# Batch Norm computational graph:-

$$X \xrightarrow[(N,D)]{} \boxed{\begin{array}{c} \text{Batch Norm} \\ BN(r, \beta) \end{array}} \xrightarrow[(N,D)]{Y} \text{Some transf.} \xrightarrow{L} (1,)$$

$$\mu = \frac{1}{N} \sum_{i=1}^{N} x_i^o \qquad \sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i^o - \mu)^2$$

$$\uparrow (D,) \qquad\qquad\qquad \uparrow (D,)$$

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \qquad y = r\hat{x} + \beta \qquad \text{where } r, \beta \text{ are learnable params of a BN layer}$$

$$\uparrow (N,D) \qquad\qquad\qquad \uparrow \quad \uparrow \\ (N,D) \;\; (D,)$$

**problem:-** $\frac{\partial L}{\partial Y}$ is known/available, of shape $(N,D)$

calculate $\underset{(N,D)}{\frac{\partial L}{\partial x}}, \underset{(D,)}{\frac{\partial L}{\partial r}}, \underset{(D,)}{\frac{\partial L}{\partial \beta}}$

$$\frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial \beta} \rightarrow \text{since } \frac{\partial Y}{\partial \beta} = 1, \quad \frac{\partial L}{\partial \beta} = \sum_{i=1}^{N} \frac{\partial L}{\partial y_i}, \text{ where } y_i \text{ is the } i\text{-th row of } Y$$
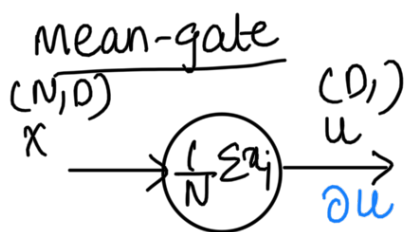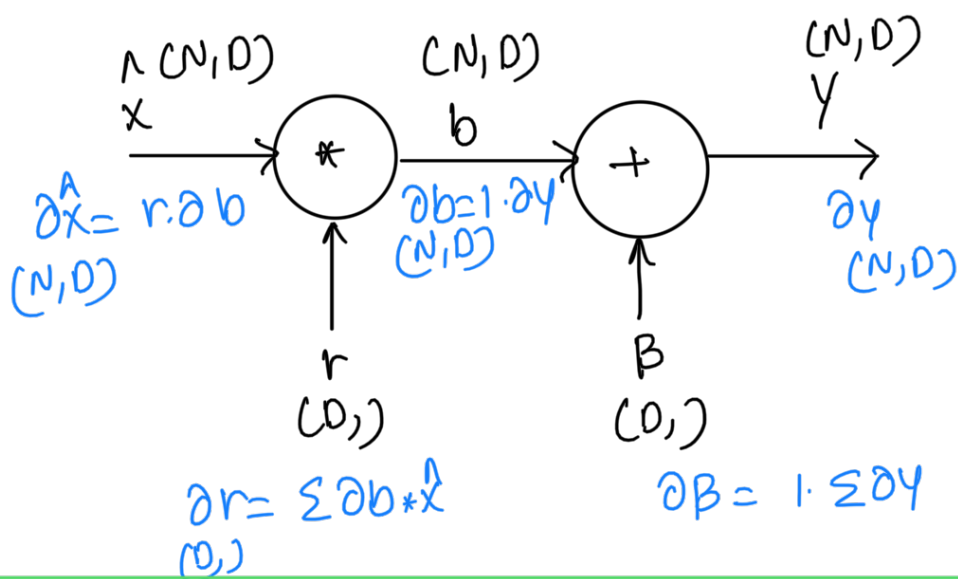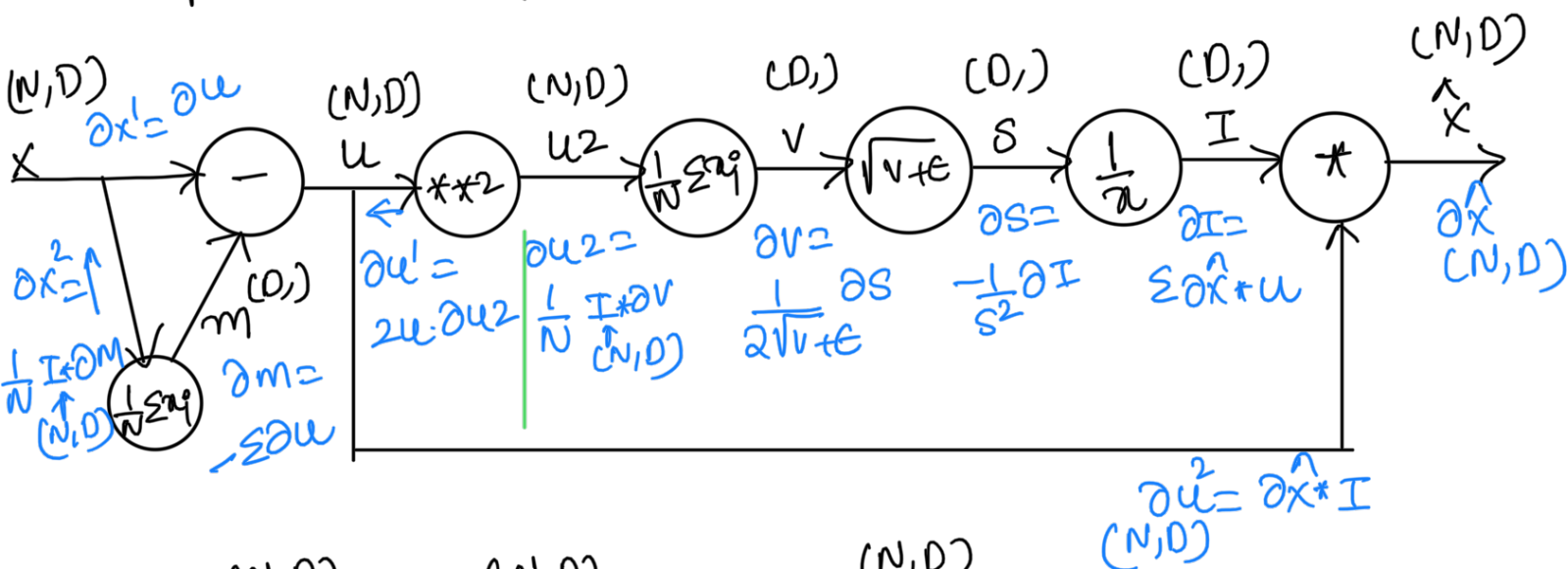
similarly, $\frac{\partial L}{\partial r} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial r}$ and since $\frac{\partial Y}{\partial r} = \hat{x}$,

$$\frac{\partial L}{\partial r} = \sum_{i=1}^{N} \frac{\partial L}{\partial y_i} * \hat{x}_i \quad \text{where } y_i \text{ \& } \hat{x}_i \text{ are } i\text{-th rows of } Y \text{ \& } \hat{X}, \text{ respectively.}$$
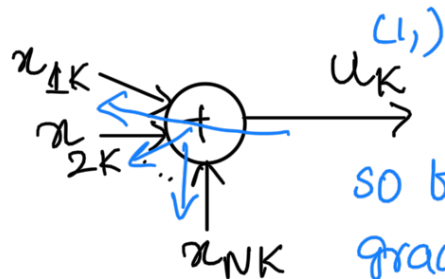
$\frac{\partial L}{\partial r}$ & $\frac{\partial L}{\partial \beta}$ are relatively straight forward to calculate.

$\frac{\partial L}{\partial x}$ is easier to calculate using a computational graph.

a computational graph is shown below :-



Top row graph:

$(N,D)$ — $x$   $\partial x' = \partial u$

$(N,D)$ — $u$

$(N,D)$ — $u2$

$(D,)$ — $v$

$(O,)$ — $s$

$(O,)$ — $I$

$(N,D)$ — $\hat{x}$

Nodes: $(-)$ , $(**2)$ , $(\frac{1}{N}\Sigma_i)$ , $(\sqrt{v+\epsilon})$ , $(\frac{1}{x})$ , $(*)$

$\partial x^2 = I$

$\frac{1}{N} I \cdot \partial M$
$(N,D) \quad \frac{1}{N}\Sigma_i$

$m \quad (D,)$

$\partial m = -\Sigma \partial u$

$\partial u' = 2u \cdot \partial u2$

$\partial u2 = \frac{1}{N} I \cdot \partial v$
$(N,D)$

$\partial v = \frac{1}{2\sqrt{v+\epsilon}} \partial s$

$\partial s = \frac{-1}{s^2} \partial I$

$\partial I = \Sigma \partial \hat{x} * u$

$\partial \hat{x}$
$(N,D)$

$\partial u2 = \partial \hat{x} * I$
$(N,D)$

$(N,D)$
$\hat{x}$
$\partial x$
$(N,D)$

Second row graph:

$\hat{x}$ $(N,D)$   $b$ $(N,D)$   $y$ $(N,D)$

Nodes: $(*)$ , $(+)$

$\partial \hat{x} = r \cdot \partial b$
$(N,D)$

$\partial b = 1 \cdot \partial y$
$(N,D)$

$\partial y$
$(N,D)$

$r$ $(D,)$

$B$ $(D,)$

$\partial r = \Sigma \partial b * \hat{x}$
$(D,)$

$\partial B = 1 \cdot \Sigma \partial y$

---

**mean-gate**

$(N,D)$ — $x$   node $(\frac{1}{N}\Sigma_i)$   $(D,)$ — $u$ , $\partial u$

mean-gate can be thought of as D sum-gates in parallel, each with N inputs corresponding to the N examples of a given feature.



$x_{1K}$ , $x_{2K}$ ... $x_{NK}$   node $(+)$   $u_K$ $(1,)$

$K = \{1 ... D\}$, for D features

so for each feature $u_k$, the gradient is copied to the N samples

$\Rightarrow$ essentially copy $u$ vector N times so it is of shape $(N,D)$ and divide by N

$\partial x = \frac{1}{N} I \cdot \partial u$
$(N,D)$