

a

$$-\sum_{w \in V_{oc}} y_w \log(\hat{y}_w) = -\log(\hat{y}_0)$$

$$-\sum_{w \in V_{oc}} y_w \log(\hat{y}_w) \stackrel{y_w \neq 0}{=} -y_0 \log(\hat{y}_0) = -\log(\hat{y}_0)$$

$$\textcircled{b} \quad \frac{\partial J(v_c, 0, U)}{\partial v_c} = ?$$

$$= \frac{\partial -\log \frac{e^{u_0^T v_c}}{\sum_w e^{u_w^T v_c}}}{\partial v_c}$$

$$= \frac{-\partial u_0^T v_c}{\partial v_c} + \frac{\partial \log \sum_w e^{u_w^T v_c}}{\partial v_c}$$

$$= -u_0 + \frac{\sum_t e^{u_t^T v_c} u_t}{\sum_w e^{u_w^T v_c}}$$

$$= -u_0 + \sum_t \frac{e^{u_t^T v_c} u_t}{\sum_w e^{u_w^T v_c}}$$

$$= -u_0 + \sum_t \hat{y}_t u_t$$

$$= -Uy + U\hat{y} = U(\hat{y} - y)$$

C

$$w = 0$$

$$\frac{\partial (-u_0^T v_c + \log \sum_w e^{u_w^T v_c})}{\partial u_0}$$

$$= -v_c + \frac{e^{u_0^T v_c} v_c}{\sum_t e^{u_t^T v_c}}$$

$$= -v_c + \hat{y}_0 v_c$$

$$w \neq 0$$

$$\frac{\partial (-u_0^T v_c + \log \sum_w e^{u_w^T v_c})}{\partial u_w}$$

$$= 0 + \frac{e^{u_w^T v_c} v_c}{\sum_t e^{u_t^T v_c}}$$

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

2

$$U: [0 \times |W|]$$

$$V'_c = V_c \begin{bmatrix} 0 & 0 & 0 & \dots & \overset{\text{0th index}}{\downarrow} & \dots & 0 & 0 & 0 \end{bmatrix}_{[1 \times |W|]}$$

$$\frac{\partial J}{\partial U} = V'_c + V_c \underset{[0 \times 1]}{\overset{[1 \times |W|]}{\hat{y}^T}}$$

②

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

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

7

$$\frac{\partial \left( -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right)}{\partial v_c}$$

$$= \frac{-\cancel{\sigma(u_0^T v_c)}(1-\cancel{\sigma(u_0^T v_c)})u_0}{\cancel{\sigma(u_0^T v_c)}} - \sum_{k=1}^K \frac{\cancel{\sigma(-u_k^T v_c)}(1-\cancel{\sigma(-u_k^T v_c)})(-u_k)}{\cancel{\sigma(-u_k^T v_c)}}$$

$$= -(1-\sigma(u_0^T v_c))u_0 + \sum_{k=1}^K (1-\sigma(-u_k^T v_c))u_k$$

\_\_\_\_\_

$$\frac{\partial \left( -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right)}{\partial u_0}$$

$$= \frac{-\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c))v_c}{\sigma(u_0^T v_c)} = 0$$

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

\_\_\_\_\_

$$\begin{aligned}
 & \frac{\partial \left( -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right)}{\partial u_k} \\
 &= 0 - \frac{\sigma(-u_k^T v_c)(1 - \sigma(-u_k^T v_c))(-v_c)}{\sigma(-u_k^T v_c)} \\
 &= (1 - \sigma(-u_k^T v_c))v_c
 \end{aligned}$$

more efficient because this approach  
 computes over  $K$  samples ( $O(K)$ ) but in  
 naive softmax computing the denominator costs  $O(N)$   
 and we know that  $K \ll N$

9

suppose that s samples are equal to  $u_k$

$$J \left( -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right)$$

---


$$\frac{\partial}{\partial u_k}$$

$$= 0 - (s) \left( \frac{\sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c)) (-v_c)}{\sigma(-u_k^T v_c)} \right) + 0$$

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



(h)

(i)

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

(ii)

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

(iii) 0

One thing that can easily be seen is analogy of king : male :: queen : female. Also proximity of similar words such as “female” and “woman” is visible.

