

Batch normalisation: Derivation for alternative backward propagation proof.

let us define

$$X = \begin{Bmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{Bmatrix} = \begin{Bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M,1} & x_{M,2} & \dots & x_{M,D} \end{Bmatrix} \in [M \times D]$$

$$\mu = [\mu_1, \mu_2, \dots, \mu_D] = \left[\frac{1}{M} \sum_{k=1}^M x_{k,1}, \frac{1}{M} \sum_{k=1}^M x_{k,2}, \dots, \frac{1}{M} \sum_{k=1}^M x_{k,D} \right] \in [1 \times D]$$

$$\sigma^2 = [\sigma_1, \sigma_2, \dots, \sigma_D] = \left[\frac{1}{M} \sum_{k=1}^M (x_{k,1} - \mu_1)^2, \frac{1}{M} \sum_{k=1}^M (x_{k,2} - \mu_2)^2, \dots, \frac{1}{M} \sum_{k=1}^M (x_{k,D} - \mu_D)^2 \right] \in [1 \times D]$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \in [1 \times D]$$

$$y_i = \gamma \hat{x}_i + \beta \in [1 \times D]$$

Taking into consideration the method from paper,
Notice that,

$$\begin{aligned} \frac{dl}{d\mu} &= \dots + \frac{dl}{d\sigma^2} \sum_{k=1}^M -\frac{2(x_i - \mu)}{M} = \dots + \frac{dl}{d\sigma^2} -\frac{2}{M} \sum_{k=1}^M (x_i - \mu) \\ &= \dots + \frac{dl}{d\sigma^2} -\frac{2}{M} \left(\sum_{k=1}^M x_i \right) - M\mu \\ &= \dots + \frac{dl}{d\sigma^2} -\frac{2}{M} \times 0 = 0 \dots + 0. \end{aligned}$$

Then we can get

$$\frac{dl}{d\hat{x}_i} = \frac{dl}{d\hat{x}_i} \cdot \frac{1}{\sqrt{\sigma^2 + \varepsilon}} + \left[\sum_{k=1}^M (\hat{x}_k - \mu) \left(\frac{-1}{2} \right) (\sigma^2 + \varepsilon)^{-3/2} \right] \cdot \frac{2(\hat{x}_i - \mu)}{M}$$

$$+ \sum_{k=1}^M \frac{dl}{d\hat{x}_k} \frac{-1}{\sqrt{\sigma^2 + \varepsilon}} \cdot \frac{1}{M}$$

$$= \frac{dl}{d\hat{x}_i} \cdot \frac{1}{\sqrt{\sigma^2 + \varepsilon}} - \frac{(\sigma^2 + \varepsilon)^{-3/2}}{M} \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} (\hat{x}_k - \mu) \right] (\hat{x}_i - \mu)$$

$$- \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \sum_{k=1}^M \frac{dl}{d\hat{x}_k} \frac{2(\hat{x}_k - \mu)}{M} \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} (\hat{x}_k - \mu) \right] (\hat{x}_i - \mu)$$

$$= \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} (\hat{x}_k - \mu) \right] (\hat{x}_i - \mu)$$

$$= \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} \hat{x}_k \right] \hat{x}_i$$

$$= \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \sum_{k=1}^M \frac{dl}{d\hat{x}_k} \hat{x}_k - \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} \right] \hat{x}_i$$

$$= \frac{1}{M \sqrt{\sigma^2 + \varepsilon}} \sum_{k=1}^M \frac{dl}{d\hat{x}_k} \hat{x}_k$$

$$= \sum_{k=1}^M \frac{dl}{d\hat{x}_k} \hat{x}_k - \left[\sum_{k=1}^M \frac{dl}{d\hat{x}_k} \right] \hat{x}_i$$