

Deep Learning101

# Softmax와 Cross Entropy 미분

Ans:  $\hat{y} - y$

안녕하세요 여러분! 신박AI입니다



오늘은 지난 영상에서 잠시 소개드렸던,  
Softmax와 Cross-entropy의 미분값에  
대해 말씀을 드리고자 합니다

# 우리는 지난 RNN영상에서 Cross-entropy를 손실함수로 쓰는 이유를,

Cross-entropy라는 손실함수를 사용하도록 하겠습니다

$$W = \begin{bmatrix} -0.011 & 0.13 \\ -0.123 & 0.014 \end{bmatrix} \quad V = \begin{bmatrix} -0.141 & 0.038 \\ 0.056 & -0.105 \\ -0.132 & 0.14 \end{bmatrix}$$
$$U = \begin{bmatrix} 0.094 & -0.02 & 0.135 \\ 0.135 & -0.069 & -0.009 \end{bmatrix}$$
$$h_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad h_1 = \begin{bmatrix} 0.094 \\ 0.341 \end{bmatrix} \quad h_2 = \begin{bmatrix} 0.15 \\ -0.022 \end{bmatrix}$$
$$o_1 = \begin{bmatrix} -0.008 \\ -0.009 \\ 0.006 \end{bmatrix} \quad \hat{y}_1 = \begin{bmatrix} 0.332 \\ 0.332 \\ 0.337 \end{bmatrix} \quad o_2 = \begin{bmatrix} 0.33 \\ 0.01 \\ -0.022 \end{bmatrix} \quad \hat{y}_2 = \begin{bmatrix} 0.33 \\ 0.341 \\ 0.33 \end{bmatrix}$$
$$\hat{y}_2 = \text{softmax}(o_2)$$
$$o_2 = Vh_2$$
$$h_2 = \tanh(W h_1 + U x_2)$$

baseball

bat

아구  
배트

baseball  
bat

# 역전파 계산에서,

Cross-entropy를 쓰는 이유는 추후에 다룰 역전파의 계산에서

$$W = \begin{bmatrix} -0.011 & 0.13 \\ -0.123 & 0.014 \end{bmatrix} \quad V = \begin{bmatrix} -0.141 & 0.038 \\ 0.056 & -0.105 \\ -0.132 & 0.14 \end{bmatrix}$$

$$U = \begin{bmatrix} 0.094 & -0.02 & 0.135 \\ 0.135 & -0.069 & -0.009 \end{bmatrix}$$

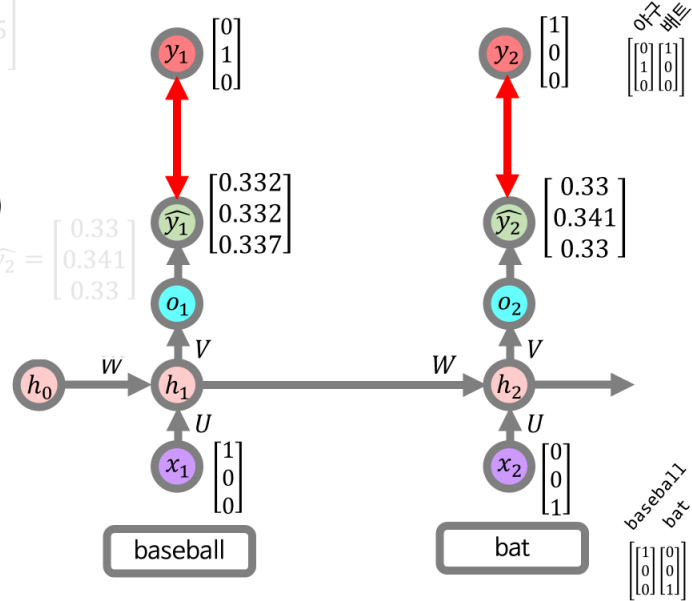
$$h_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad h_1 = \begin{bmatrix} 0.094 \\ 0.135 \end{bmatrix} \quad h_2 = \begin{bmatrix} 0.15 \\ 0.01 \end{bmatrix}$$
$$CE = - \sum_i^C t_i \log((f(s)_i)$$

$$o_1 = \begin{bmatrix} -0.008 \\ -0.009 \\ 0.006 \end{bmatrix} \quad \hat{y}_1 = \begin{bmatrix} 0.332 \\ 0.332 \\ 0.337 \end{bmatrix} \quad o_2 = \begin{bmatrix} 0.33 \\ 0.01 \\ -0.022 \end{bmatrix} \quad \hat{y}_2 = \begin{bmatrix} 0.33 \\ 0.341 \\ 0.33 \end{bmatrix}$$

$$\hat{y}_2 = \text{softmax}(o_2)$$

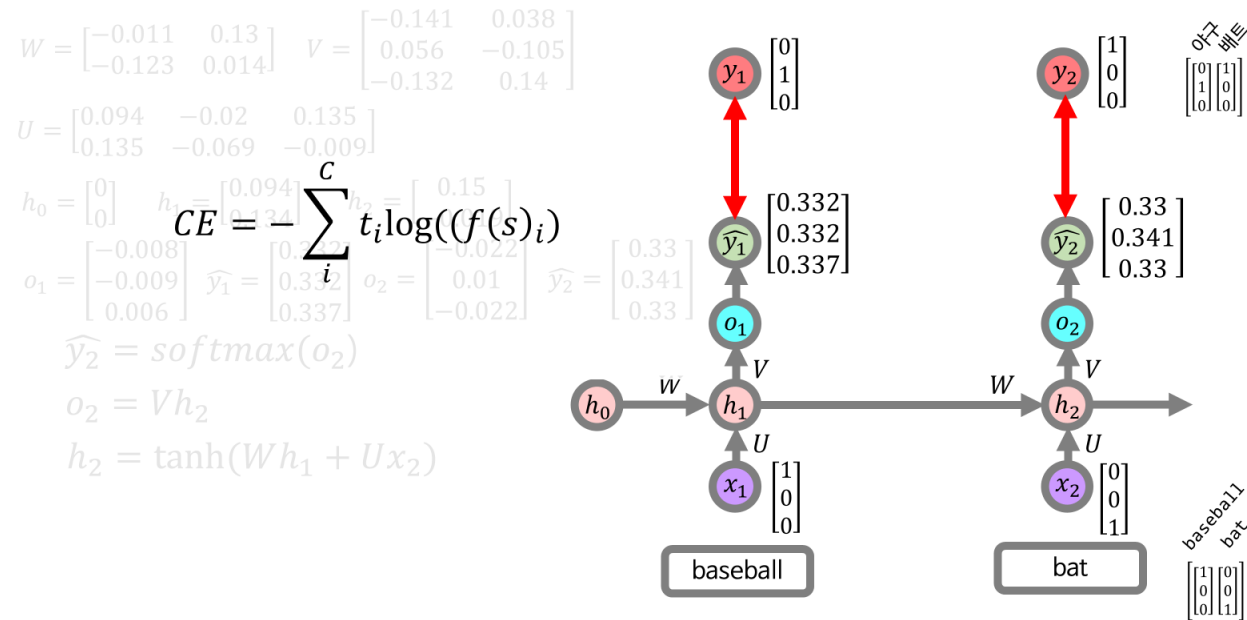
$$o_2 = Vh_2$$

$$h_2 = \tanh(W h_1 + U x_2)$$



Softmax함수와 함께 쓸 때 그 계산이 아주 용이해지기 때문이라고 말씀을 드렸습니다.

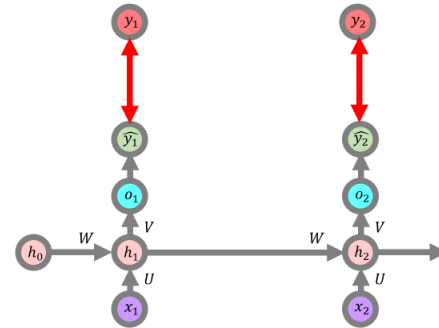
Softmax함수와 함께 쓸 때 역전파 계산이 아주 용이해지기 때문입니다



그리고 그 계산값은  $\hat{y} - y$  이 된다고만 말씀을 드렸고, 그 자세한 도출 과정은 시간 관계상 생략했었습니다.

이 부분은 제가 다른 영상에서 도출과정을 보여드리려고 합니다.  
많이 기대해주세요 ;)

$$\begin{aligned}\frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} \\ \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V} \\ \frac{\partial L}{\partial V} &= (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V}\end{aligned}$$



# LSTM 영상에서도 결과만 보여드리고 과정은 생략하였습니다.

$\hat{y}_t - y$ 는  $\partial L / \partial Z_t$  가 됩니다

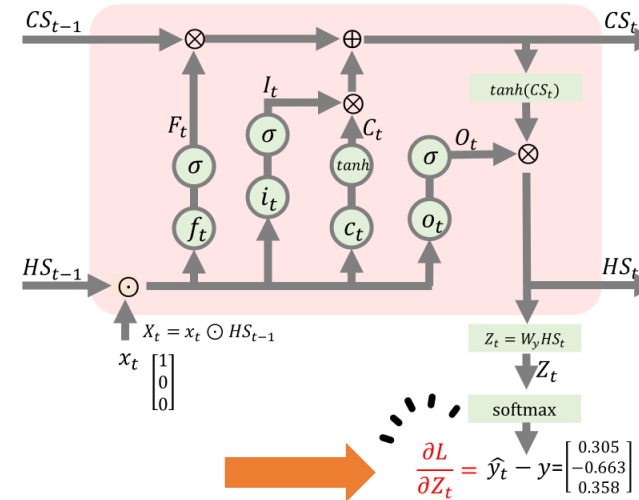
Forget Gate:  $F_t = \sigma(W_f X_t)$   
 $F_t = \sigma(f_t)$   
 $f_t = W_f X_t$

Input Gate:  $I_t = \sigma(W_i X_t)$   
 $I_t = \sigma(i_t)$   
 $i_t = W_i X_t$

Candidate Gate:  $C_t = \tanh(W_c X_t)$   
 $C_t = \tanh(c_t)$   
 $c_t = W_c X_t$

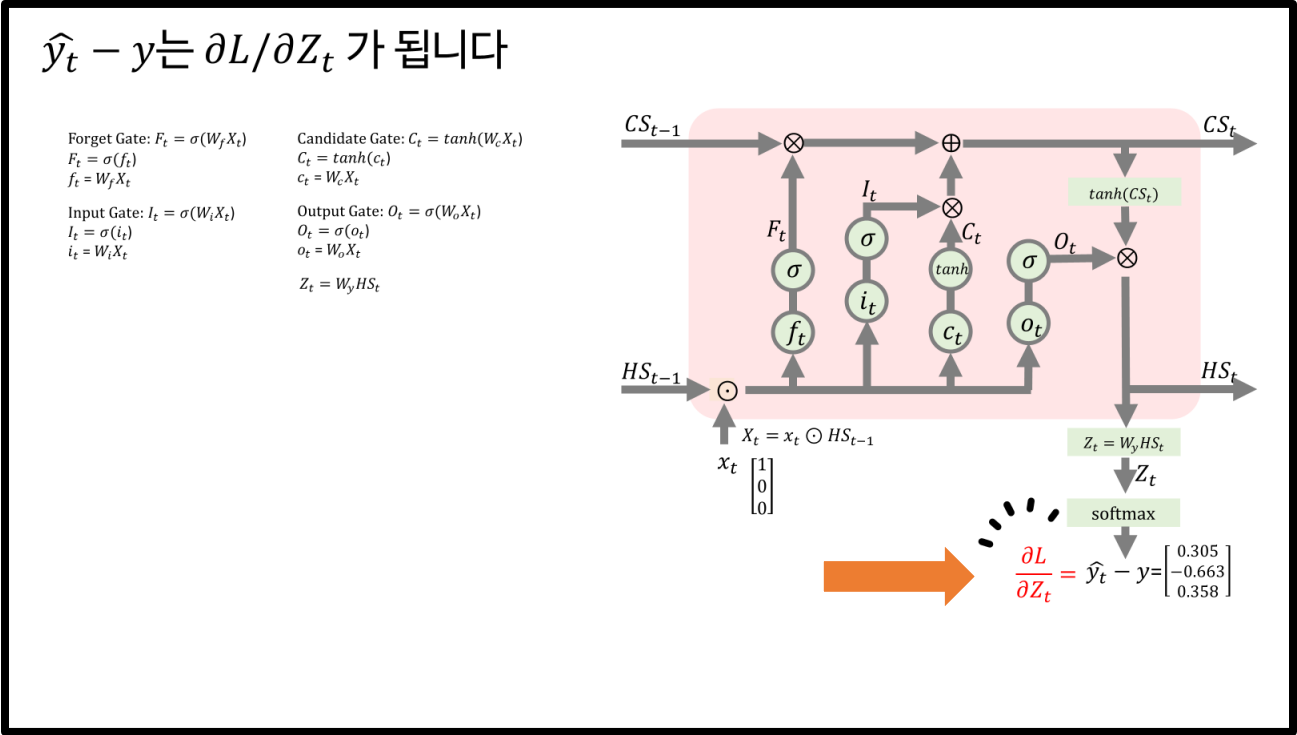
Output Gate:  $O_t = \sigma(W_o X_t)$   
 $O_t = \sigma(o_t)$   
 $o_t = W_o X_t$

$Z_t = W_y H S_t$





# 오늘 이 부분에 대해서 자세히 설명 드리고자 합니다.



Softmax와 Cross-entropy의 조합은  
RNN이나 LSTM뿐만 아니라,

다중 클래스를 분류하는 여러 신경망에서  
대표적으로 쓰이는 조합이기 때문에,

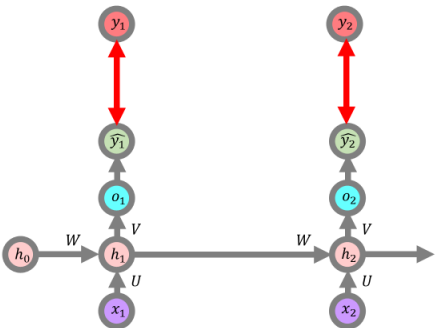
그 도출 과정을 알아두는 것은 여러모로  
유익할 수 있다고 생각이 듭니다.

그러면 바로 시작하도록 하겠습니다.

$$\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$$

이 부분은 제가 다른 영상에서 도출과정을 보여드리려고 합니다.  
많이 기대해주세요 ;)

$$\begin{aligned} \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} \\ \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V} \\ \frac{\partial L}{\partial V} &= (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V} \end{aligned}$$

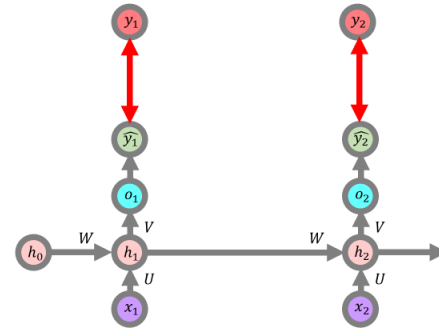


우리는  $\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$  이 편미분 값이

$$\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$$

이 부분은 제가 다른 영상에서 도출과정을 보여드리려고 합니다.  
많이 기대해주세요 ;)

$$\begin{aligned} \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} \\ \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V} \\ \frac{\partial L}{\partial V} &= (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V} \end{aligned}$$



$\hat{y} - y$ , 이 값으로 도출되는 과정을 알아보는 것입니다.

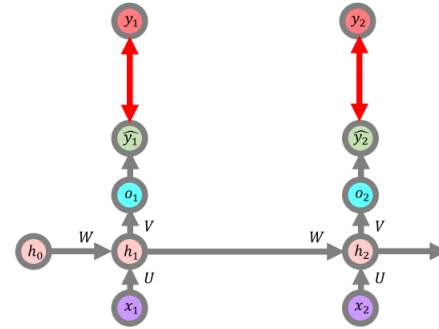
$$\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$$



$$\hat{y} - y$$

이 부분은 제가 다른 영상에서 도출과정을 보여드리려고 합니다.  
많이 기대해주세요 ;)

$$\begin{aligned}\frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} \\ \frac{\partial L}{\partial V} &= \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V} \\ \frac{\partial L}{\partial V} &= (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V}\end{aligned}$$



여기에서  $\partial L_1 / \partial \widehat{y}_1$  은 구분하자면 Cross Entropy의 편미분 값이고,

$$\frac{\partial L_1}{\partial \widehat{y}_1} \frac{\partial \widehat{y}_1}{\partial o_1}$$



이 부분은 Softmax의 편미분 값입니다.

$$\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$$

여기서 우리는 Softmax의 편미분 부터 구해보도록 하겠습니다.

$$\frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1}$$

일반적인 설명을 위해 숫자를 문자  $i, j$ 로 바꾸어 보겠습니다.

$$\frac{\partial L}{\partial \hat{y}} \boxed{\frac{\partial \hat{y}_i}{\partial o_j}}$$

Softmax의 공식은 다음과 같습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}}$$

쉬운 설명을 위해 다음과 같이 세 개의 값들로 이루어진 행렬로 가정하겠습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

Softmax 값인  $\hat{y}_i$  은 다음과 같이 계산할 수 있습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

먼저  $\partial \widehat{y}_1 / \partial o_1$  부터 계산해보도록 하겠습니다.

$$\frac{\partial \widehat{y}_i}{\partial o_j}$$

$$\widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \widehat{y}_1}{\partial o_1} =$$

$\partial \widehat{y}_1 / \partial o_1$  는 다음과 같이 바꾸어 쓸 수가 있습니다.

$$\frac{\partial \widehat{y}_i}{\partial o_j}$$

$$\widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \widehat{y}_1}{\partial o_1} = \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$



$\partial \hat{y}_1 / \partial o_1$  를 구하기 위해서는 다음 미분공식이 필요합니다

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_1}{\partial o_1} = \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

여기서  $e^{o_1}$  을  $g(x)$  로 보고,  $e^{o_1} + e^{o_2} + e^{o_3}$  을  $h(x)$ 로 본다면,

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_1}{\partial o_1} = \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

그러면, 다음과 같이 식을 대입할 수 있고,

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{(e^{o_1})'(e^{o_1}+e^{o_2}+e^{o_3}) - e^{o_1}(e^{o_1}+e^{o_2}+e^{o_3})'}{(e^{o_1}+e^{o_2}+e^{o_3})(e^{o_1}+e^{o_2}+e^{o_3})} \end{aligned}$$

$e^x$ 의 미분값은 자기 자신이기 때문에,

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$f(x) = e^x \rightarrow f'(x) = e^x$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \rightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_1}{\partial o_1} = \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \rightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

$$= \frac{(e^{o_1})'(e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1}(e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})}$$

# 이렇게 바꾸어 쓸 수가 있습니다

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$f(x) = e^x \rightarrow f'(x) = e^x$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \rightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= \frac{e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \end{aligned}$$

$$f(x) = \frac{g(x)}{h(x)} \rightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

또한 이 부분은

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) \\ &= \frac{e^{o_1} (e^{o_1}+e^{o_2}+e^{o_3}) - e^{o_1} (e^{o_1}+e^{o_2}+e^{o_3})'}{(e^{o_1}+e^{o_2}+e^{o_3})(e^{o_1}+e^{o_2}+e^{o_3})} \end{aligned}$$

$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$

$o_1$ 에 관한 편미분만 고려하기 때문에  $e^{o_2}$ 과  $e^{o_3}$ 은 그냥 사라지게 됩니다.

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} &= \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + \cancel{e^{o_2}} + \cancel{e^{o_3}})'}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$

그래서 정리하면,  $\partial \hat{y}_1 / \partial o_1$  는 다음과 같이 바꾸어 쓸 수가 있습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1}(e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1}e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$



$e^{o_1}$  을 이렇게 옮겨서 다음과 같이 항들을 새롭게 어레인지 해 보면

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1}(e^{o_1}+e^{o_2}+e^{o_3}) - e^{o_1}(e^{o_1}+e^{o_2}+e^{o_3})}{(e^{o_1}+e^{o_2}+e^{o_3})^2} \end{aligned}$$

$\partial \widehat{y}_1 / \partial o_1$  는 다음과 같이 바꾸어 쓸 수가 있습니다.

$$\frac{\partial \widehat{y}_i}{\partial o_j}$$

$$\widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \widehat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{(e^{o_1} + e^{o_2} + e^{o_3} - e^{o_1})}{(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$

그리고 이렇게 항들을 분리하면,

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) &= \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{(e^{o_1}+e^{o_2} + e^{o_3} - e^{o_1})}{(e^{o_1} + e^{o_2} + e^{o_3})} \\ &= \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \left( \frac{(e^{o_1} + e^{o_2} + e^{o_3})}{(e^{o_1} + e^{o_2} + e^{o_3})} - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \right) \end{aligned}$$

$\partial \hat{y}_1 / \partial o_1$  는 결국  $\hat{y}_1 (1 - \hat{y}_1)$  로 간단하게 바꿀 수 있습니다

$$\frac{\partial \hat{y}_i}{\partial o_j}$$

$$\hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_1} &= \frac{\partial}{\partial o_1} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) &= \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{(e^{o_1} + e^{o_2} + e^{o_3} - e^{o_1})}{(e^{o_1} + e^{o_2} + e^{o_3})} \\ &= \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \left( \frac{(e^{o_1} + e^{o_2} + e^{o_3})}{(e^{o_1} + e^{o_2} + e^{o_3})} - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \right) = \hat{y}_1 (1 - \hat{y}_1) \end{aligned}$$

이 전개는  $\partial \widehat{y}_2 / \partial o_2$  의 경우도 마찬가지입니다

$$\frac{\partial \widehat{y}_i}{\partial o_j}$$

$$\widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1}+e^{o_2}+e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \widehat{y}_2}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_2}}{e^{o_1}+e^{o_2}+e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_2}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{(e^{o_1}+e^{o_2} + e^{o_3} - e^{o_2})}{(e^{o_1} + e^{o_2} + e^{o_3})} \\ &= \frac{e^{o_2}}{(e^{o_1} + e^{o_2} + e^{o_3})} \left( \frac{(e^{o_1} + e^{o_2} + e^{o_3})}{(e^{o_1} + e^{o_2} + e^{o_3})} - \frac{e^{o_2}}{(e^{o_1} + e^{o_2} + e^{o_3})} \right) = \widehat{y}_2 (1 - \widehat{y}_2) \end{aligned}$$

$\partial \widehat{y}_3 / \partial o_3$  도 마찬가지로  $\widehat{y}_3 (1 - \widehat{y}_3)$  로 바꿀 수 있습니다

$$\frac{\partial \widehat{y}_i}{\partial o_j}$$

$$\widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \widehat{y}_3}{\partial o_3} &= \frac{\partial}{\partial o_3} \left( \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) &= \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{e^{o_3}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{(e^{o_1} + e^{o_2} + e^{o_3} - e^{o_3})}{(e^{o_1} + e^{o_2} + e^{o_3})} \\ &= \frac{e^{o_3}}{(e^{o_1} + e^{o_2} + e^{o_3})} \left( \frac{(e^{o_1} + e^{o_2} + e^{o_3})}{(e^{o_1} + e^{o_2} + e^{o_3})} - \frac{e^{o_3}}{(e^{o_1} + e^{o_2} + e^{o_3})} \right) = \widehat{y}_3 (1 - \widehat{y}_3) \end{aligned}$$



그래서 일반적으로  $\partial \hat{y}_i / \partial o_j = \hat{y}_i (1 - \hat{y}_i)$  가 됩니다 ( $i = j$ 일 경우)

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \frac{\partial}{\partial o_j} \left( \frac{e^{o_i}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

그렇다면  $i \neq j$  의 경우라면 어떨까요?

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \frac{\partial}{\partial o_j} \left( \frac{e^{o_i}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$



$\partial \hat{y}_1 / \partial o_2$  를 한번 구해보도록 하겠습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\frac{\partial \hat{y}_1}{\partial o_2} = \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right)$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

미분공식에 의해서  $\partial \hat{y}_1 / \partial o_2$  는 다음과 같이 전개할 수 있고,

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) &= \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{(e^{o_1})'(e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1}(e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$

$e^{o_1}$  는 보시는 바와 같이  $o_2$  와는 아무런 연관이 없기 때문에,

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) &= \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{(e^{o_1})'(e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1}(e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$

$(e^{o_1})'$  즉  $(= \partial e^{o_1} / \partial o_2)$  은 0가 됩니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) &= \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= \frac{0 \cdot (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1}(e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})(e^{o_1} + e^{o_2} + e^{o_3})} \end{aligned}$$

그러면 이 부분은 다 0이 되고,

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= \frac{0 \cdot (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \end{aligned}$$

$$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$$

또한 이 부분은

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= \frac{0 \cdot (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \end{aligned}$$

$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$

$o_2$  와 관련이 있는 부분은  $e^{o_2}$  밖에 없기 때문에

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= \frac{0 \cdot (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \end{aligned}$$

$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$

$o_2$  와 관련이 없는 부분들은 이렇게 없애버릴 수 있습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= \frac{0 \cdot (e^{o_1} + e^{o_2} + e^{o_3}) - e^{o_1} (e^{o_1} + e^{o_2} + e^{o_3})'}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \end{aligned}$$

$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$



그러면  $\partial \hat{y}_1 / \partial o_2$  를 다음과 같이 새로이 정리할 수 있습니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} &\longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \cdot e^{o_2} \end{aligned}$$

그러면 각각의 항들을 다음처럼  $\hat{y}_1, \hat{y}_2$ 로 바꾸어 생각하면,

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) \\ &= - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \cdot e^{o_2} \end{aligned}$$

$f(x) = \frac{g(x)}{h(x)} \longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2}$

$\partial \widehat{y}_1 / \partial o_2 = -\widehat{y}_1 \cdot \widehat{y}_2$ 가 됩니다

$$\frac{\partial \widehat{y}_i}{\partial o_j} = \begin{cases} \widehat{y}_i (1 - \widehat{y}_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad \widehat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \widehat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} &\longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})^2} \cdot e^{o_2} \\ &= -\widehat{y}_1 \cdot \widehat{y}_2 \end{aligned}$$

그래서 일반적인  $i \neq j$  의 경우,  $\partial \hat{y}_i / \partial o_j = -\hat{y}_i \cdot \hat{y}_j$ 가 됩니다

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ \hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases} \quad \hat{y}_i = \frac{e^{o_i}}{\sum_{k=1}^N e^{o_k}} \longrightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \text{softmax} \left( \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} \right)$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_2}}{e^{o_1} + e^{o_2} + e^{o_3}} \\ \frac{e^{o_3}}{e^{o_1} + e^{o_2} + e^{o_3}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial o_2} &= \frac{\partial}{\partial o_2} \left( \frac{e^{o_1}}{e^{o_1} + e^{o_2} + e^{o_3}} \right) & f(x) = \frac{g(x)}{h(x)} &\longrightarrow f'(x) = \frac{d}{dx} f(x) = \frac{d}{dx} \frac{g(x)}{h(x)} = \frac{g'(x)h(x) - g(x)h'(x)}{\{h(x)\}^2} \\ &= - \frac{e^{o_1}}{(e^{o_1} + e^{o_2} + e^{o_3})} \frac{e^{o_2}}{(e^{o_1} + e^{o_2} + e^{o_3})} \\ &= -\hat{y}_1 \cdot \hat{y}_2 \end{aligned}$$

자 여기까지가 Softmax 기울기를 구하는 과정입니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

자 여기까지가 Softmax 기울기를 구하는 과정입니다.

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

Softmax 기울기는 이렇게 구석에 잠시 두고,

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

이제는 이 결과를 바탕으로 Cross-Entropy + Softmax의 편미분을  
구해보도록 하겠습니다

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$



크로스 엔트로피 손실 함수는 다음과 같습니다

$$\frac{\partial L}{\partial o_j}$$
$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$\partial L / \partial o_j$  값은 다음과 같이 구할 수 있습니다

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$
$$\frac{\partial L}{\partial o_j} = \frac{\partial}{\partial o_j} \left( - \sum_{k=1}^N y_k \log(\hat{y}_k) \right)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

그리고 편미분 부분을 보기 좋게 안쪽으로 넣어보겠습니다

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$
$$\frac{\partial L}{\partial o_j} = \frac{\partial}{\partial o_j} \left( - \sum_{k=1}^N y_k \log(\hat{y}_k) \right)$$
$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

그리고 편미분 부분을 보기 좋게 안쪽으로 넣어보겠습니다

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$
$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$
$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{\partial \log(\hat{y}_k)}{\partial o_j}$$

그리고 체인룰을 가미하면 다음과 같이 바꿀 수 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$
$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{\partial \log(\hat{y}_k)}{\partial o_j}$$
$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

그리고 체인룰을 가미하면 다음과 같이 바꿀 수 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$
$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$
$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{\partial \log(\hat{y}_k)}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j}$$

그러면 이 부분은 log 함수의 미분 공식에 의해서..

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{\partial \log(\hat{y}_k)}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j}$$

$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

이 부분은..

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{\partial \log(\hat{y}_k)}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$



이렇게 바꿀 수 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\widehat{y}_k)$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \boxed{\frac{1}{\widehat{y}_k}} \frac{\partial \widehat{y}_k}{\partial o_j}$$

$$\frac{\partial \widehat{y}_i}{\partial o_j} = \begin{cases} \widehat{y}_i (1 - \widehat{y}_i), & \text{if } i = j \\ -\widehat{y}_i \cdot \widehat{y}_j, & \text{if } i \neq j \end{cases}$$

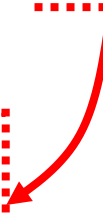
$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

그 다음 이 부분은 이미 우리가 앞에서 구한 바가 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$



그런데  $\partial \hat{y}_k / \partial o_j$  는  $k$ 값의 변화에 따라  $i = j$  가 될 수도 있고,  $i \neq j$  가 될 수도 있기 때문에,

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k$$

$$\frac{1}{\hat{y}_k}$$

$$\frac{\partial \hat{y}_k}{\partial o_j}$$

$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

Σ로만 연결 되어 있기 때문에, 이렇게 나누어 쓸 수가 있습니다

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$= -y_j \frac{1}{\hat{y}_j} \hat{y}_j (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\hat{y}_k} (-\hat{y}_j \cdot \hat{y}_k)) \quad \leftarrow \text{when } k \neq j$$

그러면 각각의 분자 분모를 상쇄하고 나면..

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \quad \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} (-\hat{y}_j \cdot \cancel{\hat{y}_k})) \quad \leftarrow \text{when } k \neq j$$

다음과 같이 정리할 수 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \quad \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} (-\hat{y}_j \cdot \cancel{\hat{y}_k})) \quad \leftarrow \text{when } k \neq j$$

$$= -y_j(1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

계속해서 정리를 다음과 같이 하면,

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} \hat{y}_j \cdot \cancel{\hat{y}_k}) \quad \leftarrow \text{when } k \neq j$$

$$= -y_j (1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

이렇게 정리를 할 수가 있습니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \quad \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} \hat{y}_j \cdot \cancel{\hat{y}_k}) \quad \leftarrow \text{when } k \neq j$$

$$= -y_j (1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k)$$



그런데 이 부분은 우리가 one-hot encoding으로 실제값을 하기로 했었기 때문에,

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \quad \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} \hat{y}_j \cdot \cancel{\hat{y}_k}) \quad \leftarrow \text{when } k \neq j$$

$$= -y_j(1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k)$$

one-hot encoding은 모든 벡터의 합이 언제나 1이 됩니다.

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

to	=	[1,0,0,0,0,0,0,0,0,...,0]
be	=	[0,1,0,0,0,0,0,0,0,...,0]
or	=	[0,0,1,0,0,0,0,0,0,...,0]
not	=	[0,0,0,1,0,0,0,0,0,...,0]
to	=	[1,0,0,0,0,0,0,0,0,...,0]
be	=	[0,1,0,0,0,0,0,0,0,...,0]
that	=	[0,0,0,0,1,0,0,0,0,...,0]
is	=	[0,0,0,0,0,1,0,0,0,...,0]
the	=	[0,0,0,0,0,0,1,0,0,...,0]
question	=	[0,0,0,0,0,0,0,1,0,...,0]

$$= -y_j \frac{1}{\hat{y}_j} \hat{y}_j (1 - \hat{y}_j)$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\hat{y}_k} \hat{y}_j \cdot \hat{y}_k)$$

$$= -y_j(1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k)$$

← when k = j

← when k ≠ j

“1”

그래서 마지막으로 정리하면  $\partial L / \partial o_j$ 은  $\hat{y}_j - y_j$  로 최종 정리할 수가 있습니다!

$$L = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

$$\frac{\partial L}{\partial o_j}$$

$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \quad \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j) \quad \leftarrow \text{when } k = j$$

$$+ (- \sum_{k \neq j}^N y_k \frac{1}{\cancel{\hat{y}_k}} \hat{y}_j \cdot \cancel{\hat{y}_k}) \quad \leftarrow \text{when } k \neq j$$

$$= -y_j (1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j \quad \therefore$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k) = \hat{y}_j - y_j$$

여기까지  $\partial L / \partial o$ 가  $\hat{y}_j - y_j$ 가 되는 과정에 대해 전개해 보았습니다

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V}$$

$$\frac{\partial L}{\partial V} = (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V}$$

$$\frac{\partial \hat{y}_i}{\partial y_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

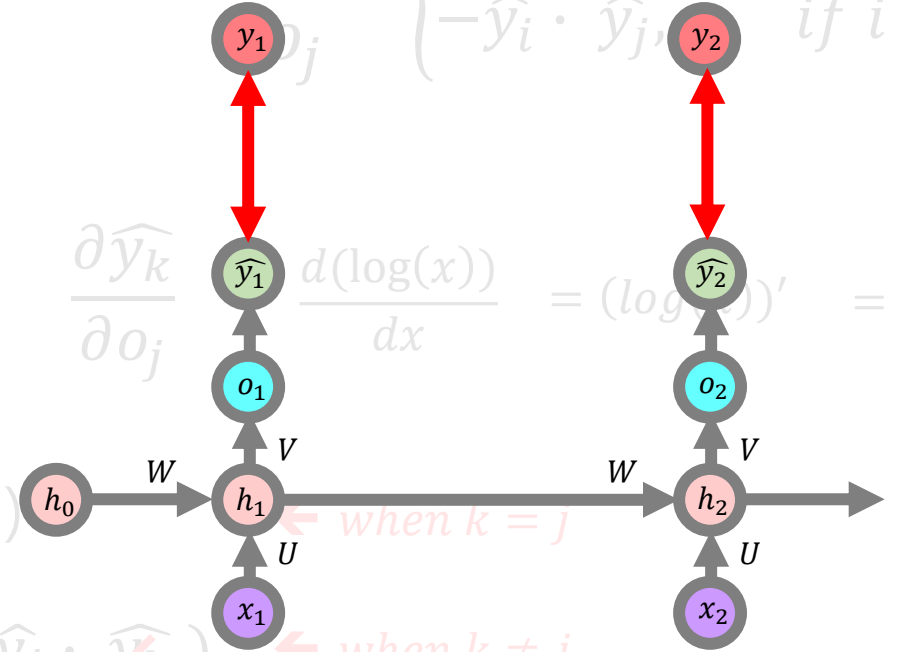
$$\frac{\partial L}{\partial y_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \frac{d(\log(x))}{dx} = (\log(x))' = \frac{1}{x}$$

$$= -y_j \frac{1}{\hat{y}_j} \hat{y}_j (1 - \hat{y}_j) + (- \sum_{k \neq j}^N y_k \frac{1}{\hat{y}_k} \hat{y}_j \cdot \hat{y}_k)$$

$$= -y_j (1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k) = \hat{y}_j - y_j$$



과정은 좀 복잡하지만 최종적인 그래디언트는 softmax 확률값에 실제값을 빼주기만 하면 되니

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V}$$

$$\frac{\partial L}{\partial V} = (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V}$$

$\frac{\partial \hat{y}_i}{\partial y_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$

$$= -y_j \frac{1}{\hat{y}_j} \hat{y}_j (1 - \hat{y}_j) + (-\sum_{k \neq j}^N y_k \frac{1}{\hat{y}_k} \hat{y}_j \cdot \hat{y}_k)$$

$$= -y_j(1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k) = \hat{y}_j - y_j$$

Softmax + Cross-entropy 조합이 많이 활용되고 있는 것 같습니다.

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial V} + \frac{\partial L_2}{\partial V} = - \sum_{k=1}^N y_k \log(\hat{y}_k)$$

$$\frac{\partial L}{\partial V} = \frac{\partial L_1}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial o_1} \frac{\partial o_1}{\partial V} + \frac{\partial L_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial o_2} \frac{\partial o_2}{\partial V}$$

$$\frac{\partial L}{\partial V} = (\hat{y}_1 - y_1) \frac{\partial o_1}{\partial V} + (\hat{y}_2 - y_2) \frac{\partial o_2}{\partial V}$$

$$\frac{\partial \hat{y}_i}{\partial y_j} = \begin{cases} \hat{y}_i (1 - \hat{y}_i), & \text{if } i = j \\ -\hat{y}_i \cdot \hat{y}_j, & \text{if } i \neq j \end{cases}$$

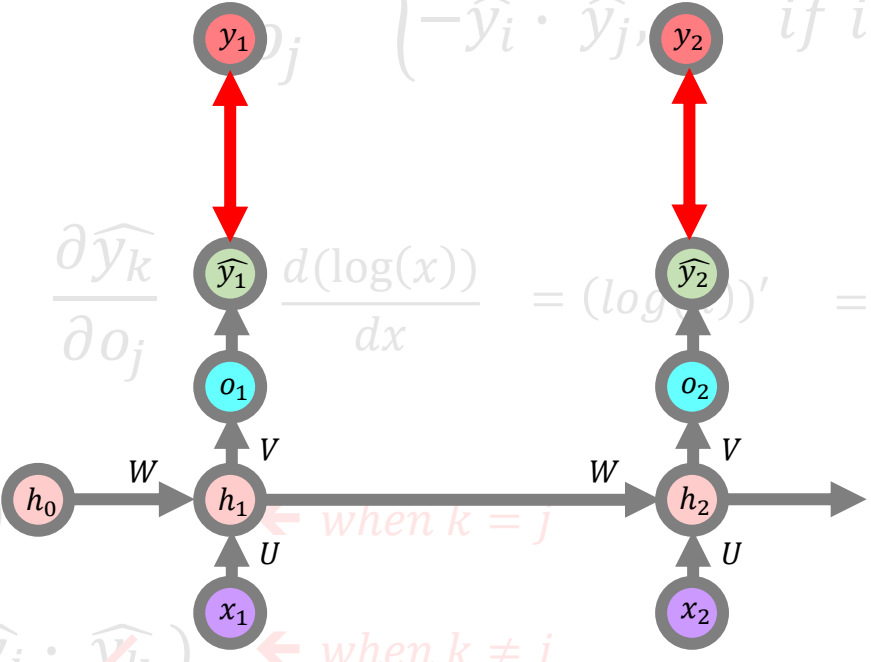
$$\frac{\partial L}{\partial o_j} = - \sum_{k=1}^N y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial o_j} \frac{d(\log(x))}{dx} = (\log(\hat{y}_j))' = \frac{1}{\hat{y}_j}$$

$$= -y_j \frac{1}{\hat{y}_j} \hat{y}_j (1 - \hat{y}_j) + (- \sum_{k \neq j}^N y_k \frac{1}{\hat{y}_k} \hat{y}_j \cdot \hat{y}_k)$$

$$= -y_j (1 - \hat{y}_j) + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + y_j \hat{y}_j + \sum_{k \neq j}^N y_k \hat{y}_j$$

$$= -y_j + \hat{y}_j (y_j + \sum_{k \neq j}^N y_k) = \hat{y}_j - y_j$$



여기까지가 오늘 제가 준비한 softmax와  
cross-entropy의 기울기에 관한  
영상입니다.

본의 아니게 변변한 그림 하나 없이  
수식만으로 설명 드리게 되어  
죄송하게 생각합니다 ☹



다음 영상에서는 보다 더 재미있는 주제로  
찾아뵙도록 하겠습니다. 계속 영상 시청  
부탁드립니다. 😊

그럼 다음 시간에 또 만나요!

감사합니다!

좋은 하루 되세요!!

이 채널은 여러분의 관심과 사랑이 필요합니다

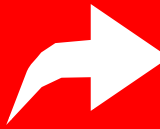
좋아요



댓글



공유



구독



‘좋아요’와 ‘구독’버튼은 강의 준비에 큰 힘이 됩니다!

좋아요



댓글



공유



구독



그리고 영상 자료를 사용하실때는  
출처 '신박AI'를 밝혀주세요





Copyright © 2024 by 신박AI

All rights reserved

본 문서(PDF)에 포함된 모든 내용과 자료는 저작권법에 의해 보호받고 있으며, 신박AI에 의해 제작되었습니다.

본 자료는 오직 개인적 학습 목적과 교육 기관 내에서의 교육용으로만 무료로 제공됩니다.

이를 위해, 사용자는 자료 내용의 출처를 명확히 밝히고,

원본 내용을 변경하지 않는 조건 하에 본 자료를 사용할 수 있습니다.

상업적 사용, 수정, 재배포, 또는 이 자료를 기반으로 한 2차적 저작물 생성은 엄격히 금지됩니다.

또한, 본 자료를 다른 유튜브 채널이나 어떠한 온라인 플랫폼에서도 무단으로 사용하는 것은 허용되지 않습니다.

본 자료의 어떠한 부분도 상업적 목적으로 사용하거나 다른 매체에 재배포하기 위해서는 신박AI의 명시적인 서면 동의가 필요합니다.

위의 조건들을 위반할 경우, 저작권법에 따른 법적 조치가 취해질 수 있음을 알려드립니다.

본 고지 사항에 동의하지 않는 경우, 본 문서의 사용을 즉시 중단해 주시기 바랍니다.

