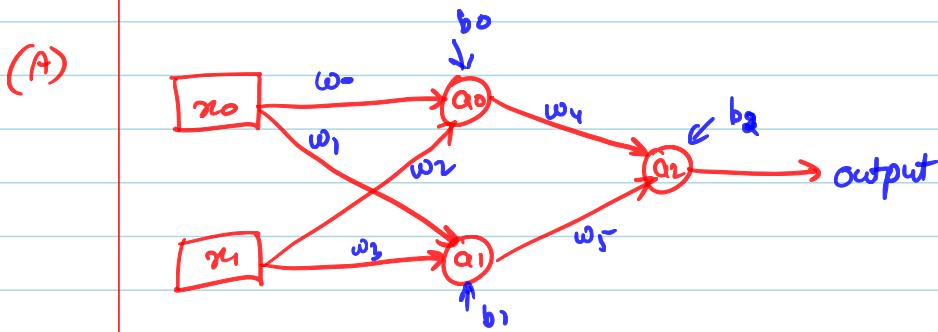


Neural network derivation for backpropagation.



input: (x_0, x_1)

output: q_2

$$\text{Error: } L = \frac{1}{\alpha} (y - q_2)^2 \quad \text{--- (1)}$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad \text{--- (2)}$$

Equations:

$$z_0 = w_0 x_0 + w_1 x_1 + b_0 \quad \text{--- (3)}$$

$$a_0 = \sigma(z_0) \quad \text{--- (4)}$$

$$z_1 = w_2 x_0 + w_3 x_1 + b_1 \quad \text{--- (5)}$$

$$a_1 = \sigma(z_1) \quad \text{--- (6)}$$

$$z_2 = w_4 a_0 + w_5 a_1 + b_2 \quad \text{--- (7)}$$

$$a_2 = \sigma(z_2) \quad \text{--- (8)}$$

For gradient descent, we want to find

$$\frac{\partial L}{\partial b_i}, \frac{\partial L}{\partial w_j} \quad \forall i, j$$

① Find $\sigma'(x)$ in terms of $\sigma(x)$.

Ans

$$\begin{aligned} \sigma(x) &= \frac{1}{1+e^{-x}} \Rightarrow \frac{d\sigma(x)}{dx} = \frac{d}{dx} (1+e^{-x})^{-1} = -\frac{1}{(1+e^{-x})^2} [0 - e^{-x}] \\ &= \left(\frac{1}{1+e^{-x}}\right) \left(\frac{e^{-x}}{1+e^{-x}}\right) = \sigma(x) \left[\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}}\right] \\ &= \sigma(x) [1 - \sigma(x)] \end{aligned}$$

② Find all derivatives related to eq ①, ⑦, ④.

Ans ① $L = \frac{1}{\alpha} (q_2 - y)^2$ where $y \rightarrow$ fixed as per given input, not a variable in the model.

$$\Rightarrow \frac{dL}{dq_2} = q_2 - y.$$

② $z_0 = w_0 x_0 + w_1 x_1 + b_0$

$(x_0, x_1) \rightarrow$ input, not a parameter of NN.

$$\Rightarrow \frac{\partial z_0}{\partial w_0} = x_0, \quad \frac{\partial z_0}{\partial w_1} = x_1, \quad \frac{\partial z_0}{\partial b_0} = 1$$

$$④ a_0 = \sigma(z_0)$$

$$\Rightarrow \frac{da_0}{dz_0}, \frac{d\sigma(z_0)}{dz_0} = \sigma(z_0) (1 - \sigma(z_0))$$

③ Find all partial derivatives: $\frac{\partial L}{\partial b_i}$, $\frac{\partial L}{\partial w_j}$

$$\text{Ans} \quad L = \frac{1}{2} (a_2 - y)^2 \Rightarrow \frac{\partial L}{\partial z_2} = a_2 - y$$

$$\rightarrow a_2 = \sigma(z_2) \Rightarrow \frac{da_2}{dz_2} = \sigma'(z_2) \Rightarrow \frac{\partial L}{\partial z_2} = \sigma'(z_2)(a_2 - y)$$

$$\rightarrow z_2 = w_4 a_0 + w_5 q_1 + b_2 \Rightarrow \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial b_2} = \sigma'(z_2)(a_2 - y) \times 1$$

$$\text{likewise, } \frac{\partial L}{\partial w_4} = \frac{\partial L}{\partial z_2} \times a_0, \quad \frac{\partial L}{\partial w_5} = \frac{\partial L}{\partial z_2} a_1$$

$$\rightarrow a_1 = \sigma(z_1) \Rightarrow \frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial z_1}$$

$$= \frac{\partial L}{\partial z_2} \frac{\partial}{\partial z_1} (w_4 a_0 + w_5 q_1 + b_2)^\circ$$

$$= \frac{\partial L}{\partial z_2} w_4 \frac{\partial a_1}{\partial z_1} = w_4 \frac{\partial L}{\partial z_2} \sigma'(z_1)$$

$$\rightarrow a_0 = \sigma(z_0) \Rightarrow \frac{\partial L}{\partial z_0} = w_4 \frac{\partial L}{\partial z_2} \sigma'(z_0)$$

$$\rightarrow z_1 = w_1 a_0 + w_3 q_1 + b_1 \Rightarrow \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial b_1} = (a_2 - y) \sigma'(z_2) w_5 \sigma'(z_1)$$

$$\Rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_1} = (a_2 - y) \sigma'(z_2) w_5 \sigma'(z_1) a_0$$

$$\Rightarrow \frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_3} = (a_2 - y) \sigma'(z_2) w_5 \sigma'(z_1) a_1$$

$$\rightarrow z_0 = w_0 a_0 + w_2 q_1 + b_0$$

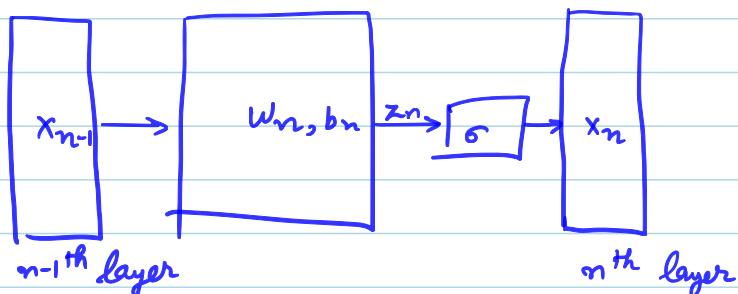
$$\Rightarrow \frac{\partial L}{\partial b_0} = \frac{\partial L}{\partial z_0} \frac{\partial z_0}{\partial b_0} = (a_2 - y) \sigma'(z_2) w_4 \sigma'(z_0)$$

$$\Rightarrow \frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial z_0} \frac{\partial z_0}{\partial w_0} = (a_2 - y) \sigma'(z_2) w_4 \sigma'(z_0) a_0$$

$$\Rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_0} \frac{\partial z_0}{\partial w_2} = (a_2 - y) \sigma'(z_2) w_4 \sigma'(z_0) a_1$$

(B) Matrix form:-

let the layers are connected as:



Equations:

$$x_n = \sigma(z_n) \quad \text{where } x_n: \text{column vector of size } m_n$$

-①

$z_n: \text{column vector of size } m_n$

σ is applied component wise on z_n .

$$z_n = W_n x_{n-1} + b_n \quad \text{where } b_n: \text{column vector of size } m_n$$

-②

$W_n: \text{matrix of size } m_n \times m_{n-1}$

or $x_n, z_n, b_n \in \mathbb{R}^{m_n \times 1}$
 $W_n \in \mathbb{R}^{m_n \times m_{n-1}}$

① $\frac{dx_n}{dz_n}$: we first need to define the meaning of vector differentiation.

let $y = f(x)$

$$y \in \mathbb{R}^{n \times 1}$$

$$x \in \mathbb{R}^{m \times 1}$$

Then $\frac{dy}{dx} := \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \vdots & & \vdots \\ \frac{\partial y_1}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}_{m \times n}$

Point to remember: if y is a scalar, $\frac{dy}{dx}$ is $[1]$ instead of $[-]$

Coming to $\frac{dx_n}{dz_n}$, $[a_{ij}] = \frac{dx_n}{dz_n}$

$$\Rightarrow a_{ij} = \frac{d x_n^{(i)}}{d z_n^{(j)}} = \frac{d \sigma(z_n^{(i)})}{d z_n^{(j)}} = \begin{cases} 0 & i \neq j \\ \sigma'(z_n^{(i)}) & i=j \end{cases}$$

$$\Rightarrow \frac{dx_n}{dz_n} = \text{diag} [\sigma'(z_n^{(1)}), \sigma'(z_n^{(2)}), \dots, \sigma'(z_n^{(m)})]$$

$m_n \times m_m$

$\in \mathbb{R}$ ————— ?

(2) $\frac{\partial z_n}{\partial b_n} :$ $Z_n = W_n x_{n-1} + b_n$
 $\Rightarrow \frac{\partial z_n}{\partial b_n}$ should be a matrix of size $m_n \times m_m$.

$$[a_{ij}] = \frac{\partial z_n}{\partial b_n} \Rightarrow a_{ij} = \frac{\partial z_n^{(j)}}{\partial b_n^{(i)}}$$

$$\text{Now, } Z_n^{(j)} = \sum_{k=1}^{m_{n-1}} W_n^{(j,k)} X_{n-1}^{(k)} + b_n^{(j)}$$

$$\Rightarrow a_{ij} = \frac{\partial z_n^{(j)}}{\partial b_n^{(i)}} = \frac{\partial}{\partial b_n^{(i)}} \left[\sum_{k=1}^{m_{n-1}} W_n^{(j,k)} X_{n-1}^{(k)} + b_n^{(j)} \right]$$

$$= \frac{\partial b_n^{(j)}}{\partial b_n^{(i)}} = \delta_{ij}$$

$$\Rightarrow a_{ij} = \delta_{ij} \text{ or } \frac{\partial z_n}{\partial b_n} = I_{m_n \times m_m}$$

(3) $\frac{\partial z_n}{\partial x_{n-1}} :$ $Z_n = W_n x_n + b_n$
 \Rightarrow Resulting matrix should have a size of $m_n \times m_{n-1}$
 (from definition)

$$[a_{ij}] = \frac{\partial z_n}{\partial x_{n-1}} \Rightarrow a_{ij} = \frac{\partial z_n^{(j)}}{\partial x_{n-1}^{(i)}}$$

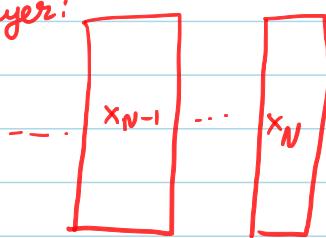
$$= \frac{\partial}{\partial x_{n-1}^{(i)}} \left[\sum_{k=1}^{m_{n-1}} W_n^{(j,k)} X_{n-1}^{(k)} + b_n^{(j)} \right]$$

$$= W_n^{(j,i)}$$

$$\Rightarrow a_{ij} = W_n^{j,i} \text{ or } \frac{\partial z_n}{\partial x_{n-1}} = W_n^T$$

(4) $\frac{\partial z_n}{\partial W_n} :$ This will be 3rd order tensor & we don't know/want to go into that. We will go by another approach later.

(c) last layer:



$$\text{loss } L = \frac{1}{2} (x_N - y)^T (x_N - y)$$

Now, we will compute all partial derivatives

① $\frac{\partial L}{\partial x_N}$ = column vector of size m_N . from definition.

$$\begin{aligned} [a_{i^*}] = \frac{\partial L}{\partial x_N} &\Rightarrow a_{i^*} = \frac{\partial L}{\partial x_N^{(i^*)}} = \frac{\partial}{\partial x_N^{(i^*)}} \frac{1}{2} \sum_{k=1}^{m_N} (x_N^{(k)} - y^{(k)})^2 \\ &= x_N^{(i^*)} - y^{(i^*)} \\ \Rightarrow \frac{\partial L}{\partial x_N} &= x_N - y \quad \text{--- } \textcircled{1} \end{aligned}$$

② We have $x_N = \sigma(z_N)$ & know what is $\sigma'(z_N)$

$$\Rightarrow \frac{\partial L}{\partial z_N} = \underbrace{\frac{\partial x_N}{\partial z_N}}_{m_N \times m_N} \underbrace{\frac{\partial L}{\partial x_N}}_{m_N \times 1} = \underbrace{\sigma'(z_N)}_{m_N \times 1} (x_N - y) \quad \text{--- } \textcircled{2}$$

Remember that chain rule is reversed here!!

③ We have $z_N = w_N x_{N-1} + b_N$.

$$\Rightarrow \frac{\partial L}{\partial b_N} = \frac{\partial z_N}{\partial b_N} \frac{\partial L}{\partial z_N} = I_{m_N \times m_N} \cdot \sigma'(z_N) \frac{\partial L}{\partial x_N}$$

$$\Rightarrow \boxed{\frac{\partial L}{\partial b_N} = \sigma'(z_N) \frac{\partial L}{\partial x_N} \quad \text{or} \quad \frac{\partial L}{\partial b_e} = \sigma'(z_e) \frac{\partial L}{\partial x_e}}$$

$$= \sigma'(z_N) (x_N - y)$$

$$\boxed{\textcircled{4} \quad \frac{\partial L}{\partial x_{N-1}} = \frac{\partial z_N}{\partial x_{N-1}} \frac{\partial L}{\partial z_N} = w_N^T \sigma'(z_N) \frac{\partial L}{\partial x_N}}$$

→ Going back one more layer to achieve generalization,

$$x_{N-1} = \sigma(z_{N-1})$$

$$z_{N-1} = w_{N-1}^T x_{N-2} + b_{N-1}$$

$$\frac{\partial L}{\partial z_{N-1}} = \frac{\partial x_{N-1}}{\partial z_{N-1}} \frac{\partial L}{\partial x_{N-1}} = \sigma'(z_{N-1}) w_N^T \frac{\partial L}{\partial x_N}$$

$$\frac{\partial L}{\partial x_{N-2}} = \frac{\partial z_{N-1}}{\partial x_{N-2}} \frac{\partial L}{\partial z_{N-1}} = w_{N-1}^T \underbrace{\sigma'(z_{N-1}) w_N^T \sigma'(z_N)}_{\text{pattern.}} \frac{\partial L}{\partial x_N}$$

$$\frac{\partial L}{\partial b_{N-1}} = \frac{\partial z_{N-1}}{\partial b_{N-1}} \cdot \frac{\partial L}{\partial z_{N-1}} = \frac{\partial L}{\partial z_{N-1}}$$

- what about $\frac{\partial L}{\partial w_N}$?

let $[a_{ij}] = \frac{\partial L}{\partial w_N}$. How it will look like?

$$a_{ij} = \frac{\partial L}{\partial w_N^{(i,j)}} . \text{ Effect of } w_N^{(i,j)} \text{ will be reflected to } x_N^{(i)} \text{ only}$$

$$\Rightarrow a_{ij} = \frac{\partial L}{\partial x_N^{(i)}} \frac{\partial x_N^{(i)}}{\partial w_N^{(i,j)}} \quad (\text{why order changed? bcz now, this is a real number})$$

$$= (x_N^{(i)} - y^{(i)}) \cdot \frac{\partial x_N^{(i)}}{\partial z_N^{(i)}} \frac{\partial z_N^{(i)}}{\partial w_N^{(i,j)}}$$

$$= (x_N^{(i)} - y^{(i)}) \sigma'(z_N^{(i)}) \cdot \frac{\partial}{\partial w_N^{(i,j)}} \left(\sum_{k=1}^{m_{N-1}} w_N^{(i,k)} x_{N-1}^{(k)} + b_N^{(i)} \right)$$

$$= \underline{\hspace{10em}} \cdot (x_N y^{(j)})$$

$$\Rightarrow a_{ij} = \alpha_i \beta_j = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} [-\beta -] \quad []$$

$$\Rightarrow \frac{\partial L}{\partial w_N} = \underbrace{[\sigma'(z_N) (x_N - y)]}_{m_N \times 1} \quad \underbrace{x_{N-1}^T}_{1 \times m_{N-1}} = \frac{\partial L}{\partial z_N} x_{N-1}^T$$

$m_N \times m_{N-1}$ as expected

let's say we know the fact that

$$\textcircled{1} \text{ if } A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \text{ then } \text{vec}(A) := \begin{bmatrix} a \\ c \\ b \\ d \end{bmatrix}$$

\textcircled{2} if $Z = ABC$, then $\text{vec}(Z) = (C^T \otimes A) \text{vec}(B)$
Can we do something?

for $\frac{\partial L}{\partial w_N}$, we can find $\frac{\partial L}{\partial w_N}$ where $w_N = \text{vec}(W_N)$

and since the result will a vector of same size, we can devectorize?

$$z_N = w_N x_{N-1} + b_N \Rightarrow z_N = \text{vec}(w_N x_{N-1}) + b_N \\ = \text{vec}(I w_N x_{N-1}) + b_N \\ = (x_{N-1}^T \otimes I) w_N + b_N$$

$$W_N \in \mathbb{R}^{m_N \times m_{N-1}} \\ x_{N-1} \in \mathbb{R}^{m_{N-1} \times 1} \Rightarrow \text{long } \mathbb{R}^{1 \times m_{N-1}} \\ I \in \mathbb{R}^{m_N \times m_N} \Rightarrow x_{N-1}^T \otimes I \in \mathbb{R}^{m_N \times (m_N \times m_{N-1})}$$

$$\Rightarrow \frac{\partial L}{\partial w_N} \stackrel{?}{=} \frac{\partial z_N}{\partial w_N} \times \frac{\partial x_{N-1}}{\partial z_N} \times \frac{\partial L}{\partial x_{N-1}} \\ \underbrace{(x_{N-1}^T \otimes I)}_{m_N \times m_N}^T \underbrace{\sigma'(z_N)}_{m_N \times 1} \underbrace{x_{N-1}}_{m_{N-1} \times 1} \rightarrow \text{well defined } \checkmark$$

$$= (x_{N-1}^T \otimes I) \sigma'(z_N) (x_{N-1})$$

will $\text{devec}(\text{RHS}) \stackrel{?}{=} \sigma'(z_N) (x_{N-1}) x_{N-1}^T$?

$$\text{Take example: } x_{N-1} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} \quad I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$\sigma'(z_N)(x_{N-1}) = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix}$$

$$\Rightarrow \sigma'(z_N)(x_{N-1}) x_{N-1}^T \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix} [m \times n] = \begin{bmatrix} m_0 y_0 & m_1 y_0 \\ m_0 y_1 & m_1 y_1 \\ m_0 y_2 & m_1 y_2 \end{bmatrix}$$

$$\& (x_{N-1} \otimes I) = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \otimes \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} x_0 & 0 & 0 \\ 0 & x_1 & 0 \\ 0 & 0 & x_1 \end{pmatrix}$$

$$\text{Product} = \begin{pmatrix} m_0 y_0 & m_1 y_0 \\ m_0 y_1 & m_1 y_1 \\ m_0 y_2 & m_1 y_2 \end{pmatrix} \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} = \underbrace{\begin{pmatrix} m_0 y_0 \\ m_1 y_1 \\ m_0 y_1 \\ m_1 y_2 \\ x_0 y_0 \\ x_1 y_1 \\ x_0 y_1 \\ x_1 y_2 \\ x_1 y_1 \\ x_1 y_2 \end{pmatrix}}$$

$$\underbrace{\begin{pmatrix} m_0 y_0 \\ m_1 y_1 \\ m_0 y_1 \\ m_1 y_2 \\ x_0 y_0 \\ x_1 y_1 \\ x_0 y_1 \\ x_1 y_2 \\ x_1 y_1 \\ x_1 y_2 \end{pmatrix}}_{m \times n}$$

→ Need to prove for general case:

$$\text{devec}(X \otimes I)Y = YX^T$$

or

$$(X \otimes I)Y = \text{vec}(YX^T) \rightarrow \text{This is truth!! So everything matches.}$$

→ From above, if $Z = WX + b$

$$\text{Then } \frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} X^T$$

$$\frac{\partial L}{\partial X} = W^T \frac{\partial L}{\partial Z}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial Z}$$

Books referred before these notes:

① Mathematics for deep learning - Ronald T. Kneusel

② Kronecker Products & Matrix Ghouls: with Applications - Alexander Graham.

