

BACKPROPAGATION

Simple neural network example
 $(\alpha^{(L-3)})_{b^{(L-3)}}^{w^{(L-3)}} (\alpha^{(L-2)})_{b^{(L-2)}}^{w^{(L-2)}} (\alpha^{(L-1)})_{b^{(L-1)}}^{w^{(L-1)}} (\alpha^{(L)})_{b^{(L)}}^{w^{(L)}} (\text{desired output})$

$\text{Cost} = C_o = (\alpha^{(L)} - \gamma)^2 \quad z^{(L)} = w^{(L)} \sigma^{(L-1)} + b^{(L)}$

$\Delta w^{(L)} \rightarrow \Delta z^{(L)} \rightarrow \Delta \alpha^{(L)} \rightarrow \Delta C_o$
 $\frac{\partial z^{(L)}}{\partial w^{(L)}} \times \frac{\partial \alpha^{(L)}}{\partial z^{(L)}} \times \frac{\partial C_o}{\partial \alpha^{(L)}} = \frac{\partial C_o}{\partial w^{(L)}} \quad (\text{chain rule})$

$\frac{\partial z^{(L)}}{\partial w^{(L)}} = \frac{\partial}{\partial w^{(L)}} (w^{(L)} \alpha^{(L-1)} + b^{(L)}) = \alpha^{(L-1)}$

$\frac{\partial \alpha^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$

$\frac{\partial C_o}{\partial \alpha^{(L)}} = \frac{\partial}{\partial \alpha^{(L)}} (\alpha^{(L)} - \gamma)^2 = 2(\alpha^{(L)} - \gamma)$

$\frac{\partial C_o}{\partial b^{(L)}} = 2 \alpha^{(L-1)} \sigma'(z^{(L)}) (\alpha^{(L)} - \gamma)$

$\frac{\partial C_o}{\partial w^{(L)}} = 2 \sigma'(z^{(L)}) (\alpha^{(L)} - \gamma) \quad \frac{\partial C_o}{\partial \alpha^{(L-1)}} = 2 w^{(L)} \sigma'(z^{(L)}) (\alpha^{(L-1)} - \gamma)$

We need to average all $\frac{\partial C_o}{\partial w^{(L)}}$ and $\frac{\partial C_o}{\partial b^{(L)}}$

From each training example to obtain $\frac{\partial C_o}{\partial w^{(L)}}$ and $\frac{\partial C_o}{\partial b^{(L)}}$, the partial derivatives of the full cost function:

$\frac{\partial C_o}{\partial w^{(L)}} = \frac{1}{n} \sum_{n=0}^{N-1} \frac{\partial C_o}{\partial w^{(L)}}$

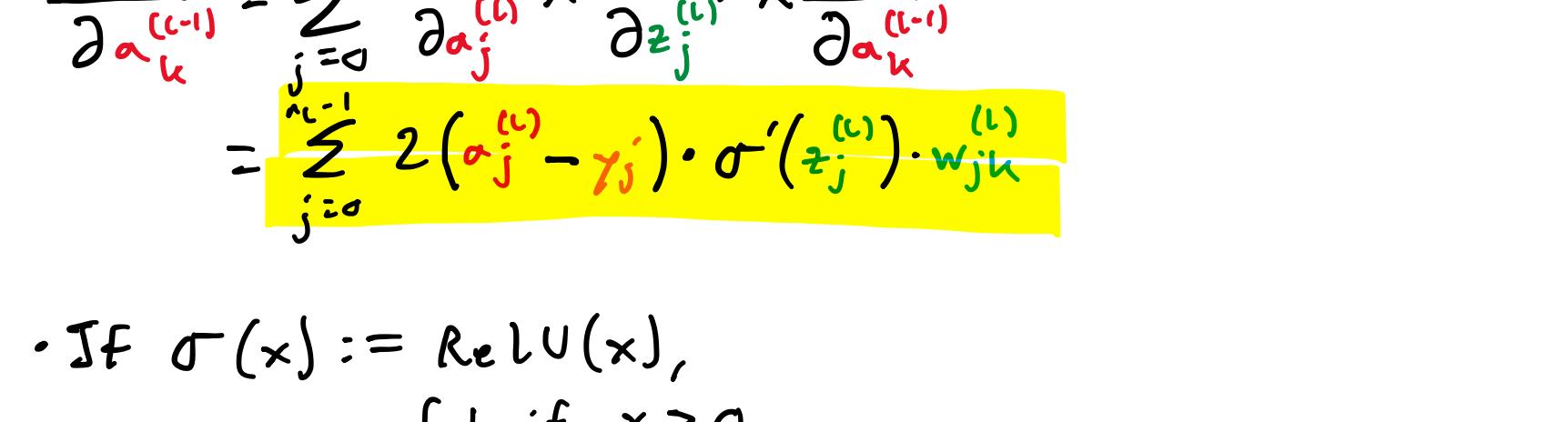
$\frac{\partial C_o}{\partial b^{(L)}} = \frac{1}{n} \sum_{n=0}^{N-1} \frac{\partial C_o}{\partial b^{(L)}}$

$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \\ \vdots \\ \frac{\partial C}{\partial b^{(L)}} \end{bmatrix}$

← gradient

Repeat process recursively.

Multineuron layer example



Index layer L neurons with j and L-1 with k.

$C_o = \sum_{j=0}^{N-1} (\alpha_j^{(L)} - \gamma_j)^2 \quad z_j^{(L)} = \sum_{k=0}^{N-1} w_{jk}^{(L)} \alpha_k^{(L-1)} + b_j^{(L)}$

Chain rule mostly stays the same

$\frac{\partial C_o}{\partial \alpha_k^{(L-1)}} = \sum_{j=0}^{N-1} \frac{\partial C_o}{\partial \alpha_j^{(L)}} \cdot \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial z_j^{(L)}}{\partial \alpha_k^{(L-1)}}$

since a single neuron from the previous layer influences multiple neurons in L.

$\frac{\partial C_o}{\partial w_{jk}^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \times \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}}$

$= 2(\alpha_j^{(L)} - \gamma_j) \cdot \sigma'(z_j^{(L)}) \cdot \alpha_k^{(L-1)}$

$\frac{\partial C_o}{\partial b_j^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \times \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial z_j^{(L)}}{\partial b_j^{(L)}}$

$= 2(\alpha_j^{(L)} - \gamma_j) \cdot \sigma'(z_j^{(L)})$

$\frac{\partial C_o}{\partial \alpha_k^{(L-1)}} = \sum_{j=0}^{N-1} \frac{\partial C_o}{\partial \alpha_j^{(L)}} \times \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial z_j^{(L)}}{\partial \alpha_k^{(L-1)}}$

$= \sum_{j=0}^{N-1} 2(\alpha_j^{(L)} - \gamma_j) \cdot \sigma'(z_j^{(L)}) \cdot w_{jk}^{(L)}$

If $\sigma(x) := \text{ReLU}(x)$,

$\sigma'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$

If $\sigma(x) := \text{sigmoid}(x)$,

$\sigma'(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}}$

$= \sigma(x) \cdot \frac{1-1+e^{-x}}{1+e^{-x}} = \sigma(x) \cdot \left(\frac{1-e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right)$

$= \sigma(x) \cdot (1-\sigma(x)) = \sigma(1-\sigma)$

If $\sigma(\vec{x}) := \text{softmax}(\vec{x})$ where \vec{x} is a vector,

$\sigma_i = \frac{e^{x_i}}{\sum e^{x_n}} \rightarrow \frac{\partial \sigma_i}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\frac{e^{x_i}}{\sum e^{x_n}} \right) = \frac{\sum e^{x_n} - e^{x_i} e^{x_j}}{(\sum e^{x_n})^2}$

$\text{Consider } i \neq j: \quad \frac{\partial \sigma_i}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\frac{e^{x_i}}{\sum e^{x_n}} - e^{x_i} e^{x_j} \right) = -\left(\frac{e^{x_i}}{\sum e^{x_n}} \right) \left(\frac{e^{x_j}}{\sum e^{x_n}} \right) + \boxed{-\sigma_i \sigma_j}$

$\text{Consider } i=j: \quad \frac{\partial \sigma_i}{\partial x_j} = \frac{e^{x_i} e^{x_n} - e^{x_i} e^{x_i}}{(\sum e^{x_n})^2} = \frac{e^{x_i} (\sum e^{x_n} - e^{x_i})}{(\sum e^{x_n})^2}$

$= \left(\frac{e^{x_i}}{\sum e^{x_n}} \right) \left(\frac{\sum e^{x_n} - e^{x_i}}{\sum e^{x_n}} \right) = \left(\frac{e^{x_i}}{\sum e^{x_n}} \right) \left(1 - \frac{e^{x_i}}{\sum e^{x_n}} \right)$

$= \boxed{\sigma_i (1 - \sigma_i)}$

In general,

$\frac{\partial \sigma_i}{\partial x_j} = \sigma_i (\delta_{ij} - \sigma_j), \text{ where } \delta_{ij} := \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i=j \end{cases}$

The Jacobian of softmax(\vec{x}) is

$J_s = \begin{bmatrix} \frac{\partial \sigma_1}{\partial x_1} & \dots & \frac{\partial \sigma_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \sigma_n}{\partial x_1} & \dots & \frac{\partial \sigma_n}{\partial x_n} \end{bmatrix} \text{ where } n \text{ is the number of components in } \vec{x}.$

$\rightarrow \begin{bmatrix} \sigma_1(1-\sigma_1) & \sigma_1 \sigma_2 & \dots \\ \sigma_1 \sigma_1 & \sigma_2(1-\sigma_2) & \dots \\ \vdots & \vdots & \vdots \\ \sigma_n(1-\sigma_n) & \dots & \sigma_n \sigma_1 \end{bmatrix}$

Let $a(w) = \sigma(z(w))$ represent the activation of the output layer wrt its weights.

$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w} = J_s \cdot \frac{\partial z}{\partial w}$

$\text{Recall: } z_j^{(L)} = \sum_k w_{jk}^{(L)} \alpha_k^{(L-1)} + b_j^{(L)} \Rightarrow \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} = \alpha_k^{(L-1)}$

$\rightarrow \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} = \alpha_k^{(L-1)} \quad \text{only nonzero for } i=k$

$\frac{\partial C_o}{\partial w_{ij}^{(L)}} = \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial w_{ij}^{(L)}}$

$= \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial z_i^{(L)}} \cdot \frac{\partial z_i^{(L)}}{\partial w_{ij}^{(L)}} = \sigma_t(\delta_{ti} - \sigma_i) \cdot \alpha_j^{(L-1)}$

$= -\frac{1}{\alpha_i} \cdot \delta_{ti} (\sigma_i - \sigma_t) \cdot \alpha_j^{(L-1)}$

$= \boxed{(\sigma_i - \delta_{ti}) \sigma_j^{(L-1)}}$

$= (\alpha_i^{(L)} - \delta_{ti}) \alpha_j^{(L-1)}$

$\frac{\partial C_o}{\partial b_i^{(L)}} = \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial b_i^{(L)}}$

$= \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial z_i^{(L)}} \cdot \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} = \sigma_t(\delta_{ti} - \sigma_i)$

$= -\frac{1}{\alpha_i} \cdot \delta_{ti} (\sigma_i - \sigma_t)$

$= \boxed{\alpha_i^{(L)} - \delta_{ti}}$

$= \alpha_i^{(L)} - \gamma_i = \boxed{\sigma_i - \gamma_i}$

$\frac{\partial C_o}{\partial w_{ik}^{(L)}} = \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial z_i^{(L)}} \cdot \frac{\partial z_i^{(L)}}{\partial w_{ik}^{(L)}} = \sigma_t(\delta_{ti} - \sigma_i) \cdot \alpha_k^{(L-1)}$

$= -\frac{1}{\alpha_i} \cdot \delta_{ti} (\sigma_i - \sigma_t) \cdot \alpha_k^{(L-1)}$

$= \boxed{\alpha_i^{(L)} - \delta_{ti} \alpha_k^{(L-1)}}$

$\frac{\partial C_o}{\partial b_k^{(L)}} = \frac{\partial C_o}{\partial \alpha_k^{(L)}} \cdot \frac{\partial \alpha_k^{(L)}}{\partial z_k^{(L)}} \cdot \frac{\partial z_k^{(L)}}{\partial b_k^{(L)}} = \sigma_t(\delta_{tk} - \sigma_t)$

$= -\frac{1}{\alpha_k} \cdot \delta_{tk} (\sigma_t - \sigma_t)$

$= \boxed{0}$

$\frac{\partial C_o}{\partial w_{jk}^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \cdot \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} = \sigma_t(\delta_{tj} - \sigma_j) \cdot \alpha_k^{(L-1)}$

$= -\frac{1}{\alpha_j} \cdot \delta_{tj} (\sigma_j - \sigma_t) \cdot \alpha_k^{(L-1)}$

$= \boxed{(\sigma_j - \delta_{tj}) \alpha_k^{(L-1)}}$

$= (\alpha_j^{(L)} - \delta_{tj}) \alpha_k^{(L-1)}$

$\frac{\partial C_o}{\partial b_j^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \cdot \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial z_j^{(L)}}{\partial b_j^{(L)}} = \sigma_t(\delta_{tj} - \sigma_j)$

$= -\frac{1}{\alpha_j} \cdot \delta_{tj} (\sigma_j - \sigma_t)$

$= \boxed{\alpha_j^{(L)} - \delta_{tj}}$

$\frac{\partial C_o}{\partial w_{ik}^{(L)}} = \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial z_i^{(L)}} \cdot \frac{\partial z_i^{(L)}}{\partial w_{ik}^{(L)}} = \sigma_t(\delta_{ti} - \sigma_i) \cdot \alpha_k^{(L-1)}$

$= -\frac{1}{\alpha_i} \cdot \delta_{ti} (\sigma_i - \sigma_t) \cdot \alpha_k^{(L-1)}$

$= \boxed{(\alpha_i^{(L)} - \delta_{ti}) \alpha_k^{(L-1)}}$

$\frac{\partial C_o}{\partial b_k^{(L)}} = \frac{\partial C_o}{\partial \alpha_k^{(L)}} \cdot \frac{\partial \alpha_k^{(L)}}{\partial z_k^{(L)}} \cdot \frac{\partial z_k^{(L)}}{\partial b_k^{(L)}} = \sigma_t(\delta_{tk} - \sigma_t)$

$= -\frac{1}{\alpha_k} \cdot \delta_{tk} (\sigma_t - \sigma_t)$

$= \boxed{0}$

$\frac{\partial C_o}{\partial w_{jk}^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \cdot \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} = \sigma_t(\delta_{tj} - \sigma_j) \cdot \alpha_k^{(L-1)}$

$= -\frac{1}{\alpha_j} \cdot \delta_{tj} (\sigma_j - \sigma_t) \cdot \alpha_k^{(L-1)}$

$= \boxed{(\sigma_j - \delta_{tj}) \alpha_k^{(L-1)}}$

$= (\alpha_j^{(L)} - \delta_{tj}) \alpha_k^{(L-1)}$

$\frac{\partial C_o}{\partial b_j^{(L)}} = \frac{\partial C_o}{\partial \alpha_j^{(L)}} \cdot \frac{\partial \alpha_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial z_j^{(L)}}{\partial b_j^{(L)}} = \sigma_t(\delta_{tj} - \sigma_j)$

$= -\frac{1}{\alpha_j} \cdot \delta_{tj} (\sigma_j - \sigma_t)$

$= \boxed{\alpha_j^{(L)} - \delta_{tj}}$

$\frac{\partial C_o}{\partial w_{ik}^{(L)}} = \frac{\partial C_o}{\partial \alpha_i^{(L)}} \cdot \frac{\partial \alpha_i^{(L)}}{\partial z_i^{(L)}} \cdot \frac{\partial z_i^{(L)}}{\partial w_{ik}^{(L)}} = \sigma_t(\delta_{ti} - \sigma_i) \cdot \alpha_k^{(L-1)}$

$= -\frac{1}{\alpha_i} \cdot \delta_{ti} (\sigma_i - \sigma_t) \cdot \alpha_k^{(L-1)}$

$= \boxed{(\alpha_i^{(L)} - \delta_{ti}) \alpha_k^{(L-1)}}$

$\frac{\partial C_o}{\partial b_k^{(L)}} = \frac{\partial C_o}{\partial \alpha_k^{(L)}} \cdot \frac{\partial \alpha_k^{(L)}}{\partial z_k^{(L)}} \cdot \frac{\partial z_k^{(L)}}{\partial b_k^{(L)}} = \sigma_t(\delta_{tk} - \sigma_t)</$