

682 - Assignment - 2

Batch Normalization \Rightarrow alternative backward

Derivation

required: $-\frac{\partial l}{\partial x_i}, \frac{\partial l}{\partial \gamma}, \frac{\partial l}{\partial \beta}$ with $\frac{\partial l}{\partial \hat{x}_i}, \frac{\partial l}{\partial \mu_B}, \frac{\partial l}{\partial \sigma_B^2}$

μ_B - Mean of $(1 \times D)$ sized batch

σ_B^2 - variance \hat{x}_i - normalized i/p

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad y = \gamma \hat{x}_i + \beta$$

$$\Rightarrow \frac{\partial y}{\partial \hat{x}_i} = \gamma \quad \frac{\partial \hat{x}_i}{\partial \mu_B} \Rightarrow \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \quad (\text{chain rule})$$

Accordingly $\sigma^2 = \text{variance} = \frac{1}{m} \sum_{i=1}^m (\hat{x}_i - \mu_B)^2$

$$\frac{\partial \sigma_B^2}{\partial \mu_B} = \frac{1}{m} \sum_{i=1}^m -2(\hat{x}_i - \mu_B)$$

By chain rule $\Rightarrow \frac{\partial l}{\partial \sigma_B^2} = \frac{\partial l}{\partial \hat{x}_i} \times \frac{\partial \hat{x}_i}{\partial \sigma_B^2}$

We know that $\hat{x}_i = \frac{(x_i - \mu_B)}{\sqrt{\sigma_B^2 + \epsilon}}$

$$\Rightarrow \frac{\partial \hat{x}_i}{\partial \sigma_B^2} \Rightarrow \frac{\partial \hat{x}_i}{\partial \sigma_B^2} (\hat{x}_i - \mu_B) (\sigma_B^2 + \epsilon)^{-3/2}$$

$$\Rightarrow -\frac{1}{2} \sum_{i=1}^m (\hat{x}_i - \mu_B) (\sigma_B^2 + \epsilon)^{-3/2}$$

①

We require $\frac{\partial l}{\partial \mu_B} = \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \mu_B} + \frac{\partial l}{\partial \sigma_B^2} \cdot \frac{\partial \sigma_B^2}{\partial \mu_B}$

$$\Rightarrow \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \left(\frac{\partial l}{\partial \sigma_B^2} \cdot \frac{1}{m} \sum_{i=1}^m -2(\hat{x}_i - \mu_B) \right)$$

$$\Rightarrow \left(\sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \left(-2 \frac{\partial \ell}{\partial \sigma_B^2} \left[\frac{1}{m} \sum_{i=1}^m \hat{x}_i - \frac{1}{m} \sum_{i=1}^m \mu_B \right] \right)$$

$$\Rightarrow \left(\sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \left(-2 \frac{\partial \ell}{\partial \sigma_B^2} \cdot \left(\mu_B - \frac{\mu_B \cdot m}{m} \right) \right)$$

$$\Rightarrow \left(\sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) \quad \text{--- (2)}$$

$$\text{Nom } \frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial \ell}{\partial \mu_B} \cdot \frac{\partial \mu_B}{\partial x_i} + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{\partial \sigma_B^2}{\partial x_i}$$

required

$$\frac{\partial \hat{x}_i}{\partial x_i} = (\sqrt{\sigma_B^2 + \epsilon})^{-1} \cdot \Delta \frac{\partial \mu_B}{\partial x_i} = \frac{1}{m}$$

$$\frac{\partial \sigma_B^2}{\partial x_i} = \frac{2}{m} (x_i - \mu_B)$$

$$\text{required } \Rightarrow \frac{\partial \ell}{\partial x_i} = \left(\frac{\partial \ell}{\partial \hat{x}_i} \cdot (\sigma_B^2 + \epsilon)^{-1/2} \right) + \left(\frac{\partial \ell}{\partial \mu_B} \cdot \frac{1}{m} \right)$$

$$\text{From (1), (2)} \quad \frac{\partial \ell}{\partial x_i} = \left(\frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{(\sigma_B^2 + \epsilon)^{1/2}} \right) + \left(\frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{2}{m} (x_i - \mu_B) \right)$$

$$\Rightarrow \frac{\partial \ell}{\partial x_i} = \left(\frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{(\sigma_B^2 + \epsilon)^{1/2}} \right) + \left(\frac{1}{m} \sum_{j=1}^m \frac{\partial \ell}{\partial \hat{x}_j} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right)$$

$$- \left(\frac{1}{2} \sum_{j=1}^m \frac{\partial \ell}{\partial \hat{x}_j} (x_j - \mu_B) (\sigma_B^2 + \epsilon)^{-1/2} \cdot \frac{2}{m} (x_i - \mu_B) \right)$$

$$\Rightarrow \left(\frac{\partial \ell}{\partial \hat{x}_i} \cdot (\sigma_B^2 + \epsilon)^{-1/2} \right) - \left((\sigma_B^2 + \epsilon)^{-1/2} \cdot \frac{1}{m} \sum_{j=1}^m \frac{\partial \ell}{\partial \hat{x}_j} \right) +$$

$$+ \left((\sigma_B^2 + \epsilon)^{-1/2} \cdot \frac{1}{m} \cdot \hat{x}_i \sum_{j=1}^m \frac{\partial l}{\partial \hat{x}_j} \cdot \hat{x}_j \right)$$

$$\Rightarrow \textcircled{1} \quad \frac{\partial l}{\partial x_i} = \frac{(\sigma_B^2 + \epsilon)^{-1/2}}{m} \left[m \frac{\partial l}{\partial \hat{x}_i} - \sum_{j=1}^m \frac{\partial l}{\partial \hat{x}_j} - \hat{x}_i \sum_{j=1}^m \frac{\partial l}{\partial \hat{x}_j} \cdot \hat{x}_j \right]$$

Hence Derived $\frac{\partial l}{\partial x_i}$

$$\frac{\partial l}{\partial \beta} = \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial \beta} \Rightarrow \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

$$\Rightarrow \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} \Rightarrow \frac{\partial l}{\partial y_i} \cdot \delta$$

$$\frac{\partial l}{\partial \gamma} = \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial \gamma} \Rightarrow \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \hat{x}_i$$

$$\textcircled{2} \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \hat{x}_i$$

$$\textcircled{3} \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

$$\textcircled{4} \quad \frac{\partial l}{\partial y_i} = \text{dout}$$