

Agenda

1) Batch Normalization .

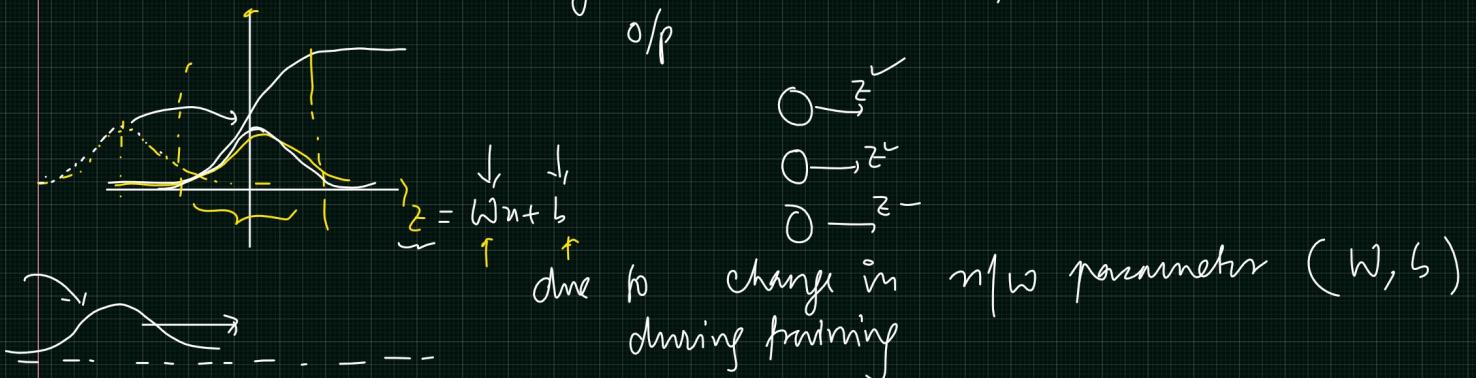
2) Transfer learning .

still NN can suffer vanishing & exploding gradient

due to ,

ICS : Internal Covariate shift

L, Change in distribution of network activation



out "f", wt. init,

→ It has been experimentally proven :-

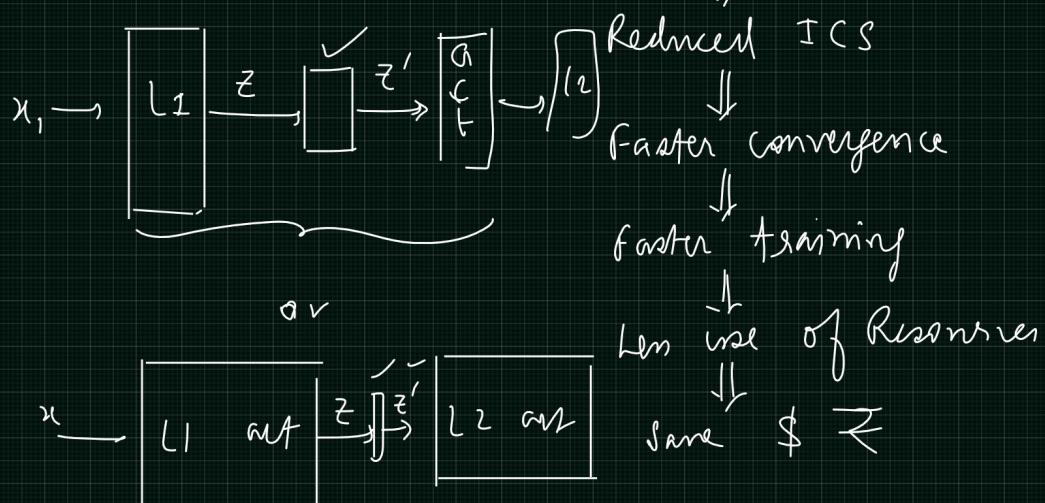
{ LeCun et.al. 1998 & Wiesler & Ning 2011 }

The network converges faster if inputs to
the layer are
whitened / normalized

linearly transformed to zero mean } std.
& unit variance } Normal
distribution

→ Expectation :

Fix the distribution for each layer { Dynamically }



Batch Normalization { BN }

↳ In 2015, Sergey Ioffe & Christian Szegedy

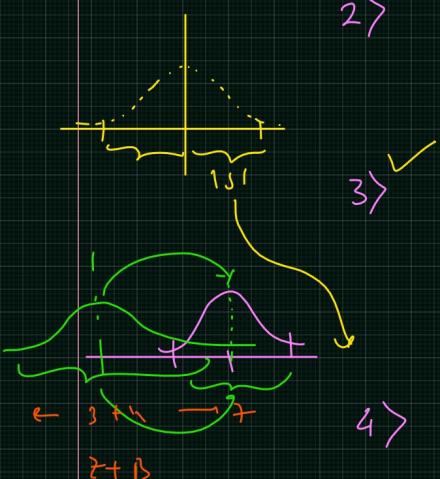
↳ proposed an extra set of operation to perform before or after activation layer.

Algorithm : - { While Training }

$$1) \quad M_B = \frac{1}{m_B} \sum_{i=1}^{m_B} x^{(i)} \quad \begin{matrix} \xrightarrow{x^{(i)}, z^{(i)}} \\ \leftarrow \text{batch mean} \end{matrix}$$

$m_B \rightarrow \text{batch size}$

$$2) \quad \sigma_B^2 = \frac{1}{m_B} \sum_{i=1}^{m_B} (x^{(i)} - M_B)^2 \quad \leftarrow \text{Batch Variance}$$

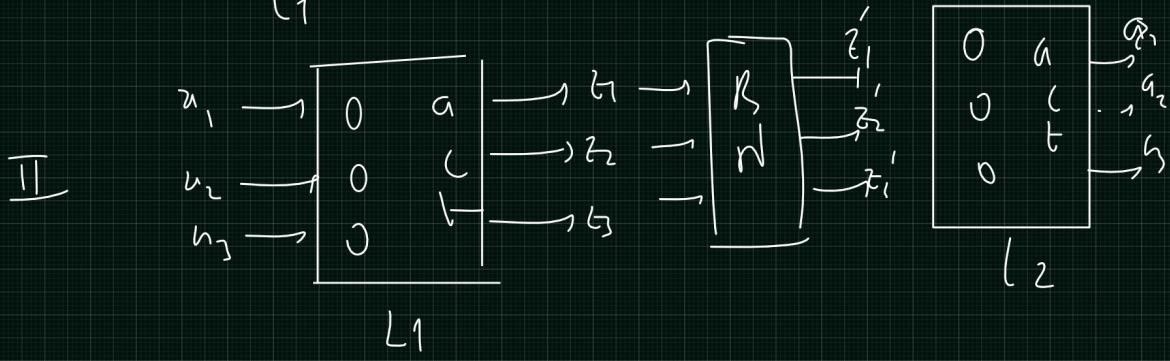
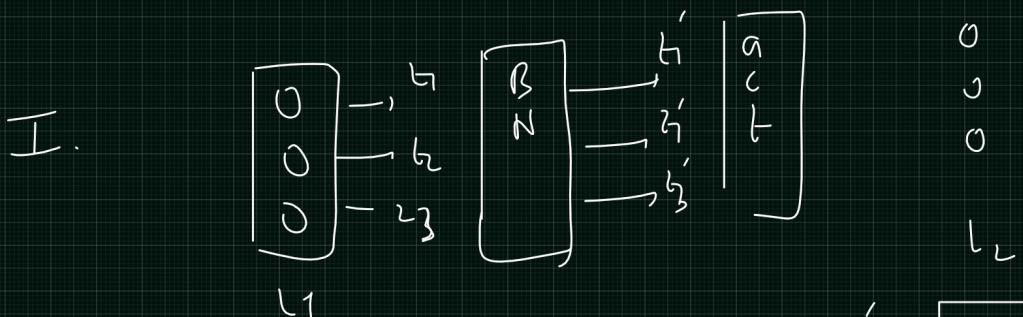


$$3) \quad \hat{x}^{(i)} = \frac{x^{(i)} - M_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \leftarrow \text{Normalizing } \hat{x}^{(i)}$$

$\epsilon = 10^{-7}$ } smoothing term } to avoid zero division error

$$4) \quad \hat{z}^{(i)} = \gamma \otimes \hat{x}^{(i)} + \beta \quad \leftarrow \text{Scaling } \hat{z}^{(i)} \quad \leftarrow \text{Shifting } \hat{z}^{(i)} \quad \leftarrow \text{Learnable parameter}$$

$$\boxed{\text{layer 1}} \xrightarrow{z} \boxed{\beta N}$$



NN (ω, b) $\frac{\partial e}{\partial \omega}, \frac{\partial e}{\partial b}$ $\omega = \omega - \eta \frac{\partial e}{\partial \omega}$	βN (γ, β) $\frac{\partial e}{\partial \gamma}, \frac{\partial e}{\partial \beta}$ $\gamma = \gamma - \eta \frac{\partial e}{\partial \gamma}$
---	--

After training

μ is calculated from all μ_B
 σ^2 is calculated from all σ_B^2

while prediction

$$\left\{ \begin{array}{l} (i) \quad \hat{x}^{(i)} = \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ (ii) \quad z^{(i)} = \gamma \otimes \hat{x}^{(i)} + \beta \end{array} \right.$$

More Observation:-

→ γ & β → 2 new trainable parameters

→ μ & σ^2 → 2 additional parameters (non trainable)
They are calculated using moving average internally.

total parameters = 4 → $\gamma, \beta, \mu, \sigma$

trainable = 2 → γ, β

non trainable = 2 → μ, σ

→ Now Cost function is :-

$$C(\underbrace{w, b, \gamma, \beta})$$

↓
updated by backpropagation

→ Advantages :-

i) You don't need scaling of data if you are using BN as a first layer.

ii) Although it adds 2 extra parameters to train still it converges to the solution faster

Without BN With BN

90%, 90%

100 steps 70 steps

iii) It solves unstable gradient problem to some extent
 get affected by

iv) It doesn't depends on choice of activation fn and wt. initialization.

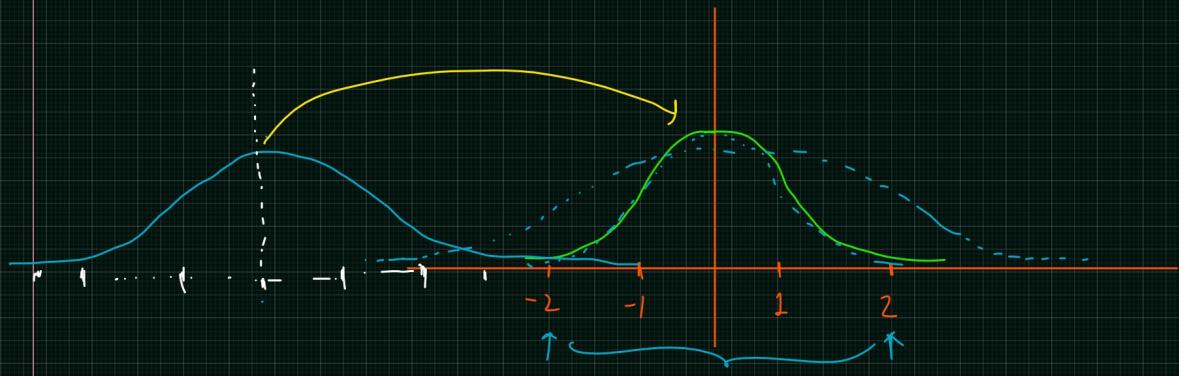
→ Disadvantage :-

i) It increases the complexity of nn.

ii) No. of trainable parameter increased.

iii) Complexity \Rightarrow runtime penalty \Rightarrow more prediction.

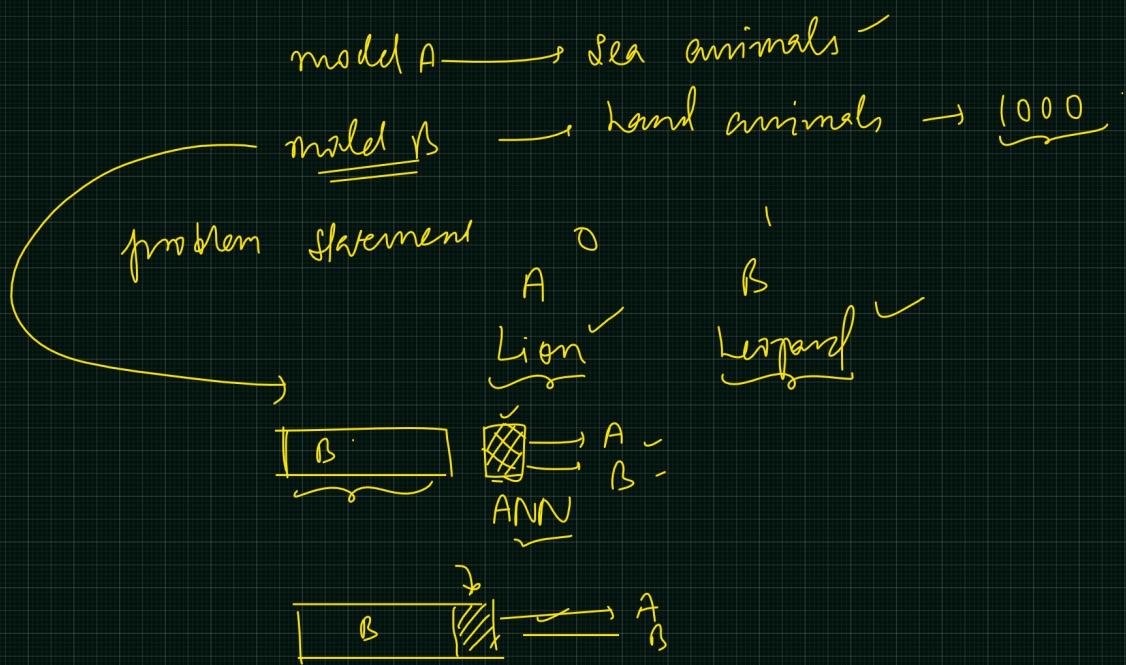
w> Training time increased but convergence is fast



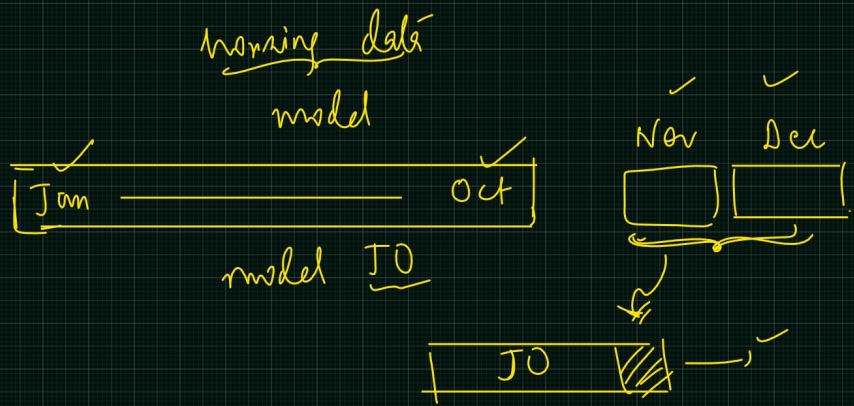
Transfer Learning :-



pre trained model



Repression



19

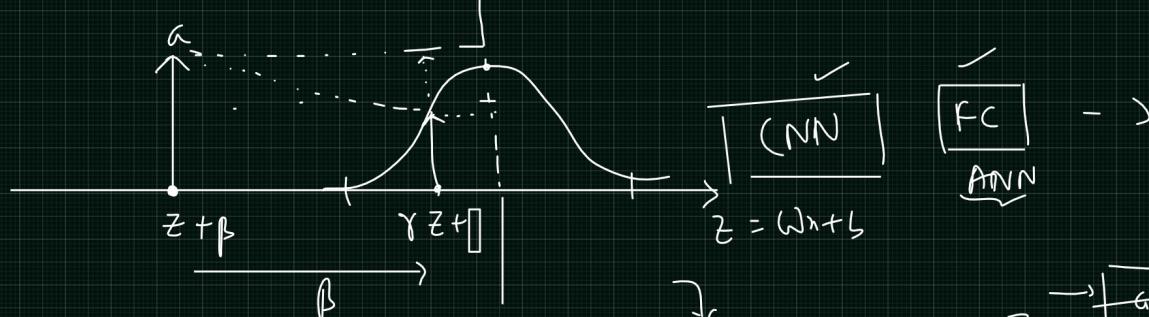
$$\begin{array}{c}
 \text{19th} \rightarrow \text{cur} \\
 \text{18m} \rightarrow 75\% \\
 \text{19m} \rightarrow 85\%
 \end{array}$$

$$\begin{array}{c}
 \text{19th} \rightarrow 90\% \\
 \text{18m} \rightarrow 9
 \end{array}$$

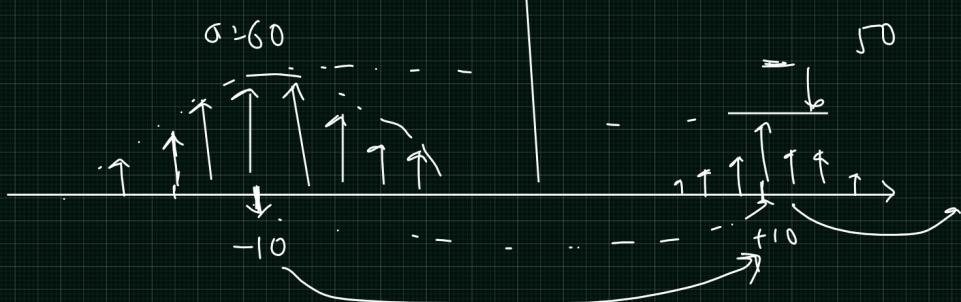
NLP

$$\begin{array}{c}
 \text{CV} \\
 \text{RNN} \\
 0
 \end{array}$$

$$\sigma(z) = A$$

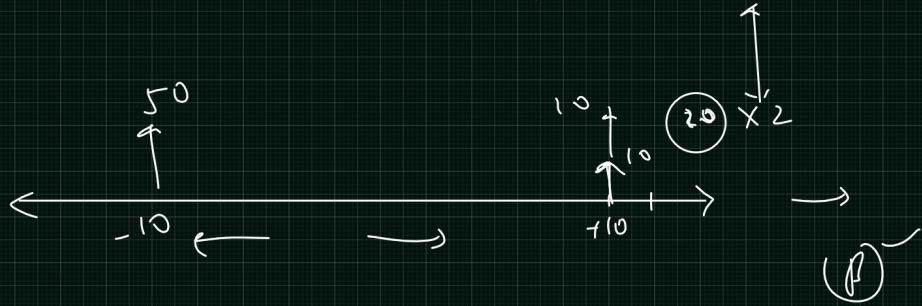


$$\begin{array}{c}
 z \rightarrow g \\
 -1 \quad 1
 \end{array}$$

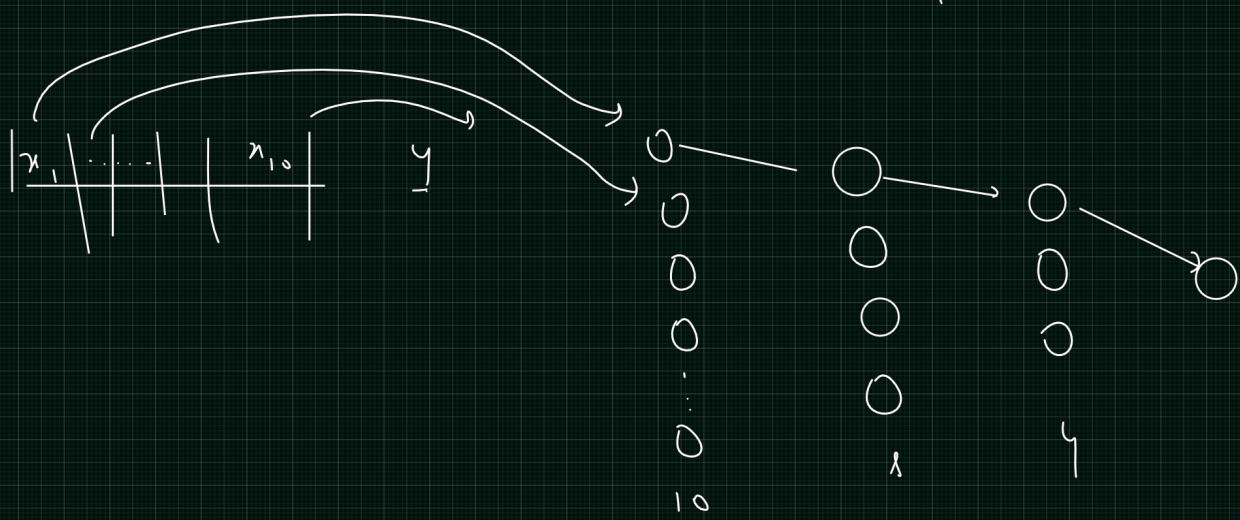
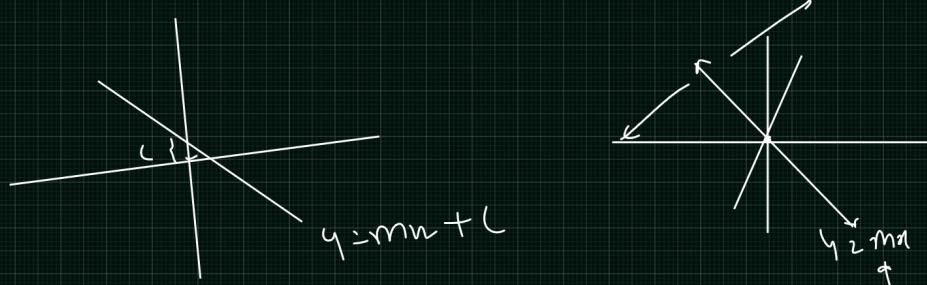


$$\begin{array}{c}
 z \rightarrow L \\
 -1 \quad 1 \\
 \omega x + b
 \end{array}$$

$$\sigma \left(\frac{1}{z} \right) = a$$



\rightarrow



$$\text{ReLU} = \max(z, 0)$$

$$= \begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$$

