

## 1. Understanding word2vec

(a) To prove:

$$-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o)$$

$y$  is an one-hot vector with 1 in context vector position and 0 in all other positions.

$$y_w = 1 \quad \text{if } w \text{ is context word} \\ = 0 \quad \text{otherwise.}$$

$$\begin{aligned} \therefore -\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) &= -[y_1 \log(\hat{y}_1) + \dots + y_o \log(\hat{y}_o) + \dots + y_v \log(\hat{y}_v)] \\ &= -[0 + \dots + 1 \cdot \log(\hat{y}_o) + \dots + 0] \\ &= -\log \hat{y}_o \end{aligned}$$

(b) To find:  $\frac{\partial J_{\text{ns}}(v_c, o, u)}{\partial v_c}$  ?

$$\begin{aligned} J_{\text{ns}}(v_c, o, u) &= -\log P(O=o|C=c) \\ &= -\frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \end{aligned}$$

$$\begin{aligned} &= -[\log \exp(u_o^T v_c) - \log \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)] \\ &= -[u_o^T v_c - \log \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)] \end{aligned}$$

$$\frac{\partial J(v_c, o, u)}{\partial v_c} = \frac{\partial}{\partial v_c} [\log \sum_{w \in \text{Vocab}} \exp(u_w^T v_c) - u_o^T v_c] = \frac{\partial}{\partial v_c} [\log \sum_w \exp(u_w^T v_c)] - \frac{\partial}{\partial v_c} [u_o^T v_c] \quad \text{①} \quad \text{②}$$

$$\text{②} \rightarrow \frac{\partial}{\partial v_c} [u_o^T v_c] = u_o.$$

$$\begin{aligned} \text{①} \rightarrow \frac{\partial}{\partial v_c} [\log \sum_w \exp(u_w^T v_c)] &= \frac{1}{\sum_w \exp(u_w^T v_c)} \sum_{x \in \text{Vocab}} \frac{\partial}{\partial v_c} \exp(u_x^T v_c) \\ &= \frac{1}{\sum_w \exp(u_w^T v_c)} \sum_{x \in \text{Vocab}} \exp(u_x^T v_c) \frac{\partial}{\partial v_c} (u_x^T v_c) \\ &= \frac{1}{\sum_w \exp(u_w^T v_c)} \sum_{x \in \text{Vocab}} \exp(u_x^T v_c) u_x \end{aligned}$$

$$= \frac{\sum_{x \in Y_{\text{ocab}}} \exp(Ux^T V_c) Ux}{\sum_{w \in Y_{\text{ocab}}} \exp(Uw^T V_c)} = \sum_{x \in Y_{\text{ocab}}} P(x|c) Ux = \sum_{x \in Y_{\text{ocab}}} \hat{y}_x \cdot Ux$$

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial V_c} = \textcircled{1} - \textcircled{2} = \cancel{\psi_d} / \cancel{\sum_{x \in Y_{\text{ocab}}} \hat{y}_x} / \cancel{U/c}$$

$$= \sum_{x \in Y_{\text{ocab}}} \hat{y}_x \cdot Ux - U_0 = \sum_{x \in Y_{\text{ocab}}} \hat{y}_x \cdot Ux - \sum_{w \in Y_{\text{ocab}}} y_w \cdot U_w$$

③  $\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_i}$  ? when  $i=0$  and  $i \neq 0$ .

When  $i=0$

$$\begin{aligned} \frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_0} &= \frac{\partial}{\partial U_0} \left[ \log \sum_{w \in Y_{\text{ocab}}} \exp(Uw^T V_c) - U_0^T V_c \right] \\ &= \frac{\partial}{\partial U_0} \left[ \log \sum_{w \in Y_{\text{ocab}}} \exp(Uw^T V_c) \right] - \frac{\partial}{\partial U_0} [U_0^T V_c] \end{aligned}$$

①                      ②

②  $\rightarrow \frac{\partial}{\partial U_0} [U_0^T V_c] = V_c$

$$\begin{aligned} \textcircled{1} \rightarrow \frac{\partial}{\partial U_0} \left[ \log \sum_{w \in Y_{\text{ocab}}} \exp(Uw^T V_c) \right] &= \frac{1}{\sum_w \exp(Uw^T V_c)} \cdot \sum_{x \in Y_{\text{ocab}}} \frac{\partial}{\partial U_0} \exp(Ux^T V_c) \\ &= \frac{1}{\sum_w \exp(Uw^T V_c)} \cdot \left[ \frac{\partial}{\partial U_0} \exp(U_1^T V_c) + \dots + \frac{\partial}{\partial U_0} \exp(U_0^T V_c) + \dots + \frac{\partial}{\partial U_0} \exp(U_{|Y|}^T V_c) \right] \\ &= \frac{1}{\sum_w \exp(Uw^T V_c)} \cdot \exp(U_0^T V_c) \cdot \frac{\partial}{\partial U_0} (U_0^T V_c) \\ &= \frac{\exp(U_0^T V_c)}{\sum_w \exp(Uw^T V_c)} \cdot V_c = P(0|c) \cdot V_c \end{aligned}$$

$$\frac{\partial}{\partial U_0} J_{\text{ns}}(V_c, 0, U) = \textcircled{1} - \textcircled{2} = P(0|c) \cdot V_c - V_c = \sum_{x \in Y_{\text{ocab}}} y_x \cdot \hat{y}_x \cdot V_c - V_c$$

When  $i \neq 0$ .

$$\begin{aligned} \frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_i} &= \frac{\partial}{\partial U_i} \left[ \log \sum_{w \in \text{vocab}} \exp(Uw^T V_c) - U_0^T V_c \right] \\ &= \frac{\partial}{\partial U_i} \left[ \log \sum_{w \in \text{vocab}} \exp(Uw^T V_c) \right] - \frac{\partial}{\partial U_i} (U_0^T V_c) \quad (2) \end{aligned}$$

$$(2) \rightarrow \frac{\partial}{\partial U_i} (U_0^T V_c) = 0$$

$$\begin{aligned} (1) \rightarrow \frac{\partial}{\partial U_i} \left[ \log \sum_{w \in \text{vocab}} \exp(Uw^T V_c) \right] &= \frac{1}{\sum_{w \in \text{vocab}} \exp(Uw^T V_c)} \cdot \sum_{x \in \text{vocab}} \frac{\partial}{\partial U_i} \exp(Ux^T V_c) \\ &= \frac{1}{\sum_{w \in \text{vocab}} \exp(Uw^T V_c)} \cdot \left[ \frac{\partial}{\partial U_i} \exp(U_1^T V_c) + \dots + \frac{\partial}{\partial U_i} \exp(U_i^T V_c) + \dots + \frac{\partial}{\partial U_i} \exp(U_v^T V_c) \right] \\ &= \frac{1}{\sum_{w \in \text{vocab}} \exp(Uw^T V_c)} \left[ \frac{\partial}{\partial U_i} \exp(U_i^T V_c) \right] = \frac{1}{\sum_{w \in \text{vocab}} \exp(Uw^T V_c)} \cdot \exp(U_i^T V_c) \cdot \frac{\partial}{\partial U_i} (U_i^T V_c) \\ &= \frac{\exp(U_i^T V_c)}{\sum_{w \in \text{vocab}} \exp(Uw^T V_c)} \cdot V_c = P(i|c) \cdot V_c \end{aligned}$$

$$\begin{aligned} \frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_i} &= (1) - (2) = P(i|c) \cdot V_c - 0 \\ &= \hat{y}_i \cdot V_c \end{aligned}$$

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

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

$$\begin{aligned} &= \frac{e^{-x}}{1+e^{-x}} \times \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-x}} \left[ 1 - \frac{1}{1+e^{-x}} \right] = \sigma(x) [1 - \sigma(x)] \\ \sigma'(x) &= \sigma(x) [1 - \sigma(x)] \end{aligned}$$



$$c) J_{\text{neg-sample}}(V_c, 0, U) = -\log(\sigma(U_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-U_k^T V_c))$$

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial V_c} = \frac{\partial}{\partial V_c} \left[ \underset{\textcircled{1}}{-\log(\sigma(U_0^T V_c))} - \frac{\partial}{\partial V_c} \left[ \sum_{k=1}^K \underset{\textcircled{2}}{\log(\sigma(-U_k^T V_c))} \right] \right]$$

$$\textcircled{1} \rightarrow = -\frac{1}{\sigma(U_0^T V_c)} \sigma(U_0^T V_c)(1-\sigma(U_0^T V_c)) \frac{\partial}{\partial V_c} (U_0^T V_c)$$

$$= -\frac{1}{\sigma(U_0^T V_c)} \sigma(U_0^T V_c)(1-\sigma(U_0^T V_c)) U_0$$

$$\textcircled{2} \rightarrow \sum_{k=1}^K \frac{\partial}{\partial V_c} \log(\sigma(-U_k^T V_c)) = \sum_{k=1}^K \frac{1}{\sigma(-U_k^T V_c)} \sigma(-U_k^T V_c)(1-\sigma(-U_k^T V_c)) \frac{\partial}{\partial V_c} (-U_k^T V_c)$$

$$= \sum_{k=1}^K \frac{1}{\sigma(-U_k^T V_c)} \sigma(-U_k^T V_c)(1-\sigma(-U_k^T V_c)) \cdot (-U_k)$$

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial V_c} = \textcircled{1} - \textcircled{2}$$

$$= -\frac{1}{\sigma(U_0^T V_c)} \sigma(U_0^T V_c)(1-\sigma(U_0^T V_c)) U_0 + \sum_{k=1}^K \frac{1}{\sigma(-U_k^T V_c)} \sigma(-U_k^T V_c)(1-\sigma(-U_k^T V_c)) U_k$$

Similarly,

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_0} = -\frac{1}{\sigma(U_0^T V_c)} \sigma(U_0^T V_c)(1-\sigma(U_0^T V_c)) \cdot V_c + \frac{1}{\sigma(-U_0^T V_c)} \sigma(-U_0^T V_c)(1-\sigma(-U_0^T V_c)) \cdot V_c$$

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U_k} = 0 + \frac{1}{\sigma(-U_k^T V_c)} \sigma(-U_k^T V_c)(1-\sigma(-U_k^T V_c)) \cdot V_c$$

$$\frac{\partial J_{\text{ns}}}{\partial V_c} = -(1-\sigma(U_0^T V_c)) U_0 + \sum_{k=1}^K (1-\sigma(-U_k^T V_c)) U_k$$

$$\frac{\partial J_{\text{ns}}}{\partial U_0} = -(1-\sigma(U_0^T V_c)) \cdot V_c + \frac{1}{\sigma(-U_0^T V_c)} \sigma(-U_0^T V_c)(1-\sigma(-U_0^T V_c)) \cdot V_c$$

$$\frac{\partial J_{\text{ns}}}{\partial U_k} = (1-\sigma(-U_k^T V_c)) \cdot V_c$$

This loss fn. is more efficient than naive softmax because in naive softmax we have to compute  $(\sum_{\text{vectors}} U_i^T V_c)$  everytime which is computationally expensive, whereas in negative sampling only the sampled 'k' vectors are used thus resulting in faster computation.

$$f) J_{\text{LSS-gram}}(V_c, w_{t-m}, \dots, w_{t+m}, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(V_c, w_{t+j}, U)$$

$$i) \frac{\partial J_{\text{LSS}}(V_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial U} J(V_c, w_{t+j}, U)$$

$$ii) \frac{\partial J_{\text{LSS}}(V_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial V_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial}{\partial V_c} J(V_c, w_{t+j}, U)$$

$$iii) \frac{\partial J_{\text{LSS}}(V_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial V_w} \text{ when } w \neq c = 0 \quad \left( \text{as the loss function depends only on } V_c \text{ and } U_1, \dots, U_K \right)$$