

Exercise 3

1. Derivative of GeLU

We know: $\text{GeLU} = x \cdot P(X \leq x) = x \cdot \phi(x)$ (1)

$$\phi(x) = \frac{1}{2} \left(1 + \text{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right) \xrightarrow{\mu=0, \sigma=1} \frac{1}{2} \left(1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right) \right) \quad (2)$$

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

By equation (1), $\frac{\partial \text{GeLU}}{\partial x} = \phi(x) + x \cdot \phi'(x)$

By equation (2), $\phi'(x) = \text{erf}'\left(\frac{x}{\sqrt{2}}\right)$, let $z = \frac{x}{\sqrt{2}}$

By equation (3), $\text{erf}'(x) = \frac{\partial \text{erf}(z)}{\partial z} \cdot \frac{\partial z}{\partial x} = \frac{2}{\sqrt{\pi}} \cdot e^{-z^2} \cdot \frac{\partial z}{\partial x} = \sqrt{\frac{2}{\pi}} \cdot e^{-\frac{x^2}{2}}$
& chain rule

Altogether, $\frac{\partial \text{GeLU}}{\partial x} = \frac{1}{2} \left(1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right) \right) + x \cdot \sqrt{\frac{2}{\pi}} \cdot e^{-\frac{x^2}{2}}$

2. Derivative of Leaky ReLU.

if $x > 0$, $\frac{\partial \text{L-ReLU}}{\partial x} = 1$

if $x = 0$, define $\frac{\partial \text{L-ReLU}}{\partial x}$ as 1.

if $x < 0$, $\frac{\partial \text{L-ReLU}}{\partial x} = \frac{\partial \alpha x}{\partial x} = \alpha$, where α is the constant defined in L-ReLU.

3. Prove the backprop through einsum.

We know: $A \in \mathbb{R}^{n \times m}$, $B \in \mathbb{R}^{m \times k}$, $C = \text{einsum}(\overset{u}{ij}, \overset{v}{jk} \rightarrow \overset{w}{ik}, A, B)$ (1)

equa. (1) implies: $C_{ik} = \sum_{j=1}^m a_{ij} b_{jk}$ (2)

$$\frac{\partial L}{\partial a_{ij}} = \sum_{i=1}^n \sum_{k=1}^k \frac{\partial L}{\partial C_{ik}} \cdot \frac{\partial C_{ik}}{\partial a_{ij}} = \sum_{i=1}^n \sum_{k=1}^k \frac{\partial L}{\partial C_{ik}} \cdot \frac{\sum_{j=1}^m \frac{\partial a_{ij} b_{jk}}{\partial a_{ij}}}{\frac{\partial a_{ij} b_{jk}}{\partial a_{ij}}} \quad (3)$$

Only when $i=i$, $j=j$, the derivative w.r.t. a_{ij} is non-zero, hence

$$D_{ij} = \frac{\partial L}{\partial a_{ij}} = \sum_{k=1}^k \frac{\partial L}{\partial C_{ik}} \cdot b_{jk} \quad \text{where } D = \text{C-grad}(C)$$

which implies $D = \text{einsum}(\overset{w}{ik}, \overset{v}{jk} \rightarrow \overset{u}{ij}, \text{C-grad}(B))$