

center word  
& window

$$\left\{ \begin{array}{l} P(O|C) = P(O=o|C=c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \\ u_o \in U = \begin{bmatrix} | & & | \\ u_w & \cdots & u_w \\ | & & | \end{bmatrix} \quad v_c \in V = \begin{bmatrix} | \\ v_w \\ \vdots \\ | \end{bmatrix} \end{array} \right.$$

$$J_{\text{naive-softmax}}(v_c, o, U) = -\log P(O|C)$$

(a)  $-\sum_{w \in V} y_w \log \hat{y}_w = -\log(\hat{y}_o)$

The mechanism of  $y_w$  (one-hot) is to choose the element of the highest probability in  $\hat{y}_w$ .

(b)

(i)  $\nabla_{V_c} J(v_c, o, U) = -\nabla_{V_c} \log \frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)}$

$$= -\nabla_{V_c} \log \exp(u_o^T v_c) + \nabla_{V_c} \log \sum_w \exp(u_w^T v_c)$$

$$= -\frac{\exp(u_o^T v_c)}{\exp(u_o^T v_c)} \cdot u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \nabla_{V_c} \exp(u_w^T v_c)$$

$$= -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) \cdot u_w$$

$$= -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) \cdot u_w$$

$$= -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \hat{y}_w$$

$$= -u_o + U^T \hat{y} = U^T (\hat{y} - y)$$

(ii)

$$-u_o + U^T \hat{y} = 0 \Rightarrow U^T \hat{y} = u_o \quad (\text{right prediction})$$

When  $\hat{y} = (U^T)^{-1} u_o$ , the gradient = 0 hold True

(iii)

$$V_c := V_c + \alpha u_o - \alpha U^T \hat{y}$$

$\downarrow$   
make the  $V_c$  more similar to  $u_o$  with LR "  $\alpha$ "

$\downarrow$   
similar to expected class  $u_w$  (soft)

(iv)  $u_c(\text{good}) \approx \alpha u_c(\text{best})$

when  $L_2$ -normalization.  $\frac{u_c(\text{good})}{\|u_c\|_2} \approx \frac{\alpha u_c(\text{best})}{\|\alpha u_c\|_2} = \frac{u_c(\text{best})}{\|u_c\|_2}$

the semantic change happens.

$$\begin{aligned}
 \text{c)} \quad \nabla_{U_c} J(V_c, O, U) &= -\nabla_{U_c} \log \frac{\exp(U_o^T V_c)}{\sum_{w \in V} \exp(U_w^T V_c)} \\
 &= -\nabla_{U_c} \log \exp(U_o^T V_c) + \nabla_{U_c} \log \sum_{w \in V} \exp(U_w^T V_c) \\
 &= -\nabla_{U_c} U_o^T V_c + \frac{1}{\sum_{w \in V} \exp(U_w^T V_c)} \cdot \sum_{w \in V} \exp(U_w^T V_c) \cdot \nabla_{U_c} U_w^T V_c \\
 &= -\nabla_{U_c} U_o^T V_c + \frac{1}{\sum_{w \in V} \hat{y}_w} \nabla_{U_c} \cdot \sum_{w \in V} U_w^T V_c \\
 &= (\frac{1}{\sum_{w \in V} \hat{y}_w} \nabla_{U_c} U_o^T - \nabla_{U_c} U_o^T) \cdot V_c \\
 &= (\hat{y} - y)^T V_c
 \end{aligned}$$

$$\begin{aligned}
 \text{d)} \quad \nabla_{U_c} J(V_c, O, U) &= \left[ \frac{\partial J(V_c, O, U)}{\partial U_1}, \frac{\partial J(V_c, O, U)}{\partial U_2}, \dots, \frac{\partial J(V_c, O, U)}{\partial U_{|\text{vocab}|}} \right] \\
 &= [\hat{y}_1^T V_c, \hat{y}_2^T V_c, \dots, \hat{y}_{|\text{vocab}|}^T V_c] \\
 &\quad \text{one-hot} \\
 &= (\hat{y} - y)^T V_c
 \end{aligned}$$

e) Leaky ReLU :  $f(x) = \max(\alpha x, x)$

$$\frac{\partial}{\partial x} f(x) = \begin{cases} \alpha & x < 0 \\ 1 & x > 0 \end{cases}$$

f) Sigmoid:

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

$$\begin{aligned}
 \sigma'(x) &= \frac{1}{(1 + e^{-x})^2} \cdot -e^{-x} \\
 &= \frac{1}{1 + e^{-x}} \cdot \frac{1 - e^{-x}}{1 + e^{-x}} \\
 &= \sigma(x) (1 - \sigma(x))
 \end{aligned}$$

$$\text{g) } i) \quad J_{\text{neg-sample}}(V_c, O, U) = -\log \sigma(U_o^T V_c) - \sum_{s=1}^k \log (\sigma(-U_{ws}^T V_c))$$

$$\begin{aligned}
 \nabla_{V_c} J(V_c, O, U) &= -\nabla_{V_c} \log \sigma(U_o^T V_c) - \sum_{s=1}^k \nabla_{V_c} \log (\sigma(-U_{ws}^T V_c)) \\
 &= -\frac{1}{\sigma(U_o^T V_c)} \cdot \sigma(U_o^T V_c) (1 - \sigma(U_o^T V_c)) U_o - \sum_{s=1}^k \frac{\sigma(-U_{ws}^T V_c)}{\sigma(-U_{ws}^T V_c)} (1 - \sigma(-U_{ws}^T V_c)) U_{ws} \\
 &= -(1 - \sigma(U_o^T V_c)) U_o + \sum_{s=1}^k (1 - \sigma(-U_{ws}^T V_c)) U_{ws}
 \end{aligned}$$

$$\nabla_{U_o} J(V_c, O, U) = -(1 - \sigma(U_o^T V_c)) V_c$$

$$\nabla_{U_{ws}} J(V_c, O, U) = \sum_{s=1}^k (1 - \sigma(-U_{ws}^T V_c)) V_c$$

$$(ii) U_{0, \{w_1, \dots, w_k\}} = \begin{bmatrix} u_0 \\ -u_{w_1} \\ \vdots \\ -u_{w_k} \end{bmatrix} \in \mathbb{R}^{(k+1) \times d} \quad V_c \in \mathbb{R}^{d \times 1}$$

$$\Rightarrow \sigma(1 - U_{0, \{w_1, \dots, w_k\}} V_c) \in \mathbb{R}^{(k+1) \times 1}$$

$$\Rightarrow \nabla_{V_c} J(V_c, 0, U) = -U_{0, \{w_1, \dots, w_k\}}^T \sigma(1 - U_{0, \{w_1, \dots, w_k\}} V_c) \quad ①$$

$$\nabla_{U_{0, \{w_1, \dots, w_k\}}} J(V_c, 0, U) = \frac{U_{0, \{w_1, \dots, w_k\}}}{|U_{0, \{w_1, \dots, w_k\}}|} \odot \underbrace{\sigma(1 - U_{0, \{w_1, \dots, w_k\}} V_c)}_{(k+1) \times 1} \cdot \underbrace{V_c^T}_{1 \times d}$$

(iii) This function only sample  $k$  negative samples, but the naive softmax loss need to calculate whole dataset.

(h)

$$J_{\text{neg-sample}}(V_c, 0, U) = -\log(\sigma(U_0^T V_c)) - \sum_{s=1}^k \log(\sigma(-U_{ws}^T V_c))$$

from g(i), we have  $\nabla_{U_{ws}} J_{\text{neg-sample}}(V_c, 0, U) = (1 - \sigma(-U_{ws}^T V_c)) V_c$

for **not distinct** case:  $\nabla = \text{sum of all } w = w_s$

$$\begin{aligned} \nabla_{U_{ws}} J_{\text{neg-sample}}(V_c, 0, U) &= \sum_{\substack{w \in V \\ w=w_s}} (1 - \sigma(U_w^T V_c)) V_c \\ &= \text{num\_non-distinct} \times (1 - \sigma(U_{ws}^T V_c)) V_c \\ &= |\{w \in \text{vocab} \mid w = w_s\}| (1 - \sigma(U_{ws}^T V_c)) V_c \end{aligned}$$

(i)

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

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

$$\frac{\partial J_{\text{skip-gram}}(V_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U_w} = 0$$