

# BUILDING OF SEMANTIC NETWORKS TO DETERMINE THE DEGREE OF TEXT SIMILARITY OR DIFFERENCE

O. O. Dmytrenko<sup>1,2, a</sup>, D. V. Lande<sup>1, b</sup>

<sup>1</sup>*National Technical University of Ukraine «Igor Sikorsky Kyiv Polytechnic Institute»*

<sup>2</sup>*Institute for Information Recording of National Academy of Sciences of Ukraine*

## Abstract

The paper presents a method for comparing text documents based on the building and comparison of the corresponding semantic networks. This method can be the basis for building systems for comparing legal documents in the framework of parliamentary control. The algorithm for building semantic networks as one of the types of ontologies is also considered. This algorithm can also be used in systems of automatic abstracting of legal information in order to generate concise information-rich reports, brief annotations or digests. The proposed method can be used in the process of processing queries during information retrieval, providing the ability to determine the degree of similarity or difference in the structure and semantics of texts.

*Keywords:* semantic network, natural language processing, legal information, horizontal visibility network, text comparison, computational linguistics

## Introduction

The rapid development of information and telecommunication technologies causes the rapid accumulation of data in various sources — text files, emails, and web pages [1] in various presentation formats. The number of legal documents submitted in electronic form, and hence the amount of information that an expert in this field has to deal with, is also constantly growing. And in order to make informed decisions based on existing legal data, it is sometimes necessary to read thousands of documents, rejecting information noise. Therefore, the task of simplifying access to the essence of the text, extracting the key highlights, ideas and pre-stated content aspects, without the need to process a huge amount of information, is relevant for the legal field. Also, the main task is to determine the degree of similarity or difference and also inconsistencies in legal documents. All these problems lead to the need to develop and improve existing technological solutions and create new ones in order to ensure prompt processing and analysis of legal information. Taking into account the huge volume of legal texts, the task of formalizing textual data and presenting them in a form that would be convenient for automatic processing is urgent [2, 3, 4]. The purpose of the paper is to present a method for determining the degree of similarity between text documents, based on the use of directed weighted networks of terms, where the nodes of such networks are key terms of the text, and edges are semantic-semantic relationships between these terms in the text.

## 1. Main material

### Building of semantic networks

An example of a subject domain model (ontology), which can be represented as a huge array of text data, and which will be convenient for computer processing, is a directed weighted network of terms. Directed Weighted Network of Terms (DWNT) is a semantic model of text representation, where the nodes of such a network are key terms (words and phrases), which are used as the names of concepts in a particular subject area, and the edges is semantic-syntactic connections between these terms. Comparing the DWNTs obtained for different texts, accordingly, allows us to determine the semantic similarity of the respective texts. Building of network of term is carried out in several stages [2], including pre-processing of text data, extraction of key terms, construction of undirected network of terms (using the algorithm of horizontal visibility graph), i.e. determining undirected connections between terms, and further determining the directions of connections and their weight values.

For the pre-processing of text data, some of the most common techniques are used, including automatic segmentation into individual sentences and subsequent tokenization of the sentences — segmentation of the input text of sentences into elementary units (tokens) [5]. After tokenization, within each sentence Part-of-Speech tagging (PoS tagging) is doing [6]. PoS tagging consists in assigning each word in the text to a certain part of the language and assigning it a corresponding tag. In addition, in order to obtain canonical, lexical forms of tokens (lemmas), the lemmatization of individual marked tokens is carried out. This step allows to further group different forms of the same word so that they can be

<sup>a</sup>dmytrenko.o@gmail.com

<sup>b</sup>dwlande.o@gmail.com

analyzed as a single element. The functions of various Python programming language packages and libraries have been used to computerize word processing, classify tokens, and assign appropriate tags to them. In particular, for the texts presented in Ukrainian and Hebrew, the Pipeline functions of the «Stanza» library and, accordingly, the Ukrainian and Hebrew language models were used. The «NLTK» library was used for the texts presented in English. Russian-language texts are processed using the «pymorphy2» library. The following link [universaldependencies.org/docs/u/pos/](https://universaldependencies.org/docs/u/pos/) contains a set of predefined tags that the above-mentioned libraries use to match each word in a sentence to a specific part of the language.

For the extracting terms, words related to parts of speech such as noun (NOUN tag), including common names (PROPN tag), adjective (ADJ tag) and conjunction (CCONJ tag) were used. To build a network of terms, individual words that belong to parts of speech such as nouns (common names with the PROPN tag have been reassigned for convenience) were used. The following templates were used to construct the phrases:

- for bigrams — «ADJ NOUN»;
- for threegrams — «NOUN CCONJ NOUN», «ADJ ADJ NOUN»;
- for fourgrams — «ADJ NOUN CCONJ NOUN», «ADJ CCONJ ADJ NOUN».

Next, the removal of individual stop words (individual articles, prepositions, conjunctions, some verbs, adverbs and pronouns), and which do not have information load is carried out. The list of stop words was formed on the basis of a combination of several stop dictionaries, ones of which for Ukrainian, Russian, English and Hebrew language are available at [github.com/stopwords-iso](https://github.com/stopwords-iso). And also, each list was expanded with the another available in the Python package — [pypi.org/project/stop-words/](https://pypi.org/project/stop-words/). It is also planned to edit the stop words dictionary by adding and removing from the list of words that have been identified by experts within the research area. Using keyword and phrase templates, the next step is to form a sequence of terms where more phrases precede the phrases and words that are part of them, with the initial order of occurrence in the sentence being taken into account for single words. Next stage is to separate the key terms from the text for each formed term of the sequence, the so-called tuple of three elements is built: the first is the term (word or formed according to the presented templates); the next is a tag that is assigned to a word depending on its belonging to a certain part of the language, or a collective tag for the corresponding template; the last element of such a set — the numerical value of *GTF* (Global Term Frequency) — a global indicator of the importance of the term [2, 4]:

$$GTF = \frac{n_i}{\sum_k n_k}$$

where  $n_i$  is a number of terms  $i$  appearances in the text; sum is a general or global number of formed terms in the whole text. Taking into account the marking of

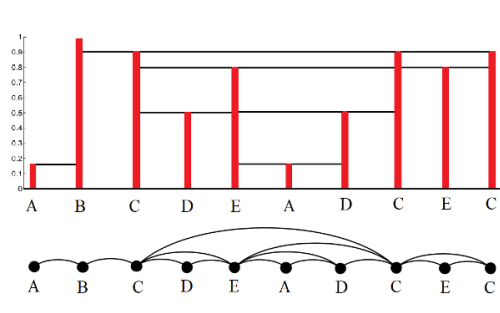


Fig. 1. Example of building the Horizontal Visibility Graph

parts of speech, *GTF* in this case is calculated taking into account the first two elements of the tuple — the term and tag. The number of such identical tuples in the whole sequence, which is normalized to the total number of generated terms, determines the value of the third element of the tuple — *GTF*. Unlike the usual *TF – IDF* statistic, *GTF* allows to more effectively find information-important elements of text when working with a text corpus of a predefined topic, when the information-important term occurs in almost every document in the corpus.

To build an undirected network of terms, as a terminological ontology of a particular subject area, this paper considers and applies an approach to building networks based on time series — Horizontal Visibility Graph algorithm (HVG) [7]. The Horizontal Visibility Graph Algorithm (HVG) [8], in turn, is an extension of the standard Visibility Graph Algorithm (VG) [9]. Horizontal visibility graphs are constructed within each individual sentence, where each term corresponds to a statistical estimate *GTF* (Global Term Frequency) — a global indicator of the importance of the term. An undirected network of terms using the Horizontal Visibility Graph Algorithm is built in two stages [7]. The first step is to mark on the horizontal axis a sequence of nodes  $t_i$ , each of which corresponds to the terms in the order in which they occur in the text; and the weighted values numerical estimates  $x_i$  that corresponded to *GTF* and intended to reflect how important a word is to a document in a collection or corpus are marked on the vertical axis. In the second stage, the horizontal visibility graph is created. It is considered, two nodes  $t_i$  and  $t_j$  corresponding to the elements of the time series  $x_i$  and  $x_j$ , are connected in a HVG if and only if,

$$x_k \leq \min(x_i; x_j)$$

for all  $t_k$  ( $t_i < t_k < t_j$ ), where  $i < k < j$  are the nodes of graph. The obtained undirected network of terms is called the horizontal visibility graph (HVG) (see figure 1).

Therefore, the considered HVG algorithm makes it possible to construct an undirected network structure from time series on the basis of texts in the case when numerical weight values (*GTF* in our case) are assigned to an individual words or phrases.

If a priori there is an undirected connection between the respective nodes in the horizontal visibility graph,

the directions of links in a undirected network of terms are established on the principle of entering a shorter term into a term, which is its an extension [10].

The weight values of the connections between the nodes in the directed network are determined by the principle proposed in [11, 12]: the vertices of the graph corresponding to the same terms of the previously constructed directed network are combined ("merged"). As a result, the weight values of the connections between the pairs of nodes are determined by the number of same directed connections between these nodes. Since any graph is determined by the adjacency matrix, the task of determining the weight values of the links is reduced to the concatenation of columns and corresponding rows, i.e. a weighted compactification of the horizontal visibility graph [7]. The resulting matrix defines an oriented weighted graph formed of vertices that correspond to unique terms in the text. The weight value of the edge, that connect the vertex  $i$  with the vertex  $j$  is determined by the number of occurrences of the term  $t_i$  before the term  $t_j$  in the text.

The resulting network can be saved in «graphml» and «json» formats. The open-source software package «Gephi» designed for network analysis and visualization is used to visualize networks presented in «graphml» format. The «json» format can be convenient for use in systems for building and visualizing semantic networks. During visualization, only the text of the term (words or phrases) is displayed as node labels, without specifying the part of the speech which was assigned to the term at the stage of PoS tagging.

### Building of semantic networks

When comparing the semantic networks considered above, the generally accepted approach is used. Matrix  $A$ , which is the difference of matrices corresponding to these semantic networks, is considered. And a norm of matrix  $A$  is evaluated as a measure of divergence. The norm of the matrix reflects the order of magnitude of the matrix elements. In this case, it is recommended to use the Frobenius norm  $\|\cdot\|_F$  that is equal to the square root of the sum of squares of all elements of the corresponding matrix:

$$\|A\|_F = \sqrt{\sum_{i,j} n_{ij}^2}$$

Of course, the dimension of the two compared matrices must coincide. In reality, the composition of terms in different semantic matrices differs. Therefore, the compared networks are mutually complementary with terms that are part of their overall composition.

## 2. Example of approbation of the method

The degree of similarity of the texts was determined by the example of biblical texts, which are well-known and translated into almost all languages (in particular, the authors researched texts in Ukrainian, Russian, English and Hebrew). The text of the sacred book of the Torah, the Pentateuch of Moses, was used

to build networks of terms and further research. In particular, the Ukrainian translation by Ivan Ogienko that available at the link [uk.wikisource.org/wiki/](http://uk.wikisource.org/wiki/) was used. The English version of the text is available at the link [www.sacred-texts.com](http://www.sacred-texts.com), the Hebrew is available at [www.ccel.org](http://www.ccel.org), The Russian version of the translation made by Archimandrite Macarius is available at the link [ru.wikisource.org/wiki/](http://ru.wikisource.org/wiki/). In general, all five books «Genesis», «Exodus», «Leviticus», «Numbers» and «Deuteronomy» were researched.

As a result of processing these texts, the networks of terms as ontological models were obtained. Figure 2 shows a fragment of the network of terms that corresponds to the fourth book «Numbers». When processing the Pentateuch of Moses, the specifics of the scriptures were taken into account. For instance, the standard list of stop words was modified at the stage of preliminary processing of texts. As a result, a separate list of exception words that in practice do not refer to stop words was formed and conversely, the list of stop words was supplemented by other words that do not have a semantic load within the researched sacred book. The most frequent synonymous words were researched separately, and as a result, a single definite token was assigned. Also, due to the presence of archaisms in similar sacred texts, some words could be assigned incorrect tags during PoS-tagging. Therefore, this issue requires manual processing.

Globality in the calculation of GTF was determined within the whole book, or within each individual chapter, depending on the text for which the network of terms was built i.e. for the entire book or individual chapter. Therefore, the same terms may have different GTF values within a single section and for the whole text, respectively. These different GTF values affect the build of the horizontal visibility graph.

In order to achieve a slight sparseness of the matrices, the links with a weight equal to 1 were also removed. Next, unconnected nodes with nodes degree equal to 0 were also removed. Such nodes could appear, in particular, after links removing.

All mentioned above affect the topology of networks and lead to the following consequences: the network of terms built for the whole book may contain nodes that do not exist in the network of terms built for a individual chapter, and vice versa — the network of terms for a individual chapter may contain nodes that do not exist in the general network built for all text.

Further comparison of the obtained semantic networks built for different texts with applying the Frobenius measure as comparable approach allows determining the semantic closeness and similarity of the corresponding texts.

The book «Numbers» (the fourth part of the Pentateuch of Moses and the Old Testament) is closest in content to a legal document. This Book contains a census of the adult people of Israel when they were on the Sinai Peninsula and the plains of Moab and regulates the rules of life of these people.





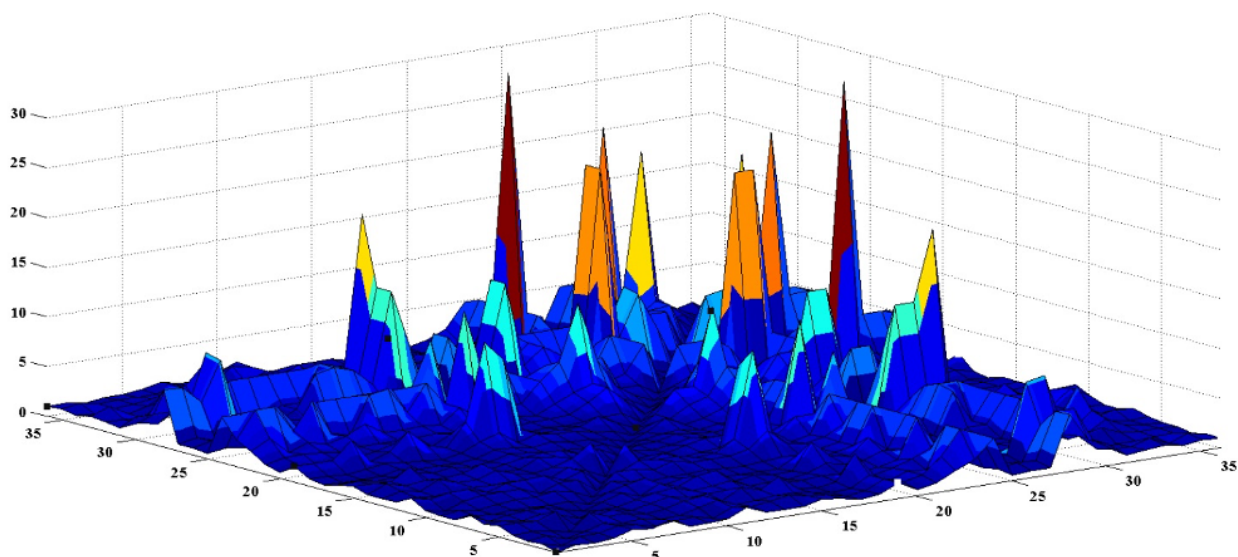


Fig. 4. Graph of semantic matrices differences corresponding to individual chapters of the book «Numbers»

text documents is also considered. This algorithm can also be used in systems of automatic abstracting of legal information in order to generate concise information-rich reports, brief annotations or digests. Also, the proposed method can be used when processing information requests during information retrieval. It allows determining the degree of semantic closeness and similarity or divergence of the text's composition to further determine the relevance of the document to the information needs of Internet users and users of information retrieval systems. As a result, it will increase the pertinency of such systems. Thus, the use of methods for building semantic networks and determining the degree of similarity of texts in modern information retrieval systems and systems of automatic information abstracting (in particular, legal documents) will contribute to the formation and improvement of conceptual and terminological apparatus in the legal field, harmonization of national and international law and, as a result, will promote national security.

## References

1. Mayer-Schönberger V., Cukier K. Big data: A revolution that will transform how we live, work, and think. — Houghton Mifflin Harcourt, 2013.
2. Lande D. V., Dmytrenko O. O., Radziivska O. H. Subject Domain Models of Jurisprudence According to Google Scholar Scientometrics Data. — 2020. — Vol. 2604. — P. 32–43.
3. Ланде Д. В., Дмитренко О. О., Радзівська О. Г. Побудова онтологій в галузі права за даними сервісу Google Scholar. — 2019. — Т. 1, № 28. — С. 74–85.
4. Lande D, Dmytrenko O. Using Part-of-Speech Tagging for Building Networks of Terms in Legal Sphere // CEUR Workshop Proceedings. — 2021. — Vol. 2870. — P. 87–97.
5. Manning C. D., Raghavan P., Schütze H. An Introduction to Information Retrieval. — Cambridge University Press, 2009. — P. 22–36.
6. Santorini Beatrice. Part-of-speech tagging guidelines for the Penn Treebank Project. — 1990.
7. The use of horizontal visibility graphs to identify the words that define the informational structure of a text / Lande D. V., Snarskii A. A., Yagunova E. V., and Pronoza E. V. // 2013 12th Mexican International Conference on Artificial Intelligence / IEEE. — 2013. — P. 209–215.
8. Horizontal visibility graphs: Exact results for random time series / Luque Bartolo, Lacasa Lucas, Ballesteros Fernando, and Luque Jordi // Physical Review E. — 2009. — Vol. 80, no. 4. — P. 046103.
9. From time series to complex networks: The visibility graph / Lacasa L., Luque B., Ballesteros F., Luque J., and Nuno J. C. // Proceedings of the National Academy of Sciences. — 2008. — Vol. 105, no. 13. — P. 4972–4975.
10. Ланде Д. В., Дмитренко О. О., Радзівська О. Г. Визначення напрямків зв'язків у мережі термінів. Інформаційні технології та безпека. — 2019. — С. 103–112.
11. Дмитренко О. О. Побудова направлених зважених мереж термінів із застосуванням Part-of-speech tagging. — 2020. — Т. 22, № 4. — С. 47–55.
12. Lande D., Dmytrenko O. Methodology for Extracting of Key Words and Phrases and Building Directed Weighted Networks of Terms with Using Part-of-speech Tagging // Selected Papers of the XX International Scientific and Practical Conference “Information Technologies and Security”(ITS 2020), CEUR Workshop Proceedings. — 2020. — Vol. 2859. — P. 168–177. — Access mode: <http://ceur-ws.org/Vol-2859/paper14.pdf>.