

Correlation network of concepts determined by the dynamics of publications

Dmytro Lande

*Institute for information recording
National academy of science of
Ukraine*

Kyiv, Ukraine

<https://orcid.org/0000-0003-3945-1178>

Leonard Strashnoy

*Institute Infectious disease department
University of California Los Angeles
Los Angeles, USA*

lstrashnoy@gmail.com

Iryna Balagura

*Institute for information recording
National academy of science of
Ukraine*

Kyiv, Ukraine

<https://orcid.org/0000-0001-9627-2091>

Abstract— A technique for forming, clustering and visualizing so-called correlation networks is herein proposed. The links between nodes of such networks correspond to the values of cross-correlations between vectors - sets of parameters corresponding to these nodes modified in a certain way. To build network structures for each node (topic), vectors are formed - arrays of numbers corresponding to a certain time series. As an example, the article considers a time series generated by the Google Books Ngram Viewer service. This approach, unlike the existing one, has advantages such as intuitive and realistic rules, the definition of the weight of nodes and links, a reliable mathematical basis for correlation analysis, an accounting of previously unused parameters of time series of publications corresponding to entities, allowing one to the group said entities according to their trends in time, and objectivity and relative simplicity. This technique can be based on data obtained, for example, from content monitoring systems, and can be used in analytical systems for various purposes in order to generalize a set of variables without explicit links between them.

Keywords — *cross-correlation network, publication dynamics, Google Books Ngram Viewer, visualization of network structures, cluster analysis*

I. INTRODUCTION

Network analysis is a tool that is widely used in scientometrics to detect scientific communication, the structure of documents, topic relations, and others [1, 2]. One of the most actual tasks of scientometrics is to identify trends in science. The task has significant value for science policy [3]. The most modern method for topic trend detection uses cumulative quantitative analysis of publications, authors, and citations, co-word networks [4-6]. But the big number of papers could reflect as growing of interest to the topic as the period of fading interest. The small amount of papers in the topic sometimes shows their uniqueness. The task of searching scientific trends is still opened because there is no exact methods and tools [7, 8]. Every research topic has own lifecycle. By comparing different topics we could predict future movements of trends and detect hot topics [9, 10]. The presence of one topic in different journals, countries, institutions makes it possible to identify opportunities and priorities for each unit, the strategy of working with data from scientometric analysis [11,12, 13]

We propose to compare scientific topics using publication activity tendencies. The idea is to cluster topics with similar movements. In the paper we propose method which help to classify topics using. The method could be

used in scientometrics and other field for clustering groups of objects.

Modern information technologies can't not be imagined without methods and tools for processing network structures, but the structural features are not always clearly expressed. There is always a question of how to build a network in order to apply a wide range of methods and tools for processing it, to obtain and interpret the results if the researcher has only certain entities – nodes -- at his disposal, but the connections between them are not clearly defined. If a single entity can be represented as a homogeneous multidimensional vector of parameters, it is possible to establish similarity relationships, and apply classification or cluster analysis methods to identify groups of similar documents.

In this paper, a method is proposed that puts the dynamics vector corresponding to the distribution of documents by dates (years) in accordance with the essence (a concept from the subject area). More specifically, each year is assigned a number-the number of times an entity appears in publications covered by the Google Books system. The dimension of this vector corresponds to the number of years and the length of the time interval during which the array of publications was analyzed.

The purpose of this paper is to present a methodology for forming, clustering, and ranking nodes and visualizing so-called directed correlation networks, graph structures, and relationships between nodes (concepts, entities) that correspond to the values of correlations between sets of parameters corresponding to these entities.

At the same time, it should be noted that correlation does not directly mean the presence of causal relationships, so correlation networks cannot be considered as causal, semantic maps. At the same time, correlation, along with other criteria, can be considered as the basis for probabilistic estimates of similarity. In other words, correlation networks can be considered as the basis for applying fuzzy semantic network technologies, for example, for further scenario analysis.

II. METHOD

We use number of publications for each year and form vectors. After constructing of vectors corresponding to individual entities, a correlation network is formed, which can be considered as a way to store and visualize entities that are objectively related to each other [14]. Indeed, it is

possible to form vectors of dynamics for various entities, the relationship between which is not always explicit.

Figure 1 shows a fragment of the interface for obtaining dynamics corresponding to the topic "Artificial Intelligence".

Below, a method is proposed for constructing a network of interconnections of entities (concepts), consisting of the following stages [15, 16]. A request to the Google Books Ngram Viewer service is generated for each entity. The analysis period is also defined - the dimension of the corresponding dynamics vectors.

As a result of performing queries, a set of dynamics vectors corresponding to the given concept is determined, similar to those shown in Fig. 1.

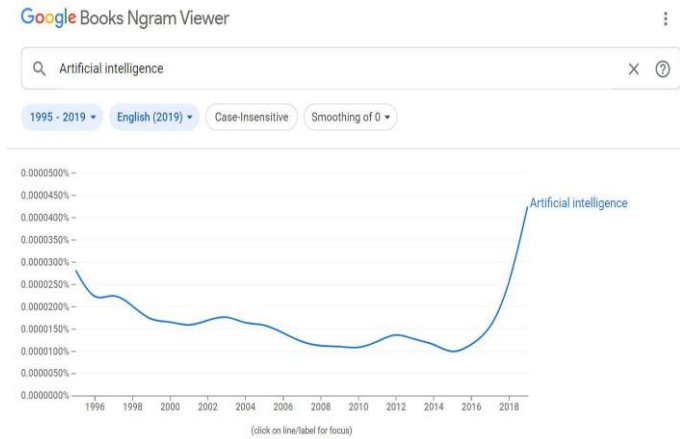


Fig. 1-a Fragment of the interface of the Google Books Ngram Viewer server, where the vector of dynamics of the emergence of the concept of Artificial Intelligence is represented as a graph

The set of maximal cross-correlations between the obtained vectors is calculated, and the corresponding correlation matrix with elements is formed:

$$a_{ij} = \max_m \frac{\sum_{k=1}^{n-m} w_{k+m}^i w_k^j}{\sqrt{\sum_{k=m+1}^n (w_k^i)^2} \sqrt{\sum_{k=1}^{n-m} (w_k^j)^2}}. \quad (1)$$

Here each entity s_k from the multitudes $S = \{s_k\}_{k=1}^{|S|}$, a vector of parameter values is assigned to each entity in the set $\overline{w}^k = (w_1^k, w_2^k, \dots, w_n^k)$, where n is the number of elements in the set of parameters.

The max function is used for the reasons that processes are similar in nature and may have similar dynamic behavior, but it is possible with a time shift.

- The adjacency matrix is formed in accordance with formula (1) and stored in a CSV file. Due to the fact that the adjacency table contains links between all nodes, according to [14], links whose value is less than some selected threshold are ignored. The choice of this threshold completely depends on the experience of analysts. In the information technology

described, a correlation matrix is formed and transmitted for processing and visualization to the network structure analysis system Gephi (<https://gephi.org/>) [17]. Gephi is a fast and simple program for visualization and analysis of network structures, provides the fast layout, efficient data research, as well as visualization of large-scale networks. At the same time, the CSV adjacency matrix for the Gephi system has some features that need to be taken into account (zeros on the diagonal, the location of the characters ";" and others).

- The values of this matrix are sent in CSV format to the Gephi system. This system has a number of modes, among which the "Data Lab" mode is used for monitoring network characteristics. In this mode, in addition to the usual degrees of matrix nodes, you can calculate their values by PageRank, Hits, modularity, and so on. In addition, there are options for ranking matrix nodes (entities) by these parameters.
- Object group modularity classes are defined and the loaded network structure is then clustered [14], [16]. Modularity is calculated as the difference between the fraction of edges within a cluster in the network under consideration and the expected fraction of edges within a cluster in a network where vertices have the same degree as in the original one, but the edges are randomly distributed. The modularity of the network can be expressed by the formula:

$$Q = \frac{1}{2m} \sum_{i,j} \left[a_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j), \quad (2)$$

where a_{ij} is the element of the adjacency matrix A , m is the number of edges in the graph, k_i , k_j are the steps of the corresponding nodes, and δ is the Kronecker Delta function (shows whether the nodes are located i and j in the same modularity class).

- Network visualization is performed in the Gephi system.
- At the last stage, an expert interpretation of the results is performed.

III. IMPROVING THE METHOD

It is proposed to take into account two points when constructing a correlation network, namely: 1) which process started first; 2) absolute values of the time series for mutual correlation, i.e. the value of the directed connection between nodes A and B, is determined in proportion to the sum of the values of the numerical series corresponding to node A.

Let each element s_k from the set of objects $S = \{s_k\}_{k=1}^{|S|}$ be matched with a vector of parameter values $\overline{w}^k = (w_1^k, w_2^k, \dots, w_n^k)$, where n is the number of elements in this set.

To implement point 1 b, the formula for determining the relationship between objects i and j (18) will take the form:

$$a_{ij} = \max_{0 < m \leq K} \frac{\sum_{k=1}^{n-m} w_{k+m}^i w_k^j}{\sqrt{\sum_{k=m+1}^n (w_k^i)^2} \sqrt{\sum_{k=1}^{n-m} (w_k^j)^2}}, \quad (3)$$

where K is the width of the window of possible time offsets.

The max function is used for the following reasons: processes that are similar in nature may have similar dynamic behavior, but it is possible with a time shift. In contrast to the method

described above, accounting m is performed not according to the range of values $[-K, K]$, but in an interval $[1, K]$.

To implement the second point, each of the matrix elements a_{ij} is multiplied by the value of the sum of the

values of the corresponding vector $v_i = C \sum_{k=1}^n w_k^i$, where C

is the normalizing constant. When further using the Gephi visualization tools, the network was defined as non-directional, the node sizes corresponded to the node degrees of the weighted directional network, clustering, if necessary, is calculated using the OpenOrd or Fruchterman Reingold algorithms, and the node modularity is calculated with Resolution = 0.5 (as an extended network, the node size is a weighted power).

IV. EXAMPLES

As a demonstration example, let's consider three entities (Node1, Node2, Node3), each of which corresponds to a time series:

Node1: (0, 1, 2, 3, 4, 5, 4, 3, 2, 1, 0, 0, 0)

Node2: (0, 0, 0, 0, 1, 2, 3, 4, 3, 2, 1, 0, 0)

Node3: (0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 2, 1, 0)

The processes corresponding to these three nodes are shown in Figure 2.

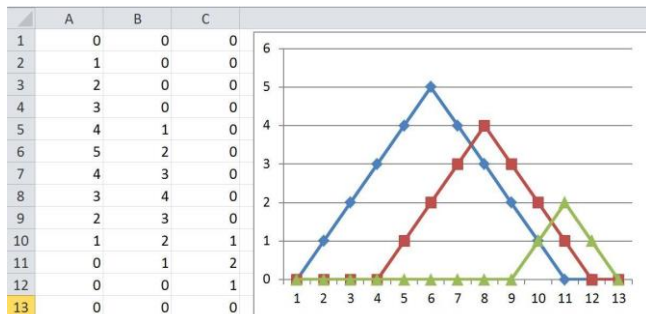


Fig. 2. Processes that match the test example

Visualization of a table corresponding to the correlation matrix calculated using the above algorithm:

```
Node1;Node2;Node3
Node1;0.000;0.818;0.623
Node2;0.818;0.000;0.766
Node3;0.623;0.766;0.000
```

it has the form shown in Fig. 3.

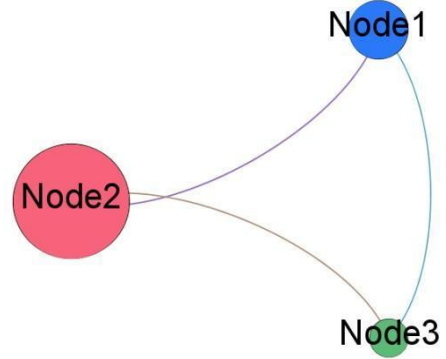


Fig. 3. The correlation network of the example calculated by the algorithm from [15]

In this matrix, node 2 is represented by the largest circle, although it is obvious that the process corresponding to node 1 started earlier and has a larger amplitude.

To correct this discrepancy, the presented improved algorithm allows one to obtain the following matrix of node relationships, the visualization of which is shown in Figure 4:

```
Node1; Node2; Node3
Node1;0.000;1.022;0.779
Node2;0.611;0.000;0.613
Node3;0.002;0.050;0.000
```

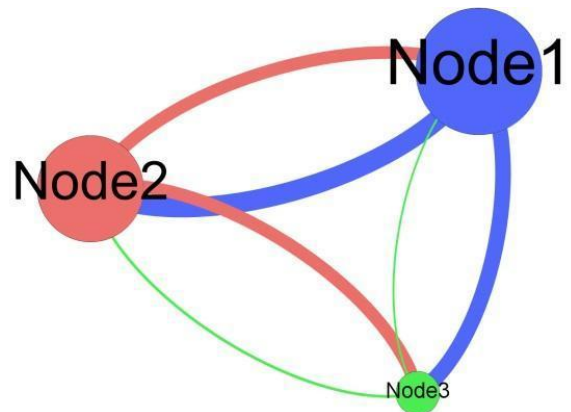


Fig. 4. Directed weighted correlation network of the example calculated using an improved algorithm

V. CONCEPT NETWORK BASED ON GOOGLE BOOKS NGRAM VIEWER

To build a network of concepts related to modern trends in Computer Science, data obtained by accessing the Google Books Ngram Viewer service was considered as an information source. As an example, we consider 20 concepts

listed in Table 1. It also defines the period of analysis (1995-2019 years).

TABLE 1. ENTITIES-REQUESTS TO THE GOOGLE BOOKS NGRAM VIEWER SERVICE

| N Entites | Entity | Abbreviati on |
|-----------|-----------------------------|---------------|
| 1 | Big data | BDT |
| 2 | Complex networks | CNT |
| 3 | Machine learning | MNL |
| 4 | Deep learning | DPL |
| 5 | Neural networks | NNT |
| 6 | Data mining | DTM |
| 7 | Semantic web | SWB |
| 8 | Pattern Recognition | PTR |
| 9 | Complex systems | CST |
| 10 | Artificial intelligence | ARI |
| 11 | Smart grids | SMG |
| 12 | Social computing | SCC |
| 13 | Natural language processing | NLP |
| 14 | Informetrics | INM |
| 15 | Social network analysis | SNA |
| 16 | Information retrieval | INR |
| 17 | Information extraction | INE |
| 18 | Computer vision | CMV |
| 19 | Digital libraries | DLB |
| 20 | Recommender Systems | RSS |

Based on the Google Books Ngram Viewer service, dynamic vectors corresponding to specified concepts are defined, represented in JSON format in the source code of the output form, for example on fig.5:

```

ngrams.data = [{"timeseries": [9.575330750521971e-09, 6.9888743681190135e-09,
1.2791656622823666e-08, 1.1356319440380958e-08, 1.2087312484254653e-08,
1.0711201703372808e-08, 1.2456910170044466e-08, 3.029423822908939e-08, 1.1776428721077536e-08,
7.028031934197543e-09, 8.192412970231544e-09, 6.390932227873236e-09,
8.478197699446355e-09, 7.651633993077667e-09, 6.700836774342633e-09, 5.9736504631757725e-09,
4.666732333902246e-09, 1.3329707115872225e-08, 8.425238284814895e-09,
1.3226189052877755e-08, 1.0391362437189855e-08, 4.0528647105020355e-08, 2.27677522070006e-08,
1.8817276625702142e-08, 1.8356137942987516e-08], "parent": "", "ngram": "Informetrics",
"type": "NGRAM"}, {"timeseries": [6.626123649766669e-08, 4.4549572919549973e-08,
4.416813581542556e-08, 3.641325463377143e-08, 3.4242845083973343e-08, 3.139854243272566e-08,
2.8454546492569554e-08, 2.8115011376858094e-08, 3.27441078695756e-08,
3.3102121932415685e-08, 3.2795835380738936e-08, 3.006990922926889e-08, 2.393094433728038e-08,
2.0940479572573167e-08, 2.411564814508438e-08, 2.3531855575242844e-08,
3.441038387563822e-08, 5.1720220994866395e-08, 3.282066529664007e-08, 2.867479231838388e-08,
3.108477386604136e-08, 3.396015557655119e-08, 4.183678115055045e-08,
4.8998536783528834e-08, 7.350192987587434e-08], "parent": "", "ngram": "Computer vision",
"type": "NGRAM"}, {"timeseries": [2.1733773891696728e-08, 2.3040740870783338e-08,
2.617509764490933e-08, 2.3740966526020202e-08, 1.967496743304764e-08, 1.9771428938497593e-08,
1.961558915297701e-08, 1.84878370396065e-08, 1.9918157789788893e-08,
2.0516173648843505e-08, 1.9428796349529875e-08, 2.2530743493121008e-08,
1.8192798378890984e-08, 1.2451295106975122e-08, 1.1564347701664701e-08,
1.6192725382779827e-08, 1.944471783588142e-08, 1.9834835995879985e-08, 1.9200374623551397e-08,
2.3690152062272318e-08, 2.5081311250119143e-08, 2.8842013221745e-08,
3.20785569984383e-08, 4.34508002911159e-08, 6.21655316562567e-08], "parent": "",
"ngram": "Natural language processing", "type": "NGRAM"}];
    
```

Fig. 5. Vectors in JSON format

As a result of the analysis of 20 concepts of existence, a corresponding weighted correlation matrix was obtained, a network was formed and its clustering was carried out (Fig. 5).

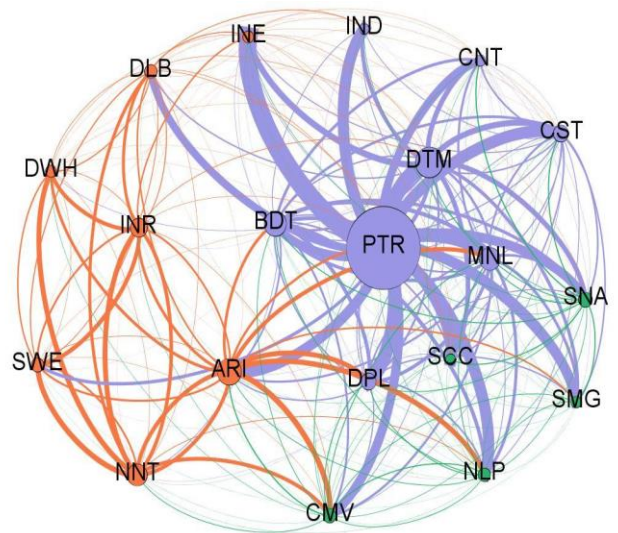


Fig. 6. Network of entities (concepts) in the Gephi system environment

Figure 7-9 shows typical dynamics corresponding to the concepts included in the various clusters shown in Figure 6.

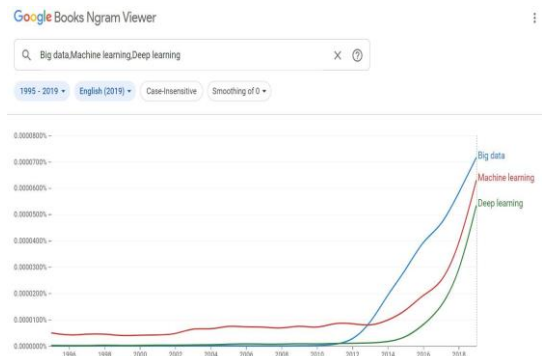


Fig. 7. Entity dynamics (Big data cluster, Machine learning, Deep learning)

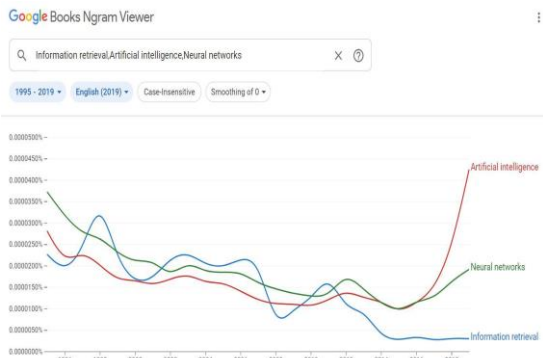


Fig. 8. Entity dynamics (cluster Artificial intelligence, Neural networks, Information retrieval)

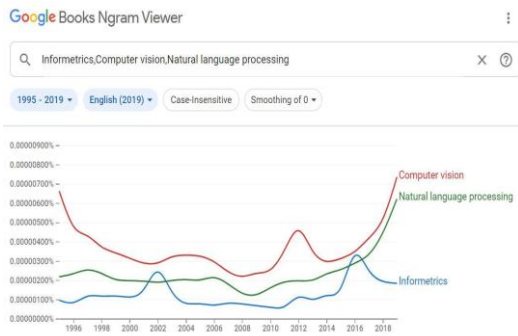


Fig. 9 the Dynamics of entities (the cluster Computer Vision, Natural language processing, Informetrics)

Examples of entities of other types that you can use the developed method for:

1. Political leaders are characterized by their attitude to various spheres of public life.
2. Consumers of products - here are options sellers, the sources of products [14].
3. Mass media as content entities, in this case, the parameters can be words in the headings of articles that are printed in these publications.

VI. CONCLUSIONS

The article proposes the concept of a correlation network determined by the dynamics of its appearance in publications and describes the methodology for its formation, clustering and visualization.

The presented approach, in contrast to the existing ones, has the following advantages:

- intuitive, close-to-reality rules for determining the weight of nodes and links;
- the reliable mathematical basis for correlation analysis;
- accounting for previously unused parameters, time series of publication dynamics that correspond to specific features (topics) and allows you to group entities by their development trends over time;
- objectivity – the dataset is responsible for the "purity" of data;
- the relative ease of implementation (ready-made software systems such as Gephi, the R language, etc. can be used).

The method is demonstrated using data obtained from the Google Books Ngram Viewer system. At the same time, it can be used on other data, for example, obtained from a content monitoring system [15], used in analytical systems for various purposes (for example, medical [16, 18]) in order to generalize a set of variables without explicit links between them.

REFERENCES

- [1] Singh, Punit Kumar & Singh, B. K. Analysis of Social Structures in Scientometrics. In B. K. Singh et al. (Eds.), Academic Libraries: Collection to Connectivity (A Collection of Essays in Honour of Dr. T. N. Dubey) New Delhi: Shree Publishers & Distributors, 2019, pp. 245-254.
- [2] D. Lande, M. Fu, W. Guo, I. Balagura, I. Gorbov & H. Yang "Link prediction of scientific collaboration networks based on information

retrieval", World Wide Web : Internet and Web Information Systems, N 23, pp. 2239-2257, 2020.

- [3] Mills, M.C., Rahal, C. A scientometric review of genome-wide association studies. *Commun Biol* 2, 9 (2019). <https://doi.org/10.1038/s42003-018-0261-x>
- [4] Mohammadreza Kamali, Dina Jahaninfard, Amid Mostafaie, Mahsa Davarazar, Ana Paula Duarte Gomes, Luis A.C. Tarelho, Raf Dewil, Tejraj M. Aminabhavi, Scientometric analysis and scientific trends on biochar application as soil amendment, *Chemical Engineering Journal*, Volume 395, 2020
- [5] Manoj K umar Verma Mapping the Research Trends on Information Literacy of Selected Countries during 2008-2017: A Scientometric Analysis 2019, *DESIDOC Journal of Library & Information Technology*, Vol. 39, No. 3, pp. 125-130
- [6] Zhang, Y., Li, C., Ji, X. et al. The knowledge domain and emerging trends in phytoremediation: a scientometric analysis with CiteSpace. *Environ Sci Pollut Res* 27, pp. 15515–15536, 2020.
- [7] V. Frehe, V. Rugaitis, F. Teuteberg "Scientometrics: How to perform a Big Data Trend Analysis with ScienceMiner", *Lecture Notes in Informatics (LNI), Proceedings - Series of the Gesellschaft fur Informatik (GI)*, 2014
- [8] Amos Darko, Albert P.C. Chan, Xiaosen Huo, De-Graft Owusu-Manu, A scientometric analysis and visualization of global green building research, *Building and Environment*, Volume 149, 2019, pp. 501-511.
- [9] Úbeda-Sánchez, Á.M., Fernández-Cano, A. and Callejas, Z. (2019), "Inferring hot topics and emerging educational research fronts", *On the Horizon*, Vol. 27 No. 2, pp. 125-134.
- [10] Ke He, Junbiao Zhang, Yangmei Zeng, Knowledge domain and emerging trends of agricultural waste management in the field of social science: A scientometric review. *Science of The Total Environment*, Volume 670, 2019, pp. 236-244
- [11] Soffer O, Geifman D. Comparing research topics in European and International Communication Association journals: Computational analysis. *International Communication Gazette*. June 2020. doi:10.1177/1748048520928334
- [12] Mikova N., Sokolova A. Comparing data sources for identifying technology trends, *Technology Analysis & Strategic Management* Volume 31, 2019 - Issue 11, pp. 1353-1367
- [13] Ruiz-Rosero, J., Ramirez-Gonzalez, G. & Viveros-Delgado, J. Software survey: ScientoPy, a scientometric tool for topics trend analysis in scientific publications. *Scientometrics* 121, 1165–1188 2019. <https://doi.org/10.1007/s11192-019-03213-w>
- [14] J. W. Foreman, "Using Data Science to Transform Information into Insight Data Smart", Wiley, 2013
- [15] D.V. Lande, A.O. Snarskii, "Networks determined by the dynamics of thematic information streams" *Data Recording, Storage Processing*, vol. 22, pp. 56-61, 2020.
- [16] D. Lande, L. Strashnoy, "Cross-Correlation of Publications Dynamics and Pandemic Statistics" Available at SSRN: <https://ssrn.com/abstract=3625725> or DOI: <https://dx.doi.org/10.2139/ssrn.3625725> (June 12, 2020). - 9 p.
- [17] K. Cherven. "Mastering Gephi Network Visualization", Packt Publishing, 2015.
- [18] D. Lande, L. Strashnoy. "Directed Correlation Networks, Determined by the Dynamics of COVID-19 Distribution in Various Countries". Available at SSRN: <http://ssrn.com/abstract=3674041>, DOI: <https://dx.doi.org/10.2139/ssrn.3674041> (Posted: 28 Aug 2020). - 7 p.