

Visualising text co-occurrence networks

HIRSCH, Laurie and ANDREWS, Simon

Available from Sheffield Hallam University Research Archive (SHURA) at:

<http://shura.shu.ac.uk/12832/>

This document is the author deposited version. You are advised to consult the publisher's version if you wish to cite from it.

Published version

HIRSCH, Laurie and ANDREWS, Simon (2016). Visualising text co-occurrence networks. CEUR Workshop Proceedings, 1637, 19-27.

Repository use policy

Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in SHURA to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

Visualising Text Co-occurrence Networks

Laurie Hirsch

Simon Andrews

Sheffield Hallam University

Abstract. We present a tool for automatically generating a visual summary of unstructured text data retrieved from documents, web sites or social media feeds. Unlike tools such as word clouds, we are able to visualise structures and topic relationships occurring in a document. These relationships are determined by a unique approach to co-occurrence analysis. The algorithm applies a decaying function to the distance between word pairs found in the original text such that words regularly occurring close to each other score highly, but even words occurring some distance apart will make a small contribution to the overall co-occurrence score. This is in contrast to other algorithms which simply count adjacent words or use a sliding window of fixed size. We show, with examples, how the network generated can be presented in tree or graph format. The tree format allows for the user to interact with the visualisation and expand or contract the data to a preferred level of detail. The tool is available as a web application and can be viewed using any modern web browser.

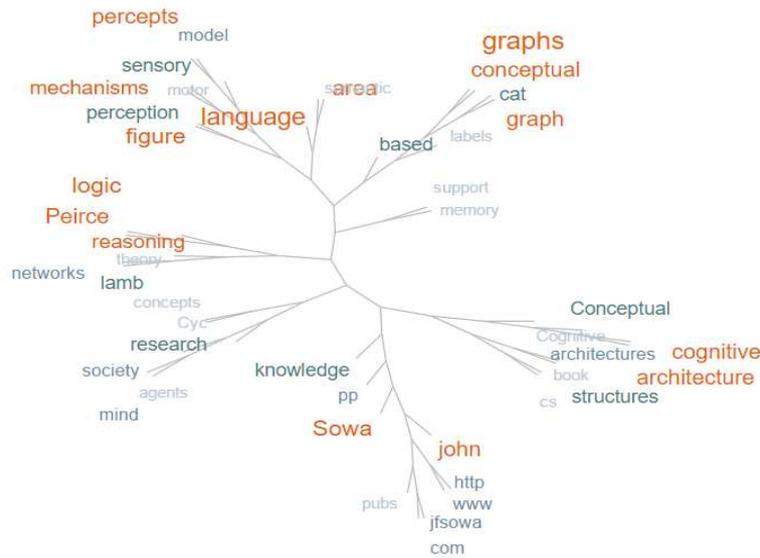
1 Background

Visual representations have proved to be useful alternatives to linear text documents. The mind mapping technique was introduced in the 1960s and is thought to encourage learning. However, creating mind maps can be a complex and time-consuming undertaking and the ability to automatically produce text visualisations has attracted significant research in recent decades. A number of possible benefits have been attributed to such tools including managing information overload, providing summaries and ‘impression formation’. Tools have been developed for identifying topics and topic correlations, displaying knowledge and generating concept clouds [1][2]. Here we will briefly outline a number of existing techniques and then show how we have developed a method based on word co-occurrence which can be used for generating both graphs and trees in various types of diagram. Here we include a number of example visualisations, all of which are based on the text of a paper concerning conceptual structures[3]¹

¹ Available at <http://www.jfsowa.com/pubs/ca4cs.pdf> It may help the reader to briefly read the article before viewing the visualisations.

available and gives the user an indication of the relationship between the key terms in the visualisation. The Sowa text produces the tree cloud shown in figure 2³

Fig. 2. Tree Cloud



The tree cloud includes colouring, font sizes and arcs to indicate relationships between topics.

2 Description of the System

In this section we will describe how our system known as txt2vz (<http://txt2vz.appspot.com/>) works and will compare visualisations produced with other text visualisation tools.

2.1 Pre-processing

To reduce dimensionality of the document(s) all words are placed in lower case, stop words are removed and stemming applied, such that only the most frequent form of a word is preserved.

2.2 Significance Measure.

We define a measure of significance for a pair (P, Q) of words, based on the number of occurrences of (P, Q), or more specifically the co-occurrences and the distance between

³ Using the tool at http://treecloud.univ-mlv.fr/cgi-bin/NuageArbore_EN.cgi#

P and Q where the distance between P and Q is defined to be the number of words between P and Q:

$$\text{significance}(P, Q) = \sum_{i=1}^M B^{\text{distance}(PQ_i)} \quad (1)$$

where M is the number of co-occurrences of P and Q; d_i is the distance between P and Q in the i th co-occurrence; $0 < B < 1$ B is between 0 and 1 and typically set to 0.9. We do not consider the significance if the distance is beyond a pre-set maximum distance which has a default of 20 words. The use of a decaying function here is in contrast to commonly used ‘sliding window’ methods of computing co-occurrence where we simply count the number of times that two words occur within a predefined distance.

2.3 Network Generation Algorithm.

The visualisations produced by txt2vz are intended for use in a web application and the visualisation should be presented to the user in a reasonable time. Even after dimension reduction there are likely to be a large number of unique words in a text and computing the co-occurrence value for each possible pair can be time consuming. We therefore sort the words according to frequency and select the top N words for the next stage where N is typically set to 200. The significance of each pair of words from the remaining set is computed and all the word pairs are sorted in descending order by their accumulated co-occurrence values. An undirected graph can then be built by selecting the top K word pairs and creating an edge between the two words of each pair. The graph is built using the d3.js software library [6]

2.4 Txt2vz network

The simplest format for txt2vz has been described previously [7] and the visualization for the Sowa text is shown in Figure 3.

Fig. 5. Radial Tree (small)

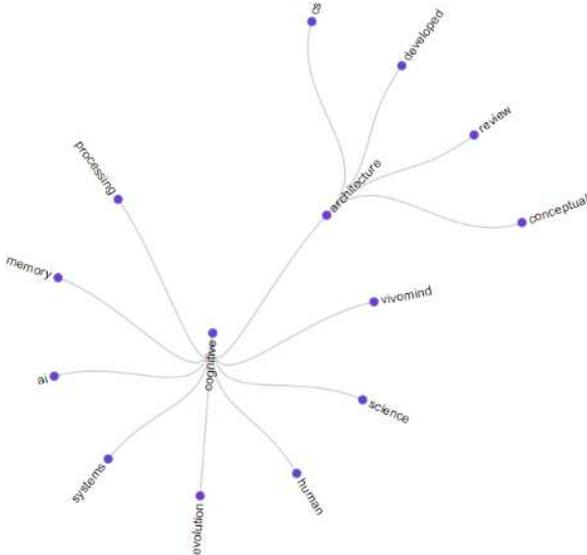
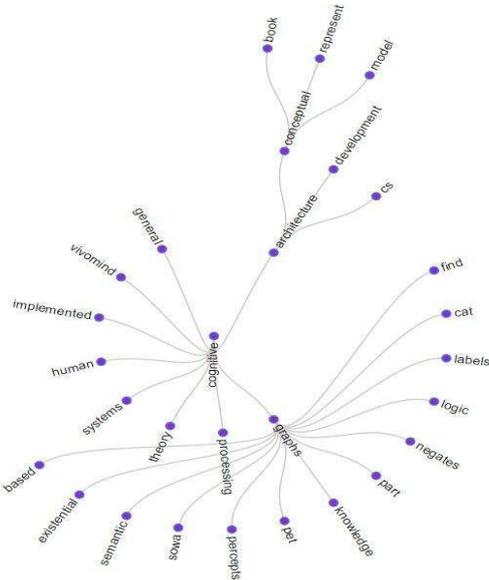


Fig. 6: Radial Tree



3 Discussion and further work

Txt2vz is work in progress. Apart from the basic graph types there are many parameters that affect the appearance of the final visualisation. We have therefore added a control panel to the web application so that the user can experiment and view different perspectives on the same piece of text. For example, reducing the number of word pairs to be analyzed will produce the smaller graph shown in figure 5, which some users may prefer. We would not want to argue that one particular graph type or cloud is the ‘best’ but we do suggest that it may be useful for the user to be able to switch between different types of visualization.

We would like to spend more time evaluating the usefulness of the tool as perceived by human subjects. We also hope to investigate the feasibility of using Txt2vz as part of a web search engine such that a user could be presented with a quick visual summary of the content of the pages pointed to by the result links. Lastly we are investigating the possibility of scaling txt2vz such that it can produce visualisations of large text datasets. We leave the reader with two further radial tree visualisations on different texts. Figure 7 shows a visualization of this article and Figure 8 shows a visualization of Darwins ‘On the Origin of Species’.

Fig. 7

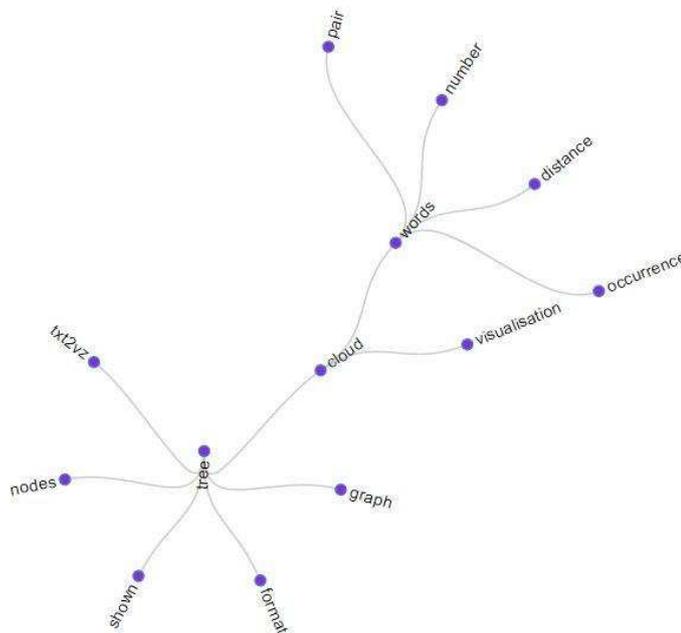


Fig. 8. The Origin of Species



4 References

1. Smith, A., Chuang, J., Hu, Y., Boyd-Graber, J. and Findlater, L., 2014. Concurrent Visualization of Relationships between Words and Topics in Topic Models. Sponsor: Idibon, 79.
2. Aga, R.T. and Wartena, C., 2015, October. Constructing concept clouds from company websites. In Proceedings of the 15th International Conference on Knowledge Technologies and Data-driven Business (p. 38). ACM
3. Sowa, J.F., 2011. Cognitive architectures for conceptual structures. In *Conceptual Structures for Discovering Knowledge* (pp. 35-49). Springer Berlin Heidelberg.
4. Viégas, F.B., Wattenberg, M., Tag Clouds and the Case for Vernacular Visualization, *ACM Interactions*, XV.4 - July/August, 2008
5. Gambette, P. and Véronis, J., 2010. Visualising a text with a tree cloud. In *Classification as a Tool for Research* (pp. 561-569). Springer Berlin Heidelberg..
6. Bostock, M., 2014. Data-Driven Documents-D3. js.
7. Hirsch, L. and Tian, D., 2013, January. Txt2vz: a new tool for generating graph clouds. In *International Conference on Conceptual Structures* (pp. 322-331). Springer Berlin Heidelberg.
8. Reingold, E.M. and Tilford, J.S., 1981. Tidier drawings of trees. *Software Engineering, IEEE Transactions on*, (2), pp.223-228.