

1 Extra Credit Challenge I: Backward pass in V in Skip-Gram

The objective now is to get the gradient so we can update our word vectors to increase the probability.

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \frac{\partial}{\partial \mathbf{v}_c} \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) \quad (1)$$

We leave the left term for the end, and keep progressing with the other part.

Applying the chain rule on $f(x) = \log$ and $g(x) = \sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{v}_c} \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) &= \frac{1}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sum_k \exp(\mathbf{u}_k^\top \mathbf{v}_c) \\ &= \frac{1}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \sum_k \frac{\partial}{\partial \mathbf{v}_c} \exp(\mathbf{u}_k^\top \mathbf{v}_c) \end{aligned}$$

Again, applying the chain rule on $f(x) = \exp$ and $g(x) = \mathbf{u}_k^\top \mathbf{v}_c$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{v}_c} \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) &= \frac{1}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \sum_k \exp(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k \\ &= \sum_k \frac{\exp(\mathbf{u}_k^\top \mathbf{v}_c)}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \mathbf{u}_k \end{aligned}$$

Note the fraction was a constant. When including it inside the summation, we end up having a probability similar to what we have started with:

$$\frac{\partial}{\partial \mathbf{v}_c} \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) = \sum_k P(K|C) \mathbf{u}_k$$

Putting it all together:

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_k P(K|C) \mathbf{u}_k = -\mathbf{u}_o + \sum_k \hat{\mathbf{y}} \cdot \mathbf{u}_k \quad (2)$$

Matrix Notation:

\mathbf{u}_o can be written as a matrix multiplication Uy , since y will only be 1 for the real output vector of U , \mathbf{u}_o

$\sum_k \hat{\mathbf{y}} \cdot \mathbf{u}_k$ is a linear transformation of vectors \mathbf{u}_k in U scaled by $\hat{\mathbf{y}}$. This can be written as $U\hat{\mathbf{y}}$.

Rewritten both of them, we show that (3) and (4) from the statement of the question were equivalent:

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_k \hat{\mathbf{y}} \cdot \mathbf{u}_k = U(\hat{\mathbf{y}} - \mathbf{y}) \quad (3)$$

Alternative 1: Back-propagate through all layers

The derivative of the error directly w.r.t $(\frac{\partial J}{\partial s_j})$ the scores is fairly simply, resulting on $\hat{y} - y$. This is derived in Appendix B.2, Eq 24.

Having obtained the equation for U first (although it is the second question on this exercise), we can back-propagate now the error into the hidden layer \mathbf{h} :

$$\frac{\partial J}{\partial h_i} = \sum_{j=1}^V \frac{\partial J}{\partial s_j} \frac{\partial s_j}{\partial h_i} = \sum_{j=1}^V (\hat{y}_j - y_j) \cdot u_{ij} = \sum_{j=1}^V e_j \cdot u_{ij} := EH_i$$

EH is an N -dim vector representing the sum of the output vectors of all the words in the vocabulary, weighted by their prediction error. Now that we have the error w.r.t the hidden layer, we can move onto V :

$$\frac{\partial J}{\partial v_{ci}} = \frac{\partial J}{\partial h_i} \frac{\partial h_i}{\partial v_{ci}} = EH_i \cdot x_c$$

This is equivalent to a tensor product:

$$\frac{\partial J}{\partial V} = x \otimes EH_i = xEH_i^T$$

Resulting on a $V \times N$ matrix, same as V , fulfilling the *shape convention*.

Only one component of \mathbf{x} is non-zero, therefore only one column of $\frac{\partial J}{\partial V}$ will be zero.

Alternative 2: Using the computation graph and matrix notation

Applying the chain rule entirely from the loss to the matrix V we are interested in:

$$\frac{\partial J}{\partial V} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial h} \frac{\partial h}{\partial V} = U(\hat{y} - y)$$

We can compute the local gradients during the forward pass (derived at Appendix A.3) such that:

$$\frac{\partial J}{\partial s} = \hat{y} - y; \frac{\partial s}{\partial h} = U; \frac{\partial s}{\partial V} = 1$$

Interpretation

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_k \hat{y} \cdot \mathbf{u}_k$$

This has an interesting interpretation:

The gradient w.r.t. the input word vector, is the difference between the real output vector of the word and what our model thinks the context of what the center words should look like.

That "thinking" is an expectation, as a **weighted average of each word in the vocabulary, weighted by their probability of occurring in the context of that word.**

$$\frac{\partial J}{\partial h_i} = EH_i$$

and,

$$\mathbf{v}_c^{new} = \mathbf{v}_c^{old} - \eta EH_i$$

EH is the sum of output vectors for all words in the vocabulary weighted by their prediction error $e_j = \hat{y}_j - y_j$. We can understand above equation as adding a portion of every output vector to the input vector.

If $\hat{y}_j > y_j$ (*overestimating*):

The input vector \mathbf{v}_c will tend to move farther from the output vector u_j .

If $\hat{y}_j < y_j$ (*underestimating*) which is only true for the real target $y_j = 1$:

The input vector \mathbf{v}_c will tend to move closer to the output vector u_j .

The movement of the input vector \mathbf{v}_c is determined by the prediction error of all vectors in the vocabulary. The larger the error, the more significant effects a word will exert on the movement of the input vector of the context word.

2 Extra Credit Challenge II: Backward pass in U Skip-Gram

Chain rule to compute the derivative of the cost w.r.t. each weight u_{ij} .

We know also that the shape of the gradients of U has to be the same shape as U .

$$\frac{\partial J}{\partial u_{ij}} = \frac{\partial J}{\partial s_j} \frac{\partial s_j}{\partial u_{ij}} = (\hat{y}_j - y_j) h_i = (\hat{y}_j - y_j) h_i = e_j \cdot h_i \quad (4)$$

As can be seen in Fig 3 on Appendix A.3 for a specific output an entire row from U is contributing. Generalizing to it:

$$\frac{\partial J}{\partial u_j} = e_j \cdot \mathbf{h} \quad (5)$$

Because the error e_j depends on the real value, or ground truth y (which is 0 for all entries except 1 for entry $j = o$, and $\mathbf{h} = \mathbf{v}_c$) we can decompose 5 into two depending on the value of j :

$$\frac{\partial J}{\partial u_w} = \begin{cases} (\hat{y}_w - 1) \mathbf{v}_c, & \text{if } w = o \\ \hat{y}_w \mathbf{v}_c, & \text{otherwise} \end{cases} \quad (6)$$

Interpretation

Note that the update equation tell us that:

$$u_{ij}^{new} = u_{ij}^{old} - \eta \cdot e_j \cdot h_i \quad (7)$$

$$\text{Or equivalently: } \mathbf{u}_j^{new} = \mathbf{u}_j^{old} - \eta \cdot e_j \cdot \mathbf{h}$$

If $\hat{y}_j > y_j$ (*overestimating*):

Then we subtract a proportion of the hidden vector (or the center word embedding \mathbf{v}_c) from the output word embedding u_j , making u_j farther away from \mathbf{v}_c .

If $\hat{y}_j < y_j$ (*underestimating*) which is only true for the real target $y_j = 1$:

Then we subtract a proportion of the center word embedding \mathbf{v}_c from the output word embedding u_o , making the two embeddings for the same word closer to each other.

A Setup, Notation and Forward Pass

A.1 Setup

We have one-hot encoded pairs of words: (input,output), (center,outside).

The goal is to create a model which is able to assign maximum probability for the outside words of a center word, given their co-occurrence alongside the text (or corpus).

The matrix W projects words from the Vocabulary space into a smaller space of a chosen dimensionality N . That is the our embedding space.

The matrix W' projects back into the Vocabulary space. These are the **scores** after a forward pass.

These scores are the input of a softmax layer, which turns them into the **predicted probabilities** p or \hat{y} .

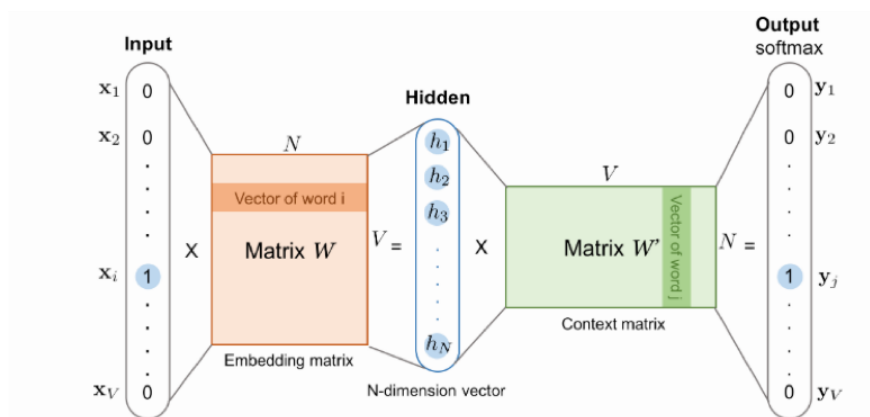


Figure 1: Skip-gram
[Liliang Weng's Blog](#)

A.2 Notation

There are many different ways to name the variables of this setup.

We will first move into the exercise notation the convention in Fig. 1, which matches with [this incredibly nice work](#) on the math details behind skip-gram:

Exercise	Figure 1	Description
x	—	One-hot encoded center (input) word
V	W^\top	Each column in V represents the input word embedding
U	W'	Each column in U represents the output word embedding
\mathbf{v}_c	W_{w_I}	Embedding for center (input) word with index C
\mathbf{u}_o	W'_{w_O}	Embedding for context (output) word with index O
\mathbf{s}_o	—	Score for word with index O . $\mathbf{s}_o = \mathbf{u}_o^\top \mathbf{v}_c$
$\hat{y}_o = \mathbf{p}_o$	—	Probability for word with index O . $\hat{y} = \text{softmax}(\mathbf{s})$

A.3 Forward Pass

Our model predicts the probability of the output word O given a center word C . $P(O|C)$. For each word in the text as a center word, and for each word lying on a window around this center word, we want our model to give a high probability.

That probability is the result of the forward pass of our neural setup for the skip gram.

Let's make explicit the forward pass by taking a look at the computing graph of our setup:

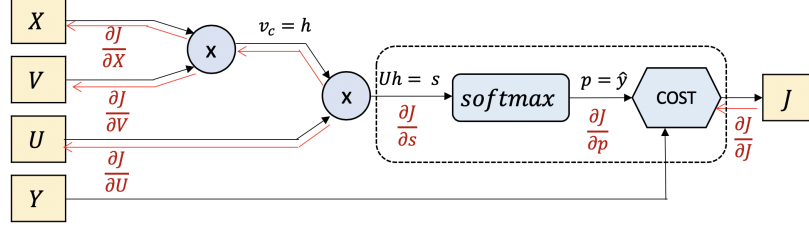


Figure 2: Skip-gram Computation Graph

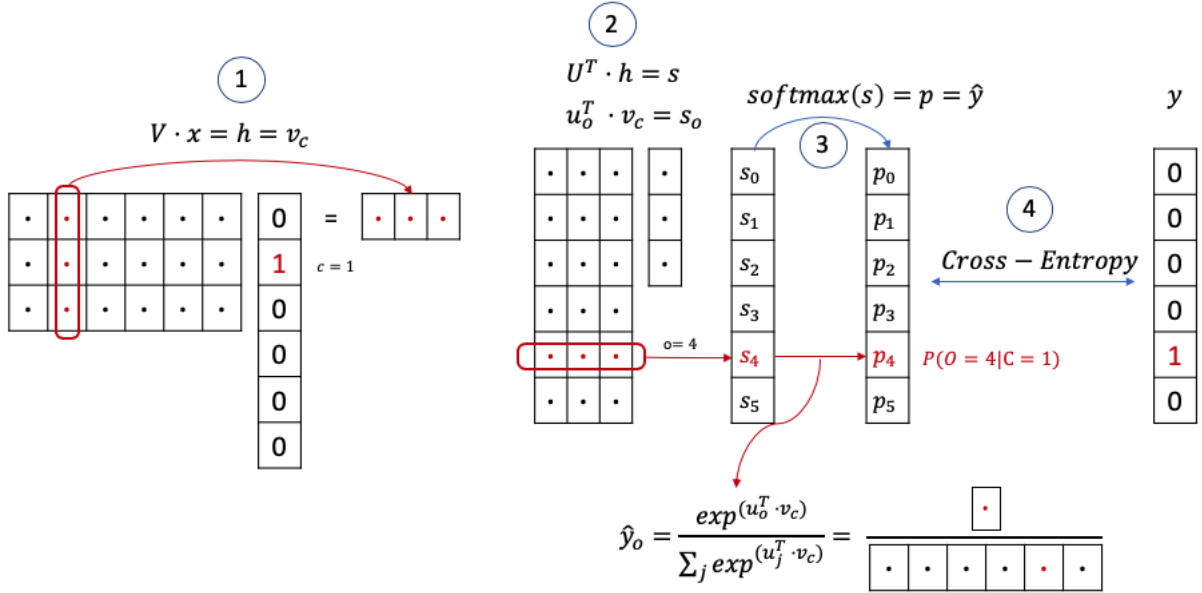


Figure 3: Skip-gram forward pass scheme for 1 prediction

$$\mathbf{h} = V\mathbf{x} := \mathbf{v}_c \quad (\text{F1})$$

$$\text{Or equivalently: } h_i = \sum_{j=1}^V x_i \cdot V_{ij}$$

$$\mathbf{s} = U^T \mathbf{v}_c \quad (\text{F2})$$

$$\text{Or equivalently: } s_i = \mathbf{u}_o^T \mathbf{v}_c =$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{s}) = \frac{\exp(\hat{\mathbf{s}}_i)}{\sum_w \exp(\hat{\mathbf{s}}_w)} = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)} \quad (\text{F3})$$

Therefore, a forward pass in our network setup outputs the probability of context (output) word to be in the context of a center (input) word:

$$\hat{y} = P(O|C) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)}$$

In order to maximize that probability, we phrase it in terms of minimizing a cost function.

That cost function is the negative log loss of the prediction, which is equivalent to the cross-entropy loss between the real distribution and the prediction, since the real distribution is a one-hot vector, yielding to the negative log likelihood for the prediction at the index of the real word.

A.4 Cost Function

$$\begin{aligned} J(\theta) &= J(u_c, o, U) = -\log \hat{y} = -\log \left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \right) \\ &= -\log(\exp(\mathbf{u}_o^\top \mathbf{v}_c)) + \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) \\ &= -\mathbf{u}_o^\top \mathbf{v}_c + \log \left(\sum_w \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right) \end{aligned} \tag{F4}$$

B Softmax & Cross-Entropy Backward Pass

In this section we derive mathematically the forward and backward passes of Softmax and Cross-Entropy:

B.1 FORWARD PASS

[1] **Softmax** \Rightarrow Transforms scores into probabilities

$$p(y_i = c|X) = p_i = \frac{\exp(\vec{s}_i)}{\sum_{j=1}^K \exp(\vec{s}_j)} \quad (8)$$

[2] **Maximum Likelihood** \Rightarrow How close are these probabilities to the ground truth?

We are interested in the derivative of the loss with respect to the parameters.

This derivative is the one we will use to update the parameters to perform better next time.

We are modeling a categorical distribution and we know that the maximum likelihood solution for a it is equivalent to minimize the negative log-likelihood of our data w.r.t. the parameters, which is equivalent to minimize the cross entropy between the true distribution y and the distribution generated by our model .

Let's derive it. The likelihood function for a categorical distribution is:

$$P(Y|X, \theta) \stackrel{iid}{=} \prod_{n=1}^N p(y_n|X_n, \theta) = \prod_{n=1}^N \prod_{k=1}^K p_{nk}^{y_{nk}} \quad (9)$$

Where N is the number of examples in the data set, and K the number of possible classes.

Taking the negative logarithm:

$$-\log p(y_i|X, \theta) = -\sum_{n=1}^N \sum_{k=1}^K y_{nk} \log p_{nk} \quad (10)$$

Where y_{nk} is a one-hot encoded vector, where every entry is zero except the k -th entry, representing the k class. That correspond to the cross-entropy.

B.2 BACKWARD PASS

We are looking for $\frac{\partial L}{\partial s}$ to back-propagate the error. We can apply the chain rule to decompose:

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial s} \frac{\partial s}{\partial \theta} = \frac{\partial L}{\partial s} \frac{\partial s}{\partial \theta} \quad (11)$$

where the first term is the derivative of the cross-entropy loss and the second term is the derivative of the softmax function. This is the reason why normally the derivative of the loss with respect to the scores is directly taken involving both softmax and cross-entropy together. This simplifies greatly the calculations

[1] **DERIVATIVE OF SOFTMAX**

$$\frac{\partial p_i}{\partial s_j} = \frac{\partial \left(\frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)} \right)}{\partial s_j} \quad (12)$$

To differentiate the numerator, we can use:

$$f(x) = \frac{g(x)}{h(x)} \Rightarrow f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{h^2(x)} \quad (13)$$

Therefore, we have different solutions if $i = j$ or not, that we need to derive independently:

- Solution for $i = j$

$$\frac{\partial p_i}{\partial s_j} = \frac{\exp(s_i) \sum_{j=1}^K \exp(s_j) - \exp(s_j) \exp(s_i)}{(\sum_{j=1}^K \exp(s_j))^2} = \frac{\exp(s_i) (\sum_{j=1}^K \exp(s_j) - \exp(s_j))}{\sum_{j=1}^K \exp(s_j) \sum_{j=1}^K \exp(s_j)} = -p_i(1 - p_j) \quad (14)$$

- Solution for $i \neq j$

$$\frac{\partial p_i}{\partial s_j} = \frac{0 \cdot \sum_{j=1}^K \exp(s_j) - \exp(s_j) \cdot \exp(s_i)}{(\sum_{j=1}^K \exp(s_j))^2} = \frac{\exp(s_j) \exp(s_i)}{\sum_{j=1}^K \exp(s_j) \sum_{j=1}^K \exp(s_j)} = -p_j(p_i) \quad (15)$$

where in both cases we have made use of Equation 12.

If we combine both solutions making use of Kronecker Delta:

$$\frac{\partial p_i}{\partial s_j} = p_i \cdot (\delta_{ij} - p_j) \quad (16)$$

[2] DERIVATIVE OF CROSS ENTROPY

We have derived already the first component of the chain rule, the local gradient. We need now to compute the upcoming gradient to back-propagate it into the network so we can update their weights accordingly to their contribution to the error.

Starting with the NLL cost function for a single prediction, Equation 10:

$$J(\theta|y_i, X) = - \sum_{j=1}^K y_j \log p_j \quad (17)$$

$$\frac{\partial J}{\partial s_i} = - \sum_{j=1}^K \frac{\partial (y_j \log p_j)}{\partial s_i} = - \sum_{j=1}^K y_j \frac{\partial (\log p_j)}{\partial s_i} = - \sum_{j=1}^K y_j \frac{1}{p_i} \frac{\partial (p_j)}{\partial s_i} \quad (18)$$

Substituting 16:

$$\frac{\partial J}{\partial s_i} = y_j \frac{1}{p_i} \cdot (p_i \cdot (\delta_{ij} - p_j)) = y_j (\delta_{ij} - p_j) \quad (19)$$

If we decompose the sum now on terms where $i = j$ and $i \neq j$:

$$\frac{\partial J}{\partial s_i} = -y_j(1 - p_i) - \sum_{j \neq i}^K y_i(0 - p_i) = \quad (20)$$

$$-y_j(1 - p_i) + \sum_{j \neq i}^K y_j p_i = \quad (21)$$

$$y_i p_i - y_i + \sum_{j \neq i}^K y_j p_i = \quad (22)$$

$$p_i(y_i + \sum_{j \neq i}^K y_j) - y_i \quad (23)$$

Note how the above purple selection is 1 as it is the sum of all possible values of a one-hot encoded vector.

$$\frac{\partial J}{\partial s_i} = p_i - y_i = \hat{y}_i - y_i \quad (24)$$