

3 Layer Normalization

↳ It's applied to each input individually rather than to one feature across all inputs.

Normalization per data point

① Normalization:

1- Data Normalization:

Also referred to as "Standardization". It's a transformation applied to the input data. It is a calculation performed once, before training starts. It transforms the data so that:

$$\text{mean} = 0, \text{standard deviation} = 1$$

↳ the most common formula is:

$$x_{\text{norm}} = \frac{x - \text{mean}}{\text{std}}$$

Original value

Std = 1 means the average distance of data points from the mean is 1.
Spread is moderate. Not too tight and not too wide.

Why do we do it?

↳ if input features have different ranges (one in $[0, 1]$, another in $[0, 1000]$), the model gives more importance to the large feature because it contributes more to the loss → causes the weights connected to that feature to dominate the training while the others are barely updated.

↳ training unstable → Normalization fixes this by putting all features on a similar scale (mean 0, std 1) so they contribute equally to the learning process.

Example:

1) ImageNet and Image Normalization

Most pre-trained models (ViT..) are trained using the ImageNet normalization.

$$\begin{aligned} \text{IMAGENET_MEAN} &= [0, 485, 0, 456, 0, 406] \\ \text{STD} &= [0.229, 0.224, 0.225] \end{aligned}$$

↳ these values were computed over millions of images.

A pixel with RGB values:

$$[R, G, B] = [120, 100, 90]$$

feature

Step 1: Scale to $[0, 1]$

$$R' = \frac{120}{255} = 0.47, G' = \frac{100}{255} = 0.392, B' = \frac{90}{255} = 0.353$$

Step 2: Normalize using ImageNet std and mean

$$R - \text{mean} = \frac{0.47 - 0.485}{0.229} = -0.065$$

⋮

2) Let's say we have:

5 vectors (samples)

each vector has 3 features (dimensions)

↳ we'll do a per feature normalization

$$X = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \\ 4 & 7 & 10 \\ 5 & 8 & 11 \end{bmatrix}$$

Step 1: Compute [mean] per feature = column:

$$\mu = \begin{bmatrix} \frac{1+2+3+4+5}{5} \\ \frac{4+5+6+7+8}{5} \\ \frac{7+8+9+10+11}{5} \end{bmatrix} = \begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix}$$

Step 2: Compute [std] per feature:

$$\sigma_j = \sqrt{\frac{1}{5} \sum_{i=1}^5 (x_{ij} - \mu_j)^2} = \begin{bmatrix} 1,014 \\ 1,914 \\ 2,814 \end{bmatrix}$$

Step 3: Normalize each value?

$$x_{ij} - \text{mean} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

3 Layer Normalization

↳ It is a technique used to normalize the pre-activations (as it comes before the activation function) instead of normalizing across all feature dimensions.

Normalization per data point

↳ It is a technique used to normalize the pre-activations (as it comes before the activation function) instead of normalizing across all feature dimensions.

Why is it used instead of BatchNorm?

RNNs & sequence length varies → hard to use batch stats.

↳ for BatchNorm, we need shape = (B, D) but in RNNs, shape = (B, T, D)

time step (e.g., words)

Suppose we have 3 sequences with different lengths

Sample	Seq Length	
A	4	we need to pad
B	3	→ them all to
C	2	T = 4. Let's assume D = 1

Input shape = $(3, 4, 1)$

$$X = \begin{bmatrix} [1][2][3][4] \\ [1][2][3][PAD] \\ [1][2][PAD][PAD] \end{bmatrix}$$

↳ = 0 [common practice]

Now let's apply BatchNorm per time step to have the desired shape = (B, D)

so at t=2, we have

$$x[:, 2, :] = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$

↳ so the calculus will be noisy

→ So, with LayerNorm, for a given sample

$$x = [x_1, x_2, \dots, x_d], \text{ we:}$$

1- Compute the mean and std per sample

$$\mu = \frac{1}{d} \sum_{j=1}^d x_j \quad \sigma^2 = \frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2$$

2- Normalize

$$\hat{x}_j = \frac{x_j - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

3) Apply learnable shift and scale

$$y_{ij} = \gamma \hat{x}_{ij} + \beta$$

Learnable parameters per feature

2) Batch Normalization

It is a regularization or optimization technique. It is designed to normalize the activation of the hidden layers using statistics computed over a mini-batch. By normalizing the output of a layer to have a zero mean and unit std, BatchNorm allows to mitigate the problem of internal covariate shift.

↳ this is when the distribution of activations shifts during training.

BatchNorm is applied after the linear or convolutional layer and before the activation function (e.g., ReLU) which ensures that the inputs to the non-linearity are well conditioned.

Example:

① Let's say we have:

- a batch of 3 samples
- each sample has 3 features

$$X = \begin{bmatrix} 3 & 1 \\ 4 & 2 \\ 5 & 3 \end{bmatrix} \text{ shape} = (3, 2)$$

↳ we apply batch norm → normalize each column independently across the batch.

$$\text{Step 1: } \begin{cases} \mu_1 = \frac{1}{3} [3+4+5] = 4 ; \quad \sigma_1^2 = \frac{1}{3} [(3-4)^2 + (4-4)^2 + (5-4)^2] = \frac{2}{3} \\ \mu_2 = 2 \quad \sigma_2^2 = \frac{2}{3} \end{cases}$$

Step 2: Normalize

• We add a small $\epsilon = 1e-5 = 10^{-5}$ for stability

$$\sigma = \sqrt{\sigma^2 + \epsilon} = 0.8165$$

• Normalization formula

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

Linear/conv layer →

BatchNorm →

Activation function

$$X = \begin{bmatrix} -1,225 & -1,225 \\ 0 & 0 \\ +1,225 & +1,225 \end{bmatrix}$$

Step 3: Apply Learnable Scale and Shift

↳ Let's assume $\gamma = [2, 0.5]$

$$\beta = [1, 0]$$

$$y_{ij} = \gamma \hat{x}_{ij} + \beta$$

Why use scale and shift?

γ, β restore flexibility. They allow the model to "undo" the standardization if needed.

↳ Normalization always squashes everything between 0 and 1 but maybe some neurons should be more active (= features more emphasized)

• If the network wants the feature to behave like before it can learn $\gamma = 8, \beta = 0$

• If the network wants to suppress a feature $\gamma = 0$