#1-(a). Because only $o$ is the position that the value is $1$, and $0$ everywhere else.

#1-(b).

Substituting (1) into (2), then

$$J = \log \sum_w \left( \exp(U_w^T V_c) \right) - U_o^T V_c$$

Here, $\dfrac{\partial J}{\partial V_c} = -U_o + \dfrac{\overset{\hat{y}_w}{\overbrace{\sum_w \exp(U_w^T V_c) \cdot U_w}}}{\sum_w \exp(U_w^T V_c)} = -U_o + \sum_w \hat{y}_w \cdot U_w$

$$[U] = \begin{bmatrix} | & | & & | \\ U_1 & U_2 & \cdots & U_w & \cdots \\ | & | & & | \end{bmatrix}$$ and $y$ is one-hot vector, so $[U]\,y$ indicates $U_o$

Also, $[U]\hat{y}$ indicates $\sum_w \hat{y}_w \cdot U_w$. Therefore, $\underline{\underline{\dfrac{\partial J}{\partial V_c} = [U](\hat{y} - y)}}$

(1) The gradient will be zero if $\hat{y} = y$.

(2) By subtracting the gradient from $V_c$, $V_c$ will be adjusted to the direction that reduce the difference between predicted value and true value. In other words, $V_c$ will be closer to the outside word vectors in its context.

# 1-(c).

i) $w=0$ $(U_w = U_o)$

$$\dfrac{\partial J}{\partial U_w} = -V_c + \dfrac{\overset{\hat{y}_w = \hat{y}_o}{\overbrace{\exp(U_w^T V_c) \cdot V_c}}}{\sum_w \exp(U_w^T V_c)} = \underline{\underline{(\hat{y}_o - 1)V_c}}$$

ii) $w \neq 0$

$$\dfrac{\partial J}{\partial U_w} = \dfrac{\overset{\hat{y}_w}{\overbrace{\exp(U_w^T V_c) \cdot V_c}}}{\sum_w \exp(U_w^T V_c)} = \underline{\underline{\hat{y}_w V_c}}$$

# 1-(d). With the answer of 1-(c), and assuming $\forall = |Vocab|$,

$$\dfrac{\partial J}{\partial U_1} = \begin{cases} (\hat{y}_1 - 1)V_c & \text{if } w = 1 = 0 \\ \hat{y}_1 V_c & \text{if } w = 1 \neq 0 \end{cases}$$

$$\dfrac{\partial J}{\partial U_2} = \begin{cases} (\hat{y}_2 - 1)V_c & \text{if } w = 2 = 0 \\ \hat{y}_2 V_c & \text{if } w = 2 \neq 0 \end{cases} \cdots$$

$$\dfrac{\partial J}{\partial U_\forall} = \begin{cases} (\hat{y}_\forall - 1)V_c & \text{if } w = \forall = 0 \\ \hat{y}_\forall V_c & \text{if } w = \forall \neq 0 \end{cases}$$

# 1 -(e).

$$\frac{d\phi}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

# 1 - (f).

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

# 1 -(g).

$$J = -\log\left(\sigma(U_o^T V_c)\right) - \sum_{s=1}^{K} \log\left(\sigma(-U_{w_s}^T V_c)\right)$$

(i). For $V_c$,

$$\frac{\partial J}{\partial V_c} = - \frac{\sigma(U_o^T V_c)(1 - \sigma(U_o^T V_c)) \cdot U_o}{\sigma(U_o^T V_c)} - \sum_s \frac{\sigma(-U_{w_s}^T V_c) \cdot (1 - \sigma(-U_{w_s}^T V_c)) \cdot (-U_{w_s})}{\sigma(-U_{w_s}^T V_c)}$$

$$= -\left(1 - \sigma(U_o^T V_c)\right) U_o + \sum_s \left(1 - \sigma(-U_{w_s}^T V_c)\right) U_{w_s}$$

For $U_o$,

$$\frac{\partial J}{\partial U_o} = -\left(1 - \sigma(U_o^T V_c)\right) V_c$$

For $U_{w_s}$,

$$\frac{\partial J}{\partial U_{w_s}} = \left(1 - \sigma(-U_{w_s}^T V_c)\right) V_c$$

(ii)

$$U_{o, \{w_1, \cdots, w_k\}}^T \cdot V_c = \left[ U_o^T V_c, -U_{w_1}^T V_c, \cdots, -U_{w_k}^T V_c \right]$$

Therefore, we can reuse $\sigma\left(U_{o, \{w_1, \cdots, w_k\}}^T \cdot V_c\right) - \mathbb{1}$ .

(iii)

We need lesser vectors than the naive-softmax loss.

# 1 -(h).

$$J = -\log\left(\sigma(U_o^T V_c)\right) - \sum_{w_j = w_s} \log\left(\sigma(-U_{w_j}^T V_c)\right) - \sum_{w_j \neq w_s} \log\left(\sigma(-U_{w_j}^T V_c)\right)$$

Assuming there are $m$ $W_s$ words, $\quad \frac{\partial J}{\partial U_s} = m\left(1 - \sigma(-U_{w_s}^T V_c)\right) V_c$

(i)

$$\frac{\partial J\left(v_c, w_{t-m}, \cdots, w_{t+m}, U\right)}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J\left(v_c, w_{t+j}, U\right)}{\partial U}$$
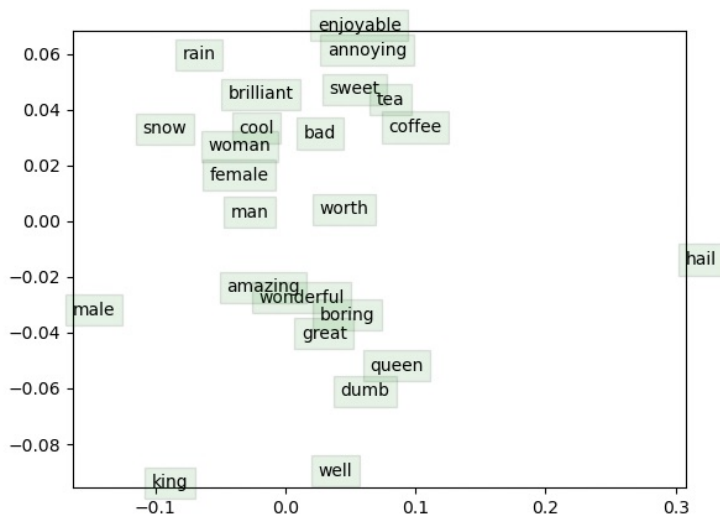
(ii)

$$\frac{\partial J\left(v_c, w_{t-m}, \cdots, w_{t+m}, U\right)}{\partial v_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J\left(v_c, w_{t+j}, U\right)}{\partial v_c}$$

(iii)

$$\frac{\partial J\left(v_c, w_{t-m}, \cdots, w_{t+m}, U\right)}{\partial v_w} = 0$$

# #2 - (C).



There are many well-clustered words such as (amazing, wonderful, boring, great). However, there are also not well-clustered words like "hail", which should have clustered with (rain, snow). Also, in evaluation, antonyms are considered as similar, and it is shown in the cluster "great (pos) ↔ boring (neg)"