# Batch Normalization (element-wise operation)

why use it: backpropagation tends to training more on the top layers, once bottom layer's weight changed, what we learned before need to be learned again. Hard to converge, internal covariate shift

batch norm allows stablize training process and allows faster convergence

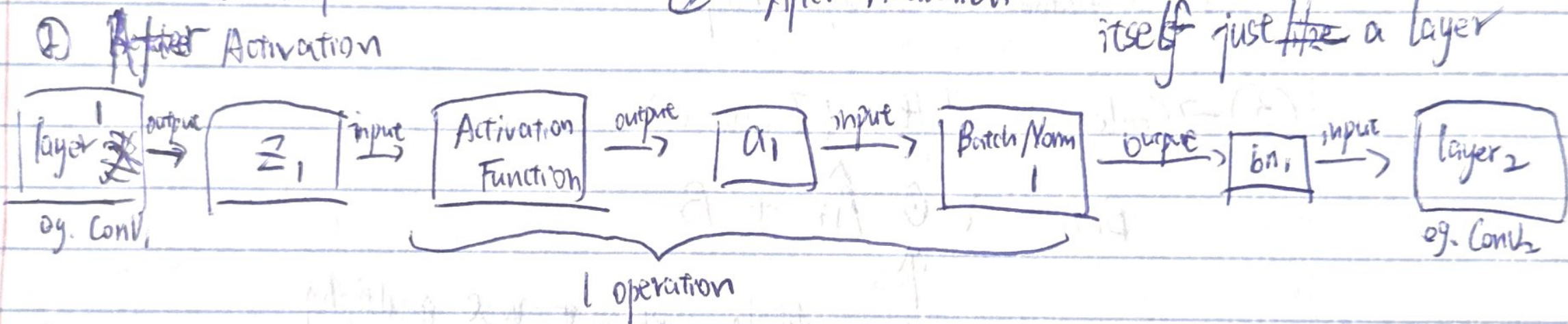benefit: allows larger learning rate, loss exploding/vanishing gradient

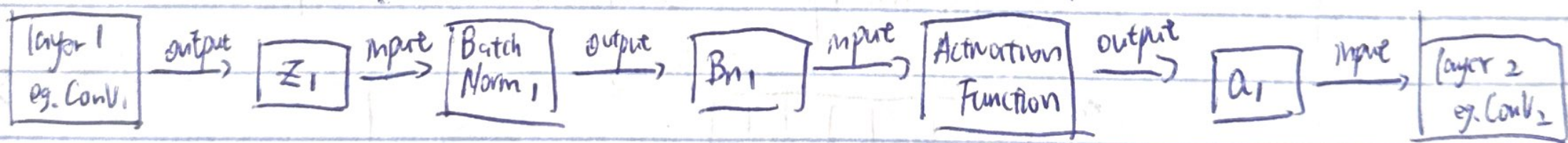where to use: normalize the input for fist layer and each (hidden) layers

① before activation   ② After Activation

itself just like a layer

① After Activation

layer 1 $\xrightarrow{output}$ $Z_1$ $\xrightarrow{input}$ Activation Function $\xrightarrow{output}$ $a_1$ $\xrightarrow{input}$ Batch Norm 1 $\xrightarrow{output}$ $bn_1$ $\xrightarrow{input}$ layer 2

eg. Conv₁                                                                                                                    eg. Conv₂

1 operation

② Before Activation

layer 1 $\xrightarrow{output}$ $Z_1$ $\xrightarrow{input}$ Batch Norm 1 $\xrightarrow{output}$ $Bn_1$ $\xrightarrow{input}$ Activation Function $\xrightarrow{output}$ $a_1$ $\xrightarrow{input}$ layer 2

eg. Conv₁                                                                                                                    eg. Conv₂

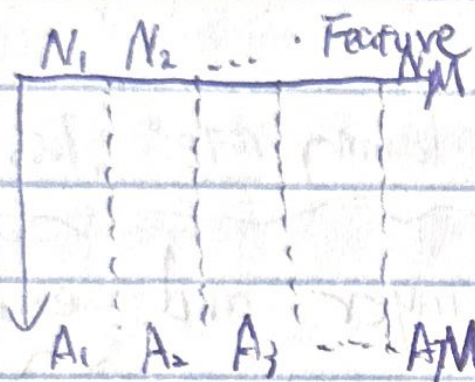Parameters ( Each Batch Norm Layer has its own copy of parameters )

learnable : $\beta$ , $\gamma$

saved : mean $\mu$ , var $\sigma$

Input

channels

CNN ( N , C , H , W )   N = # of image in one mini-batch,   C = # of feature maps (~~filter~~)

## Application Process

**①** Mini-Batch $\longrightarrow$ **②** Normalize

$(N, C, H, W)$

$\mu_i = \frac{1}{M} \sum A_i = \frac{1}{N \times H \times W} \sum_{N} \sum_{H} \sum_{W} x_{N,H,W}$

each sample is a feature

$\sigma_i = \sqrt{\frac{1}{M} \sum (A_i - \mu)^2}$

$\hat{A}_i = \frac{A_i - \mu_i}{\sigma_i}$

$\longrightarrow$ Feature

$\mu \quad \mu_1, \mu_2, \cdots, \mu_M$

$\sigma \quad \sigma_1, \sigma_2, \cdots, \sigma_M$

| $N_1$ | $N_2$ | $\cdots$ | Feature $M$ |
|---|---|---|---|

$\downarrow$

$A_1 \quad A_2 \quad A_3 \quad \cdots \quad A_M$

each feature is a vector/matrix map.

| $\hat{A}_1$ | $\hat{A}_2$ | | $\hat{A}_M$ |
|---|---|---|---|

Normalized value Now have 0 mean and unit variance / std (1)

**③** $\longrightarrow$ Scale and Shift (Innovation)

$$BN_i = \gamma \odot \hat{A}_i + \beta$$

$\uparrow$

⊘ not matrix multiply, element-wise multiply

⊘ not hyperparameters, but trainable parameters

$\longrightarrow$ Features

$\gamma$ | | | | | | | |

$\beta$ | | | | | | | |

| $N_1$ | $N_2$ | $N_3$ | $\cdots$ | $N_M$ | Features |
|---|---|---|---|---|---|

$\downarrow$

$BN_1 \quad BN_2 \quad \cdots \quad BN_M$