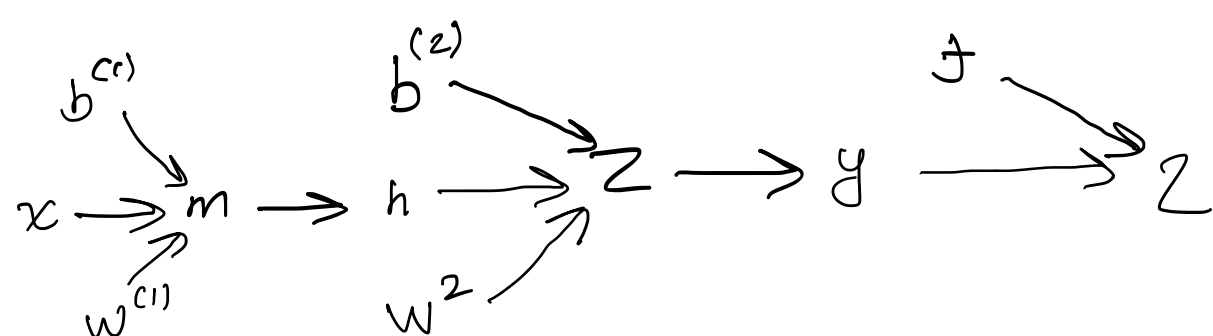


→ Find gradient descent update rule for $w^{(1)}$ and $b^{(1)}$.

Computation graph:



$$m_j = w_{ij}^{(1)} x_j + b_j^{(1)}$$

$$z_j = w_{ij}^{(2)} h_j + b_j^{(2)}$$

→ what we have derived from part (c):

$$\bar{z} = 1$$

$$\bar{y}_k = \frac{-(1+y_k) + z_k}{y_k(1-y_k)} \quad (\text{non-vectorized}).$$

$$\bar{z}_i = \begin{cases} \bar{y}_i \cdot \frac{e^{z_i}(\bar{z}_i) - e^{z_i}e^{z_i}}{(\bar{z}_i)^2} & \text{if } i=j \\ \bar{y}_i \cdot -\frac{e^{z_i}e^{z_i}}{(\bar{z}_i)^2} & \text{if } i \neq j \end{cases}$$

→ vectorized:

$$\bar{z} = 1$$

$$\bar{y} = \frac{-(1+y) + z}{y(1-y)}$$

→ division element wise.

$$\bar{z} = \bar{y} \cdot (\sigma(\text{softmax}))^T$$

→ with σ denoting Jacobian.

→ continuing our computations:

$$\bar{h}_j = \bar{z}_j \cdot \frac{\partial z_j}{\partial h_j} = w_{ij}^{(2)} \Rightarrow$$

$$\bar{h}_j = \bar{z}_j - w_{ij}^{(2)}$$

$$\bar{m}_j = \bar{h}_j \cdot \frac{\partial h_j}{\partial m_j} = \frac{\partial}{\partial m_j} (\text{ReLU}(m_j))$$

$$= \begin{cases} 1 & \text{if } m_j > 0 \\ 0 & \text{if } m_j \leq 0 \end{cases}$$

$$\bar{m}_j = \begin{cases} \bar{h}_j & \text{if } m_j > 0 \\ 0 & \text{if } m_j \leq 0 \end{cases}$$

$$w_{ij}^{(1)} = \bar{m}_j \cdot \frac{\partial m_j}{\partial w_{ij}^{(1)}} = x_j$$

$$b_j^{(1)} = \bar{m}_j \cdot \frac{\partial m_j}{\partial b_j^{(1)}} = \bar{m}_j$$

$$= \bar{m}_j \cdot x_j$$

→ non-vectorized, we have the following update rules:

$$\bar{h}_j = \bar{z}_j \cdot w_{ij}^{(2)}$$

$$\bar{w}_{ij}^{(1)} = \bar{m}_j \cdot x_j$$

$$\bar{m}_j = \begin{cases} \bar{h}_j & \text{if } m_j > 0 \\ 0 & \text{if } m_j \leq 0 \end{cases}$$

$$\bar{b}_j^{(1)} = \bar{m}_j$$

→ vectorized:

$$\bar{h} = w^{(2)} \cdot (\bar{z})^T$$

$$\bar{w}^{(1)} = \bar{m} \cdot x^T$$

$$\bar{b}^{(1)} = \bar{m}$$

$$\bar{m} = \max(\bar{h}, 0)$$