

(a).

$$\begin{aligned}
-\sum_{w \in Vocab} y_w \log(\hat{y}_w) &= - \left[y_o \log(\hat{y}_o) + \sum_{\substack{w \in Vocab \\ w \neq o}} y_w \log(\hat{y}_w) \right] \\
&= - \left[1 \cdot \log(\hat{y}_o) + \sum_{\substack{w \in Vocab \\ w \neq o}} 0 \cdot \log(\hat{y}_w) \right] \\
&= -\log(\hat{y}_o)
\end{aligned}$$

(b).

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{v}_c} &= \frac{\partial(-\log(\hat{y}_o))}{\partial \mathbf{v}_c} \\
&= - \frac{\partial(\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)})}{\partial \mathbf{v}_c} \\
&= - \frac{\partial(\mathbf{u}_o^\top \mathbf{v}_c) - \partial(\log \sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c))}{\partial \mathbf{v}_c} \\
&= -\mathbf{u}_o + \frac{1}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \cdot \sum_{w \in Vocab} \mathbf{u}_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \\
&= -\mathbf{u}_o + \sum_{w \in Vocab} \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{k \in Vocab} \exp(\mathbf{u}_k^\top \mathbf{v}_c)} \cdot \mathbf{u}_w \\
&= -\mathbf{u}_o + \sum_{w \in Vocab} p(\mathbf{u}_w | \mathbf{v}_c) \cdot \mathbf{u}_w \\
&= -\mathbf{u}_o + \sum_{w \in Vocab} \hat{y}_w \mathbf{u}_w \\
&= -\mathbf{u}_o + \mathbf{U} \hat{\mathbf{y}} \\
&= -\mathbf{U} \mathbf{y} + \mathbf{U} \hat{\mathbf{y}} \\
&= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y})
\end{aligned}$$

(c).

case1: $w \neq o$:

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{u}_w} &= \frac{\partial(-\log(\hat{y}_o))}{\partial \mathbf{u}_w} \\
&= -\frac{\partial(\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)})}{\partial \mathbf{u}_w} \\
&= -\frac{\partial(\mathbf{u}_o^\top \mathbf{v}_c) - \partial(\log \sum_{k \in Vocab} \exp(\mathbf{u}_k^\top \mathbf{v}_c))}{\partial \mathbf{u}_w} \\
&= \frac{1}{\sum_{k \in Vocab} \exp(\mathbf{u}_k^\top \mathbf{v}_c)} \cdot \mathbf{v}_c \exp(\mathbf{u}_w^\top \mathbf{v}_c) \\
&= \mathbf{v}_c \cdot p(\mathbf{u}_w | \mathbf{v}_c) \\
&= \mathbf{v}_c \cdot \hat{y}_w
\end{aligned}$$

case 2: $w = o$:

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{u}_w} &= \frac{\partial(-\log(\hat{y}_o))}{\partial \mathbf{u}_o} \\
&= -\frac{\partial(\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)})}{\partial \mathbf{u}_o} \\
&= -\frac{\partial(\mathbf{u}_o^\top \mathbf{v}_c) - \partial(\log \sum_{k \in Vocab} \exp(\mathbf{u}_k^\top \mathbf{v}_c))}{\partial \mathbf{u}_o} \\
&= -\mathbf{v}_c + \frac{1}{\sum_{k \in Vocab} \exp(\mathbf{u}_k^\top \mathbf{v}_c)} \cdot \mathbf{v}_c \exp(\mathbf{u}_o^\top \mathbf{v}_c) \\
&= \mathbf{v}_c (p(\mathbf{u}_o | \mathbf{v}_c) - 1) \\
&= \mathbf{v}_c \cdot (\hat{y}_o - 1)
\end{aligned}$$

(d).

$$\frac{\partial \mathbf{J}}{\partial \mathbf{U}} = \left[\frac{\partial \mathbf{J}}{\partial \mathbf{u}_1}, \frac{\partial \mathbf{J}}{\partial \mathbf{u}_2}, \frac{\partial \mathbf{J}}{\partial \mathbf{u}_3}, \dots, \frac{\partial \mathbf{J}}{\partial \mathbf{u}_{|Vocab|}} \right]$$

(e).

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

(f).

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{v}_c} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))}{\partial \mathbf{v}_c} \\
&= \frac{\partial}{\partial \mathbf{v}_c} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))] - \frac{\partial}{\partial \mathbf{v}_c} \left[\sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\
&= \left[-\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)^2) \mathbf{u}_o \right] - \\
&\quad \left[\sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c)^2) (-\mathbf{u}_k) \right] \\
&= -\mathbf{u}_o (1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) + \sum_{k=1}^K \mathbf{u}_k (1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{u}_o} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))}{\partial \mathbf{u}_o} \\
&= \frac{\partial}{\partial \mathbf{u}_o} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))] - \frac{\partial}{\partial \mathbf{u}_o} \left[\sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\
&= \left[-\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)^2) \mathbf{v}_c \right] - 0 \\
&= -\mathbf{v}_c (1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))}{\partial \mathbf{u}_k} \\
&= \frac{\partial}{\partial \mathbf{u}_k} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))] - \frac{\partial}{\partial \mathbf{u}_k} \left[\sum_{t=1}^K \log(\sigma(-\mathbf{u}_t^\top \mathbf{v}_c)) \right] \\
&= 0 - \left[\frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c)^2) (-\mathbf{v}_c) \right] \\
&= \mathbf{v}_c (1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))
\end{aligned}$$

(g).

$$\begin{aligned}
\frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))}{\partial \mathbf{u}_k} \\
&= \frac{\partial}{\partial \mathbf{u}_k} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))] - \frac{\partial}{\partial \mathbf{u}_k} \left[\sum_{t=1}^K \log(\sigma(-\mathbf{u}_t^\top \mathbf{v}_c)) \right] \\
&= 0 - \left[\sum_{\mathbf{u}_t = \mathbf{u}_k} \frac{1}{\sigma(-\mathbf{u}_t^\top \mathbf{v}_c)} (\sigma(-\mathbf{u}_t^\top \mathbf{v}_c) - \sigma(-\mathbf{u}_t^\top \mathbf{v}_c)^2) (-\mathbf{v}_c) \right] \\
&= \sum_{\mathbf{u}_t = \mathbf{u}_k} \mathbf{v}_c (1 - \sigma(-\mathbf{u}_t^\top \mathbf{v}_c))
\end{aligned}$$

(h).

$$\begin{aligned}
\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U} &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U} \\
\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c
\end{aligned}$$

$w \neq c$:

$$\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_w$$