

# XCS224N Assignment #2: word2vec (20 Points + 5 Point Extra Credit Challenge)

## 1 Understanding word2vec (Refresher)

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  $c$ .

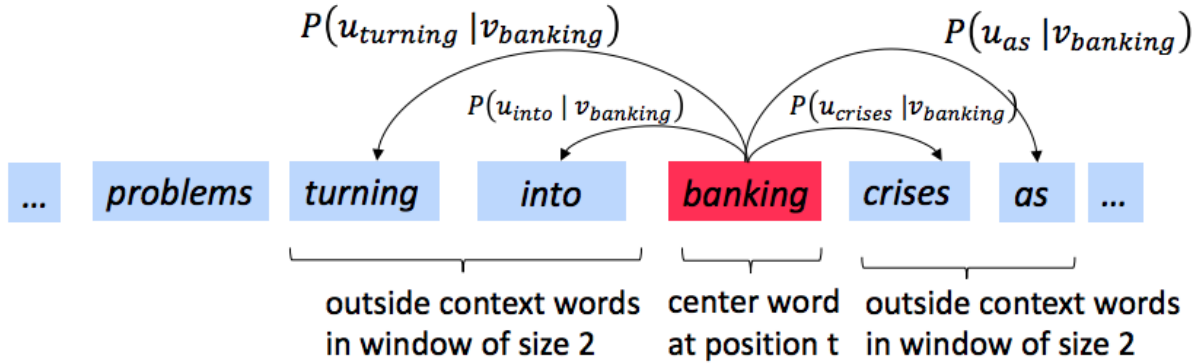


Figure 1: The word2vec skip-gram prediction model with window size 2

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:

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

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

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

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). There are 2 optional questions below which you may attempt (*Note: Bonus points will be awarded for the optional assignments*)

## 2 Extra Credit Challenge I (2.5 Points)

The partial derivative of  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$  in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{U}$  is given below:

$$\frac{\partial J}{\partial \mathbf{v}_c} = \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \quad (3)$$

or equivalently,

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{w=1}^V \hat{y}_w \mathbf{u}_w \quad (4)$$

The naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ ; This is equivalent to:

$$- \sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o). \quad (5)$$

**Reasoning:** Since  $\mathbf{y}$  is a one-hot vector, all  $y_k = 0$  where  $k \neq o$ .  $y_o = 1$ , so we are left with  $-\log(\hat{y}_o)$ .

**Write the steps** to arrive at equation 3 or 4, the partial derivative of  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$ , starting from equation 5. Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{U}$ . The first few steps have been provided below (loss function  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ ) and the rest of the proof may take 4 or 5 steps.

$$\begin{aligned} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log(\hat{y}_o) && \text{(refer to equation 5)} \\ &= -\log\left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}\right) \\ &= -\left(\log(\exp(\mathbf{u}_o^\top \mathbf{v}_c)) - \log\left(\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)\right)\right) \\ &= -\mathbf{u}_o^\top \mathbf{v}_c + \log\left(\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)\right) \end{aligned}$$

### 3 Extra Credit Challenge II (2.5 Points)

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 is given below:

$$\frac{\partial J}{\partial \mathbf{U}} = \mathbf{v}_c(\hat{\mathbf{y}} - \mathbf{y})^\top \quad (6)$$

or equivalently:

$$\frac{\partial J}{\partial \mathbf{u}_w} = \begin{cases} (\hat{y}_w - 1)\mathbf{v}_c & \text{if } w = o \\ \hat{y}_w\mathbf{v}_c & \text{otherwise} \end{cases} \quad (7)$$

**Write the steps** required to arrive at the partial derivative 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 are two cases you need to consider: 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  $\mathbf{v}_c$ . The proof may take 4 or 5 steps. The loss function  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  is:

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

## 4 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 virtual environment. This guarantees that you have all the necessary packages to complete the assignment. Make sure to take a look at the equations presented in section 2 and 3, especially **equations (3) for gradient (of loss function  $J_{\text{naive-softmax}}(v_c, o, U)$ ) with respect to  $v_c$  and (6) for gradient with respect to each of the ‘outside’ word vectors,  $u_w$ ’s** which we’ll be using in our gradient calculation code.

```
source venv/bin/activate
```

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

```
deactivate
```

**Note:** A2 code can be found here : <https://github.com/scpd-proed/XCS224N-A2>

- (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 naive softmax loss and gradient function. 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`. To complete these steps you will need to refer to Sections 2 and 3, Equations (3) and (6).

### Pseudo code for skip-gram

---

#### Algorithm 1 Skipgram

---

```
procedure SKIPGRAM(outside_words, loss_and_gradient_func, ** kwargs)    ▷ The center word vector
    loss ← 0
    grad_center ← np.zeros(*args)
    grad_outside ← np.zeros(*args)
    for word in outside_words do                                     ▷ Iterate over outside words
        loss_current, gradc, grado ← loss_and_gradient_func(** kwargs)
        loss ← loss + loss_current                                   ▷ Loss is accumulated
        grad_center ← grad_center + gradc
        grad_outside ← grad_outside + grado
    end for
    return loss, grad_center, grad_outside                            ▷ Return loss and gradients
end procedure
```

---

- (b) (5 points) Complete the implementation for your SGD optimizer in `sgd.py`. Test your implementation by running `python sgd.py`.
- (c) (3 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. There is no additional code to write for this part; just run `python run.py`.

**Note:** *Please do not use external python modules which are not specified in the requirements.txt file. This will cause the autograder to fail.*

**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` and the corresponding wordvectors as `sample_vectors.json` in your project directory.

## Submission Instructions

1. **Please** do not use external python modules which are not specified in the `requirements.txt` file. This will cause the autograder to fail.
2. Run the `collect_submission.sh` script to produce your `assignment2.zip` file or check `README.md` to see the necessary list of files to submit.
3. Upload your `assignment2.zip` file or the individual files via the **Assignment 2 Coding Submission Link** on the Assignment 2 page of your SCPD learning portal.
4. If you completed the extra credit derivation prompts, upload your responses as a PDF via the **Assignment 2 Extra Credit Submission Link** on the Assignment 2 page of your SCPD learning portal.

## 5 Appendix

Negative sampling is briefly introduced in Lecture 2: Word2Vec: Model Variants and an implementation is provided in the Assignment 2 coding assignment. For detailed notes on the math behind negative sampling, see below.

1. Negative Sampling loss 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:

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (8)$$

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

Below we compute 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  $\mathbf{u}_k$ .

**Solution:** Firstly:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \times \frac{\partial}{\partial \mathbf{v}_c} \sigma(\mathbf{u}_o^\top \mathbf{v}_c) - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \times \frac{\partial}{\partial \mathbf{v}_c} \sigma(-\mathbf{u}_k^\top \mathbf{v}_c) && \text{(chain rule on log)} \\ &= -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{u}_o}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} - \sum_{k=1}^K -\frac{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))\mathbf{u}_k}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} && \text{(chain rule on } \sigma) \\ &= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{u}_o - \sum_{k=1}^K -(1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))\mathbf{u}_k && \text{(cancel)} \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{u}_o - \sum_{k=1}^K (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1)\mathbf{u}_k && \text{(rearrange)} \end{aligned}$$

Secondly:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_o} &= \frac{\partial}{\partial \mathbf{u}_o} \left( -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \right) \\ &= -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \times \frac{\partial}{\partial \mathbf{u}_o} \left( \sigma(\mathbf{u}_o^\top \mathbf{v}_c) \right) && \text{(chain rule on log)} \\ &= -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{v}_c}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} && \text{(chain rule on } \sigma) \\ &= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{v}_c && \text{(cancel)} \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{v}_c && \text{(rearrange)} \end{aligned}$$

---

<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.

Thirdly, for all  $k = 1, 2, \dots, K$ :

$$\begin{aligned}
 \frac{\partial J}{\partial \mathbf{u}_k} &= \frac{\partial}{\partial \mathbf{u}_k} \left( -\log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right) \\
 &= \frac{1}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \times \frac{\partial}{\partial \mathbf{u}_k} \left( \sigma(-\mathbf{u}_k^\top \mathbf{v}_c) \right) \quad (\text{chain rule on log}) \\
 &= \frac{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))\mathbf{v}_c}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \quad (\text{chain rule on } \sigma) \\
 &= (1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c))\mathbf{v}_c \quad (\text{cancel})
 \end{aligned}$$

The naive-softmax loss contains a summation over the entire vocabulary as part of computing the  $P(O = o \mid C = c)$  term. Here, we don't do that calculation, approximating it with  $K$  samples (where  $K$  is much smaller than the vocabulary size).

2. 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:

$$\mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (9)$$

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 the implementation.

The proofs for the three partial derivatives are given below:

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

*Note that the final derivatives of  $\mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$  with respect to all the model parameters  $\mathbf{U}$  and  $\mathbf{V}$  are provided in the appendix part 1 derivation*

**Solution:** Given a loss function  $J$ , we already know how to obtain the following derivatives

$$\frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}} \quad \text{and} \quad \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

Therefore, for skip-gram, the gradients for the loss of one context window can be expressed in terms of these:

$$\begin{aligned}
 \frac{\partial \mathbf{J}_{\text{skip-gram}}(w_{t-m} \dots w_{t+m})}{\partial \mathbf{U}} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}, \\
 \frac{\partial \mathbf{J}_{\text{skip-gram}}(w_{t-m} \dots w_{t+m})}{\partial \mathbf{v}_c} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}, \\
 \frac{\partial \mathbf{J}_{\text{skip-gram}}(w_{t-m} \dots w_{t+m})}{\partial \mathbf{v}_w} &= \mathbf{0}, \text{ when } w \neq c.
 \end{aligned}$$