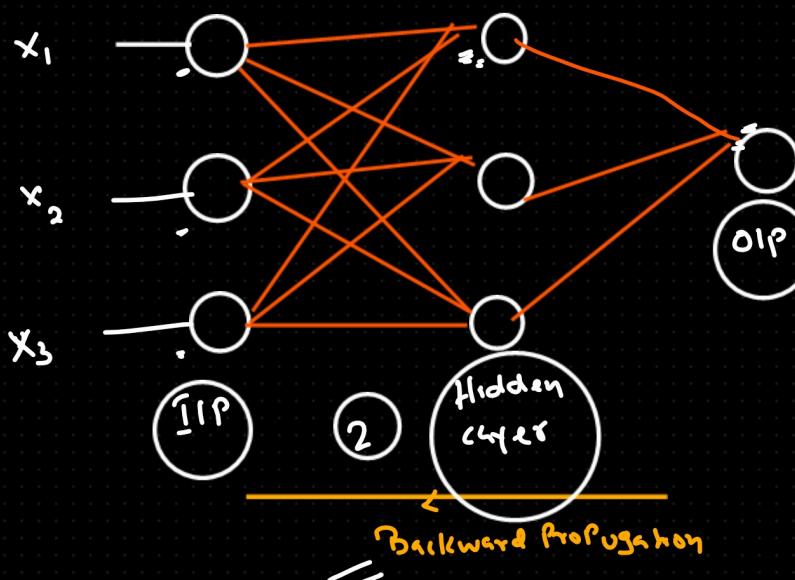


- ~~1~~ Vanishing grad | exploding grad. ✓
- 2 Data Scaling → Batch Norm. ✓
- 3 Dropout ↘
- 4 Regularization ↗
- 5 Weight Initialization ↗

- 1 Batch norm
- 2 Weight init.
- 3 Vanishing grad
- 4 PyTorch
- 5 Hyperparameter tuning
(Keras tuner)

Vanishing grad Prob

1 forward propagation



epoch

= train

$$\gamma \rightarrow \text{loss}(\hat{\gamma} - \gamma)$$

$$\left(\sum + \frac{\Delta(\gamma)}{\gamma} \right)$$

Sigmoid, tanh, ReLU

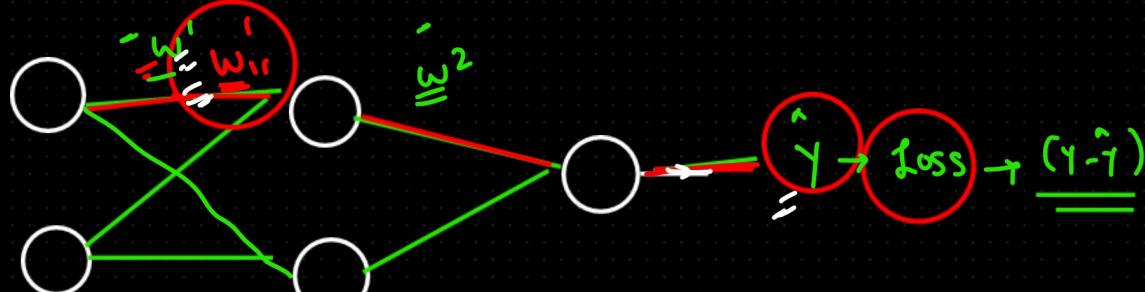
Softmax

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



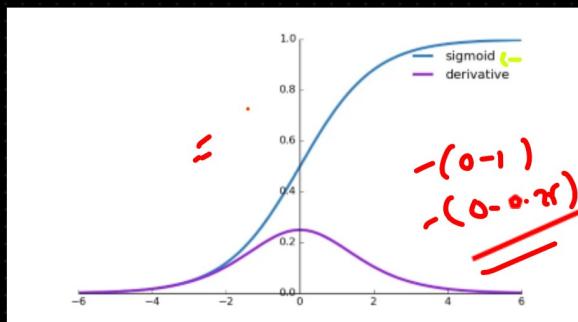
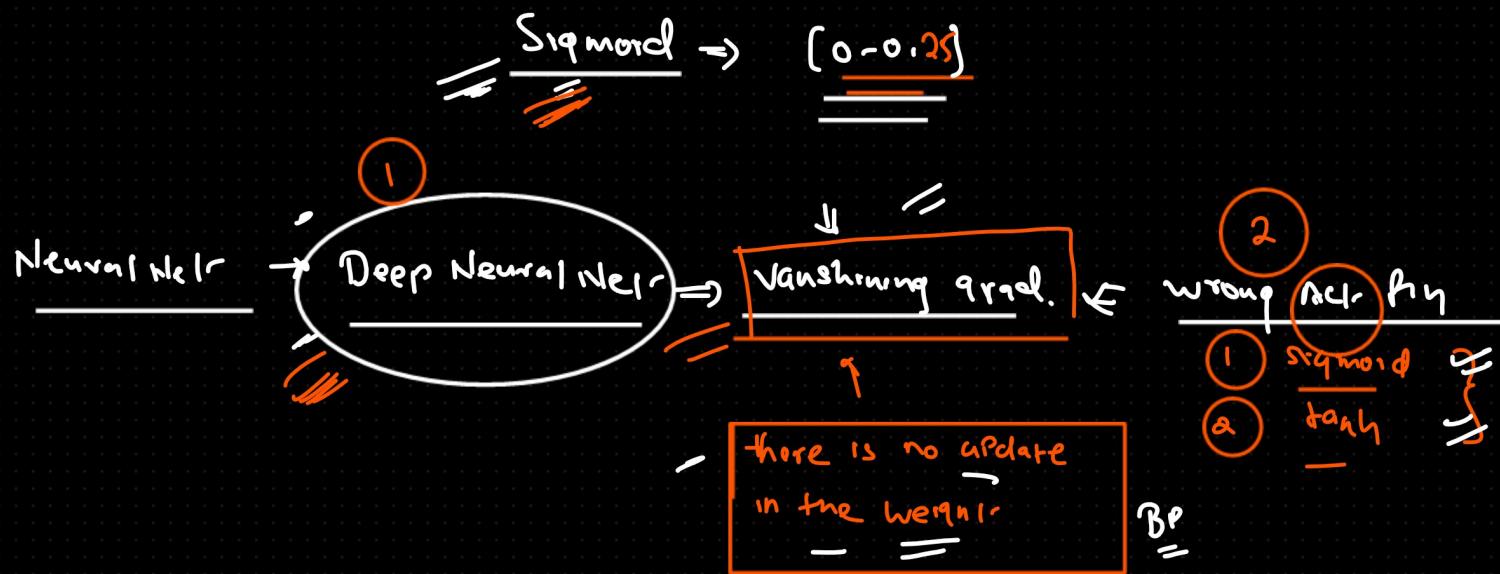
$$\boxed{w_n = w_0 - \eta \frac{\partial L}{\partial w}}$$

$$\boxed{\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial I} \times \frac{\partial I}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial w_i}}$$

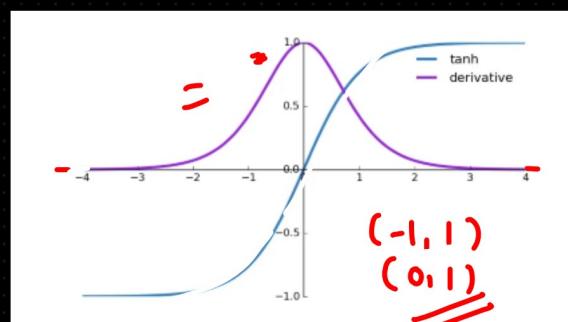
$$= \boxed{\text{curr}(I/p \times \text{weights}) = 0.12} \Rightarrow \boxed{\text{Backpropagation}}$$

(Derivation / Partial Derivation)

[0-1]

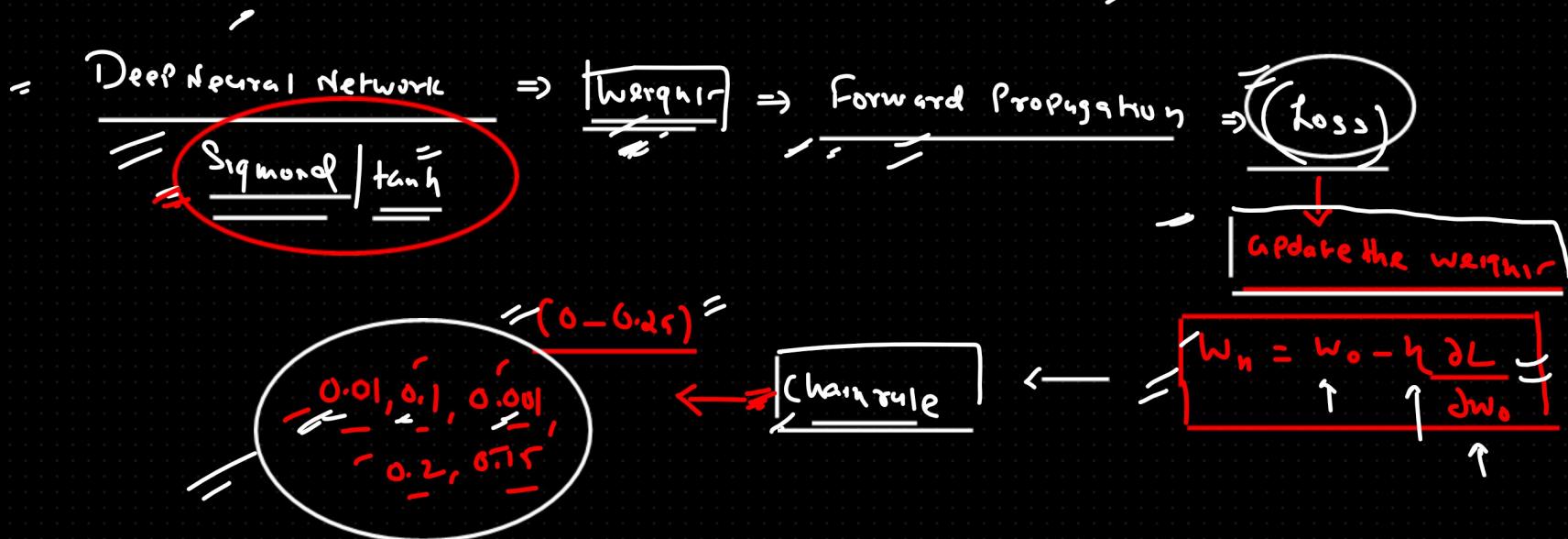


$$\begin{aligned}
 & \text{Sigmoid} \\
 & = [0, 1] \\
 & = \\
 & \text{B.P. derivation} \quad -[0, 0.25]
 \end{aligned}$$



$$\begin{aligned}
 & \text{tanh.} \\
 & = [0, 1]
 \end{aligned}$$

$$0.1 \times 0.1 \times 0.1 \times 0.1 = \underline{\underline{0.0001}}$$



Small \times Small \times Small ... = Very small value

$L \rightarrow \gamma \rightarrow w_2 \rightarrow o \rightarrow w_1$

$\frac{\partial L}{\partial w}$

$1 \quad 0.01 \quad 0.0001 \quad 0.9999$

$W_{new} = w_{old} - \eta \frac{\partial L}{\partial w}$

$= 1 - 0.000018$

$= 0.9999$ negligible

$0.9999 - 0.00001$ gradient value

which is going to be vanish

Vanishing grad.

→ Deep Neural Net



Loss will not change



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

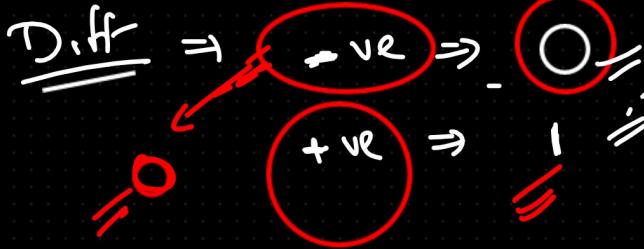
$$\max(0, z)$$

$$(-\infty, +\infty)$$

$$(0, +\infty)$$

{

It will remove negative value from your data }



Diving Relu ← 0 or 1

chain rule

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



L-ReLu

$$= \max(0, x) \leftarrow (-\infty, +\infty)$$

$$\max(0.1, x)$$

3

Proper weight initialization

$$\begin{array}{ll} 1 & \text{Xavier} \\ 2 & \text{Global} \end{array}$$

4

Batch Normalization

Weight init

Exploding

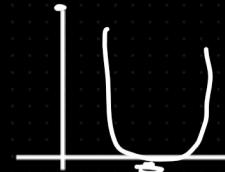
$$(2 \times 0.5 \times 1 \times 3 \times 5 =)$$

RMS

Derivation

$$\begin{array}{l} 0, 1 \\ \text{Chain rule} \end{array}$$

Chain rule



How to improve a Performance of Neural Network? -

1

Handle OR TC vanishing gradient

- ↳ Activation f_f
- ↳ Correct weight initializing

2

Overfitting

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

3

Normalization (fast training)

- ↳ Normalization / Std.
- ↳ Batch Normalization

Exploding grad.

9

Optimizer

- ↳ Momentum
- ↳ AdaGrad
- ↳ Adadelta | RMSProp
- ↳ Adam

5

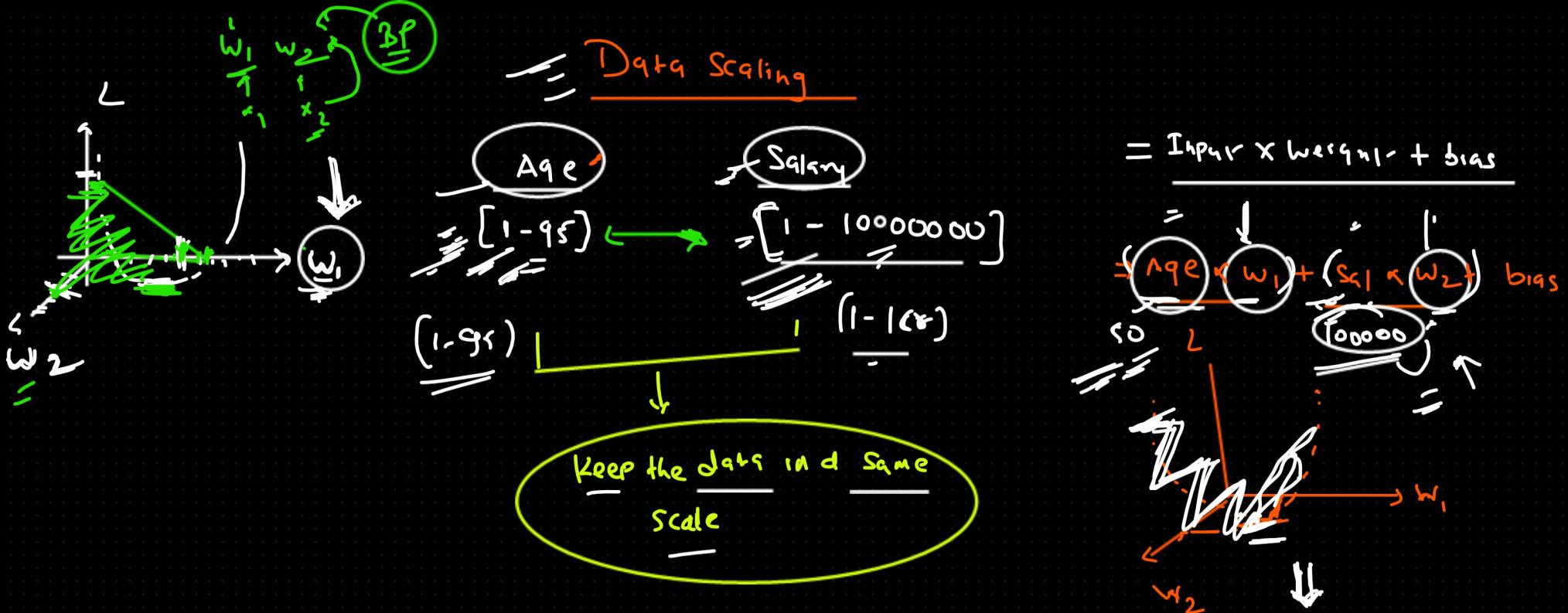
hyperParameter

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

6

Activation

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



1 Standardization

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

$x_i - \mu$

σ

$(-3, +3)$

Z_{table}

S/N

Salary

2 Normalization

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

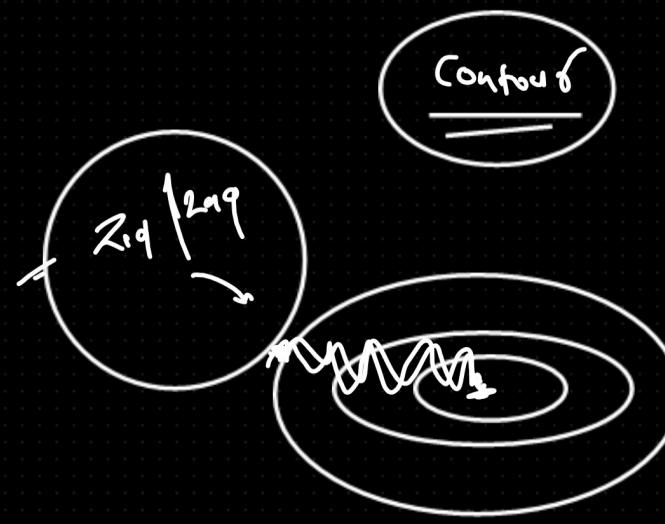
$x_i - x_{\min}$

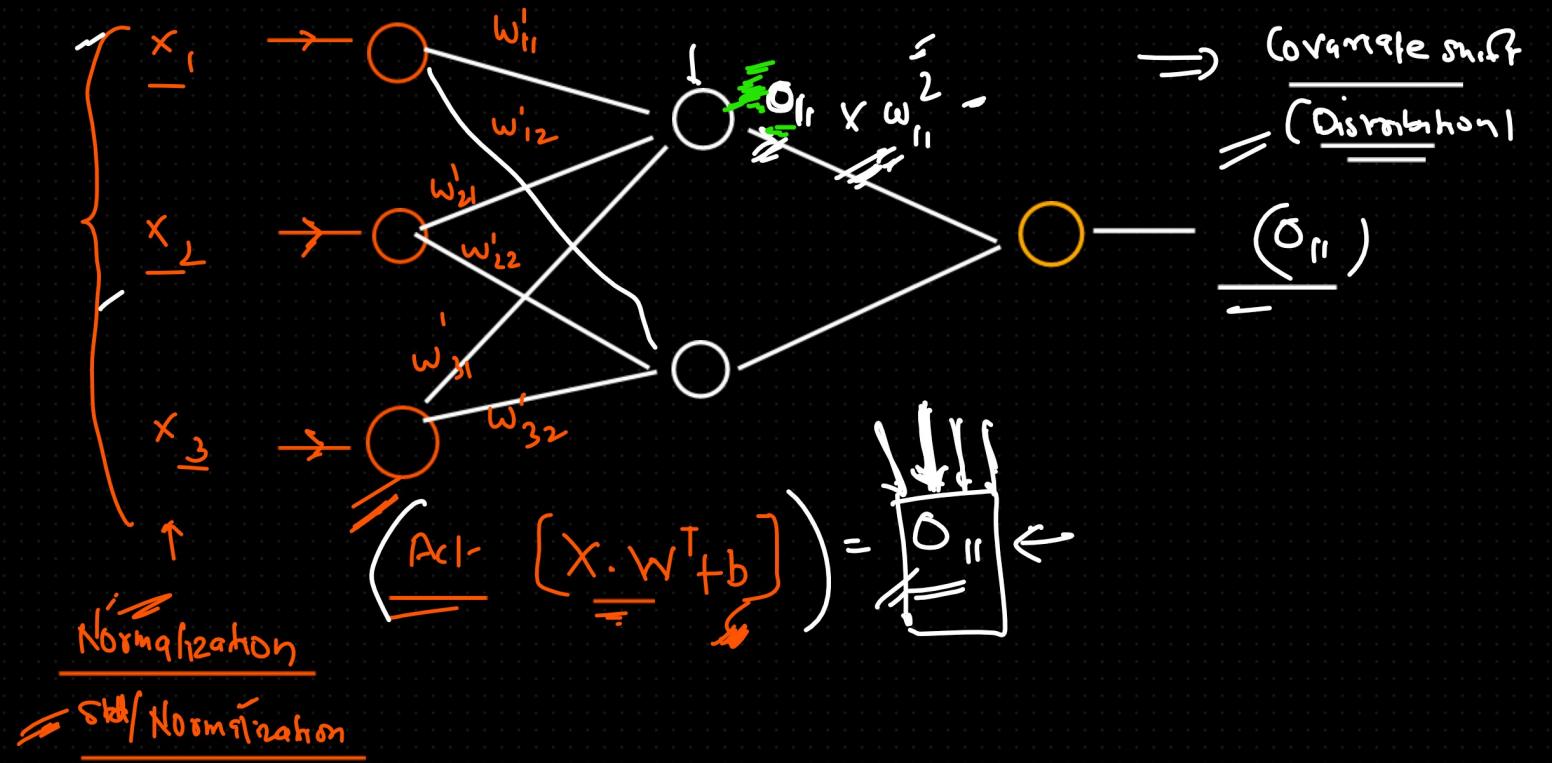
$x_{\max} - x_{\min}$

$[0-1]$

C GPA

$(0-10)$

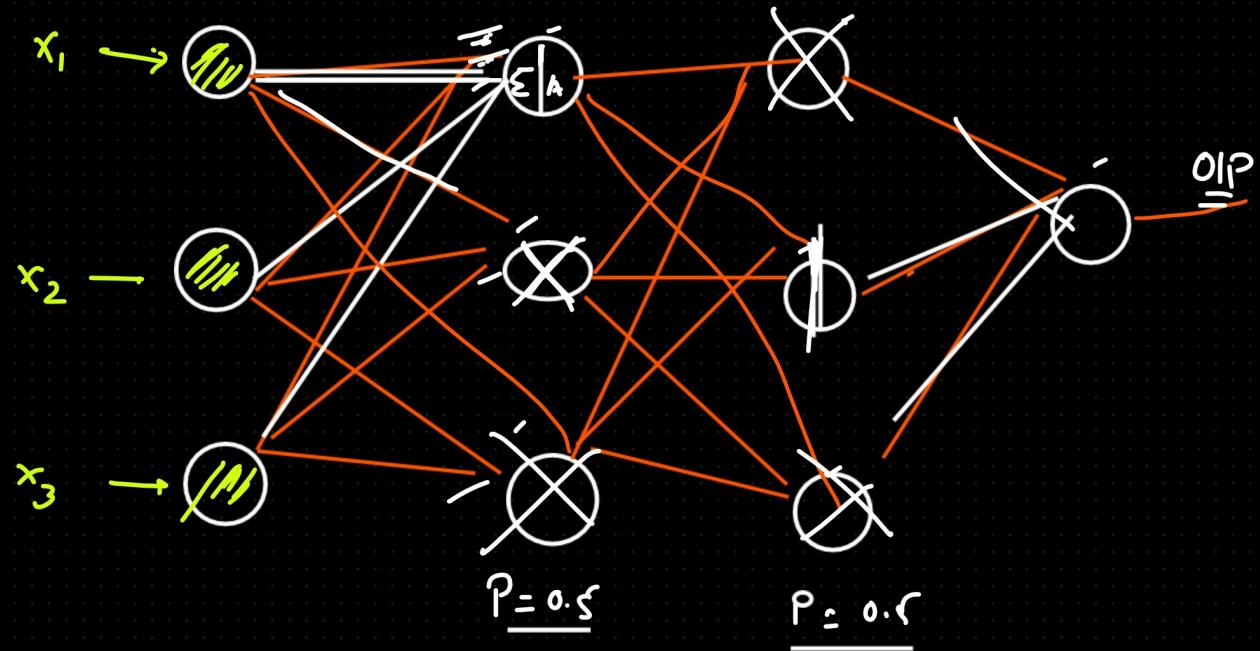




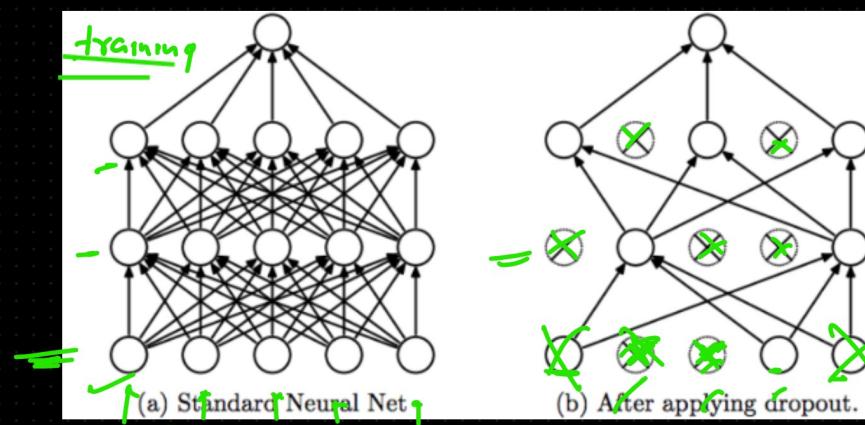
Drop out

Training

(1 epoch)



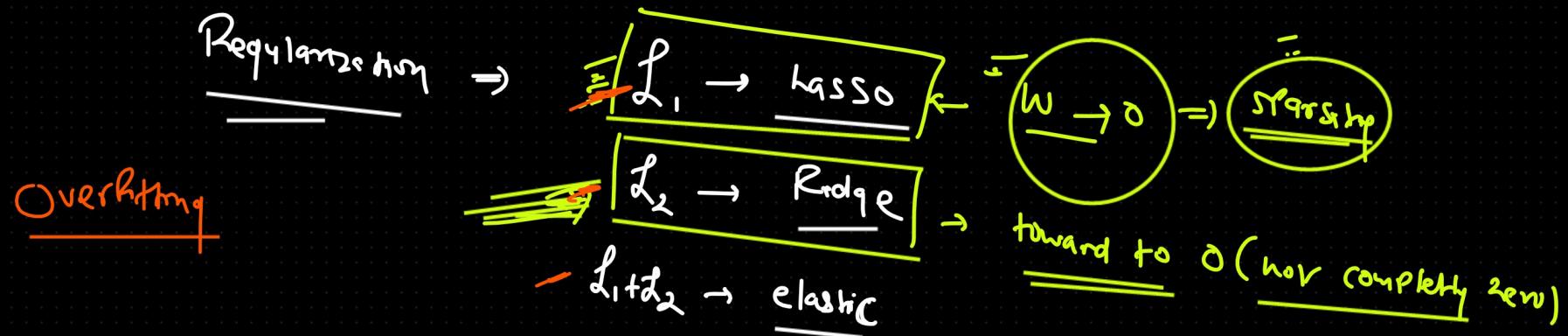
5 epoch



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

$\approx 50\%$
Shutdown

ff



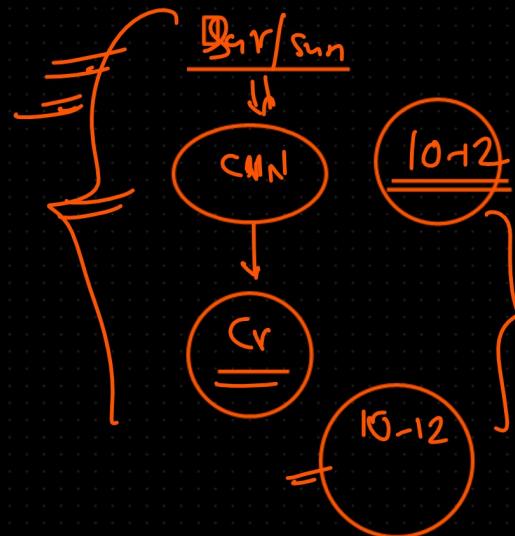
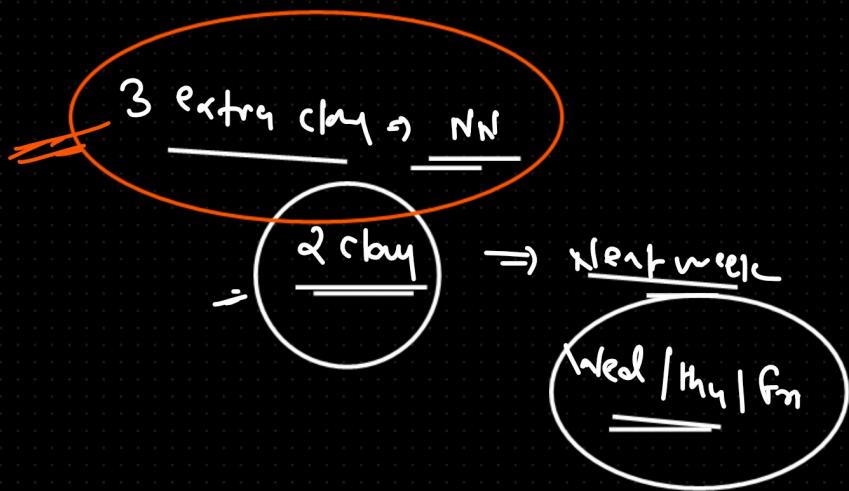
$$\text{Loss} + \frac{\text{Penalty}}{\mathcal{L}_1, \mathcal{L}_2 / \mathcal{L}_1 + \mathcal{L}_2}$$

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

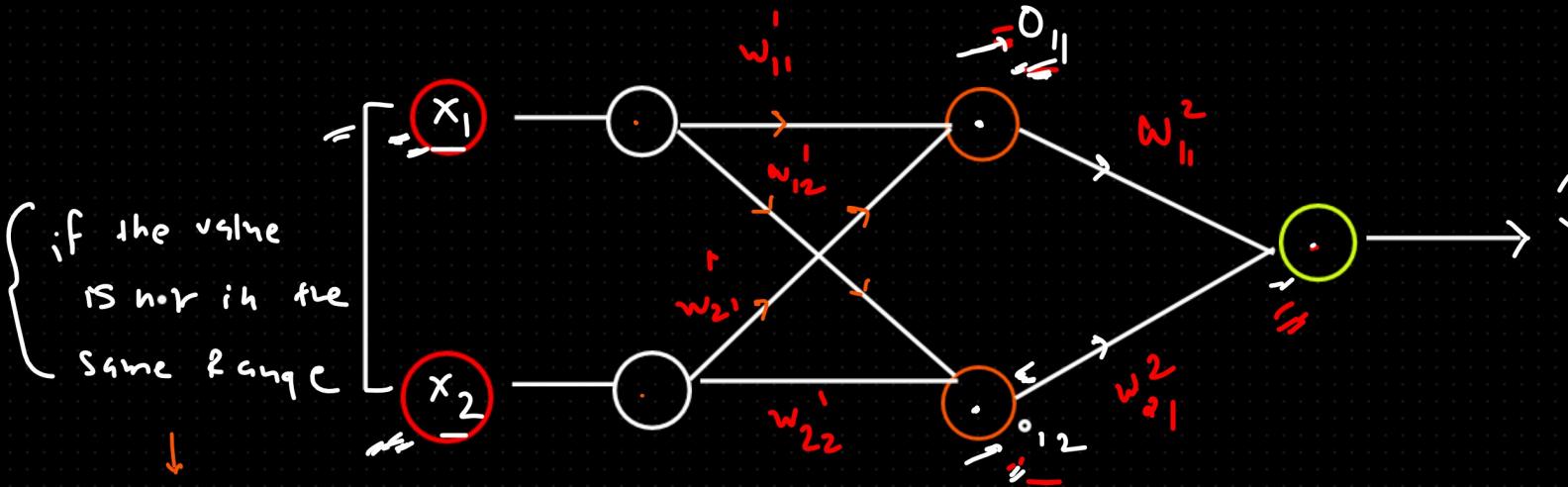
$$\mathcal{L}_1 + \mathcal{L}_2 \rightarrow \text{Loss} + \lambda(|w|) + \lambda(w^2)$$

Cost

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



Batch normalization

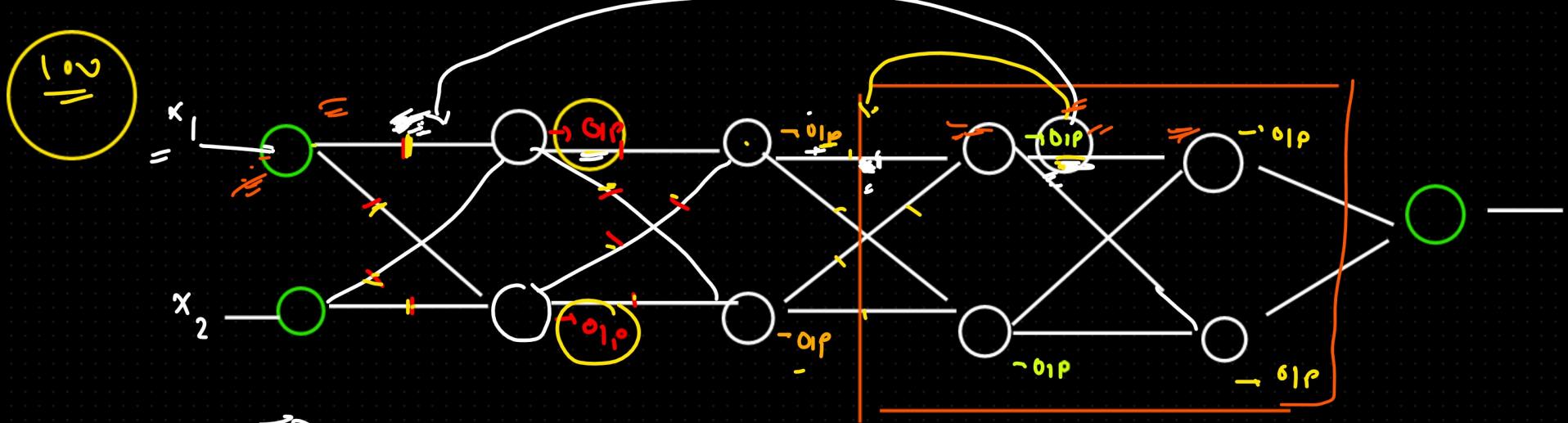


$$z = \frac{x_1 \times w_{11}' + x_2 \times w_{21}' + b_{11}}{}$$

$$\hat{o}_{11} = \text{Act}(z)$$

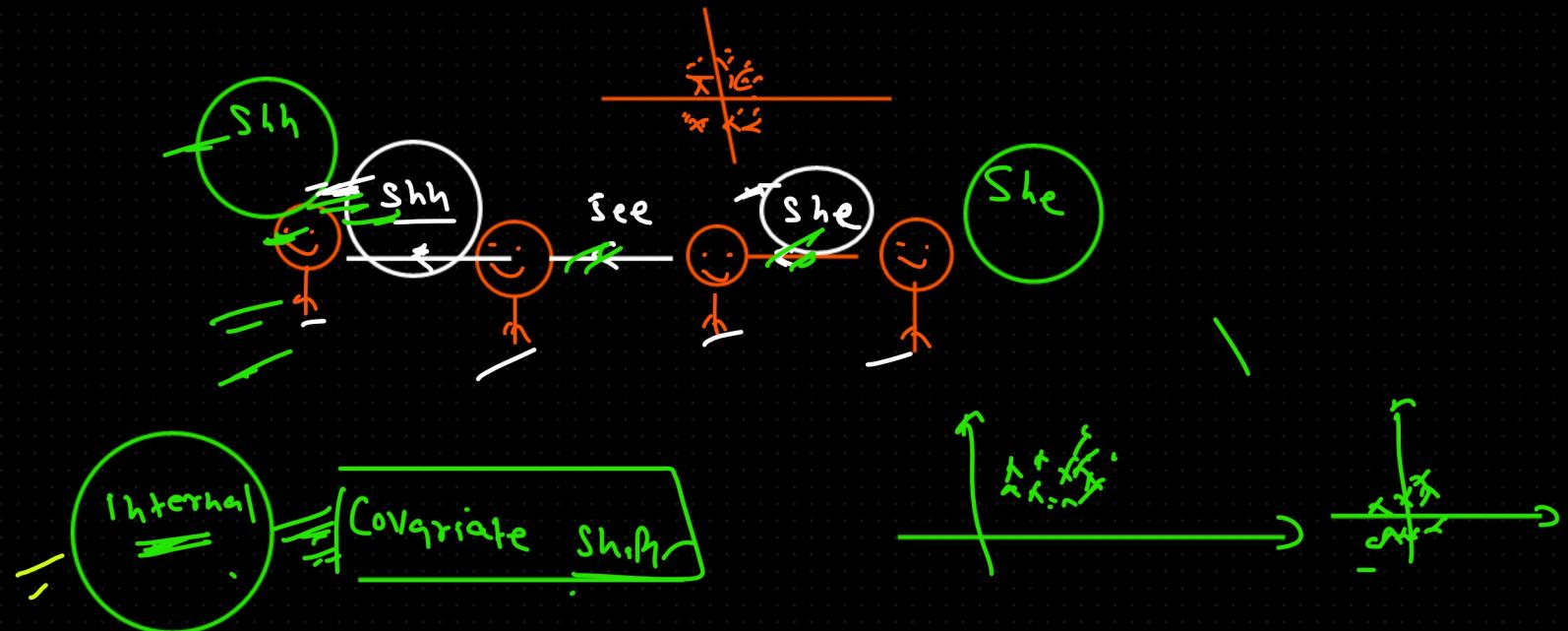
= Batch normalization → Batch Gradient descent

$\hat{o}_{11}, \hat{o}_{12} \rightarrow$ Near Layer

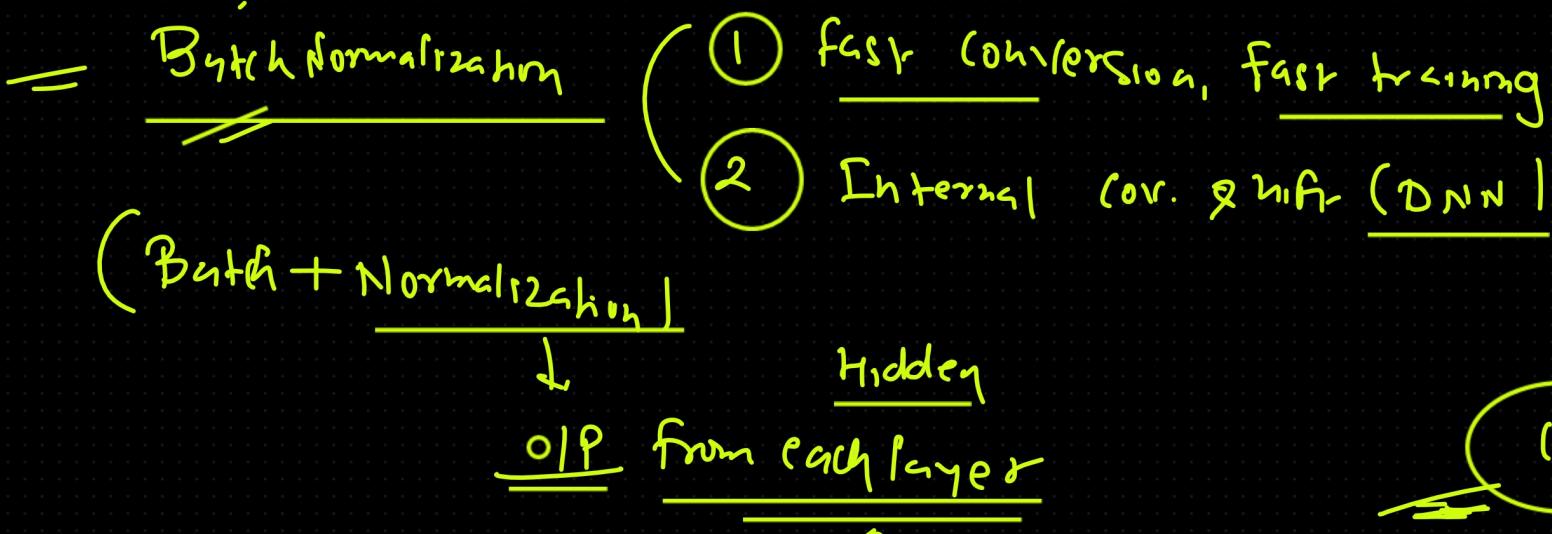


$=$ Batch Normalization
 $-$ Deep Layer

Deep NN



Normalization



$$\underline{\text{Data}} \rightarrow \underline{\text{Scaling}}$$

Scaling = $\underline{\text{ICR}}$

Shallow

$$= \underline{65\%} \underline{70\%}$$

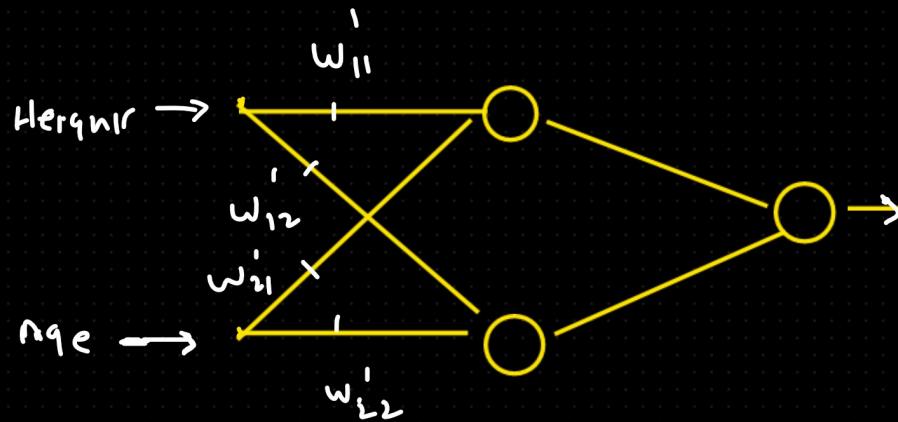
Scaling
Batch = $\underline{15 \text{ hr.}}$

$$= \underline{3.5 \text{ hr.}} \underline{72\%}$$

= Increase \rightarrow Batch \rightarrow 15 hr. \rightarrow 85%.

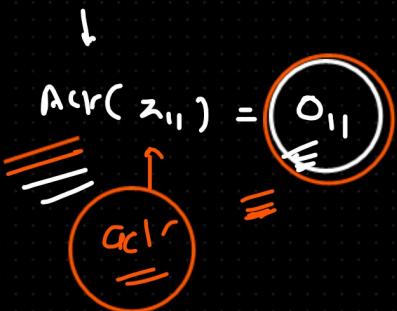
4.5

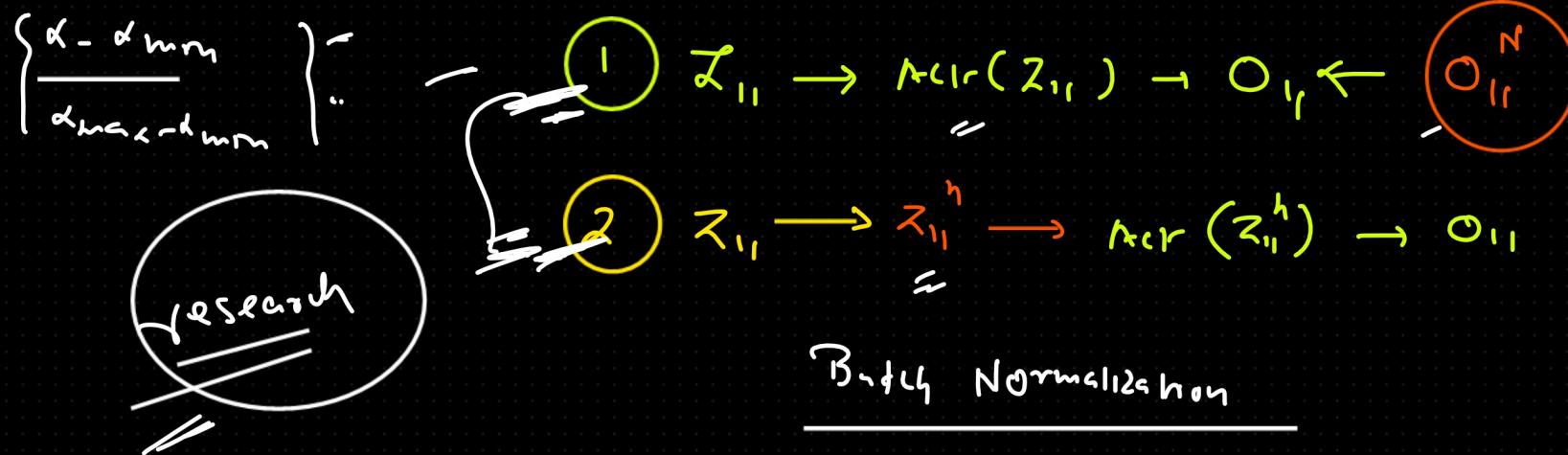
90%



(1) $\underline{z}_{11} = w_{11}^T \alpha_{\text{hergnir}} + nqe \cdot w_{12}^T$

$$\left\{ \frac{\text{std}}{o_{11} - m} \right\}$$

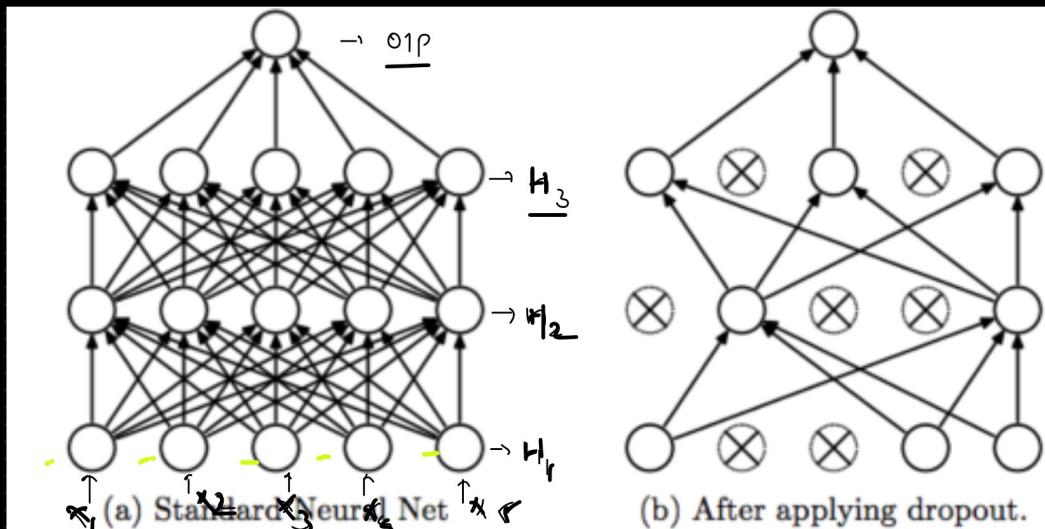




Dropout \Rightarrow Recent research in the field of NN

(2014)

Concept



$\{1 \text{ epoch} \Rightarrow \text{forward } (O1P)$
 +
 Backward
 $= (\text{Parameter})$

Shutdown, Shutdown (\equiv nothing)
 \equiv Randomly

~~↓~~

- ~~○ ○~~ \Rightarrow Leprosy

- $\rightarrow \circ \circ$ Leprosy

- ~~○ X~~ \Rightarrow What is beneficial?

- $\rightarrow \circ X$

- $\circ \circ$ Dropout?

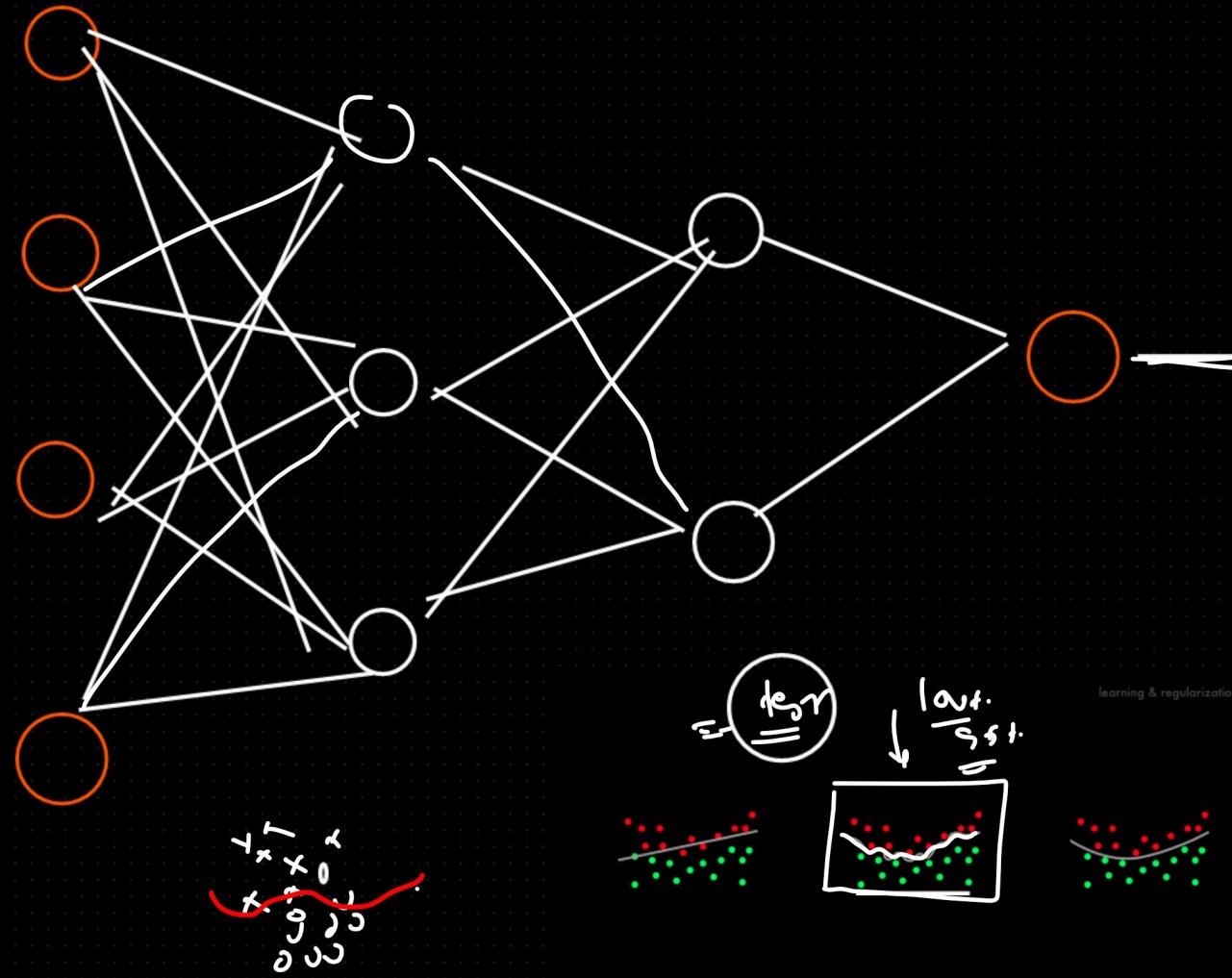
Shutdown \rightarrow Training
(deactivation)

Testing \rightarrow New neurons will be activated

P = 0.1 50%
Shutdown

Deep ≈
fully connected

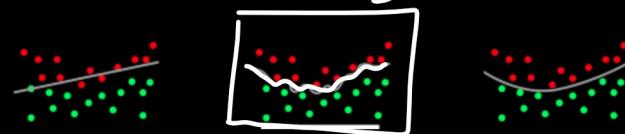
Precise



= for

lavr.
sst.

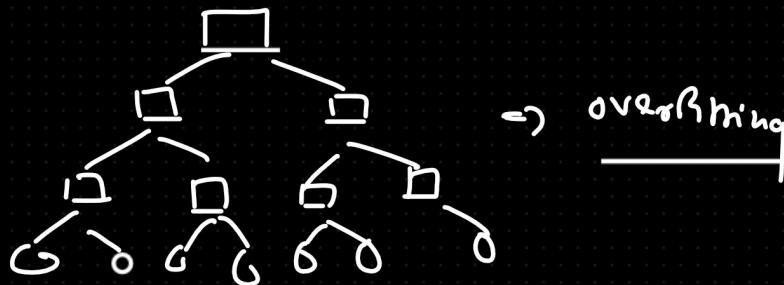
learning & regularization

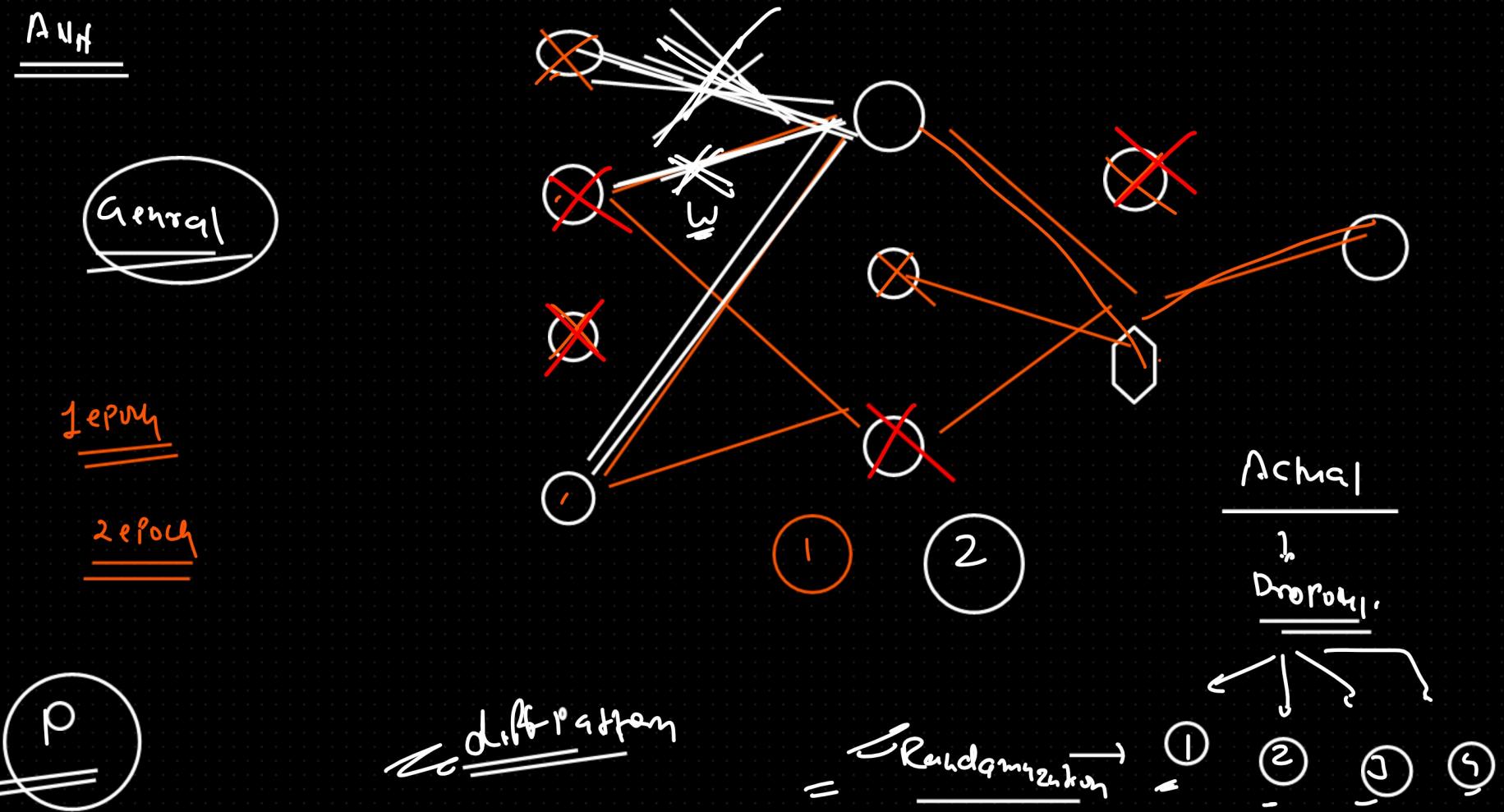
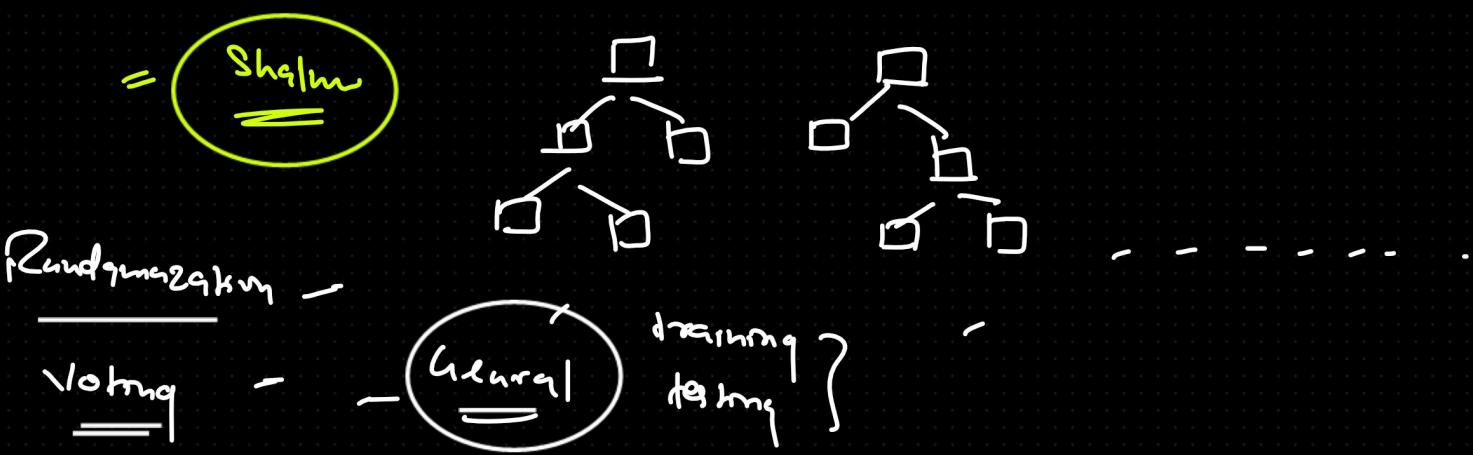


Decision tree → Random forest

(Bootstrap)

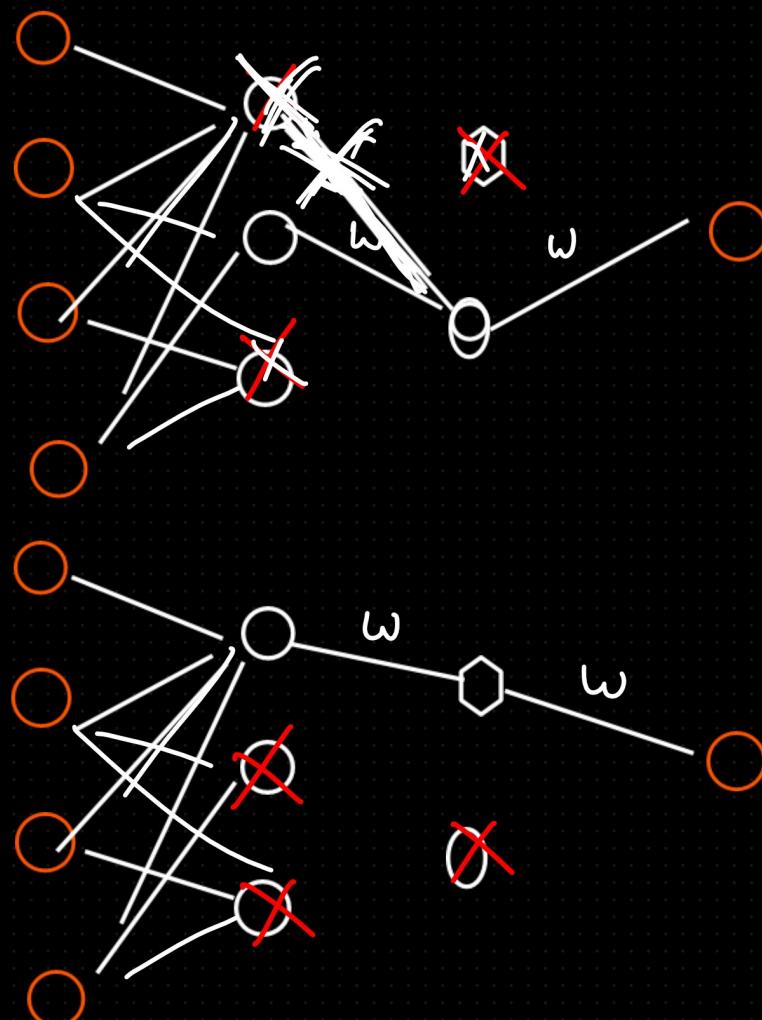
- (R+C)



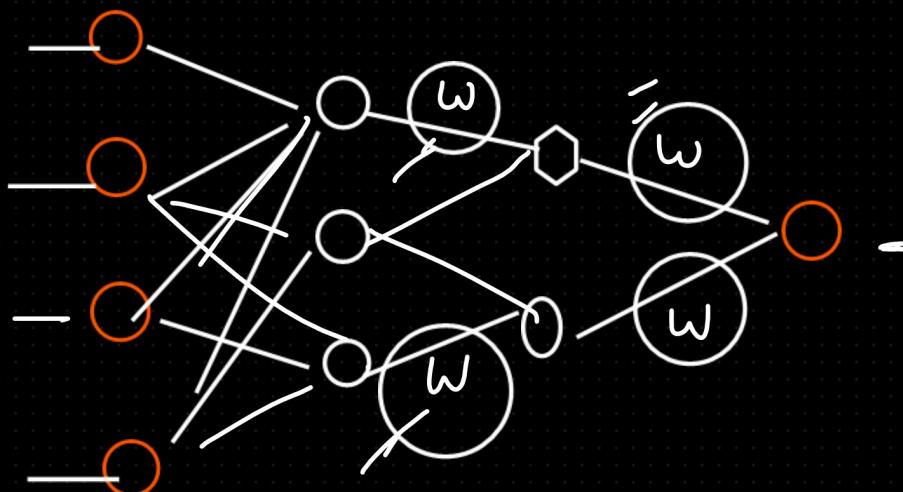


$$\cancel{P_c = \frac{0.25}{\text{---}} = \text{deactivation}}$$

$$\cancel{\underline{1-P} = \frac{0.75}{\text{---}} = \text{activation}}$$

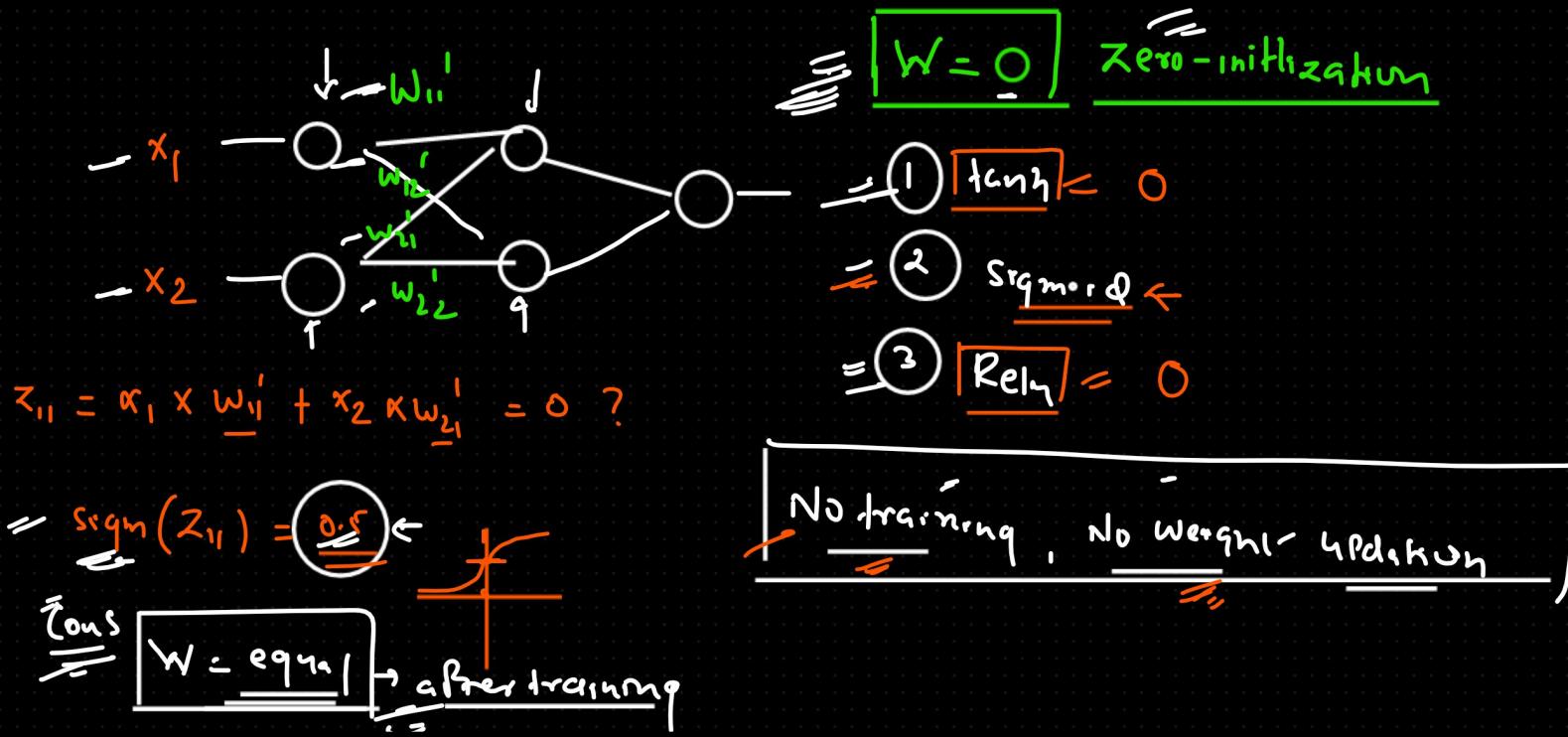


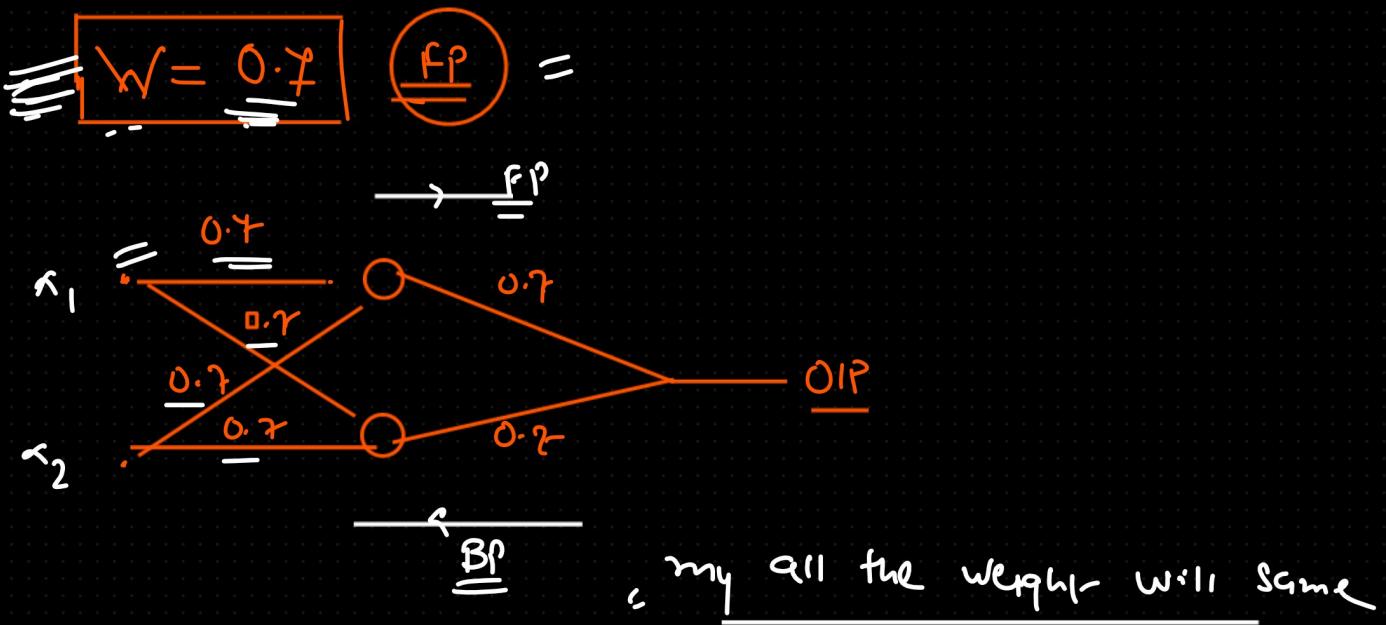
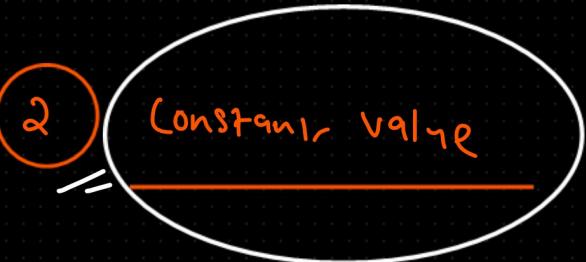
testing



Weight Initialization

- 1 Zero initialization
- 2 Initialization with constant value :=
- 3 Random initialization
- 4 Xavier / glorot - init
- 5 He init

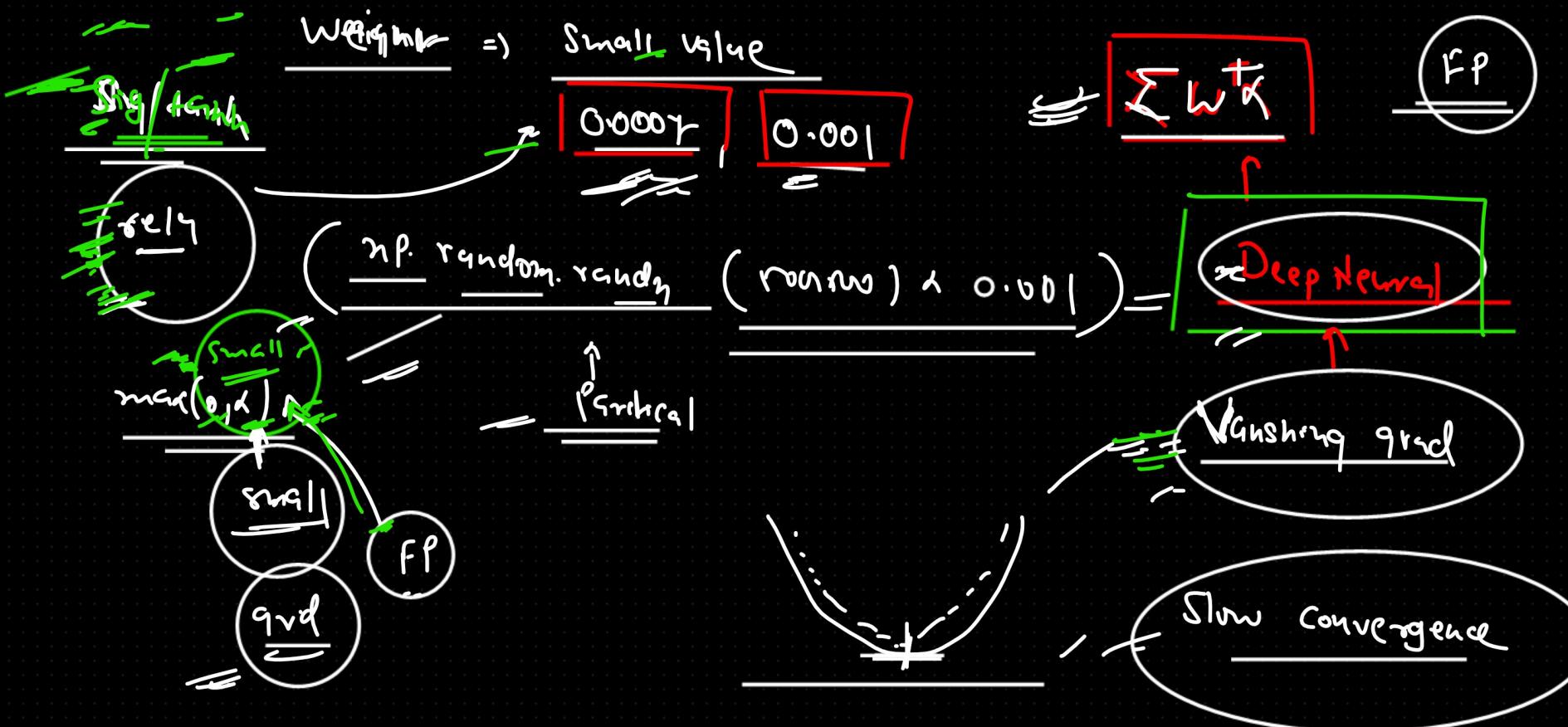




3 Random initialization

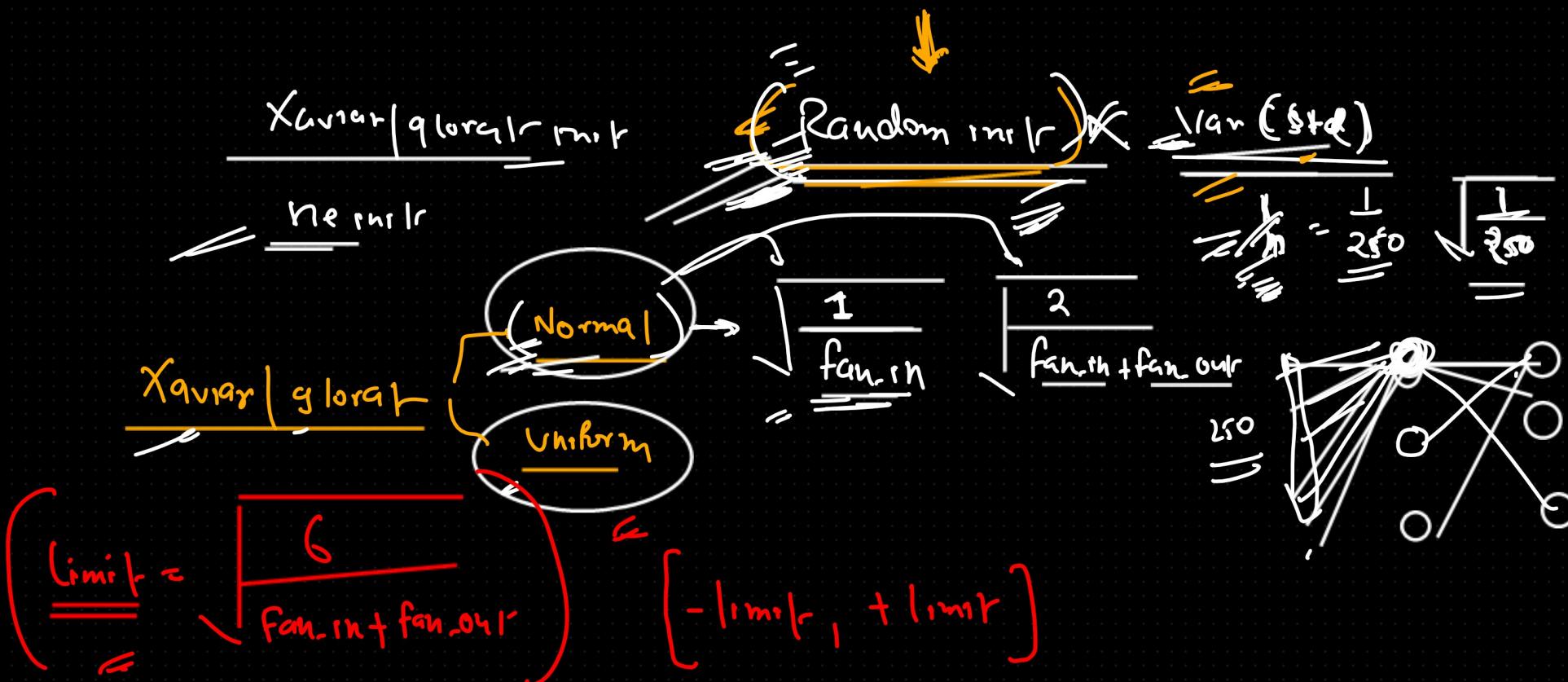
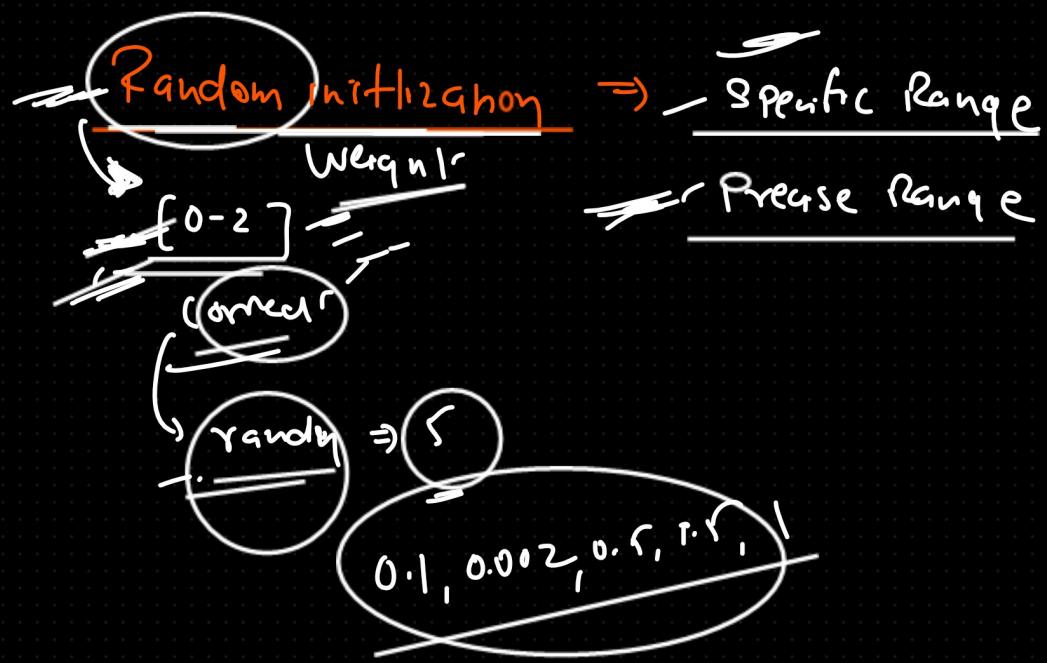
- ↳ 1 Initialization with small value
- ↳ 2 Initialization with larger value

tanh, sigmoid, relu



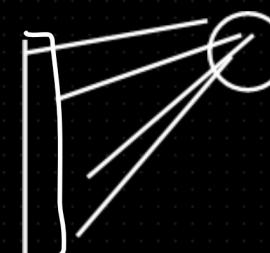
Large weight initialization





He init → Uniform → $\left[-\text{limit}_{\text{g}} + \text{limit}_{\text{t}} \right] \rightarrow \frac{6}{\text{fan-in}}$

Normal → Random init \times n.p random $\rightarrow \frac{2}{\text{fan-in}}$



1 Zero initialization

2 Initialization with constant

3 Random initialization $\begin{cases} \text{Small} \\ \text{Large} \end{cases} \rightarrow ? \times$ (we don't the specific Range)

4 Xavier / Glorot init

my weight can misslead to me

↪ Normal (NP-random) $\times \sqrt{\frac{1}{\text{f-in}}} \text{ or } \sqrt{\frac{2}{\text{f-in} + \text{f-out}}}$

↪ Uniform $[-\text{limit}, +\text{limit}] \Rightarrow \text{limit} \left\{ \sqrt{\frac{6}{\text{f-in} \text{ four}}}, \sqrt{\frac{6}{\text{f-out} \text{ four}}} \right\}$

5 he init

↪ Normal (NP-random) $\times \sqrt{\frac{2}{\text{f-in}}}$

↪ Uniform $[-\text{limit}_{\text{t}}, +\text{limit}_{\text{t}}] \rightarrow \left[-\sqrt{\frac{6}{\text{f-in}}}, +\sqrt{\frac{6}{\text{f-in}}} \right]$

= Random initialization

$\equiv []$

~~$\left(\frac{1}{h} \text{ or } \sqrt{\frac{1}{n}} \right)$~~

(var) (std)

Normal $\propto \frac{1}{\sqrt{n}}$

$h = \frac{\text{input to the node}}$

