

Batch Normalisation

①

Even if we do proper weight initialisation & choose a right activation function \rightarrow still we may encounter gradient problem

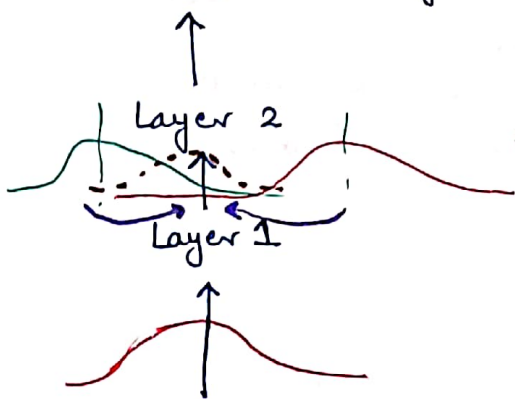
\downarrow

- ICS : Internal Covariate Shift

\rightarrow change in distribution of network

- activation o/p due to change in network parameter during training

Let's say we have a bell shaped distribution in Layer 1. If we move to Layer 2 there may be a shift in the curve or also in the distribution range of the curve due to the changes in the value of z . Hence we need to bring the distribution in a similar form in Layer 2 as well.



$$z = wx + b$$

Why?

\rightarrow It has been experimentally proven [LeCun et al, 1998 & Wiese and May 2011]

"The network converges faster if inputs to the layer are ~~narrowed~~ i.e. if they are widened normalised \Rightarrow linearly transformed to zero mean ($\mu=0$) and unit variance" ($\sigma=1$)

Our Expectation:

- fix distribution for each layer \Rightarrow reduction in Internal Covariate Shift (ICS)

\downarrow

Save dollars. \Leftarrow efficient use of resources \Leftarrow fast training \Leftarrow faster convergence

\downarrow

Batch Normalisation solves this (BN) \because In 2015 Sergey Ioffe & Christian Szegedy have

- published this paper

\downarrow
[proposed extra set of operations which can be performed before or after the activation layer]

BN Algorithm (While Training)

(2)

1) Calculate the batch mean

$$\mu_B = \frac{1}{m_B} \sum_{i=1}^{m_B} x^i$$

\downarrow batch mean \downarrow number of samples in the batch \downarrow input

2) Calculate the batch variance

$$\sigma_B^2 = \frac{1}{m_B} \sum_{i=1}^{m_B} (x^i - \mu_B)^2$$

Batch variance

3)

$$\hat{x}^i = \frac{x^i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

\downarrow Normalising x^i

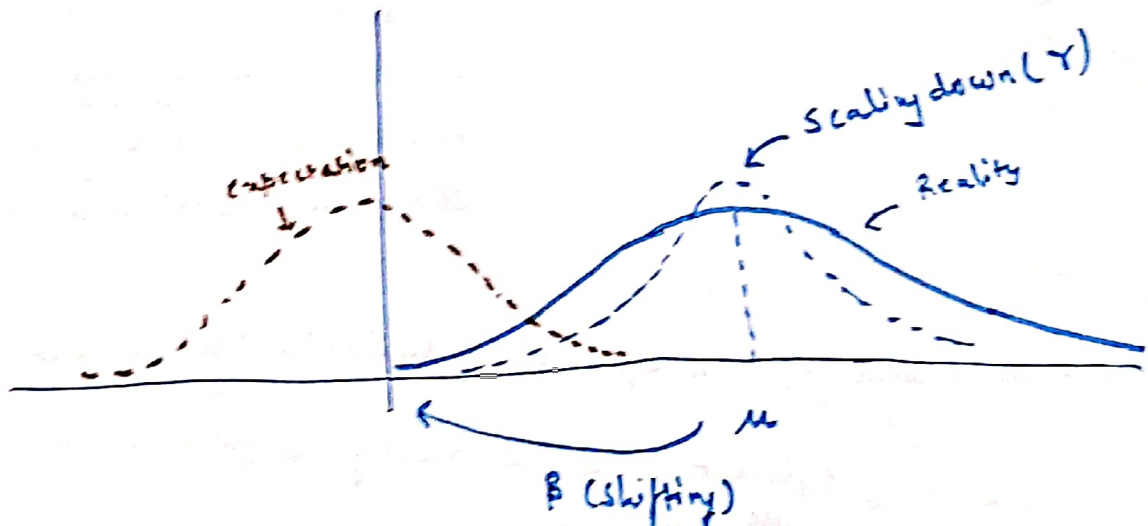
This has been introduced because σ_B can be zero.
 Smoothing Term $\epsilon = 10^{-7}$ To avoid zero division error.

4)

$$z^i = \gamma \otimes \hat{x}^i + \beta$$

\downarrow scaling parameter \downarrow Learnable Parameters \downarrow shifting parameter

shift & scale the distribution



Calculate overall mean (μ) and standard deviation (σ) by using moving average on μ_{BA} & σ_{BA}

while prediction:-

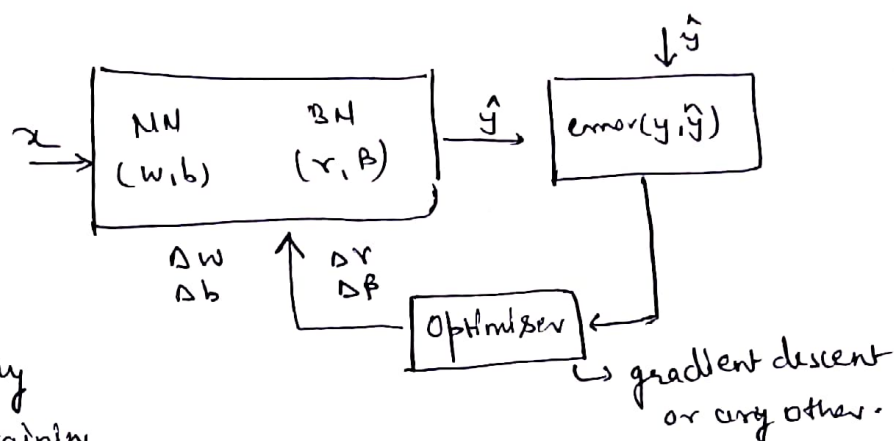
(i) $x^{(i)} = \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$ (ii) $z^{(i)} = \gamma \hat{x}^{(i)} + \beta$

① $\rightarrow \text{layer 1} \xrightarrow{z} \text{BN} \xrightarrow{\hat{z}} \text{activation function} \xrightarrow{a} \text{layer 2}$

② $\rightarrow \text{layer 1} \rightarrow \text{act}^n \xrightarrow{a} \text{BN} \xrightarrow{\hat{a}} \text{layer 2}$

③ Two approaches to apply Batch Normalisation.

In normal neural network trainable parameters are just weights and biases where as in a neural network with batch normalisation we have two extra parameters named γ (gamma) and β (Beta) that needs training



σ^2, μ will be Calculated Internally after completion of training

- Extra parameters but not trainable
our cost function is dependent on 4 parameters.

$C(w, b, \gamma, \beta) \rightarrow$ updated by backpropagation

Disadvantages

- ① It increases the complexity of the network.
- ② Number of learnable parameters increased
- ③ Runtime penalty due to complex network \Rightarrow slow prediction.
- ④ Training time is increased but convergence will be faster.

Advantages

- ① You don't need scaling of data if you are using Batch Normalisation as a 1st layer.
- ② It converges faster despite having two extra learnable parameters.
- ③ It helps to reduce the vanishing and exploding gradient issue drastically
- ④ It doesn't get affected by choice of activation function and weight initialisation technique.