# Batch Normalization

Normalizing inputs to speed up learning

Given some intermediate values in NN, for layer $l$ unit $i$

given as $Z^{[l](i)} = \{ z^{(1)}, z^{(2)}, \ldots, z^{(m)} \}$

where each values are from some inputs of a mini batch.

so $\mu = \frac{1}{m} \sum_i z^{(i)}$   $\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$

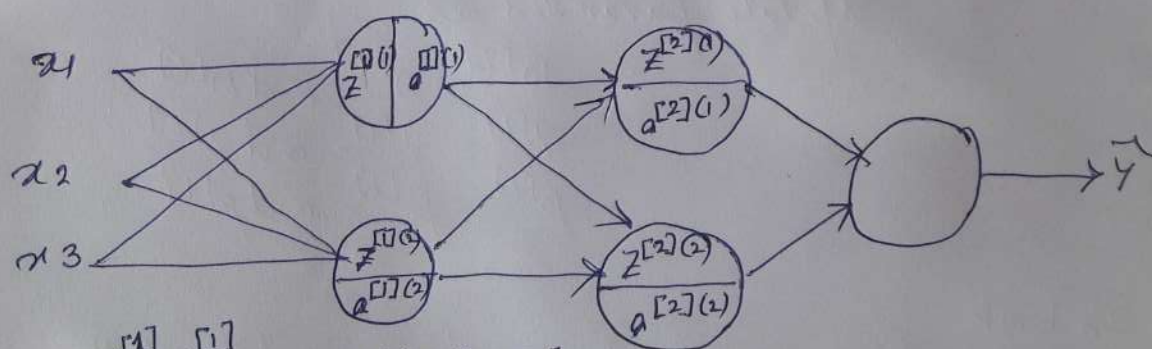$Z^{(i)}_{norm} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$

$\hat{z}^{(i)} = \gamma z^{(i)}_{norm} + \beta$.

$\gamma$ & $\beta$ are learnable.

If $\gamma = 1$, $\beta = 0$

$\hat{z}^{(i)} = z^{(i)}_{norm}$.

If $\gamma = \sqrt{\sigma^2 + \epsilon}$   $\beta = \mu$

then $\hat{z}^{(i)} = z^{(i)}$.

Filling Batch norm in a neural network



$X \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \underset{\text{Batch norm}}{\xrightarrow{\beta^{[1]}, \gamma^{[1]}}} \tilde{z}^{[1]} \longrightarrow a = g(\tilde{z}^{[1]})$

$\Big] w^{[2]}, b^{[2]}$

$\cdots \longleftarrow a^{[2]} \longleftarrow \tilde{z}^{[2]} \underset{\text{Batch norm}}{\xleftarrow{\beta^{[2]}, \gamma^{[2]}}} z^{[2]} \longleftarrow$

Parameters $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]} \cdots$
$\gamma^{[1]}, \beta^{[1]}, \gamma^{[2]}, \beta^{[2]} \cdots$

The bias term may be removed from calculation

as $Z^{[L]}$ $= W^{[L]} a^{[L-1]} + b^{[L]}$

so $\mu^{[L](i)} = \dfrac{1}{m} \sum_{L=1}^{m} Z^{[L](i)} = P^{(i)} + b^{[L](i)}$

The bias term remains in $\mu$

so in $Z^{(i)} - \mu$, the bias term gets removed.

Implementing Gradient descent with BN

for $k = 1$ to $m$ (# of minibatches)

    do forward propagation using $X^{(k)}$

    In each hidden layer use BN to replace
        $Z^{[L]}$ with $\tilde{Z}^{[L]}$

    use back prop to compute
        $dW^{[L]}, d\gamma^{[L]}$ & $d\beta^{[L]}$

    update parameters as
$$W^{[L]} = W^{[L]} - \alpha \, dW^{[L]}$$
$$\gamma^{[L]} = \gamma^{[L]} - \alpha \, d\gamma^{[L]}$$
$$\beta^{[L]} = \beta^{[L]} - \alpha \, d\beta^{[L]}$$

In book

H is considered as a design matrix of the activation of a minibatch

$$H = \begin{vmatrix} x^{(i)} \\ \downarrow \\ i \end{vmatrix} \xrightarrow{\text{minibatch } k} = \begin{array}{|c|c|c|} x^{\{1\}(1)} & x^{\{2\}(1)} & x^{\{3\}(1)} \\ \hline x^{\{1\}(2)} & x^{\{2\}(2)} & x^{\{3\}(2)} \end{array}$$

$$\mu = \begin{bmatrix} \text{mean} \\ \text{of} \\ \text{rows} \end{bmatrix} \qquad \sigma = \begin{bmatrix} \text{variance} \\ \text{of} \\ \text{rows} \end{bmatrix}$$

$$H' = \dfrac{H - \mu}{\sigma}$$