## Batch Normalization →

$$\mu_\ell = \sum_p^M x_{p\ell} \cdot \frac{1}{N} \quad ; \quad \sigma^2 = \frac{1}{N} \sum_p^M (x_{p\ell} - \mu_\ell)^2$$

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \varepsilon}} \longrightarrow \hat{x}_{k\ell} = \frac{x_{k\ell} - \mu_\ell}{\sqrt{\sigma_\ell^2 + \varepsilon}}$$

$$y = \gamma \hat{x} + \beta \longrightarrow \begin{pmatrix} y_{11} & y_{12} & \cdots \end{pmatrix} = \begin{pmatrix} \gamma_1 \hat{x}_{11} + \beta_1 & \gamma_2 \hat{x}_{12} + \beta_2 & \cdots \end{pmatrix}$$

$$y_{k\ell} = \gamma_\ell \hat{x}_{k\ell} + \beta_\ell \longrightarrow \frac{\partial y_{k\ell}}{\partial x_{ij}} = \gamma_\ell (\delta_{i,k} \, \delta_{j,k}) \longrightarrow \text{Element wise.}$$

$$\frac{\partial h}{\partial x_{ij}} = \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \frac{\partial y_{k\ell}}{\partial x_{ij}} = \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \frac{\partial y_{k\ell}}{\partial \hat{x}_{k\ell}} \cdot \frac{\partial \hat{x}_{k\ell}}{\partial x_{ij}}$$

* just translation →

$$\hat{x}_{k\ell} = x_{k\ell} - \mu_\ell = x_{k\ell} - \frac{1}{N} \sum_p^M x_{p\ell}$$

$$\frac{\partial \hat{x}_{k\ell}}{\partial x_{ij}} = \begin{vmatrix} 1 - \frac{1}{N} & \text{if} \quad k=i, \, \ell \neq j \\ -1/N & \text{if} \quad k \neq i, \, \ell = j \\ 0 & \text{Otherwise} \end{vmatrix} = \delta_{\ell j} \left( \delta_{i,k} - \frac{1}{N} \right)$$

* Variance $\longrightarrow$

$$\sigma_\ell^2 = \sum_p^M (x_{p\ell} - \mu_\ell)^2 \cdot \frac{1}{N}$$

$$\frac{\partial \sigma_\ell^2}{\partial x_{i'j}} = \frac{1}{N} \cdot 2 \sum_p^M \frac{\partial (x_{p\ell} - \mu_\ell)}{\partial x_{ij}} \cdot (x_{p\ell} - \mu_\ell) = \frac{2 \, \delta_{\ell,j}}{N} \sum_p^1 (x_{p\ell} - \mu_\ell)\left(\delta_{p,i} - \frac{1}{N}\right)$$

$$= \frac{2 \, \delta_{\ell,j}}{N}\left[ (x_{i\ell} - \mu_\ell) + \frac{1}{N}\sum_p^M \mu_\ell - \frac{1}{N}\sum_p^M x_{p\ell} \right] =$$

$$\underbrace{\frac{N}{N}\mu_\ell - \mu_\ell}$$

$$= \frac{2}{N} \, \delta_{\ell,j} \, (x_{i\ell} - \mu_\ell)$$

* Norm $\longrightarrow$

$$\frac{\partial \hat{x}_{\kappa\ell}}{\partial x_{ij}} = \frac{\partial (x_{\kappa\ell} - \mu_\ell)}{\partial x_{ij}}(\sigma^2 + \varepsilon)^{-1/2} - \frac{1}{2}\frac{\partial \sigma^2}{\partial x_{ij}}(x_{\kappa\ell} - \mu_\ell)(\sigma_\ell^2 + \varepsilon)^{-3/2}$$

$$= \delta_{\ell,j}\left(\delta_{\kappa,i} - \frac{1}{N}\right)(\sigma_\ell^2 + \varepsilon)^{-1/2} - \frac{1}{N}\delta_{\ell,j}(x_{i\ell} - \mu_\ell)(x_{\kappa\ell} - \mu_\ell)(\sigma_\ell^2 + \varepsilon)^{-3/2}$$

* Cost function $\longrightarrow$

$$\frac{\partial h}{\partial x_{ij}} = \sum_{\kappa=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{\kappa\ell}} \cdot \frac{\partial y_{\kappa\ell}}{\partial \hat{x}_{\kappa\ell}} \cdot \frac{\partial \hat{x}_{\kappa\ell}}{\partial x_{ij}} =$$

$$= \sum_{\kappa=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{\kappa\ell}} \cdot \gamma_\ell \cdot \delta_{\ell,j}(\sigma^2 + \varepsilon)^{-1/2}\left[ \delta_{\kappa,i} - \frac{1}{N} - \frac{1}{N}(x_{i\ell} - \mu_\ell)(x_{\kappa\ell} - \mu_\ell)(\sigma_\ell^2 + \varepsilon)^{-1/2} \right]$$

$$= \frac{\gamma_j}{N} (\tau_j^2 + \varepsilon)^{-1/2} \left[ \frac{\partial h}{\partial y_{ij}} - \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}} - \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}} (x_{ij} - \mu_j)(x_{kj} - \mu_j)(\tau_j^2 + \varepsilon)^{1/2} \right]$$

$$= \frac{\gamma_j}{N} (\tau_j^2 + \varepsilon)^{-1/2} \left[ \frac{\partial h}{\partial y_{ij}} - \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}} - (x_{ij} - \mu_j)(\tau_j^2 + \varepsilon)^{-1/2} \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}} (x_{kj} - \mu_j) \right]$$

$$\frac{\partial h}{\partial \gamma_{ij}} = \sum_{k=1}^{M} \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \frac{\partial y_{k\ell}}{\partial \gamma_{ij}} = \sum_{k=1}^{M} \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \delta_{\ell ij} \cdot \hat{x}_{k\ell} =$$

$$= \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}} \cdot \hat{x}_{kj}$$

$$\frac{\partial h}{\partial \beta_{ij}} = \sum_{k=1}^{M} \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \frac{\partial y_{k\ell}}{\partial \beta_{ij}} = \sum_{k=1}^{M} \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \delta_{\ell j} = \sum_{k=1}^{M} \frac{\partial h}{\partial y_{kj}}$$