

Word and Document Embeddings

Atanas Denkov¹

Abstract: TODO

Keywords: Machine Learning; Distributed Computing; Self-organizing Maps; Unsupervised Learning

¹ Hochschule für Technik und Wirtschaft, Angewandte Informatik, Wilhelminenhofstr., 12459 Berlin, Germany
s0559025@htw-berlin.de

1 Continuous Bag of Words

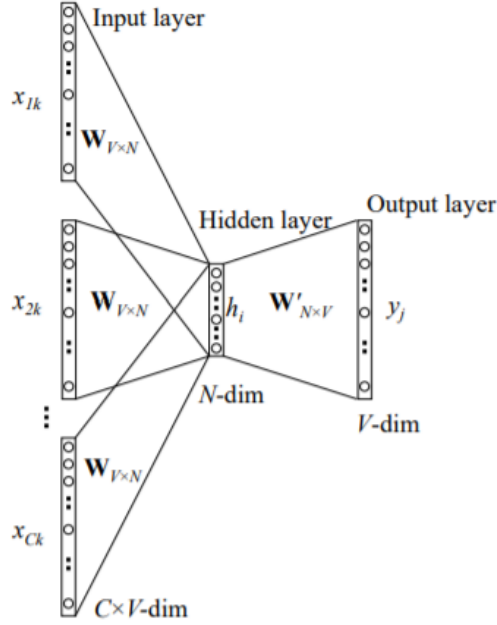


Fig. 1: The Continuous Bag of Words model

The Continuous Bag of Words model is a shallow neural network used for creating word vectors, also referred to as word *embeddings*. The set of word vectors forms a vector space that captures the similarity between words. The objective of the model is to predict the true center word w_o given the context words w_1, w_2, \dots, w_c that surround it. For example, in the sentence *The big brown fox jumped over the puddle*, we would like to predict the word *jumped* given the context words *[the, big, brown, over, the, puddle]*. The algorithm takes a probabilistic approach and tries to model the probability $p(w_o | w_1, w_2, \dots, w_c)$. The objective is then formulated as maximizing that probability.

1.1 Forward pass

The input to our model is a set of context words represented as one-hot encoded vectors x_1, x_2, \dots, x_c where $x_c \in \mathbb{R}^V$. V corresponds to the size of the vocabulary that the model can perceive. Each one-hot vector has a 1 at the index that is unique to the word it represents ($x_k = 1$) and 0's everywhere else ($\{x_{k'} = 0 | k' \in [1, 2, \dots, V], k' \neq k\}$). The neural network has a hidden matrix $\mathbf{W} \in \mathbb{R}^{V \times N}$ that holds a word vector $v_w \in \mathbb{R}^N$ for each word w when

it appears as a context word. Given that the context words are the input to our model, this matrix is called the *input matrix*. The hidden layer of the network functions as a lookup table to extract the input vectors v_1, v_2, \dots, v_c by computing the following dot products

$$\begin{aligned} v_1 &= \mathbf{W}^T x_1 \\ v_2 &= \mathbf{W}^T x_2 \\ &\vdots \\ v_c &= \mathbf{W}^T x_c. \end{aligned} \tag{1}$$

These vectors are then averaged to form the hidden state

$$\mathbf{h} = \frac{1}{C}(v_1 + v_2 + \dots + v_c), \mathbf{h} \in \mathbb{R}^N. \tag{2}$$

N denotes the size of the embedding space and is typically a hyperparameter. The model has an additional weight matrix $\mathbf{W}' \in \mathbb{R}^{N \times V}$ that holds a word vector $v'_w \in \mathbb{R}^N$ for each word w when it appears as a center word. We refer to this matrix as the *output matrix*. If we compute the dot product between the hidden state \mathbf{h} and a word vector v'_j , we obtain a score that measures the similarity between the context and the word represented by that word vector.

$$u_j = v'_j{}^T \mathbf{h} \tag{3}$$

We convert the scores into probabilities via the softmax function and obtain our prediction

$$p(w_j|w_I) = \hat{y}_j = \frac{\exp u_j}{\sum_{j'=1}^V \exp u_{j'}} = \frac{\exp v'_j{}^T \mathbf{h}}{\sum_{j'=1}^V \exp v'_{j'}{}^T \mathbf{h}}. \tag{4}$$

1.2 Backward pass

The objective of our model is to maximize the probability $p(w_j|w_I)$ defined in Equation 4 when w_j really is the center word given the context w_I . That is, we want to compare the target probability $y_j = 1$ with our estimate \hat{y}_j . We do this by using a simplified version of the cross-entropy loss function

$$\begin{aligned} H(\hat{y}_j, y_j) &= -y_j \log \hat{y}_j \\ &= -\log \hat{y}_j \\ &= -\log \frac{\exp u_j}{\sum_{j'=1}^V \exp u_{j'}} \end{aligned} \tag{5}$$

$$= -u_j + \log \sum_{j'=1}^V \exp u_{j'}. \tag{6}$$

The transition from (5) to (6) follows the logarithm rule for fractions

$$\log_b \left(\frac{A}{B} \right) = \log_b (A) - \log_b (B). \quad (7)$$

We use gradient descent to update v'_j and v_1, v_2, \dots, v_c such that $H(\hat{y}_j, y_j) \approx 0$. The chain rule tells us that

$$\frac{\partial H}{\partial v'_j} = \frac{\partial H}{\partial u_j} \cdot \frac{\partial u_j}{\partial v'_j}. \quad (8)$$

Calculating the derivative of the loss yields

$$\begin{aligned} \frac{\partial H}{\partial u_j} &= \frac{\partial}{\partial u_j} - u_j + \log \sum_{j'=1}^V \exp u_{j'} \\ &= \frac{\partial}{\partial u_j} - u_j + \frac{\partial}{\partial u_j} \log \sum_{j'=1}^V \exp u_{j'} \\ &= -1 + \frac{1}{\sum_{j'=1}^V \exp u_{j'}} \frac{\partial}{\partial u_j} \sum_{j'=1}^V \exp u_{j'} \\ &= -1 + \frac{1}{\sum_{j'=1}^V \exp u_{j'}} \frac{\partial}{\partial u_j} \exp u_j \\ &= -1 + \frac{\exp u_j}{\sum_{j'=1}^V \exp u_{j'}} \\ &= \hat{y}_j - 1. \end{aligned} \quad (9)$$

Next we take the derivative of u_j with respect to v'_j .

$$\frac{\partial u_j}{\partial v'_j} = \frac{\partial}{\partial v'_j} v'_j \cdot \mathbf{h} = \mathbf{h}. \quad (10)$$

By substituting our calculations in Equation 8 and considering that the index of the true center word is j , then our derivative is

$$\frac{\partial H}{\partial v'_{j'}} = \begin{cases} (\hat{y}_{j'} - 1) \cdot \mathbf{h}, & \text{if } j' = j \\ (\hat{y}_{j'}) \cdot \mathbf{h} & \text{otherwise.} \end{cases} \quad (11)$$

Since we are using gradient descent, our update rule for $v'_{j'}$ becomes

$$v'_{j'}^{(new)} = \begin{cases} v'_{j'}^{(old)} - \alpha \cdot ((\hat{y}_{j'} - 1) \cdot \mathbf{h}), & \text{if } j' = j \\ v'_{j'}^{(old)} - \alpha \cdot \hat{y}_{j'} \cdot \mathbf{h} & \text{otherwise.} \end{cases} \quad (12)$$

Note that when $j' = j$, the negative sign of the expression $\hat{y}_{j'} - 1$ cancels out the negative sign of the update rule, thus pushing the value of v'_j towards the context \mathbf{h} . Conversely,

when $j' \neq j$, the output vector is pushed away from the context. In other words, vectors that have similar contexts are clustered together by being pushed away from unrelated contexts and at the same time being pushed towards contexts that are meaningful. This notion is referred to as *distributional semantics*.

The next step is to compute the partial derivatives of the loss function with respect to the input vectors v_1, v_2, \dots, v_c . The chain rule tells us that

$$\frac{\partial H}{\partial v_i} = \sum_{j=1}^V \frac{\partial H}{\partial u_j} \cdot \frac{\partial u_j}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}}{\partial v_i}, \text{ for } i \in [1, 2, \dots, C]. \quad (13)$$

The first derivative in the summation has already been computed. The second one yields

$$\frac{\partial u_j}{\partial \mathbf{h}} = \frac{\partial}{\partial \mathbf{h}} v_j^T \mathbf{h} = v_j', \quad (14)$$

and the third one:

$$\frac{\partial \mathbf{h}}{\partial v_i} = \frac{\partial}{\partial v_i} \frac{1}{C} (v_1 + v_2 + \dots + v_c) = \frac{1}{C}. \quad (15)$$

Substituting (14) and (15) into (13) gives us

$$\frac{\partial H}{\partial v_i} = \begin{cases} \frac{1}{C} \cdot \sum_{j'=1}^V (\hat{y}_{j'} - 1) \cdot v_{j'}, & \text{if } j' = j \\ \frac{1}{C} \cdot \sum_{j'=1}^V \hat{y}_{j'} \cdot v_{j'} & \text{otherwise.} \end{cases} \quad \text{for } i \in [1, 2, \dots, C] \quad (16)$$

We now derive the update rule

$$v_i^{(new)} = \begin{cases} v_i^{(old)} - \alpha \cdot \frac{1}{C} \cdot \sum_{j'=1}^V (\hat{y}_{j'} - 1) \cdot v_{j'}, & \text{if } j' = j \\ v_i^{(old)} - \alpha \cdot \frac{1}{C} \cdot \sum_{j'=1}^V \hat{y}_{j'} \cdot v_{j'} & \text{otherwise.} \end{cases} \quad \text{for } i \in [1, 2, \dots, C] \quad (17)$$

1.3 Optimization

Notice that Equation 4 requires iterating over the whole vocabulary to compute the normalization factor. The vocabulary usually exceeds hundreds of thousands of words, which renders this computation rather inefficient. A technique called *negative sampling* deals with this by reformulating the objective function, the weights and the update equations. Instead of going over the whole vocabulary, we can sample a set of words from it that are unlikely to occur near the context and update their weights in such a way that they are

pushed further away from the context, while the true center word is pushed towards it. First, we remodel the objective function to be

$$H(\hat{y}_j, y_j) = -\log \sigma(u_j) - \sum_{k=1}^K \log \sigma(-u_k). \quad (18)$$

The words $\{u_k | k \in [1, 2, \dots, K]\}$ are sampled from a noise distribution $P(w)$, which is defined in terms of word frequencies. A good heuristic value for K is 10. Note that K is a lot smaller than V . If $f(w)$ defines the relative frequency of the word w in our corpus, then the probability of using w as a negative sample is defined by

$$P(w) = \frac{f(w)^{3/4}}{\sum_{i=1}^V f(w_i)^{3/4}}. \quad (19)$$

The power factor $3/4$ is determined heuristically and tends to favor words that appear less frequently.

Now that we have changed the objective, we need to recalculate the partial derivatives for our weights. Backpropagating the update signal to the output vector v'_j is done via the following derivative chain:

$$\frac{\partial H}{\partial v'_j} = \frac{\partial H}{\partial \sigma(u_j)} \cdot \frac{\partial \sigma(u_j)}{\partial u_j} \cdot \frac{\partial u_j}{\partial v'_j}. \quad (20)$$

The first derivative yields

$$\begin{aligned} \frac{\partial H}{\partial \sigma(u_j)} &= \frac{\partial}{\partial \sigma(u_j)} - \log \sigma(u_j) - \sum_{k=1}^K \log \sigma(-u_k) \\ &= \frac{\partial}{\partial \sigma(u_j)} - \log \sigma(u_j) - \frac{\partial}{\partial \sigma(u_j)} \sum_{k=1}^K \log \sigma(-u_k) \\ &= -\frac{1}{\sigma(u_j)} \frac{\partial}{\partial \sigma(u_j)} \sigma(u_j) \\ &= -\frac{1}{\sigma(u_j)}. \end{aligned} \quad (21)$$

The derivative of the sigmoid function is defined as

$$\frac{\partial \sigma(u_j)}{\partial u_j} = \sigma(u_j) \cdot (1 - \sigma(u_j)). \quad (22)$$

The third derivative has been calculated in Equation 10. Chaining these produces

$$\frac{\partial H}{\partial v'_{j'}} = \begin{cases} (\sigma(u_{j'}) - 1) \cdot \mathbf{h}, & \text{if } j' = j \\ \sigma(u_{j'}) \cdot \mathbf{h} & \text{otherwise.} \end{cases} \quad (23)$$

Note that we need to compute this just for the center word w_j and for all negative samples w_k instead of for the whole vocabulary. We substitute this equation in (12) to obtain the new update rule.

We now need to backpropagate the update signal to the context vectors v_1, v_2, \dots, v_c . Again, we determine the derivative chain and then compute each derivative. First, we define the set S as the intersection between the index of the current center word j and the indices of the negative samples k , $k \in [1, 2, \dots, K]$.

$$\frac{\partial H}{\partial v_i} = \sum_{j' \in S} \frac{\partial H}{\partial \sigma(u_{j'})} \cdot \frac{\partial \sigma(u_{j'})}{\partial u_{j'}} \cdot \frac{\partial u_{j'}}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}}{\partial v_i} \quad (24)$$

$$\frac{\partial H}{\partial v_i} = \begin{cases} \frac{1}{C} \cdot \sum_{j' \in S} \sigma(u_{j'}) - 1 \cdot v'_{j'}, & \text{if } j' = j \\ \frac{1}{C} \cdot \sum_{j' \in S} \sigma(u_{j'}) \cdot v'_{j'}, & \text{otherwise. } (j' = k) \end{cases} \quad (25)$$

We substitute the gradient in (17) to obtain the new update rule.

2 Distributed Memory Model of Paragraph Vectors

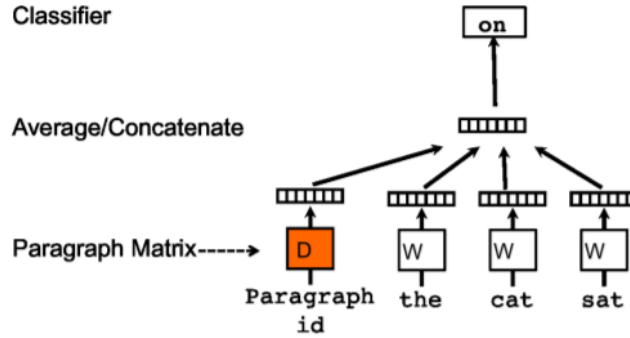


Fig. 2: The DM-PV model.

The distributed memory model of paragraph vectors, also referred to as DM-PV, attempts to generalize the abilities of CBOW to variable-length documents. In DM-PV, every document is mapped to a unique vector represented by a row in a matrix $\mathbf{D} \in \mathbb{R}^{M \times N}$, where M denotes the number of documents. As before, every word is also mapped to a vector in an input matrix $\mathbf{W} \in \mathbb{R}^{V \times N}$ and an output matrix $\mathbf{W}' \in \mathbb{R}^{N \times V}$, respectively. The only difference between DM-PV and CBOW is that the document vector d_j that contains the words is asked

to contribute to the prediction task by being integrated into the hidden state in Equation (2). See Figure 2 for a visual representation. Thus, the new hidden state is computed as follows:

$$\mathbf{h} = d_j + \frac{1}{C}(v_1 + v_2 + \dots + v_c), \mathbf{h} \in \mathbb{R}^N. \quad (26)$$

The document vector can be thought of as another word. It acts as a memory that remembers what is missing from the current context, i.e the document's topic. Note that the document vector is shared across all contexts $v_1 + v_2 + \dots + v_c$ generated from the same document, but not across documents. Conversely, the word vector matrices \mathbf{W} and \mathbf{W}' are shared across documents. During backpropagation, the derivative of the cost function with respect to the target document vector

$$\frac{\partial H}{\partial d_j} \quad (27)$$

is computed as well, thus pushing documents with similar words towards those words and therefore, towards each other.

2.1 Inference

So far, we've only considered the case of training the algorithm in order to create numeric representations for **previously seen** documents. But how do we generate such representations for unseen documents? This problem is addressed in the **inference stage**.

Given is an unseen document d . A randomly-initialized vector is added to our model's weight matrix \mathbf{D} . That vector represents the document vector for d . The vector for d is updated by gradient descending on \mathbf{D} for i epochs while keeping \mathbf{W} , \mathbf{W}' and their gradients fixed. The number of epochs to descent on d should match the number of epochs used during training.