The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

# The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map

[1] Yu-Hsiang Yang, [2] Rua-Huan Tsaih, [3] Huimin Bhikshu
[*1,] *Dept. of Management Information Systems, National Chengchi University,*
*yuxiang1001@gmail.com*
[2,] *Dept. of Management Information Systems, National Chengchi University,*
*tsaih@mis.nccu.edu.tw*
[3] *Dharma Drum Buddhist College, General Education, Taipei National University of the Arts,*
*huimin@ddbc.edu.tw*

## *Abstract*

*The purpose of this study was to propose a multi-layer topic map analysis using co-word analysis of informetrics with Growing Hierarchical Self-Organizing Map (GHSOM). The topic map illustrated the delicate intertwining of subject areas and provided a more explicit illustration of the concepts within each subject area. We applied GHSOM, a text-mining Neural Networks tool, to obtain a hierarchical topic map. After taking up one example of altruism in evaluation, we suggest that topic map may disclose some important facts from a whole bunch of data.*

**Keywords***: Topic-map, Co-word, Growing Hierarchical Self-Organizing Map, GHSOM, Altruism*

## 1. Introduction

This study p roposed a h ierarchical mapping model usin g co -word analysis and Growing Hierarchical Sel f-Organizing Map ( GHSOM) [1,2]. Since Price [ 3] first sugg ested t he possibility o f dynamic mapping usi ng the s cientific method, r esearch in bib liometrics and sciento metrics has developed techniques to analyze data sets fro m within publications [4]. Most early work in this field focused on identifying networks (or clusters) of authors, papers, or references. Based on the nature of words, which are important carriers of scientific concepts, ideas and knowledge [5], co-word analysis was also adopted to iden tify se mantic themes [6]. Co-word analysis simplifies and pr ojects data into specific visual representations while maintaining the essential information contained within it.

Noyons [7] suggested that bibliometric mapping of science appeared to have experienced a revival, due to increased i nterest in in formation t echnology, since t he mid-1990s. Man y st udies, su ch as [7,8,9,10,11,12] ha ve a pplied bi bliometric maps us ing c o-word an alysis t o v isualize c ognitive structures, based o n scientific to pics, as well as th e relationships l inking them. In p articular, Noyons and van Raan [11] adopted the Self-Organizing Map (SOM) technique [13], to apply co-word approach to scientific mapping (i.e. the organization of science based topics). Furthermore, Shih et al [14] and Li and Chang [15 ] propose a layer ed knowledge-map using the clustering of k eyterms through GHSO M [1,2]. This is an updated version of SOM, enabling the visualization of hierarchical topic maps.

The objectives of this study were to reveal the major topics or conceptual interrelations of research related to altruis m as an example, in order to gain a better understanding of the quantit ative aspects of recorded data and di scover features of r esearch r elevant t o altruism e mbedded in the SSCI dat abase. Altruistic behavior is a selfless prosocial behavior for the welfare of others. It is also a traditional virtue in many cultures and a core aspect of various religions such as Bu ddhism, Islam, and Christianity and so on. Thus, we adop ted GHSOM in c o-word analysis to cluster the conceptual top ics into a representation of dynamic 2-dimentional interrelated structures within the data.

## 2. Dataset and Method

The dataset used in this study was derived from the SSCI database of the Web of Science, created by the Institute for Scientific Information. It comprehensively indexes ov er 1,9 50 journals across 50

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

social sciences disciplines. It also indexes individually selected, relevant items from over 3,300 of the world's leading scientific and technical journals[1].

An empirical search command was used by "Topic=(altruism) OR Topic=("altruist* behavio*") OR Topic=("helping beh avio*") OR Topic=("prosocial behav io*") r efined by Do cument Ty pe= (ARTICLE O R R EVIEW) "to retrieve d ata related to al truism. The documents spe cifically in cluded articles or reviews in th e study. Book reviews, papers of proceeding, letters, notes, meeting abstracts were not t aken in to con sideration. A total o f 4,271 pap ers pub lished between 1956 and 20 09 we re found.

The study applied co-word analysis with GHSOM to cluster the major topics of a large collection of documents based on research related to altruism, and provide a topical landscape of the field. As with co-citation analysis, co-word analysis has been used to determine the strength of relationships among textual containers, w hether t he co ntainers are full-text documents, their surr ogates, fields w ithin documents (e.g. titles, descriptors), or queries submitted to information retrieval systems. Techniques for th e a nalysis of wo rd co- occurrence are gen erally si milar to those u sed for co-citation a nalysis, consisting of cluster analyses, multidimensional scaling methods [16].

Co-occurrence an alysis of document co ntent is usuall y p erformed on s ubstantive k eywords appearing in a bibliographic database record field such as the title, descriptors, or abstract. These fields encapsulate the topicality o f a document, although keywords from the body of te xt could be used as well [16]. Th e b enefits of co-word analysis can he mixed depen ding on th e app lication such as clustering major top ics o f a lar ge coll ection o f docu ments based on their con tent and prov iding a topical landscape of a field. Many studies, such as [7,8,9,10,11,12] had applied informetric maps using co-word ana lysis to visual ize cogn itive structures, b ased o n scient ific topics, as well as t he relationships linking them.

Co-word an alysis embraces a large number of different methods to determine the clusters of word co-occurrence. For t he purposes of the present study, we choose GHSOM used successfully before in comparable studies to identify distinctive clusters of papers [14,15].

Self-Organizing Map was design ed according to t he concept o f u nsupervised ar tificial neura l networks to process high-dimensional data and provided visual results [11,13,17,18]. However, SOM requires a predefined number of nodes (neural processing units) and implements a static architecture. These nodes result in a representation o f hierar chical r elations with limited capa bility. GHSOM approach w as dev eloped to o vercome these limitations, an d is often a pplied in fi eld the i nformation extraction [1,2,14,15,19]. GHSOM is based on the characteristic of SOM, but it can automatically grow its ow n multi-layer hi erarchical structure, i n wh ich each la yer en compasses a nu mber of SOMs, as shown in Figure 1.

The process o f applying GHSOM to topic analysis is illustrated in Figure 2. The three phases are: the data preprocessing phase; the clustering phase; and the interpreting phase.
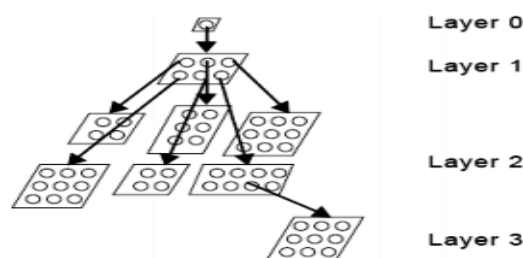


**Fig 1.** Structures of GHSOM [2]

```
+-------------------------------------------+
|          Data preprocessing:              |
|          Determine key-terms              |
+-------------------------------------------+
                    |
                    v
+-------------------------------------------+
|              Clustering:                  |
|      Obtain an acceptable GHSOM result     |
+-------------------------------------------+
                    |
                    v
+-------------------------------------------+
|             Interpreting:                 |
|    Identify the topic categories represented |
|               in GHSOM                    |
+-------------------------------------------+
```

**Fig 2.** The three phases of the topic analysis process

In the data preprocessing phase, key-terms such as titles, keywords, and subject categories are used to represent the contents of the documents. Meaningful key-terms describing the articles are extracted directly from the documents without any manual intervention. These key-terms are weighted according to a *tf* x *idf* the state-of-the-art weighting scheme shown in equation (1) [2,14,16,20].

$$w_i(d) = tf_i(d) \times \log(N / df_i) \tag{1}$$

In equation (1), $w_i(d)$ represents the weight of the $i$th term in document (d), $tf_i(d)$ represents the number of times the ith term appears in document (d), N represents the total number of documents, and $df_i$ represents how many documents contain the ith term. The weighted value for a term will always be greater than or equal to zero. This weighting scheme assigns high values to terms considered important for describing the contents of a document and discriminating between various documents. A high weight is earned by frequent appearances of a term in a given document, with infrequent appearance of terms within the entire collection of documents. In this manner, weight assignment tends to filter out common terms. Based upon weighting values, we selected the top order distinct key-terms for document representation [16,20]. The resulting key-term vectors were used for GHSOM training.

In the clustering phase, the GHSOM experiment[2] was conducted through the trial and error method, using various values for breadth and depth and different normalizations to gain an acceptable GHSOM model for the analysis. The results of GHSOM are shown as Figure 2.

In the interpreting phase, for each node of GHSOM of the first-layer and some nodes of the second-layer which will be re-grouped into the layer 3, we counted the $df_i$ value of each key-term in all articles cluster them into a particular node and assigned a key-term with the highest $df_i$ value (or several key-terms if their $df_i$ values were very close) as the topic category. If there were more than five topics, we would denote it as multidisciplinary. For the remaining nodes, the utmost five important key-terms would be automatically assigned by the GHSOM using the *tf* x *idf* weighting scheme.

## 3. Results

### 3.1. Overview of Productivity

A total of 4,271 papers related to altruism were retrieved from the SSCI database. Figure 2 shows the number of published papers on the topic of altruism, between 1956 and 2009. According to numerical data, a large number of research papers published in recent years (2007-2009) have been catalogued in the SSCI database, with distribution rates of 326 (7.6 %), 416 (9.7 %), and 406 (9.5 %) against the total number of papers, respectively. It has also been observed that a trend in the growth of these numbers appears to have begun in 1991. Figure 3 shows the number of citations of published

---

[2] We used GHSOM toolbox in the Matlab R2007a® package to conduct the GHSOM experiment.

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

papers related to altruism made each year. The figures suggest that the number of these citations has also been increasing. Clearly, the topic of altruism has received a great deal of attention from researchers in the fields of social sciences.

The ten countries ranked as the top publishers of catalogues in the SSCI database are illustrated in Figure 4. The figure shows how the USA is the dominant country, followed by England, Canada and so on.



**Fig 3.** Number of published papers



**Fig 4.** Number of citation (source: ISI Web of Science)

Figure 5 provides the top ten subject areas in which altruism was most widely studied, within the social sciences. The most highly ranked subject area was economics, followed by social psychology and developmental psychology related to altruism. It was also observed that the main part of studies was related to psychology, accounting for over 40 % of total.

Table 3 shows the 10 articles receiving the most citations. The results show how Trivers [21] was an icon in altruism; however, if we take into account the average number of citations per year, the work of Goodman [22] was more influential than that of Trivers [21]. Most of the articles were in the fields of psychology and biology. In addition, Robert Goodman and Ernst Fehr had the two most cited articles, shown in bold text.
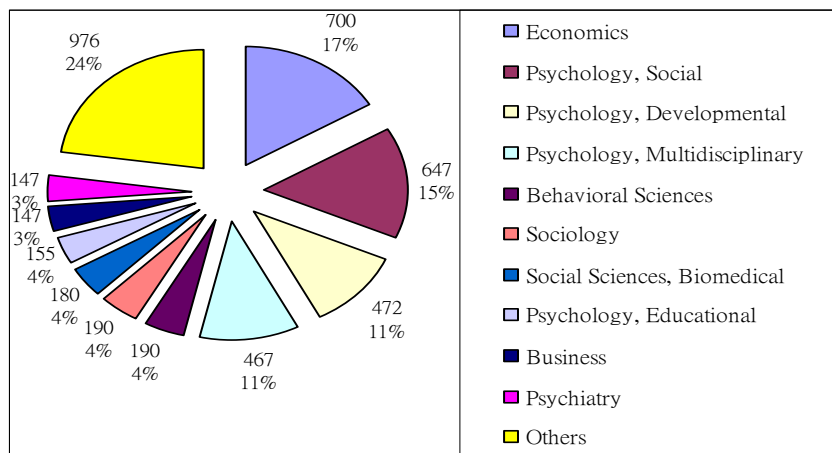
The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

**Fig 5.** Top 10 subject areas for articles related to altruism

**Table 3.** The10 most cited articles (Data retrieved on August 23, 2010)

| Authors | Article | Year | TC | ACPY |
|---|---|---|---|---|
| Trivers, R. L. | Evolution of reciprocal altruism | 1971 | 2,410 | 60 |
| **Goodman, R.** | **The strengths and difficulties questionnaire: A research note** | **1997** | **1,025** | **73** |
| **Fehr, E. and Gachter, S.** | **Altruistic punishment in humans** | **2002** | **591** | **65** |
| Pratto, F., Sidanius, J., Stallworth, L. M. and Malle, B. F. | Social-dominance orientation - a personality variable predicting social and political-attitudes | 1994 | 529 | 31 |
| Andreoni, J. | Impure altruism and donations to public-goods - a theory of warm-glow giving | 1990 | 455 | 21 |
| **Goodman, R.** | **Psychometric properties of the strengths and difficulties questionnaire** | **2001** | **434** | **43** |
| Conner, M. and Armitage, C. J. | Extending the theory of planned behavior: A review and avenues for further research | 1998 | 381 | 29 |
| Krebs, D. L. | Altruism - examination of concept and a review of research | 1970 | 370 | 9 |
| **Fehr, E. and Fischbacher, U.** | **The nature of human altruism** | **2003** | **361** | **45** |
| Colquitt, J. A. | On the dimensionality of organizational justice: A construct validation of a measure | 2001 | 359 | 36 |

*TC: times cited; ACPY: Average Citations per Year*

### 3.2. GHSOM and Topic Analysis

Through the process of applying GHSOM to topic analysis as showed in Figure 2, we obtained the result as showed in Figure 7 in the clustering phase. The model comprised three layers and 56 nodes. All 4,271 articles were clustered into a SOM of 2 x 3 nodes in layer 1, where all articles that had been clustered into the six nodes were further re-grouped into a SOM of 2 x 2 (i.e. node 1, 2, 3, 5, and 6) or 2 x 3 (i.e. node 4) nodes in layer 2, respectively. The articles clustered into nodes 4.1, 4.3 and 6.2 were further re-grouped into a SOM of 2 x 2 nodes in layer 3. The articles clustered into node 4.1.4 were further re-grouped into a SOM of 2 x 2 nodes in layer 4.

In the interpreting phase, for each node of GHSOM, we count the $df_i$ value of each key-term in all articles cluster them into a particular node and assigned a key-term with the highest $df_i$ value (or several key-terms if their $df_i$ values were very close), as the topic category. If there were more than five

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

topics, we wo uld den ote it a s multidisciplinary. The res ults ar e presented in Figures 8, 9, and 10 , in which the number in the parenthesis refers to the number of clustered articles. For instance, there were 283 articles clustered into node 1, and based upo n the inter pretation, it was na med th e "psychology applied in management and busin ess categor y"; 434 articles in node 2 as "the econom ics categ ory", 1308 ar ticles i n no de 3 as th e "psychology, multidiscipline category", 6 03 ar ticles in node 4 as the "economics & multidiscipline cat egory", 51 8 articles in node 5 as t he " psychology, developmental category", 351 articles in node 6 as the "evolution category". Based on these dominant topical clusters in the col lection of articl es, further speci fic topics were o btained i n layer 2, (Figure 9 ). For instance, articles in the "psychology applied in management and business category" were further re-grouped into sub-category topics includin g "organizational", "work predictors", "perfor mance", and "applied psychology" in nod e 1.1; th e sub -category to pics in cluding " management", " business, and applied psychology" in no de 1.2; the sub-catego ry t opics incl uding " management", "performance", "applied psychology", "work", and "job" in node 1.3; and sub-category topics including "donation", "attitudes", "biomedical social sciences", "PEOH (public, en vironmental, an d o ccupational he al th)", an d "transplantation" in node 1.4. Articles in a number of nodes of layer 2 (that is, nodes 3.1, 3.3, 3.5, 4.1, 5.2, and 5.4) were further re-grouped into more specific subcategories in layer 3, as shown in Figure 10.

The interpretation results for the second- and third-layer of GHSOM shown in Figure 9 and 10 were more delicate than those in Figure 5 were. It was observed that the interpretation results for the second-layer were more specific than in the first-layer. For instance, articles in nodes 1.1 and 1.3 belonged to the ca tegory of "psychology a pplied i n management a nd b usiness" in n ode 1 , bu t th ey b oth ha ve further differentiations. Nod e 1.3 focuses on organizational, work, pre dictors, and performance, while node 1.1 focuses on the overall aspect of management and bus iness. An other i nteresting observat ion shown in Figure 9 is that the two n eighboring nodes a re much more closely relate d t han the re mote nodes. For example, articles clustered in node 6.4 at the bottom-right corner of Figure 9 are obviously very different from those cluster ed in node 1.1 in the top-left corner o f Figure 9, bu t they are more closely related to those in nodes 6.1, 6.2 and 6.3.
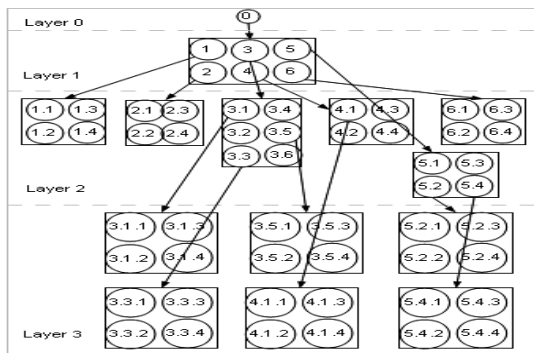

Fig 7. The GHSOM result


Fig 8. First-layer interpretation results of GHSOM.

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

**Fig 9.** Second-layer interpretation result of GHSOM. PSY is the abbreviation for psychology; PEOH refers to public, environmental, and occupational health; SCI is science; MATH is mathematics; ENV is environment.



**Fig. 10.** Third-layer interpretation result of GHSOM. HCSS is the abbreviation for health care sciences and services; SOC refers to social; SCI is science; EDU is education.

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

## 4. Discussion

The results of the GHSOM co mplied with the s ubject area rankings in the first layer, and provided more explicit topics implying the interrelationship of the different subject areas in the second or t hird layers. For example, the sociology in Fi gure 5 is in the node 4.1.2 of Figure 9, indicating that research regarding altruism related to sociology was relevant to economics, social sciences with mathematical methods, and int erdisciplinary ap plied mathematics. The first-layer interpretation results g ive t he disciplinary map while the second- and thir d-layer interpretation results present topic maps indicating the relationship among different disciplines.

Furthermore, the evolution category in node 6 with 351 papers did not appear in Figure 5, where the top ten subject ar eas a re li sted. The category of ev olution in node 6 co-exists with a nu mber of disciplines suc h as anthropol ogy and biolo gy in node 6.1, zool ogy in 6.3, and biologi cal psy chology and biomedical social scienc es i n 6.4 , which i mplies t hat these stu dies w ere i nterdisciplinary an d focused o n evolution. T o be more precise, th e topi cs i n nodes 6 .1, 6.2, 6.3 and 6.4 ex plained wh y evolution had become one of the major clusters in Figure 8. For example, node 6.1 tells us that groups of research associated with anthropology were strongly related to cooperation and reciprocal altruism, which was b ased on the resear ch of biological ev olutionary findings. At the sa me ti me node 6 .2 illustrates how the groups of coo peration and r eciprocity were corr elated wit h the r esearch of group selection i n b iology. I n addition, node 6.4 shows that th e group o f bi omedical soci al s cience researchers targeted cooperation, which is based o n the research of behavioral sciences and biological psychology. Node 6.3 gives us a hint that the beh avioral science s group applies the id eas of cooperation and r eciprocal altruis m in zool ogy. T he above f our gr oups o f research are all b ased on evolutionary c oncepts wh ich have h ad a long and sound histo ry o f sin ce D arwin's g reat di scovery. More specifically, the works such as Trivers [21], Fehr and Fischbacher [23], and Fehr and Gachter [24] in Table 3 could explain the above suggestion, because their articles were prominently cited in research related to altruism. This highlights the need for co-citation analysis in the future.

## 5. Conclusions

To su m up, we observed a steady growth in the number of papers r elated to al truism between the years of 1956 and 2009 in this study. The three most influential authors were Trivers, R., Goodman, R., and Fehr, E . wi th regar d to the n umber of times ci ted. The Good man's pap er is a standard psychometric to ol that has become w idely u sed in psy chology m easuring p rosocial or alt ruistic behaviours. Obviously, Trivers, R is a giant of the field, so it was reassuring he comes out on top!
The study also shows that the variety of research appeared to be scattered across a wide range of subject ar eas, and tha t th e th ree main sub ject ar eas were pri marily wi thin t he f ields o f economics, psycho logy, an d soc iology. However, the GH SOM to ol had all o f t he b enefit o f SOM, i n provid ing a map from a higher dimensional i nput spac e t o a l ower di mensional map space, as well as providing a global orientation of independently growing maps in the individual layers of the hi erarchy, w hich facilitated navigation across b ranches. The t opic map using GHSOM in co-word analysis illustrated the delicate intertwining of subject areas and provided a more explicit il lustration of th e con cepts within each su bject area. T he result of the topic map may indic ate that the co ncept o f ev olution pl ayed an i mportance role in multidiscipline within the research related to altruism. This suggests that topic map may disclose some important facts from a whole bunch of data.

## 6. References

1. Dittenbach M, Rau ber A, Merkl D (2002) Unco vering hierarchical struct ure in data u sing t he growing hierarchical self-organizing map. Neurocomputing 48: 199-216.
2. Rauber A, Merkl D, Dittenbach M (2002) The growing hierarchical self-organizing map: exploratory analysis of high-dimensional data. IEEE Transactions on Neural Networks 13: 1331.
3. Price DJS (1965) Networks of scientific papers. Science 149: 510-515.
4. Leydesdorff L (1987) Various methods for the mapping of science. Scientometrics 11: 295-324.

The Research of Multi-Layer Topic Map Analysis using Co-word Analysis with Growing Hierarchical Self-organizing Map
Yu-Hsiang Yang, Rua-Huan Tsaih, Huimin Bhikshu
International Journal of Digital Content Technology and its Applications. Volume 5, Number 3, March 2011

5. Van Raan AFJ, Tijssen RJW (1993) The neural net of neural network research. Scientometrics 26: 169-192.

6. Boyack K W, Klavans R, Bo rner K ( 2005) Mapp ing t he backb one of s cience. Scient ometrics 6 4: 351-374.

7. Noyons E (2001) Bibliometric mapping of science in a policy context. Scientometrics 50: 83-98.

8. Hassan E (200   3) Si multaneous  mapping  of in teractions betw een s cientific and   technological knowledge  bases: th e cas e  of space c  ommunications. Journal   of t he  American  Society  for Information Science and Technology 54: 462-468.

9. Chau M, Huang Z, Q in J, Zhou Y, Chen H (20 06) Building a scientific knowledge web portal: The NanoPort experience. Decision Support Systems 42: 1216-1238.

10. Gr upp H , Schmoch U, ed itors ( 1992) Percept ion of  scien tification as  measured by  referencing between patents and papers. Heidelberg: Spring Verlag. 73-128 p.

11. Noyons E, van  Raan A (19 98) Moni toring scien tific developments f rom a  dynamic perspectiv e: self-organized  structuring t o  map neural n etwork r esearch. Jou rnal of  the A merican Society f or Information Science 49: 68-81.

12. Ding Y, Chowdhury G, Foo S (2001) Bibliometric cartography of information retrieval research by using co-word analysis. Information Processing & Management 37: 817-842.

13. Kohonen  T ( 1982) Sel f-organized  formation  of  topologically corr ect  feature  maps. Bio logical cybernetics 43: 59-69.

14. Shih J, Chang Y, Che n W (200 8) Using GHSOM to  construct leg al maps fo r Taiwan's securities and futures markets. Expert Systems with Applications 34: 850-858.

15. Li ST, Chang WC (2009) Design And Evaluation Of A Layered Thematic Knowledge Map System. Journal of Computer Information Systems 49.

16. Wolfram D (2003) Applied informetrics for information retrieval research. Westport, Connecticut: Greenwood Publishing Group.

17. Campanario J   (1995) Usin g  neural  networks to study    networks of  scienti fic  journals. Scientometrics 33: 23-40.

18. Kohonen T, K aski S, Lagus K, Salojarvi J, H onkela J, et al. ( 2000) Self organization of a massive document collection. IEEE Transactions on Neural Networks 11: 574-585.

19. Tsaih  R, Lin  W, Hu ang  S (2009) Exp loring Fr audulent Fina ncial  Reporting wit h  GHSOM. Intelligence and Security Informatics: 31-41.

20. Salton G (198  9)  Automatic tex t p rocessing: th e tr ansformation, analys is, a nd retr ieval  of information by computer. Reading, MA: Addison-Wesley.

21. Trivers RL (19 71) EVO LUTION O F RECIPROCAL ALTRUISM. Qu arterly Rev iew of B iology 46: 35-57.

22. Goodman R (1997) The strengths and difficulties questionnaire: A r esearch note. Journal of Child Psychology and Psychiatry and Allied Disciplines 38: 581-586.

23. Fehr E, Fischbacher U (2003) The nature of human altruism. Nature 425: 785-791.

24. Fehr E, Gachter S (2002) Altruistic punishment in humans. Nature 415: 137-140.