

Batch Normalization : Back-propagation

$$X = \begin{matrix} & x_{11} & x_{12} & x_{13} \\ x_{21} & & x_{22} & x_{23} \end{matrix}$$

$$\mu = \frac{x_{11} + x_{21}}{2}, \quad \frac{x_{12} + x_{22}}{2}, \quad \frac{x_{13} + x_{23}}{2}$$

$$= \frac{1}{N} \sum_{i=1}^N x_{i1}, \quad \frac{1}{N} \sum_{i=1}^N x_{i2}, \quad \frac{1}{N} \sum_{i=1}^N x_{i3}$$

$$\sigma^2 = \frac{1}{2} ((x_{11} - \mu_1)^2 + (x_{21} - \mu_1)^2), \quad \frac{1}{2} [(x_{12} - \mu_2)^2 + (x_{22} - \mu_2)^2],$$

$$\frac{1}{2} [(x_{13} - \mu_3)^2 + (x_{23} - \mu_3)^2]$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_{i1} - \mu_1)^2, \quad \frac{1}{N} \sum_{i=1}^N (x_{i2} - \mu_2)^2, \quad \dots$$

$$\hat{x} = \frac{x_{11} - \mu_1}{\sqrt{\sigma_1^2 + \epsilon}}, \quad \frac{x_{12} - \mu_2}{\sqrt{\sigma_2^2 + \epsilon}}, \quad \frac{x_{13} - \mu_3}{\sqrt{\sigma_3^2 + \epsilon}}$$

$$\frac{x_{21} - \mu_1}{\sqrt{\sigma_1^2 + \epsilon}}, \quad \frac{x_{22} - \mu_2}{\sqrt{\sigma_2^2 + \epsilon}}, \quad \frac{x_{23} - \mu_3}{\sqrt{\sigma_3^2 + \epsilon}}$$

$$y = r_1 \hat{x}_{11} + \beta_1, \quad r_2 \hat{x}_{12} + \beta_2, \quad r_3 \hat{x}_{13} + \beta_3$$

$$r_2 \hat{x}_{21} + \beta_1, \quad r_2 \hat{x}_{22} + \beta_2, \quad r_3 \hat{x}_{23} + \beta_3$$

$$\frac{dl}{dy} = \frac{dl}{dy_{11}}, \quad \frac{dl}{dy_{12}}, \quad \frac{dl}{dy_{13}}$$

$$\frac{dl}{dy_{21}}, \quad \frac{dl}{dy_{22}}, \quad \frac{dl}{dy_{23}}$$

upstream x local gradient

$$\frac{dL}{dr_i} = \sum_{i,j} \frac{dL}{dy_{ij}} \times \frac{dy_{ij}}{dr_i} = \frac{dL}{dy_{11}} \cdot \hat{x}_{11} + \frac{dL}{dy_{21}} \cdot \hat{x}_{21} = \frac{dL}{dx_i}$$

$$\frac{dL}{dr} = \left(\frac{dL}{dy} \odot \hat{x} \right) \rightarrow \text{sum along axis} = 0$$

$$\frac{dL}{dr_j} = \sum_{i=1}^N \hat{x}_{ij} \cdot \frac{dL}{dy_{ij}}$$

\odot elementwise multiplication

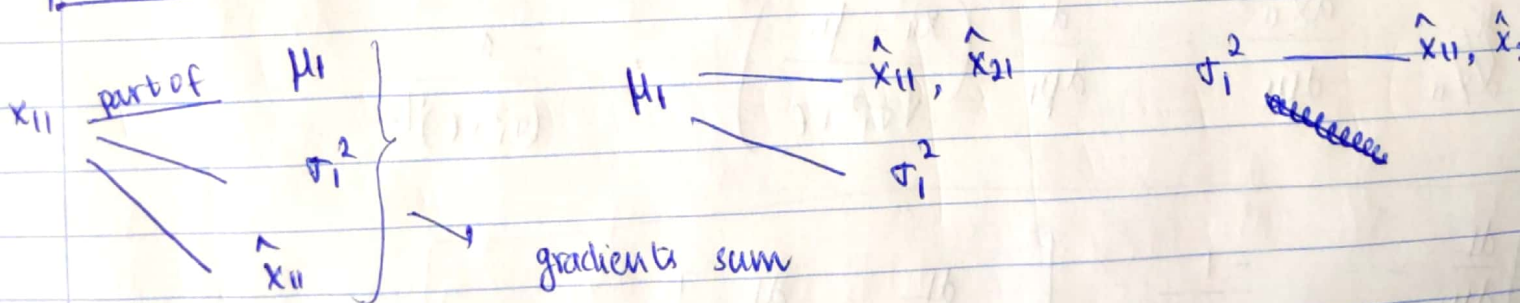
$$\frac{dL}{db_i} = \sum_{i,j} \frac{dL}{dy_{ij}} \frac{dy_{ij}}{db_i} = \frac{dL}{dy_{11}} + \frac{dL}{dy_{12}}$$

$$\therefore \frac{dL}{db_j} = \sum_{i=1}^N \frac{dL}{dy_{ij}}$$

$$\frac{dL}{db} = \frac{dL}{dy} \rightarrow \text{sum along axis} = 0$$

$$\frac{dL}{d\hat{x}_{11}} = \sum \frac{dL}{dy_{ij}} \frac{dy_{ij}}{d\hat{x}_{11}} = \frac{dL}{dy_{11}} \times r_i$$

$$\therefore \frac{dL}{d\hat{x}} = \frac{dL}{dy} \odot r$$



$$\frac{dL}{d\sigma_1^2} = \sum \frac{dL}{d\hat{x}_{ij}} \frac{d\hat{x}_{ij}}{d\sigma_1^2} = \frac{dL}{d\hat{x}_{11}} \cdot \frac{d\hat{x}_{11}}{d\sigma_1^2} + \frac{dL}{d\hat{x}_{21}} \cdot \frac{d\hat{x}_{21}}{d\sigma_1^2}$$

$$\frac{dL}{d\hat{x}_{11}} \cdot \frac{d\hat{x}_{11}}{d\sigma_1^2} = \frac{d}{d\sigma_1^2} \frac{x_{11} - \mu_1}{(\sigma_1^2 + \epsilon)^{0.5}} = -\frac{1}{2} (\sigma_1^2 + \epsilon)^{-3/2} (x_{11} - \mu_1)$$

$$\boxed{\frac{dL}{d\sigma_1^2} = -\frac{1}{2} (\sigma_1^2 + \epsilon)^{-3/2} \left[\frac{dL}{d\hat{x}_{11}} (x_{11} - \mu_1) + \frac{dL}{d\hat{x}_{21}} (x_{21} - \mu_1) \right]} \quad \text{--- A}_1$$

$$= -\frac{1}{2} (\sigma_1^2 + \epsilon)^{-1} \left[\frac{dL}{d\hat{x}_{11}} \frac{x_{11} - \mu_1}{(\sigma_1^2 + \epsilon)^{1/2}} + \frac{dL}{d\hat{x}_{21}} \frac{(x_{21} - \mu_1)}{\sqrt{\sigma_1^2 + \epsilon}} \right]$$

$$\frac{dL}{d\sigma_1^2} = -\frac{1}{2} (\sigma_1^2 + \epsilon)^{-1} \left[\frac{dL}{d\hat{x}_{11}} \hat{x}_{11} + \frac{dL}{d\hat{x}_{21}} \hat{x}_{21} \right] \quad \text{--- A}_2$$

$$\boxed{\frac{dL}{d\sigma_1^2} = -\frac{1}{2} (\sigma_1^2 + \epsilon)^{-1} \left[\sum_{i=1}^N \frac{dL}{d\hat{x}_{ii}} \hat{x}_{ii} \right]}$$

$$\frac{dL}{d\mu_1} = \sum \frac{dL}{d\hat{x}_{ij}} \frac{d\hat{x}_{ij}}{d\mu_1} = \frac{dL}{d\hat{x}_{11}} \frac{d\hat{x}_{11}}{d\mu_1} + \frac{dL}{d\hat{x}_{21}} \frac{d\hat{x}_{21}}{d\mu_1}$$

$$\frac{dL}{d\hat{x}_{11}} \cdot \frac{d\hat{x}_{11}}{d\mu_1} = \frac{d}{d\mu_1} \left(\frac{x_{11} - \mu_1}{\sqrt{\sigma_1^2 + \epsilon}} \right) = \frac{-1}{(\sigma_1^2 + \epsilon)^{1/2}}$$

$$\boxed{\frac{dL}{d\mu_1} = \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \left[\frac{dL}{d\hat{x}_{11}} + \frac{dL}{d\hat{x}_{21}} \right]} \quad \text{--- b}_1$$

$$\frac{d\sigma_1^2}{d\mu_1} = \frac{d}{d\mu_1} \left[\frac{x_{11} - \mu_1}{(\sigma_1^2 + \epsilon)^{0.5}} + \frac{x_{21} - \mu_1}{(\sigma_1^2 + \epsilon)^{0.5}} \right]$$

$$= \frac{d}{d\mu} \left(\frac{1}{2} \left[(x_{11} - \mu_1)^2 + (x_{21} - \mu_1)^2 \right] \right)$$

$$= -\frac{2}{N} \left[(x_{11} - \mu_1) + (x_{21} - \mu_1) \right]$$

$$\boxed{\frac{d\sigma_1^2}{d\mu_1} = -\frac{2}{N} \sum_{i=1}^N (x_{i1} - \mu_1)}$$

$$\begin{aligned} \therefore \frac{dL}{d\mu_1} &= b_1 + \frac{dL}{d\sigma_1^2} \cdot \frac{d\sigma_1^2}{d\mu_1} + \frac{-1}{2(\sigma_1^2 + \epsilon)} \left[\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \right] \cdot \frac{-2}{N} \sum_{i=1}^N (x_{i1} - \mu_1) \\ &= \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \left[\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \right] + \left[\frac{-1}{2(\sigma_1^2 + \epsilon)} \left[\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \right] \cdot \frac{-2}{N} \sum_{i=1}^N (x_{i1} - \mu_1) \right] \end{aligned}$$

$$\frac{dL}{d\mu_1} = \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \right) + \frac{1}{N} \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \hat{x}_{i1} \right) \cdot \sum_{i=1}^N (x_{i1} - \mu_1)$$

$$\boxed{\frac{dL}{d\mu_1} = \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} + \frac{1}{N} \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} \hat{x}_{i1} \right) \sum_{i=1}^N (x_{i1} - \mu_1) - B_1}$$

Consider $\sum_{i=1}^N (x_{i1} - \mu_1) = N \cdot \mu_1 - N\mu_1 = 0$

$$\boxed{\frac{dL}{d\mu_1} = \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{i1}} - B_2}$$

$$\frac{dL}{dx_{ii}} = \sum_{i \neq j} \frac{dx_{ij}}{d\hat{x}_{ij}} \frac{d\hat{x}_{ij}}{dx_{ii}} = \frac{dL}{d\hat{x}_{ii}} \cdot \frac{d\hat{x}_{ii}}{dx_{ii}}$$

$$\boxed{\frac{dL}{d\hat{x}_{ii}} = \frac{dL}{d\hat{x}_{ii}} \cdot \frac{1}{(\sigma_i^2 + \epsilon)^{0.5}}}$$

$$\boxed{\frac{d\mu_i}{dx_{ii}} = \frac{1}{N}}$$

$$\frac{d\sigma_i^2}{dx_{ii}} = \frac{2}{N} \underbrace{\sum_{i=1}^N (x_{ii} - \mu_i)}_0 + \frac{-2}{N} \sum_{i=1}^N (x_{ii} - \mu_i) \cdot \frac{d\mu_i}{dx_{ii}} = \frac{-2}{N} \sum_{i=1}^N (x_{ii} - \mu_i) \frac{d\mu_i}{dx_{ii}}$$

$$\frac{dL}{dx_{ii}} = \frac{dL}{d\hat{x}_{ii}} \cdot \frac{1}{(\sigma_i^2 + \epsilon)^{0.5}} + \frac{dL}{d\mu_i} \frac{d\mu_i}{dx_{ii}} + \frac{dL}{d\sigma_i^2} \times \frac{d\sigma_i^2}{dx_{ii}}$$

$$= \frac{dL}{d\hat{x}_{ii}} \cdot \frac{1}{(\sigma_i^2 + \epsilon)^{0.5}} + \frac{-1}{(\sigma_i^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{ii}} \cdot \frac{1}{N} \quad \text{and} \quad \frac{0}{N} \sum_{i=1}^N$$

$$+ \left(-\frac{1}{2}\right)(\sigma_i^2 + \epsilon)^{-1} \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{ii}} \hat{x}_{ii} \right) \cdot \frac{-2}{N} \sum_{i=1}^N (x_{ii} - \mu_i) \frac{d\mu_i}{dx_{ii}} \left(\frac{1}{N} \right)$$

$$\frac{dL}{dx_{ii}} = \frac{dL}{d\hat{x}_{ii}} \frac{1}{(\sigma_i^2 + \epsilon)^{0.5}} - \frac{1}{N(\sigma_i^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{ii}} + \frac{1}{N^2(\sigma_i^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{ii}} \hat{x}_{ii} \cdot \sum_{i=1}^N \hat{x}_{ii}$$

$$= \frac{1}{N(\sigma_1^2 + \epsilon)^{0.5}} \left[\cancel{N \cdot \frac{dL}{d\hat{x}_{11}}} - \cancel{\sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}}} + \frac{1}{N} \sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}} \hat{x}_{1i} - \cancel{\sum_{i=1}^N \hat{x}_{1i}} \right]$$

$$\frac{d\sigma_1^2}{dx_{11}} = \frac{2}{N} \sum_{i=1}^N (x_{1i} - \mu_1) + \frac{-2}{N} \sum_{i=1}^N (x_{1i} - \mu_1) \cdot \left(\frac{d\mu_1}{dx_{11}} \right) = \frac{1}{N}$$

$$\frac{d\sigma_1^2}{dx_{11}} = 0 + \frac{-2}{N^2} \sum_{i=1}^N (x_{1i} - \mu_1) = 0$$

$$\therefore \frac{dL}{dx_{11}} = \frac{dL}{d\hat{x}_{11}} \cdot \frac{1}{(\sigma_1^2 + \epsilon)^{0.5}} + \frac{1}{N} \left[\right]$$

$$\frac{d\sigma_1^2}{dx_{11}} = \frac{d}{dx_{11}} \frac{1}{N} \sum_{i=1}^N (x_{1i} - \mu_1)^2$$

$$= \frac{2}{N} (x_{11} - \mu_1) + \left[\cancel{\frac{-2}{N} \sum_{i=1}^N (x_{1i} - \mu_1) \cdot \frac{d\mu_1}{dx_{11}}} \right]$$

$$= \frac{2}{N} (x_{11} - \mu_1) + \frac{-2}{N^2} \left(\sum_{i=1}^N x_{1i} - \mu_1 \right) \quad \because \frac{d\mu_1}{dx_{11}} = \frac{1}{N}$$

we don't take $\frac{d\sigma_1^2}{dx_{11}} \times \frac{d\mu_1}{dx_{11}} - \frac{dx_{11}}{d\mu_1}$ since we already computed $\frac{dL}{dx_{11}}$ wrt μ_1

$$\boxed{\frac{d\sigma_1^2}{dx_{11}} = \frac{2}{N} (x_{11} - \mu_1)}$$

$$\therefore \frac{dL}{dx_{11}} = \frac{dL}{d\hat{x}_{11}} \cdot \frac{1}{(\sigma_1^2 + \epsilon)^{0.5}} + \frac{-1}{(\sigma_1^2 + \epsilon)^{0.5}} \sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}} \cdot \frac{1}{N}$$

$$= \frac{-1}{2(\sigma_1^2 + \epsilon)} \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}} \hat{x}_{1i} \right) \cdot \frac{2}{N} (x_{11} - \mu_1)$$

$$\boxed{\frac{dL}{dx_{11}} = \frac{1}{N(\sigma_1^2 + \epsilon)^{0.5}} \left[\frac{dL}{d\hat{x}_{11}} - \sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}} - \left(\sum_{i=1}^N \frac{dL}{d\hat{x}_{1i}} \hat{x}_{1i} \right) \cdot \hat{x}_{11} \right] - C}$$