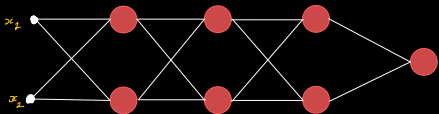# Batch Normalization

→ An algorithmic method which makes training of NN more faster & better.

→ It contains normalizing activation vectors from hidden layer using mean & variance of current batch.

→ This step is applied right before (or right after) the non-linear function.

→ In short we will normalize activation of each neuron.



→ Generally our inputs $(x_1, x_2)$ are normalized $(u = 0, \sigma = 1)$

→ In Batch normalization, o/p of each node will be normalized. $(u = 0, \sigma = 1)$.

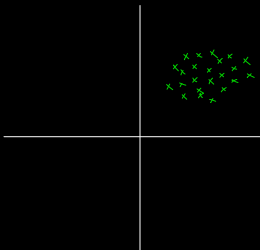→ This process will be for each hidden layers. makes training fast & stable.

## Why use batch normalization..?

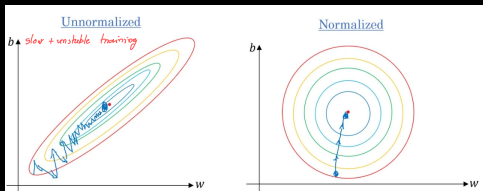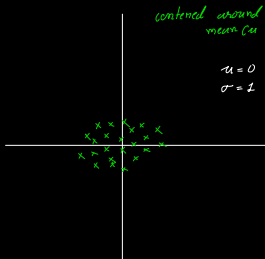→ In NN, it is advisable to normalize data before giving to model.

Example

Original Data

Normalized Data

centered around
mean (≈

$u = 0$
$\sigma = 1$



Unnormalized

slow & unstable training

$b$

$w$

Normalized

$b$

$w$

→ Due to unnormalized data, the cost function will look like first image, stretched in one direction.

→ in this scenario, we can't go with large learning-rate, because it may over-shoot in one direction.

→ So if normalized i/p are make training smoother & faster, what if we normalize the o/p of each activation in NN.

Covariate shift

→ A situation where the distribution of i/p changes between the training & testing data, while relationship between i/p & o/p remains same.

| Training Distribution | | Testing Distribution | | |
| --- | --- | --- | --- | --- |
| dogs in water | cats in grass | seen classes | | unseen classes |
| | | dogs in grass | cats in water | bicycle, boat... |

→ This image shows covariate shift
→ As image shows, the model is trained to detect dog or cat in grass or water.
→ But at testing time, model get some unseen objects located in grass or water.

In this situation we need to re-train our model.

## Internal covariate shift

"We define internal covariate shift as change in the distribution of network activations due to the change in network parameters during training."



network-1                    network-2

→ Here, i/p of network-2 is output of last layer of network-1.
→ The o/p of network-1 is depending on weights & bias of layers.
→ Since weights of network-1 are constantly changing the o/p is also constantly changing.
→ Due to this scenario, i/p distribution of network-2 is constantly changing due to this problem network-2 faces issue in training (unstable training).
→ This problem is called as Internal covariate shift.
→ In case of ICS without data normalization, learning-rate should be ...

# How batch normalization works...?

=> Works with mini-batch gradient descent.
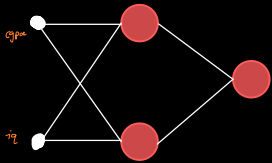=> Applies layer by layer, optional for each layer.

| cgpa | iq | placed |
|------|-----|--------|
| 6.2  | 100 | 1      |
| 6.2  | 84  | 0      |
| 9.7  | 92  | 0      |
| 7.7  | 76  | 1      |



# Now process will goes....

$$z_{11} = w_{11} \, cgpa + w_{21} \, iq + b_{11}$$

Normalizing $z_{11}$ before further process.

$$z_{11}^N \quad (\mu = 0, \sigma = 1)$$

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

OR

$$z_{11} = w_{11} \, cgpa + w_{21} \, iq + b_{11}$$

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

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

most popular method

$$\boxed{a_{11}}$$

Now we will normalize this $a_{11}$

$$a_{11}^N$$

$$\boxed{a_{11}^N}$$

for $z_{11}^N = \dfrac{z_{11} - \mu}{\sigma}$ , 'μ' will be mean of batch.

$$\mu_B = \frac{1}{m} \sum_{i=1}^{m} z_{11}^i \qquad \text{where}$$

$$m = \text{batch size}$$

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

$$m = \text{batch size}$$

$$z_{11}^i = \frac{z_{11}^i - u_B}{\sigma_B + \epsilon} \quad \text{where}$$

$\epsilon = \text{Error term}$

Now we have $Z_{11}^N$, still not done yet, we will do one-more operation here.

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

(in keras)

where:

$\gamma$ and $\beta$ are learnable parameters.

By default:

$\gamma = 1$

$\beta = 0$

$g(Z_{11}^{BN}) = a_{11}$

→ This process, we will do for each neurons of a layer...

→ Each neuron will have their own $\gamma$ and $\beta$ parameters...

Q. Why we're using $\gamma$ and $\beta$.

→ Sometimes we don't need normalized data on our NN doesn't want data to be normalized.

→ In this kinds of situations, using $\gamma$ and $\beta$ value we can change our distribution.

during Back-propagation, these 2 params will also get updated..

$$\gamma = \gamma - \alpha \frac{\partial L}{\partial \gamma} \qquad \beta = \beta - \alpha \frac{\partial L}{\partial \beta}$$

## Advantages

→ Fast & stable training
→ Regularization effect.
→ Reduces weight initialization impact.