

# 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  $c$ .

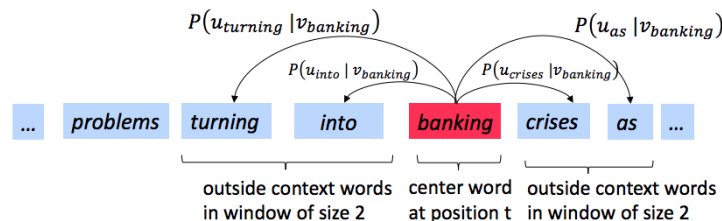


Figure 1: The word2vec skip-gram prediction model with window size 2

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  $w \in \text{Vocabulary}$ .<sup>1</sup>

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

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

Another way to view this loss is as the cross-entropy<sup>2</sup> 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  $k^{\text{th}}$  entry in these vectors indicates the conditional probability of the  $k^{\text{th}}$  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).

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

<sup>1</sup>Assume that every word in our vocabulary is matched to an integer number  $k$ .  $\mathbf{u}_k$  is both the  $k^{\text{th}}$  column of  $\mathbf{U}$  and the 'outside' word vector for the word indexed by  $k$ .  $\mathbf{v}_k$  is both the  $k^{\text{th}}$  column of  $\mathbf{V}$  and the 'center' word vector for the word indexed by  $k$ . **In order to simplify notation we shall interchangeably use  $k$  to refer to the word and the index-of-the-word.**

<sup>2</sup>The Cross Entropy Loss between the true (discrete) probability distribution  $p$  and another distribution  $q$  is  $-\sum_i p_i \log(q_i)$ .

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

Your answer should be one line.

$$\begin{aligned} (a) \quad -\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) &= -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - \dots - y_o \log(\hat{y}_o) - \dots - y_w \log(\hat{y}_w) \\ &= 0 \cdot \log \sim - 0 \cdot \log \sim \dots - 1 \cdot \log(\hat{y}_o) - \dots - 0 \cdot \log \sim = -\log(\hat{y}_o) \end{aligned}$$

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

$$\begin{aligned} (b) \quad J &= -\log \frac{\exp(\mathbf{U}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)} \\ \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial}{\partial \mathbf{v}_c} \left( -\log \frac{\exp(\mathbf{U}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)} \right) = \frac{\partial}{\partial \mathbf{v}_c} (-\mathbf{U}_o^T \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)) \\ &= -\mathbf{U}_o^T + \frac{\exp(\mathbf{U}_w^T \mathbf{v}_c) \cdot \mathbf{U}_w^T}{\sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)} = -\mathbf{U}_o^T + P(w|c) \cdot \mathbf{U}_w^T \quad \mathbf{U}_o = \mathbf{U} \cdot \mathbf{y} \rightarrow P(w|c) \cdot \mathbf{U}_w^T = \mathbf{U} \cdot \hat{\mathbf{y}}^T \\ \therefore \frac{\partial J}{\partial \mathbf{v}_c} &= \mathbf{U}^T (\hat{\mathbf{y}} - \mathbf{y})^T \end{aligned}$$

(c) (5 points) Compute the partial derivatives of  $J_{\text{naive-softmax}}(\mathbf{v}_c, \mathbf{o}, \mathbf{U})$  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  $\mathbf{v}_c$ .

$$\begin{aligned} (c) \quad \frac{\partial J}{\partial \mathbf{u}_o} &= \frac{\partial}{\partial \mathbf{u}_o} \left( -\log \frac{\exp(\mathbf{U}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)} \right) = -\mathbf{v}_c + \frac{\partial}{\partial \mathbf{u}_o} (\log \sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)) \\ &= -\mathbf{v}_c + \frac{\exp(\mathbf{U}_w^T \mathbf{v}_c) \cdot \mathbf{v}_c}{\sum_{w \in \text{Vocab}} \exp(\mathbf{U}_w^T \mathbf{v}_c)} = -\mathbf{v}_c + P(w|c) \cdot \mathbf{v}_c \\ \therefore \frac{\partial J}{\partial \mathbf{u}_o} &= \begin{cases} (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{v}_c & (w = o) \\ \hat{\mathbf{y}}_w \cdot \mathbf{v}_c & (w \neq o) \end{cases} \end{aligned}$$

(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 vector.

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

- (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.<sup>3</sup>

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.

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

$$\text{ii) } \frac{\partial \mathbf{J}}{\partial \mathbf{u}_o} = \frac{\partial}{\partial \mathbf{u}_o} \left( -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right) = -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1-\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \cdot \mathbf{v}_c}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} = -(1-\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \cdot \mathbf{v}_c$$

$\mathbf{u}_o$  포함  $\times \rightarrow$  연산에 영향  $\times$

$$\text{iii) } \frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} = \frac{\partial}{\partial \mathbf{u}_k} \left( -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right) = \sum_{k=1}^K \frac{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)(1-\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \cdot \mathbf{v}_c}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} = \sum_{k=1}^K (1-\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \cdot \mathbf{v}_c$$

$\mathbf{u}_k$  포함  $\times \rightarrow$  연산에 영향  $\times$

iv) Naive Softmax Loss에 비해 Negative Sampling Loss의 연산량이 더 작기 때문에 more efficient!

- (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:

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

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

Write down three partial derivatives:

- $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
- $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$
- $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w$  when  $w \neq c$

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

**Once you're done:** Given that you computed the derivatives of  $\mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$  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  $\mathbf{J}_{\text{skip-gram}}$  with respect to all parameters. You're ready to implement word2vec!

$$\begin{aligned} \text{(i) } \frac{\partial \mathbf{J}}{\partial \mathbf{U}} &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}} & \text{(ii) } \frac{\partial \mathbf{J}}{\partial \mathbf{v}_c} &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c} & \text{(iii) } \frac{\partial \mathbf{J}}{\partial \mathbf{v}_w} &= 0 \quad (w \neq c) \end{aligned}$$

center word가 아니면 파라미터 업데이트 되지 않음  
 $\therefore \frac{\partial \mathbf{J}}{\partial \mathbf{v}_w} = 0 \quad (w \neq c)$