

CS 224n Assignment #2: word2vec (43 Points)

1 Written: Understanding word2vec (23 points)

Let's have a quick refresher on the word2vec algorithm. The key insight behind word2vec is that '*a word is known by the company it keeps*'. Concretely, suppose we have a 'center' word c and a contextual window surrounding c . We shall refer to words that lie in this contextual window as 'outside words'. For example, in Figure 1 we see that the center word c is 'banking'. Since the context window size is 2, the outside words are 'turning', 'into', 'crises', and 'as'.

The goal of the skip-gram word2vec algorithm is to accurately learn the probability distribution $P(O|C)$. Given a specific word o and a specific word c , we want to calculate $P(O = o|C = c)$, which is the probability that word o is an 'outside' word for c , i.e., the probability that o falls within the contextual window of c .

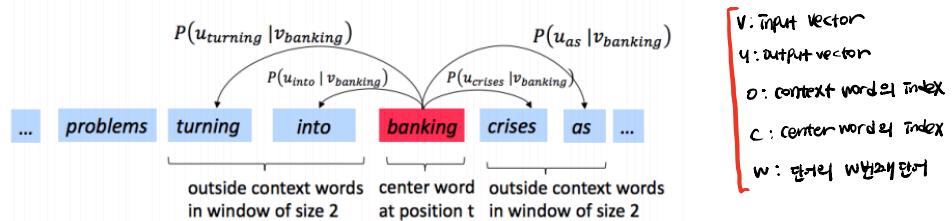


Figure 1: The word2vec skip-gram prediction model with window size 2

In word2vec, the conditional probability distribution is given by taking vector dot-products and applying the softmax function:

$$P(O = o | C = c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (1)$$

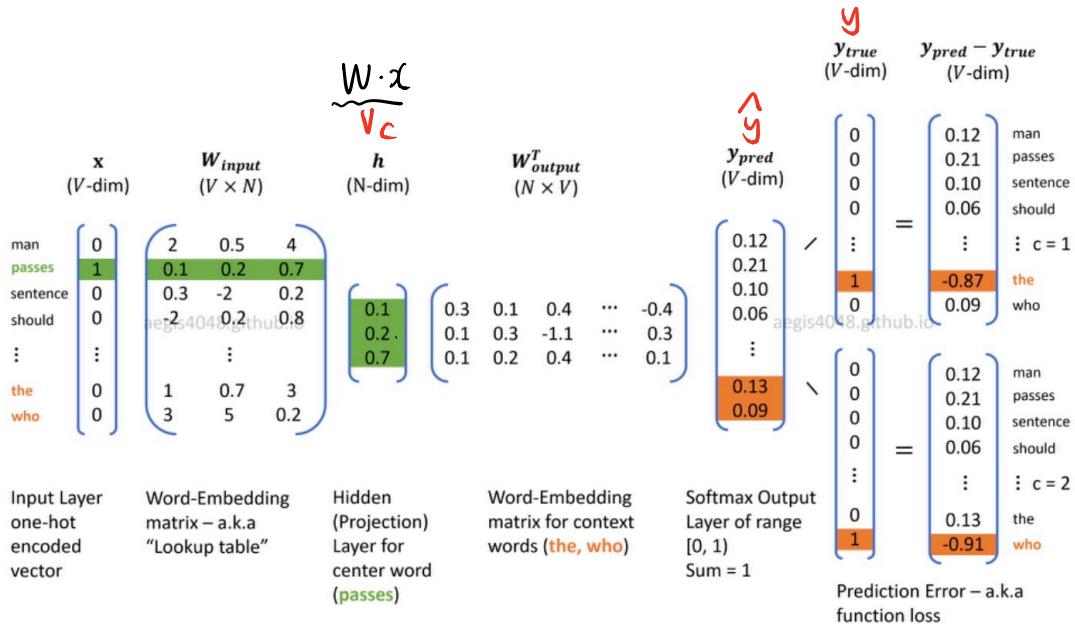
Here, \mathbf{u}_o is the 'outside' vector representing outside word o , and \mathbf{v}_c is the 'center' vector representing center word c . To contain these parameters, we have two matrices, \mathbf{U} and \mathbf{V} . The columns of \mathbf{U} are all the 'outside' vectors \mathbf{u}_w . The columns of \mathbf{V} are all of the 'center' vectors \mathbf{v}_w . Both \mathbf{U} and \mathbf{V} contain a vector for every $w \in \text{Vocabulary}$.¹

Recall from lectures that, for a single pair of words c and o , the loss is given by:

$$J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c). \quad (2)$$

Another way to view this loss is as the cross-entropy² between the true distribution \mathbf{y} and the predicted distribution $\hat{\mathbf{y}}$. Here, both \mathbf{y} and $\hat{\mathbf{y}}$ are vectors with length equal to the number of words in the vocabulary. Furthermore, the k^{th} entry in these vectors indicates the conditional probability of the k^{th} word being an 'outside word' for the given c . The true empirical distribution \mathbf{y} is a one-hot vector with a 1 for the true outside word o , and 0 everywhere else. The predicted distribution $\hat{\mathbf{y}}$ is the probability distribution $P(O|C = c)$ given by our model in equation (1).

(Onehot vector)
단어의 outside Word의 1이 2, 나머지는 0



- (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

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

Your answer should be one line.

y : 실제 확률분포, \hat{y} : 모델에서 구한 확률분포 (둘다 단어개수와 동일한 길이)

v : input vector
 y : output vector
 o : outside word의 index
 c : center word의 index
 w : 단어의 word index

$$\hat{y}_o = p(u_o | v_c) = \frac{e^{u_o^T v_c}}{\sum_{w \in Vocab} e^{u_w^T v_c}}$$

y 는 one-hot vector로 실제 outside vector인 0에 대한만 1이고, 나머지는 0

$$-\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -(\underbrace{y_1 \log(\hat{y}_1)}_{\downarrow} + \underbrace{y_2 \log(\hat{y}_2)}_{\downarrow} + \dots + \underbrace{y_o \log(\hat{y}_o)}_{\downarrow} + \dots + \underbrace{y_{|V|} \log(\hat{y}_{|V|})}_{\downarrow})$$

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

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

$$\begin{aligned}
 \frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{v}_c} &= -\frac{\partial}{\partial \mathbf{v}_c} [-\log(\hat{y}_o)] \\
 &= -\frac{\partial}{\partial \mathbf{v}_c} \left[-\log \left(\frac{e^{\mathbf{u}_o^\top \mathbf{v}_c}}{\sum_{w \in \text{Vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \right) \right] \\
 &= -\frac{\partial}{\partial \mathbf{v}_c} \left[\log(e^{\mathbf{u}_o^\top \mathbf{v}_c}) - \log \left(\sum_{w \in \text{Vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) \right] \\
 &= -\frac{\partial}{\partial \mathbf{v}_c} [\mathbf{u}_o^\top \mathbf{v}_c] + \frac{\partial}{\partial \mathbf{v}_c} \left[\log \left(\sum_{w \in \text{Vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) \right] \\
 &\quad \text{Softmax처럼 계산된다} \quad \text{인덱스는 다르게} \\
 &\quad \text{log e 떠라...} \quad \text{인덱스는 다르게} \\
 &\quad \mathbf{a}^\top \mathbf{x} \stackrel{?}{=} (\mathbf{a}, \mathbf{x}; \text{vector}) \quad \text{인덱스는 다르게} \\
 &= -(\mathbf{u}_o) + \frac{\sum_{x \in \text{Vocab}} \mathbf{u}_x \cdot e^{\mathbf{u}_x^\top \mathbf{v}_c}}{\sum_{w \in \text{Vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \\
 &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \frac{e^{\mathbf{u}_x^\top \mathbf{v}_c}}{\sum_{w \in \text{Vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \cdot \mathbf{u}_x \\
 &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} p(u_x | \mathbf{v}_c) \cdot \mathbf{u}_x \quad P(O=o | C=c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \\
 &\quad \text{인덱스는 다르게} \\
 &\quad \hat{y}_o = p(u_o | \mathbf{v}_c) \\
 &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \hat{y}_x \cdot \mathbf{u}_x \quad \Rightarrow (\mathbf{y}, \hat{\mathbf{y}}, \mathbf{v}) \text{로의 표현}
 \end{aligned}$$

- (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. ① when $w = o$, the true 'outside' word vector, and ② for all other words. Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c .

① $w \neq o$ 일 때 실제 outside vector 인 경우

② $w \neq o$ 일 때 실제 outside vector가 아닐 경우

$$\begin{aligned} \textcircled{1} \quad \frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{u}_{w=0}} &= \frac{\partial}{\partial \mathbf{u}_{w=0}} \left[-\log(\hat{y}_o) \right] \\ &= \frac{\partial}{\partial \mathbf{u}_{w=0}} \left[-\log \left(\frac{e^{\mathbf{u}_o^\top \mathbf{v}_c}}{\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \right) \right] \\ &= -\frac{\partial}{\partial \mathbf{u}_{w=0}} \left[\mathbf{u}_o^\top \mathbf{v}_c \right] + \frac{\partial}{\partial \mathbf{u}_{w=0}} \left[\log \left(\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) \right] \\ &= -\mathbf{v}_c + \left(\frac{\frac{\partial}{\partial \mathbf{u}_{w=0}} \left(\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) \cdot \mathbf{v}_c}{\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \right) \underbrace{e^{\mathbf{u}_o^\top \mathbf{v}_c}}_{\text{(여기서 } w=o \text{인 경우)} \text{ 놓친다.)}} \\ &= -\mathbf{v}_c + \hat{y}_o \cdot \mathbf{v}_c = \mathbf{v}_c (\hat{y}_o - 1) \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{u}_{w \neq o}} &= \frac{\partial}{\partial \mathbf{u}_{w \neq o}} \left[-\log(\hat{y}_o) \right] \\ &= \frac{\partial}{\partial \mathbf{u}_{w \neq o}} \left[-\log \left(\frac{e^{\mathbf{u}_o^\top \mathbf{v}_c}}{\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \right) \right] \\ &= -\frac{\partial}{\partial \mathbf{u}_{w \neq o}} \left[\mathbf{u}_o^\top \mathbf{v}_c \right] + \frac{\partial}{\partial \mathbf{u}_{w \neq o}} \left[\log \left(\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) \right] \\ &= 0 + -\frac{e^{\mathbf{u}_{w \neq o}^\top \mathbf{v}_c} \cdot \mathbf{v}_c}{\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^\top \mathbf{v}_c}} \underbrace{\hat{y}_{w \neq o}}_{\text{여기서 } w \neq o \text{인 경우)} \\ &= \mathbf{v}_c \cdot \hat{y}_{w \neq o} \end{aligned}$$

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

Sigmoid 편미분

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

Please compute the derivative of $\sigma(x)$ with respect to x , where x is a vector.

$$\begin{aligned}\frac{d\sigma}{dx} &= \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] = \frac{d}{dx} \left[(1+e^{-x})^{-1} \right] \\ &= \left[-(1+e^{-x})^{-2} \right] [-e^{-x}] \\ &= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{(1+e^{-x})} \left(\frac{e^{-x}}{1+e^{-x}} \right) \\ &= \left(\frac{1}{(1+e^{-x})} \right) \cdot \left(\frac{e^{-x} + -1}{1+e^{-x}} \right) = \frac{1}{1+e^{-x}} \cdot \left(1 - \frac{1}{1+e^{-x}} \right) \\ &= \sigma(x) (1 - \sigma(x))\end{aligned}$$

- (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:

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

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 $\mathbf{J}_{\text{neg-sample}}$ 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.

CASE 1 : Center Word에 대해서 미분

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{neg-sample}}}{\partial \mathbf{v}_c} &= \frac{\partial}{\partial \mathbf{v}_c} \left[-\log(6(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(6(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= \frac{\partial}{\partial \mathbf{v}_c} \left[-\log(6(\mathbf{u}_o^\top \mathbf{v}_c)) \right] - \frac{\partial}{\partial \mathbf{v}_c} \left[\sum_{k=1}^K \log(6(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &\stackrel{\text{의미}}{=} -\frac{\mathbf{u}_o 6(\mathbf{u}_o^\top \mathbf{v}_c)(1 - 6(\mathbf{u}_o^\top \mathbf{v}_c))}{6(\mathbf{u}_o^\top \mathbf{v}_c)} - \sum_{k=1}^K \frac{-\mathbf{u}_k 6(-\mathbf{u}_k^\top \mathbf{v}_c)(1 - 6(-\mathbf{u}_k^\top \mathbf{v}_c))}{6(-\mathbf{u}_k^\top \mathbf{v}_c)} \\ &\stackrel{\text{sigmoid 미분}}{=} -\mathbf{u}_o (1 - 6(\mathbf{u}_o^\top \mathbf{v}_c)) + \sum_{k=1}^K \mathbf{u}_k (1 - 6(-\mathbf{u}_k^\top \mathbf{v}_c)) \end{aligned}$$

CASE 2 : outside word에 대해서 미분

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{negative sample}}}{\partial \mathbf{u}_o} &= \frac{\partial}{\partial \mathbf{u}_o} \left[-\log(6(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(6(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= \frac{\partial}{\partial \mathbf{u}_o} \left[-\log(6(\mathbf{u}_o^\top \mathbf{v}_c)) \right] - \frac{\partial}{\partial \mathbf{u}_o} \left[\sum_{k=1}^K \log(6(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= \frac{\partial}{\partial \mathbf{u}_o} \left[-\log(6(\mathbf{u}_o^\top \mathbf{v}_c)) \right] \end{aligned}$$

\$\mathbf{u}_o\$가 \$\mathbf{u}_k\$에
대한 미분

$$= - \left[\frac{V_c \cdot 6(u_0^\top v_c)(1 - 6(u_0^\top v_c))}{6(u_0^\top v_c))} \right] = -V_c(1 - 6(u_0^\top v_c))$$

CASE 3 : K negative samples 와 대비 미분

$$\begin{aligned}\frac{\partial J_{\text{neg-samp}}}{\partial u_k} &= \frac{\partial}{\partial u_k} \left[-\log(6(u_0^\top v_c)) - \sum_{k=1}^K \log(6(-u_k^\top v_c)) \right] \\ O &= \frac{\partial}{\partial u_k} \left[-\log(6(u_0^\top v_c)) \right] - \frac{\partial}{\partial u_k} \left[\sum_{k=1}^K \log(6(-u_k^\top v_c)) \right] \\ &= -\frac{\partial}{\partial u_k} \left[\sum_{k=1}^K \log(6(-u_k^\top v_c)) \right] \\ &= -\frac{-V_c 6(-u_k^\top v_c)(1 - 6(-u_k^\top v_c))}{6(-u_k^\top v_c)} = V_c(1 - 6(-u_k^\top v_c))\end{aligned}$$

- (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) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
- (ii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$
- (iii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w$ when $w \neq c \rightarrow$ 다른 뜻의 weight $\Rightarrow 0$

Write your answers in terms of $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U}$ and $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c$. This is very simple – each solution should be one line.

Once you're done: Given that you computed the derivatives of $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ with respect to all the model parameters \mathbf{U} and \mathbf{V} in parts (a) to (c), you have now computed the derivatives of the full loss function $J_{\text{skip-gram}}$ with respect to all parameters. You're ready to implement word2vec!

$$\begin{aligned} \textcircled{1} \quad & \partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U} \\ & \sum_{-m \leq j \leq m, j \neq 0} \partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad & \partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c \\ & \sum_{-m \leq j \leq m, j \neq 0} \partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c \end{aligned}$$

$$\begin{aligned} \textcircled{3} \quad & \partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w \quad (w \neq c) \\ & = 0 \end{aligned}$$