$$\text{normalize} = \begin{bmatrix} \text{mean} \rightarrow 0 \\ sd \rightarrow 1 \end{bmatrix}$$

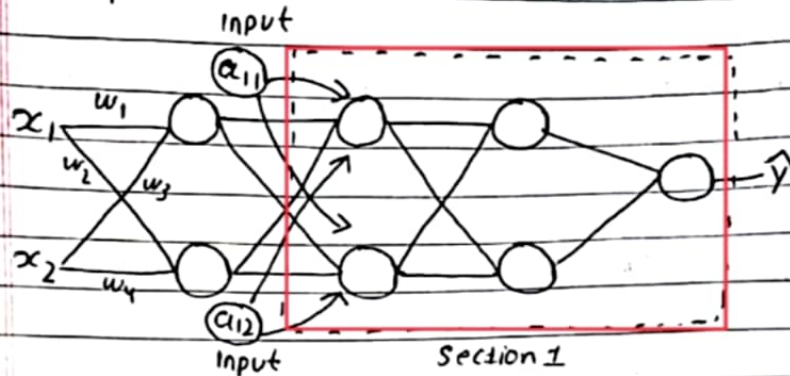**Batch Normalization** – It is used to speed up Neural network training

It normalizes the activation vectors from hidden layer using mean and variance of current batch

Example:



We normalize the input $x_1, x_2$ (It's normal case) but In Batch normalization process we also normalize the $[a_{11}, a_{12}]$.
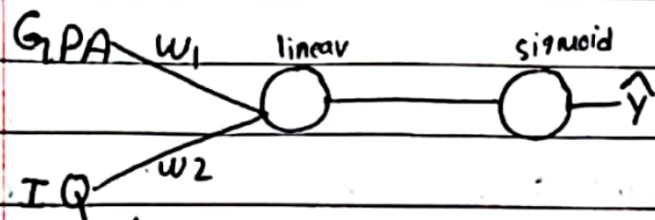
**Concept:**



Input          Section 1

⇒ Section 1 lai xutai individual neural network socham. Section 1 ko first layer ko node) ley input $a_{11}$ and $a_{12}$ linxa. $a_{11}$ and $a_{12}$ values after one batch of training change vaisakxa ra yesko distribution ne change hunxa kina vane. Euta batch ko training vaye pari weight update hunxa. $a_{11}$ ra $a_{12}$ ($w_1, w_2, w_3$ $w_4$) maa depend xa ra       harek batch ko training vaye pari ($w_1, w_2, w_3$ re $w_4$) pani change hunxa

so, section 1           ko node) ley different distribution           Input linxa Taile training stable hunna ra model ley na chaini kura pani learn garxa ra slow hunxa so, we use concept of batch normalization.

Jaba $a_{11}$ ra $a_{12}$ normalize hunxa yesko value) between 0 and 1 ku range ma hunxa so, Section 1 ko layer node) leg stable data pauxan ra training faster hunxa.

[Internal covariance shift- change of distribution of activation due to change of network parameters during training]

| GPA | IQ | Placement |
|-----|-----|-----------|
| 8.9 | 100 | 1 |
| 6.2 | 25 | 0 |
| 8.3 | 2 | 0 |
| 9.2 | 70 | 1 |

GPA $w_1$ — linear — sigmoid — $\hat{y}$

IQ $w_2$

$z'_{11} = 8.9 \times w_1 + 100 \times w_2 + b_{11}$   $\big[ z'_{11} \rightarrow$ first row batch ko

$\qquad = M'$

$z'_{11} \rightarrow$ first node

First layer $\big]$

[ $M' \rightarrow$ linear output of first row]

Similarly we will have $M^2, M^3, M^4$

Our batch size is 4 so, in one batch we have
4 rows so, we will get 4 outputs [ $M', M^2, M^3, M^4$ ]

Now, we will normalize them using formula

$\qquad Z'_{nor} = \dfrac{m^i - \mu}{6}$   [ $i$ is row of batch ]

$\therefore \mu = \dfrac{1}{n} \sum\limits_{i=1}^{n} m^i$

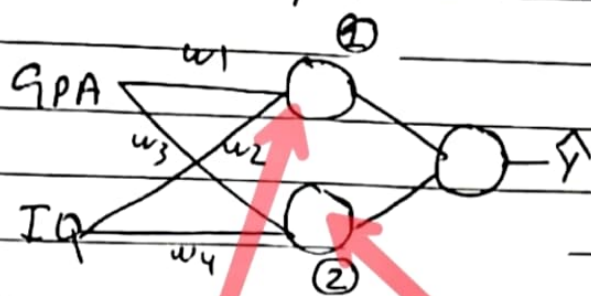$\Rightarrow \dfrac{[M' + M^2 + M^3 + M^4]}{4 \,(batch\,size)}$   [n is batch size]

$6 = \sqrt{\dfrac{1}{n} \sum\limits_{i=1}^{n} (m^i - \mu)^2}$

Now, let's do for $z'_{nor} = \dfrac{m^1 - \mu}{6 + \varepsilon}$   ( $\varepsilon$ is small value
if 6 is 0 the $\varepsilon$
will prevent it from
undefined )

Similary we will get ( $z'_{nor}, z^2_{nor}, z^3_{nor}, z^4_{nor}$ )

## If It has multiple nodes like below figure



| GPA | IQ | Plac |
|-----|-----|------|
| 1 | 2 | 1 |
| 10 | 11 | 0 |

we have to calculate $u$ and $6$ seperately for $1, 2$ nodes and do normalize.

| Node 1 Outputs: | Node 2 outputs |
|---|---|
| $\rightarrow 1 \cdot w_1 + 2 \times w_3 + b_{11}$ | $1 \cdot w_2 + 2 \times w_4 + b_{12}$ |
| $\rightarrow 10 \cdot w_1 + 11 \times w_3 + b_{11}$ | $10 \cdot w_2 + 11 \times w_4 + b_{12}$ |
| (calculate seperate $u$ and $6$ and normalize) | (calculate seperate $u$ and $6$ and normalize) |

## Now,

### In depth Intution

$$Z'_{11} \rightarrow Z'_{11 nor} \rightarrow P'_{11} \rightarrow q(P'_{11}) \rightarrow a'_{11}$$

( $Z'_{11}$  first layer ko first node ley    batch ko first row ko output)

Here $\boxed{P = Y \cdot Z + \beta}$ ($Y$ and $\beta$ are learnable parameters whose values can be find during training)

Jaile normalize garda ramro accuraly na aauna sakxa tei vara $Y$ and $\beta$ learnable parameter rakheko yo pani $w$ and $b$ update vako jasari update hunxa during back propogation like $\boxed{Y_{new} = Y_{old} - \dfrac{d\ Loss}{d\ Y_{old}}}$

Generally $Y = 1$ and $\beta = 0$ during first time.