

1. Written

a) Given: $\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$

Similarly, $\text{softmax}(x+c)_i = \frac{e^{(x+c)_i}}{\sum_j e^{(x+c)_j}} = \frac{e^{x_i} e^c}{\sum_j e^{x_j} e^c} = \frac{e^c e^{x_i}}{e^c \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}}$

Therefore, $\text{softmax}(x) = \text{softmax}(x+c)$.

$= \text{softmax}(x)_i$.

b) Given: $\sigma(x) = \frac{1}{1+e^{-x}}$

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

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

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

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

Using the given definition of $\sigma(x)$, this can be expressed as:

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x) (1 - \sigma(x)).$$

c) i. $J = -\sum_i y_i \log(\hat{y}_i) = -\sum_i y_i \log(p(x_i|c)) = -\log \left(\frac{\exp(u_k^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \right)$ (Because $y_i = 1$).

$$= -\left(\log(\exp(u_k^T v_c)) - \log\left(\sum_{w=1}^W \exp(u_w^T v_c)\right) \right) = -(u_k^T v_c) - \log\left(\sum_{w=1}^W \exp(u_w^T v_c)\right)$$

$$\frac{\partial J}{\partial v_c} = -\left(\frac{\partial}{\partial v_c} (u_k^T v_c) - \frac{\partial}{\partial v_c} \log\left(\sum_{w=1}^W \exp(u_w^T v_c)\right) \right)$$

$$\frac{\partial}{\partial v_c} (u_k^T v_c) = u_k$$

Continued on next page

$$\begin{aligned}
\frac{\partial}{\partial v_c} \log \left(\sum_{w=1}^W \exp(u_w^T v_c) \right) &= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \frac{\partial}{\partial v_c} \exp(u_x^T v_c) \\
&= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \exp(u_x^T v_c) \cdot \frac{\partial}{\partial v_c} (u_x^T v_c) \\
&= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \exp(u_x^T v_c) \cdot u_x \\
&= \sum_{x=1}^W \frac{\exp(u_x^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot u_x = \sum_{x=1}^W P(u_x | v_c) \cdot u_x = \sum_{x=1}^W \hat{y}_x u_x
\end{aligned}$$

therefore,

$$\begin{aligned}
\frac{\partial J}{\partial v_c} &= - \left(u_k - \sum_{x=1}^W \hat{y}_x u_x \right) = -u_k + \sum_{x=1}^W \hat{y}_x u_x \\
&= -u \cdot y + u \cdot \hat{y} = u(-y + \hat{y}) = u(\hat{y} - y).
\end{aligned}$$

c) ii. From part i,

$$\begin{aligned}
J &= -(\log(\exp(u_0^T v_c)) - \log(\sum_{w=1}^W \exp(u_w^T v_c))) \\
&= -(u_0^T v_c) + \log \left(\sum_{w=1}^W \exp(u_w^T v_c) \right)
\end{aligned}$$

consider two cases for the value of u_w :

① $w=0$

$$\begin{aligned}
\frac{\partial}{\partial u_w} -(u_0^T v_c) &= -v_c & \frac{\partial}{\partial u_w} \log \left(\sum_{w=1}^W \exp(u_w^T v_c) \right) &= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \frac{\partial}{\partial u_w} \exp(u_x^T v_c) \\
& & &= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \exp(u_x^T v_c) \cdot v_c \\
& & &= \sum_{x=1}^W \frac{\exp(u_x^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot v_c = \sum_{x=1}^W \hat{y}_x v_c
\end{aligned}$$

$$\frac{\partial J}{\partial u_w} \Big|_{w=0} = -v_c + \sum_{x=1}^W \hat{y}_x v_c = v_c(1 - \hat{y}_w)$$

continued on
next page

② $w \neq 0$

$$\frac{\partial}{\partial u_w} - (u_0^T v_c) = 0$$

$$\begin{aligned} \frac{\partial}{\partial u_w} \log \left(\sum_{w=1}^W \exp(u_w^T v_c) \right) &= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \frac{\partial}{\partial u_w} \exp(u_x^T v_c) \\ &= \frac{1}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot \sum_{x=1}^W \exp(u_x^T v_c) \cdot v_c \\ &= \sum_{x=1}^W \frac{\exp(u_x^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot v_c = \sum_{x=1}^W \hat{y}_x v_c \end{aligned}$$

$$\left. \frac{\partial J}{\partial u_w} \right|_{w \neq 0} = 0 + \sum_{x=1}^W \hat{y}_x v_c = \hat{y}_w v_c$$

therefore,

$$\frac{\partial J}{\partial u_w} = \begin{cases} (1 - \hat{y}_w) v_c, & w=1 \\ \hat{y}_w v_c, & w \neq 0 \end{cases}$$

$$c) \text{ iii. } J = -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

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

$$\frac{\partial}{\partial v_c} (-\log(\sigma(u_0^T v_c))) = -\frac{1}{\sigma(u_0^T v_c)} \cdot \frac{\partial}{\partial v_c} (\sigma(u_0^T v_c)) = -\frac{1}{\sigma(u_0^T v_c)} \cdot \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c)) \cdot \frac{\partial}{\partial v_c} (u_0^T v_c)$$

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

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

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

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

therefore,

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

continued on
next page.

ii) Consider two cases for the value of μ_w :

① $w=0$

$$\begin{aligned} \frac{\partial}{\partial \mu_w} (-\log(\sigma(\mu_0^T v_c))) &= -\frac{1}{\sigma(\mu_0^T v_c)} \cdot \sigma(\mu_0^T v_c)(1-\sigma(\mu_0^T v_c)) \cdot \frac{\partial}{\partial \mu_w} (\mu_0^T v_c) \\ &= -(1-\sigma(\mu_0^T v_c)) v_c = (\sigma(\mu_0^T v_c) - 1) v_c \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mu_w} \left(\sum_{k=1}^K \log(\sigma(-\mu_k^T v_c)) \right) &= \sum_{k=1}^K \frac{1}{\sigma(-\mu_k^T v_c)} \cdot \sigma(-\mu_k^T v_c)(1-\sigma(-\mu_k^T v_c)) \cdot \underbrace{\frac{\partial}{\partial \mu_w} (-\mu_k^T v_c)}_{=0} \\ \frac{\partial J}{\partial \mu_w} \Big|_{w=0} &= (\sigma(\mu_0^T v_c) - 1) v_c \end{aligned}$$

② $w \neq 0$ (or, $w=k$)

$$\frac{\partial}{\partial \mu_w} (-\log(\sigma(\mu_0^T v_c))) = -\frac{1}{\sigma(\mu_0^T v_c)} \cdot \sigma(\mu_0^T v_c)(1-\sigma(\mu_0^T v_c)) \cdot \underbrace{\frac{\partial}{\partial \mu_w} (\mu_0^T v_c)}_{=0}$$

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

$$\frac{\partial J}{\partial \mu_w} \Big|_{w \neq 0} = \sum_{k=1}^K (1-\sigma(-\mu_k^T v_c)) v_c$$

Therefore,

$$\frac{\partial J}{\partial \mu_w} = \begin{cases} (\sigma(\mu_0^T v_c) - 1) v_c, & w=0 \\ \sum_{k=1}^K (1-\sigma(-\mu_k^T v_c)) v_c, & w \neq 0. \end{cases}$$

v. $J_{\text{skip-gram}}(\text{word}_c, \dots, \text{word}_e, \dots, \text{word}_m) = \sum_{-m \leq j \leq m, j \neq 0} F(w_{c+j}, v_c)$

Use $F(0, v_c)$ as placeholder for JCE or Jneg-sample.

$$\frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(0, v_c)}{\partial v_c}$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mu_k} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(0, v_c)}{\partial \mu_k}$$