# Deriving Layer Norm.



$$\mu = \frac{\sum_n x_n}{n}$$

$$\sigma^2 = \left(\frac{\sum_n (x_n - \mu)^2}{n-1}\right)$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$y_i = \hat{x}_i \cdot \lambda + \beta$$

$$\frac{\partial L}{\partial \beta} = \left(\frac{\partial L}{\partial \text{out}}\right)$$

$$\frac{\partial L}{\partial \lambda} = \left(\frac{\partial L}{\partial \text{out}} \cdot \text{raw}\right) \cdot \text{sum}$$

$$\frac{\partial L}{\partial \hat{x}_i} = \frac{\partial L}{\partial \text{out}_i} \cdot \lambda$$

$$\frac{\partial L}{\partial \text{raw}_i} = \lambda \cdot \frac{\partial L}{\partial \text{out}_i}$$

$$\frac{\partial L}{\partial \sigma^2} = \sum_i \frac{\partial L}{\partial \text{raw}_i} \cdot \frac{\partial \text{raw}_i}{\partial \sigma^2} \rightarrow \sum_i \lambda \cdot \frac{\partial L}{\partial \text{out}_i} \cdot \frac{d}{d\sigma^2}\left(\frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}\right) \rightarrow \sum_i \lambda \cdot \frac{\partial L}{\partial \text{out}_i} \cdot -\frac{1}{2}\left(\frac{x_i - \mu}{(\sigma^2 + \epsilon)^{3/2}}\right)$$

$$\frac{\hat{x}_i}{(\sigma^2 + \epsilon)'}$$

$$= \sum_i \beta \cdot \frac{\partial L}{\partial \text{out}_i} \cdot -\frac{1}{2} \frac{\text{raw}_i}{\sigma^2 + \epsilon}$$

$$\frac{\partial L}{\partial \mu} = \frac{\partial L}{\partial \sigma^2} \cdot \frac{\partial \sigma^2}{\partial \mu^2} + \sum_i \frac{\partial L}{\partial \text{raw}_i} \cdot \frac{\partial \text{raw}_i}{\partial L} \rightarrow \sum_i \frac{-1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial L}{\partial \text{out}_i} \cdot \lambda = \frac{\partial L}{\partial \mu}$$

$$\frac{\partial}{\partial \sigma^2} \cdot \frac{-2\sum_i x_i - \mu}{n-1} = 0$$

• Note:
$$\sum \frac{x_i - \mu}{\sigma + \epsilon} \rightarrow \frac{1}{\sqrt{\sigma^2 + \epsilon}} \sum_i \frac{x_i}{...}$$
and,
$$\sum x_i - \mu = 0$$
and,
$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial \mu} \cdot \frac{\partial \mu}{\partial x_i} + \frac{\partial L}{\partial \sigma^2} \cdot \frac{\partial \sigma^2}{\partial x_i} + \frac{\partial L}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i}$$

$$\frac{1}{n} \sum_i \frac{-1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial L}{\partial \text{out}_i} \cdot \lambda$$

$$\frac{1}{n} \cdot \frac{\partial L}{\partial \mu} + 2 \cdot \frac{x_i - \mu}{n-1} \cdot \frac{\partial L}{\partial \sigma^2} + \frac{1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial L}{\partial \hat{x}_i}$$

$$+ \frac{x_i - \mu}{n-1 \sqrt{\sigma^2 + \epsilon}} \sum_i \frac{\hat{x}_i}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial L}{\partial \text{out}_i} \cdot -\frac{1}{2} \cdot \lambda$$

$$+ \frac{1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial L}{\partial \text{out}_i} \cdot \lambda$$

$$= \frac{\lambda \cdot n^{-1}}{\sqrt{\sigma^2 + \epsilon}}\left(\frac{\partial L}{\partial \text{out}_i} \cdot n - \sum_i \frac{\partial L}{\partial \text{out}_i} - \frac{\hat{x}_i \cdot n}{n-1} \sum_i \hat{x}_i \cdot \frac{\partial L}{\partial \text{out}_i}\right) = \frac{\partial L}{\partial x}$$