

a) we know that y is a one-hot vector where $y_w = 1$ when $w = 0$ & $y_w = 0$ when $w \neq 0$. So:

$$\begin{aligned}
 - \sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) &= -[y_1 \log(\hat{y}_1) + \dots + y_0 \log(\hat{y}_0) + \dots] \\
 &= -[0 + \dots + y_0 \log(\hat{y}_0) + \dots] \\
 &= -y_0 \log(\hat{y}_0) = -\log(\hat{y}_0)
 \end{aligned}$$

b) $\frac{\partial J_{\text{naive-softmax}}}{\partial v_c} = \frac{\partial (-\log P(0|c))}{\partial v_c}$

$$\begin{aligned}
 &= \frac{\partial \left(-\log \frac{\exp(u_0^T v_c)}{\sum_w \exp(u_w^T v_c)} \right)}{\partial v_c} = \frac{\partial}{\partial v_c} \left(-\log(\exp(u_0^T v_c)) + \log\left(\sum_w \exp(u_w^T v_c)\right) \right)
 \end{aligned}$$

chain rule
↑

$$\frac{\partial}{\partial v_c} (-u_0^T v_c) + \log\left(\sum_w \exp(u_w^T v_c)\right) =$$

$$-u_0 + \sum_{w \in V} \frac{\exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} = (\hat{y} - y)u$$

$$\underline{\underline{c)}} \quad \frac{\partial \text{intra-softmax}}{\partial v_w} = \frac{\partial}{\partial v_w} \frac{-\exp(u_0^T v_c)}{\sum_w \exp(u_w^T v_c)}$$

$$= + \frac{\partial}{\partial v_w} \left(-\log \exp(u_0^T v_c) + \log \sum_w \exp(u_w^T v_c) \right) \quad (2)$$

$$\text{if } w=0: \frac{\partial}{\partial u_0} (-u_0^T v_c) = -v_c$$

$$(2): \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} v_c = p(o|c) v_c$$

$$\Rightarrow v_c (p(o|c) - 1)$$

$$\text{if } w \neq 0: \frac{\partial}{\partial u_w} (-u_0^T v_c) = 0$$

$$(2): \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} v_c = p(o|c) v_c$$

$$\Rightarrow p(o|c) v_c$$

$$\underline{\underline{d)}} \quad \frac{\partial \text{intra-softmax}}{\partial u} = \left[\frac{\partial j(v_c, o, u)}{\partial u_1}, \frac{\partial j(v_c, o, u)}{\partial u_r}, \dots, \frac{\partial j(v_c, o, u)}{\partial u_{\text{vocab}}} \right]$$

a matrix with each column has derivative

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

$$\frac{\partial \sigma(x)}{\partial x} = \frac{e^x(e^x + 1) - e^x e^x}{(e^x + 1)^2} = \frac{e^x(e^x - e^x + 1)}{(e^x + 1)^2} = \frac{e^x}{(e^x + 1)^2}$$

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

$$\Rightarrow \frac{e^x}{(e^x + 1)^2} = \sigma(x)(1 - \sigma(x))$$

f)

$$\textcircled{1} \frac{\partial j}{\partial v_c} = \frac{\partial}{\partial v_c} (-\log(\sigma(u_0^T v_c))) - \sum_k \log(\sigma(-u_k^T v_c)) =$$

$$-(1 - \sigma(u_0^T v_c))u_0 + \sum_k (1 - \sigma(-u_k^T v_c))u_k$$

$$\textcircled{2} \frac{\partial j}{\partial u_0} = -(1 - \sigma(u_0^T v_c))v_c$$

$$\textcircled{3} \frac{\partial j}{\partial u_k} = (1 - \sigma(-u_k^T v_c))v_c$$

in naive softmax function the whole U is computed but in negative sampling loss, only k of U is computed. so it is more efficient.

$$\underline{h)} \quad ① \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{j(\psi_c, \omega_t + j, u)}{\delta u}$$

$$② \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{j(\psi_c, \omega_t + j, u)}{\delta \psi_c}$$

$$③ = 0$$

