

Let's consider the following two layer architecture:

$$h_1 = \text{Relu}(w_1 x + b_1)$$

$$z = w_2 h_1 + b_2$$

$$\mathcal{L} = \text{CE}(z)$$

Where,

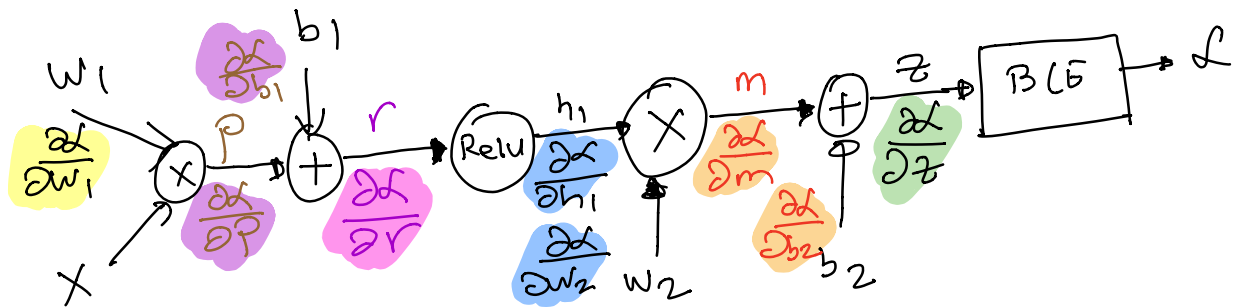
$$x \in \mathbb{R}^D, \quad h_1 \in \mathbb{R}^H, \quad b_1 \in \mathbb{R}^H$$

$$w_1 \in \mathbb{R}^{H \times D}, \quad z \in \mathbb{R}^C, \quad b_2 \in \mathbb{R}^C$$

$$w_2 \in \mathbb{R}^{C \times H}, \quad \mathcal{L} \in \mathbb{R}$$

and CE stands for cross entropy loss.

The first step in backprop is to draw the computational graph



We know, $\frac{\partial \mathcal{L}}{\partial z}$ from HW2, so

Starting with $\frac{\partial \mathcal{L}}{\partial z}$ and backpropagating:

- Since \oplus gate distributes the gradient so,

$$\frac{\partial \mathcal{L}}{\partial m} = \frac{\partial \mathcal{L}}{\partial z}$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z}$$

- Now using the tensor derivative derived in class,

$$\frac{\partial \mathcal{L}}{\partial h_1} = w_2^T \frac{\partial \mathcal{L}}{\partial m}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial m} h_1^T$$

- Since relu gate routes the gradient,
so

$$\frac{\partial \mathcal{L}}{\partial r} = \frac{\partial \mathcal{L}}{\partial h_1} \odot \mathbb{I}(r > 0)$$

- Since \oplus gate distributes the gradient,

$$\frac{\partial \mathcal{L}}{\partial p} = \frac{\partial \mathcal{L}}{\partial r}$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial r}$$

- Now using the tensor derivative, derived in class,

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial p} x^T$$