

$$(b) \quad J_{\text{naive-softmax}}(\vec{v}_0, o, \vec{u}) = -\log P(o=o | C=c)$$

$$= -\log \frac{\exp(\vec{u}_o^T \vec{v}_c)}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)}$$

$$\frac{\partial}{\partial \vec{v}_c} J_{\text{naive-softmax}}$$

$$= \frac{\partial}{\partial \vec{v}_c} \left[-\log \exp(\vec{u}_o^T \vec{v}_c) + \log \sum_w \exp(\vec{u}_w^T \vec{v}_c) \right]$$

$$= -\frac{\partial}{\partial \vec{v}_c} \vec{u}_o^T \vec{v}_c + \frac{1}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \cdot \frac{\partial}{\partial \vec{v}_c} \sum_w \exp(\vec{u}_w^T \vec{v}_c)$$

$$= -\vec{u}_o + \frac{1}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \sum_w \frac{\partial}{\partial \vec{v}_c} \exp(\vec{u}_w^T \vec{v}_c)$$

$$= -\vec{u}_o + \frac{1}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \sum_w \exp(\vec{u}_w^T \vec{v}_c) \cdot \vec{u}_w$$

$$= -\vec{u}_o + \sum_w \frac{\exp(\vec{u}_w^T \vec{v}_c) \vec{u}_w}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)}$$

$$= -\vec{u}_o + \underbrace{\sum_w P(o=w | C=c)}_{\text{predicted distribution}} \vec{u}_w$$

$$= -\vec{u} \cdot \vec{y} + \vec{u} \cdot \hat{\vec{y}} = \vec{u}(\hat{\vec{y}} - \vec{y})$$

\vec{y} : true distribution

$\hat{\vec{y}}$: predicted distribution

$$(c) J_{\text{naive-softmax}}(\vec{v}_c, 0, 1) = -\log \frac{\exp(\vec{u}_0^T \vec{v}_c)}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)}$$

① $w \neq 0$

$$\begin{aligned} \frac{\partial}{\partial \vec{u}_w} J &= -\cancel{\frac{\partial}{\partial \vec{u}_w} \vec{u}_0^T \vec{v}_c} + \frac{\partial}{\partial \vec{u}_w} \log \sum_w \exp(\vec{u}_w^T \vec{v}_c) \\ &= 0 + \frac{1}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \sum_w \frac{\partial}{\partial \vec{u}_w} \exp(\vec{u}_w^T \vec{v}_c) \\ &= \frac{1}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \frac{\partial}{\partial \vec{u}_w} \exp(\vec{u}_w^T \vec{v}_c) \\ &= \frac{\exp(\vec{u}_w^T \vec{v}_c)}{\sum_w \exp(\vec{u}_w^T \vec{v}_c)} \cdot \vec{v}_c \\ &= P(0=w | C=c) \cdot \vec{v}_c \\ &= \hat{y}_w \cdot \vec{v}_c \end{aligned}$$

② $w=0$

$$\begin{aligned} \frac{\partial}{\partial \vec{u}_0} J &= -\vec{v}_c + \hat{y}_0 \cdot \vec{v}_c \\ &= -y_0 \vec{v}_c + \hat{y}_0 \vec{v}_c \\ &= (\hat{y}_0 - y_0) \vec{v}_c \end{aligned}$$

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

重要性质: $\sigma(x) + \sigma(-x) = 1$

(e)

$$\begin{aligned}
 \textcircled{1} \quad \frac{\partial}{\partial \vec{v}_c} J &= - \frac{\partial}{\partial \vec{v}_c} \log(\sigma(\vec{u}_0^T \vec{v}_c)) - \frac{\partial}{\partial \vec{v}_c} \sum_{k=1}^K \log(\sigma(-\vec{u}_k^T \vec{v}_c)) \\
 &= - \frac{1}{\sigma(\vec{u}_0^T \vec{v}_c)} \cdot \frac{\partial}{\partial \vec{v}_c} \sigma(\vec{u}_0^T \vec{v}_c) - \sum_{k=1}^K \frac{1}{\sigma(-\vec{u}_k^T \vec{v}_c)} \cdot \frac{\partial}{\partial \vec{v}_c} \sigma(-\vec{u}_k^T \vec{v}_c) \\
 &= (\sigma(\vec{u}_0^T \vec{v}_c) - 1) \cdot \frac{\partial}{\partial \vec{v}_c} \vec{u}_0^T \vec{v}_c + \sum_{k=1}^K (\sigma(-\vec{u}_k^T \vec{v}_c) - 1) \frac{\partial}{\partial \vec{v}_c} (-\vec{u}_k^T \vec{v}_c) \\
 &= (\sigma(\vec{u}_0^T \vec{v}_c) - 1) \cdot \vec{u}_0^T + \sum_{k=1}^K \vec{u}_k (1 - \sigma(-\vec{u}_k^T \vec{v}_c)) \\
 &= -\sigma(-\vec{u}_0^T \vec{v}_c) \vec{u}_0 + \sum_{k=1}^K \sigma(\vec{u}_k^T \vec{v}_c) \vec{u}_k
 \end{aligned}$$

$$\textcircled{2} \quad \frac{\partial}{\partial \vec{u}_0} J_{\text{neg sample}} = - \frac{\partial}{\partial \vec{u}_0} \log(\sigma(\vec{u}_0^T \vec{v}_c)) - \sum_{k=1}^K \frac{\partial}{\partial \vec{u}_0} \log(\sigma(-\vec{u}_k^T \vec{v}_c))$$

$$= - \frac{1}{\sigma(\vec{u}_0^T \vec{v}_c)} \frac{\partial}{\partial \vec{u}_0} \sigma(\vec{u}_0^T \vec{v}_c)$$

$$= - \sigma(-\vec{u}_0^T \vec{v}_c) \vec{v}_c$$

$$\textcircled{3} \quad \frac{\partial}{\partial \vec{u}_k} J_{\text{neg sample}}$$

$$= - \frac{\partial}{\partial \vec{u}_k} \log(\sigma(\vec{u}_0^T \vec{v}_c)) - \sum_{k=1}^K \frac{\partial}{\partial \vec{u}_k} \log(\sigma(-\vec{u}_k^T \vec{v}_c))$$

$$= - \frac{\partial}{\partial \vec{u}_k} \log(\sigma(-\vec{u}_k^T \vec{v}_c))$$

$$= \sigma(\vec{u}_k^T \vec{v}_c) \vec{v}_c$$

Only need to compute for at most K words, but the naive softmax has to go through the whole vocabulary, computing for $|V|$ words!

4)

$$(i) \frac{\partial}{\partial U} J_{\text{skip-gram}}(\vec{v}_c, w_{t-m}, \dots, w_{t+m}, U)$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial U} J(\vec{v}_c, w_{t+j}, U)$$

$$(ii) \frac{\partial}{\partial \vec{v}_c} J_{\text{skip-gram}}(\vec{v}_c, w_{t-m}, \dots, w_{t+m}, U)$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial \vec{v}_c} J(\vec{v}_c, w_{t+j}, U)$$

$$(iii) \frac{\partial}{\partial \vec{v}_w} J_{\text{skip-gram}}(\vec{v}_c, w_{t-m}, \dots, w_{t+m}, U)$$

$$= 0$$