

Midterm Group Project : word2vec

1 Written: Understanding word2vec

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 .

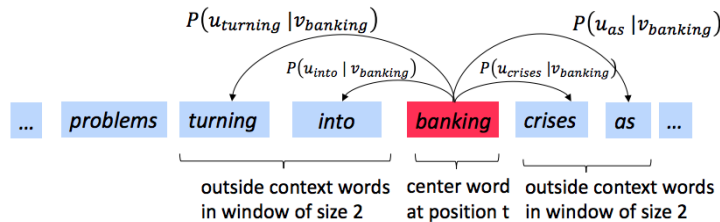


Figure 1: The word2vec skip-gram prediction model with window size 2

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:

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

Another way to view this loss is as the cross-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).

- (a) 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

¹Assume that every word in our vocabulary is matched to an integer number k . \mathbf{u}_k is both the k^{th} column of \mathbf{U} and the ‘outside’ word vector for the word indexed by k . \mathbf{v}_k is both the k^{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.**

²The 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) 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) 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 your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c .
- (d) 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 $\sigma(\mathbf{x})$ with respect to \mathbf{x} , where \mathbf{x} is a vector.

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

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

In this part you will implement the word2vec model and train your own word vectors with stochastic gradient descent (SGD). This guarantees that you have all the necessary packages to complete the assignment.

- (a) First, When you think about writing the code for word2vec,
- (a.1) Write the sigmoid function in word2vec with your Python function (or Python class) code.
 - (a.2) Write the stochastic gradient descent in word2vec with your Python function (or Python class) code.
- Only the sigmoid function and stochastic gradient descent can be represented by code.

- (b) Show time! Now we are going to load some real data and train word vectors with everything. Use **word2vec** to represent the "**Finance - related Words**". You can use `gensim` and `Konlpy`.

- You can do homework with only `gensim` dataset. But if you want to do a higher level of work, **you can crawl newspaper articles to show finance-related words**. If you use newspaper articles rather than just using `gensim` datasets, you can display a lot more words with similarities to finance.
- The script will finish and a visualization for your word vectors will appear. It will also be saved as **word2vec_plot_group number.png** in your project file. Include the plot in your report write up. And Explain the plot.

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!

3 Submission Instructions

You shall submit this assignment on GradeScope as two submissions – one for “Midterm [coding]” and another for “Midterm [solving]”:

This project is a midterm project.

- (a) Submission method : 1) YSCEC upload + 2) PAPER report directly submit (only one copy from GROUP)
- (b) Deadline for submission : **April 26, 2019 Upload and submit by 5 pm.** (There is a deduction for late submission. If you submit after May 1, we will consider that you did NOT submit your homework.)
- (c) Direct Submission place : Industry-Academic Cooperation 514 (산학협동관 514호), Front door box
- (d) For problem # 1, please solve the problem by hand.

(If you upload YSECE, please scan and convert it to pdf. Also, you can use tex file, word and Hangul. If you use them , you convert it to pdf, too.)

Please submit handwritten reports(problem #1) when submitting your report(with problem #2))

- (e) For problem # 2(coding),

Please submit the file name as follows.

- sigmoid function_group number.py,
- stochastic gradient descent_group number.py,
- word2vec_plot_group number.png
- report_group numbr.pdf (The report should include the analysis of your code and plot. You will also need to write in the report a detailed description of the idea and code that led to the project.)

If we do not have plot, we will assume that we did not submit the coding homework.

Please modify the group number that is written in bold. For example, if group1 is submitted, modify it with sigmoid function_group1.py.