Updated on: 2$^{nd}$ March, 2023

# Question (a)

Note that

$$
\boldsymbol{y}_w = \begin{cases} 1 & \text{if } w = o \\ 0 & \text{if } w \neq o \end{cases}
$$

Therefore,

$$
-\sum_{w \in \text{Vocab}} \boldsymbol{y}_w \log\left(\hat{\boldsymbol{y}}_w\right) = -(1) \log\left(\hat{\boldsymbol{y}}_o\right)
$$

$$
= -\log\left(\hat{\boldsymbol{y}}_o\right)
$$

# Question (b)

**Part (i)**

$$
J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U}) = -\log\left(\hat{\boldsymbol{y}}_o\right)
$$

$$
= -\log\left(\frac{\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)}\right)
$$

$$
= -\log\left(\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)\right) + \log\left(\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)\right)
$$

$$
= -\boldsymbol{u}_o^\top \boldsymbol{v}_c + \log\left(\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)\right)
$$

$$
\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{v}_c} = \frac{\partial}{\partial \boldsymbol{v}_c}\left[-\boldsymbol{u}_o^\top \boldsymbol{v}_c + \log\left(\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)\right)\right]
$$

$$
= -\boldsymbol{u}_o + \frac{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right) \boldsymbol{u}_k}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)}
$$

$$
= -\boldsymbol{u}_o + \sum_{k \in |\text{Vocab}|} \left(\frac{\exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)}\right) \boldsymbol{u}_k
$$

$$
= -\boldsymbol{u}_o + \sum_{k \in |\text{Vocab}|} \left(\hat{\boldsymbol{y}}_k\right) \boldsymbol{u}_k
$$

$$
= \sum_{k \in |\text{Vocab}|} \boldsymbol{y}_k(-\boldsymbol{u}_k) + \sum_{k \in |\text{Vocab}|} \left(\hat{\boldsymbol{y}}_k\right) \boldsymbol{u}_k
$$

$$
= \sum_{k \in |\text{Vocab}|} \left(\hat{\boldsymbol{y}}_k - \boldsymbol{y}_k\right) \boldsymbol{u}_k
$$

$$
= (\hat{\boldsymbol{y}}^\top - \boldsymbol{y}^\top)\boldsymbol{U}^\top \quad \text{(a row vector)}
$$

$$
= \boldsymbol{U}(\hat{\boldsymbol{y}} - \boldsymbol{y}) \quad \text{(a colume vector)}
$$

Note that a column vector is left as the answer to follow the shape convention.

**Part (ii)**

If the gradient computed in (b)(i) is equal to zero, i.e. the equation $U(\hat{y} - y) = 0$ is satisfied, then the vector $(\hat{y} - y)$ lies in the nullspace of the matrix $U$. Since the zero vector is always in the nullspace, a trival solution to the equation is $\hat{y} - y = 0$, i.e. the predicted word vector is correct.

**Part (iii)**

$$\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{v_c}} = \boldsymbol{U}(\hat{\boldsymbol{y}} - \boldsymbol{y})$$
$$= \boldsymbol{U}\hat{\boldsymbol{y}} - \boldsymbol{U}\boldsymbol{y}$$
$$= \text{predicted word vector} - \text{true word vector}$$

When the gradient is subtractedf from the word vecter $\boldsymbol{v_c}$ , $\boldsymbol{v_c}$ is improved towards the true word vector since the true word vector is added to $\boldsymbol{v_c}$. On the other hand, the updated vector is moved away from the predicted word vector since the predicted word vector is subtracted from $\boldsymbol{v_c}$.

**Part (iv)**

Consider two word vectors $\mathbf{u_x}$ and $\mathbf{u_y}$ of the two different words $\mathbf{x}$ and $\mathbf{y}$, where $\mathbf{u_x} = \alpha \mathbf{u_y}$ and $\alpha$ is some scalar.

If the two words are opposite in meanings, it is expected that these two word vectors are in opposite direction, i.e. $\alpha < 0$. In such case, since the L2 normalization preserves the direction of the word vectors, one can still distinguish the two word vectors with their directions.

In the case where the two words have the same sign ($\alpha > 0$) but with different extent (e.g. $\mathbf{x}$ is much positive than $\mathbf{y}$), one would expect such information is revealed from the vector magnitudes. However, the L2 normalization produces vectors with the same magnitudes. Information of word vectors carried by the vector magnitude is lost after performing L2 normalization.

# Question (c)

$$J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U}) = -\boldsymbol{u}_o^\top \boldsymbol{v_c} + \log\left(\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v_c}\right)\right)$$

When $w \neq o$,

$$\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_{w \neq o}}} = \frac{\partial}{\partial \boldsymbol{u_{w \neq o}}}\left[-\boldsymbol{u}_o^\top \boldsymbol{v_c} + \log\left(\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v_c}\right)\right)\right]$$
$$= 0 + \frac{\exp\left(\boldsymbol{u}_{w \neq o}^\top \boldsymbol{v_c}\right)\boldsymbol{v_c}}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v_c}\right)}$$
$$= \left(\frac{\exp\left(\boldsymbol{u}_{w \neq o}^\top \boldsymbol{v_c}\right)}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v_c}\right)}\right)\boldsymbol{v_c}$$
$$= \hat{y}_{w \neq o}\boldsymbol{v_c}$$

When $w = o$,

$$\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u}_{w=o}} = \frac{\partial}{\partial \boldsymbol{u}_{w=o}} \left[ -\boldsymbol{u}_o^\top \boldsymbol{v}_c + \log \left( \sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right) \right) \right]$$

$$= -\boldsymbol{v}_c + \frac{\exp\left(\boldsymbol{u}_{w=o}^\top \boldsymbol{v}_c\right)\boldsymbol{v}_c}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)}$$

$$= -\boldsymbol{v}_c + \left( \frac{\exp\left(\boldsymbol{u}_{w=o}^\top \boldsymbol{v}_c\right)}{\sum_{k \in |\text{Vocab}|} \exp\left(\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)} \right) \boldsymbol{v}_c$$

$$= -\boldsymbol{v}_c + \hat{\boldsymbol{y}}_{w=o}\boldsymbol{v}_c$$

$$= (\hat{\boldsymbol{y}}_{w=o} - 1)\boldsymbol{v}_c$$

Therefore,

$$\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u}_w} = \begin{cases} \hat{\boldsymbol{y}}_w \boldsymbol{v}_c & \text{if } w \neq o \\ (\hat{\boldsymbol{y}}_w - 1)\boldsymbol{v}_c & \text{if } w = o \end{cases}$$

# Question (d)

$$\frac{\partial J_{naive-softmax}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{U}} = \begin{pmatrix} \frac{\partial J(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_1}} & \frac{\partial J(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_2}} & \cdots & \frac{\partial J(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u}_{|\text{Vocab}|}} \end{pmatrix}$$

# Question (e)

$$f(x) = \max(\alpha x, x)$$

$$= \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x > 0 \end{cases}$$

$$\frac{df(x)}{dx} = \begin{cases} \alpha & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

# Question (f)

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

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

$$= \sigma(x) - \sigma(x)^2$$

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

# Question (g)

**Part (i)**

$$\frac{d}{dx} \log \sigma(x) = \frac{1}{\sigma(x)} \left( \frac{d}{dx} \sigma(x) \right)$$

$$= \frac{1}{\sigma(x)} \left( \sigma(x) \left( 1 - \sigma(x) \right) \right)$$

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

$$J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U}) = -\log(\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c})) - \sum_{s=1}^{K} \log(\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}))$$

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{v_c}} = (\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c}) - 1) \frac{\partial}{\partial \boldsymbol{v_c}} (\boldsymbol{u}_o^\top \boldsymbol{v_c}) + \sum_{s=1}^{K} (\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}) - 1) \frac{\partial}{\partial \boldsymbol{v_c}} (-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c})$$

$$= (\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c}) - 1) \boldsymbol{u}_o - \sum_{s=1}^{K} (\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}) - 1) \boldsymbol{u}_{w_s}$$

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_o}} = (\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c}) - 1) \frac{\partial}{\partial \boldsymbol{u_o}} (\boldsymbol{u}_o^\top \boldsymbol{v_c}) + 0$$

$$= (\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c}) - 1) \boldsymbol{v}_c$$

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_{w_s}}} = 0 + (\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}) - 1) \frac{\partial}{\partial \boldsymbol{u_{w_s}}} (-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c})$$

$$= -(\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}) - 1) \boldsymbol{v}_c$$

**Part (ii)**

$$\boldsymbol{U}_{o,\{w_1,\ldots,w_K\}} = [\boldsymbol{u_o}, -\boldsymbol{u_{w_1}}, \cdots, -\boldsymbol{u_{w_K}}]$$

Denote

$$\boldsymbol{\lambda} = \sigma \left( \boldsymbol{U}_{o,\{w_1,\cdots,w_K\}}^\top \boldsymbol{v}_c \right) - \mathbf{1}$$

$$= \begin{pmatrix} \sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c}) - 1 \\ \sigma(-\boldsymbol{u}_{w_1}^\top \boldsymbol{v_c}) - 1 \\ \vdots \\ \sigma(-\boldsymbol{u}_{w_K}^\top \boldsymbol{v_c}) - 1 \end{pmatrix}$$

$$= \begin{pmatrix} \lambda_o \\ \lambda_1 \\ \vdots \\ \lambda_K \end{pmatrix}$$

Each element of $\boldsymbol{\lambda}$ is used in computing the partial derivatives in part (g)(i). The derivatives are simplified as follows:

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{v_c}} = \lambda_o \boldsymbol{u}_o - \sum_{s=1}^{K} \lambda_s \boldsymbol{u}_{w_s}$$

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_o}} = \lambda_o \boldsymbol{v}_c$$

$$\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u_{w_s}}} = -\lambda_s \boldsymbol{v}_c$$

**Part (iii)**

The computation of $J_{\text{naive-softmax}}$ requires computing the dot products between all word vectors with the input vector. The time complexity of the computation is $O(|\text{Vocab}|)$. On the other hand, the computation of $J_{\text{negative-sample}}$ involves only the dot products between K negative sample word vectors and the input vectors, where the time complexity of computation is $O(K)$. Hence, computing the negative sampling loss function is much more efficient than the naive-softmax function.
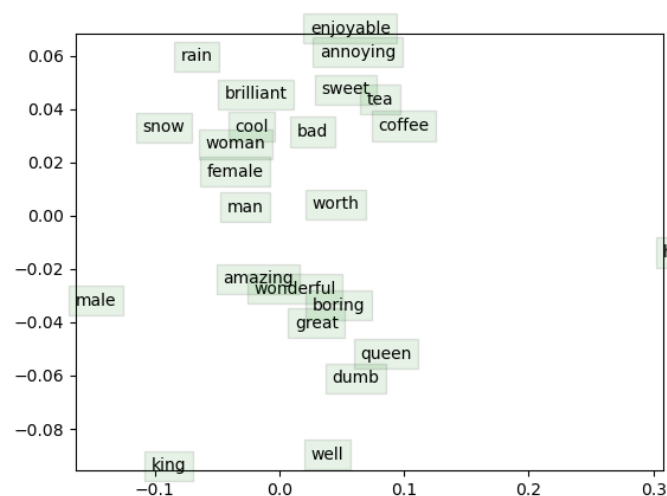
# Question (h)

$$
\begin{aligned}
\frac{\partial J_{\text{negative-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U})}{\partial \boldsymbol{u}_{\boldsymbol{w}_s}} &= \frac{\partial}{\partial \boldsymbol{u}_{\boldsymbol{w}_s}} \left( -\log(\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c})) - \sum_{k=1}^{K} \log(\sigma(-\boldsymbol{u}_{w_k}^\top \boldsymbol{v_c})) \right) \\
&= \frac{\partial}{\partial \boldsymbol{u}_{\boldsymbol{w}_s}} \left( -\log(\sigma(\boldsymbol{u}_o^\top \boldsymbol{v_c})) - \sum_{\substack{k=1 \\ w_k \neq w_s}}^{K} \log(\sigma(-\boldsymbol{u}_{w_k}^\top \boldsymbol{v_c})) - \sum_{\substack{k=1 \\ w_k = w_s}}^{K} \log(\sigma(-\boldsymbol{u}_{w_k}^\top \boldsymbol{v_c})) \right) \\
&= 0 - 0 - \frac{\partial}{\partial \boldsymbol{u}_{\boldsymbol{w}_s}} \left( \sum_{\substack{k=1 \\ w_k = w_s}}^{K} \log(\sigma(-\boldsymbol{u}_{w_k}^\top \boldsymbol{v_c})) \right) \\
&= - \sum_{\substack{k=1 \\ w_k = w_s}}^{K} (\sigma(-\boldsymbol{u}_{w_s}^\top \boldsymbol{v_c}) - 1)\boldsymbol{v}_c
\end{aligned}
$$

# Question (i)

$$
\begin{aligned}
\frac{\partial J_{\text{skip-gram}}(\boldsymbol{v}_c, w_{t-m}, \cdots, w_{t+m}, \boldsymbol{U})}{\partial \boldsymbol{U}} &= \sum_{\substack{-m \leq j < m \\ j \neq 0}} \frac{\partial J(\boldsymbol{v}_c, w_{t+j}, \boldsymbol{U})}{\partial \boldsymbol{U}} \\
\frac{\partial J_{\text{skip-gram}}(\boldsymbol{v}_c, w_{t-m}, \cdots, w_{t+m}, \boldsymbol{U})}{\partial \boldsymbol{v_c}} &= \sum_{\substack{-m \leq j < m \\ j \neq 0}} \frac{\partial J(\boldsymbol{v}_c, w_{t+j}, \boldsymbol{U})}{\partial \boldsymbol{v_c}} \\
\frac{\partial J_{\text{skip-gram}}(\boldsymbol{v}_c, w_{t-m}, \cdots, w_{t+m}, \boldsymbol{U})}{\partial \boldsymbol{v_w}} &= 0 \qquad \text{where} \quad w \neq c
\end{aligned}
$$

# Coding: Implementing word2vec Part (c)



In the plot, some synonyms ("amazing", "wonderful", "great") and antonyms ("enjoyable", "annoying") are cluster together. Terms with similar usage such as ("tea", "coffee"), ("rain", "snow") and ("female", "woman") also group together. Lastly, the linear relationships of word vectors, i.e. vec("king") − vec("male") + vec("female") ≈ vec("queen"), is observed.

## Coding: Implementing word2vec Part (c)