

# 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  $c$ .

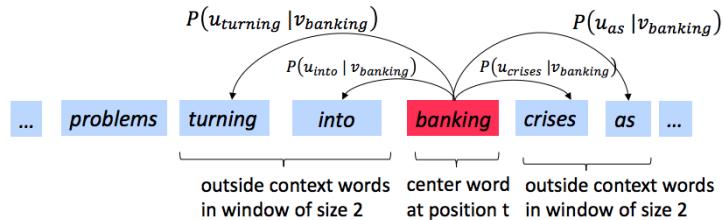


Figure 1: The word2vec skip-gram prediction model with window size 2

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  $-\sum_i p_i \log(q_i)$ .

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

Your answer should be one line.

- (b) (5 points) Compute the partial derivative of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$ . Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{U}$ .
- (c) (5 points) Compute the partial derivatives of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  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 you answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{v}_c$ .
- (d) (3 Points) The sigmoid function is given by Equation 4:

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

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

- (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.

- (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:

- (i)  $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
- (ii)  $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$

---

<sup>3</sup>Note: the loss function here is the negative of what Mikolov et al. had in their original paper, because we are doing a minimization instead of maximization in our assignment code. Ultimately, this is the same objective function.

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

## 2 Coding: Implementing word2vec (20 points)

In this part you will implement the word2vec model and train your own word vectors with stochastic gradient descent (SGD). Before you begin, first run the following commands within the assignment directory in order to create the appropriate conda virtual environment. This guarantees that you have all the necessary packages to complete the assignment.

```
conda env create -f env.yml  
conda activate a2
```

Once you are done with the assignment you can deactivate this environment by running:

```
conda deactivate
```

- (a) (12 points) First, implement the `sigmoid` function in `word2vec.py` to apply the sigmoid function to an input vector. In the same file, fill in the implementation for the softmax and negative sampling loss and gradient functions. Then, fill in the implementation of the loss and gradient functions for the skip-gram model. When you are done, test your implementation by running `python word2vec.py`.
- (b) (4 points) Complete the implementation for your SGD optimizer in `sgd.py`. Test your implementation by running `python sgd.py`.
- (c) (4 points) Show time! Now we are going to load some real data and train word vectors with everything you just implemented! We are going to use the Stanford Sentiment Treebank (SST) dataset to train word vectors, and later apply them to a simple sentiment analysis task. You will need to fetch the datasets first. To do this, run `sh get_datasets.sh`. There is no additional code to write for this part; just run `python run.py`.

*Note: The training process may take a long time depending on the efficiency of your implementation (an efficient implementation takes approximately an hour). Plan accordingly!*

After 40,000 iterations, the script will finish and a visualization for your word vectors will appear. It will also be saved as `word_vectors.png` in your project directory. **Include the plot in your homework write up.** Briefly explain in at most three sentences what you see in the plot.

## 3 Submission Instructions

You shall submit this assignment on GradeScope as two submissions – one for “Assignment 2 [coding]” and another for ‘Assignment 2 [written]’:

- (a) Run the `collect_submission.sh` script to produce your `assignment2.zip` file.
- (b) Upload your `assignment2.zip` file to GradeScope to “Assignment 2 [coding]”.
- (c) Upload your written solutions to GradeScope to “Assignment 2 [written]”.

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

- (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

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

Your answer should be one line.

先来厘清一下思路，对于 skip-gram 模型：

输入：随机化的  $\mathbf{U}, \mathbf{V}$ .      输出：结果  $\mathbf{U}, \mathbf{V}$

对某个单词  $c$ :  $P(O=o|C=c) = \frac{\exp(\mathbf{U}_o^T \mathbf{V}_c)}{\sum_{w \in Vocab} \exp(\mathbf{U}_w^T \mathbf{V}_c)}$

对词表中每个单词的  $P(O=o|C=c)$  组成  $\hat{\mathbf{y}}$ .

真正的  $\mathbf{y}$  应为一个只有 outside word 为 1, 其余为 0 的向量.

那么对一对单词  $c$  和  $o$ :

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

$$J_{\text{cross-entropy}}(\mathbf{v}_c, o, \mathbf{U}) = -\sum_{w \in Vocab} y \log P(O=w|C=c)$$

只有当  $w$  是真实的 outside word 时  $y=1$ , 其余情况  $y=0$ .

$$\begin{aligned} \therefore J_{\text{cross-entropy}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log P(O=o|C=c) \\ &= J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) \end{aligned}$$

$\because$  对每对单词  $c$  和  $o$  两者都相等,  $\therefore$  两个 loss function 相等

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

- (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  $\mathbf{U}$ .

对一对单词  $O$  和  $C$ :

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} \text{ 在课上推导过, 在这里直接给出结论:}$$

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{x \in V} P(x|c) \cdot \mathbf{v}_x$$

$$= -\mathbf{u}_o + \sum_{x \in V} \hat{\mathbf{y}}_x \cdot \mathbf{v}_x = \mathbf{U}^T (\hat{\mathbf{y}} - \mathbf{y})$$

- (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  $\mathbf{v}_c$ .

当  $w=o$  时:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \frac{\partial -\log P(O=o|C=c)}{\partial \mathbf{u}_o}$$

$$= \frac{\partial -\log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)}}{\partial \mathbf{u}_o}$$

$$= -\frac{\partial \log \exp(\mathbf{u}_o^T \mathbf{v}_c)}{\partial \mathbf{u}_o} + \frac{\partial \log \sum_{w \in \text{vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\partial \mathbf{u}_o}$$

这个数值只有一项与  $\mathbf{u}_o$  有关, 其余看作常数, 记为 A

$$= -\frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{u}_o} + \frac{\partial \log(A + \exp(\mathbf{u}_o^T \mathbf{v}_c))}{\partial \mathbf{u}_o}$$

$$= -\mathbf{v}_c^T + \frac{1}{A + \exp(\mathbf{u}_o^T \mathbf{v}_c)} \cdot \frac{\partial (A + \exp(\mathbf{u}_o^T \mathbf{v}_c))}{\partial \mathbf{u}_o}$$

$$= -\mathbf{v}_c^T + \frac{1}{A + \exp(\mathbf{u}_o^T \mathbf{v}_c)} \cdot \frac{\partial \exp(\mathbf{u}_o^T \mathbf{v}_c)}{\partial \mathbf{u}_o}$$

$$= -\mathbf{v}_c^T + \frac{1}{A + \exp(\mathbf{u}_o^T \mathbf{v}_c)} \cdot \exp(\mathbf{u}_o^T \mathbf{v}_c) \cdot \mathbf{v}_c^T$$

$$= -V_C^T + \frac{\exp(U_o^T V_C)}{\sum_{w \in \text{vocab}} \exp(U_w^T V_C)} \cdot V_C^T$$

$$= -V_C^T + \hat{y}_o \cdot V_C^T = (\hat{y}_o - 1) V_C^T$$

当  $w \neq o$  时：

$$\frac{\partial J_{\text{naive-softmax}}(V_C, o, U)}{\partial U_X} = \frac{\partial -\log p(o=o|C=c)}{\partial U_X}$$

$$\text{常数} = \frac{-\partial \log \exp(U_o^T V_C) + \partial \log \sum_{w \in \text{vocab}} \exp(U_w^T V_C)}{\partial U_X}$$

只有一项不是常数，  
其他常数认为 A.

$$= \frac{\partial \log(A + \exp(U_X^T V_C))}{\partial U_X} = \frac{1}{A + \exp(U_X^T V_C)} \cdot \frac{\partial (A + \exp(U_X^T V_C))}{\partial U_X}$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(U_w^T V_C)} \cdot \frac{\partial \exp(U_X^T V_C)}{\partial U_X}$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(U_w^T V_C)} \cdot \exp(U_X^T V_C) \cdot V_C^T$$

$$= \hat{y}_X \cdot V_C^T = \hat{y}_w \cdot V_C^T$$

综上：

$$\frac{\partial J_{\text{naive-softmax}}(V_C, o, U)}{\partial U}$$

$$= (\hat{y} - y)^T V_C$$

$1 \times n \quad 1 \times d$

技巧：

$$(1) \frac{\partial U^T V}{\partial U} = V^T$$

$$(2) \frac{\partial U^T V}{\partial V} = U$$

(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.

我们讨论向量  $x$  中的一个元素  $x_i$ :

$$\frac{\partial \sigma(x_i)}{\partial x_i} = \frac{\partial \frac{e^{x_i}}{e^{x_i} + 1}}{\partial x_i} = \frac{\partial \frac{e^{x_i}}{e^{x_i} + 1}}{\partial e^{x_i}} \cdot \frac{\partial e^{x_i}}{\partial x_i} = \frac{\partial (1 - \frac{1}{e^{x_i} + 1})}{\partial e^{x_i}} \cdot e^{x_i}$$

$$\frac{\partial e^{x_i}}{\partial x_i} = u \quad \frac{\partial (1 - \frac{1}{u+1})}{\partial u} \cdot u = \frac{1}{(u+1)^2} \cdot u = \frac{1}{(1+e^{x_i})^2} \cdot e^{x_i}$$

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

(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  $u_1, \dots, u_K$ . 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:

$$J_{\text{neg-sample}}(v_c, o, U) = -\log(\sigma(u_o^\top v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^\top v_c)) \quad (5)$$

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  $J_{\text{neg-sample}}$  with respect to  $v_c$ , with respect to  $u_o$ , and with respect to a negative sample  $u_k$ . Please write your answers in terms of the vectors  $u_o$ ,  $v_c$ , and  $u_k$ , 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.

$$\begin{aligned} (1) \quad \frac{\partial J_{\text{neg-sample}}(v_c, o, U)}{\partial v_c} &= \frac{\partial -\log(\sigma(u_o^\top v_c))}{\partial v_c} - \frac{\partial \sum_{k=1}^K \log(\sigma(-u_k^\top v_c))}{\partial v_c} \\ &= \frac{-1}{\sigma(u_o^\top v_c)} \cdot \frac{\partial \sigma(u_o^\top v_c)}{\partial v_c} - \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^\top v_c))}{\partial v_c} \\ &= -\frac{1}{\sigma(u_o^\top v_c)} \cdot \sigma(u_o^\top v_c) \cdot [1 - \sigma(u_o^\top v_c)] \cdot \frac{\partial u_o^\top v_c}{\partial v_c} \\ &\quad - \sum_{k=1}^K \frac{1}{\sigma(-u_k^\top v_c)} \cdot \frac{\partial \sigma(-u_k^\top v_c)}{\partial v_c} \end{aligned}$$

$$= [g(u_0^T v_c) - 1] \cdot u_0 - \sum_{k=1}^K \frac{1}{g(-u_k^T v_c)} \cdot g(-u_k^T v_c) \cdot [1 - g(-u_k^T v_c)]$$

$$\cdot \frac{\partial (-u_k^T v_c)}{\partial v_c}$$

$$= [g(u_0^T v_c) - 1] \cdot u_0 - \sum_{k=1}^K [1 - g(-u_k^T v_c)] \cdot (-u_k)$$

$$= [g(u_0^T v_c) - 1] \cdot u_0 + \sum_{k=1}^K [1 - g(-u_k^T v_c)] \cdot u_k$$

$$(2) \frac{\partial J_{\text{neg-sample}}(v_c, o, u)}{\partial u_0} = \frac{\partial -\log(g(u_0^T v_c))}{\partial u_0} - \boxed{\frac{\partial \sum_{k=1}^K \log(g(-u_k^T v_c))}{\partial u_0}}$$

常数

$$= \frac{\partial -\log(g(u_0^T v_c))}{\partial u_0} = \frac{-1}{g(u_0^T v_c)} \cdot g(u_0^T v_c) \cdot [1 - g(u_0^T v_c)] \cdot v_c$$

$$= [g(u_0^T v_c) - 1] \cdot v_c$$

$$(3) \frac{\partial J_{\text{neg-sample}}(v_c, o, u)}{\partial u_k} = \boxed{\frac{\partial -\log(g(u_0^T v_c))}{\partial u_k}} - \frac{\partial \sum_{k=1}^K \log(g(-u_k^T v_c))}{\partial u_k}$$

常数 除了  $u_k$  其余项常数

$$= - \frac{\partial \log(g(-u_k^T v_c))}{\partial u_k} = \frac{1}{g(-u_k^T v_c)} \cdot g(-u_k^T v_c) \cdot$$

$$[1 - g(-u_k^T v_c)] \cdot v_c^T$$

$$= [1 - g(-u_k^T v_c)] \cdot v_c^T$$

注意到这3问的结果中存在大量可以复用的项，这使得在计算梯度时，更加高效。

这个损失函数从  $V$  的多分类变成 0, 1 二分类，每次输出概率从  $V$  减小到  $2 \times K$ ，从  $V$  个向量相乘减小到  $K$  个向量相乘。（同时也可以提升词向量的效果）。

(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  $w \neq c$

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!

$$(i) \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

$$(ii) \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

$$(iii) \sum_{\substack{-m \leq j \leq m \\ j \neq 0 \\ w \neq c}} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_w} = 0$$

## 2. Coding

按照注释 U 中的每一行作为一个 outside vector, 假设一共有 n 个单词.

$$U = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} \quad \text{其中 } u_i \text{ 为 } d \text{ 维向量 } (1 \times d) \quad U \text{ 的 shape: } (n \times d)$$

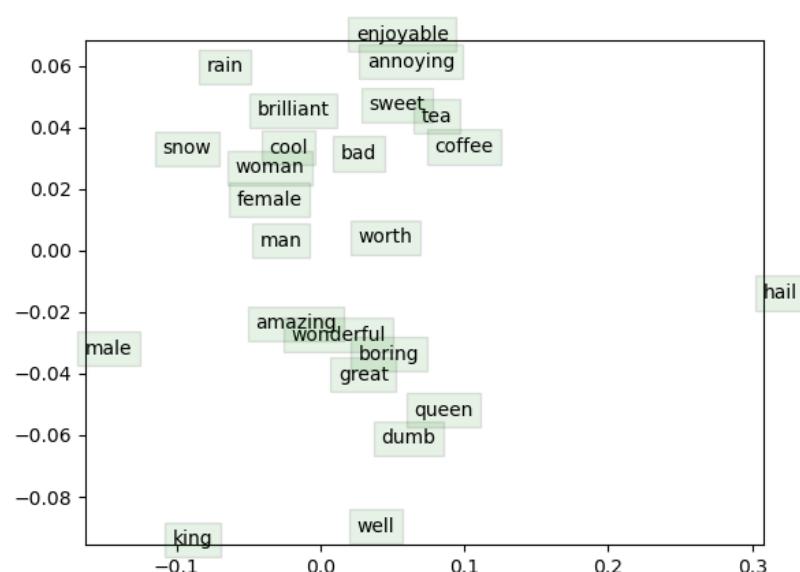
$$\frac{\partial J}{\partial U} = \frac{\partial J}{\partial \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}} = \left( \frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_n} \right) \quad d \times n \text{ 维矩阵}$$

$\frac{\partial J}{\partial u_i}$  为  $d \times 1$  维

在前面的 writing 部分, 我们已经求出了  $\frac{\partial J}{\partial u_i}$ , 现在只要代入就可以完成,

总结: 在矩阵求导和编写代码的过程中, 可以领悟到在深度学习中, 得到的结果是  $n \times k$ , 还是  $k \times n$  真的影响不大, 主要看问题 / 输出想要什么维度, 只要保证乘法时维度能 match 上就可以了.

coding 部分的结果:



从图上来看, "woman" 和 "female" 很接近, 这点是符合认知的 (除此还有, "tea", "coffee"; "amazing", "wonderful" 等等.)

male - King,  
female - queen 有对应关系.