

# Word2Vec Tutorial Part 2 - Negative Sampling

11 Jan 2017

In part 2 of the word2vec tutorial (here's [part 1](#)), I'll cover a few additional modifications to the basic skip-gram model which are important for actually [making it feasible to train](#).

When you read the tutorial on the skip-gram model for Word2Vec, you may have noticed something—it's a huge neural network!

In the example I gave, we had word vectors with 300 components, and a vocabulary of 10,000 words. Recall that the neural network had two weight matrices—a hidden layer and output layer. Both of these layers would have a weight matrix with  $300 \times 10,000 = 3$  million weights each!

Running gradient descent on a neural network that large is going to be slow. And to make matters worse, you [need a huge amount of training data](#) in order to [tune](#) that many weights and [avoid over-fitting](#). millions of weights times billions of training samples means that training this model is going to be a beast.

The authors of Word2Vec addressed these issues in their second [paper](#).

There are three innovations in this second paper:

1. Treating common word pairs or phrases as single "words" in their model.
2. Subsampling frequent words to decrease the number of training examples.
3. Modifying the optimization objective with a technique they called "Negative Sampling", which causes each training sample to update only [a small percentage of the model's weights](#).

It's worth noting that [subsampling frequent words](#) and applying Negative Sampling not only reduced the compute burden of the training process, but [also improved the quality of their resulting word vectors as well](#).

# Word Pairs and "Phrases"

The authors pointed out that a word pair like "Boston Globe" (a newspaper) has a much different meaning than the individual words "Boston" and "Globe". So it makes sense to treat "Boston Globe", wherever it occurs in the text, as a single word with its own word vector representation.

You can see the results in their published model, which was trained on 100 billion words from a Google News dataset. The addition of phrases to the model **swelled** the vocabulary size to 3 million words!

If you're interested in their resulting vocabulary, I poked around it a bit and published a post on it [here](#). You can also just browse their vocabulary [here](#).

**Phrase detection** is covered in the "Learning Phrases" section of their [paper](#). They shared their implementation in `word2phrase.c`—I've shared a commented (but otherwise unaltered) copy of this code [here](#).

I don't think their phrase detection approach is a key contribution of their paper, but I'll share a little about it anyway since it's pretty straightforward.

Each pass of their tool only looks at combinations of 2 words, but you can run it multiple times to get longer phrases. **So, the first pass will pick up the phrase "New\_York", and then running it again will pick up "New\_York\_City" as a combination of "New\_York" and "City".**

The tool counts the number of times each combination of two words appears in the training text, and then these counts are used in an equation to determine which word combinations to turn into phrases. The equation is designed **to make phrases out of words which occur together often relative to the number of individual occurrences**. It also favors phrases made of infrequent words in order to **avoid making phrases out of common words** like "and the" or "this is".

You can see more details about their equation in my code comments [here](#).

One thought I had for an alternate phrase recognition strategy would be to **use the titles of all Wikipedia articles** as your vocabulary.

## Subsampling Frequent Words

In part 1 of this tutorial, I showed how training samples were created from the source text, but I'll repeat it here. The below example shows some of the training samples (word pairs) we would take from the sentence "The quick brown fox jumps over the lazy dog." I've used a small window size of 2 just for the example. The word highlighted in blue is the input word.

Source Text	Training Samples			
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)
The	quick	brown		
The <table><tr><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)
quick	brown	fox		
The quick <table><tr><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
brown	fox	jumps		
The quick brown <table><tr><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
fox	jumps	over		

There are two "problems" with common words like "the":

1. When looking at word pairs, ("fox", "the") doesn't tell us much about the meaning of "fox". "the" appears in the context of pretty much every word.
2. We will have many more samples of ("the", ...) than we need to learn a good vector for "the".

Word2Vec implements a "subsampling" scheme to address this. For each word we encounter in our training text, there is a chance that we will effectively delete it from the text. The probability that we cut the word is related to the word's frequency.

If we have a window size of 10, and we remove a specific instance of "the" from our text:

1. As we train on the remaining words, "the" will not appear in any of their context windows.
2. We'll have 10 fewer training samples where "the" is the input word.

Note how these two effects help address the two problems stated above.

# Sampling rate

The word2vec C code implements an equation for calculating a probability with which to keep a given word in the vocabulary.

$w_i$  is the word,  $z(w_i)$  is the fraction of the total words in the corpus that are that word. For example, if the word "peanut" occurs 1,000 times in a 1 billion word corpus, then  $z(\text{'peanut'}) = 1\text{E-}6$ .

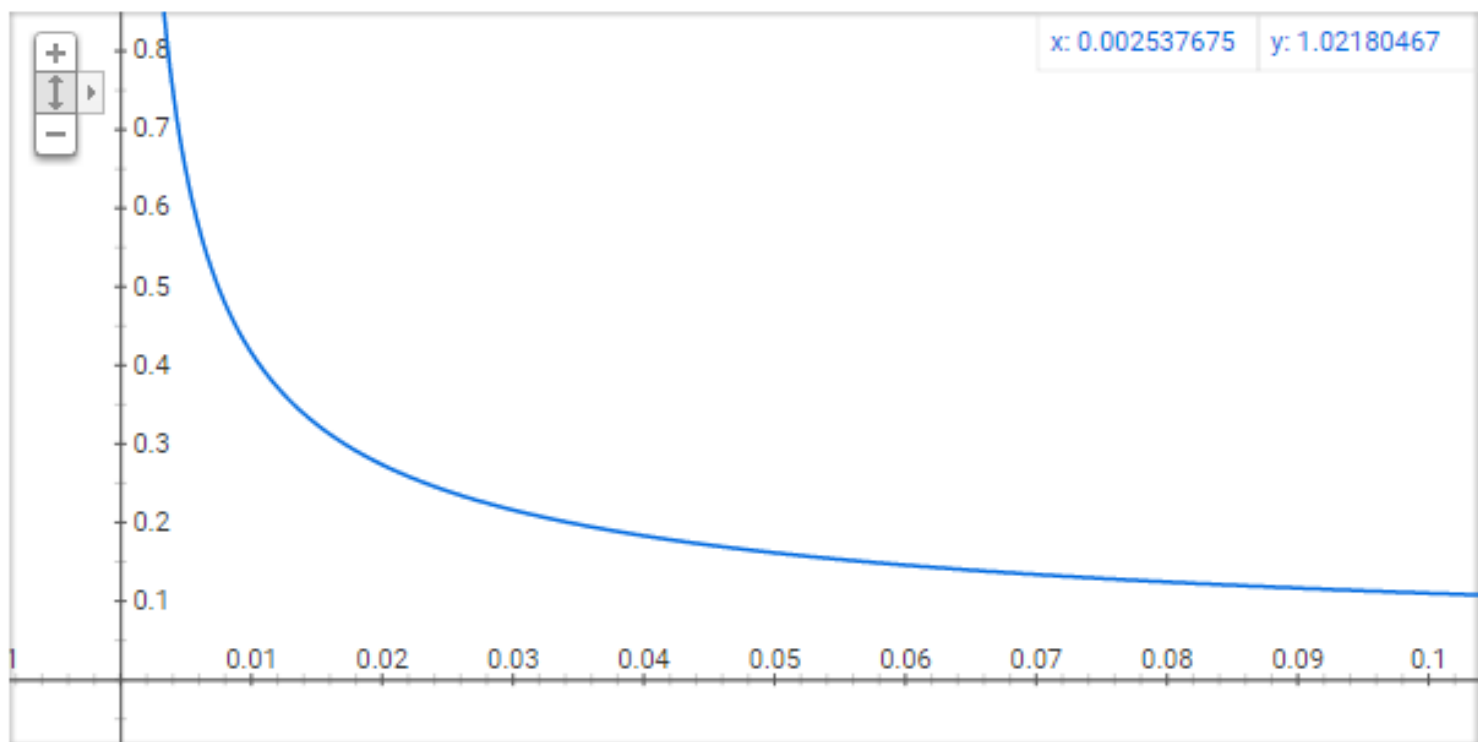
There is also a parameter in the code named 'sample' which controls how much subsampling occurs, and the default value is 0.001. Smaller values of 'sample' mean words are less likely to be kept.

$P(w_i)$  is the probability of keeping the word:

$$P(w_i) = \left( \sqrt{\frac{z(w_i)}{0.001}} + 1 \right) \cdot \frac{0.001}{z(w_i)}$$

You can plot this quickly in Google to see the shape.

Graph for  $(\sqrt{x/0.001}+1)*0.001/x$



No single word should be a very large percentage of the corpus, so we want to look at pretty small values on the x-axis.

Here are some interesting points in this function (again this is using the default sample value of 0.001).

- $P(w_i) = 1.0$  (100% chance of being kept) when  $z(w_i) \leq 0.0026$ 
  - This means that only words which represent more than 0.26% of the total words will be subsampled.

- $P(w_i) = 0.5$  (50% chance of being kept) when  $z(w_i) = 0.00746$
- $P(w_i) = 0.033$  (3.3% chance of being kept) when  $z(w_i) = 1.0$ .
  - That is, if the corpus consisted entirely of word  $w_i$ , which of course is ridiculous.

You may notice that the paper defines this function a little differently than what's implemented in the C code, but I figure the C implementation is the more authoritative version.

## Negative Sampling

Training a neural network means taking a training example and adjusting all of the neuron weights slightly so that it predicts that training sample more accurately. In other words, each training sample will tweak *all* of the weights in the neural network.

As we discussed above, the size of our word vocabulary means that our skip-gram neural network has a tremendous number of weights, all of which would be updated slightly by every one of our billions of training samples!

Negative sampling addresses this by having each training sample only modify a small percentage of the weights, rather than all of them. Here's how it works.

When training the network on the word pair ("fox", "quick"), recall that the "label" or "correct output" of the network is a one-hot vector. That is, for the output neuron corresponding to "quick" to output a 1, and for *all* of the other thousands of output neurons to output a 0.

With negative sampling, we are instead going to randomly select just a small number of "negative" words (let's say 5) to update the weights for. (In this context, a "negative" word is one for which we want the network to output a 0 for). We will also still update the weights for our "positive" word (which is the word "quick" in our current example).

The paper says that selecting 5-20 words works well for smaller datasets, and you can get away with only 2-5 words for large datasets.

Recall that the output layer of our model has a weight matrix that's 300 x 10,000. So we will just be updating the weights for our positive word

("quick"), plus the weights for 5 other words that we want to output 0. That's a total of 6 output neurons, and 1,800 weight values total. That's only 0.06% of the 3M weights in the output layer!

In the hidden layer, only the weights for the input word are updated (this is true whether you're using Negative Sampling or not).

## Selecting Negative Samples

The "negative samples" (that is, the 5 output words that we'll train to output 0) are chosen using a "unigram distribution".

Essentially, the probability for selecting a word as a negative sample is related to its frequency, with more frequent words being more likely to be selected as negative samples.

In the word2vec C implementation, you can see the equation for this probability. Each word is given a weight equal to its frequency (word count) raised to the 3/4 power. The probability for selecting a word is just its weight divided by the sum of weights for all words.

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n (f(w_j)^{3/4})}$$

The decision to raise the frequency to the 3/4 power appears to be empirical; in their paper they say it outperformed other functions. You can look at the shape of the function—just type this into Google: "plot  $y = x^{3/4}$  and  $y = x$ " and then zoom in on the range  $x = [0, 1]$ . It has a slight curve that increases the value a little.

The way this selection is implemented in the C code is interesting. They have a large array with 100M elements (which they refer to as the unigram table). They fill this table with the index of each word in the vocabulary multiple times, and the number of times a word's index appears in the table is given by  $P(w_i) * \text{table\_size}$ . Then, to actually select a negative sample, you just generate a random integer between 0 and 100M, and use the word at that index in the table. Since the higher probability words occur more times in the table, you're more likely to pick those.

## Other Resources

For the most detailed and accurate explanation of word2vec, you should check out the C code. I've published an extensively commented (but

otherwise unaltered) version of the code [here](#).

I've also created a [post](#) with links to and descriptions of other word2vec tutorials, papers, and implementations.

## Cite

McCormick, C. (2017, January 11). *Word2Vec Tutorial Part 2 - Negative Sampling*. Retrieved from <http://www.mccormickml.com>

---

## Related posts

[Concept Search on Wikipedia](#) 22 Feb 2017

[Getting Started with mlpack](#) 01 Feb 2017

[DBSCAN Clustering](#) 08 Nov 2016

---

© 2017. All rights reserved.