

# 1 Understanding word2vec

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 of the word  $o$  is an ‘outside’ word for  $c$ , i.e., the probability that  $o$  falls within the contextual windows of  $c$ .

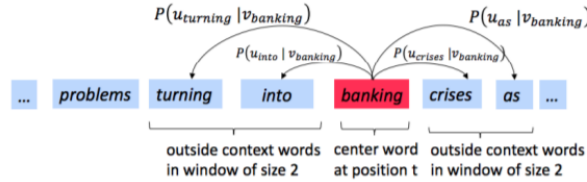


Figure 1: The word2vec skip-gram prediction model with windows size 2

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

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

Here,  $\mathbf{u}_o$  is the ‘outside’ vector representing outside word  $o$ , and  $\mathbf{v}_c$  is the center ‘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 the ‘center’ vectors  $\mathbf{v}_w$ . Both  $\mathbf{U}$  and  $\mathbf{V}$  contain a vector for every  $w \in Vocabulary$ .<sup>1</sup>

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

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

Another way to view this loss is 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 the vectors with length equal to the number of words in the vocabulary. Furthermore, the  $k^{th}$

<sup>1</sup>Assume 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 ‘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.

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

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 the Equation 1.

- (a) Naive softmax loss in the Equation 2 is the same as the cross-entropy loss between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ ; i.e.,

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

Reason: For two words  $c$  and  $o$ ,  $y_w = 1$  when  $w = o$ ,  $y_w = 0$  otherwise.

- (b) Partial derivative of  $\mathbf{J}_{naive-softmax}(\mathbf{v}_c, \mathbf{o}, \mathbf{U}) = J$  with respect to  $(v_c)$ .<sup>3</sup>  
Let us rewrite  $J$  as:

$$\begin{aligned} J &= -\log \hat{y}_o = -\log\left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}\right) \\ &= -[\log(\exp(\mathbf{u}_o^\top \mathbf{v}_c)) - \log(\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c))] \\ &= -\mathbf{u}_o^\top \mathbf{v}_c + \log(\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)) \end{aligned} \quad (4)$$

Now, let us find  $\frac{\partial J}{\partial \mathbf{v}_c}$ .<sup>4</sup>

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \frac{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \cdot \mathbf{u}_w}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (5)$$

Re-arranging,

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{w \in Vocab} \left( \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{x \in Vocab} \exp(\mathbf{u}_x^\top \mathbf{v}_c)} \mathbf{u}_w \right) \quad (6)$$

From Equation 1,

---

<sup>3</sup>For more details visit: <https://stats.stackexchange.com/questions/253244/gradients-for-skipgram-word2vec>

<sup>4</sup>If you are confused about why  $\frac{\partial \mathbf{u}_o^\top \mathbf{v}_c}{\partial \mathbf{v}_c} = \mathbf{u}_o$ , please read about numerator layout notation and denominator layout notation. You may start here: <https://www.comp.nus.edu.sg/cs5240/lecture/matrix-differentiation.pdf>

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{w \in Vocab} (\hat{y}_w \mathbf{u}_w) \quad (7)$$

Writing in the matrix form,

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= -Uy + U\hat{y} \\ &= U[\hat{y} - y] \end{aligned} \quad (8)$$