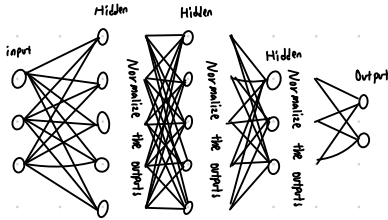


Normalization

: We can design our models such that they are easier to optimize

Types of Normalization

- Batch normalization : Used in ResNets and many other models
- Instance normalization
- Layer normalization : Used in Transformer
- Group normalization



→ the quality or condition of being in good condition

Normalization layers are often used to increase training robustness of neural network

Why we need normalization? : Main issue is "Covariate shift"

- Training assumes the training data are similarly distributed. However, in practice, each minibatch may have different distribution : "Covariate shift". It may occur in each layer of the network.



: Think MNIST dataset.

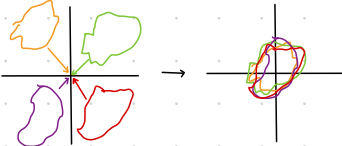
Training assumes that distributions of digit '2' are all similarly distributed, and it's very naive approach.

We need to apply normalization to make each digits' distributions as similar as possible.

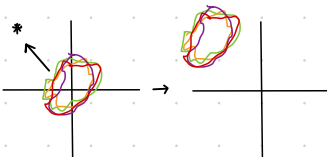
Solution: Move all minibatches to a "standard" location

General concept of Normalization

① Move all batches to have a mean of 0 and unit standard deviation : This process eliminates covariate shift between batches



② Move the entire collection to the appropriate location



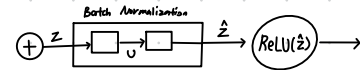
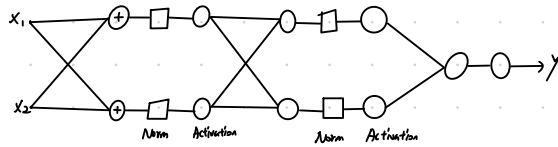
Batch normalization is a "shift-adjustment unit" that happens after the weighted addition of inputs but before the application of activation.

- Is done independently for each unit, to simplify computation.

Ex) self basis = nn.Sequential(

```
nn.Linear(X_dim, 100),
nn.BatchNorm(100),
nn.ReLU(),
...
)
```

In pictorial depiction



$$Z = \sum_i W_i X_i + b, \quad W_i X_i + W_j X_j + b = w^T X + b.$$

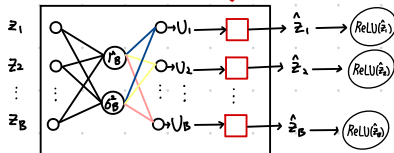
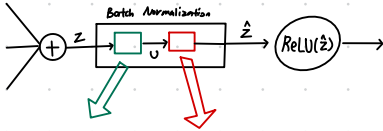
$$U = \frac{Z - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \text{Centering}$$

$$\sqrt{\sigma_B^2 + \epsilon} \quad \text{Scaling, where } \mu_B = \frac{1}{B} \sum_{i=1}^B Z_i \quad \text{and } \sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (Z_i - \mu_B)^2$$

$$\hat{Z} = \gamma U_i + \beta, \quad \text{where } \gamma, \beta \text{ are learnable features, and they're responsible for shifting distributions into appropriate locations.}$$

★ Tip: Batch Normalization prefers larger batches.

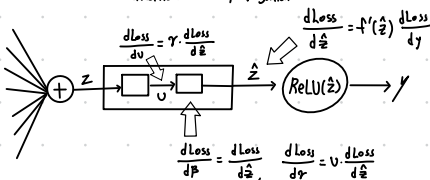
A better picture for batch norm more explicitly



$$\text{where } U_i = \frac{Z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, \quad \hat{Z}_i = \gamma U_i + \beta$$

We can compute $\frac{d\text{Loss}}{dU_i}$ individually for each U_i because the processing after the computation of U_i is independent for each U_i .

Batch normalization: Backpropagation



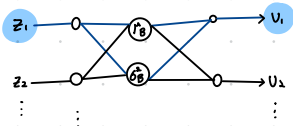
Complete derivative of the batch-loss w.r.t. z_i

$$\frac{d \text{Loss}}{dz_i} = \sum_j \frac{d \text{Loss}}{dU_j} \frac{dU_j}{dz_i}$$

We know $\frac{d \text{Loss}}{dU_i} = \gamma \cdot \frac{d \text{Loss}}{d\hat{z}_i} = \gamma \cdot f'(\hat{z}_i) \frac{d \text{Loss}}{dy}$

$$\frac{dU_j}{dz_i} ?$$

① The derivative for the "through" line ($i=j$)



$$\frac{dU_i}{dz_i} = \frac{\partial U_i}{\partial z_i} + \frac{\partial U_i}{\partial I_B} \frac{\partial I_B}{\partial z_i} + \frac{\partial U_i}{\partial O_B} \frac{\partial O_B}{\partial z_i}$$

1.1: $U_i = \frac{z_i - I_B}{\sqrt{\sigma_B^2 + \epsilon}}$, then $\frac{\partial U_i}{\partial z_i} = \frac{\partial}{\partial z_i} \left(\frac{1}{\sqrt{\sigma_B^2 + \epsilon}} (z_i - I_B) \right) = \frac{1}{\sqrt{\sigma_B^2 + \epsilon}}$

1.2: $U_i = \frac{z_i - I_B}{\sqrt{\sigma_B^2 + \epsilon}}$, then $\frac{\partial U_i}{\partial I_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}}$

1.3: $I_B = \frac{1}{B} \sum_{i=1}^B z_i$, then $\frac{\partial I_B}{\partial z_i} = \frac{1}{B}$

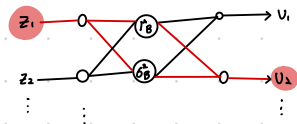
1.4: $U_i = \frac{z_i - I_B}{\sqrt{\sigma_B^2 + \epsilon}}$, then $\frac{\partial U_i}{\partial \sigma_B^2} = \frac{\partial}{\partial \sigma_B^2} \left((z_i - I_B) \cdot (\sigma_B^2 + \epsilon)^{-\frac{1}{2}} \right) = -\frac{(z_i - I_B)}{2} \cdot (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} = \frac{-(z_i - I_B)}{2(\sigma_B^2 + \epsilon)^{\frac{3}{2}}}$

1.5: $\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - I_B)^2$, then $\frac{d\sigma_B^2}{dz_i} = \frac{\partial \sigma_B^2}{\partial z_i} + \frac{\partial \sigma_B^2}{\partial I_B} \frac{\partial I_B}{\partial z_i} = \frac{2(z_i - I_B)}{B} + 0$

Hence, the derivative for the through line ($i=j$)

$$\frac{dU_i}{dz_i} = \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} - \frac{1}{B\sqrt{\sigma_B^2 + \epsilon}} + \frac{-(z_i - I_B)^2}{B(\sigma_B^2 + \epsilon)^{\frac{3}{2}}}$$

② The derivative for the "cross" lines ($i \neq j$)



$$\frac{dU_j}{dz_i} = \frac{\partial U_j}{\partial I_B} \frac{\partial I_B}{\partial z_i} + \frac{\partial U_j}{\partial \sigma_B^2} \frac{d\sigma_B^2}{dz_i}$$

$$= \frac{-1}{B\sqrt{\sigma_B^2 + \epsilon}} - \frac{(z_i - I_B)^2}{B(\sigma_B^2 + \epsilon)^{\frac{3}{2}}}$$

①+②

$$\frac{dU_i}{dz_i} = \begin{cases} \frac{1}{\sqrt{\sigma_0^2 + \epsilon}} - \frac{1}{B\sqrt{\sigma_0^2 + \epsilon}} + \frac{-(z_i - \mu_0)^2}{B(\sigma_0^2 + \epsilon)^{\frac{3}{2}}} & , \text{ if } i=j \\ \frac{-1}{B\sqrt{\sigma_0^2 + \epsilon}} - \frac{(z_i - \mu_0)^2}{B(\sigma_0^2 + \epsilon)^{\frac{3}{2}}} & , \text{ if } i \neq j \end{cases}$$

And, recall the complete derivative of the mini-batch loss w.r.t. z_i

$$\frac{dL_{\text{loss}}}{dz_i} = \sum_j \frac{dL_{\text{loss}}}{dU_j} \frac{dU_j}{dz_i}$$