

13.9.21

Neural Networks - HW03

① Prove that,

$$D_{h^{(k)}} C = D_{\hat{y}} L_y(\hat{y})$$

From forward propagation, we know that

$$C(\theta) = L_y(\hat{y}) + \lambda \cancel{\phi(\theta)} \quad (\text{since, } \lambda = 0)$$

therefore  $C = L_y(\hat{y})$

We need to prove that

$$D_{h^{(k)}} L_y(\hat{y}) = D_{\hat{y}} L_y(\hat{y})$$

Again by definition, from forward propagation

$$a^{(k)} = b^{(k)} + w^{(k)} h^{(k-1)}$$

$$h^{(k)} \leftarrow g(a^{(k)})$$

$$\hat{y} \leftarrow h^{(k)}$$

then,

$$D_{h^{(k)}} C = D_{h^{(k)}} L_y(\hat{y}) = D_{\hat{y}} L_y(\hat{y}) =$$

where the first equality follows from (1),  
and the second follows from (2).

$$C \circ L_y$$

$$D_{a^{(k)}} C = D_{h^{(k)}} C \odot g'(a^{(k)})^T$$

$$C = L_y(\hat{y}) + \lambda \Omega$$

$$\text{Since } \lambda = 0$$

$$C = L_y(\hat{y}), \quad \hat{y} = h^{(k)}$$

$$g(a^{(k)}) = h^{(k)} = \sigma a^{(k)}$$

By applying chain rule,

$$\begin{aligned} D_{a^{(k)}} (L_y(\hat{y})) &= \frac{d L_y(\hat{y})}{d a^{(k)}} = \frac{d L_y(h^{(k)})}{d a^{(k)}} \\ &= D_{\hat{y}} L_y(\hat{y}) \cdot D_{\sigma a^{(k)}} \cdot D_{a^{(k)}} \cdot \overset{1}{a^{(k)}} \\ &= D_{\hat{y}} L_y(\hat{y}) \cdot D_{\sigma a^{(k)}} = D_{\hat{y}} L_y(\hat{y}) \cdot g'(a^{(k)}) \quad \text{--- ①} \end{aligned}$$

$$g' = D_{\sigma} \quad (\text{given})$$

$$\begin{aligned} \frac{\partial L_y(\hat{y})}{\partial a_i^{(k)}} &= \sum_j \frac{d L_y(\hat{y})}{d h_j^{(k)}} \cdot \frac{\partial h^{(k)}}{\partial a_i^{(k)}} \\ &= \frac{\partial L_y(\hat{y})}{\partial h_j^{(k)}} \sum_i \frac{\partial h^{(k)}}{\partial a_i^{(k)}} \end{aligned}$$

$$\text{We also know that, } \forall i \neq j, \frac{\partial h^{(k)}}{\partial a_i^{(k)}} = 0.$$

$$\frac{\partial L_y(\hat{y})}{\partial a_i^{(k)}} = \frac{\partial L_y}{\partial h_j^{(k)}} \cdot \sum_i \frac{\partial h_i^{(k)}}{\partial a_i^{(k)}}$$

From what i understand,  $\frac{\partial h_i^{(k)}}{\partial a_i^{(k)}}$  is a diagonal element matrix which is nothing but the activation function that's being multiplied to  $D_{h^{(k)}} C$

$$\Rightarrow \frac{\partial L_y}{\partial h_j^{(k)}} \cdot \begin{bmatrix} \frac{\partial h_0^{(k)}}{\partial a_0} & \dots & \frac{\partial h_n^{(k)}}{\partial a_0} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_0^{(k)}}{\partial a_n} & \dots & \frac{\partial h_n^{(k)}}{\partial a_n} \end{bmatrix}$$

Only diagonal elements will be having values.

Since for a symmetric matrix,  $A = A^T$ ,

we can write this as a transpose.

we also know that  $\text{diag} \left\{ \frac{\partial h_0^{(k)}}{\partial a_0}, \dots, \frac{\partial h_n^{(k)}}{\partial a_n} \right\} = g(a^{(k)})$

we use this result now in ①

which gives us;

$$\begin{aligned} \frac{\partial L_y(\hat{y})}{\partial a_i^{(k)}} &= [g'(a^{(k)})]^T \cdot \frac{\partial L_y(\hat{y})}{\partial h_j^{(k)}} \\ &= D_{h^{(k)}} \cdot L_y(\hat{y}) \odot [g'(a^{(k)})]^T \\ &= D_{h^{(k)}} \cdot C \odot [g'(a^{(k)})]^T \end{aligned}$$

Hence proved.

③ Given ;  $D_{w^{(k)}} C = h^{(k-1)} \cdot D_{a^{(k)}} \cdot C$

Also given,  $[D_{w^{(k)}} C]_{ij} = \frac{\partial C}{\partial w_{ji}^{(k)}}$

This means that the matrix  $\frac{\partial C}{\partial w^{(k)}}$  is a symmetric matrix.

We know,  $\hat{y} = h^{(k)}$

$$c = L_y(\hat{y})$$

$$\frac{\partial c}{\partial w^{(k)}} = \frac{\partial L_y(\hat{y})}{\partial w^{(k)}} = \frac{\partial L_y(h^{(k)})}{\partial w^{(k)}} \quad \text{--- (i)}$$

$$\frac{\partial L_y(\hat{y})}{\partial w_{ij}^{(k)}} = \sum_m \frac{\partial L_y(\hat{y})}{\partial a_m^{(k)}} \cdot \frac{\partial a_m^{(k)}}{\partial w_{ij}^{(k)}} \quad \text{--- (1)}$$

$$a^k = b^k + \sum_i w_i^k h_i^{k-1}$$

for  $m^{\text{th}}$  column of vector  $a^k$

$$a_m^k = b_m^k + \sum_i w_{im}^k h_i^{k-1}$$

from (1), we can understand that all terms except  $\frac{\partial L_y(\hat{y})}{\partial w_{ij}^{(k)}} = 0$

$$\frac{\partial L_y(\hat{y})}{\partial w_{ij}^{(k)}} = \sum_m \frac{\partial L_y(\hat{y})}{\partial a_m^{(k)}} \cdot \sum_l h_l$$

~~$\frac{\partial L_y(\hat{y})}{\partial w} = h^{(k-1)} \cdot \frac{\partial L_y(\hat{y})}{\partial a}$~~

$h^{(k)} = \sigma a^{(k)}$  where  $\sigma \rightarrow$  activation function.

$$\frac{\partial L_y(\sigma(a^{(k)}))}{\partial w^{(k)}}$$

We know,  $a^{(k)} = w^k h^{(k-1)} + b^{(k)}$

$$\frac{\partial C}{\partial w^{(k)}} = \frac{\partial L_y}{\partial w^{(k)}} \left[ - (w^{(k)} h^{(k-1)} + b^{(k)}) \right]$$

$$= D_{\hat{y}} L_y(\hat{y}) \cdot D_{-a^{(k)}} \cdot D w^{(k)} [w^{(k)} \cdot h^{(k-1)} + b^{(k)}]$$

$$= D_{h^{(k)}} C \cdot D_{-a^{(k)}} \cdot h^{(k-1)}$$

$$g' = D_{-a^{(k)}} \rightarrow \text{Symmetric matrix, as from previous question.}$$

$$= D_{h^{(k)}} C \cdot \boxed{g'(a^{(k)})} (h^{(k-1)})$$

$$\frac{dC}{dw^{(k)}} = D_{h^{(k)}} C \cdot x \cdot g'(a^{(k)})^T \cdot h^{(k-1)}$$

from previous question again,

$$D_{a^{(k)}} \cdot C = D_{h^{(k)}} C \odot g'(a^{(k)})^T \quad \text{--- (3)}$$

$$\frac{dC}{dw^{(k)}} = D_{a^{(k)}} C \cdot h^{(k-1)} \quad \text{--- Hence proved}$$

for (3),

$$\frac{\partial L_y(\hat{y})}{\partial h^{(k)}} \odot g'(a^{(k)})^T$$

→ The above multiplication has

to result in a scalar quantity. since it is a gradient.



from the results in previous question,  
 $[g'(a^{(k)})]^T$  will yield a square matrix with  
 $\frac{\partial h_i^{(k)}}{\partial a_i} \rightarrow$  as activation.

$$\textcircled{4} \quad D_h^{(k-1)} \cdot C = (D_a^{(k)} \cdot C) \cdot w^{(k)}$$

$$D_{h^{(k-1)}} C = \frac{dC}{dh^{(k-1)}} = \frac{dL_y(\hat{y})}{dh^{(k-1)}}$$

$$\hat{y} = h^{(k)} ; \quad D_{h^{(k-1)}} C = \frac{dL_y(h^{(k)})}{dh^{(k-1)}}$$

$$\text{then } h^{(k)} = \sigma a^{(k)}$$

$$D_{h^{(k-1)}} C = \frac{dL_y}{dh^{(k-1)}} (\sigma(a^{(k)}))$$

$$a^{(k)} = w^{(k)} h^{(k-1)} + b^{(k)}$$

$$D_{h^{(k-1)}} C = \frac{dL_y}{dh^{(k-1)}} \cdot (\sigma[w^{(k)} h^{(k-1)} + b^{(k)}])$$

$$= D_{\hat{y}} L_y(\hat{y}) D\sigma(a^{(k)}) \cdot D_{h^{(k-1)}} a^{(k)}$$

$$= D_{\hat{y}} L_y(\hat{y}) \cdot D\sigma(a^{(k)}) w^{(k)}$$

$$= D_{h^{(k)}} C \cdot D\sigma(a^{(k)}) w^{(k)}$$

from the previous results,

$$D\sigma = g' \text{ applied element-wise.}$$

therefore,

$$= D_{n(k)} C \odot [g'(a^{(k)})]^T \cdot w^{(k)}$$

$$= [D_{a^{(k)}} C] w^{(k)} \equiv$$