

CS 224N Assignment #2: Word2Vec

1. Understanding Word2Vec

- (a) (3 points) Show that the naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between y and \hat{y} ; i.e., show that

$$-\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -\log(\hat{y}_o). \quad (3)$$

$$-\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -\{y_1 \log(\hat{y}_1) + \dots + y_o \log(\hat{y}_o) + \dots + y_k \log(\hat{y}_k)\} = -y_o \log(\hat{y}_o)$$

$y=0$ 이 되는 주변 단어 값들은 다 지워지고, $y=1$ 인 값을 갖는 중심 단어만 남게 되므로 식은 생략한다.

- (b) (5 points) Compute the partial derivative of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to v_c . Please write your answer in terms of y , \hat{y} , and U .

$$\begin{aligned} J_{\text{naive-softmax}}(v_c, o, U) &= -\log P(O=o | C=c) \quad \text{center word가 주어졌을 때, outside word가 나올 확률} \\ &= -\log\left(\frac{\exp(u_o^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)}\right) \\ &= -u_o^T v_c + \log\left(\sum_w \exp(u_w^T v_c)\right) \end{aligned}$$

$$\frac{\partial J}{\partial v_c} = -u_o^T + \sum_w \frac{\exp(u_w^T v_c) \cdot u_w^T}{\sum_w \exp(u_w^T v_c)}$$

$$= -u_o + \sum_w P(O=w | C=c) \cdot u_w$$

$$= -u_o + \sum_w \hat{y}_w u_w \quad \because \text{answer (a): 답에 대해서만 1이고, 나머지 값은 모두 0}$$

$$= -U(\hat{y} - y)$$

- (c) (5 points) Compute the partial derivatives of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to each of the 'outside' word vectors, \mathbf{u}_w 's. There will be two cases: when $w = o$, the true 'outside' word vector, and $w \neq o$, for all other words. Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c .

$$J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\mathbf{u}_o^T \mathbf{v}_c + \log \left(\sum_{\mathbf{w}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \right)$$

i) $w = o$ outside word 가 실제 정답인 경우 (correct class)

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_w} &= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{\mathbf{w}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \cdot \mathbf{v}_c \\ &= -\mathbf{v}_c + P(o=w | c=c) \cdot \mathbf{v}_c \\ &= (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{v}_c \end{aligned}$$

ii) $w \neq o$ outside word 가 실제 정답이 아닌 경우 (incorrect class)

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_w} &= 0 + \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{\mathbf{w}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \cdot \mathbf{v}_c \\ &= P(o=w | c=c) \cdot \mathbf{v}_c \\ &= \hat{\mathbf{y}} \mathbf{v}_c \end{aligned}$$

- (d) (3 Points) The sigmoid function is given by Equation 4:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (4)$$

Please compute the derivative of $\sigma(x)$ with respect to x , where x is a scalar. Hint: you may want to write your answer in terms of $\sigma(x)$.

$$\begin{aligned} \frac{\partial \sigma}{\partial x} &= \frac{-(-e^{-x})}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \times \frac{e^{-x}}{1+e^{-x}} \\ &= \sigma(x) \{1 - \sigma(x)\} \end{aligned}$$

- (e) (4 points) Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as w_1, w_2, \dots, w_K and their outside vectors as $\mathbf{u}_1, \dots, \mathbf{u}_K$. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, U) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (5)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.³

Please repeat parts (b) and (c), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to \mathbf{v}_c , with respect to \mathbf{u}_o , and with respect to a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{u}_o , \mathbf{v}_c , and \mathbf{u}_k , where $k \in [1, K]$. After you've done this, describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss. Note, you should be able to use your solution to part (d) to help compute the necessary gradients here.

$$\begin{aligned} \text{i) } \frac{\partial J}{\partial \mathbf{v}_c} &= -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \cdot \sigma(\mathbf{u}_o^\top \mathbf{v}_c) \cdot [1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \cdot \mathbf{u}_o^\top \\ &\quad - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \cdot \sigma(-\mathbf{u}_k^\top \mathbf{v}_c) \cdot [1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c)] \cdot (-\mathbf{u}_k)^\top \\ &= -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o + \sum_{k=1}^K [1 - \sigma(\mathbf{u}_k^\top \mathbf{v}_c)] \mathbf{u}_k \end{aligned}$$

$$\text{ii) } \frac{\partial}{\partial \mathbf{u}_o} \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) = 0 \quad (\because o \notin \{w_1, \dots, w_K\})$$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_o} &= -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \cdot \sigma(\mathbf{u}_o^\top \mathbf{v}_c) \cdot [1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \cdot \mathbf{v}_c \\ &= -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \cdot \mathbf{v}_c \end{aligned}$$

$$\text{iii) } \frac{\partial}{\partial \mathbf{u}_k} \log(\sigma(\mathbf{u}_k^\top \mathbf{v}_c)) = 0 \quad (\because o \notin \{w_1, \dots, w_K\})$$

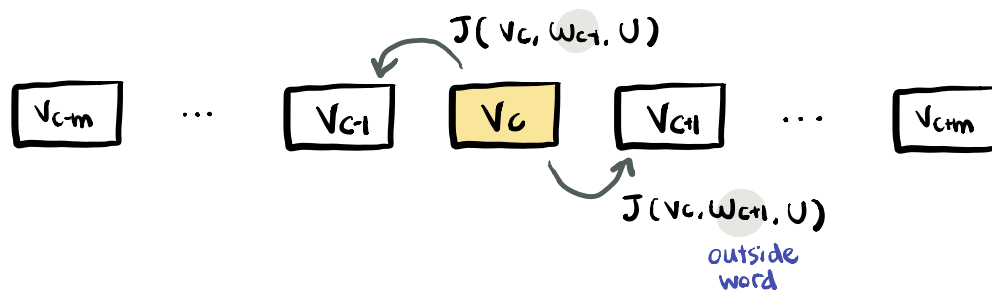
$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_k} &= -\sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \cdot \sigma(-\mathbf{u}_k^\top \mathbf{v}_c) \cdot [1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c)] \cdot (-\mathbf{v}_c) \\ &= \sum_{k=1}^K [1 - \sigma(\mathbf{u}_k^\top \mathbf{v}_c)] \cdot \mathbf{v}_c \end{aligned}$$

- (f) (3 points) Suppose the center word is $c = w_t$ and the context window is $[w_{t-m}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+m}]$, where m is the context window size. Recall that for the skip-gram version of word2vec, the total loss for the context window is:

$$J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (6)$$

Here, $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ represents an arbitrary loss term for the center word $c = w_t$ and outside word w_{t+j} . $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ could be $J_{\text{naive-softmax}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ or $J_{\text{neg-sample}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$, depending on your implementation.

Write down three partial derivatives:



$$(i) \frac{\partial J}{\partial \mathbf{U}} = \sum \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

$$(ii) \frac{\partial J}{\partial \mathbf{v}_c} = \sum \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

$$(iii) \frac{\partial J}{\partial w_c} = 0 \quad \text{center word 가 아닌 경우 update 되지 않기 때문}$$