

1) Forward.

$$a_1(x) = h(\alpha^{(0)}(x; w^{(0)}, b^{(0)})), \text{ where } \alpha^{(0)}(x; w^{(0)}, b^{(0)}) = (w^{(0)})^T x + b^{(0)}$$

↳ (hidden) second layer

$$a_2(x) = o(w^{(1)T} a_1(x) + b^{(1)})$$

↳ output softmax activation function.

Like this, 일반화 한다면 $L+1$ 개 층을 가진. $\{ \text{Input} = 1, \text{hidden} = L-1, \text{output} = 1 \}$

$$\textcircled{1} a_k(x) = h_k(x) = h^{(k)}(\alpha^{(k)}(x; w^{(k)}, b^{(k)})) \quad (\text{from } k=1 \text{ to } L-1)$$

$$\text{where, } \alpha^{(k)}(x; w^{(k)}, b^{(k)}) = (w^{(k)})^T h^{(k-1)}(x) + b^{(k)}$$

$$\textcircled{2} h_L(x) = o(\alpha^{(L)}(x)) = a_L(x)$$

Loss function.

$$l(a_L(x), y) = - \sum_c \underset{\substack{\uparrow \\ \text{predicted}}}{I_{(y=c)}} \ln \underset{\substack{\uparrow \\ \text{real}}}{a_L(x)_c} = - \ln a_L(x)_y$$

∴ with m sample \Rightarrow negative log-likelihood loss

$$J = -\frac{1}{m} \sum_{i=1}^m -\ln a_L(x_i)_y = -\frac{1}{m} \sum_{i=1}^m \sum_c I_{(y=c)} \ln a_L(x_i)_c$$

partial derivative)

$$\frac{\partial J}{\partial a_L(x)_c} = \frac{\partial (-\ln a_L(x)_y)}{\partial a_L(x)_c} = \frac{-I_{(y=c)}}{a_L(x)_c}$$

$$\therefore \frac{\partial J}{\partial a_L} = - \begin{bmatrix} \frac{I_{(y=0)}}{a_L(x)_0} \\ \vdots \\ \frac{I_{(y=L-1)}}{a_L(x)_{L-1}} \end{bmatrix} = -e(y) \oslash a_L(x)$$

$$e(y) = \begin{bmatrix} I_{(y=0)} \\ \vdots \\ I_{(y=L-1)} \end{bmatrix}$$

\oslash : element wise division.

$$\therefore \text{loss} = l = -\ln a_L(x)y$$

$$a_L(x) = \text{Softmax}(d^{(L)}(x))$$

$$\therefore \frac{\partial l}{\partial d^{(L)}(x)} = \frac{\partial l}{\partial a_L(x)y} \cdot \frac{\partial a_L(x)y}{\partial a_L(x)c} = \frac{-1}{a_L(x)y} \cdot \frac{\partial a_L(x)y}{\partial a_L(x)c}$$

$\sigma(x)(1-\sigma(x))$
for sigmoid

$$= \frac{-1}{a_L(x)y} a_L(x)y (1_{(y=c)} - a_L(x)c) = a_L(x)c - 1_{(y=c)}$$

\therefore to be general

$$\frac{\partial l}{\partial d^{(L)}} = a_L(x) - \begin{bmatrix} 1_{(y=c_1)} \\ \vdots \\ 1_{(y=c_N)} \end{bmatrix} = a_L(x) - e(y)$$

at hidden layer (need W, b)

$$\text{from } d^{(L)}(x) = (W^{(L)})^T a_{L-1}(x) + b^{(L)}$$

$$i) \frac{\partial l}{\partial w_{ij}^{(L)}} = \frac{\partial l}{\partial d^{(L)}(x)_j} \cdot \frac{\partial d^{(L)}(x)_j}{\partial w_{ij}^{(L)}} = \frac{\partial l}{\partial d^{(L)}(x)_j} \cdot a_{L-1}(x)_i$$

$$\Rightarrow \frac{\partial l}{\partial W^{(L)}} = \frac{\partial d^{(L)}}{\partial W^{(L)}} \left(\frac{\partial l}{\partial d^{(L)}} \right)^T = a_{L-1}(x) \left(\frac{\partial l}{\partial d^{(L)}} \right)^T$$

dimension matching.

$$ii) \frac{\partial l}{\partial b_i^{(L)}} = \frac{\partial l}{\partial d^{(L)}(x)_i} \cdot \frac{\partial d^{(L)}(x)_i}{\partial b_i^{(L)}} = \frac{\partial l}{\partial d^{(L)}(x)_i} \cdot 1$$

$$\Rightarrow \text{to general} \quad \frac{\partial l}{\partial b^{(L)}} = \left(\frac{\partial d^{(L)}}{\partial b^{(L)}} \right)^T \frac{\partial l}{\partial d^{(L)}} = 1 \cdot \frac{\partial l}{\partial d^{(L)}} = \frac{\partial l}{\partial d^{(L)}}$$

Then, what is $\frac{\partial \mathcal{L}}{\partial a_k}$?

From $a^{(l+1)}(x) = (w^{(l+1)})^T \cdot a_l(x) + b^{(l+1)}$

$$i) \frac{\partial \mathcal{L}}{\partial a_k(x_j)} = \sum_i \frac{\partial \mathcal{L}}{\partial a^{(l+1)}(x_i)} \cdot \frac{\partial a^{(l+1)}(x_i)}{\partial a_k(x_j)}$$

$$= \sum_i \frac{\partial \mathcal{L}}{\partial a^{(l+1)}(x_i)} w_{ij}^{(l+1)}$$

\therefore to general

$$\frac{\partial \mathcal{L}}{\partial a_k} = \left(\frac{\partial a^{(l+1)}}{\partial a_k} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial a^{(l+1)}} = w^{(l+1)} \cdot \frac{\partial \mathcal{L}}{\partial a^{(l+1)}}$$

dimension matching.

$$\frac{\partial \mathcal{L}}{\partial a^{(l)}(x)_j} \leftarrow \frac{\partial \mathcal{L}}{\partial a_k(x)_j} \cdot \frac{\partial a_k(x)_j}{\partial a^{(l)}(x)_j} = \frac{\partial \mathcal{L}}{\partial a_k(x)_i} \cdot h'(a^{(l)}(x)_j)$$

to general

$$\frac{\partial \mathcal{L}}{\partial a^{(l)}} = \left(\frac{\partial a_k}{\partial a^{(l)}} \right)^T \left(\frac{\partial \mathcal{L}}{\partial a_k} \right) = \begin{bmatrix} h'(a^{(l)}(x)_1) \\ \vdots \\ h'(a^{(l)}(x)_j) \end{bmatrix} \left(\frac{\partial \mathcal{L}}{\partial a_k} \right)$$

$$= h'(a^{(l)}(x)) \odot \left(\frac{\partial \mathcal{L}}{\partial a_k} \right)$$

$$= h'(a^{(l)}(x)) \odot w^{(l)} \frac{\partial \mathcal{L}}{\partial a^{(l)}}$$

Summarise:

$$\left\{ \begin{aligned} \text{Loss: } \mathcal{L} &= -\ln a_k(x)_y \\ a_k &= \text{Softmax}(a^{(k)}(x)) \\ a_{k+1}(x) &= h(a^{(k+1)}(x)) \\ a^{(k)}(x) &= b^{(k)} + (w^{(k)})^T a_{k+1}(x) \end{aligned} \right.$$

$$\frac{dL}{da_L} = - \begin{bmatrix} 1_{(y=y_0)/a_L(x)_0} \\ \vdots \\ 1_{(y=y_0)/a_L(x)_{L-1}} \end{bmatrix}$$

$$\therefore \text{at output level} \Rightarrow d\alpha^{(L)} = \frac{dL}{d\alpha^{(L)}} = a_L(x) - \text{eqy}$$

From layer L to 1

$$\Rightarrow dW^{(k)} = \frac{dL}{dW^{(k)}} = a_{k+1}(x) \left(\frac{dL}{d\alpha^{(k)}} \right)^T = a_{k+1}(x) (d\alpha^{(k)})^T$$

$$db^{(k)} = \frac{dL}{db^{(k)}} = \frac{dL}{d\alpha^{(k)}} = d\alpha^{(k)}$$

$$\text{where } d\alpha^{(k+1)} = \frac{dL}{d\alpha^{(k+1)}} = h'(\alpha^{(k+1)}(x)) \odot \left(\frac{dL}{da_{k+1}} \right)$$

$$= h'(\alpha^{(k+1)}(x)) \odot W^{(k)} d\alpha^{(k)}$$