

1 Derivation of Batchnorm gradients

Let X of size (m, D) contain the data for a single batch of training. Let μ_B be the mean of that batch of size $(1 \times D)$ and σ^2 of size $(1 \times D)$ be the variance. \bar{X} is the batch normalized output.

Notation: Data_k refers to the k^{th} row of Data matrix.

$$\bar{X} = \frac{X - \mu_B}{\sqrt{\sigma^2 + \epsilon}}$$

$$y = \gamma \bar{X} + \beta$$

$$\frac{\partial y}{\partial \bar{X}} = \gamma$$

We know that:

$$\bar{X}_i = \frac{X_i - \mu_B}{\sqrt{\sigma^2 + \epsilon}}$$

$$\frac{\partial \bar{X}_i}{\partial X_k} = \frac{\frac{\partial X_i}{\partial X_k} - \frac{\partial \mu_B}{\partial X_k}}{\sqrt{\sigma^2 + \epsilon}} - \frac{X_i - \mu_B}{2^{3/2} \sqrt{\sigma^2 + \epsilon}} \frac{\partial \sigma^2}{\partial X_k} \quad (1)$$

$$\frac{\partial \mu_B}{\partial X_k} = \frac{1}{m} \quad (2)$$

We can express the variance in terms of mean using the following equation.

$$\sigma^2 = \sum_{k=0}^{m-1} \frac{X_k^2}{m} - \mu_B^2$$

$$\frac{\partial \sigma^2}{\partial X_k} = \frac{2X_k}{m} - 2\mu_B \frac{\partial \mu_B}{\partial X_k}$$

From (2), we get

$$\frac{\partial \sigma^2}{\partial X_k} = \frac{2(X_k - \mu_B)}{m}$$

Substituting this into (1) we have:

$$\frac{\partial \bar{X}_i}{\partial X_k} = \frac{\frac{\partial X_i}{\partial X_k} - \frac{1}{m}}{\sqrt{\sigma^2 + \epsilon}} - \frac{(X_i - \mu_B)(X_k - \mu_B)}{m^{3/2} \sqrt{\sigma^2 + \epsilon}}$$

The above equation can be expressed in terms of \bar{X} as

$$\begin{aligned} \frac{\partial \bar{X}_i}{\partial X_k} &= \frac{m \frac{\partial X_i}{\partial X_k} - 1}{m \sqrt{\sigma^2 + \epsilon}} - \frac{\bar{X}_i \bar{X}_k}{m \sqrt{\sigma^2 + \epsilon}} \\ \frac{\partial \bar{X}_i}{\partial X_k} &= \frac{m \frac{\partial X_i}{\partial X_k} - 1 - \bar{X}_i \bar{X}_k}{m \sqrt{\sigma^2 + \epsilon}} \end{aligned} \quad (3)$$

We need to find $\frac{\partial L}{\partial \bar{X}_j}$ where L is the loss. From chain rule, we have:

$$\frac{\partial L}{\partial X_j} = \frac{\partial L}{\partial \bar{X}_0} \frac{\partial \bar{X}_0}{\partial X_j} + \frac{\partial L}{\partial \bar{X}_1} \frac{\partial \bar{X}_1}{\partial X_j} + \dots + \frac{\partial L}{\partial \bar{X}_{m-1}} \frac{\partial \bar{X}_{m-1}}{\partial X_j}$$

Which can be reduced as:

$$\frac{\partial L}{\partial X_j} = \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} \frac{\partial \bar{X}_k}{\partial X_j}$$

Substituting result from eq (3), we have:

$$\begin{aligned} \frac{\partial L}{\partial X_j} &= \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} \left(\frac{m \frac{\partial X_k}{\partial X_j} - 1 - \bar{X}_j \bar{X}_k}{m \sqrt{\sigma^2 + \epsilon}} \right) \\ \frac{\partial L}{\partial X_j} &= \frac{1}{m \sqrt{\sigma^2 + \epsilon}} \sum_{k=0}^{m-1} \left(\frac{\partial L}{\partial \bar{X}_k} (m \frac{\partial X_k}{\partial X_j} - 1 - \bar{X}_j \bar{X}_k) \right) \\ \frac{\partial L}{\partial X_j} &= \frac{1}{m \sqrt{\sigma^2 + \epsilon}} \left(m \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} \frac{\partial X_k}{\partial X_j} - \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} - \bar{X}_j \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} \bar{X}_k \right) \end{aligned}$$

We see that $\frac{\partial X_k}{\partial X_j}$ is 0 everywhere except when $j=k$. So that part of the summation reduces to summing all rows where $j=k$, which is just one row corresponding to $\frac{\partial L}{\partial \bar{X}_j}$. So the equation reduces to:

$$\frac{\partial L}{\partial X_j} = \frac{1}{m \sqrt{\sigma^2 + \epsilon}} \left(m \frac{\partial L}{\partial \bar{X}_j} - \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} - \bar{X}_j \sum_{k=0}^{m-1} \frac{\partial L}{\partial \bar{X}_k} \bar{X}_k \right) \quad (4)$$

We receive the gradient of y from the function as dout. so

$$\frac{\partial L}{\partial y} = \text{dout}$$

$$\frac{\partial L}{\partial \bar{X}} = \gamma \frac{\partial L}{\partial y}$$

We use the matrix computed from above to substitute in equation (4) to get the loss w.r.t X for all its rows.