

A Lexicographer-Friendly Association Score

Pavel Rychlý

Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
pary@fi.muni.cz

Abstract. Finding collocation candidates is one of the most important and widely used feature of corpus linguistics tools. There are many statistical association measures used to identify good collocations. Most of these measures define a formula of a association score which indicates amount of statistical association between two words. The score is computed for all possible word pairs and the word pairs with the highest score are presented as collocation candidates. The same scores are used in many other algorithms in corpus linguistics.

The score values are usually meaningless and corpus specific, they cannot be used to compare words (or word pairs) of different corpora. But end-users want an interpretation of such scores and want a score's stability. This paper present a modification of a well known association score which has a reasonable interpretation and other good features.

1 Introduction

Finding collocation candidates is one of the most important and widely used feature of corpus linguistics tools [1]. There are many statistical association measures used to identify good collocations. Most of these measures define a formula of a association score which indicates amount of statistical association between two words. The score is computed for all possible word pairs and the word pairs with the highest score are presented as collocation candidates. The same scores are used in many other algorithms in corpus linguistics, for example to compute collocations in grammatical relations and an importance of grammatical relations in the Sketch Engine [2].

There are two general problems of most association scores:

1. A score is fine-tuned to one particular corpus size and/or key word frequency. If we use a score for a corpus with very different number of tokens the resulting list is not satisfying enough or is completely wrong.
2. The score values are usually meaningless and corpus specific, they cannot be used to compare words (or word pairs) of different corpora. But end-users want an interpretation of such scores and want a score's stability. They want to compare collocation scores of different words and on different corpora or subcorpora.

The article is organized as follows. The following section describe notation and the most widely used association scores. The Section 3 illustrates these two problems on real examples. The next section defines a new score *logDice*, which is a modification of the well known association score *Dice* [3]. The *logDice* score has a reasonable interpretation, scales well on a different corpus size, is stable on subcorpora, and the values are in reasonable range.

2 Association Scores for Collocations

Almost all association score formulas use frequency characteristics from a contingency table, which records the relationship between two words (W_1, W_2). Table 1 shows an example of a contingency table. The numbers in the right-hand column and the bottom row are called marginal frequencies and the number in the bottom right-hand corner is the size of the corpus.

In the rest of this paper we will use the following symbols (the meaning is also summarized in Table 1):

- f_x = number of occurrences of word X
- f_y = number of occurrences of word Y
- f_{xy} = number of co-occurrences of words X and Y
- $R_x = \frac{f_{xy}}{f_x}$ = relative frequency of word X
- $R_y = \frac{f_{xy}}{f_y}$ = relative frequency of word Y

Table 1. Notation of frequencies of words X and Y

| | $W_1 = X$ | $W_1 \neq X$ | |
|--------------|----------------|----------------|-----------|
| $W_2 = Y$ | f_{xy} | $f_y - f_{xy}$ | f_y |
| $W_2 \neq Y$ | $f_x - f_{xy}$ | $N - f_{xy}$ | $N - f_y$ |
| | f_x | $N - f_x$ | N |

3 Widely Used Association Scores

This section summarize formulas of some association scores and gives its main characteristics. More scores, motivations, discussion of their mathematical background and full references could be find in [4].

$$\text{T-score: } \frac{f_{xy} - \frac{f_x f_y}{N}}{\sqrt{f_{xy}}}$$

$$\text{MI-score: } \log_2 \frac{f_{xy} N}{f_x f_y}$$

$$\text{MI}^3\text{-score: } \log_2 \frac{f_{xy}^3 N}{f_x f_y}$$

Minimum Sensitivity: $\min R_x, R_y$

Dice coefficient: $D = \frac{2f_{xy}}{f_x + f_y}$

MI log Freq: $MI\text{-score} \times \log f_{xy}$, used as salience in the first version of Word Sketches [2].

Table 2 lists the collocation candidates on lemmas to the verb *break* in the window from 5 tokens to the left to 5 tokens to the right. They were computed on the British National Corpus by the Manatee system [5].

Table 2. Collocation lists for different association scores

| | F_{xy} | T-score | | F_{xy} | MI-score | | F_{xy} | MI³-score |
|------|----------|----------------|---------------|----------|-----------------|------|----------|-----------------------------|
| the | 11781 | 99.223 | spell-wall | 5 | 11.698 | the | 11781 | 30.591 |
| . | 8545 | 83.897 | deadlock | 84 | 10.559 | down | 2472 | 29.882 |
| , | 8020 | 80.169 | hoodoo | 3 | 10.430 | . | 8545 | 29.558 |
| be | 6122 | 69.439 | scapulum | 3 | 10.324 | , | 8020 | 29.193 |
| and | 5183 | 65.918 | Yasa | 7 | 10.266 | be | 6122 | 28.311 |
| to | 5131 | 65.798 | intervenien | 4 | 10.224 | to | 5131 | 28.268 |
| a | 3404 | 52.214 | preparedness | 21 | 10.183 | and | 5183 | 28.246 |
| of | 3382 | 49.851 | stranglehold | 18 | 10.177 | into | 1856 | 27.854 |
| down | 2472 | 49.412 | logjam | 3 | 10.131 | up | 1584 | 26.967 |
| have | 2813 | 48.891 | irretrievably | 12 | 10.043 | a | 3404 | 26.717 |
| in | 2807 | 47.157 | Andernesse | 3 | 10.043 | have | 2813 | 26.593 |
| it | 2215 | 43.314 | irreparably | 4 | 10.022 | of | 3382 | 26.255 |
| into | 1856 | 42.469 | Thief | 37 | 9.994 | in | 2807 | 26.095 |
| he | 1811 | 39.434 | THIEf | 4 | 9.902 | it | 2215 | 25.876 |
| up | 1584 | 39.038 | non-work | 3 | 9.809 | out | 1141 | 25.821 |

| | F_{xy} | Min. Sens. | | F_{xy} | MI log Freq | | F_{xy} | Dice |
|-----------|----------|-------------------|----------|----------|--------------------|-----------|----------|-------------|
| down | 2472 | 0.027 | down | 2472 | 57.340 | down | 2472 | 0.0449 |
| silence | 327 | 0.018 | silence | 327 | 48.589 | silence | 327 | 0.0267 |
| leg | 304 | 0.016 | deadlock | 84 | 46.909 | into | 1856 | 0.0210 |
| law | 437 | 0.014 | barrier | 207 | 46.389 | leg | 304 | 0.0203 |
| heart | 259 | 0.014 | into | 1856 | 46.197 | off | 869 | 0.0201 |
| rule | 292 | 0.013 | off | 869 | 42.411 | barrier | 207 | 0.0191 |
| off | 869 | 0.013 | up | 1584 | 42.060 | law | 437 | 0.0174 |
| news | 236 | 0.013 | leg | 304 | 41.980 | up | 1584 | 0.0158 |
| into | 1856 | 0.012 | neck | 180 | 39.336 | heart | 259 | 0.0155 |
| barrier | 207 | 0.011 | law | 437 | 38.805 | neck | 180 | 0.0148 |
| away from | 202 | 0.011 | out | 1141 | 38.783 | news | 236 | 0.0144 |
| war | 294 | 0.010 | bone | 151 | 38.263 | rule | 292 | 0.0142 |
| ground | 182 | 0.010 | heart | 259 | 37.327 | out | 1141 | 0.0135 |
| record | 287 | 0.010 | Thief | 37 | 36.353 | away from | 202 | 0.0135 |
| neck | 180 | 0.010 | news | 236 | 36.296 | bone | 151 | 0.0130 |

4 logDice

As one can see from the previous section, *Dice* score gives very good results of collocation candidates. The only problem is that the values of the *Dice* score are usually very small numbers. We have defined *logDice* to fix this problem.

$$\logDice = 14 + \log_2 D = 14 + \log_2 \frac{2f_{xy}}{f_x + f_y}$$

Values of the *logDice* have the following features:

- Theoretical maximum is 14, in case when all occurrences of *X* co-occur with *Y* and all occurrences of *Y* co-occur with *X*. Usually the value is less than 10.
- Value 0 means there is less than 1 co-occurrence of *XY* per 16,000 *X* or 16,000 *Y*. We can say that negative values means there is no statistical significance of *XY* collocation.
- Comparing two scores, plus 1 point means twice as often collocation, plus 7 points means roughly 100 times frequent collocation.
- The score does not depend on the total size of a corpus. The score combine relative frequencies of *XY* in relation to *X* and *Y*.

All these characteristics are useful orientation points for any field linguist working with collocation candidate lists.

5 Conclusion

In this paper, we have presented the new association score *logDice*. The *logDice* score has a reasonable interpretation, scales well on a different corpus size, is stable on subcorpora, and the values are in reasonable range.

Acknowledgments. This work has been partly supported by the Academy of Sciences of Czech Republic under the projects 1ET200610406, 1ET100300419 and by the Ministry of Education of CR within the Centre of basic research LC536 and National Research Programme 2C06009.

References

1. Smadja, F.: Retrieving Collocations from Text: Xtract. *Computational Linguistics* **19**(1) (1994) 143–177.
2. Kilgarriff, A., Rychlý, P., Smrž, P., Tugwell, D.: The Sketch Engine. *Proceedings of Euralex (2004)* 105–116.
3. Dice, L.: Measures of the Amount of Ecologic Association Between Species. *Ecology* **26**(3) (1945) 297–302.
4. Evert, S.: The Statistics of Word Cooccurrences: Word Pairs and Collocations. Unpublished Ph.D. dissertation, University of Stuttgart (2004).
5. Rychlý, P.: Manatee/Bonito – A Modular Corpus Manager. In: P. Sojka, A. Horák (Eds.): *RASLAN 2007 Proceedings (2007)*, pp. 65–70.