

# Inverse Transform Sampling for Bibliometric Literature Analysis

Nikolaos P. Bakas

*Research and Development Dpt.  
RDC Informatics  
Athens, Greece  
n.bakas@rdc.gr*

Dionisios Koutsantonis

*Research and Development Dpt.  
RDC Informatics  
Athens, Greece  
d.koutsantonis@rdc.gr*

Vagelis Plevris

*Dept. of Civil and Architectural Engineering  
Qatar University  
Doha, Qatar  
vplevris@qu.edu.qa*

Andreas Langousis

*Department of Civil Engineering  
University of Patras  
Patras, Greece  
andlag@alum.mit.edu*

Savvas A. Chatzichristofis

*Intelligent Systems Lab.  
Neapolis University Pafos  
Pafos, Cyprus  
s.chatzichristofis@nup.ac.cy*

**Abstract**—Scientific literature is prosperously evolving, exhibiting exponential growth in the last decades. For a wide range of scientific thematic areas, it is hard or even impossible for individual researchers to analyze in detail the available published works. For this purpose, we utilize a robust multidimensional scaling procedure, to construct the bibliometric maps of the literature, for keywords, authors and references. Particularly, we propose a generic machine learning algorithm for multidimensional scaling, and describe the algorithmic procedure for the generation of the bibliometric maps.

**Index Terms**—Bibliometrics, Multidimensional Scaling, Optimization

## I. INTRODUCTION

The growth of scientific output is equivalent to a doubling every nine years, as highlighted in [1], [2]. Particularly, using publications obtained from the Web of Science, it was demonstrated that the growth rate of publications was 1% in the 18th century, increasing to 2% to 3% in the period between the two world wars, and 8% to 9% up to 2010. This means that today the total number of published works is twice what it was back in 2013. Henceforth, the writing of a review paper, either a narrative or a systematic one [3]–[7], as well as writing the literature review as a significant part of any research paper, demands from the researcher to analyze a vast amount of research papers rigorously. Hence, new terminology and relevant techniques appear, called bibliometric analysis, bibliometrics, scientometrics, scientific mapping, etc. As a result of applying advanced computer algorithms, a systematic analysis of a wide range of research papers is nowadays possible. This is a modern approach, which can significantly contribute to increasing the quality and broadening the scope of the typical review papers made by individual researchers [8], [9].

The main part of a Bibliometric study is to construct bibliometric maps of the studied topic. A bibliometric map is a visual representation of the solution of the multidimensional

scaling problem [10], [11], which is based on the assembly and further processing of the co-occurrence matrix. Bibliometric maps regard associations of how Bibliometric Objects (BO) are interrelated and appear simultaneously in research papers. This is attained graphically, through their distances in a two-dimensional map which reveals important information about how the studied BO are conceptually linked. By utilizing these computation-based analyses, the conclusions regarding the studied scientific field can be obtained based on an extended database of papers. Furthermore, through a robust computational procedure, the outcomes are documented rigorously.

Other methods such as clustering could also be applied for such a purpose; however, in [12], it is described how it may be misleading, identifying patterns in data that don't actually exist. Utilizing the proposed approach, the multidimensional scaling error is directly computed, as it is the objective function of the problem statement, offering a clear metric of the algorithmic performance. Hence, the researcher can evaluate the accuracy of the final map for the particular problem studied. Earlier versions of the proposed methodology have been utilized in [4]–[7], for the computational analysis of literature in other fields, such as Engineering, Finance, and Management.

The code for the optimization part, is based on the Inverse Transform Sampling (ITSO) optimization algorithm [13], it is written in Julia Language [14], it is open-source and generic, and the user can modify or add modules for serving particular purposes.

## II. MULTIDIMENSIONAL SCALING FOR BIBLIOMETRIC MAPPING

The rendering of a bibliometric map requires the identification of the location of each BO on a two-dimensional space, such that the distances among the objects represent the inter-item dissimilarities. Particularly, a bibliometric map has the following characteristic attributes:

- Each object (e.g. Keyword, Author, or Reference) is represented as a point on the 2D map, with its coordinates on the Cartesian plane.
- The objects with co-occurrences are connected with a line.
- The thickness of the line represents the link strength, which is proportional to the similarity (or co-occurrence) between the objects.
- The distances between the objects are indicators of their dissimilarity.

### A. Baseline formulation

We define as  $\mathbf{c} = c_{ij}$  the elements of the contingency matrix  $\mathbf{c}$ , with

$$[N] = \{1, 2, \dots, N\},$$

the iterator for the number of Bibliometric Objects  $N$ , and  $i, j \in [N]$ . The contingency table  $c_{ij}$  corresponds to the counting of the co-occurrence of objects. Henceforth, it is rational to define the similarity matrix  $\mathbf{s}$ , with elements

$$s_{ij} := \frac{c_{ij}}{\max \mathbf{c}}, \quad (1)$$

and the corresponding dissimilarity

$$ds_{ij} := 1 - s_{ij}, \quad (2)$$

both comprising values in  $[0, 1]$ . Accordingly, we have the pairwise distances of two elements on the map

$$d_{ij} := \|\mathbf{x}_i - \mathbf{x}_j\|_2^2, \quad (3)$$

where  $\mathbf{x}_i$  stands for the  $t$ -dimensional vector, defining the position of the  $i^{th}$  point, on the map. For the case of two-dimensional maps, we use  $t = 2$ , however, the procedure can be generalized for  $t > 2$  as well, for 3-Dimensional maps, or dimensionality reduction.

Ideally, we should have a map, such that

$$ds_{ij} \equiv d_{ij} \quad \forall i, j \in [N], \quad (4)$$

or alternatively,

$$ds_{ij} \propto d_{ij} \quad \forall i, j \in [N]. \quad (5)$$

However, this cannot be attained exactly, for any  $N > 3$ . Henceforth, the problem of finding the bibliometric objects' location on the map, should be formulated as an optimization problem, that is to find the **best possible** coordinates of the bibliometric objects on the map. Furthermore, the formulation of a representative objective function for such a purpose is not unique, and vastly affects the obtained shape of the map (see next section II-B). Moreover, the optimization algorithm used to solve this problem is still an open issue (see below II-C).

### B. Objective Functions

By defining  $f$  as the objective function, the problem of finding the best possible topology of the objects is formulated as finding the optimal  $\mathbf{x}$ , that is to say

$$\mathbf{x}_{opti} = \arg \min f(\mathbf{x}) := \{\mathbf{x} \in A \subseteq \mathbb{R}^t \mid \forall \mathbf{y} \in A : f(\mathbf{y}) \geq f(\mathbf{x})\}, \quad (6)$$

for a  $t$ -dimensional map, and a given domain  $A \subseteq \mathbb{R}^t$ . A number of objective functions were investigated in order to construct a map, such that each  $d_{ij}$  corresponds to the reciprocal  $ds_{ij}$ . One approach to this multidimensional scaling problem [11], can be formulated as finding a  $t$ -dimensional map, where the order of distances of the  $N$  objects, when ranked monotonically at a descending order, corresponds to the monotonically descending order of their dissimilarities [10], [15]. That is to say, higher distances resemble lower similarities. This is achieved through the minimization of the corresponding errors, called stresses. However, this approach does not always offer a consistent representation of the dissimilarities, as it does not utilize a mathematical association for the mapping of the dissimilarities to the distances, further than the monotonicity. Thus, two pairs of points on the map with comparable distances may correspond to dissimilarities with a high difference.

Another approach often used in bibliometric mapping [8], is to minimize the objective function

$$f(\mathbf{x}) = \sum_{i,j \in [N]} s_{ij} d_{ij}, \quad (7)$$

or, equivalently,

$$f(\mathbf{x}) = \sum_{i,j \in [N]} \frac{d_{ij}}{ds_{ij}}. \quad (8)$$

However, this approach does not exploit the significant information where the similarities are equal to zero, because they will be omitted in the summation. Additionally, it does not offer a measurement of the efficiency of the algorithm, as the summation is not comprehensible.

Additionally, we investigated as an objective function, a combined metric, obtained by fitting a regression curve between the dissimilarities and the distances. The regression was nonlinear, considering the dissimilarities to an  $r^{th}$  power, as the independent variable and the corresponding distances as the dependent ones. The regression equation is written as

$$d_{ij} = \alpha ds_{ij}^r + \beta + \varepsilon_{ij},$$

where  $\alpha$  is the regression coefficient,  $\beta$  is the constant term, and  $\varepsilon$  is the regression residuals. Hence, by varying the positions on map  $\mathbf{x}_i$ , the regression coefficients  $a, b$ , as well as the coefficient of determination  $R^2$  also varied. Accordingly, the objective function to be minimized can be written as

$$f(\mathbf{x}) = |\beta(\mathbf{x})| + |1 - R(\mathbf{x})^2| + \frac{|\alpha(\mathbf{x}) - 1|}{3},$$

using the  $\frac{1}{3}$  factor for  $\alpha$ , for normalization purposes. This way, the optimization algorithm is utilized to find a solution satisfying the fundamental regression requirements, for low or zero constant term, high  $R^2$  and close to a unit slope of the regression line. The number of design variables was  $2N + 1$ ,  $2N$  for the points and 1 for  $r$ , as we varied  $r$  during the optimization process, as well. The regression coefficient  $\alpha$  that was divided by 3 relaxed the variation of the slope, as it is important to avoid extreme values, for example, slopes close to zero or tending to infinity. However, this formulation is actually a double optimization problem; one for the map and one for the regression, which is nested into the former. Hence, the execution time was very slow. Additionally, the results were not satisfactory, as the errors, although small, disorientated the mapping procedure.

Accordingly, another idea was to minimize directly, the absolute value of the difference between the distances, and a nonlinear function of the dissimilarities

$$f(\mathbf{x}, \alpha, \beta, r) = |\alpha + \beta ds_{ij}^r - d(\mathbf{x})_{ij}|.$$

Therefore, by including  $\alpha, \beta$  in the optimization problem, the number of variables becomes  $2N + 3$ . However, this approach gave good results only for extreme (high or small) values of  $\alpha, \beta$ , thus the corresponding maps were not supplied information regarding the efficiency of the algorithm.

Additional approaches were investigated. In particular, instead of using the co-occurrence map, the matrix of the chi-squared statistic, as well as the corresponding p-values, were investigated as candidates for similarity matrices. However, the p-values exhibit vastly low values (e.g.  $10^{-100}$  or less) and thus the corresponding distances could not be coherent. Similarly, for the  $\chi^2$ , the values were either equal to zero, or close to 1, not supplying the necessary information to construct the map. Furthermore, instead of minimizing the summation of the errors between the distances and the dissimilarities, the maximum difference, as well as specific quantiles of the distribution of the differences, were also investigated. However, they did not perform well enough, as the minimization problem does not have a solution equal to zero, because of the error always existing in dimensionality reduction problems. Thus, optimizing the maximum or the 75% percentile, does not necessarily lead to the minimization of the other distances, as well.

Another idea was also investigated, inspired by the dropout technique [16], which is common in the training of artificial neural networks, targeting the avoidance of over-fitting. In particular, the technique was used each time the optimizer did not optimize the  $a\%$  percentile (i.e. 75%) but the  $a + \varepsilon$ , where  $\varepsilon$  stands for a randomly generated error number. However, this technique did not give satisfactory results, for the reasons stated above. Other specific techniques, such as linkage or hierarchical clustering, exist for the representation of multiple correlations between  $N$  objects, however, they do not usually offer a coherent, easily and quickly interpretable way of representing the associations between  $N$  objects simultaneously, such as the multidimensional scaling offers.

Finally, after many trials, we decided to utilize the formula

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i,j \in [N]} |ds_{ij} - d_{ij}|. \quad (9)$$

We may also add the constraint

$$\min d_{ij} \geq l, \quad (10)$$

where  $l$  is a specified threshold, to avoid the coincidence of pairs of objects in the final map.

Furthermore, we may easily weigh the elements of the objective function, by multiplying them by the summation of the occurrences for each pair  $i, j$ . This way, we can control the map creation by giving more attention to the most important elements.

### C. Optimization Algorithms

A number of optimization algorithms were investigated for the minimization of the objective function of the optimization problem, including genetic algorithms [17], pattern search [18], particle swarm optimization [19], [20], trust region [21], global search [22], and the Levenberg-Marquardt Algorithm [15], [23], [24], among others. For the trust region and global search methods, the supply of the gradient was investigated, as well. Finally, we utilized the Pure Random Orthogonal Search (PROS), [13], [25] optimization Algorithm, modified to fit the Bibliometric problem, as presented in Algorithm 1.

---

#### Algorithm 1: Bibliometric map generation

---

**Data:** Vector of Strings of the Bibliometric Objects  
**Result:** optimal positions  $\mathbf{x}_{opti} = \mathbf{x}_{i,opti} \forall i \in [N]$  on the Bibliometric Map

Compute co-occurrences  $c_{ij}$  of the studied BO, and maximum iterations  $f_e$ ;

Compute  $s, ds$  from  $c_{ij}$  (Equations 1, 2);

Initialize  $\mathbf{x} \mid \mathbf{x}_{i \in [N]} \in A$  randomly;

Assign  $\mathbf{x}_{opti} \leftarrow \mathbf{x}$ ;

Compute  $f_{opti} = f(\mathbf{x}_{opti})$ ;

**for**  $i \in f_e$  **do**

$i \leftarrow \mathcal{U}(1, N)$ ;

$j \leftarrow \mathcal{U}(1, 2)$ ;

$ra \leftarrow \mathcal{U}(0, 1)$ ;

$x_{ij} = lb_j + ra \times (ub_j - lb_j)$ ;

    Compute  $d_{ij} \forall i, j \in [N]$ ;

    Compute  $f_{run} = f(\mathbf{x})$  (Equation 9);

**if**  $f_{run} \leq f_{opti}$  **then**

$f_{opti} \leftarrow f_{run}$ ;

$\mathbf{x}_{opti} \leftarrow \mathbf{x}$ ;

**else**

$\mathbf{x} \leftarrow \mathbf{x}_{opti}$ ;

**end**

**end**

Plot Bibliometric Map;

**return**  $\mathbf{x}_{opti}$  positions

---

The procedure is generic, for all BO, denoting either keywords, authors, references, or any other object. All the

matrices  $s, c, ds, d$  are symmetric. The vector  $\mathbf{x}$  comprises the design variables of the optimization problem. For a 2-Dimensional map,  $\mathbf{x}$  denotes a set of points in the  $\mathbf{R}^2$  space, that is

$$\mathbf{x}_{i \in [N]} = (x_{i1}, x_{i2}). \quad (11)$$

Ideally and only theoretically, all the errors  $e_{ij} = ds_{ij} - d_{ij}$  should become equal to zero, but this cannot be attained in practice, because the map represents multi-dimensional relationships in the 2D space, leading in dimensionality reduction. Hence, the aim is to find all the  $\mathbf{x}_i$ , corresponding to the minimum errors  $e_{ij}$ .

The algorithm starts with the initial calculation of the co-occurrence table of the studied objects (keywords, authors or references). Accordingly, the similarity and dissimilarity matrices are obtained, as per section II-A. The optimization algorithm initializes randomly the positions of each object, and computes the distances between the objects on the bibliometric map. In the end, the optimal values of the positions are utilized to visualize the results on the bibliometric map.

Furthermore, we may specify a-priori the locations of the top frequent keywords  $[t] = \{1, 2, \dots, t\}$ ,  $\mathbf{x}_{[t]}$ . For this work, we positioned the most frequent keyword in the center of the map (0, 0), and the next 4 most frequent on a unit rectangle around zero. The optimization algorithm, runs after this step accordingly.

### III. RESULTS

In this section, using the proposed framework, we analyze the top 200 cited papers of the recent 5 years (2017-2021), retrieved by a search in Scopus digital library on 24 March 2022. As search keywords, we used the query “Artificial Intelligence”, in the title, abstract and keywords. In order to group similar keywords, we replace the initial ones, with the corrected, as presented in Table I. The results for keywords, authors and references are presented in Figures 1, 2, 3 accordingly. We may see that the procedure yields meaningful results, for keywords, as we see “artificial intelligence”, “deep learning”, and “machine learning” on one edge of the map, along with two other clusters of specific keywords, highly related. Furthermore, in Figure 2, we automatically identify the four major cooperating groups of authors. Similarly, in Figure 3 we see three groups of references occurring simultaneously in many papers, and LeCun 2015 as a major reference associated with almost all the others.

### IV. CONCLUSIONS

We presented a novel mathematical framework, for the construction of a bibliometric map, based on the Inverse Transform Sampling Optimization Algorithm. By using the contingency matrix, we may easily compute the objects’ location and draw the map. Due to the convergence properties of ITSO, the algorithm is not prone to local minima. Henceforth, the resulting Bibliometric Maps, are representative of the actual associations of the bibliometric problems. Furthermore, with the proposed algorithm we may easily control the objective function and thus the shape of the map. An analytic example

for Bibliometric Review on artificial intelligence is presented, for the keywords, authors, and references, along with the computed maps.

As described in the mathematical formulation along with the algorithmic implementation of the suggested framework, the bibliometric mapping procedure is generic, and can be applied in any particular thematic area of research. Henceforth, we foresee applications of this algorithm in other topics, apart from computer science, such as in Engineering, Finance, Learning Systems, etc. Accordingly, we aim to collaborate with researchers from other fields and analyze the corresponding research output, stemming from extended databases of papers, in the rigorous process presented.

### V. ACKNOWLEDGMENTS

The contribution of Andreas Langousis has been conducted within the project PerManeNt, which has been co-financed by the European Regional Development Fund of the European Union and Greek National Funds through the Operational Program Competitiveness, Entrepreneurship and Innovation, under the call RESEARCH – CREATE – INNOVATE (project code: T2EDK-04177).

### NOMENCLATURE

$[N]$	= $\{1, 2, \dots, N\}$ iterator for the N BO. $i, j, \in [N]$ .
$\alpha$	regression coefficient
$\beta$	regression constant term
$\mathbf{x}_i$	$t$ -dimensional vector, defining the position of the $i^{th}$ point
$\varepsilon$	regression residuals
$c_{ij}$	contingency table (co-occurrence of objects)
$f$	objective function to be minimized
$N$	number of Bibliometric Objects studied
$r$	exponent for the nonlinear regression
$t$	number of dimensions of the map. For 2D maps, $t = 2$

### REFERENCES

- [1] R. Van Noorden, “Global scientific output doubles every nine years,” *Nature News Blog*, 2014.
- [2] L. Bornmann and R. Mutz, “Growth rates of modern science: A bibliometric analysis based on the number of publications and cited references,” *Journal of the Association for Information Science and Technology*, vol. 66, no. 11, pp. 2215–2222, 2015.
- [3] C. M. Faggion, N. P. Bakas, and J. Wasiak, “A survey of prevalence of narrative and systematic reviews in five major medical journals,” *BMC Medical Research Methodology*, vol. 17, no. 1, p. 176, dec 2017. [Online]. Available: <https://bmcmedresmethodol.biomedcentral.com/articles/10.1186/s12874-017-0453-y>
- [4] V. Plevris, G. Solorzano, and N. Bakas, “Literature review of historical masonry structures with machine learning,” in *7th ECCOMAS Thematic Conference on Computational Methods in Structural Dynamics and Earthquake Engineering*. Crete, Greece: ECCOMAS, 2019, pp. 1547–1562. [Online]. Available: <https://doi.org/10.7712/120119.7018.21053>
- [5] V. Plevris, N. Bakas, G. Markeset, and J. Bellos, “Literature review of masonry structures under earthquake excitation utilizing machine learning algorithms,” in *6th ECCOMAS Thematic Conference on Computational Methods in Structural Dynamics and Earthquake Engineering*. Rhodes Island, Greece: ECCOMAS, 2017, pp. 2685–2694. [Online]. Available: <https://doi.org/10.7712/120117.5598.18688>

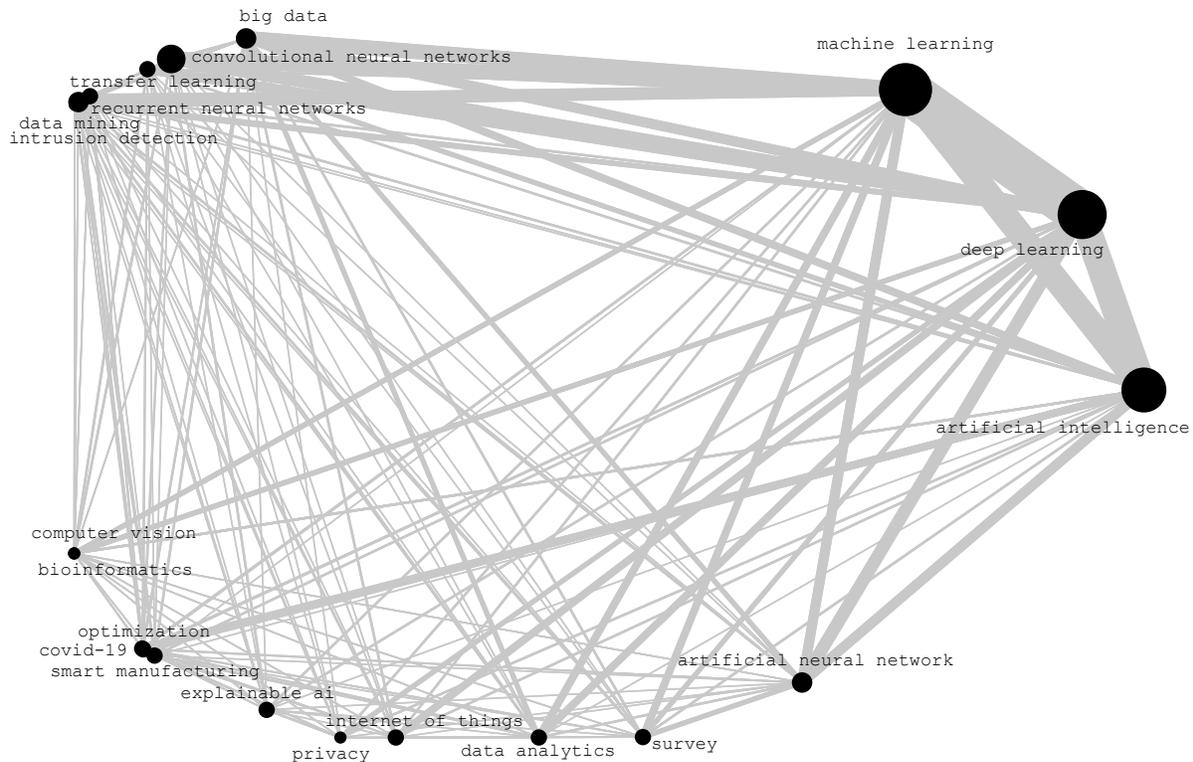


Fig. 1: Association Map of Top 20 Keywords, for Artificial Intelligence, after their Grouping

- [6] M. Papadaki, N. Bakas, E. Ochieng, I. Karamitsos, and R. Kirkham, "Big data from social media and scientific literature databases reveals relationships among risk management, project management and project success," in *6th Annual University of Maryland PM Symposium in May 2019*, pmworldjournal, 2019. [Online]. Available: <https://pmworldjournal.com/article/big-data-from-social-media-and-scientific-literature-databases>
- [7] T. Dimopoulos and N. Bakas, "An artificial intelligence algorithm analyzing 30 years of research in mass appraisals," *RELAND: International Journal of Real Estate & Land Planning*, vol. 2, no. 0, pp. 10–27, 2019. [Online]. Available: <http://ejournals.lib.auth.gr/reland/article/view/6749>
- [8] N. van Eck and L. Waltman, "Software survey: Vosviewer, a computer program for bibliometric mapping," *Scientometrics*, vol. 84, no. 2, pp. 523–538, 2009.
- [9] M. J. Cobo, A. G. López-Herrera, E. Herrera-Viedma, and F. Herrera, "Scimat: A new science mapping analysis software tool," *Journal of the American Society for Information Science and Technology*, vol. 63, no. 8, pp. 1609–1630, 2012.
- [10] R. N. Shepard, "The analysis of proximities: multidimensional scaling with an unknown distance function. i," *Psychometrika*, vol. 27, no. 2, pp. 125–140, 1962.
- [11] J. B. Kruskal, "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis," *Psychometrika*, vol. 29, no. 1, pp. 1–27, 1964.
- [12] N. Altman and M. Krzywinski, "Points of significance: clustering," 2017.
- [13] N. P. Bakas, V. Plevris, A. Langousis, and S. A. Chatzichristofis, "Itso: a novel inverse transform sampling-based optimization algorithm for stochastic search," *Stochastic Environmental Research and Risk Assessment*, vol. 36, pp. 67–76, 2022. [Online]. Available: <https://doi.org/10.1007/s00477-021-02025-w>
- [14] J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah, "Julia: A fresh approach to numerical computing," *SIAM review*, vol. 59, no. 1, pp. 65–98, 2017.
- [15] J. J. Moré, "The levenberg-marquardt algorithm: implementation and theory," in *Numerical analysis*. Springer, 1978, pp. 105–116.
- [16] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [17] A. R. Conn, N. I. Gould, and P. Toint, "A globally convergent augmented lagrangian algorithm for optimization with general constraints and simple bounds," *SIAM Journal on Numerical Analysis*, vol. 28, no. 2, pp. 545–572, 1991.
- [18] C. Audet and J. E. Dennis Jr, "Analysis of generalized pattern searches," *SIAM Journal on optimization*, vol. 13, no. 3, pp. 889–903, 2002.
- [19] Y. Shi and R. Eberhart, "A modified particle swarm optimizer," in *Evolutionary Computation Proceedings, 1998. IEEE World Congress on Computational Intelligence., The 1998 IEEE International Conference on*. IEEE, 1998, pp. 69–73.
- [20] V. Plevris, "Innovative computational techniques for the optimum structural design considering uncertainties," 2009.
- [21] R. H. Byrd, J. C. Gilbert, and J. Nocedal, "A trust region method based on interior point techniques for nonlinear programming," *Mathematical Programming*, vol. 89, no. 1, pp. 149–185, 2000.
- [22] Z. Ugray, L. Lasdon, J. Plummer, F. Glover, J. Kelly, and R. Martí, "Scatter search and local nlp solvers: A multistart framework for global optimization," *INFORMS Journal on Computing*, vol. 19, no. 3, pp. 328–340, 2007.
- [23] K. Levenberg, "A method for the solution of certain non-linear problems in least squares," *Quarterly of applied mathematics*, vol. 2, no. 2, pp. 164–168, 1944.
- [24] D. W. Marquardt, "An algorithm for least-squares estimation of non-linear parameters," *Journal of the society for Industrial and Applied Mathematics*, vol. 11, no. 2, pp. 431–441, 1963.
- [25] V. Plevris, N. P. Bakas, and G. Solorzano, "Pure random orthogonal search (pros): A plain and elegant parameterless algorithm for global optimization," *Applied Sciences*, vol. 11, no. 11, 2021. [Online]. Available: <https://doi.org/10.3390/app11115053>

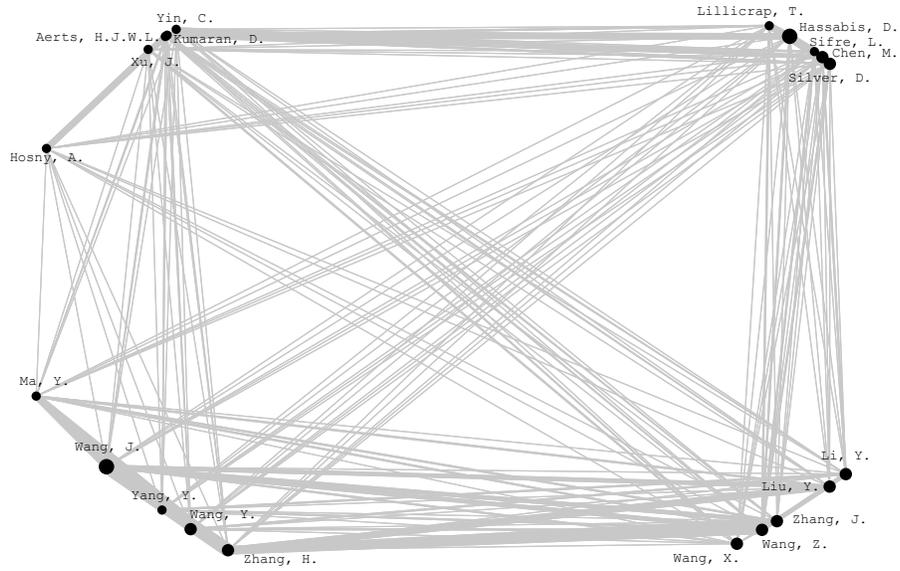


Fig. 2: Authors' cooperating groups. By using Bibliometric Maps, we distinguish 4 main clusters.

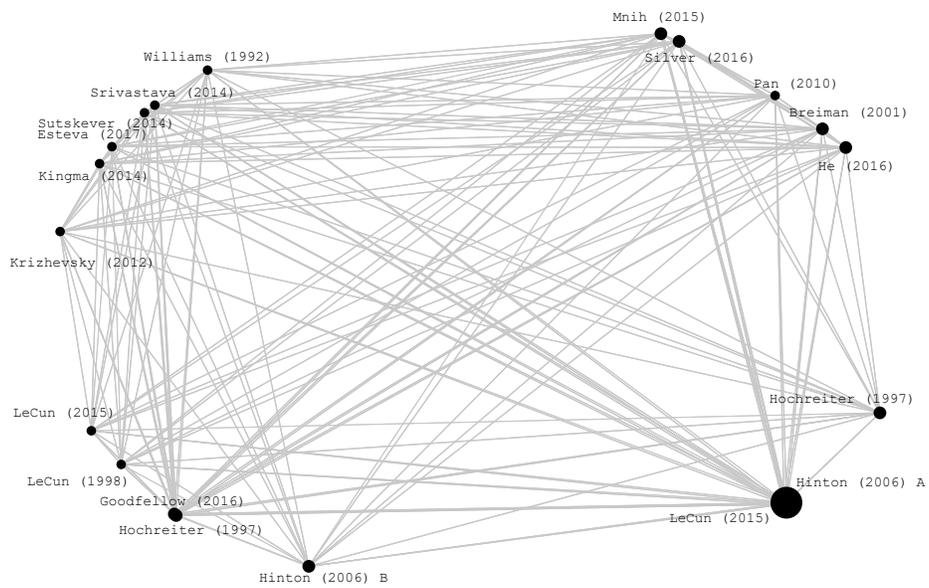


Fig. 3: Association map of top 20 references, after their Grouping.

APPENDIX

TABLE I: 100 Most frequent keywords and their grouping.

Initial String	Corrected String	Frequency	Initial String	Corrected String	Frequency
deep learning	deep learning	39	robots	robots	2
machine learning	machine learning	39	sentiment analysis	sentiment analysis	2
artificial intelligence	artificial intelligence	28	support vector machine	support vector machine	2
big data	big data	8	transparency	explainable ai	2
convolutional neural networks	convolutional neural networks	7	triboelectric nanogenerators	triboelectric nanogenerators	2
data mining	data mining	6	accountability	accountability	1
convolutional neural network	convolutional neural networks	5	age-related macular degeneration	age-related macular degeneration	1
internet of things	internet of things	5	agriculture	agriculture	1
transfer learning	transfer learning	5	asic	asic	1
data analytics	data analytics	4	black-box models	black-box models	1
optimization	optimization	4	business modelling	business modelling	1
smart manufacturing	smart manufacturing	4	chest x-ray images	chest x-ray images	1
survey	survey	4	choroidal neovascularization	choroidal neovascularization	1
artificial neural network	artificial neural network	3	comprehensibility	comprehensibility	1
bioinformatics	bioinformatics	3	computer architecture	computer architecture	1
computer vision	computer vision	3	continual learning	continual learning	1
covid-19	covid-19	3	coronavirus (covid-19)	covid-19	1
interpretability	explainable ai	3	data	data	1
intrusion detection	intrusion detection	3	data fusion	data fusion	1
neural network	artificial neural network	3	data infrastructure	data infrastructure	1
privacy	privacy	3	dataflow processing	dataflow processing	1
adversarial learning	adversarial learning	2	diabetic macular edema	diabetic macular edema	1
ai	artificial intelligence	2	diabetic retinopathy	diabetic retinopathy	1
artificial intelligence (ai)	artificial intelligence	2	energy-efficient accelerators	energy-efficient accelerators	1
automation	automation	2	explainability	explainable ai	1
backpropagation	backpropagation	2	explainable ai	explainable ai	1
big data analytics	big data analytics	2	explanation	explainable ai	1
black-box attack	black-box attack	2	explanations	explainable ai	1
building energy	building energy	2	fairness	fairness	1
classification	classification	2	functional analysis of variance (fanova)	functional analysis of variance (fanova)	1
cognitive computing	cognitive computing	2	gdpr	gdpr	1
communications	communications	2	imbalanced data	imbalanced data	1
computational intelligence	computational intelligence	2	information and communication technology	information and communication technology	1
coronavirus	covid-19	2	long short-term memory (lstm)	long short-term memory (lstm)	1
cyber-physical systems	cyber-physical systems	2	low power	low power	1
cybersecurity	cybersecurity	2	neuromorphic computing	neuromorphic computing	1
decision making	decision making	2	open the black box	black-box models	1
deep neural networks	deep learning	2	optical coherence tomography	optical coherence tomography	1
distributed computing	distributed computing	2	pneumonia	pneumonia	1
explainable artificial intelligence	explainable ai	2	radiology images	radiology images	1
federated learning	federated learning	2	random search	random search	1
governance	governance	2	rare events	rare events	1
interpretable machine learning	explainable ai	2	responsible artificial intelligence	responsible artificial intelligence	1
machine-learning	machine learning	2	screening	screening	1
medical imaging	medical imaging	2	sequence learning	sequence learning	1
neural networks	artificial neural network	2	spatial architectures	spatial architectures	1
omics	omics	2	stability plasticity	stability plasticity	1
pandemic	covid-19	2	synaptic consolidation	synaptic consolidation	1
recurrent neural network	recurrent neural networks	2	transparent models	explainable ai	1
recurrent neural networks	recurrent neural networks	2	vlsi	vlsi	1