

(a) Because y is one-hot vector $\rightarrow y_w = \begin{cases} 0 & w \neq 0 \\ 1 & w = 0 \end{cases}$

So we can remove all $y_w = 0$.

$$(b) i) \frac{\partial J(r_c, o, U)}{\partial r_c} = - \frac{\partial}{\partial r_c} \log p(o=o | C=c)$$

$$= - \frac{\partial}{\partial r_c} \log \frac{\exp(u_o^T r_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} = - \frac{\partial}{\partial r_c} \log \exp(u_o^T r_c)$$

$$+ \frac{\partial}{\partial r_c} \log \sum_{w \in \text{vocab}} \exp(u_w^T r_c) = - \frac{\partial}{\partial r_c} u_o^T r_c +$$

$$\frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} \frac{\partial}{\partial r_c} \sum_{w \in \text{vocab}} \exp(u_w^T r_c) = -u_o +$$

$$\frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} \sum_{w \in \text{vocab}} \frac{\partial}{\partial r_c} \exp(u_w^T r_c) = -u_o +$$

$$+ \frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} \sum_{w \in \text{vocab}} \exp(u_w^T r_c) \frac{\partial}{\partial r_c} u_w^T r_c$$

$$= -u_o + \frac{\sum_{w \in \text{vocab}} \exp(u_w^T r_c) u_w}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} = -Uy + \sum_{w \in \text{vocab}} \frac{\exp(u_w^T r_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)} u_w$$

$$= -Uy + \sum_{w \in \text{vocab}} \hat{y}_w u_w = -Uy + U\hat{y} = U(\hat{y} - y)$$

$$\text{w/ } r_c, u_w \in \mathbb{R}^N \quad U \in \mathbb{R}^{N \times V} \quad y, \hat{y} \in \mathbb{R}^{V \times 1} \quad \frac{\partial J}{\partial r_c} \in \mathbb{R}^N$$

$$\text{ii) } U(y - \hat{y}) = 0 \rightarrow \underbrace{Uy}_{u_0} = U\hat{y}$$

The gradient is equal to zero when predicted word is equal to outside word.

iii) We subtract this difference from v_0 to get closer to optimum point. It means that we want to have minimum difference between y and \hat{y} .

iv) In some cases, L2 normalization might take away useful information that could be relevant for the task.

For example, in sentiment analysis where the goal is to classify phrases as positive or negative, it may not be good, because L2 scales the embeddings to have a unit norm & remove information about the original magnitude.

Well in tasks such as text similarity or clustering it can be helpful in reducing the influence of outliers.

$$(c) -\frac{\partial}{\partial u_w} \log P(O=o | C=c) = -\frac{\partial}{\partial u_w} \log \frac{\exp(u_o^T r_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T r_c)}$$

$$= -\frac{\partial}{\partial u_w} \log \exp(u_o^T r_c) + \frac{\partial}{\partial u_w} \log \sum_{w \in \text{vocab}} \exp(u_w^T r_c)$$

$$= -\frac{\partial}{\partial u_w} u_o r_c + \frac{1}{\sum_w \exp(u_w r_c)} \cdot \frac{\partial}{\partial u_w} \sum_w \exp(u_w r_c)$$

$$= -\frac{\partial}{\partial u_w} u_o r_c + \frac{1}{\sum_w \exp(u_w r_c)} \cdot \frac{\partial}{\partial u_w} \exp(u_w r_c)$$

$$= -\frac{\partial}{\partial u_w} u_o r_c + \frac{\exp(u_w r_c)}{\sum_w \exp(u_w r_c)} \cdot r_c = -y_w r_c + \hat{y}_w r_c$$

$$= (\hat{y}_w - y_w) r_c^T \quad y_w = \begin{cases} 1 & w=0 \\ 0 & o, w \end{cases}$$

$$(d) \frac{\partial \mathcal{L}}{\partial \mathbf{u}} = \left[\frac{\partial \mathcal{L}}{\partial u_1}, \frac{\partial \mathcal{L}}{\partial u_2}, \dots, \frac{\partial \mathcal{L}}{\partial u_{|\text{vocab}|}} \right]$$

$$(e) f(x) = \begin{cases} x & x \geq 0 \\ \alpha x & x < 0 \end{cases} \quad f'(x) = \begin{cases} 1 & x \geq 0 \\ \alpha & x < 0 \end{cases}$$

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

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

$$\begin{aligned}
 \textcircled{9} \text{ i) } \frac{\partial J}{\partial v_c} &= \frac{\partial}{\partial v_c} \log(\sigma(u_0^T v_c)) - \frac{\partial}{\partial v_c} \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \\
 &= -\frac{1}{\sigma(u_0^T v_c)} \frac{\partial}{\partial v_c} \sigma(u_0^T v_c) - \sum_{k=1}^K \frac{\partial}{\partial v_c} \log(\sigma(-u_k^T v_c)) = -\frac{1}{\sigma(u_0^T v_c)} \\
 &\quad \times \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c)) u_0 - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \frac{\partial}{\partial v_c} \sigma(-u_k^T v_c) \\
 &= (\sigma(u_0^T v_c) - 1) u_0 - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c))
 \end{aligned}$$

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

$$\begin{aligned}
 * \frac{\partial J}{\partial u_0} &= -\frac{\partial}{\partial u_0} \log(\sigma(u_0^T v_c)) - \frac{\partial}{\partial u_0} \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \\
 &= -\frac{1}{\sigma(u_0^T v_c)} \frac{\partial}{\partial u_0} \sigma(u_0^T v_c) = (\sigma(u_0^T v_c) - 1) v_c
 \end{aligned}$$

$$\begin{aligned}
 * \frac{\partial J}{\partial u_k} &= -\frac{\partial}{\partial u_k} \log(\sigma(u_0^T v_c)) - \frac{\partial}{\partial u_k} \sum_{x=1}^K \log(\sigma(-u_x^T v_c)) \\
 &= -\frac{\partial}{\partial u_k} \log(\sigma(-u_k^T v_c)) = (1 - \sigma(-u_k^T v_c)) v_c
 \end{aligned}$$

ii)

$$u_{0, \{w_1, \dots, w_k\}} = [u_0, -u_{w_1}, \dots, -u_{w_k}]$$

$$u_0^T v_c = \begin{bmatrix} u_0^T v_c \\ \vdots \\ -u_{w_k}^T v_c \end{bmatrix} \quad 1 - \sigma(u_0^T v_c) = \begin{bmatrix} 1 - \sigma(u_0^T v_c) \\ \vdots \\ 1 - \sigma(u_{w_k}^T v_c) \end{bmatrix}$$

iii) Rather than computing the probabilities for all possible output words, negative sampling only considers a small number of negative samples per positive sample, greatly reducing the number of computations.

$$(h) \frac{\partial J}{\partial u_k} = -\frac{\partial}{\partial u_k} \log(\sigma(u_0^T v_c)) - \frac{\partial}{\partial u_k} \sum_{n=1}^k \log(\sigma(-u_n^T v_c))$$

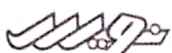
$$= -\frac{\partial}{\partial u_k} \sum_{\substack{u_n = u_k \\ n=1, \dots, k}} \log(\sigma(-u_n^T v_c)) - \frac{\partial}{\partial u_k} \sum_{\substack{u_n \neq u_k \\ n=1, \dots, k}} \log(\sigma(-u_n^T v_c))$$

$$= -\sum_{\substack{u_n = u_k \\ n=1, \dots, k}} \frac{\partial}{\partial u_k} \log(\sigma(-u_k^T v_c)) = -\left(\sum_{n=1}^k \mathbb{I}\{u_n = u_k\} \right)$$

$$\frac{\partial}{\partial u_k} \log(\sigma(-u_k^T v_c)) = \left(\sum_{n=1}^k \mathbb{I}\{u_n = u_k\} \right) (1 - \sigma(-u_k^T v_c)) v_c$$

$$(i) \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, u)}{\partial u} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, u)}{\partial u}$$

$$ii) \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, u)}{\partial v_c} = \sum_{-m \leq j \leq m} \frac{\partial J(v_c, w_{t+j}, u)}{\partial v_c}$$



$$\text{iii)} \quad \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_w} = 0 \quad w \neq c$$