

# 1 Written: Understanding word2vec (23 points)

a)  $y$  is one-hot vector, when  $w \neq o$ ,  $y_w = 0$ . So  $-\sum_{w \in V} y_w \log \hat{y}_w = \log \hat{y}_o$

b)  $U_o - \sum_{x=1}^V P(x|c) U_x = U(\hat{y} - y)$

c)

$$\begin{aligned} & \frac{\partial}{\partial u_w} \log \frac{\exp(u_o^T v_c)}{\sum_{x=1}^V \exp(u_x^T v_c)} \\ &= \frac{\frac{\partial}{\partial u_w} (\log \exp(u_o^T v_c) - \log \sum_{x=1}^V \exp(u_x^T v_c))}{v_c / 0} \\ &= \frac{\frac{\partial}{\partial u_w} \log \sum_{x=1}^V \exp(u_x^T v_c)}{\frac{1}{\sum_{x=1}^V \exp(u_x^T v_c)} \cdot \frac{\partial}{\partial u_w} \sum_{x=1}^V \exp(u_x^T v_c)} \\ &= \frac{\exp(u_w^T v_c)}{\sum_{x=1}^V \exp(u_x^T v_c)} v_c \quad \because \frac{\partial J}{\partial u_w} = \begin{cases} \hat{y}_w v_c, w \neq o \\ (\hat{y}_w - 1) v_c, w = o \end{cases} \\ &= \hat{y}_w v_c \end{aligned}$$

d)

$$\sigma(x) = \sigma(x)(1 - \sigma(x))$$

e)

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

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

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

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

With naïve-softmax loss, we don't need to compute through the whole vocabulary.

f)

$$(i) \quad \partial J_{\text{skip-gram}}(v_c, w_{t-m}), \dots, w_{t+m}, U / \partial U = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_o, w_{w+j}, U)}{\partial U}$$

$$(ii) \quad \partial J_{\text{skip-gram}}(v_c, w_{t-m}), \dots, w_{t+m}, U / \partial v_c = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_o, w_{w+j}, U)}{\partial v_c}$$

$$(iii) \quad \partial J_{\text{skip-gram}}(v_c, w_{t-m}), \dots, w_{t+m}, U / \partial v_w (\text{when } w \neq c) = 0$$