient value

which is going to be vanish

Vanishing grad.

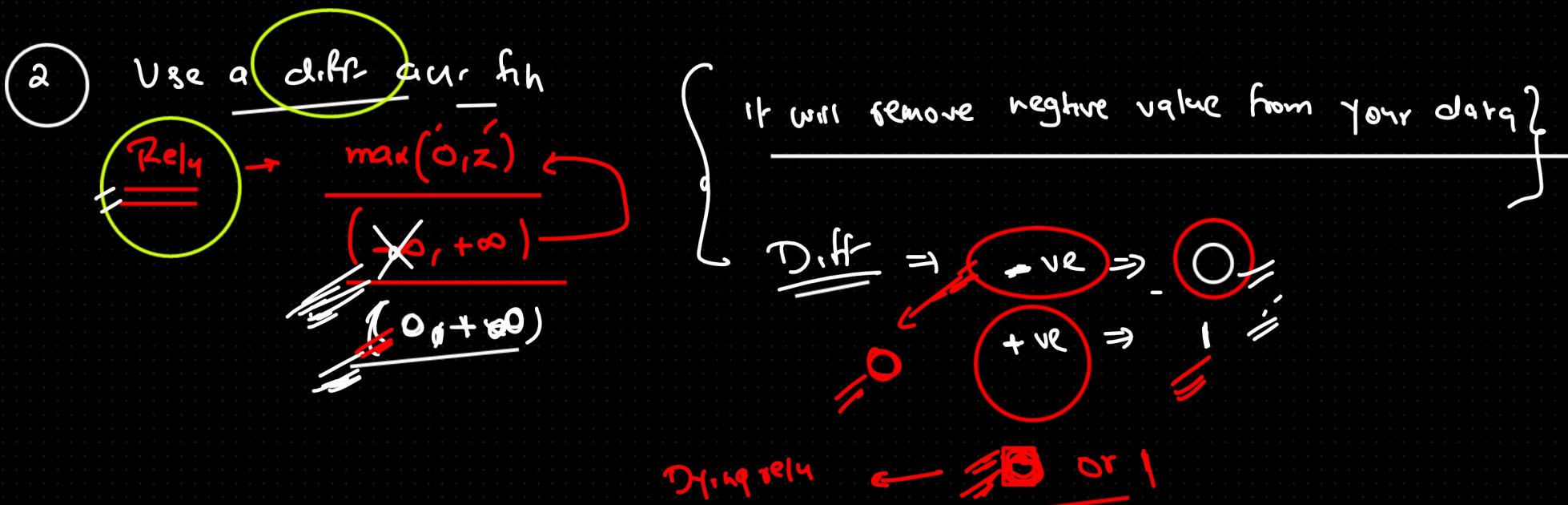
Deep Neural Net

- 1 Loss will not change
- 2 Weight will not change

Vanishing grad (How you gonna resolve it)

- 1 Reduce the Complexity of your NN:
 - ↳ Reduce the No. of layer
 - ↳ Reduce the No. of Nodes

(Shallow Neural Net)



chain rule

$$\frac{\partial L}{\partial w} = \underbrace{f(\dots)}_{\text{function}} \leftarrow \underbrace{w}_{\text{weight}}$$

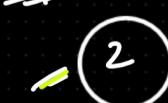


3

Proper weight initialization



Xavier ←



Glorot ←

4

Batch Normalization ←

How to Improve a Performance of Neural Network?

1

Handle OR TC Vanishing gradient

- ↳ Activation f_g - Refer the Bappy's notes
- ↳ Correct weight initializing -

2

Overfitting

- ↳ Early stopping
- ↳ Dropout -
- ↳ Regularization
- ↳ Reduce complexity (layer, Node)
- ↳ Increase data

3

Normalization (fast training)

- ↳ Normalization / Std.
- ↳ Batch Normalization -

Can refer the Bappy's notes

4

Optimizer -

- ↳ Momentum
- ↳ AdaGrad
- ↳ Adadelta / RMSprop
- ↳ Adam

5

HyperParameter

- ↳ No. of hidden layer -
- ↳ Nodes / Neuron -
- ↳ Batch size -
- ↳ Learning rate -

6

Activation

- ↳ sigmoid
- ↳ tanh
- ↳ type Relu
- ↳ softmax
- ↳ linear / step

Hyperparameter tuning
using Keras Tuner

Can refer Bappy's notes

X

Data Scaling

$$\frac{\text{Age}}{[1-95]} = \frac{\text{Salary}}{[1 - 10000000]}$$

!

Keep the data in same scale

$$= \text{Input} \times \text{weight} + \text{bias}$$

$$= \frac{\text{Age} \times w_1 + \text{Salary} \times w_2 + \text{bias}}{1000000}$$

To ensure distribution of data is same ie; normally distributed with mean=0 and std dev=1

1 Std

$$\frac{x_i - \mu}{\sigma}$$

$$\begin{aligned} & (-3, +3) \\ & \xrightarrow{\text{Z-table}} \xrightarrow{\text{STD}} \end{aligned}$$

Salary

To compress the values in the range [0,1] so that to eliminate any bias that may arise due to different scale or magnitude of data(Recall Avnish's example)

2 Normalization

$$\frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

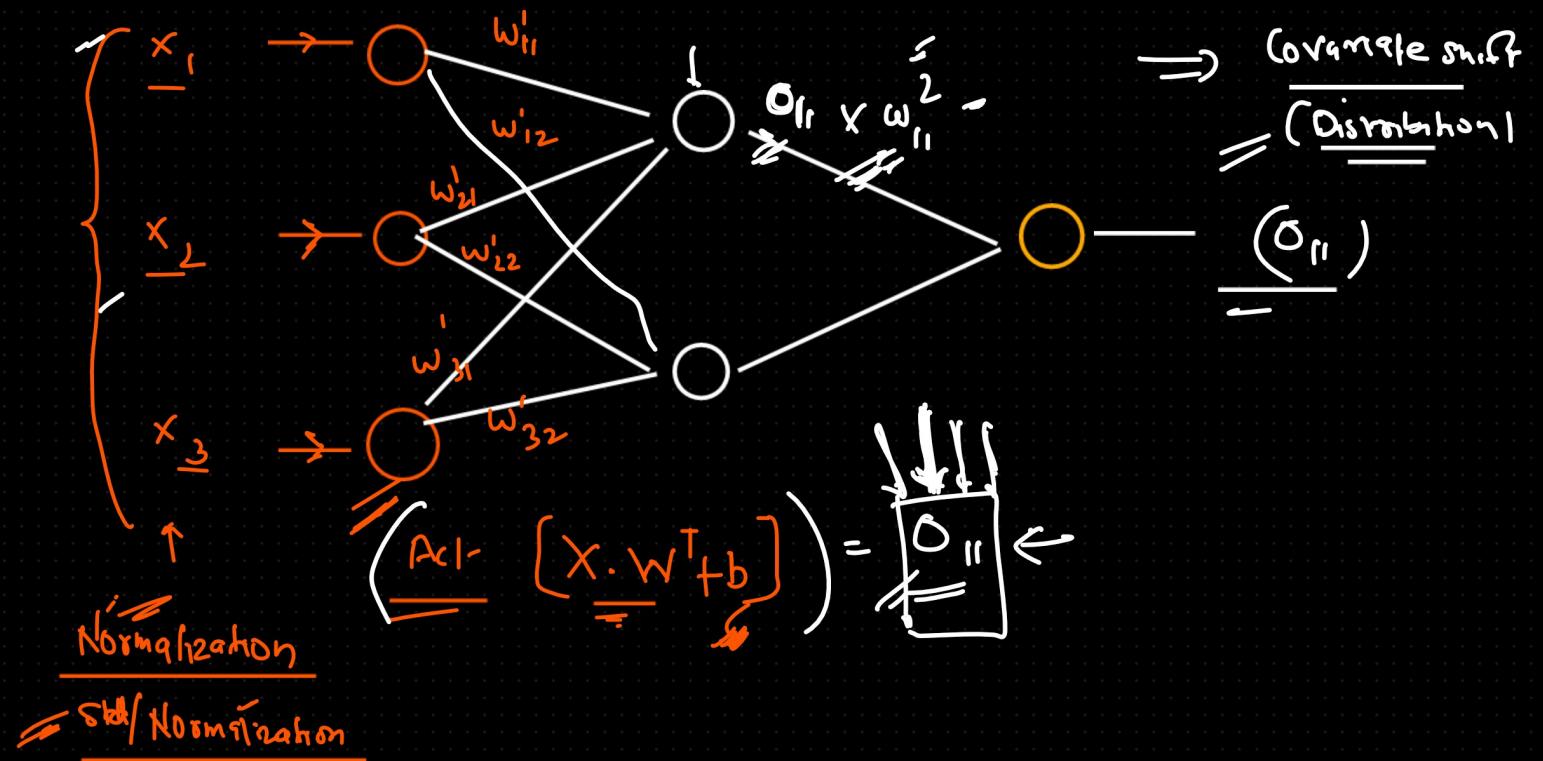
[0-1]

GPA
(0-10)

Zigzag

Wavy

Contours

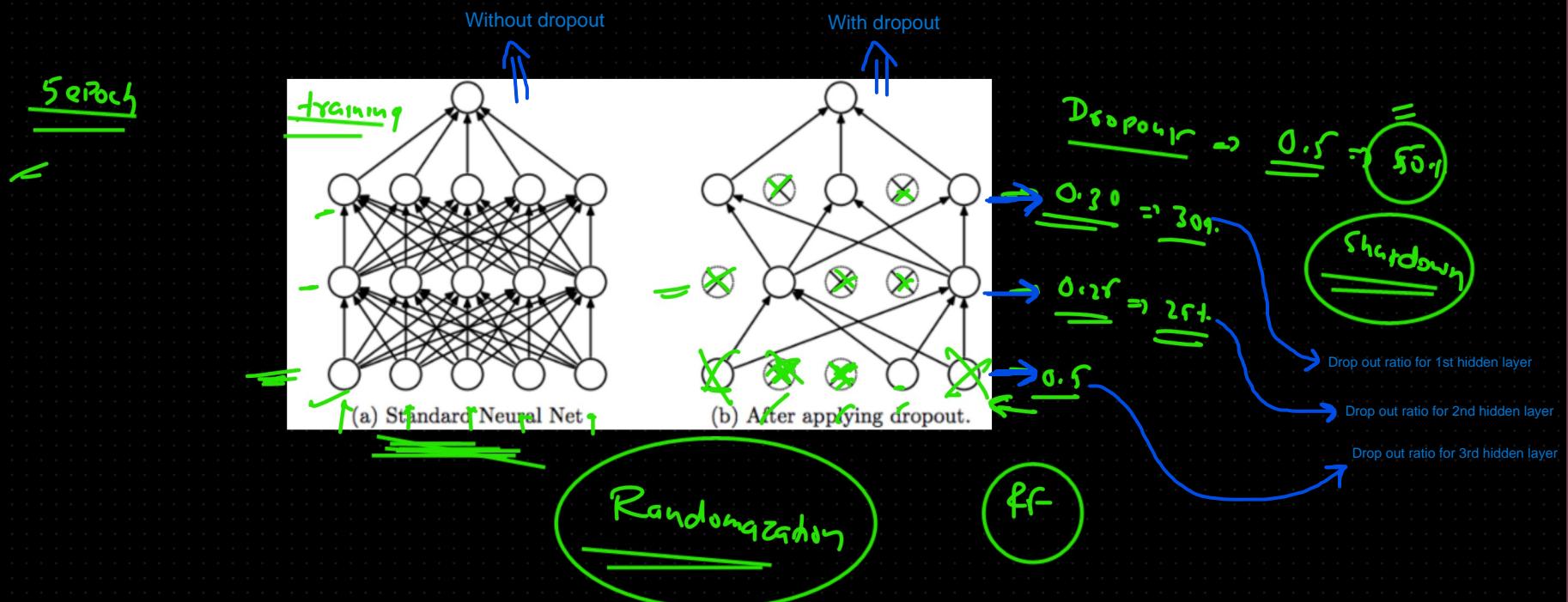
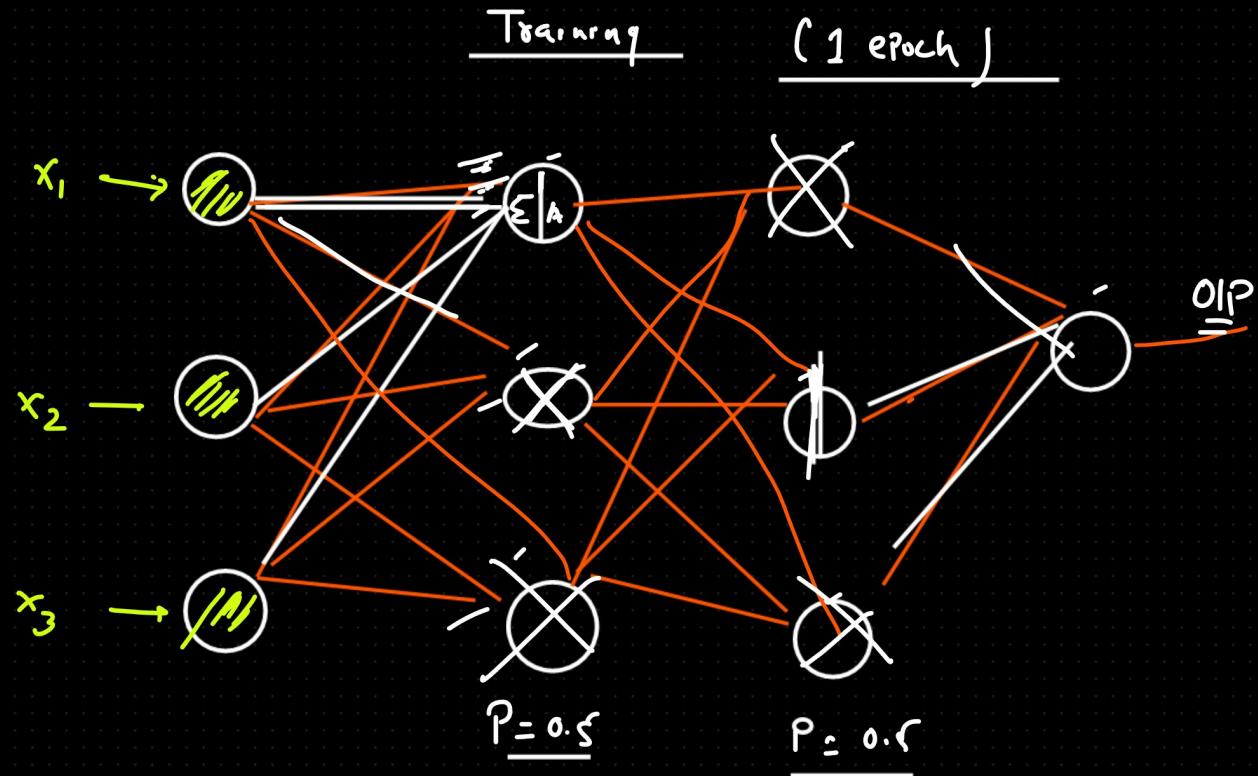


~~Drop out~~

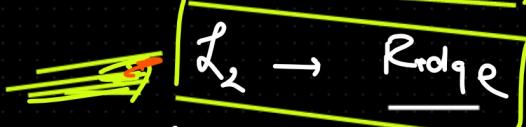
* So, basically in Drop out we are shutting down few neurons in each hidden layer randomly. Due to this the data is passed only through the neurons which are not shut down and overfitting is addressed.

* Shutting down of the neuron is done randomly in each epoch. Therefore in each epoch we are going to get the different set of neurons through which data is basically getting passed. This is similar to what we see in the random forest. This is how we are enabling randomization in this.

* No. of neurons to be picked up randomly in each epoch for shutdown is determined using Dropout ratio. We assign Dropout ratio for each hidden layer



#

Overfitting

$L_1 + L_2 \rightarrow$ elastic

toward to 0 (more completely zero)

* Regularization can be used to add the penalty (regularization parameter) in the loss/cost function which addresses the overfitting issue.

* L1 generates the sparsity since it makes weight 0 hence in DL we prefer using L2 or Ridge regularization where weights are tending to 0.

Loss functions $\{L_1, L_2\}$

$$L_1 \rightarrow \text{Loss} + \frac{\text{Penalty}}{\lambda |w|}$$

$$L_2 \rightarrow \text{Loss} + \frac{\text{Penalty}}{\lambda |w|^2}$$

$$L_1 + L_2 \rightarrow \text{Loss} + \lambda(|w_1| + |w_2|)$$

In case of multiple weights may be using the determinant which will yield the single value single value

Cost

$$\text{Loss} + \lambda \sum_{i=1}^n |w_i|^2$$

