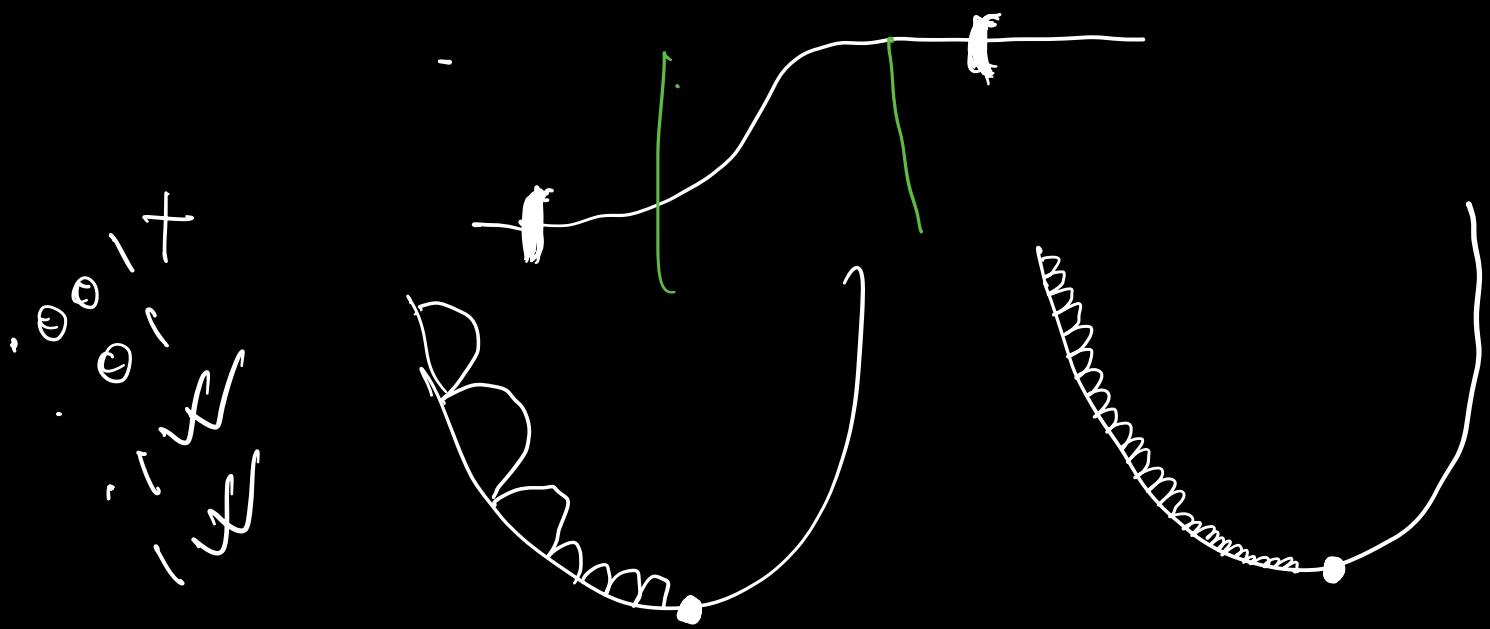
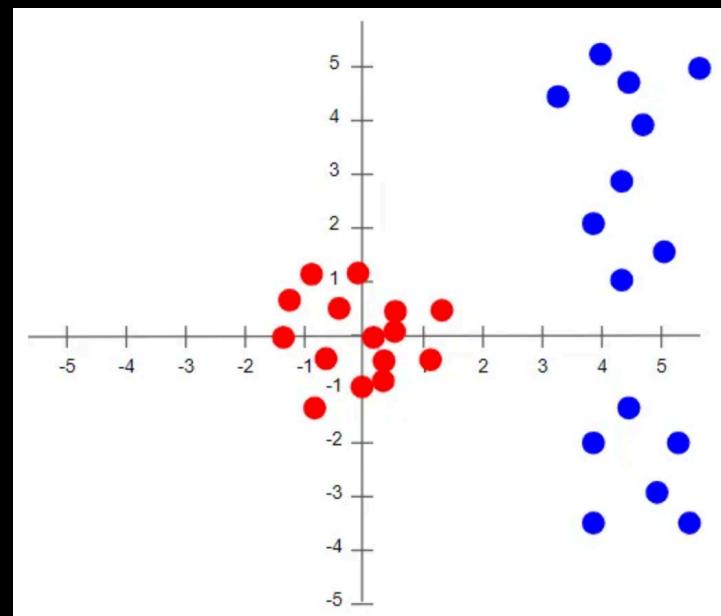
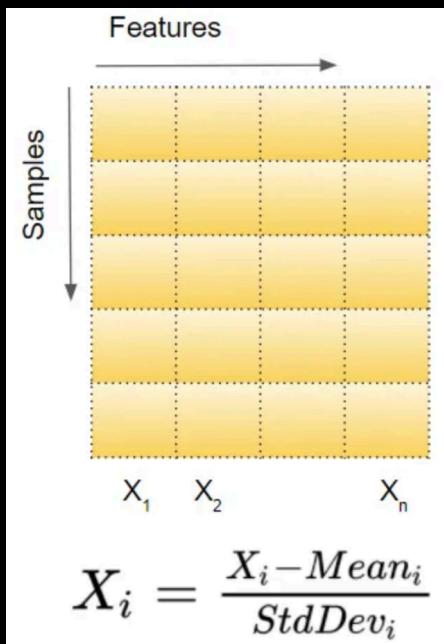


# Today's Agenda

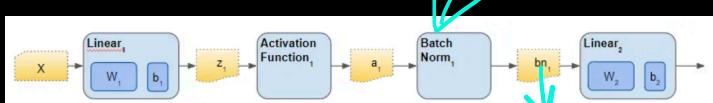
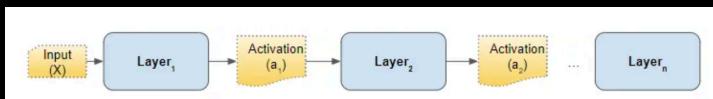
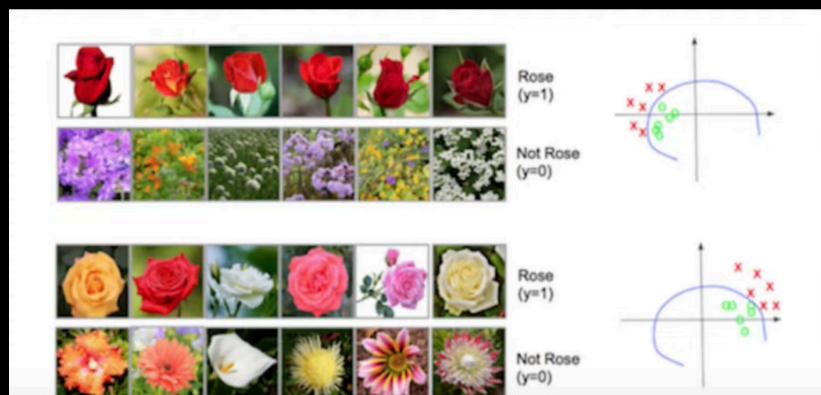
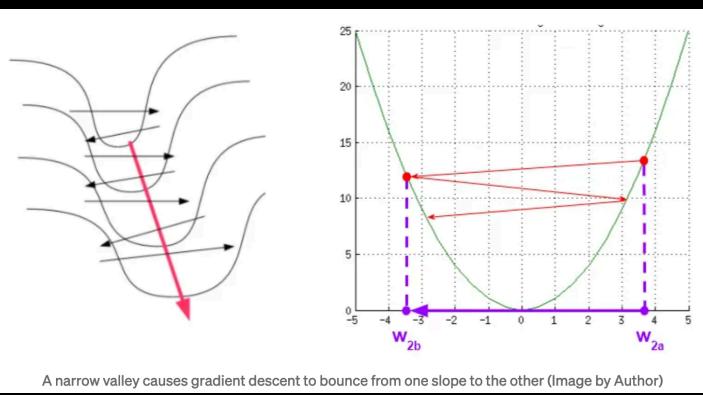
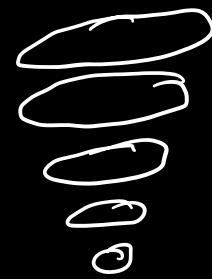
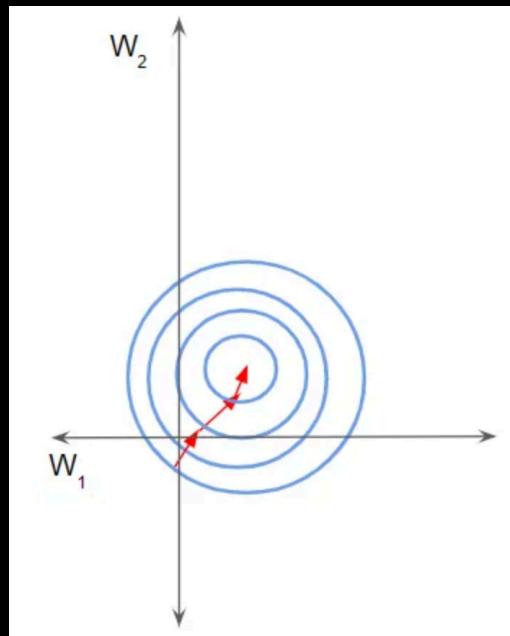
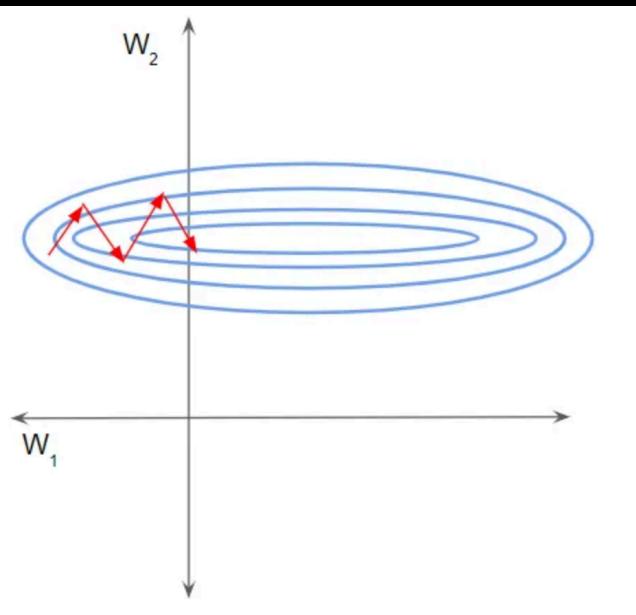
1) Batch Normalization

2) Weight Initialization



Without BN CF Curve

BN CF Curve



$\frac{N^2}{\text{Less Epochs}}$   $\frac{N^1}{\text{More Epochs}}$   $\frac{\text{Less Time}}{\text{More Time}}$   $\frac{100}{98\%}$ .  $\frac{\text{No Epochs}}{\text{With BN}}$

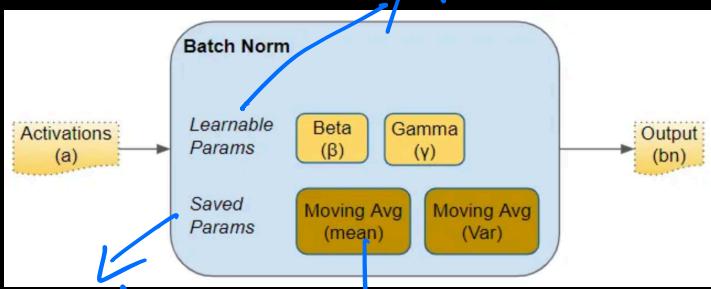
Normalized Output

$$\text{act}(\omega x + b) = a''$$

Terminable

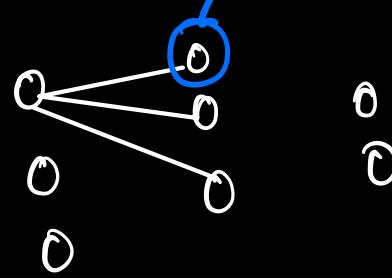
$$\text{Normal}(a'') \rightarrow \text{BN}(a'')$$

$$= \text{act}(\omega x + b)$$

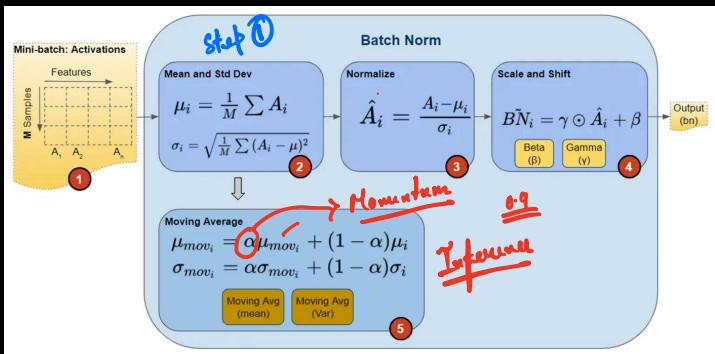


Now Derivable by Parameters  
Inference

$\beta$  &  $\gamma$  Exponential Weighted



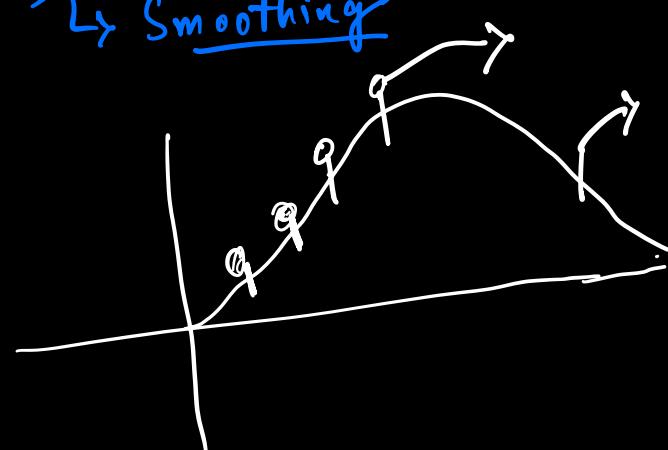
In a layer for every node activation is calculated and then BN is applied



Moving Avg

(EWMA)

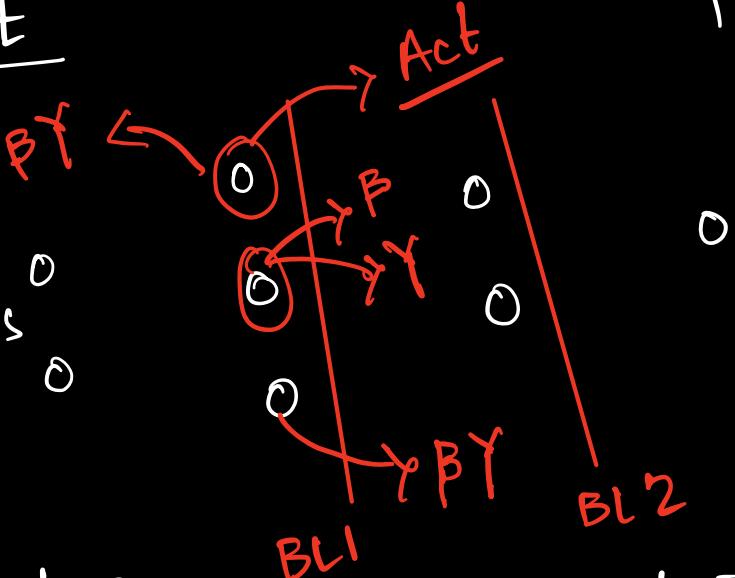
Smoothed



$\gamma \rightarrow$  Scaling factor

$\beta \rightarrow$  shift

Total Parameters  
20



Learnable Parameters

$$5 \times 4 = 20$$

Activations  $\rightarrow$

Mean  $\rightarrow$  Std-Dev  $\rightarrow$  Normalization

$\downarrow$

B.N  
(Scale & Shift)

$\{z(wx+b)\}$

Case 1

Activation  $\rightarrow$  Batch Normalization

Case 2

$wx + b \rightarrow$  Batch Normalization  $\rightarrow$  Activation  
MBxD

Feed Forward

$$\begin{array}{ll} \checkmark \quad Y = 1_x \\ \checkmark \quad \beta = 0_x \end{array}$$

$$\frac{\partial L}{\partial Y}$$

$$\begin{array}{r} 100 \\ \hline 4 - \underline{\underline{25}} \end{array}$$

Back prop

Weights, Bias,  $\beta$ ,  $Y$

Moving Average Mean / Variance

$\rightarrow$  Saved during Training

100 data points

Batch size  $\rightarrow 4$

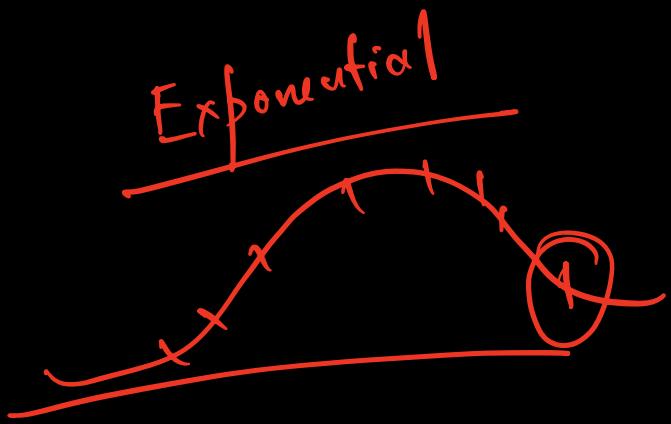
Iterations  $\rightarrow$  25

$\Sigma_1$  Moving Average & Mean

$\Sigma_1$  Moving Avg Std. Dev

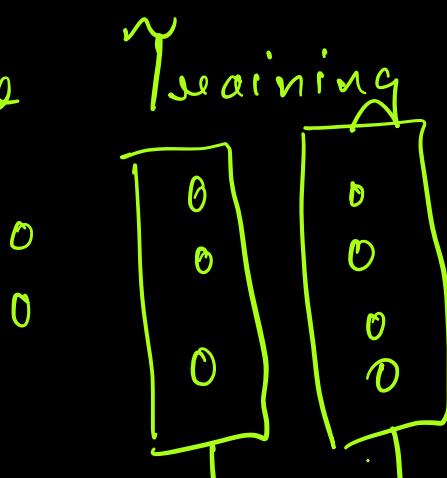
⋮  
⋮  
⋮

$\Sigma_{25}$  NAM  
 $\Sigma_{25}$  MAV



### Pros

- 1) Faster Training Time
- 2) Regularized
- 3) Not to worry about weight initialization
- 4) Stable Training



N.D N.D

$\mu_1$

$\sigma_1$

:

$\mu_{2S}$

$\sigma_{2S}$

Influencing

ESD

EWMA R

Smoothing

\*

$B_{\text{avg}} \times V_{\text{avg}}$

$(1-\beta)^n$

Input Data

skip standard

MNIST  
Dataset

$\frac{1}{255}$

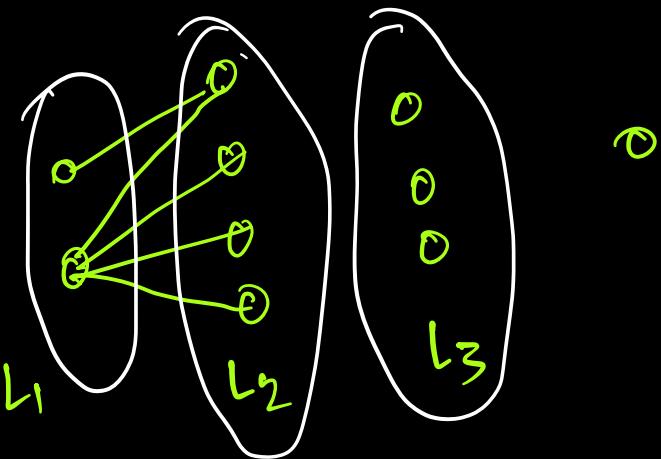
Input Layer

Batch Normalization

Data

Prelparation

Modelling



Different Distributions  
 ↳ keeps on changing  
 Training slow } bit unstable

- | <u>Weight</u> | <u>Initialization</u> |
|---------------|-----------------------|
| 1) Zero       | Initialization X      |
| 2) Non - Zero | Initialization X      |
| 3) Random     | Initialization        |

$$Wx + b$$

$$\Rightarrow 0 \cdot x + 0$$

$$\Rightarrow 0$$

## Intuition

$$z_{ii} = w_i x + b$$

$$a_{ii} = \max(0, z_{ii}) \rightarrow$$

$$w_{ii}' = w_{ii} - \frac{\partial L}{\partial w_{ii}}$$

$$z_{12} = w_1 x + b$$

$$a_{12} = \max(0, z_{12}) \rightarrow$$

$$a_{11} = a_{12} \quad (\text{No Training})$$

equal to 0

## Random Initialization

① Small values  $\rightarrow$  Vanishing gradient

② Large values  $\rightarrow$  Exploding gradient

## According to requirement

1) Mid values (not small / large)

2) All weights should be different

### 3) Range Defined (Limit)

$$\text{Variance} = \sqrt{\frac{1}{n}}$$

Experientially

previous layer nodes, Fan In

$$= \frac{1}{3} = \frac{1}{5}$$

$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

### Heuristics

(tanh) Xavier Initialization (Normal Distribution)

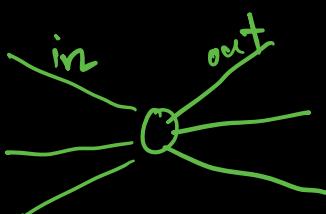
① glorot

$$W_{ini} = \sqrt{\frac{1}{\text{fan-in}}}$$

(ReLU) He Initialization (Normal Distribution)

②

$$W_{ini} = \sqrt{\frac{2}{\text{fan-in}}}$$



③ Xavier Initialization (Uniform Distribution)

$[-\text{limit}, \text{limit}]$

$$\text{limit} = \sqrt{\frac{6}{\text{fan-in} + \text{fan-out}}}$$

④ He Initialization (Uniform Distribution)

$[\text{-limit}, \text{limit}]$

$$\text{limit} = \sqrt{\frac{6}{\text{fan-in}}}$$