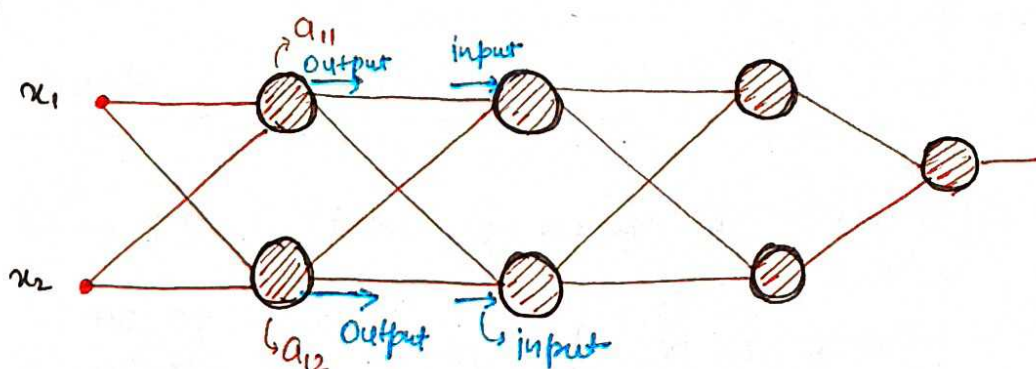


# Batch Normalization

What is Batch Norm?

Batch Normalization (BN) is an algorithmic method which makes the training of Deep Neural Network (DNN) faster and more stable.

It consists of normalizing activation vectors from hidden layers using the mean and variance of the current batch. This normalization step is applied right before (or right after) the non-linear function.



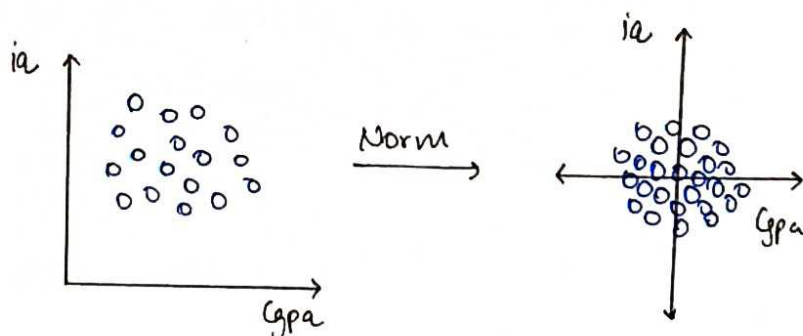
Batch Norm  $\rightarrow$  Normalize the output of the activation fn.

mean = 0      Std = 1

Why use Batch Norm?

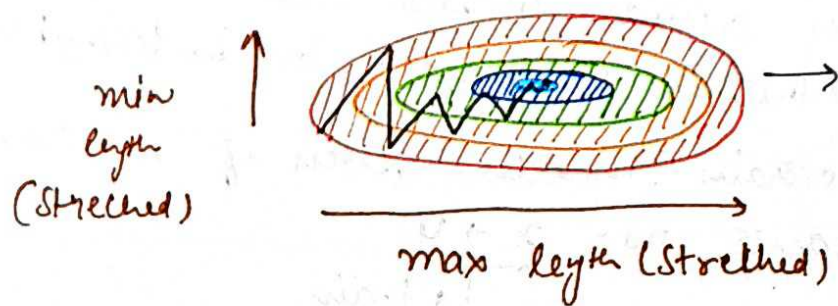
Advice  $\rightarrow$  1. Scaling data  
2. Normalization

Gpa	iq	placed
7	70	1
8	80	0
9	90	1
6	60	0



# Contour plot of unnormalized data.

(63)

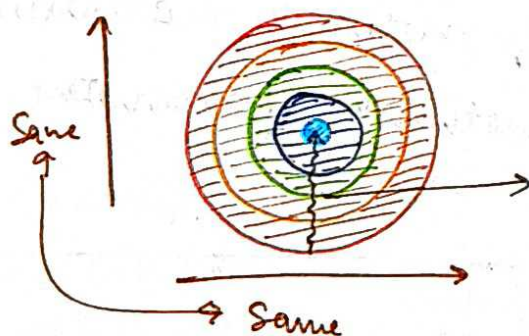


We cannot approach with higher learning rate. Highly chance to overshoot in min length direction.

That's why we use small learning rate and training will be slow.

Un-normalized data  $\rightarrow$  training would be slow.

## Contour plot after Normalize the data



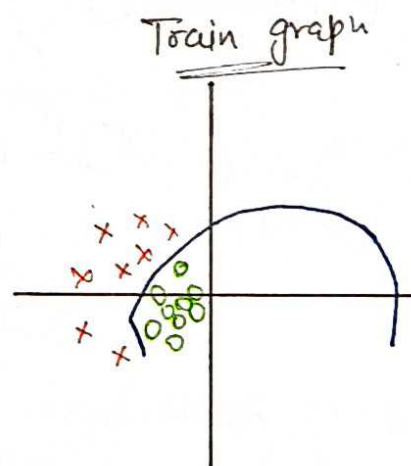
Training is faster and stable.

## Covariate Shift

Train data  $\downarrow$

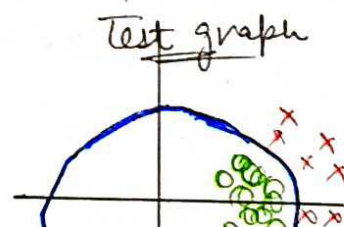
Red colour  $\rightarrow$  Rose Image  $\Rightarrow$  Rose ( $y=1$ )

different flower image  $\Rightarrow$  Not Rose ( $y=0$ )



Test Data  $\downarrow$

different color  $\rightarrow$  Rose Image  $\Rightarrow$  Rose ( $y=1$ )





In input coln of test data  $\rightarrow$  distribution changed

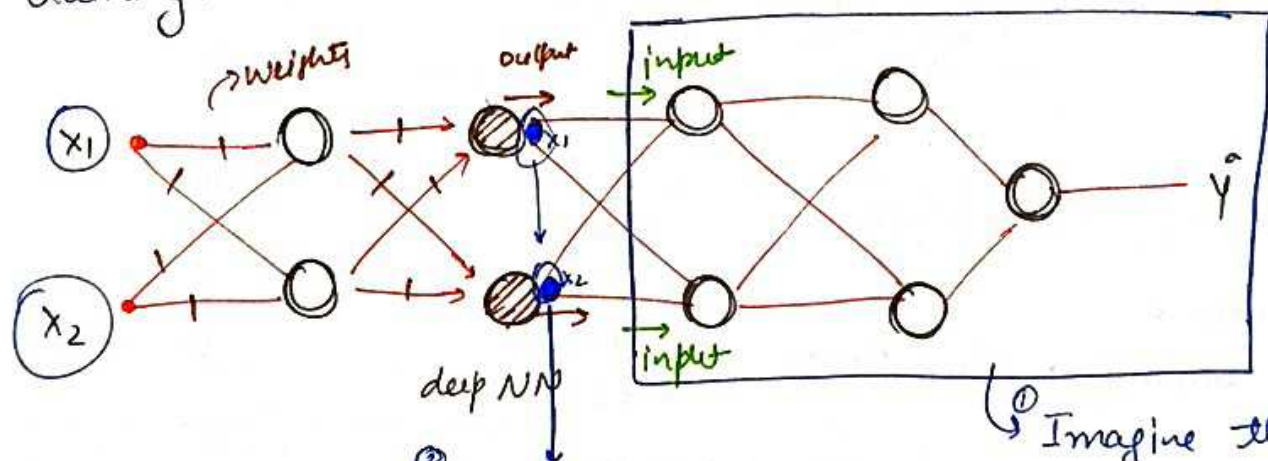
Same decision can divide between red and green points but distribution changed. ~~bad in testing position~~

we have to retain model even if relation bet<sup>n</sup>  $x \rightarrow y$  is same as  $\frac{x \rightarrow y}{\text{train}}$

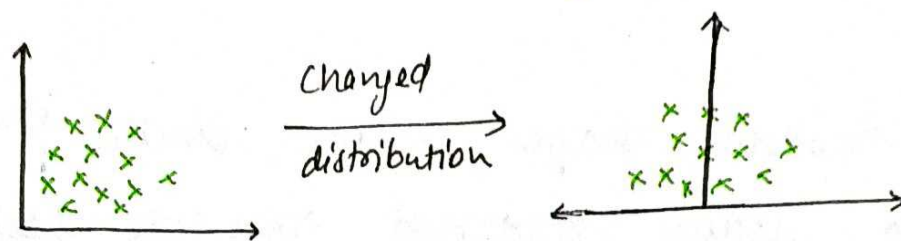
or relation is same  $x \rightarrow y$  still model not perform well. This is called covariate shift.

### Internal Covariate

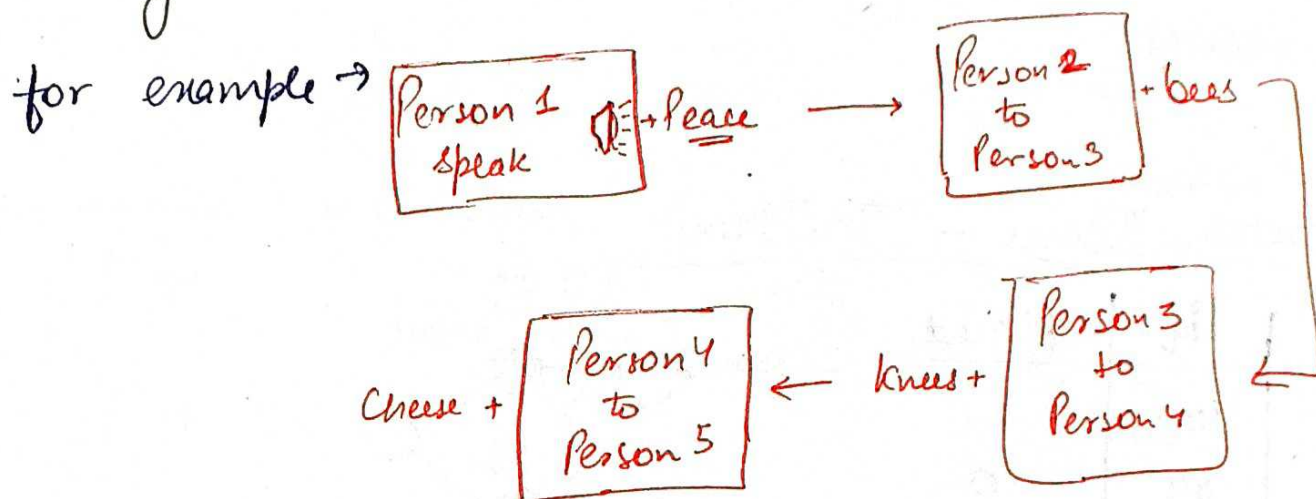
We define Internal covariate shift as the change in the distribution of network activation due to the change in network parameters during training.



- ③ Output of Hidden layer depend on previous layer of weights. And weights are constantly changed
- ② Imagine these point are working as  $x_1$  and  $x_2$ . But actually these are not input  $\rightarrow$  these are output of hidden layer.
- ① Imagine this is independent NN.



So, Independence NN face difficulties to train and training are unstable.



In this example first person speak peace and information travel and changed. Last person got information which is cheese.

Same in Neural Network, Input  $(x_1, x_2, \dots)$  pass the <sup>correct</sup> information but again and again change in weight and information are also changed (distribution <sub>correct</sub>) and wrong information to  $\hat{y}$ . Our model not trained correctly. Model not stable.

Here's Batch normalization play important role.

Batch Norm ensure after every activation function (output of activation function) is normally distributed.

$$\mu = 0, \quad \sigma^2 = 1$$



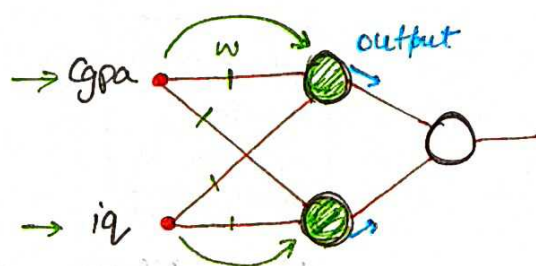
With the help of Batch Norm, Next layer will get stable ground to do their work. And training also improve.

Output of every hidden layer convert into Gaussian so this help to reduce internal covariate shift.

→ If you don't use Batch norm then use small learning rate with internal covariate shift.

### Batch Norm - The How

gpa	iq	placed
8.9	100	1
6.2	89	0
9.1	91	0
7.7	76	1
⋮	⋮	⋮
6.7	91	0



#### Bullet points for Batch Norm

- Mini-batch gradient descent
- layer by layer

$$z_{11} = w_1 \text{gpa} + w_2 \text{iq} + b$$

$$g(z_{11}) = a_{11}$$

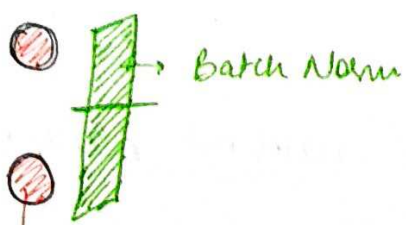
batch-Normz Normalize the output of hidden layer.

How? ↓

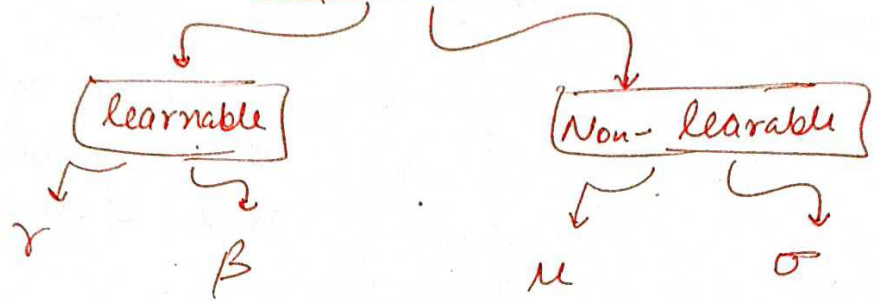
first method ↓

$$z_{11} \rightarrow z_{11}^N \rightarrow g(z_{11}^N) = a_{11}$$

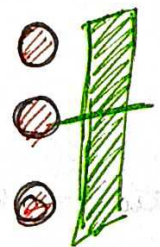
first find the  $z_{11}$  then Normalize the  $z_{11}$  value of every node in hidden layer.



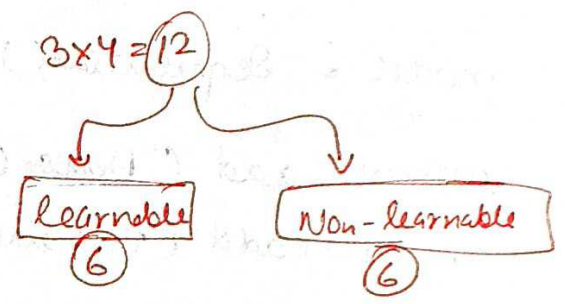
for every neuron that you have store 4 parameters



→ let assume, we have 3 neuron in Hidden layer



Total Parameter is

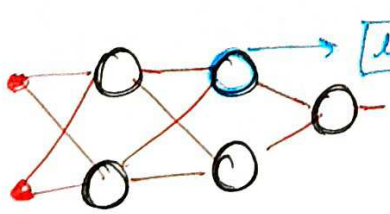


### Advantages

1. Batch Norm makes training more stable.

when you use hyperparameters, you have "wider range of values" that you can choose.   
 ← mean

2. Training become faster → learning rate value can be set to higher value.
3. Batch Norm act as Regularization (helps in overfitting)   
 But no much powerful like Regularization method.

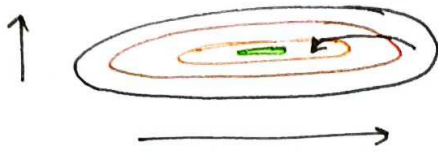


$[\mu, \sigma]$  depend on Batch

if batch value change then  $\mu, \sigma$  change.   
 And if  $\mu, \sigma$  change then activation change and randomness add → reduce overfitting.

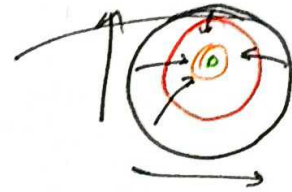
4) Weight init impact reduce

Without Batch Norm



face difficulties  
to reach optimal  
solution.

Without Batch Norm



use diff weight  
Still easily you can  
reach optimal solution

### Keras Implementation

```
model = Sequential()
```

```
model.add(Dense(3, activation='relu', input_dim=2)),
```

```
model.add(BatchNormalization()),
```

```
model.add(Dense(2, activation='relu')),
```

```
model.add(BatchNormalization())
```

```
model.add(Dense(1, activation='sigmoid'))
```



## Second Method

65

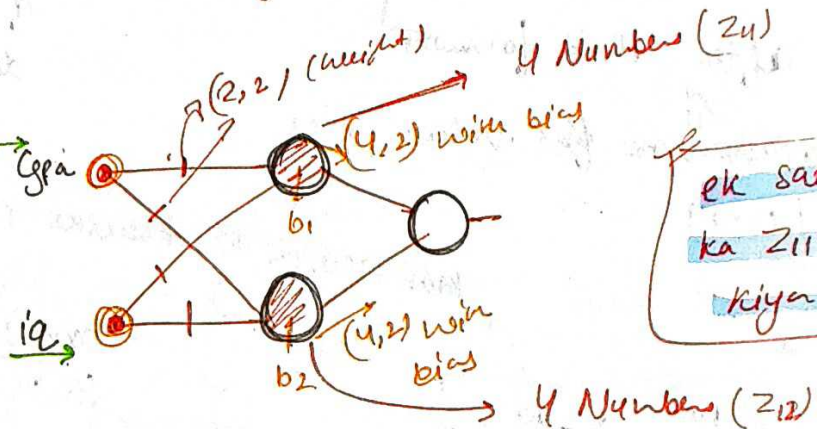
$$Z_{11} \rightarrow g(Z_{11}) = a_{11} \rightarrow a_{11}^N$$

first, find the  $Z_{11}$  and put into  $g(Z_{11})$  and find the  $a_{11}$  value of every node of hidden layer apply normalization.

$$\frac{Z_{11} - \mu}{\sigma} = Z_{11}^N$$

how to find  $\mu$ ?

Cgpa	iq	placed
6.2 ①	89	0
9.1 ②	91	0
7.7 ③	76	1
8.9 ④	150	1
6.7	91	0



ek sath 4 points ka  $Z_{11}$  find kiya

Input shape

$$\rightarrow (4 \times 2) \text{ and multiply with } (2 \times 2)$$

$$(4 \times 2) * (2 \times 2) = (4 \times 2)$$

$$\rightarrow (4, 2) + (1, 2) = (4, 2)$$

bias

after adding bias we get  $(4, 2)$  of every node.

Now calculate  $\mu$ ,

$$\mu = \frac{1}{m} \sum_{i=1}^m Z_{11}^i$$

where  $m = 4$  i.e. of batch size.

[ek sath 4  $Z_{11}$  find kiye the wahi use karke karne hai]

let, batch-size = 4

Take 4 points and Normalize it. send to NN (ek sath 4 points send kiye)



$$\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_{11}^i - \mu_B)^2}$$

we have to find  $\mu_B$  and  $\sigma_B$  for two times because NN have 2 node in hidden layer.

Step 1:  $\rightarrow$

So, we have 4 activation as per neuron (node).

Now we have to calculate as per activation of per node.

$$z_{11}^i = \frac{z_{11}^i - \mu_B}{\sigma_B + \epsilon}$$

4  $z_{11}$  the, unka individual  $z_{11}^i$  calculate hogae

add this error term if std is 0 then denominator not get 0.

applying after this formula all activation  $z_{11}^i$  value lies betn (0-1)

har neuron ka khudke  $\gamma, \beta$  parameter hogae

Step 2:  $\rightarrow$

$$z_{11}^{BN} = \gamma z_{11}^N + \beta$$

learnable parameter

correct value of  $\gamma$  and  $\beta$  find during training, but initial values is  $\gamma=1, \beta=0$  in keras. after back propagation values ~~not~~ changed.

Step 3 =  $g(z_{11}^{BN}) = a_{11}$

Step 4 =  $a_{11} \rightarrow$  activation func<sup>n</sup>

Why we use  $\gamma$  and  $\beta$ ?

Because in Neural Network some dataset not need to Normalization. So, these values ( $\gamma, \beta$ ) help to reverse the Normalization.

eg:- let assume

$$\gamma = \sigma + E \quad \beta = \mu \implies$$

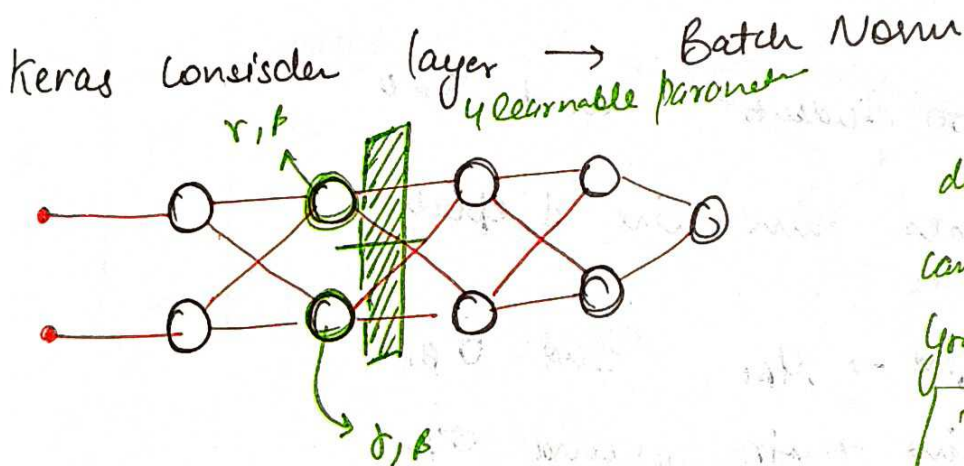
these value come

these value come when dataset not need normalization.

$$Z_{ii} \Rightarrow Z_{ii}^N \rightarrow Z_{ii}^{BN}$$

both are same

So, reverse the Normalization and back to  $Z_{ii}$  value.



during learning, we can update  $\gamma$  using gradient descent.

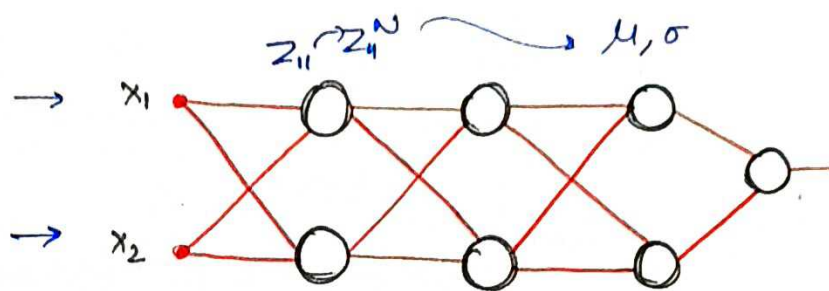
$$\gamma = \gamma - \eta \frac{\partial L}{\partial \gamma}$$



## Batch Norm During test

cgpa	iq	plaid
8	80	1
7	70	0
6	60	1
1		
1		
1		
9	90	1

100 students



In training, we use batch or minibatch like  $b=3$  then <sup>train</sup> first 3 datapoint together. for Normalizing process we find  $\mu$  and  $\sigma$  with the help of 3 datapoints.

But in test, we have only 1 query point. So, we use Exp weighted moving avg.   
 Exponential

let assume,

data  $\rightarrow$  100 students and  $b=4$

So, 25 times data run in 1 epoch

$\Rightarrow$  At batch 1  $\rightarrow$  find  $\rightarrow \mu_{\beta_1}$  and  $\sigma_{\beta_1}$

$\Rightarrow$  At batch 2  $\rightarrow$  find  $\rightarrow \mu_{\beta_2}$  and  $\sigma_{\beta_2}$

⋮

$\Rightarrow$  A batch<sub>25</sub>  $\rightarrow$  find  $\rightarrow \mu_{\beta_{25}}$  and  $\sigma_{\beta_{25}}$

EWMA Maintain  $\rightarrow \mu$  and  $\beta$

$\rightarrow$  After ~~test~~ training complete, last updated  $\mu$  and  $\beta$  use in test.