

(a) $y \rightarrow$ true distribution $\hat{y} \rightarrow$ predicted distribution

$$-\sum_{w \in V} y_w \log(\hat{y}_w) = -\left(y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2) + \dots + y_v \log(\hat{y}_v)\right) = -\log(\hat{y}_0)$$

all these values are zero except when $w=0$. In that case, $y_0 = 1$

$$(b) J_{\text{nsn}} = -\log(e^{u_0^T v_c} \times (\sum_{w \in V} e^{u_w^T v_c})) = \log(\sum_{w \in V} e^{u_w^T v_c}) - \log(e^{u_0^T v_c})$$

$(J_A) \qquad (J_B)$

$$J = J_A - J_B$$

$$\Rightarrow \frac{\partial J}{\partial v_c} = \frac{\partial J_A}{\partial v_c} - \frac{\partial J_B}{\partial v_c}$$

$$\frac{\partial J_A}{\partial v_c} = \frac{\partial (u_0^T v_c)}{\partial v_c} = u_0$$

$$\frac{\partial J_B}{\partial v_c} = \frac{\partial (\log \sum_{w \in V} e^{u_w^T v_c})}{\partial v_c} = \frac{1}{\sum_{w \in V} e^{u_w^T v_c}} \sum_{s \in V} u_s e^{u_s^T v_c} \rightarrow \hat{y}_s$$

chain rule

$$\Rightarrow \frac{\partial J_B}{\partial v_c} = \sum_{s \in V} \hat{y}_s u_s$$

$$\Rightarrow \frac{\partial J_{\text{nsn}}}{\partial v_c} = \sum_{w \in V} \hat{y}_w u_w - u_0$$

$$(c.1) \quad w=0 \quad \frac{\partial J_{\text{nsn}}}{\partial u_w} = \frac{\partial (-\log(\hat{y}_0))}{\partial u_w} = -\frac{\partial}{\partial u_w} \left(\log \frac{e^{u_0^T v_c}}{\sum_{w \in V} e^{u_w^T v_c}} \right)$$

$$= -\frac{\partial}{\partial u_w} (u_0^T v_c) + \frac{\partial}{\partial u_w} \left(\log \left(\sum_{w \in V} e^{u_w^T v_c} \right) \right) \quad \Leftarrow \text{generalized}$$

now we use the fact that $w=0$

$$\Rightarrow -v_c + \frac{1}{\sum_{w \in V} e^{u_w^T v_c}} (e^{u_w^T v_c} v_c)$$

$$= -v_c + v_c \hat{y}_w = v_c (\hat{y}_0 - 1)$$

(c.2) $w \neq 0$

$$\Rightarrow 0 + v_c \hat{y}_w = v_c \hat{y}_w$$

$$(d) \quad \frac{\partial J_{\text{softmax}}}{\partial u} = \left[\frac{\partial J(u_1, 0, u)}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_r} \right]$$

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

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

$$\Rightarrow \sigma'(x) = \left(\frac{1}{\sigma(x)} - 1 \right) \times \sigma(x) \times \sigma(x) = \sigma(x) (1 - \sigma(x))$$

$$(f) \quad \frac{\partial J_{\text{ns}}}{\partial v_c} = - \frac{\partial}{\partial v_c} \log(\sigma(u_0^T v_c)) - \frac{\partial}{\partial v_c} \left(\sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right) \quad (2)$$

$$\frac{-1}{\sigma(u_0^T v_c)} (u_0 \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c))) = u_0 (\sigma(u_0^T v_c) - 1)$$

$$(2): \quad \sum_{k=1}^K \frac{-u_k \sigma(A^T) (1 - \sigma(A))}{\sigma(A)} = \sum_{k=1}^K -u_k (1 - \sigma(-u_k^T v_c))$$

$$\frac{\partial J_{\text{ns}}}{\partial v_c} = u_0 (\sigma(u_0^T v_c) - 1) + \sum_{k=1}^K u_k (1 - \sigma(-u_k^T v_c))$$

• $\frac{\partial J_{\text{ns}}}{\partial u_0} \Rightarrow$ the Σ doesn't have any u_0 in it! just like previous part
The first part is
 $\rightarrow v_c (\sigma(u_0^T v_c) - 1)$

• $\frac{\partial J_{\text{ns}}}{\partial u_k} \Rightarrow$ First part has no u_k . In the second part, since there's only one match, Σ gets out of the way.

$$\frac{\partial J_{\text{ns}}}{\partial u_k} = v_c (1 - \sigma(-u_k^T v_c))$$

• The naive softmax considers the whole vocabulary for computation but the negative-sampling loss function only needs k samples and $K \ll V$ & is more efficient.

$$(g) \quad J_{ns} = -\log(\sigma(u_0^T v_c)) - \sum_{k \in A} \log(\sigma(-u_k^T v_c)) - \sum_{k \in B} \log(\sigma(-u_k^T v_c))$$

$\hookrightarrow \text{equal to } u_k$
 $\hookrightarrow \text{unequal to } u_k$

$$\frac{\partial J_{ns}}{\partial u_k} = 0 - \sum_{k \in A} \log(\sigma(-u_k^T v_c)) - 0$$

$$= n v_c (1 - \sigma(-u_k^T v_c))$$

$$(h) \quad (i) \quad \left[\frac{\partial J}{\partial w_{t-m}}, \frac{\partial J}{\partial w_{t-m+1}}, \dots, \frac{\partial J}{\partial w_{t+m}} \right]$$

$$(ii) \quad \sum_{-m \leq j \leq m} \frac{\partial J_{sg}(v_c, w_{t+j}, v)}{\partial v_c}$$

$$(iii) \quad \text{Zero}$$