

Written part

1.

(a). Since the true empirical distribution \vec{y} is by one-hot encoding in the vector \vec{y} , only $y_o = 1$, all other $y_w = 0$, $w \in \text{Vocabulary}$

Hence, $-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o)$

$$\begin{aligned}
 (b) \frac{\partial J(v_c, o, U)}{\partial v_c} &= \frac{\partial -u_o^T v_c + \log\left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right)}{\partial v_c} \\
 &= -u_o + \frac{\frac{\partial}{\partial v_c} \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \\
 &= -u_o + \frac{\sum_{w \in \text{Vocab}} u_w \cdot \exp(u_w^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \\
 &= -u_o + \sum_{w' \in \text{Vocab}} \frac{\exp(u_{w'}^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} u_{w'} \\
 &= -u_o + \sum_{w' \in \text{Vocab}} P(u_{w'} | v_c) \cdot u_{w'} \\
 &= -U y + U \hat{y} \\
 &= U(\hat{y} - y)
 \end{aligned}$$

(c) 2 cases can be merged into 1

$$\begin{aligned}
 1) \text{ When } w = o \\
 \frac{\partial J(v_c, o, U)}{\partial u_w} &= \frac{\partial -u_o^T v_c + \log\left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right)}{\partial u_w} \\
 &= -v_c \cdot \mathbb{I}_{\{w=o\}} + P(w|c) \cdot v_c
 \end{aligned}$$

$$\begin{aligned}
 &= V_c (P(w|c) - \mathbb{I}_{\{w=0\}}) \\
 &= V_c (\hat{y}_w - y_w)
 \end{aligned}$$

(d) for $\vec{x}^T = [x_1, x_2, \dots, x_n]$
 the sigmoid function σ acts as

$$\sigma(\vec{x})^T = \left[\frac{e^{x_1}}{e^{x_1} + 1}, \frac{e^{x_2}}{e^{x_2} + 1}, \dots, \frac{e^{x_n}}{e^{x_n} + 1} \right]$$

$$\begin{aligned}
 \text{So } \frac{\partial \sigma(\vec{x})_i}{\partial x_j} &= \mathbb{I}_{\{i=j\}} \sigma(x_i) \cdot (1 - \sigma(x_i)) \\
 &= \text{diag}[\sigma(x) \cdot (1 - \sigma(x))^T]
 \end{aligned}$$

notice the dimension here, by taking transpose to get a matrix.

$$(e) J_{\text{neg-sample}}(V_c, 0, U) = -\log(\sigma(u_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T V_c))$$

$$\begin{aligned}
 \textcircled{1} \frac{\partial J_{\text{neg-sample}}(V_c, 0, U)}{\partial V_c} &= \frac{\partial -\log(\sigma(u_0^T V_c))}{\partial V_c} - \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^T V_c))}{\partial V_c} \\
 &= \frac{\text{diag}[\sigma(u_0^T V_c) \cdot (1 - \sigma(u_0^T V_c))] \cdot u_0}{-\sigma(u_0^T V_c)} \\
 &\quad - \sum_{k=1}^K \frac{\text{diag}[\sigma(-u_k^T V_c) \cdot (1 - \sigma(-u_k^T V_c))] \cdot -u_k}{\sigma(-u_k^T V_c)} \\
 &= -(1 - \sigma(u_0^T V_c)) \cdot u_0 + \sum_{k=1}^K (1 - \sigma(-u_k^T V_c)) \cdot u_k
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \frac{\partial J_{\text{neg-sample}}(V_c, 0, U)}{\partial u_0} &= \frac{\partial -\log(\sigma(u_0^T V_c))}{\partial u_0} \\
 &= \frac{\sigma(u_0^T V_c) \cdot (1 - \sigma(u_0^T V_c)) \cdot V_c}{-\sigma(u_0^T V_c)} \\
 &= -(1 - \sigma(u_0^T V_c)) \cdot V_c
 \end{aligned}$$

$$(3) \frac{\partial J_{\text{neg-sample}}(V_c, 0, U)}{\partial U_k} = (1 - \sigma(U_k^T V_c)) \cdot V_c$$

(f) define $J_{\text{skip-gram}}(V_c, W_{t-m}, \dots, W_{t+m}, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(V_c, W_{t+j}, U)$

$$(i) \frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(V_c, W_{t+j}, U)}{\partial U}$$

$$(ii) \frac{\partial J_{\text{skip-gram}}}{\partial V_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(V_c, W_{t+j}, U)}{\partial V_c}$$

(iii) when $w \neq c$

$$\frac{\partial J_{\text{skip-gram}}}{\partial V_w} = 0$$