

1.

a) Since  $y$  is a one-hot vector, the summation over all words is a summation of  $w=0$  only

$$-\sum_{w \in \text{voc}} y_w \log(\hat{y}_w) = -[y_0 \log(\hat{y}_0), \dots, y_0 \log(\hat{y}_0), \dots y_w \log(\hat{y}_w)] = -1 \cdot \log(\hat{y}_0) \\ = -\log(\hat{y}_0)$$

$$\begin{aligned} \text{b) } \frac{\partial J}{\partial v_c} - \log\left(\frac{\exp(U_0^T v_c)}{\sum_{w \in \text{voc}} \exp(U_w^T v_c)}\right) &= -\frac{\partial}{\partial v_c} \left( \frac{U_0^T v_c}{\log \sum \exp(U_w^T v_c)} \right) \\ &= \frac{\partial}{\partial v_c} (-U_0^T v_c) + \frac{1}{\sum \exp(U_w^T v_c)} \cdot \sum \exp(U_w^T v_c) \frac{\partial}{\partial v_c} U_w^T v_c \\ &= -U_0 + \sum_{x \in \text{voc}} \frac{\exp(U_x^T v_c)}{\sum \exp(U_w^T v_c)} \cdot U_x \\ &= -U_0 + \sum_x U_x \frac{\exp(U_x^T v_c)}{\sum \exp(U_w^T v_c)} \\ &= -U_0 + \sum_x U_x p(x|c) \\ &= U(\hat{y} - y) \end{aligned}$$

c) When  $w \neq 0$

$$\begin{aligned} \frac{\partial}{\partial U_w} - \log\left(\frac{\exp(U_0^T v_c)}{\sum_{w \in \text{voc}} \exp(U_w^T v_c)}\right) &= \frac{\log(\sum \exp(U_w^T v_c))}{\partial U_w} \\ &= \frac{\sum \exp(U_w^T v_c) \cdot \sum_x \exp(U_w^T v_c) \cdot v_c^T}{\sum_x \exp(U_w^T v_c)} \\ &= \sum_x \frac{\exp(U_w^T v_c)}{\exp(U_w^T v_c)} v_c^T \\ &= \hat{y} v_c^T \end{aligned}$$

$$\begin{aligned} w=0 &\Rightarrow -v_c^T + \hat{y} v_c^T \\ &= (\hat{y} - y) v_c^T \end{aligned}$$

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

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

$$\begin{aligned} \text{LHS: } & \frac{1}{\sigma(u_0^T v_c)} \frac{\partial}{\partial v_c} \sigma(u_0^T v_c) \\ &= \frac{1}{\sigma(u_0^T v_c)} \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c)) u_0 \\ &= (1 - \sigma(u_0^T v_c)) u_0 \end{aligned}$$

$$\begin{aligned} \text{RHS: } & \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \\ &= \sum_k \frac{1}{\sigma(-u_k^T v_c)} \frac{\partial}{\partial v_c} \sigma(-u_k^T v_c) \\ &= \sum_k \frac{1}{\sigma(-u_k^T v_c)} \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c)) (-u_k) \\ &= \sum_k (1 - \sigma(-u_k^T v_c)) (-u_k) \end{aligned}$$

$$\frac{\partial J}{\partial v_c} = - (1 - \sigma(u_0^T v_c)) u_0 + \sum_k (1 - \sigma(-u_k^T v_c)) u_k$$

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

$$\begin{aligned} \frac{\partial J}{\partial u_k} &= \log(\sigma(-u_k^T v_c)) \\ &= \frac{1}{\sigma(-u_k^T v_c)} \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c)) v_c^T \\ &= (1 - \sigma(-u_k^T v_c)) v_c^T \end{aligned}$$

We can use cached activation values to speed up the computation and there is no need to sum over the entire corpus to get probabilities.

$$f) \frac{\partial J_{sg}}{\partial u} (v_c, w_{-m}, \dots, w_{+m}, u) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{+j}, u)}{\partial u}$$

$$\frac{\partial J_{sg}}{\partial v_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{+j}, u)}{\partial v_c}$$

$$\frac{\partial J_{sg}}{\partial v_w} = 0 \quad (w \neq c)$$