

# Word Embedding

L.J. Brown

September 30, 2017

## 1 Introduction

This paper outlines a few methods used in this repository to map all unique words found in a corpus to vectors in a vector space. The idea, and hope, is that some relationships between words found in the corpus will be preserved through this mapping and will manifest as characteristics of the word vectors. Successful implementations and strides in this area include [Latent Semantic Analysis](#), ["Word2vec"](#), and ["GloVe"](#). At its heart, the goal of word embedding is to quantify the conceptual relationships that exist within our vocabulary and to create a mathematical system to describe and analyze them. Abstractly word embedding seeks to build a sort of concept space where different regions and axes can represent ideas, and to embed word vectors inside this space where their lengths and directions place them in meaningful locations. Models search for word vectors whose distances apart represent essential relationships between them. The key idea (and assumption) that is exploited by these models is that the more frequently words appear together in natural language the more related these words are conceptually. Most methods exploit this property by first counting the cooccurrences of word pairs and then by building word vectors accordingly. Models search for vectors whose inner product or cosine similarity or euclidean distance apart can be traced back to this list of cooccurrences. Many different techniques have been employed to find these word vectors, but the power of these representations is probably best illustrated by a commonly cited example from Word2vec's results:

$$\vec{king} - \vec{queen} + \vec{woman} = \vec{man}$$

## 2 General Strategy

The essential steps taken by each implementation within this repository are as follows:

1. **Construct a Cooccurrence Matrix**

Build a square "Cooccurrence Matrix",  $\mathbf{A}$ , where each element roughly represents some function,  $f$ , of the frequency or probability that two

words,  $w_i$  and  $w_j$ , occur together in the corpus (the diagonal elements being associated with the same word twice).

$$\mathbf{A}_{n \times n} = \begin{matrix} & \begin{matrix} w_1 & w_2 & \dots & w_n \end{matrix} \\ \begin{matrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{matrix} & \begin{pmatrix} f(w_1, w_1) & f(w_1, w_2) & \dots & f(w_1, w_n) \\ f(w_2, w_1) & f(w_2, w_2) & \dots & f(w_2, w_n) \\ \vdots & \vdots & \ddots & \vdots \\ f(w_n, w_1) & f(w_n, w_2) & \dots & f(w_n, w_n) \end{pmatrix} \end{matrix}$$

## 2. Decompose the Cooccurrence Matrix into unique Word Vectors

Search for word vectors with the soft constraint that given any word vector pair, their inner product will yield a value close to the two values in the co-occurrence matrix associated with those two words. This constraint can be written as

$$\vec{w}_i \cdot \vec{w}_j^T \approx f(w_i, w_j)$$

If the desired word vectors,  $(\vec{w}_1, \dots, \vec{w}_n)$ , are written as the rows of some "Word Vector Matrix",  $\mathbf{W}$ , then this constraint can be written as

$$\mathbf{W}\mathbf{W}^T \approx \mathbf{A}$$

Several methods to perform this decomposition step are explored within this repository.

## 3 Cooccurrence Matrix Decomposition Methods

The methods implimented in this repository to decompose the Cooccurrence Matrix,  $\mathbf{A}$ , into  $\mathbf{W}\mathbf{W}^T$  are

### 1. Stochastic Gradient Descent

An iterative optimization method used to minimize an objective function. Stochastic Gradient Descent (SGD) is often used over Batch Gradient Decent (BGD) for large noisy datasets. Our implementation of SGD Draws heavily on the implementations by "Word2vec" and "GloVe".

### 2. Eigen Decomposition of a forced Symmetric Positive Matrix

First the Cooccurrence Matrix,  $\mathbf{A}$ , is updated to have positive Eigen values. This is achieved by modifying the unused diagonal entries to ensure diagonal dominance. The new real symmetric matrix,  $\hat{\mathbf{A}}$ , can be written as  $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ , and because of the positive Eigen values,  $\mathbf{W} = \mathbf{Q}\sqrt{\mathbf{\Lambda}}$  will satisfy the constraint.

### 3. Singular Value Decomposition of a forced Symmetric Positive Definite Matrix

Similarly to the EVD method, the Cooccurrence Matrix,  $\mathbf{A}$ , is updated, but now to be positive definite. This is again achieved by modifying the unused diagonal entries to ensure diagonal dominance with the additional constraint that all diagonal entries are greater than zero. The new positive definite matrix,  $\tilde{\mathbf{A}}$ , can be written as  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , and because it is symmetric and positive definite,  $\mathbf{U} = \mathbf{V}$ , and  $\mathbf{W} = \mathbf{V}\sqrt{\mathbf{\Sigma}}$  will satisfy the constraint.

## 4 Details of Cooccurrence Matrix

This section details how this repository defines its  $n \times n$  cooccurrence matrix,  $\mathbf{A}_{n \times n}$ , whose rows and columns correspond to the set of unique words,  $W$ , and whose elements are the output of a cooccurrence function,  $f(w_i, w_j)$ , which takes words as inputs.

For clarity it is useful to define a function,  $m(w_i) = i$ , that converts unique words to unique integers in the range 1 to  $n$  and its inverse,  $m^{-1}(i) = w_i$ , which does the opposite.

$n \equiv$  Number of unique words in corpus,  $W \equiv$  Set of unique words in corpus

$$W = \{w_1, \dots, w_n\}, I = \{1, \dots, n\}$$

$$m : W \rightarrow I, m^{-1} : I \leftarrow W$$

$f(w_i, w_j) \equiv$  Cooccurrence function defined below

$$\mathbf{A}_{n \times n} = \begin{matrix} & \begin{matrix} c_1 & c_2 & \dots & c_n \end{matrix} \\ \begin{matrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{matrix} & \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} \end{matrix}$$

$$\text{where } a_{ij} = f(m^{-1}(r_i), m^{-1}(c_j)) = f(w_i, w_j)$$

For each element,  $a_{ij}$ , the words used as inputs for  $f$  are the words whose integer mappings correspond to the elements row,  $r_i$ , and column,  $c_j$ .

In order to define the cooccurrence function,  $f$ , it is useful to introduce some new variables. For each of the  $m$  sentences in the corpus we define  $s_i$  as the sequence of words in that particular sentence. Where  $S$  is a **sequence** containing the  $m$  **sequences**,  $(s_1, \dots, s_m)$ , corresponding to the  $m$  sentences of the corpus.

$$S = (s_1, \dots, s_m)$$

For example if the  $i^{\text{th}}$  sentence in the corpus is: “God made mud.”, then the corresponding sequence,  $s_i$ , in  $S$  would be

$$s_i = (\tilde{w}_1, \tilde{w}_2, \tilde{w}_3)$$

where  $\tilde{w}_1, \tilde{w}_2, \tilde{w}_3 \in W$

Also used in the construction of the concurrence matrix is a distance function,  $d(s_i, \{\tilde{w}_j, \tilde{w}_k\})$ , whose parameters are a sequence of words,  $s_i$ , in  $S$  (corresponding to a sentence of the corpus) and a set or an unordered pair of words that are members of the sequence  $s_i$ . For example using the  $i^{\text{th}}$  sentence in the corpus again and a random pair from that sequence  $\{\tilde{w}_2, \tilde{w}_3\} = \{\text{"made"}, \text{"mud"}\}$

$$d(s_i, \{\tilde{w}_2, \tilde{w}_3\}) = d(\{\text{"God"}, \text{"made"}, \text{"mud"}\}, \{\text{"made"}, \text{"mud"}\})$$

It is important to note that although  $\tilde{w}_2, \tilde{w}_3 \in W$

it is not necessarily true that  $\tilde{w}_2 = w_2$  or  $\tilde{w}_3 = w_3$ .

The distance function,  $d(s_i, \{\tilde{w}_j, \tilde{w}_k\})$ , returns a value one more than the number of words between the word pair,  $\{\tilde{w}_j, \tilde{w}_k\}$ , in sentence corresponding to the sequence,  $s_i$ .

$$d(s_i, \{\tilde{w}_j, \tilde{w}_k\}) = |j - k|$$

One more useful definition due to obscure notation,  $\{[s_i]^2\}$  is the set of all unordered pairs of words in the sequence  $s_i$ . Using the example above for  $s_i$

$$\{[s_i]^2\} = \{\{\text{"God"}, \text{"made"}\}, \{\text{"God"}, \text{"mud"}\}, \{\text{"made"}, \text{"mud"}\}\}$$

Finally the cooccurrence function,  $f(w_i, w_j)$ , is defined as

$$f(w_i, w_j) = \sum_{s \in S} \sum_{p \in \{[s]^2\}} \begin{cases} \log\left(\frac{1}{d(s,p)}\right), & \text{if } w_i, w_j \in p \text{ and if } w_i \neq w_j \\ 0, & \text{otherwise} \end{cases}$$

Stated simply,  $f(w_i, w_j)$  finds all times that the words  $w_i$  and  $w_j$  appear together in sentences of the corpus, and sums the  $\log\left(\frac{1}{\text{their distance apart}}\right)$ . Note that the value of the cooccurrence function,  $f(w_i, w_j)$ , when the two inputs are the same word is zero. This corresponds to the diagonal entries of the cooccurrence matrix  $\mathbf{A}$ . The matrix  $\mathbf{A}$  will also be symmetric due to the property of  $f(w_i, w_j)$  that,

$$f(w_i, w_j) = f(w_j, w_i)$$

therefore,

$$a_{ij} = a_{ji},$$

$$\mathbf{A} = \mathbf{A}^T$$

## 5 Details of Word Vectors and Word Vector Matrix

**Defining Constraints** The soft constraints imposed states that for any given word vector pair,  $\vec{w}_i$  and  $\vec{w}_j$ , excluding duplicate pairs,  $\vec{w}_i$  and  $\vec{w}_i$ , their inner product will yield a value close to the two values in the co-occurrence matrix associated with those two words. This constraint can be written as

$$\vec{w}_i \cdot \vec{w}_j^T : \Leftrightarrow f(w_i, w_j) = a_{ij}$$

Since  $\mathbf{A}$  is symmetric order does not matter when taking the inner product of two word vectors

$$\vec{w}_i^T \cdot \vec{w}_j = \vec{w}_j^T \cdot \vec{w}_i$$

Therefore,

$$\vec{w}_i \cdot \vec{w}_j^T = \vec{w}_i^T \cdot \vec{w}_j : \Leftrightarrow a_{ij} = a_{ji}$$

If the desired word vectors,  $(\vec{w}_1, \dots, \vec{w}_n)$ , are written as the rows of some "Word Vector Matrix",  $\mathbf{W}$ , then this constraint can be written as

$$\mathbf{W}\mathbf{W}^T : \Leftrightarrow \mathbf{A}$$

where  $\mathbf{W}$  is written as,

$$\mathbf{W}_{n \times v} = \begin{pmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_n \end{pmatrix}$$

and where  $v$  is the dimensionality of the word vectors (not yet specified).

## 6 Decomposition using Stochastic Gradient Descent

An iterative optimization method used to minimize an objective function. Stochastic Gradient Descent (SGD) is often used over Batch Gradient Descent (BGD) for large noisy datasets. Our implementation of SGD Draws heavily on the implementations by "Word2vec" and "GloVe".

The definition of an error function,  $\varepsilon$ , between two word vectors,  $\vec{w}_i$  and  $\vec{w}_j$ , based on the soft constraints chosen is

$$\varepsilon(i, j) = \vec{w}_i^T \cdot \vec{w}_j - a_{ij}$$

The chosen objective function,  $J$ , to minimize during the optimization process is

$$J = \frac{1}{\binom{n}{2}} \sum_{r=1}^n \sum_{c=r+1}^n \begin{cases} \varepsilon(r, c)^2, & \text{if } r \neq c \\ 0, & \text{otherwise} \end{cases}$$

where  $\frac{n^2 - n}{2} = \binom{n}{2}$  is the number of upper-triangular non-diagonal elements in the matrix  $\mathbf{A}$ .

When  $J$  is at a minimum then the  $\mathbf{W}$  will best meet our soft constraint.

Next is the definition an individual objective function for each word vector,  $J_{\vec{w}_i}$ . During the optimization process, and during each iteration, we randomly select a column vector of  $\mathbf{W}$ ,  $\vec{w}_i$ , to update by computing the gradient of the chosen individual word vector,  $\frac{\partial J_{\vec{w}_i}}{\partial \vec{w}_i}$ . We define the objective function for an individual word vector,  $\vec{w}_i$ , as

$$J_{\vec{w}_i} = \frac{1}{2} \sum_{j=1}^n \begin{cases} \varepsilon(i, j)^2, & \text{if } i \neq j \\ 0, & \text{otherwise} \end{cases}$$

And compute the gradient of that specific objective function,  $\frac{\partial J_{\vec{w}_i}}{\partial \vec{w}_i}$ , as

$$\frac{\partial J_{\vec{w}_i}}{\partial \vec{w}_i} = \sum_{j=1}^n \begin{cases} \vec{w}_i \odot (\vec{w}_i^T \cdot \vec{w}_j - y_{ij}), & \text{if } i \neq j \\ 0, & \text{otherwise} \end{cases}$$

The optimization process in pseudo code is outlined bellow:

$n \equiv$  Number of unique words in corpus

$m \equiv$  Number of training iterations

$\eta \equiv$  Learning rate

For  $i = 1, 2, \dots, m$ , do:

$x :=$  random number range[1,  $n$ ]

$$\vec{w}_x := \vec{w}_x - \eta \frac{\partial J_{\vec{w}_x}}{\partial \vec{w}_x}$$

## 7 Eigen Decomposition of a forced Symmetric Positive Matrix

Our goal is to decompose the square, real, symmetric, cooccurrence matrix  $\mathbf{A}_{n \times n}$  into a word vector matrix,  $\mathbf{W}$ , multiplied by its transpose. This would satisfies the conditions that the inner product between any two word vectors is equal to the two elements in matrix  $\mathbf{A}$  corresponding to the two words intersection.

$$\mathbf{A} = \mathbf{W}\mathbf{W}^T$$

Real symmetric matrices can be decomposed into the form

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$$

Where  $\mathbf{Q}$  is an orthogonal matrix and  $\mathbf{\Lambda}$  is a diagonal matrix whose entries are the eigenvalues of  $\mathbf{A}$ . If the eigenvalues of  $\mathbf{A}$  are all positive then we can write as

$$\mathbf{W} = \mathbf{Q}\sqrt{\mathbf{\Lambda}}$$

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T = \mathbf{Q}\sqrt{\mathbf{\Lambda}}\sqrt{\mathbf{\Lambda}}^T\mathbf{Q}^T = \mathbf{Q}\sqrt{\mathbf{\Lambda}}\left(\mathbf{Q}\sqrt{\mathbf{\Lambda}}\right)^T = \mathbf{W}\mathbf{W}^T$$

$$\mathbf{\Lambda} = \mathbf{\Lambda}^T \text{ and } \sqrt{\mathbf{\Lambda}} = \sqrt{\mathbf{\Lambda}}^T$$

However, it is not guaranteed that the eigenvalues of  $\mathbf{A}$  will all be positive. But we can force this to be the case by making use of the free diagonals of the cooccurrence matrix  $\mathbf{A}$  and making it diagonally dominant.

A matrix is diagonally dominant if  $|\mathbf{a}_{ii}| \geq \sum_{j \neq i} |\mathbf{a}_{ij}|$  for all  $i$ .

## 8 Singular Value Decomposition of a forced Symmetric Positive Definite Matrix

If we add the condition that  $\mathbf{a}_{ii} > 0$  for all  $i$ , then this matrix will also be positive definite and we can use singular value decomposition to compute  $\tilde{\mathbf{A}} = \mathbf{W}\mathbf{W}^T$ . If  $\tilde{\mathbf{A}}$  is a symmetric positive definite matrix then by the spectral theorem we know that  $\mathbf{\Sigma} = \mathbf{\Lambda}$  and  $\mathbf{U} = \mathbf{V} = \mathbf{Q}$ .

$$\tilde{\mathbf{A}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

$$\tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T, \text{ therefore } \mathbf{U} = \mathbf{V}$$

And we can write  $\tilde{\mathbf{A}}$  as  $\mathbf{W}\mathbf{W}^T$  where  $\mathbf{W} = \mathbf{V}\sqrt{\mathbf{\Lambda}}$