

# Efficient Calculation of Bigram Frequencies in a Corpus of Short Texts

Melvyn Drag, Gauthaman Vasudevan

*Avlino, Inc. Holdmel, NJ*

---

## Abstract

We show that an efficient and popular method for calculating bigram frequencies is unsuitable for bodies of short texts and offer a simple alternative. Our method has the same computational complexity as the old method and offers an exact count instead of an approximation.

*Keywords:* NLP, Machine Learning

---

## 1. Acknowledgements

This short note is the result of a brief conversation between the authors and Joel Nothman. We came across a potential problem, he gave a sketch of a fix, and we worked out the details of a solution.

## 2. Calculating Bigram Frequencies

A common task in natural language processing is to find the most frequently occurring word pairs in a text(s) in the expectation that these pairs will shed some light on the main ideas of the text, or offer insight into the structure of the language. One might be interested in pairings of adjacent words, but in some cases one is also interested in pairs of words in some small neighborhood. The neighborhood is usually referred to as a window, and to illustrate the concept consider the following text and bigram set:

Text: "I like kitties and doggies"  
 Window: 2  
 Bigrams: {(I like), (like kitties), (kitties and), (and doggies)}

and this one:

Text: "I like kitties and doggies"  
 Window: 4  
 Bigrams: {(I like), (I kitties), (I and), (like kitties), (like and), (like doggies), (kitties and), (kitties doggies), (and doggies)}.

## 3. The Popular Approximation

Bigram frequencies are often calculated using the approximation

$$freq(*, word) = freq(word, *) = freq(word) \quad (1)$$

In a much cited paper, Church and Hanks [2] use '=' in place of '≈' because the approximation is so good. Indeed, this approximation will only cause errors for the very few words which occur near the beginning or the end of the text.

Take for example the text appearing above - the bigram (doggies, \*) does not occur once, but the approximation says it does.

An efficient method for computing the contingency matrix for a bigram (word1, word2) is suggested by the approximation. Store  $freq(w1, w2)$  for all bigrams (w1, w2) and the frequencies of all words. Then,

- $freq(word1, word2)$  is known,
- $freq(\sim word1, word2) \approx freq(word2) - freq(word1, word2)$ ,
- $freq(word1, \sim word2) \approx freq(word1) - freq(word1, word2)$ ,
- and  $freq(\sim word1, \sim word2)$  is easily computed.

The statistical importance of miscalculations due to this method diminishes as our text grows larger and larger. Interest is growing in the analysis of small texts, however, and a means of computing bigrams for this type of corpus must be employed. This approximation is implemented in popular NLP libraries and can be seen in many tutorials across the internet. People who use this code, or write their own software, must know when it is appropriate.

#### 4. An Alternative Method

We propose an alternative. As before, store the frequencies of words and the frequencies of bigrams, but this time store two additional maps called **too\_far\_left** and **too\_far\_right**, of the form {word : list of offending indices of word}. The offending indices are those that are either too far to the left or too far to the right for approximation (1) to hold. All four of these structures are built during the construction of a bigram finder, and do not cripple performance when computing statistical measures since maps are queried in  $O(1)$  time.

As an example of the contents of the new maps, in “Dogs are better than cats”, **too\_far\_left**['dog'] = [0] for all windows. In “eight mice eat eight cheese sticks” with window 5, **too\_far\_left**['eight'] = [0,3]. For ease of computation the indices stored in **too\_far\_right** are transformed before storage using:

$$\widehat{id}x = length - idx - 1 = g(idx) \quad (2)$$

where  $length$  is the length of the small piece of text being analyzed. Then, **too\_far\_right**['cats'] = [ $g(4) = id x$ ] = [0 =  $\widehat{id}x$ ].

Now, to compute the exact number of occurrences of a bigram we do the computation:

$$freq(*, word) = (w - 1) * wordfd[word] - \sum_{i=1}^N (w - tfl[word][i] - 1) \quad (3)$$

where  $w$  is the window size being searched for bigrams,  $wfd$  is a frequency distribution of all words in the corpus,  $tfl$  is the map **too\_far\_left** and  $N$  is the number of occurrences of the  $word$  in a position too far left. The computation of  $freq(word, *)$  can now be performed in the same way by simply substituting  $tfl$  with  $tfr$  thanks to transformation  $g$ , which reverses the indexing.

#### References

- [1] S. Bird, E. Klein, and E. Loper, *Natural Language Processing With Python*, 1st ed., O'Reilly Media, Inc., 2009.
- [2] K. W. Church and P. Hanks, *Word Association Norms, Mutual Information, And Lexicography*, Computational Linguistics **16** (1990), no. 1, 22–29.