

# On the role of words in the network structure of texts: application to authorship attribution

Camilo Akimushkin<sup>a,\*</sup>, Diego R. Amancio<sup>b</sup>, Osvaldo N. Oliveira Jr.<sup>a</sup>

<sup>a</sup>*São Carlos Institute of Physics, University of São Paulo, Avenida Trabalhador  
São-carlense 400, São Carlos, São Paulo, Brazil*

<sup>b</sup>*Institute of Mathematics and Computer Science, University of São Paulo, Avenida  
Trabalhador São-carlense 400, São Carlos, São Paulo, Brazil*

---

## Abstract

Well-established automatic analyses of texts mainly consider frequencies of linguistic units, e.g. letters, words and bigrams, while methods based on co-occurrence networks consider the structure of texts regardless of the nodes label (i.e. the words semantics). In this paper, we reconcile these distinct viewpoints by introducing a generalized similarity measure to compare texts which accounts for both the network structure of texts and the role of individual words in the networks. We use the similarity measure for authorship attribution of three collections of books, each composed of 8 authors and 10 books per author. High accuracy rates were obtained with typical values from 90% to 98.75%, much higher than with the traditional the TF-IDF approach for the same collections. These accuracies are also higher than taking only the topology of networks into account. We conclude that the different properties of specific words on the macroscopic scale structure of a whole text are as relevant as their frequency of appearance; conversely, considering the identity of nodes brings further knowledge about a piece of text represented as a network.

*Keywords:* complex networks, word semantics, authorship attribution, similarity measures

*2010 MSC:* 00-01, 99-00

---



---

\*Corresponding author

Email address: `diego.rafael@gmail.com` (Diego R. Amancio)

## 1. Introduction

The huge volume of written text produced everyday makes it imperative to use automatic tools to retrieve relevant information, e.g. with text summarization, information retrieval methods, polarity analysis, citation analysis, and document classification [1, 2, 3, 4, 5, 6, 7]. An essential step in many of these tasks is to compare pieces of texts, as in classification of texts into categories [4] and in search engines where typically a list of texts relevant to a given query is retrieved. A special case is the pairwise comparison, where one searches for similarities between pairs of texts, which is actually a typical subtask in the authorship attribution process [8]. Automatic authorship attribution has been made with varied strategies [9], from the use of first-order statistics of linguistic elements to the processing of text represented as networks [10, 11]. For example, the frequency of characters [12, 13], phonemes [14], and morphemes [15, 16] has been explored, with texts normally modelled as lists of individual words, i.e. word order is disregarded. The archetype of such models is the so-called bag-of-words (BoW) model [17], where the text is represented as the set of its constitutive words by counting the number of appearances for each word. Word frequencies, which follow Zipf's law [18, 19], can then be used straightforwardly as attributes in a machine learning scheme [20] or to further build specific similarity measures.

Variations of the BoW model have been developed to address possible biases, e.g. the tendency of larger texts of being more likely to be considered similar to any other. These variations include the use of the term frequency-inverse document frequency (TF-IDF) statistic [21, 4], where lower relevance is assigned to words frequent in the document as well as in the whole collection. The model has also been modified to incorporate other kinds of data, such as in the bag-of-features model used for image analysis [22]. Another important modification is to consider  $n$ -grams, i.e., groups of  $n$  adjacent words [23, 24], in an attempt to take syntactic information into account, since the BoW model disregards word ordering. In other types of work, the syntactic roles of the words in sentences

are used for authorship attribution [25, 26]. It must be noted, nevertheless, that all of these approaches are based on the counting of features, even if some consider small-scale structural relationships.

An alternative perspective has been developed in recent years from the discovery that language features may be best described by complex networks models [27]. The structure of a text, for instance, can be mapped onto a co-occurrence network [10], which is characterized by power-law distributions [19, 28], and core-periphery structures [29]. Even though the general features of these complex networks remain analogous for texts in the same language, the network representation can also be used for classification tasks, particularly for authorship attribution [10, 30, 31].

While the frequency-based methods overlook all structural relationships among words farther than in the same sentence, the methods based on co-occurrence networks ignore the identity of the words (i.e. which actual word corresponds to a given node), thus characterizing the texts only on the basis of the network topology. In this study, we reconcile both viewpoints to show that, from a network perspective, words can play relevant roles in the structure of a text besides their frequencies.

## 2. Methods

The methodology proposed to address the authorship attribution task consists of four steps: i) construct a co-occurrence network for each text; ii) obtain various distance matrices for the collection using the proposed similarity metrics (see below); iii) join the various distance matrices with multi-dimensional scaling [32]; and iv) analyze the resulting data with standard supervised learning algorithms [20]. These steps are described in detail below. The model was applied to three collections of 80 literary texts. Each collection contains 10 texts per author for 8 authors from the 19th century, with 22 of the 24 authors being native English writers (details of the collections are included in the Supporting Information).

### 2.1. Network construction and characterization

Texts are pre-processed for constructing the networks, with stopwords, such as articles and prepositions, being removed, and lemmatization being applied to reduce different forms to a common base form. Lemmatization is assisted by a part-of-speech tagger based on entropy maximization [33], in order to solve ambiguities in mapping words to their lemmatized form. From the resulting pre-processed text, a co-occurrence (or word adjacency) network is built, where each distinct word is a node and two nodes are connected if the words appear consecutively in the text. The link is directed according to the natural reading order. For instance, the title of this paper generates the network: role  $\rightarrow$  word  $\rightarrow$  network  $\rightarrow$  structure  $\rightarrow$  text  $\rightarrow$  application  $\rightarrow$  author  $\rightarrow$  attribution. Each link has a default weight equal to one, which is increased by one unit each time the pair of words appears in the text.

Networks were characterized in this study by four well-known node-local metrics:

1. Degree ( $k_i$ ): this metric corresponds to the number of links attached to a node. As a consequence of the construction rules imposed by co-occurrence networks, there is strong correlation between this metric and the word frequency.
2. Average shortest path length ( $l_i$ ): this is the typical distance between two nodes of the network, given by:

$$l_i = N^{-1} \sum_j d_{ij}, \quad (1)$$

where  $d_{ij}$  is the shortest path length between nodes  $i$  and  $j$ , and  $N$  is the number of nodes. This metric is useful to identify keywords in written texts, irrespectively of the word frequency [34]. Low values of  $l$  are not only associated to the frequent words, but also to the words appearing close to other relevant words in the text.

3. Betweenness centrality ( $B_i$ ): the betweenness is the fraction of all shortest

paths that pass through the node, i.e.

$$B_i = \sum_{i \neq j \neq k} \frac{n_{jk}^{(i)}}{n_{jk}}, \quad (2)$$

where  $n_{jk}^{(i)}$  is the number of shortest paths from  $j$  to  $k$  passing through  $i$  and  $n_{jk}$  is the total number of shortest paths from  $j$  to  $k$ . In text analysis, the betweenness can be interpreted as a measure to quantify the ability of a word to appear in restrict or wider contexts [10].

4. Intermittency ( $I_i$ ): the intermittency is a measure that quantifies the spatial distribution of a given word along a text. To define this measure, consider the text as a sequence of tokens. This sequence generates, for each word  $i$ , a time series  $T^{(i)} = \{t_1^{(i)}, t_2^{(i)}, \dots, t_{f_i}^{(i)}\}$ , where  $t_j^{(i)}$  corresponds to the position of the  $j$ -th occurrence of the word  $i$ . The interval recurrence ( $\tau$ ) for word  $i$  is defined as the spatial difference between two occurrences, i.e.  $\tau_j^{(i)} = t_j^{(i)} - t_{j-1}^{(i)}$ . The set of all values of  $\tau_j^{(i)}$ , i.e.  $\mathcal{T}^{(i)} = \{\tau_1^{(i)}, \tau_2^{(i)}, \dots\}$  is used to quantify the regularity of the appearance of  $i$  along the sequence of tokens. More specifically, this regularity is computed using the intermittency defined as:

$$I_i = \sigma_{\mathcal{T}} / \langle \mathcal{T} \rangle = \left[ \frac{\langle \mathcal{T}^2 \rangle}{\langle \mathcal{T} \rangle^2} - 1 \right]^{1/2}, \quad (3)$$

where  $\sigma_{\mathcal{T}}$  and  $\langle \mathcal{T} \rangle$  are the standard deviation and average of  $\mathcal{T}$ , respectively. In text networks, the intermittency also measures the relevance of words, since it has been shown that intermittent (i.e. bursty) words are the ones most related to the subject being approached [34].

## 2.2. Similarity metrics

The novelty introduced in this work is to compare the words representing the most relevant nodes in the network topology, in contrast to previous approaches where only the statistics of topological metrics were taken into account [10]. We consider as the most relevant the nodes possessing the highest degree and

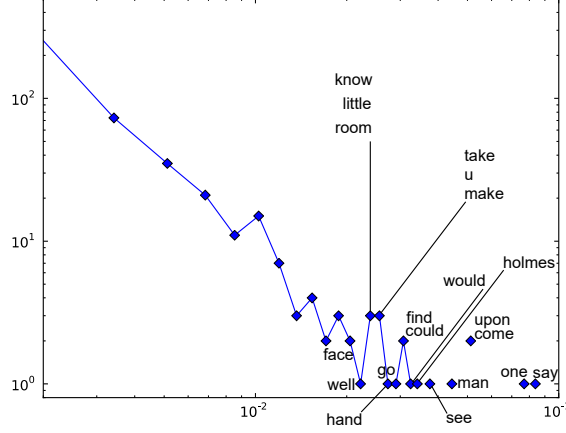


Figure 1: Betweenness centrality distribution for “The Memoirs of Sherlock Holmes”. The top 20 nodes are labeled by the corresponding words.

betweenness. As for the other metrics, namely average shortest paths and intermittency, we chose the nodes with lowest values. We tested the largest shortest paths but the results were not as good. Intermittency obeys a power law with positive exponent. We hypothesize that two pieces of text will be similar if there is significant overlap in the words (nodes) considered most relevant in both texts.

Figures 1 and 2 show the distributions of betweenness centralities for two books from Arthur Conan Doyle: The Memoirs of Sherlock Holmes and The Return of Sherlock Holmes, which could be expected to be similar since these books were written by the same author in the same series of novels. In both figures, the highest centralities belong to the same words (19 out of 20) which also occupy almost the same relative positions. We shall therefore test the hypothesis that not only the frequency of usage but also the long-scale topological metrics of the organization of words may be reliable signatures of authorship.

To quantify this we introduce a similarity measure between pairs of texts as follows: for each network metric considered we give a rank  $R$  to each word  $w$  from a subset  $V$  of top words with unique properties. In our approach the

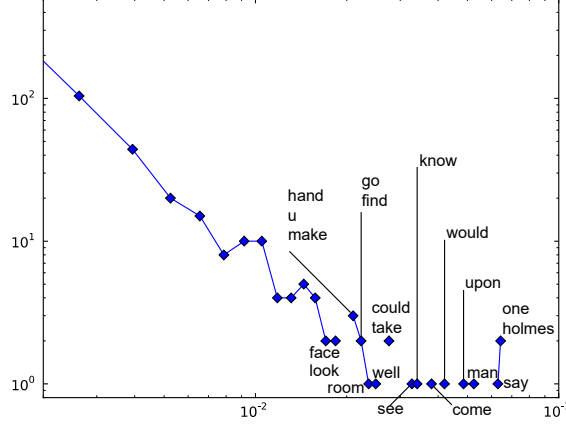


Figure 2: Betweenness centrality distribution for “The Return of Sherlock Holmes”. The top 20 nodes are labeled with the corresponding words.

importance of a word depends on the particular network metric considered. As mentioned above, we select the words with the highest connectivity and betweenness centrality, and with the smallest values of shortest path and intermittency. For these two latter metrics the smallest values were chosen because their distributions present power laws with positive coefficients. We choose sets of 100 top words since in subsidiary experiments we observed that the interval between 50 – 150 words gave the best results. A ranking is assigned to each word, starting with the maximum value (100 in this case) for the word with the most extreme value (e.g. “say” for the betweenness centrality of “The Memoirs of Sherlock Holmes”), and decreasing in one unit for each consecutive word until reaching the last of the top words which receives a ranking value of one. With these rankings, the similarity between two texts  $A$  and  $B$  for a given network metric is given by

$$A \cdot B = \sum_{w \in V_A \cap V_B} R_A(w) R_B(w), \quad (4)$$

that is, if a word is present in the top words subsets of both texts, the product of its rankings adds to their similarity.

This similarity metric is guaranteed to be high only if the same words occupy

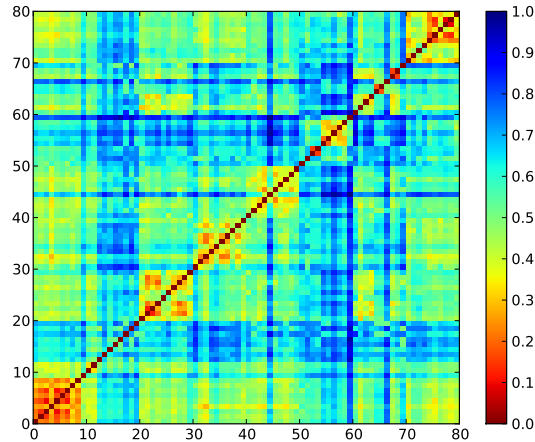


Figure 3: Betweenness centrality distance matrix for the second collection, each decade corresponds to texts from the same author.

similar positions in the distributions such as in figures 1 and 2, with higher influence from the highest-ranked words. Equation 4 implies that the norm or similarity of a text with itself is always the same, that is,  $A \cdot A = B \cdot B = n(n+1)(2n+1)/6$ , where  $n$  is the size of  $V$ . We therefore normalize all similarities for this value to be one and the minimum value to be zero, and define the distance  $D_{AB}$  between two texts as being one minus this normalized value. It is worth noting that other similarity metrics could be used to compare two pairs of texts, but the dot product adopted here appears to be the most straightforward, as it is done in bag-of-words methods [4].

With all the values  $D_{AB}$  we produce a distance matrix for each metric. The distance matrix for the betweenness centrality of one of the collections is shown in figure 3, where the indices 0 to 9 correspond to texts from the first author, texts 10 to 19 to the second author and so on. Note that, in general, texts from the same author appear to be closer among themselves compared to texts from different authors even if they are relatively separated (e.g. texts 10 – 19 and 50 – 59).



### 2.3. Combining Distance Matrices

One strength of the approach is the ability of observing different aspects of the network structure simultaneously. Each metric yields a different distance matrix; hence, we can observe the similarity between texts at different scales. We now combine information from the distinct metrics in order to have useful data for the classification algorithms.

In this study we employed two strategies for the input into the classification algorithms. In the first, we simply used the whole of the distance matrices for the different metrics, i.e. with distances as attributes. In the second strategy, we reduced the dimensionality of the distance metrics with Multi-dimensional scaling (MDS) [35], with the aim of capturing the highest similarities while eliminating possible unnecessary information that may harm the classification task. MDS was conceived to map distances into positions in a space so that the distances between these positions reproduce as well as possible the original input distances. The space obtained is usually intended to have a small dimensionality and the algorithm is largely used for visualization purposes. The positions obtained when applying the algorithm to map one of the distance matrices to a two-dimensional space is presented in figure 4 reflecting the similarities between same-author texts already observed with the distance matrices. We use MDS to map the four distance matrices of each collection into four subspaces and then join these subspaces into a space of bigger dimension: if we write the positions in each subspace as a matrix  $M \times N_i$  where  $M$  is the number of points (80 texts per collection in our case) and  $N_i$  is the dimensionality of the subspace, then the positions in the total space are given by a matrix  $M \times (N_1 + N_2 + N_3 + N_4)$  where each row is composed joining head to tail the corresponding rows of the positions on the subspaces.

Instead of a bi-dimensional mapping such as that of figure 4, the dimensionalities  $N_i$  are calculated based on the stress or cost function (the difference between the actual and the obtained distances). As stress is a monotonically decreasing function of the number of dimensions we set a threshold of 10% of the value for one dimension (also known as the elbow method) which was usually

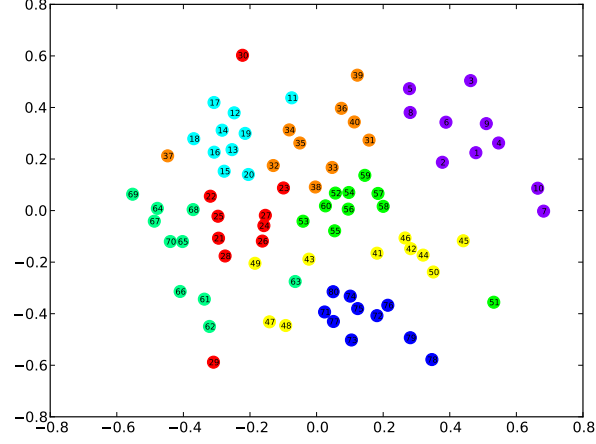


Figure 4: Bi-dimensional MDS mapping of the betweenness centrality distances for the third collection. Numbers correspond to texts indices, colors correspond to authors.

found to be reached at  $N_i = 6$ .

#### 2.4. Data analysis

The final positions on the composed space are the attributes for the data analysis algorithms. Analysis is done with supervised learning algorithms from the main types currently in use: tree-based J48; K-Nearest Neighbors (KNN); Naive Bayesian (NB); and Radial Basis Function Network (RBFN). For all cases 10-fold cross validation is applied and the parameters are set to their default values [20]. For KNN the number of neighbors is set to three which is the smallest odd non-trivial value. For RBFN the number of clusters is set to eight which is the number of authors. Authorship is also addressed using the standard TF-IDF model. Since TF-IDF returns a distance matrix, we also use MDS in this single matrix in order to apply the same classification algorithms on both approaches.

### 3. Results and Discussion

The approach based on distance matrices as input for the classification algorithms was applied for the three collections, for which the success score for classification by chance is  $1/8 = 12.5\%$ . The results are outstanding as shown in table 1, especially when MDS was used. It seems therefore that reducing the dimensionality actually amounted to an efficient feature selection, probably eliminating data that brought noise to the analysis. With MDS, typical accuracy rates were above 90% and the maximum value was 98.75% obtained with KNN for the third collection which corresponds to only one text (out of 80) not correctly classified. These scores greatly surpassed the values obtained by applying the TF-IDF method, for which the mean scores among collections were 36.67% for J48, 66.25% for KNN, 63.75% for NB, and 65% for RBFN, as shown in figure 5. These scores demonstrate the added value of using the network structure over relying only on the frequency of appearance of features. Significantly, the higher scores for the approach introduced here are maintained when changing the classification algorithm (KNN, NB, and RBFN), which indicates the robustness of the proposed metrics. Also worth noting is that the present approach outperforms a previous one where the topology of networks was taken without considering the labels of the nodes (words) [31], for which the accuracy rates for the second collection studied here were 63.75% with J48, 88.75% with KNN, 81.25% with NB, and 83.75% with RBFN. In addition, the approach presented is less demanding, both computationally and conceptually, than the previous one.

Taken together, the results indicate that, apart from the frequency of appearance and syntactical relations, certain words are essential to the structure of a text as a whole. The procedure to identify such words using complex networks has been successful, since utilization of these words is author-dependent. The co-occurrence network procedure allows one to observe the features of a word at different scales in the text. For instance, words with low intermittency, i.e. whose appearance in the text is highly periodic, had a high relevance for

	J48	KNN	NB	RBFN
Without MDS				
Collection 1	73.75	85.00	85.00	62.50
Collection 2	73.75	83.75	82.50	78.75
Collection 3	75.00	92.50	80.00	71.25
Using MDS				
Collection 1	72.50	87.50	92.50	90.00
Collection 2	63.75	97.50	93.75	96.25
Collection 3	73.75	98.75	92.50	95.00

Table 1: Accuracy rates (in percentage) in identifying the authors in the three collections, using several machine learning algorithms. Results are shown with the input comprising the whole distance matrices (without MDS) and applying MDS on the matrices. For the three last algorithms MDS improved accuracy in all cases.

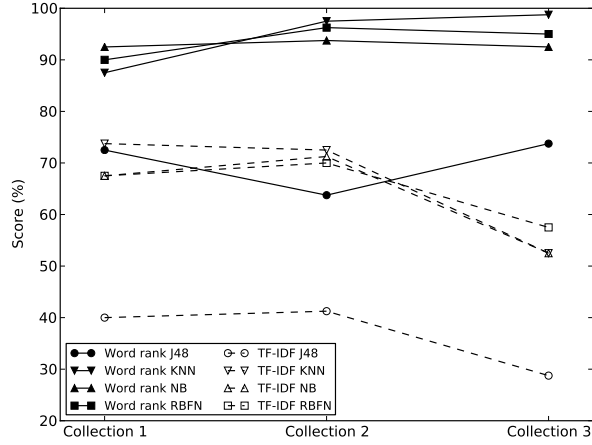


Figure 5: Scores obtained with the method introduced here, also using MDS, and with TF-IDF for the three collections of texts.

the authorship attribution, even though a word can have a low intermittency and not appear in most paragraphs.

Ours is a unified framework for multivariate analysis of texts. In contrast to other multivariate approaches, the generalized similarity measure 4 allows to easily introduce new features to the scheme using some of the many node-local network metrics in existence. Care must be taken, however, because not all metrics are useful and some can even lower the performance. For example, we tested the clustering coefficient (characterized by a bell-shaped distribution) and eigenvalue centrality without success. Even though the computation of the similarity measures between documents resembles that of TF-IDF (i.e. cosine similarity) there are significant differences, mainly the fact that the vectors are of a much smaller size and that there are no repeated values. A question to be further studied is the optimal ranking procedure: we chose a ranking *a la* Zipf because of the presence of power law distributions, but other rankings could be possible. While both TF-IDF and our method account for the heterogeneity of sizes of texts, our ranking procedure has two principal advantages: computation is faster and most importantly, the ranking does not have to be repeated every time the collection is modified, which is especially advantageous with big collections.

#### **4. Conclusions**

We have introduced an approach by which the representation of text with complex networks is enhanced by considering the words corresponding to the nodes. This is done with a similarity metric to compare two pieces of text where the presence of the most relevant words, according to network metrics, is taken into account. When the distance matrices obtained with the similarity metrics were used as input into machine learning algorithms, a high accuracy was achieved which reached 98.75% for one of the book collections. Significantly, the accuracy was considerably higher than for traditional methods based on TF-IDF, being also higher than other network approaches that did not consider the

label of the nodes. Also relevant is that the performance was improved with dimensionality reduction with MDS, which is advantageous owing to the lower computational cost.

With regard to the limitations, one should emphasize that the present approach is not useful for very short texts (such as a summary of an article). The method can be extended to employ other metrics and multi-node distributions. As some authors have pointed out [36], it is likely that every person has a characteristic writing fingerprint owing to their particular way to learn a language. If this is the case, the traits that define such fingerprint are probably complex and not bounded to one single measure. Finally, the approach proposed could be used for such other applications as part-of-speech analysis of network distributions and resolution of word polysemy.

## Acknowledgments

This work was supported by CNPq (Brazil) and FAPESP (grants 2014/20830-0, 2013/14262-7 and 2016/19069-9).

## References

## References

- [1] W. Liang, Physica A 468 (2017) 802 – 808.
- [2] T. C. Silva, D. R. Amancio, EPL (Europhysics Letters) 98 (2012) 58001.
- [3] X. Zhong, J. Liu, Y. Gao, L. Wu, Physica A 466 (2017) 462 – 475.
- [4] C. D. Manning, P. Raghavan, H. Schütze, et al., Introduction to information retrieval, volume 1, Cambridge university press Cambridge, 2008.
- [5] M. Z. Asghar, A. Khan, S. Ahmad, M. Qasim, I. A. Khan, PLoS One 12 (2017) e0171649.
- [6] D. R. Amancio, O. N. Oliveira Jr., L. F. Costa, Journal of Informetrics 6 (2012) 427 – 434.

- [7] M. P. Viana, D. R. Amancio, L. F. Costa, *Journal of Informetrics* 7 (2013) 371 – 378.
- [8] P. Juola, *Foundations and Trends in information Retrieval* 1 (2006) 233–334.
- [9] E. Stamatatos, *J. Am. Soc. Inf. Sci. Technol.* 60 (2009) 538–556.
- [10] D. R. Amancio, E. G. Altmann, O. N. Oliveira Jr, L. F. Costa, *New Journal of Physics* 13 (2011) 123024.
- [11] D. R. Amancio, *Journal of Statistical Mechanics: Theory and Experiment* 2015 (2015) P03005.
- [12] F. Peng, D. Schuurmans, S. Wang, V. Keselj, in: *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics-Volume 1*, Association for Computational Linguistics, pp. 267–274.
- [13] H. J. Escalante, T. Solorio, M. Montes-y Gómez, in: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, Association for Computational Linguistics, pp. 288–298.
- [14] C. Forstall, W. Scheirer, in: *Journal of the Chicago Colloquium on Digital Humanities and Computer Science*, volume 1.
- [15] O. V. Kukushkina, A. Polikarpov, D. V. Khmelev, *Problems of Information Transmission* 37 (2001) 172–184.
- [16] C. E. Chaski, *International journal of digital evidence* 4 (2005) 1–13.
- [17] Z. Harris, *Word* 10 (1954) 146–62.
- [18] G. K. Zipf, *The psycho-biology of language*, Houghton, Mifflin, 1935.
- [19] R. Ferrer-i Cancho, R. V. Solé, *Journal of Quantitative Linguistics* 8 (2001) 165–173.

- [20] D. R. Amancio, C. H. Comin, D. Casanova, G. Travieso, O. M. Bruno, F. A. Rodrigues, L. F. Costa, PLoS One 9 (2014) e94137.
- [21] K. Sparck Jones, Journal of documentation 28 (1972) 11–21.
- [22] J. Zhang, M. Marszałek, S. Lazebnik, C. Schmid, International journal of computer vision 73 (2007) 213–238.
- [23] V. Kešelj, F. Peng, N. Cercone, C. Thomas, in: Proceedings of the conference pacific association for computational linguistics, PACLING, volume 3, pp. 255–264.
- [24] R. Clement, D. Sharp, Literary and linguistic computing 18 (2003) 423–447.
- [25] H. Baayen, H. Van Halteren, F. Tweedie, Literary and Linguistic Computing 11 (1996) 121–132.
- [26] J. Rygl, K. Zemková, V. Kovár, in: Proceedings of Sixth Workshop on Recent Advances in Slavonic Natural Language Processing, RASLAN, pp. 111–119.
- [27] J. Cong, H. Liu, Physics of Life Reviews 11 (2014) 598 – 618.
- [28] S. N. Dorogovtsev, J. F. F. Mendes, Proceedings of the Royal Society of London. Series B: Biological Sciences 268 (2001) 2603–2606.
- [29] M. Choudhury, D. Chatterjee, A. Mukherjee, in: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, Association for Computational Linguistics, pp. 162–170.
- [30] A. Mehri, A. H. Darooneh, A. Shariati, Physica A 391 (2012) 2429–2437.
- [31] C. Akimushkin, D. R. Amancio, O. N. Oliveira Jr., PLoS One 12 (2017) e0170527.
- [32] I. Borg, P. Groenen, Modern Multidimensional Scaling: Theory and Applications, Springer, 2005.



- [33] B. B. Greene, G. M. Rubin, Automatic grammatical tagging of english, 1971. Department of Linguistics, Brown University, Providence, Rhode Island.
- [34] M. Ortuno, P. Carpena, P. Bernaola-Galvn, E. Muoz, A. M. Somoza, EPL (Europhysics Letters) 57 (2002) 759.
- [35] J. B. Kruskal, Psychometrika 29 (1964) 1–27.
- [36] H. Van Halteren, H. Baayen, F. Tweedie, M. Haverkort, A. Neijt, Journal of Quantitative Linguistics 12 (2005) 65–77.