

1 Vanishing grad | Exploding grad. ✓

2 Data Scaling → Batch Norm. ✓

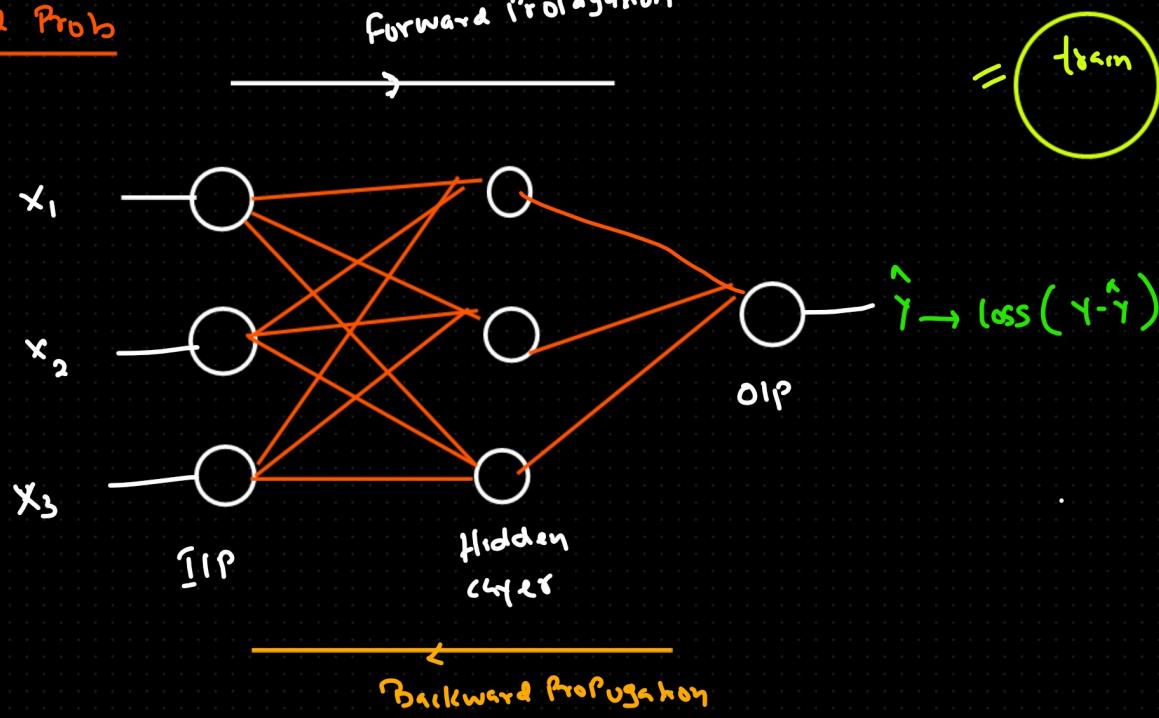
3 Drop out ↵

4 Regularization ↵

5 Weight Initialization ✓

Vanishing grad Prob
"

Forward Propagation

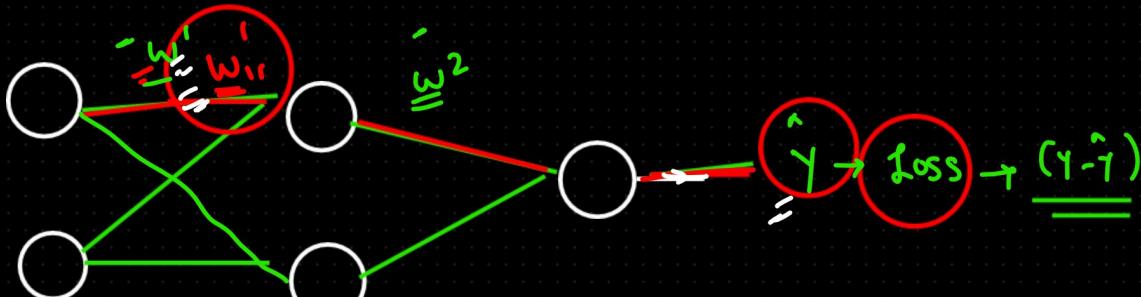


Gradient will be very very small \rightarrow There won't be any sort of a weight update

Vanishing

Vanishing gradient

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



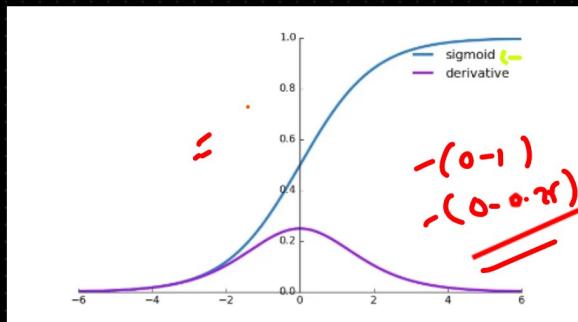
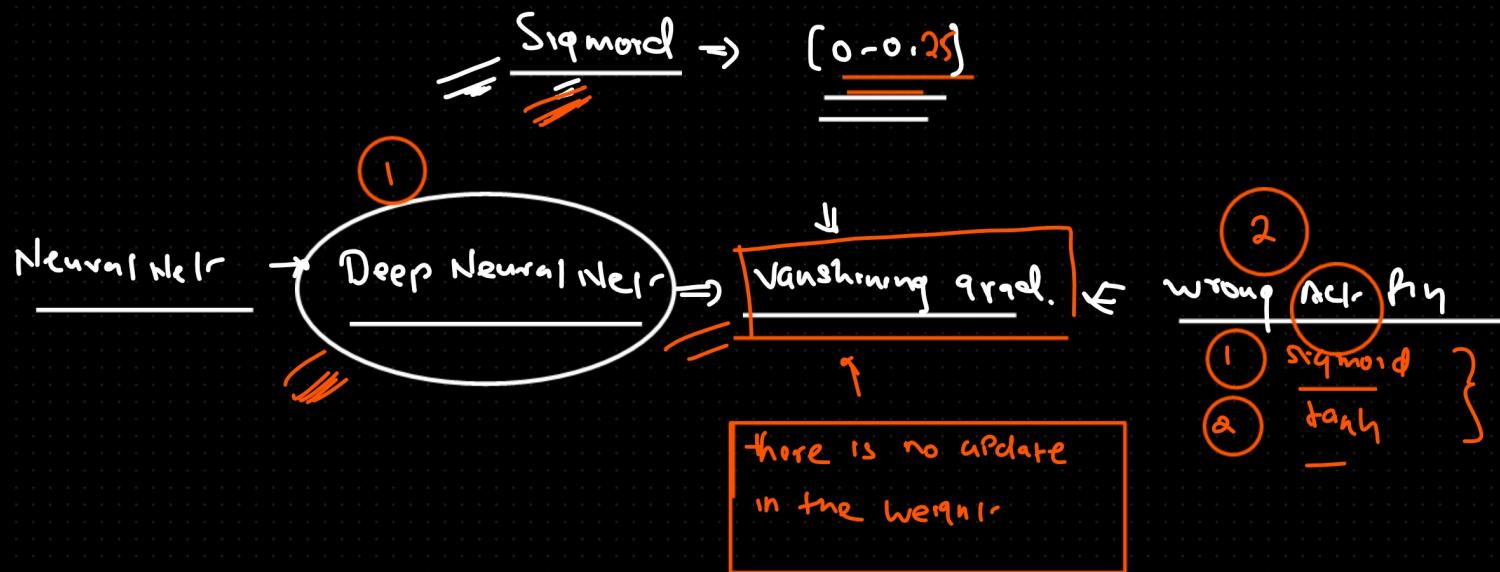
$$\hat{w}_n = w_0 - \eta \frac{\partial L}{\partial w}$$

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

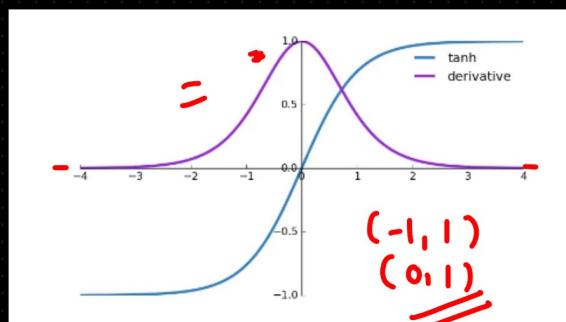
$$= \text{curr}(I/p \times \text{weights}) = 0.12 \quad \Rightarrow \text{Backpropagation}$$

Derivation / Partial Derivation

[0-1]

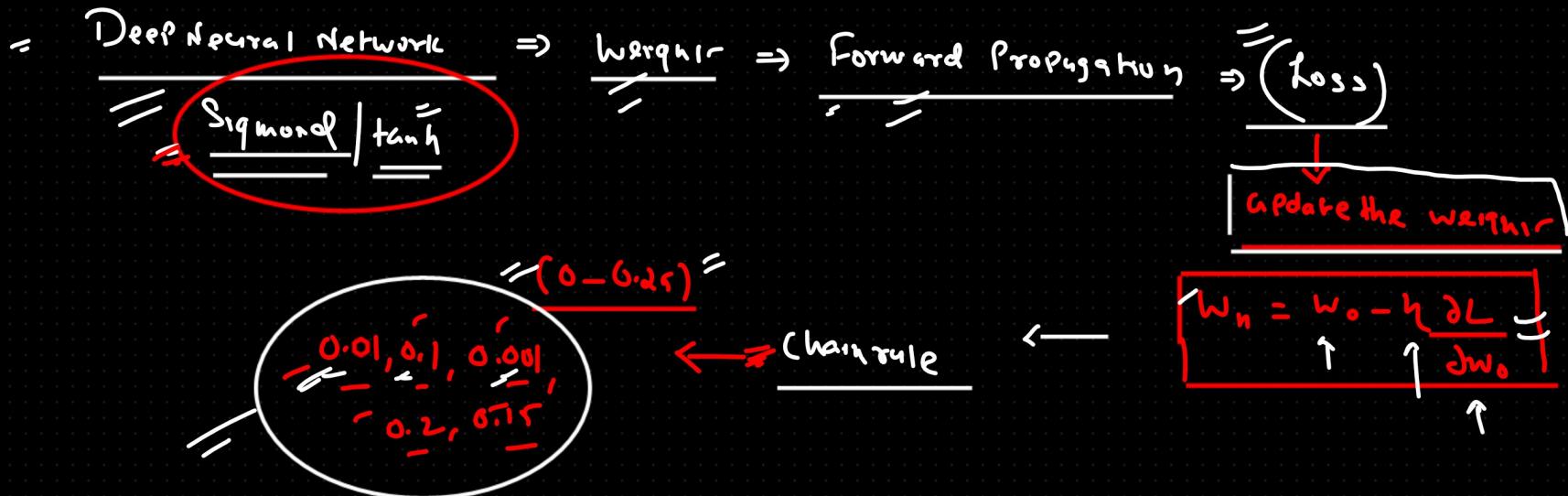


Sigmoid



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)

- 2 Use a cliff-cut fn

Rely

$$\max(0, z) \quad (-\infty, +\infty)$$

$$0 \quad +\infty$$

If will remove negative value from your data

Dif \Rightarrow -VR \Rightarrow 0

+VR \Rightarrow 1

Dif \Rightarrow 0 or 1

chain rule

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

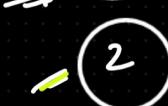


3

Proper weight initialization



Xavier ←



Glorot ←

4

Batch Normalization ←

How to improve a Performance of Neural Network ?



Handle OR TC vanishing gradient

- ↳ Activation f_g -

- ↳ Correct weight initializing -



Overfitting

- ↳ Early stopping

- ↳ Dropout -

- ↳ Regularization

- ↳ Reduce complexity (layer, Node)

- ↳ Increase data



Normalization (fast training)

- ↳ Normalization / Std. -

- ↳ Batch Normalization -



Optimizer -

- ↳ Momentum

- ↳ AdaGrad

- ↳ Adadelta / RMSProp

- ↳ Adam



HyperParameter

- ↳ No. of hidden layer -

- ↳ Nodes / Neuron -

- ↳ Batch size -

- ↳ Learning rate -



Activation

- ↳ sigmoid

- ↳ tanh

- ↳ type Relu

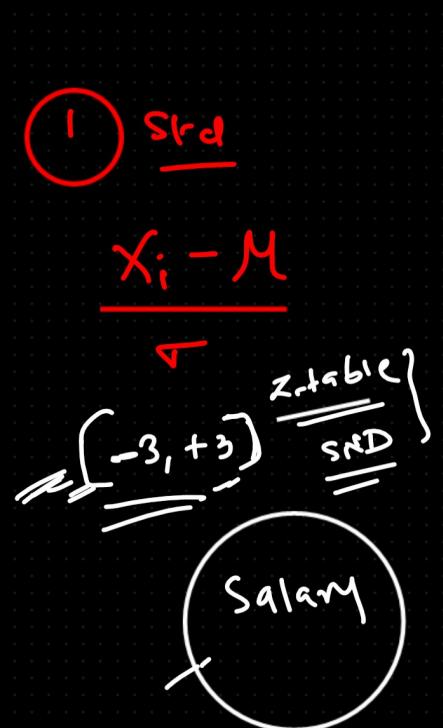
- ↳ softmax

- ↳ Linear / Step

Data Scaling

$$\frac{\text{Age}}{[1 - q_5]} = \frac{\text{Salary}}{[1 - 10000000]}$$

Keep the data in d same
scale

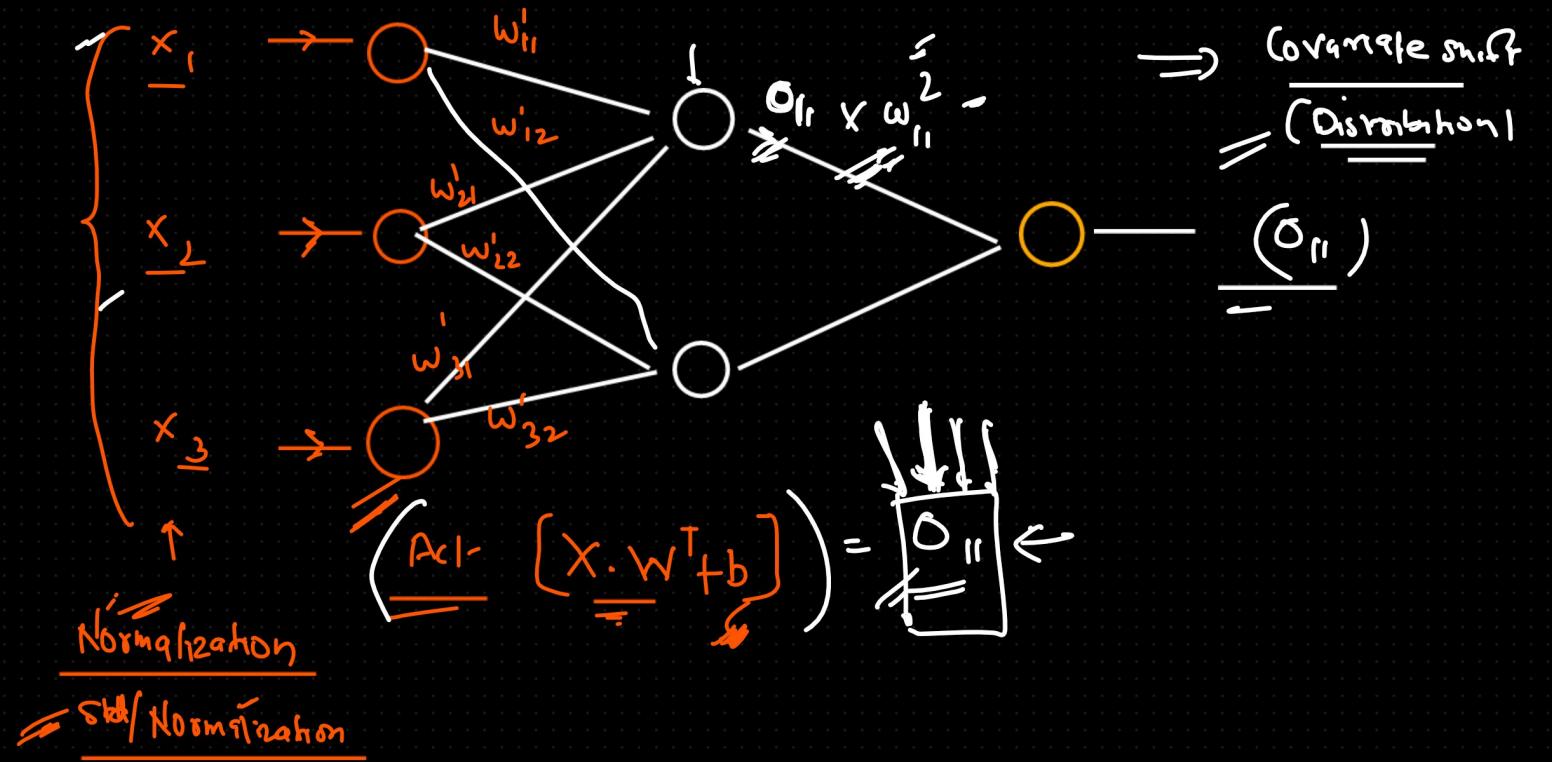


2 Normalization

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

[0-1]

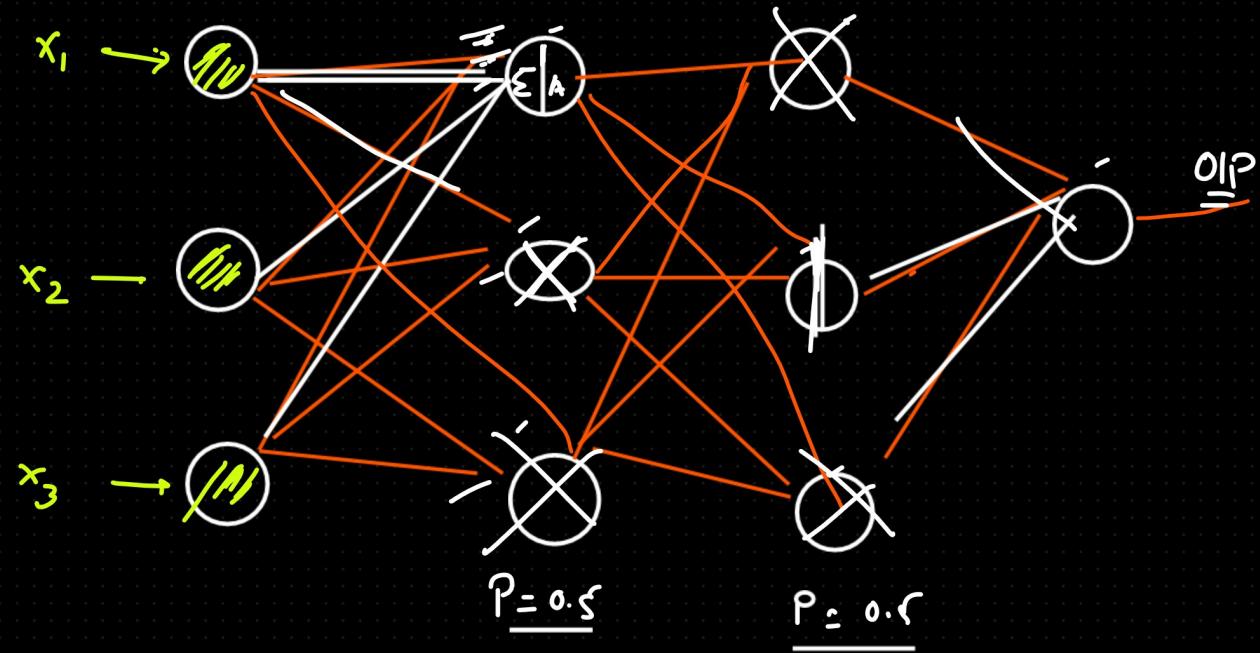
(0-10)



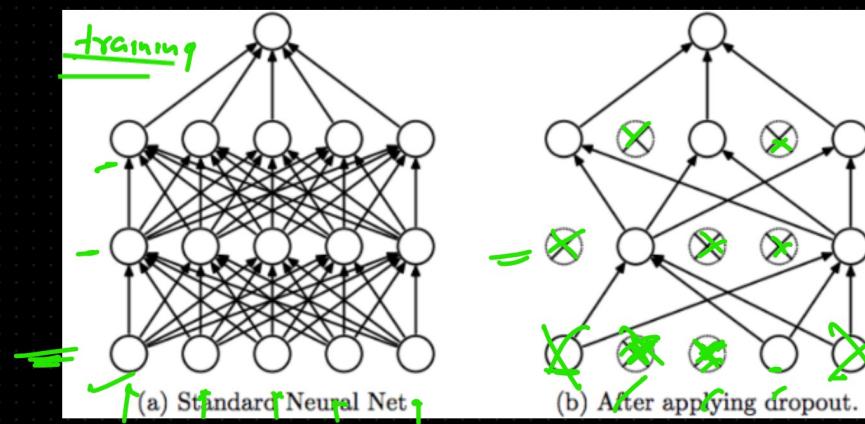
Drop out

Training

(1 epoch)



5 epoch

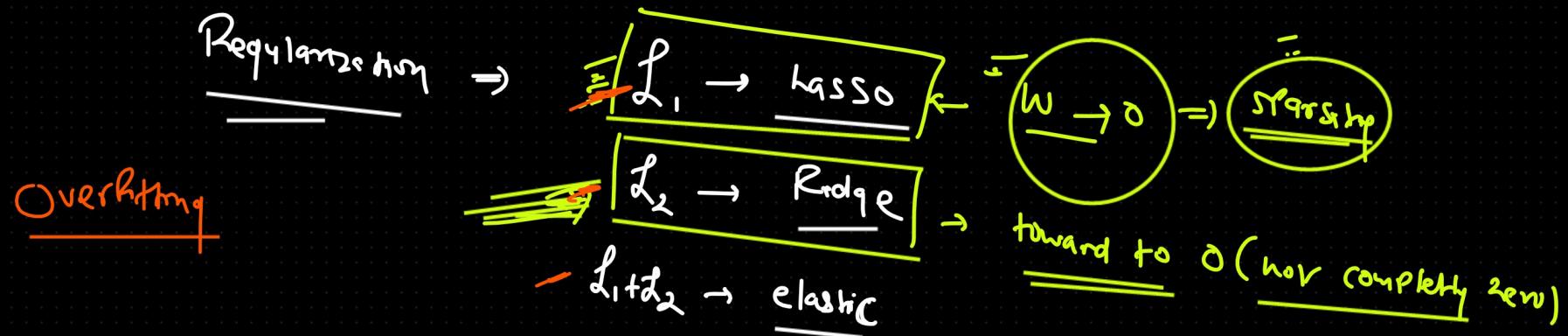


$$\begin{aligned} \text{Dropout} &\rightarrow 0.5 \Rightarrow 50\% \\ \rightarrow 0.30 &= 30\% \\ \rightarrow 0.28 &= 28\% \\ \rightarrow 0.5 & \end{aligned}$$

$\approx 50\%$
shutdown

Randomization

ff



$$\text{Loss} + \frac{\text{Penalty}}{\mathcal{L}_1, \mathcal{L}_2 / \mathcal{L}_1 + \mathcal{L}_2}$$

$$\begin{aligned} \mathcal{L}_1 &\rightarrow \text{Loss} + \lambda |w| \\ \mathcal{L}_2 &\rightarrow \text{Loss} + \lambda |w|^2 \end{aligned} \quad \left. \right\} \text{Single value}$$

$$\mathcal{L}_1 + \mathcal{L}_2 \rightarrow \text{Loss} + \lambda(|w_1| + |w_2|)$$

Cost

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

