

# Training a Word Embedding Model

## Setup

Phrase: \_\_\_\_\_

# of words: 4

Context size: 1

# of dimensions: 2

Learning rate: 0.1

### Training pairs

Context words

(first) (centre) : (left) (right)

(second) : :

(third) : :

(fourth) : :

Random numbers  
between 0 and 2

### Centre word embeddings ( $v$ )

Word	Dimension 1 (X)	Dimension 2 (Y)

Random numbers  
between 0 and 2

### Context embeddings ( $u$ )

Word	Dimension 1 (X)	Dimension 2 (Y)

## Training

### 1. Choose Training Example

Pick one centre word and an associated context word.

(centre)      (context)

\_\_\_\_\_

Make sure that the words are next to each other in the beginning phrase!

Now, write down the initial embeddings for each word.

$v_w$  → (centre): \_\_\_\_\_

$u_c$  → (context): \_\_\_\_\_

Centre embeddings come from the first matrix, while context embeddings come from the second.

### 2. Compute Similarity

Next, take the **dot product** between the centre words and all adjacent context words.

$$\text{similarity (dot product)} = (v_{w\chi} * u_{c\chi}) + (v_{wy} * u_{cy})$$

Use the following formula, where  $u_{c\chi}$  and  $u_{cy}$  are the X and Y values of the context word, and  $v_{w\chi}$  and  $v_{wy}$  are the X and Y values of the centre word. Write the output below:

\_\_\_\_\_

If the word is on either of the edges, it will only have one adjacent context word.

Then, take the exponential of each of those values ( $e^{\text{similarity}}$ )

\_\_\_\_\_

### 3. Compute Softmax

Take the sum of the previously calculated values:

\_\_\_\_\_

And then divide the similarity values by the sum you just took:

\_\_\_\_\_

Congratulations! You now have your model's predicted probabilities, in decimal form. You can check the percentage predictions by multiplying that by 100.

### 4. Improve the Embeddings

First, we want to compute the loss. The loss measures how wrong the model's predictions are.

$$\text{Cross-entropy} = (p(\chi) * -\log(q(\chi)))$$

That looks very scary, but all it's doing is:

- taking the logarithm of the previous, decimal prediction
- multiplying it by 0 if it's not the chosen ("true) context, or multiplying it by 1, if it is the true context word.

The logarithm is used in this scenario, in order to punish incorrect answers harder if they're more sure of themselves. Perform the loss function on the softmax value belonging to the "true" context word.

\_\_\_\_\_

The next part of training is figuring out what changes need to be made to the embeddings. Use the following formula:

$$\text{true context embedding} - \left[ \left( \frac{\text{true context word softmax}}{\text{true context word softmax}} * \frac{\text{Context X 1}}{\text{Context X 1}} + \frac{\text{true context word softmax}}{\text{true context word softmax}} * \frac{\text{Context X 2}}{\text{Context X 2}} \right) \right]$$

Finally, multiply the centre embedding by the softmax of the true context word, multiply it by the learning rate(0.1), and add it to the embedding.

$$\text{New centre embedding} = \text{Centre embedding} + \left( \frac{\text{Context word softmax}}{\text{Context word softmax}} * \text{Centre embedding} * 0.1 \right)$$