

(a)

$$-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -[0 \cdot \log(\hat{y}_0) + 0 \cdot \log(\hat{y}_1) + \dots + 1 \cdot \log(\hat{y}_0) + \dots + 0 \cdot \log(\hat{y}_{\text{Vocab}})]$$

$$= -\log(\hat{y}_0) \neq$$

(c)

$$J_{\text{naive-softmax}}(v_c, 0, U) = -\log \frac{\exp(u_0^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}$$

$$= \underbrace{-u_0^T v_c}_{\textcircled{1}} + \underbrace{\log \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}_{\textcircled{2}}$$

(b)

$$\frac{\partial \textcircled{1}}{\partial v_c} = -u_0 = -U y$$

$$\frac{\partial \textcircled{2}}{\partial v_c} = \frac{1}{\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T v_c)} \cdot \frac{\partial}{\partial v_c} \left[ \sum_{w \in \text{Vocab}} \exp(u_w^T v_c) \right]$$

$$= \frac{1}{\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T v_c)} \cdot \sum_{w \in \text{Vocab}} [\exp(u_w^T v_c) \cdot u_w]$$

$$= \sum_{w \in \text{Vocab}} \left[ \frac{\exp(u_w^T v_c)}{\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T v_c)} u_w \right] = \sum_{w \in \text{Vocab}} \hat{y}_w \cdot u_w$$

$$= U \hat{y}$$

$$\Rightarrow \frac{\partial J_{\text{naive-softmax}}}{\partial v_c} = U(\hat{y} - y) \neq$$

① the gradient is zero when  $\hat{y} = y$   
 ② subtracting the gradient is to update  $v_c$  to let  $y$  be closest to the one-hot vector whose 0 entry is 1. It is equivalent to increase  $P(O=0|C=)$

(c)

if  $w=0$ ,

$$\frac{\partial \textcircled{1}}{\partial u_w} = -v_c$$

non-zero only when  $w''=w$ 

$$\frac{\partial \textcircled{2}}{\partial u_w} = \frac{1}{\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T v_c)} \cdot \frac{\partial}{\partial u_w} \left( \sum_{w'' \in \text{Vocab}} \exp(u_{w''}^T v_c) \right)$$

$$= \frac{1}{\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T v_c)} \cdot \exp(u_w^T v_c) \cdot v_c$$

$$= \hat{y}_w \cdot v_c$$

$$\Rightarrow \frac{\partial J_{\text{naive-softmax}}}{\partial u_w} = (\hat{y}_w - 1) v_c \stackrel{w=0}{=} (\hat{y}_0 - 1) v_c \quad \#$$

if  $w \neq 0$ 

$$\frac{\partial \textcircled{1}}{\partial u_w} = 0$$

$$\frac{\partial \textcircled{2}}{\partial u_w} = \hat{y}_w \cdot v_c$$

$$\Rightarrow \frac{\partial J_{\text{naive-softmax}}}{\partial u_w} = \hat{y}_w \cdot v_c \quad \#$$

$$(d) \frac{\partial J_{\text{naive-softmax}}}{\partial u} = \begin{bmatrix} \frac{\partial J}{\partial u_1} & \frac{\partial J}{\partial u_2} & \dots & \frac{\partial J}{\partial u_{|\text{Vocab}|}} \end{bmatrix} \in \mathbb{R}^{d \times |\text{Vocab}|} \quad \#$$

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

$$\sigma'(x) = \frac{(e^x+1) \cdot e^x - e^x \cdot e^x}{(e^x+1)^2} = \frac{e^x}{e^x+1} \cdot \frac{1}{e^x+1} = \sigma(x) \cdot \left(1 - \frac{e^x}{e^x+1}\right) = \sigma(x) [1 - \sigma(x)] \quad \#$$

$$J_{\text{neg-sample}}(v_c, 0, U) = \underbrace{-\log(\sigma(u_0^T v_c))}_{\textcircled{1}} - \underbrace{\sum_{k=1}^K \log(\sigma(-u_k^T v_c))}_{\textcircled{2}}$$

f)

$$\frac{\partial \textcircled{1}}{\partial v_c} = -\frac{1}{\sigma(u_0^T v_c)} \cdot \sigma'(u_0^T v_c) \cdot u_0$$

$$= [\sigma(u_0^T v_c) - 1] u_0$$

$$\frac{\partial \textcircled{2}}{\partial v_c} = -\sum_{k=1}^K \left[ \frac{1}{\sigma(-u_k^T v_c)} \cdot \sigma'(-u_k^T v_c) \cdot -u_k \right]$$

$$= \sum_{k=1}^K [1 - \sigma(-u_k^T v_c)] u_k$$

$$\Rightarrow \frac{\partial J_{\text{neg-sample}}}{\partial v_c} = [\sigma(u_0^T v_c) - 1] u_0 - \sum_{k=1}^K [\sigma(-u_k^T v_c) - 1] u_k \quad \#$$

$$\frac{\partial \textcircled{1}}{\partial u_0} = [1 - \sigma(u_0^T v_c)] v_c$$

$$\frac{\partial \textcircled{2}}{\partial u_0} = 0 \quad \text{because } 0 \notin \{w_1, w_2, \dots, w_K\}$$

$$\Rightarrow \frac{\partial J_{\text{neg-sample}}}{\partial u_0} = [1 - \sigma(u_0^T v_c)] v_c \quad \#$$

$$\frac{\partial \textcircled{1}}{\partial u_k} = 0 \quad \text{because } 0 \notin \{w_1, w_2, \dots, w_K\}$$

$$\begin{aligned} \frac{\partial \textcircled{2}}{\partial u_k} &= -\frac{\partial}{\partial u_k} \sum_{k'=1}^K \log(\sigma(-u_{k'}^T v_c)) = -\frac{\partial}{\partial u_k} \log(\sigma(-u_k^T v_c)) \\ &= \frac{-1}{\sigma(-u_k^T v_c)} \cdot \sigma'(-u_k^T v_c) \cdot (-v_c) = [1 - \sigma(-u_k^T v_c)] v_c \\ &\Rightarrow \frac{\partial J_{\text{neg-sample}}}{\partial u_k} = [1 - \sigma(-u_k^T v_c)] v_c \quad \# \end{aligned}$$

(g) When  $K$  negative samples may not be distinct

$$\frac{\partial \textcircled{1}}{\partial u_k} = 0 \quad \text{because } 0 \notin \{w_1, w_2, \dots, w_K\}$$

$$\frac{\partial \textcircled{2}}{\partial u_k} = -\frac{\partial}{\partial u_k} \left( \sum_{\substack{1 \leq k' \leq K \\ w_{k'} = w_k}} \log \sigma(-u_{k'}^T v_c) + \sum_{\substack{1 \leq k' \leq K \\ w_{k'} \neq w_k}} \log \sigma(-u_{k'}^T v_c) \right)$$

$$= - \sum_{\substack{1 \leq k' \leq K \\ w_{k'} = w_k}} \frac{1}{\sigma(-u_{k'}^T v_c)} \cdot \sigma'(-u_{k'}^T v_c) \cdot (-v_c)$$

$$= \sum_{\substack{1 \leq k' \leq K \\ w_{k'} = w_k}} [1 - \sigma(-u_{k'}^T v_c)] v_c \quad \#$$

(h)  $J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U) = \sum_{\substack{-m \leq \bar{j} \leq m \\ \bar{j} \neq 0}} J(v_c, w_{t+\bar{j}}, U)$

$$(i) \frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{\substack{-m \leq \bar{j} \leq m \\ \bar{j} \neq 0}} \frac{\partial J(v_c, w_{t+\bar{j}}, U)}{\partial U} \quad \#$$

$$(ii) \frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum_{\substack{-m \leq \bar{j} \leq m \\ \bar{j} \neq 0}} \frac{\partial J(v_c, w_{t+\bar{j}}, U)}{\partial v_c} \quad \#$$

$$(iii) \frac{\partial J_{\text{skip-gram}}}{\partial v_w} = 0 \quad \text{when } w \neq c \quad \#$$