We are IntechOpen, the world's leading publisher of Open Access books Built by scientists, for scientists



186,000

200M



Our authors are among the

TOP 1%





WEB OF SCIENCE

Selection of our books indexed in the Book Citation Index in Web of Science™ Core Collection (BKCI)

Interested in publishing with us? Contact book.department@intechopen.com

Numbers displayed above are based on latest data collected. For more information visit www.intechopen.com



Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes

Frizo Janssens^{1,2,3}, Lin Zhang^{1,4} and Wolfgang Glänzel^{1,5} ¹K.U. Leuven, Steunpunt O&O Indicatoren, Dept. MSI, Leuven ²Attentio SA/NV, StudioTROPE building, Bloemenstraat 32, B-1000 Brussels, ³K.U. Leuven, ESAT-SCD, Leuven ⁴WISE Lab, Dalian University of Technology, Dalian ⁵Hungarian Academy of Sciences, IRPS, Budapest ^{1,2,3}Belgium ⁴China ⁵Hungary

1. Introduction

The history of cognitive mapping of science is as long as the history of computerised scientometrics itself. While the first visualisations of the structure of science were considered part of information services, i.e., an extension of scientific review literature (Garfield, 1975, 1988), bibliometricians soon recognised the potential value of structural science studies for science policy and research evaluation as well. At present, the identification of emerging and converging fields and the improvement of subject delineation are in the foreground.

The main bibliometric techniques are characterised by three major approaches, particularly the analysis of citation links (cross-citations, bibliographic coupling, cocitations), the lexical approach (text mining), and their combination. The widely used method of co-citation clustering was introduced independently by Small (1973, 1978) and Marshakova (1973). Although the principle of bibliographic coupling had already been discovered earlier by Fano (1956) and Kessler (1963), coupling-based techniques have been used for mapping the structure of science only decades after co-citation analysis had become a standard tool in visualising the structure of science (e.g., Glänzel & Czerwon, 1996; Small, 1998). Cross-citation based cluster analysis for science mapping has to be distinguished from the previous two methods; while the former two types can be - and usually are - based on links connecting individual documents, the latter approach requires aggregation of documents to units like journals, subject categories, etc., among which cross-citation links are established. The obvious advantages of this method (e.g., the possibility to analyse directed information flows among these units or the assignment/aggregation of units to larger structures) are contrasted by some limitations and shortcomings such as possible biases caused by the use of predefined units. Thus, for instance, Leydesdorff (2006), Leydesdorff and Rafols (2008), and Boyack et al. (2008) used journal cross-citation matrices, while Moya-Anegon (2007) used subject co-citation analysis to visualise the structure of science and its dynamics.

Source: Data Mining and Knowledge Discovery in Real Life Applications, Book edited by: Julio Ponce and Adem Karahoca, ISBN 978-3-902613-53-0, pp. 438, February 2009, I-Tech, Vienna, Austria

Earlier, a completely different approach was introduced by Callon et al., (1983) and Callon, Law and Rip (1986). Their mapping and visualisation tool Leximappe was based on a lexical approach, particularly, co-word analysis. The notion of lexical approach, which was originally based on extracting keywords from records in indexing databases, was later on deepened and extended by using advanced text-mining techniques in full texts (cf. Kostoff et al., 2001, 2005; Glenisson et al., 2005a,b).

Whatever method is used to study the structure of science, cluster algorithms have beyond doubt become the most popular technique in science mapping. The sudden, large interest the application of these techniques has found in the community is contrasted by objections and criticism from the viewpoint of information use in the framework of research evaluation (e.g., Noyons, 2001; Jarneving, 2005). For instance, clustering based on co-citation and bibliographic coupling has to cope with several severe methodological problems. This has been reported, among others by Hicks (1987) in the context of co-citation analysis and by Janssens et al. (2008) with regard to bibliographic coupling. One promising solution is to combine these techniques with other methods such as text mining (e.g., combined cocitation and word analysis: Braam et al., 1991; combination of coupling and co-word analysis: Small (1998); hybrid coupling-lexical approach: Janssens et al., 2007b, 2008). Most applications were designed to map and visualise the cognitive structure of science and its change in time, and, from a policy-relevant perspective, to detect new, emerging disciplines. Improvement of subject-classification schemes was in most cases not intended. Jarneving (2005) proposed a combination of bibliometric structure-analytical techniques with statistical methods to generate and visualise subject coherent and meaningful clusters. His conclusions drawn from the comparison with 'intellectual' classification were rather sceptical. Despite several limitations, which will be discussed further in the course of the present study, cognitive maps proved useful tools in visualising the structure of science and can be used to adjust existing subject classification schemes even on the large scale as we will demonstrate in the following.

The main objective of this study is to compare (hybrid) cluster techniques for cognitive mapping with traditional 'intellectual' subject-classifications schemes. The most popular subject classification schemes created by Thomson Scientific (Philadelphia, PA, USA) are based on journal assignment. Therefore journal cross-citation analysis puts itself forward as underlying method and we will cluster the document space using journals as predefined units of aggregation. In contrast to the method applied by Leydesdorff (2006), who uses the Journal Citation Reports (JCR), we calculate citations on a paper-by-paper basis and then assign individual papers indexed in the Web of Science (WoS) database to the journals in which they have been published. The use of the JCR would confine us to data as available in the JCR and prevent us from combining cross-citation analysis with a textual approach. What is more, proceeding from the document level allows us to control for document types and citation windows, and to combine bibliometrics-based techniques with other methods like text mining. This results in a higher precision since irrelevant document types and 'lowweight journals' can be excluded. This way we can present the results of a hybrid (i.e., combined/integrated) citation-textual cluster analysis to compare those with the structure of an existing 'intellectual' subject classification scheme created and used by Thomson Scientific. The aim of this comparison is exploring the possibility of using the results of the cluster analysis to improve the subject classification scheme in question.

90

1.1 Cognitive mapping vs. subject classification

The objective of the present study is two-fold. The first task is not merely visualising the field structure of science by presenting yet another map based on an alternative approach, but to validate and improve existing subject classifications used for research evaluation. In particular, the question arises of in how far observed 'migration' of journals among science fields can be adopted to improve classification. The second issue is, however, a methodological one, namely to evaluate improved methods of hybrid clustering techniques.

The 22-field subject classification scheme of the Essential Science Indicators (ESI) of Thomson Scientific, which actually forms a partition of the Web of Science universe with practically unique subject assignment, is used as the "control structure". In particular, we propose the following approach in seven steps to solve the integration of cluster analysis and cognitive mapping into subject classification.

- 1. Evaluation of existing subject-classification schemes and visualisation of their crosscitation graph
- 2. Labelling subject fields using cognitive characteristics
- 3. Studying the cognitive structure based on hybrid cluster analysis and visualisation of the cross-citation graph
- 4. Evaluation of science areas resulting from cluster analysis
- 5. Labelling clusters using cognitive characteristics and representative journals suggested by the PageRank algorithm
- 6. Comparison of subject fields and cluster structure
- 7. Migration of journals among subject fields

2. Data sources and data processing

In order to accomplish the above objectives, more than six million papers of the type article, letter, note and review indexed in the Web of Science (WoS) in the period 2002–2006 have been taken into consideration. Citations to individual papers have been aggregated from the publication year till 2006. The complete database has been indexed and all terms extracted from titles, abstracts and keywords have been used for "labelling" the obtained clusters.

Citations received by these papers have been determined for a variable citation window beginning with the publication year, up to 2006, on the basis of an item-by-item procedure using special identification-keys made up of bibliographic data elements extracted from first-author names, journal title, publication year, volume and first page.

In a first step, journals had to be checked for name changes, merging or splitting and identified accordingly. Journals which were not covered in the entire period have been omitted. Furthermore, only journals that have published at least 50 papers in the period under study were considered. A second threshold was used afterwards to remove all journals for which the sum of references and citations was lower than 30. The resulting number of retained journals was 8,305. Most of the subsequent analyses were performed in Java and MATLAB. We also made use of the MATLAB Tensor Toolbox (Bader, 2006).

3. Methods

In this section we briefly describe the methodological background and the algorithms and procedures that have been applied. The first subsection refers to the outlines of the textual approach; this is followed by the description of the cross-citation analysis. The journal

clustering techniques described in the subsequent paragraphs are applied to the textual and citation data separately and used for combined (hybrid) clustering as well. This procedure is described in the following step by step.

3.1 Text analysis

All textual content was indexed with the Jakarta Lucene platform (Hatcher, 2004) and encoded in the Vector Space Model using the TF-IDF weighting scheme (Baeza-Yates, 1999). Stop words were neglected during indexing and the Porter stemmer was applied to all remaining terms from titles, abstracts, and keyword fields. The resulting term-by-document matrix contained nine and a half million term dimensions (9,473,061), but by ignoring all tokens that occurred in one sole document, only 669,860 term dimensions were retained. Those ignored terms with a document frequency equal to one are useless for clustering purposes. The dimensionality was further reduced from 669,860 term dimensions to 200 factors by Latent Semantic Indexing (LSI) (Deerwester, 1990; Berry, 1995), which is based on the Singular Value Decomposition (SVD). The reduction of the number of features in a vector space by application of LSI improves the performance of retrieval, clustering, and classification algorithms. Text-based similarities were calculated as the cosine of the angle between the vector representations of two papers (Salton, 1986).

3.2 Citation analysis

Since the present study analyses the structure of science on the level of journals, all local citations between papers are aggregated to form a journal cross-citation graph. For cluster analysis we ignored the direction of citations by symmetrising the journal cross-citation matrix. At the level of journal clusters, the journal cross-citations can be further aggregated into inter-cluster citations.

From the raw number of cross-citations between two journals (or clusters, respectively), a normalised similarity can be calculated by dividing it by the square root of the product of the total number of citations to or from the first journal (cluster), and the total number of citations to or from the second. Intra-cluster 'self-citations' are counted only once.

For visualisation of the networks we use the similarities just described as edge weights between two clusters or fields (see Figure 2 for an example). For clustering, however, we calculated the similarity of two journals somewhat differently because we didn't want to ignore, for instance, that both journals could be highly cited by a third one. That's why we opted to use "second order" journal cross-citation similarities for clustering. The journal cross-citation numbers are usually stored in a square, symmetric matrix. With "second-order similarities" we mean that the cross-citation values between a journal and all other journals (i.e., row or column of the matrix with cross-citation numbers) are used as input for another step of pairwise similarity calculation. The second-order similarities are found by calculating the cosine of the angle between pairs of vectors containing all symmetric journal cross-citation values between the two respective journals and all other journals. Hence, the ultimate similarity of two journals is based on their respective similarities with all other journals.

The journal cross-citation graph is also analysed to identify important high-impact journals. We use the PageRank algorithm (Brin, 1998) to determine representative journals in each cluster. Besides, the graph can also be used to evaluate the quality of a clustering outcome.

3.3 Clustering

In order to subdivide the journal set into clusters we used the agglomerative hierarchical cluster algorithm with Ward's method (Jain, 1988). It is a hard clustering algorithm, which means that each individual journal is assigned to exactly one cluster.

3.3.1 Number of clusters

Determination of the optimal number of clusters in a data set is a difficult issue and depends on the adopted validation and chosen similarity measures, as well as on data representation. In general, the number of clusters is determined by comparing the quality of different clustering solutions based on various numbers of clusters. Cluster quality can be assessed by internal or external validation measures. Internal validation solely considers the statistical properties of the data and clusters, whereas external validation compares the clustering result to a known gold standard partition. Halkidi, Batistakis and Vazirgiannis (2001) gave an overview of quality assessment of clustering results and cluster validation measures. The strategy that we adopted to determine the number of clusters is a combination of distancebased and graph-based methods. This compound strategy encompasses observation of a dendrogram, text- and citation-based mean Silhouette curves, and modularity curves. Besides, the Jaccard similarity coefficient and the Rand index are used to compare the obtained results with an intellectual classification scheme.

3.3.2 Dendrogram

A preliminary judgment is offered by a dendrogram, which provides a visualisation of the distances between (sub-) clusters (see Figure 4 for an example). It shows the iterative grouping or splitting of clusters in a hierarchical tree. A candidate number of clusters can be determined visually by looking for a cut-off point where an imaginary vertical line would cut the tree such that resulting clusters are well separated. Because of the difficulty to define the optimal cut-off point on a dendrogram (Jain, 1988), we complement this method with other techniques.

3.3.3 Silhouette curves

A second appraise for the number of clusters is given by the curve with mean *Silhouette values*. The Silhouette value for a document ranges from –1 to +1 and measures how similar it is to documents in its own cluster vs. documents in other clusters (Rousseeuw, 1987). The average Silhouette value for all clustered objects (e.g., journals) is an intrinsic measurement of the overall quality of a clustering solution with a specific number of clusters. Since Silhouette values are based on distances, depending on the chosen distance measure and reference data different Silhouette values can be calculated. For instance, we use the complement of cosine similarity applied to text and citation data.

The quality of a specific partition can be visualised in a *Silhouette plot*. In a Silhouette plot (see Figures 1 & 5), the sorted Silhouette values of all members of each cluster (or field) are indicated with horizontal lines. The more the Silhouette profile of a cluster (field) is to the right of the vertical line at the value 0, the more coherent the cluster (field) is, whereas negative values indicate that the corresponding objects should rather belong to another cluster (field).

3.3.4 Modularity curves

The quality of a clustering can also be evaluated by calculating the modularity of the corresponding partition of the cross-journal citation graph (Newman & Girvan, 2004; Newman, 2006). Up to a multiplicative constant, modularity measures the number of intracluster citations minus the expected number in an equivalent network with the same clusters but with citations given at random. Intuitively, in a good clustering there are more citations within (and fewer citations between) clusters than could be expected from random citing. The expected number of citations between two journals is based on their respective degrees and on the total number of citations in the network.

For an additional 'external validation' of clustering results, we also use modularity curves computed from a network containing all journals as nodes, but with edge weights equal to the number of ISI Subject Categories commonly assigned to both journals by Thomson Scientific (out of the total of 254).

3.3.5 Jaccard similarity coefficient and Rand index

The Jaccard index is the ratio of the cardinality of the intersection of two sets and the cardinality of their union. The Jaccard similarity coefficient is an extension of the Jaccard index and can be used as a measure for external cluster validation. The Rand index is another external validation measure to quantify the correspondence between a clustering outcome and a ground-truth categorisation (Jain, 1988). In contrast to the Jaccard coefficient, the Rand index does take into account negative matches as well. Both measures result in a value between 0 and 1, with 1 indicating identical partitions. In Figure 8, we use the Jaccard index to compare each cluster with every field from the intellectual ESI classification, in order to detect the best matching fields for each cluster.

3.3.6 Hybrid clustering

As mentioned at the outset, in general four major approaches are used for clustering sets of scientific papers, particularly, the lexical approach and three citation-based methods, namely cross-citation, bibliographic coupling, and co-citation analysis. Each of the methods alone suffers from severe shortcomings. For example, typical problems with bibliographic coupling and co-citations are sparse matrices, the lack of consensual referencing in some areas (Braam et al., 1991b; Jarneving, 2007), document types with insufficient number of references (e.g., letters) that have to be excluded (bibliographic coupling), the incompleteness due to missing citations to recent years (co-citation analysis), the missing 'critical mass' for emerging field detection (co-citation analysis, cf. Hicks, 1987), and the bias towards high-impact journals (co-citation analysis). If strict citation-based criteria are applied, then the resulting citations-by-document matrix is extremely sparse. In this case, rejection of relationship between two entities (e.g., journals or documents) tends to be unreliable. On the other hand, any lexical (text-based) approach is usually based on rather rich vocabularies and peculiarities of natural language. The result is, according to our observations, a rather 'smooth' or gradual transition between what is related and what is not. Therefore, the relationship is somewhat fuzzy and not always reliable. Hence, both the textual and citation-based approaches provide different perceptions of similarities among the same data. Textual information might indicate similarities that are not visible to bibliometric techniques, but true document similarity can also be obscured by differences in

94

vocabulary use, or spurious similarities might be introduced as a result of textual preprocessing, or because of polysemous words or words with little semantic value. The combination of the two worlds helps to improve the reliability of relationship and therefore of the clustering algorithm as well.

Therefore, the present study combines cross-citation analysis with text mining. The former can be applied to directed links as well as to the symmetrised transaction matrix. Symmetrisation also compensates for the incompleteness caused by the lack of citations to recent years and allows links between journals to be considered strong and subject-relevant even if these are asymmetric or even unidirectional. In order to reduce noise caused by 'small' journals and extremely weak citation links, thresholds have been applied to both citation links and number of papers (see previous section).

The text mining analysis supplements the citation analysis. In particular, the textual information is integrated with the bibliometric information before the clustering algorithm is applied. In the present study, the actual integration is achieved by weighted linear combination of the corresponding distance matrices. The methodology and advantages of hybrid clustering have been substantiated in more detail in earlier studies devoted to the analysis of different research fields (see Glenisson et al., 2005; Janssens et al., 2007a, 2007b, 2008). In addition, the lexical approach allows to 'label' clusters using automatically detected salient terms.

In Section 4.3, Silhouette and modularity curves will be used to compare results of textbased, citation-based and hybrid clustering, and we will substantiate that the hybrid method in general outperforms the other two.

3.4 Multidimensional scaling

Multidimensional scaling (MDS) can be used to represent high-dimensional vectors (for example, the centroids of journal clusters) in a lower dimensional space by explicitly requiring that the pairwise distances between the points approximate the original high-dimensional distances as precisely as possible (Mardia, 1979). If the dimensionality is reduced to two or three dimensions, these mutual distances can directly be visualised. It should, however, be stressed that interpretations concerning such a low-dimensional approximation of very high-dimensional distances must be handled with care.

4. Results

4.1 Evaluation of existing 'intellectual' subject-classification schemes

The multidisciplinary databases *Science Citation Index Expanded* (SCIE) and *Social Sciences Citation Index* (SSCI) of Thomson-Reuters (formerly Institute for Scientific Information, ISI, Philadelphia, PA, USA) traditionally did not provide a direct subject assignment for indexed papers. The annual Science Citations Index Guides, the Journal Citation Reports (JCR) and more recently the Website of Thomson Scientific, however, contain regularly updated lists of (S)SCI journals assigned to one or more subject matters (ISI Subject Categories) each. For lack of an appropriate subject-heading system, more or less modified versions of this Subject Category scheme were often used in bibliometric studies too, namely as an indirect subject assignment to individual papers based on the journals in which they had been published. Such assignment systems based on journal classification have been developed among others

by Narin and Pinski (see, for instance, Narin, 1976; Pinski & Narin, 1976). This was followed by classification schemes developed by other institutes as well. Nowadays two ISI systems are widely used, in particular, the ISI Subject Categories, which are available in the JCR and through journal assignment in the Web of Science as well, and the Essential Science Indicators (ESI).

Field #	ESI Field	Field #	ESI Field
1	Agricultural Sciences	12	Mathematics
2	Biology & Biochemistry	13	Microbiology
3	Chemistry	14	Molecular Biology & Genetics
4	Clinical Medicine	15	Multidisciplinary
5	Computer Science	16	Neuroscience & Behavior
6	Economics & Business	17	Pharmacology & Toxicology
7	Engineering	18	Physics
8	Environment/Ecology	19	Plant & Animal Science
9	Geosciences	20	Psychology/Psychiatry
10	Immunology	21	Social Sciences
11	Materials Sciences	22	Space Science

Table 1. The 22 broad science fields according to the Essential Science Indicators (ESI)

While the first system assigns multiple categories to each journal and is too fine grained (254 categories) for comparison with cluster analysis, the ESI scheme is forming a partition (with practically unique journal assignment) and the 22 fields are large enough. Therefore the ESI classification seems to be a good choice for our analysis.

Subject fields will be considered like automatically generated clusters. One precondition for easy comparison with results from hard clustering is that the reference classification system must form a partition of the WoS universe, while most schemes allow multiple assignments (e.g., the above-mentioned ISI Subject Categories). The only commonly known subject scheme for ISI products that meets the criterion is the ESI classification system. This subject classification scheme is in principle based on unique assignment; only about 0.6% of all journals were assigned to more than one field over a five-year period. For the present exercise, assignment has to be de-duplicated in the case of journals which merged or split up during the period of 5 years, declaredly a somewhat arbitrary procedure. Nonetheless, the assignment remains correct and results in no more than a slightly narrower scope for several journals. The field structure of the ESI scheme is presented in Table 1.

The question arises whether field classification according to the ESI scheme could still be improved. In particular, we will analyse whether journal assignments to fields can be considered optimum. Figure 1 presents the evaluation of the 22 ESI fields based on the cross-citation- (left) and text-based (right) Silhouette values (see Section 3.3.3). Several fields seem not to be consistent enough from both perspectives. Above all, the Silhouette values of field #2 (Biology & Biochemistry), #4 (Clinical Medicine), #7 (Engineering), #19 (Plant & Animal Science) and #21 (Social Sciences) substantiate that at least five of the 22 fields are not sufficiently consistent.

96



Fig. 1. Silhouette plot for 22 ESI fields based on journal cross-citations (left) and based on text (right)

4.2 Labelling subject fields using cognitive characteristics and visualization of the cross-citation network



Fig. 2. Network of the 22 ESI fields based on cross-citation links

Simultaneously to the above validation, the textual approach also provides the best TF-IDF terms – out of a vocabulary of 669,860 terms – describing the individual fields. These terms are presented in Table 2. Although these terms already provide an acceptable characterisation of the topics covered by the 22 fields, considerable overlaps are apparent between pairs of fields, respectively: Engineering (#7) and Computer Science (#5), Chemistry (#3) and Materials Science (#11), Plant & Animal Science (#19) and Environment/Ecology (#8), as well as Biology & Biochemistry (#2), Molecular Biology & Genetics (#14) and Clinical Medicine (#4). In addition, the terms characterising the social sciences (#21) reflect a pronounced heterogeneity of the field. The structural map of the 22 ESI fields based on cross-citation links is presented in Figure 2. For the visualisation we used Pajek (Batagelj & Mrvar, 2002). The network map confirms the strong links we have found based on the best terms between fields #2 & #14, #3 & #11, #5 & #7, and #8 & #19, respectively.

Field	Best 50 terms
1	soil; crop; milk; fruit; seed; cultivar; wheat; dry; rice; ha; chees; diet; fat; ferment; nutrit; meat; farm; grain; starch; fertil; irrig; agricultur; dietari; intak; wine; flour; antioxid; sensori; fatti; sugar; juic; nutrient; moistur; harvest; maiz; veget; cook; leaf; soybean; nitrogen; farmer; season; vitamin; potato; weed; textur; dairi; bacteria; fresh; corn;
2	enzym; dna; receptor; rat; peptid; metabol; lipid; genom; insulin; muscl; transcript; ca2; amino; glucos; mutat; rna; molecul; diabet; kinas; inhibitor; hormon; mice; mrna; neuron; fluoresc; mutant; cancer; assai; serum; vitro; secret; bone; recombin; mitochondri; coli; brain; tumor; ligand; liver; antibodi; subunit; ion; apoptosi; yeast; intracellular; vivo; cholectorol; biologi; offin; calcium;
3	polym; catalyst; crystal; ion; bond; molecul; solvent; atom; ligand; hydrogen; film; polymer; adsorpt; aqueou; poli; nmr; methyl; spectroscopi; thermal; chemistri; bi; electrod; spectra; cu; catalyt; cation; mol; copolym; anion; angstrom; amino; chiral; nm; ir; electrochem; salt; reactor; copper; chlorid; ionic; surfact; aromat; ni; h2o; fluoresc; column; chromatographi; alkyl; cl; alcohol;
4	cancer; therapi; tumor; infect; surgeri; pain; hospit; arteri; syndrom; diabet; injuri; bone; lesion; chronic; symptom; surgic; renal; breast; carcinoma; serum; transplant; lung; mortal; muscl; liver; coronari; cardiac; physician; rat; hypertens; recurr; malign; pulmonari; receptor; oral; men; therapeut; postop; ci; hiv; vascular; mutat; ct; hepat; infant; diagnos; tumour; pregnanc; antibodi; il;
5	web; queri; internet; graph; schedul; wireless; semant; logic; node; busi; video; processor; traffic; execut; fuzzi; server; machin; packet; finit; fault; ltd; grid; hardwar; messag; cach; mesh; xml; multimedia; qo; bandwidth; custom; scalabl; bit; multicast; 3d; iter; java; ip; onlin; metric; platform; polynomi; retriev; neural; circuit; heurist; algebra; robot; topolog; broadcast;
6	firm; price; trade; economi; busi; capit; invest; wage; tax; financi; organiz; incom; bank; compani; sector; corpor; employ; stock; monetari; custom; labor; privat; strateg; welfar; incent; asset; profit; employe; polit; household; game; worker; inflat; job; union; foreign; brand; earn; forecast; labour; reform; export; unemploy; insur; retail; volatil; team; credit; pai; financ;
7	nonlinear; fuzzi; finit; machin; robot; sensor; motion; veloc; nois; crack; thermal; ltd; circuit; vehicl; neural; fuel; voltag; vibrat; elast; beam; shear; turbul; schedul; fault; deform; film; plane; stochast; iter; steel; compress; custom; wind; friction; actuat; concret; logic; soil; geometr; laser; graph; antenna; cylind; traffic; oscil; calibr; autom; geometri; grid; reactor;
8	soil; forest; habitat; river; sediment; ecolog; lake; pollut; land; ecosystem; climat; season; veget; fish; seed; landscap; biomass; nutrient; predat; agricultur; sludg; toxic; groundwat; bird; stream; wast; sea; island; wastewat; wetland; nitrogen; fire; ha; emiss; urban; coastal; flood; biodivers; reproduct; basin; nest; pesticid; seedl; crop; dry; microbi; watersh; graze; winter; rainfal;
9	rock; basin; sediment; sea; fault; ocean; miner; seismic; climat; isotop; earthquak; ic; tecton; ma; soil; southern; volcan; atmospher; mantl; geolog; wind; northern; reservoir; metamorph; precipit; river; cretac; lake; faci; eastern; assemblag; veloc; sedimentari; crust; melt; marin; continent; magma; or; deform; east; flux; granit; belt; fractur; shallow; earth; slope; cloud; clai;

98

Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes

Field	Best 50 terms
10	immun; il; infect; antigen; antibodi; mice; vaccin; receptor; cytokin; hiv; cd4; lymphocyt; ifn; autoimmun; dc; cd8; macrophag; viru; inflammatori; peptid; hla; mhc; tnf; nk; ig; molecul; tumor; lp; serum; tcr; pathogen; innat; assai; chemokin; dendrit; allergen; viral; igg; interleukin; monocyt; apoptosi; neutrophil; epitop; allerg; immunolog; secret; inflamm; dna; vitro; th2;
11	alloi; steel; film; coat; corros; glass; crack; microstructur; ceram; powder; fiber; grain; thermal; sinter; polym; crystal; deform; fabric; weld; fibr; fatigu; concret; fractur; si; specimen; cast; tensil; melt; cement; ni; silicon; shear; bond; microscopi; fe; ion; wear; adhes; cu; copper; nanoparticl; lamin; nanotub; aluminum; compress; roll; elast; creep; atom; al2o3;
12	algebra; theorem; finit; asymptot; infin; manifold; let; polynomi; graph; nonlinear; invari; omega; inequ; singular; lambda; convex; proof; compact; ellipt; conjectur; bar; epsilon; infinit; sigma; phi; symmetr; stochast; hyperbol; banach; topolog; metric; integ; matric; lie; exponenti; markov; curvatur; norm; eigenvalu; kernel; hilbert; cohomolog; geometr; quadrat; covari; dirichlet; semigroup; iter; parabol; theta;
13	infect; bacteria; viru; bacteri; pathogen; dna; genom; pcr; parasit; coli; enzym; mutant; yeast; microbi; viral; hiv; rna; vaccin; immun; encod; virul; antibiot; transcript; sp; assai; escherichia; virus; plasmid; clone; candida; 16; soil; biofilm; antibodi; microorgan; fungal; amino; antigen; bacillu; recombin; fungi; albican; gram; mutat; phylogenet; mice; pseudomona; ferment; rrna; genotyp;
14	dna; chromosom; genom; transcript; mutat; receptor; kinas; mous; mice; rna; allel; mutant; apoptosi; cancer; mrna; rat; phenotyp; muscl; polymorph; embryo; tumor; drosophila; phosphoryl; ca2; neuron; actin; clone; encod; prolifer; mitochondri; enzym; genotyp; vitro; assai; vivo; il; embryon; epitheli; recombin; pcr; chromatin; mammalian; regulatori; linkag; transgen; loci; delet; haplotyp; homolog; yeast;
15	dna; genom; scientist; receptor; brain; soil; climat; earth; molecul; neuron; rna; chromosom; mice; mutat; africa; transcript; biologi; ocean; infect; fossil; india; sea; evolutionari; rock; fuel; logic; southern; island; enzym; marin; insect; fluoresc; cancer; quantum; sediment; scienc; bone; viru; australia; immun; ecolog; fish; china; atmospher; your; mind; rat; bird; ic; colour;
16	neuron; brain; rat; receptor; cortex; motor; cognit; cortic; cerebr; mice; neural; stroke; sleep; nerv; lesion; synapt; seizur; epilepsi; axon; schizophrenia; hippocamp; spinal; symptom; pain; alzheim; hippocampu; dopamin; injuri; parkinson; neurolog; deficit; syndrom; eeg; nervou; sensori; stimuli; dementia; ms; stimulu; glutam; muscl; nucleu; astrocyt; chronic; gaba; frontal; sclerosi; auditori; cord; alcohol;
17	rat; receptor; inhibitor; toxic; therapeut; cancer; metabol; vitro; mice; liver; pharmacokinet; oral; therapi; pharmaceut; enzym; antagonist; assai; vivo; pharmacolog; dna; tablet; inflammatori; tumor; metabolit; lipid; brain; agonist; diabet; cytotox; antioxid; kinas; lung; peptid; apoptosi; ca2; serum; administ; molecul; potent; chronic; insulin; mug; mum; liposom; p450; renal; hepat; inhibitori; immune; ligand;

Data Mining and Knowledge Discovery in Real Life Applications

Field	Best 50 terms
18	quantum; laser; film; beam; spin; atom; scatter; crystal; ion; nonlinear; excit;
	photon; lattic; nois; thermal; oscil; dope; symmetri; veloc; emiss; finit; decai;
	spectra; wavelength; si; diffract; neutron; nm; plane; acoust; fiber; hole;
	superconduct; motion; spectral; dielectr; collis; coher; glass; semiconductor;
	neutrino; perturb; detector; algebra; elast; soliton; waveguid; relativist; amplitud;
	alloi;
19	fish; dog; egg; forest; genu; breed; habitat; seed; infect; diet; sp; season; larva;
	reproduct; leaf; bird; nest; hors; cow; soil; predat; sea; cat; taxa; flower; fruit;
	veget; parasit; pig; milk; seedl; prei; mate; shoot; cattl; southern; trait; genera; fed;
	island; nov; ecolog; lake; insect; pollen; viru; river; juvenil; farm; pathogen;
20	psycholog; cognit; emot; student; mental; adolesc; anxieti; symptom; school;
	item; child; psychiatr; gender; sexual; attitud; cope; mother; interview;
	schizophrenia; suicid; skill; questionnair; belief; abus; therapi; men; word;
	psychotherapi; aggress; mood; verbal; teacher; cue; stimuli; satisfact; judgment;
	job; infant; development; violenc; trait; ptsd; stimulu; style; interperson; peer;
	prime; esteem; distress; recal;
21	polit; student; school; teacher; gender; urban; nurs; court; reform; war; legal;
	discours; profession; parti; disabl; interview; capit; rural; attitud; child; ethnic;
	privat; welfar; democraci; democrat; ethic; employ; justic; feder; violenc; worker;
	agenc; teach; sexual; economi; incom; academ; immigr; sociolog; moral; african;
	skill; mental; librari; men; sector; land; crime; china; civil;
22	star; galaxi; solar; orbit; radio; telescop; emiss; stellar; veloc; disk; galact; earth;
	planet; flux; atmospher; satellit; wind; mar; cosmic; binari; cloud; flare; dust;
	spectral; luminos; redshift; jet; accret; dwarf; planetari; cosmolog; mission;
	motion; observatori; burst; spectra; photometr; gravit; comet; sun; bright; infrar;
	grb; shock; ngc; dark; supernova; spacecraft; radial; halo;

Table 2. The best 50 TF-IDF terms describing the 22 ESI fields

4.3 Cluster analysis: text-based, citation-based and hybrid

Figure 3 compares the performance of text-based, cross-citation and hybrid clustering by several evaluation methods, for various numbers of clusters. For each of the three clustering types, Figure 3(1) presents for various cluster numbers (2 to 30) the modularity calculated from the journal cross-citation graph. Since this evaluation is based on cross-citation data, it is not a surprise that the text-only clustering provides worse results than cross-citation clustering, which performs best here. However, very interesting to note is that the hybrid clustering (integrated text and cross-citation information) provides results highly comparable to those from cross-citation clustering, especially for 7 or more than 12 clusters. The modularity scores for cross-citation clustering indicate that any number of clusters larger than 9 is acceptable. On the other hand, the modularity curve for text-only clustering contains a maximum for eight clusters.

In Figure 3(2), Silhouette curves based on (the complement of) cross-citation values show the somewhat counter-intuitive but beneficial result that hybrid clustering always performs better than cross-citation clustering, although the evaluation only considers citations here. This again demonstrates the power of hybrid clustering: the combined heterogeneous

100

citation-textual approach is superior to both methods applied separately. Nevertheless, this figure does not provide a clear clue with respect to the number of clusters to choose.



Fig. 3. Performance evaluation of text-based, citation-based and hybrid clustering based on (1) modularity calculated from the journal cross-citation graph, and based on Silhouette curves calculated from (2) journal cross-citations, (3) second-order journal cross-citations, (4) text-based distances, and (5) linearly combined distances. For an additional 'external validation' of clustering results compared to ISI Subject Categories, the lower-right figure (6) uses modularity computed from a network containing all journals as nodes, but with edge weights equal to the number of ISI Subject Categories commonly assigned to the corresponding journals by ISI/Thomson Scientific (out of the total of 254 categories).

Silhouette curves based on the complement of second-order cross-citations are shown in Figure 3(3). Again, the hybrid clustering almost always performs best.

In Figure 3(4), the Silhouette values are computed only from textual distances. Naturally, the citation-based clustering performs worst here, while the integrated clustering scores almost as good as the text-only clustering and for some cluster numbers even better.

Figure 3(5) shows Silhouette curves based on linearly combined text-based and citationbased distances (with equal weight). Here, combined data and mere citations give comparable results, which might be an indication that there is a preponderance of citation over text data in the combined Silhouette values.

Finally, Figure 3(6) provides an *external validation* of clustering results by expert knowledge available in the ISI Subject Categories assigned to journals by ISI/Thomson Scientific. The modularity curves are computed from a network containing all journals as nodes, but with edge weights equal to the number of ISI Subject Categories in common (out of the total of 254 categories). Again very interesting to see is that hybrid clustering outperforms text-only and citation-based clustering. The optimal number of clusters according to this type of evaluation is 7.

	Modularity based on journal cross- citation graph	Modularity based on common ISI Subject Categories	MSV based on textual distances	MSV based on 2 nd order journal cross- citations	MSV based on linearly combined distances	Rand index with 22 ESI fields as reference classification
22 ESI fields	0.47533	(0.52604)	0.057237	0.016017	0.062807	(1)
22 citation- based clusters	0.54676	0.44244	0.09319	0.057337	0.18938	0.90463
22 text- based clusters	0.50451	0.45091	0.11829	0.035447	0.12987	0.90582
22 Hybrid clusters	0.54677	0.48839	0.1206	0.05453	0.18951	0.90867

Table 3. Evaluation of 22 ESI fields and 22 citation-based, text-based and hybrid clusters by modularities and mean Silhouette values (MSV). Highest values in each column are shown in bold.

In Table 3 we compare the quality of the partition of 22 ESI fields with the quality of the 22 clusters resulting from citation-based, text-based and hybrid clustering. The only evaluation measure for which the 22 human-made ESI fields score best is modularity based on ISI Subject Categories. As already explained before, this evaluation type computes modularity from a network containing all journals as nodes and with edge weights equal to the number of ISI Subject Categories commonly assigned to the corresponding journals by ISI/Thomson Scientific (out of the total of 254 categories). Since there is a direct correspondence between the 22 ESI fields and these 254 Subject Categories (a field is an aggregation of multiple subject categories), it is not at all surprising (not to mention unfair) that the ESI fields outperform the clusters for this type of evaluation. For all other data-driven evaluation types it is clear that automatic clustering does better than human expert classification.

Hybrid clustering always performs at least as good as text-based or citation-based clustering, except for evaluation by second order cross-citations. However, small the difference, the last column shows that the 22 hybrid clusters correspond best to the 22 ESI fields. It should be noted that the values in Table 3 can differ somewhat from the values in Figure 3 because, for the sake of a fair comparison with ESI fields, in the table only 7729 journals were considered for which a field assignment was available.

4.4 Evaluation of hybrid clusters

The cluster dendrogram shows the structure in a hierarchical order (see Figure 4). We visually find a first clear cut-off point at three clusters, a second one around seven, and 22 clusters also seemed to be an acceptable/ appropriate number. This value coincides with the number of fields according to the ESI classification scheme. The Silhouette plots in Figure 5 and the mean Silhouette values in Table 3 substantiate that the 22 hybrid clusters are furthermore acceptable for both the citation and the text-mining approach. The same conclusion can be drawn from computed modularity scores.

The number of three clusters results in an almost trivial classification. Intuitively, these three high-level clusters should comprise natural and applied sciences, medical sciences, and social sciences and humanities. The solutions with 3 and 22 clusters will be analysed in more detail in Section 4.5. The solution comprising of seven clusters results in a non-trivial classification. The best TF-IDF terms (see Table 5) show that three of these clusters represent the natural/applied sciences, whereas two classes each stand for the life sciences and the social sciences and humanities. This situation is also reflected by the cluster dendrogram in Figure 4. A closer look at the best TF-IDF terms reveals that social sciences cluster (#1 of the 3-cluster solution) is split into the cluster #1 (economics, business and political science) and #6 (psychology, sociology, education), the life-science cluster (#3 in the 3-cluster scheme) is split into clusters #3 (biosciences and biomedical research) and #7 (clinical, experimental medicine and neurosciences) and, finally, the sciences cluster #2 of the 3-cluster scheme is distributed over three clusters in the 7-cluster solution, particularly, the cluster comprising biology, agriculture and environmental sciences (#2), physics, chemistry and engineering (#4) as well as mathematics and computer science (#5).



Fig. 4. Cluster dendrogram for hybrid hierarchical clustering of 8305 journals, cut off at 22 clusters on the left-hand side. Two other vertical lines indicate the cut-off points for 7 and 3 clusters.

The hybrid, i.e. the combined citation-textual based clustering yields acceptable results (see Figure 5), and is distinctly superior to both methods applied separately. Nonetheless, we must not conceal that we can also find clusters of lesser quality, notably cluster #1, in the hybrid classification.



Fig. 5. Evaluation of the hybrid clustering solution with 22 clusters by citation based Silhouette plot (left), text based Silhouette plot (centre) and the plot with Silhouette values based on combined data (right).

4.5 Cognitive characteristics of clusters

As already mentioned in the previous section, another nice point to cut off the dendrogram is at three clusters (cf. the right-most vertical line in Figure 4). Although this refers to a rather trivial case, it might be worthwhile to have a look at term representation of this structure before we deal with 'labelling' the 22 clusters that we have obtained from the hybrid algorithm. This will also help us to understand the hierarchical architecture of the subject structure of science. Table 4 lists the best 50 terms for each of the three top-level clusters which definitely confirm the presence of the expected clusters. Indeed, cluster #1 comprises the social sciences, cluster #2 the natural and applied sciences and cluster #3 the medical sciences. The distribution of journals over clusters is surprisingly well-balanced.

Cluster (# journals)	Best 50 terms
1	polit; student; school; firm; cognit; psycholog; war; gender; price; emot;
(n=2144)	mental; capit; teacher; trade; economi; reform; adolesc; child; busi; discours;
	attitud; urban; skill; court; organiz; moral; text; employ; privat; interview;
	narr; profession; sexual; parti; legal; incom; english; job; music; anxieti;
	invest; german; welfar; academ; belief; write; sector; violenc; religi; teach
2	soil; finit; film; nonlinear; thermal; ion; crystal; algebra; polym; ltd; forest;
(n=3447)	atom; veloc; sediment; laser; quantum; motion; graph; theorem; seed; alloi;
	asymptot; deform; sea; fish; bond; coat; grain; sensor; beam; polynomi;
	hydrogen; fiber; fault; machin; season; emiss; crack; fuzzi; shear; habitat; nois;
	steel; dry; plane; fe; catalyst; elast; sp; glass
3	cancer; infect; therapi; tumor; receptor; rat; dna; pain; diabet; mice; bone;
(n=2714)	brain; muscl; hospit; syndrom; chronic; injuri; mutat; surgeri; serum; lesion;
	arteri; neuron; immun; liver; hiv; il; symptom; antibodi; metabol; inhibitor;
	renal; enzym; breast; surgic; lung; therapeut; mortal; vaccin; genom;
	transcript; nurs; assai; transplant; inflammatori; peptid; insulin; cardiac;
	carcinoma; oral

Table 4. Best 50 TF-IDF terms describing the 3 top-level clusters

According to the terms, economics, business and psychology are the dominant issues in the first cluster which represents the social sciences. The most characteristic terms of the second cluster represent the full spectrum of the sciences including mathematics, geosciences and engineering. Also some subfields of agriculture & environment are covered. Cluster #3, finally, covers biosciences, biomedical research, clinical & experimental medicine and neurosciences.

Cluster (# journals)	Best 50 terms
1 (n=1384)	polit; firm; war; price; trade; economi; capit; busi; reform; urban; court; parti; gender; privat; invest; organiz; sector; corpor; employ; moral; labor; legal; incom; financi; discours; tax; music; compani; contemporari; welfar; essai; union; foraign; domograpi; job; land; uage; givil; chine; labour; back; parre
2 (n=1264)	union; foreigh; democraci; job; fand; wage; civit; china; fabour; book; harr; worker; democrat; german; school; liber; internet; text; religi soil; forest; sediment; fish; seed; habitat; sea; season; river; lake; sp; basin; rock; genu; veget; crop; leaf; climat; southern; ecolog; egg; land; ocean; fruit; dry; island; biomass; northern; miner; nutrient; predat; marin; reproduct; nest; larva; bacteria; taxa; winter; cultivar; ha; nitrogen; ecosystem; seedl; eastern; ic; atmospher; flower; brood; wheat; bird
3 (n=1558)	cancer; infect; tumor; receptor; dna; rat; therapi; mice; mutat; immun; il; antibodi; liver; serum; genom; enzym; transcript; hiv; diabet; assai; inhibitor; viru; antigen; vaccin; peptid; apoptosi; metabol; carcinoma; lung; renal; chromosom; bone; kinas; breast; vitro; chronic; muscl; mrna; therapeut; transplant; syndrom; insulin; dog; inflammatori; hepat; lesion; rna; pcr; diet; molecul
4 (n=1334)	film; ion; crystal; polym; thermal; atom; alloi; laser; bond; coat; quantum; beam; steel; hydrogen; catalyst; crack; glass; fiber; molecul; nm; spectroscopi; spectra; veloc; ltd; finit; cu; vibrat; solvent; deform; electrod; shear; powder; spin; elast; fabric; adsorpt; si; nonlinear; excit; sensor; fuel; fe; poli; polymer; diffract; emiss; aqueou; ni; nmr; corros
5 (n=849)	algebra; finit; nonlinear; graph; theorem; asymptot; polynomi; fuzzi; infin; manifold; let; invari; stochast; schedul; inequ; convex; robot; singular; proof; logic; omega; machin; iter; topolog; nois; traffic; infinit; metric; motion; lambda; web; compact; epsilon; neural; integ; circuit; symmetr; ellipt; bar; fault; node; matric; geometr; markov; sigma; exponenti; queri; custom;
6 (n=760)	wireless; video student; school; cognit; psycholog; teacher; mental; adolesc; emot; child; symptom; anxieti; gender; psychiatr; skill; attitud; abus; teach; item; word; interview; disabl; mother; schizophrenia; sexual; alcohol; speech; instruct; belief; cope; english; profession; questionnair; suicid; violenc; classroom; verbal; youth; academ; peer; therapi; men; development; semant; stimuli; discours; linguist; phonolog; deficit; infant; offend
7 (n=1156)	pain; therapi; hospit; injuri; arteri; nurs; brain; surgeri; neuron; symptom; physician; syndrom; muscl; bone; diabet; rat; lesion; coronari; chronic; stroke; cancer; mortal; cardiac; surgic; receptor; infect; nerv; hypertens; men; infant; implant; cognit; ct; ey; cerebr; smoke; pregnanc; fractur; tumor; mri; cardiovascular; elderli; ci; motor; spinal; sleep; oral; questionnair; myocardi; vascular

Table 5. Best 50 TF-IDF terms describing the 7 top-level clusters

The 50 best TF-IDF terms describing the 22 hybrid citation–lexical clusters are listed in Table 6. Cluster #1 of the 3-cluster scheme is split up in seven clusters, particularly, in #1, #6, #9, #11, #14, #21 and #22. However, this sub-classification of the social sciences is less straightforward. Cluster #6 represents economics and business and political science, cluster #9 stands for psychology and linguistics, cluster #21 covers psychology and psychiatry, #11 comprises sociology and education, and cluster #1 is rather focussed on the humanities. Cluster #14 and #22 seem to have more heterogeneous profiles among these 'social and humanity clusters'. Although cluster #14 largely covers information and library science, the terms reflect a large overlap with other clusters. The same applies to cluster #22, which has obviously an even fuzzier structure. On the other hand, #9 and #21 are both covering psychology but focussing on different aspects, namely cognitive (#9) and medical (#21) issues.

Similarly to the structural analysis in Section 4.2, we use now a network with citation links among clusters to study the relationship of clusters based on hybrid cross-citation/textual information (see Figure 6). The observations made on the basis of most characteristic terms are confirmed by the link structure. The social-sciences and humanities clusters form two groups that are each strongly interlinked; one consists of clusters #1, #6, #14 and #22 with focus on humanities, economics, business, political and library science, the other one comprises #9, #11 and #21 with sociology, education and psychology. This is in line with the hierarchical structure shown in Figure 4. These two groups correspond to the two social-sciences clusters in the 7-cluster solution (cf. Section 4.4).

In the natural and applied sciences, we have found eight clusters, particularly, #2, #4, #5, #8, #15, #18, #19 and #20. On the basis of the most important TF-IDF terms (see Table 6) we can assign clusters #2, #15 and #19 to geosciences, environmental science, biology and agriculture, which, in turn, form a larger group corresponding to the first of the three "mega-clusters" in the 7-cluster solution. The graphic network presentation in Figure 6 confirms this interpretation. These three clusters form a group at the bottom. The other sciences clusters are more clearly recognisable, and distinctly separate fields. Thus clusters #4 represents chemistry, #20 physics, #5 engineering, #8 mathematics, and cluster #18 computer science. These science clusters form two groups, #4, #20 and #5 form one group of chemistry, physics and engineering, while #8 and #18 form the third group comprising mathematics and computer science. The network presentation and the hierarchical architecture in the dendrogram confirm the term characterisation.

The interpretation of the most characteristic terms of the life-science clusters is somewhat more complicated. Here we have a biomedical and a clinical group. These two groups are in line with the hierarchical structure of the dendrogram in Figure 4 but less clearly distinguished in the graphical network presentation (Figure 6). Nonetheless, the terms provide an excellent description for at least some of the medical clusters: cluster #7 stands for the neuro- and behavioural sciences, #3 for bioscience, #10 for the clinical and social medicine, #13 microbiology and veterinary science, #12 non-internal medicine, #16 hematology and oncology and #17 cardiovascular and respiratory medicine. According to the dendrogram clusters 3, 13, 16 and clusters 7, 10, 12, 17 form one larger cluster each. On the basis of the best terms, we can characterise these groups as the bioscience-biomedical and the clinical and neuroscience group, respectively.



Fig. 6. Network structure of hybrid clusters represented by the three most important TF-IDF terms

Cluster	50 best terms
1	polit; war; court; music; moral; essai; legal; philosophi; narr; text; literari; book; contemporari; french; religi; write; german; discours; ethic; civil; reform; christian; philosoph; justic; fiction; coloni; nineteenth; archaeolog; religion; aesthet; british; english; poetri; stori; feder; truth; church; russian; artist; liber; revolut; historian; gender; roman; america; militari; god; democraci; ideolog; china;
2	rock; basin; sediment; fault; sea; climat; soil; ic; miner; ocean; seismic; atmospher; river; wind; isotop; veloc; earth; star; tecton; earthquak; solar; precipit; ma; volcan; mantl; southern; lake; satellit; geolog; cloud; land; northern; groundwat; metamorph; flux; rainfal; shear; deform; forecast; weather; melt; crust; slope; faci; flood; sedimentari; clai; galaxi; season; magma;
3	receptor; rat; dna; genom; enzym; transcript; mutat; mice; metabol; peptid; diabet; cancer; insulin; chromosom; kinas; inhibitor; lipid; ca2; muscl; mrna; rna; neuron; molecul; vitro; apoptosi; mous; liver; tumor; glucos; assai; brain; hormon; mutant; amino; vivo; serum; mitochondri; embryo; fluoresc; diet; secret; phosphoryl; therapeut; bone; phenotyp; polymorph; prolifer; toxic; therapi; antibodi;
4	polym; ion; catalyst; crystal; bond; molecul; film; solvent; atom; hydrogen; ligand; nmr; polymer; poli; aqueou; adsorpt; thermal; methyl; spectroscopi; spectra; copolym; cation; fiber; cu; nm; bi; coat; mol; blend; nanoparticl; anion; chemistri; catalyt; ms; electrod; resin; chromatographi; ir; surfact; silica; copper; gel; amino; salt; column; uv; spectrometri; chiral; angstrom; fluoresc;

Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes

Data Mining and Knowledge Discovery in Real Life Applications

Cluster	50 best terms
5	crack; finit; concret; veloc; elast; turbul; vibrat; shear; thermal; nonlinear; beam; deform; fuel; motion; ltd; steel; cylind; combust; convect; flame; fatigu; compress; fractur; vehicl; jet; reynold; plane; wind; stiff; pipe; buckl; shell; friction; damp; vortex; cool; turbin; coal; fire; blade; bend; porou; lamin; axial; reservoir; rotor; specimen; cement; actuat; mesh; price; firm; trade; tax; economi; capit; incom; wage; invest; bank; financi; monetari; stock; welfar; labor; inflat; sector; privat; incent; household; earn; game; asset; employ; insur; forecast; reform; foreign; unemploy; volatil; profit;
	worker; polit; labour; fiscal; corpor; debt; monei; credit; investor; busi; export;
7	financ; shock; macroeconom; fund; currenc; equiti; school; agricultur; brain; neuron; rat; stroke; lesion; receptor; pain; cerebr; ct; mri; motor; cognit; injuri; spinal; nerv; mr; seizur; epilepsi; cortex; neurolog; tumor; arteri; dementia; sleep; cortic; syndrom; muscl; therapi; symptom; parkinson; neural; cord; alzheim; mice; synapt; chronic; axon; aneurysm; patholog; pet; eeg; nervou; surgeri; elderli; surgic; sclerosi; ms; deficit; rehabilit; sensori;
8	algebra; finit; theorem; manifold; infin; nonlinear; polynomi; let; graph; asymptot; singular; omega; invari; inequ; lambda; ellipt; convex; compact; conjectur; epsilon; hyperbol; infinit; proof; bar; symmetr; phi; sigma; topolog; banach; lie; eigenvalu; metric; matric; curvatur; perturb; integ; norm; parabol; cohomolog; hilbert; geometr; plane; lattic; semigroup; explicit; stochast; iter; dirichlet: exponenti; holomorph:
9	word; cognit; speech; semant; english; linguist; phonolog; stimulu; stimuli; lexic; cue; sentenc; speaker; verb; prime; perceptu; student; acoust; text; item; discours; verbal; auditori; emot; recal; syntact; brain; hear; skill; deficit; write; motor; judgment; listen; nois; letter; noun; mental; learner; psycholog; neuropsycholog; voic; instruct; vowel; motion; german; execut; ev; grammar; grammat;
10	nurs; hospit; physician; therapi; cancer; pain; mortal; pregnanc; infant; symptom; ethic; smoke; diabet; infect; birth; ci; men; interview; hiv; injuri; profession; chronic; syndrom; questionnair; worker; school; adolesc; student; surgeri; mental; child; neonat; sexual; breast; healthcar; caregiv; visit; matern; elderli; mother; satisfact; hypertens; vaccin; attitud; rural; cohort; doctor; cardiac; staff; gender;
11	student; school; teacher; teach; classroom; instruct; skill; academ; curriculum; literaci; disabl; learner; profession; colleg; cognit; peer; child; faculti; gender; reform; write; psycholog; pupil; graduat; attitud; undergradu; text; emot; interview; belief; discours; pedagog; think; gift; adolesc; preschool; english; inquiri; elementari; girl; boi; development; leadership; pedagogi; polit; web; tutor; team; item; intellig;
12	bone; ey; muscl; sport; athlet; pain; implant; surgeri; fractur; injuri; knee; dental; hip; surgic; nerv; anterior; retin; postop; oral; tendon; corneal; teeth; flap; periodont; graft; lesion; ocular; posterior; therapi; radiograph; ligament; neck; nasal; fixat; cartilag; dentin; glaucoma; laser; femor; shoulder; cari; player; ankl; cement; acuiti; tooth; syndrom; symptom; arthroplasti; rehabilit;

108

Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes

Cluste	er 50 best terms
13 14	infect; dog; hiv; vaccin; viru; hors; cow; milk; parasit; cat; pig; cattl; antibodi; diet; immun; pcr; calv; breed; viral; herd; sheep; pathogen; serum; antigen; therapi; farm; malaria; antibiot; veterinari; hospit; assai; egg; dairi; dna; bird; genotyp; chicken; pneumonia; hepat; fed; tuberculosi; bovin; goat; mortal; lesion; epidemiolog; outbreak; canin; virus; respiratori; firm; organiz; busi; librari; web; internet; custom; compani; employe; job; onlin; brand; team; strateg; journal; career; corpor; satisfact; student; price; trust; advertis; academ; profession; librarian; attitud; organis; leadership; cognit;
15	 enterpris; commerc; invest; sector; manageri; financi; psycholog; polit; retail; skill; commit; interview; emot; ventur; capit; purchas; intent; citat; book; retriev; text; soil; seed; crop; cultivar; leaf; fruit; bacteria; wheat; dry; rice; enzym; pathogen; shoot; nitrogen; microbi; ferment; bacteri; ha; pollut; sediment; milk; nutrient; dna; fertil; germin; biomass; seedl; season; agricultur; coli; sludg; irrig; veget;
16	infect; wast; grain; flower; co2; yeast; toxic; fungi; starch; genom; maiz; sp; grown; sugar; inocul; mutant; forest; cancer; tumor; therapi; il; carcinoma; transplant; breast; immun; infect; lung; liver; renal; antibodi; receptor; antigen; mice; malign; serum; chronic; prostat; lesion; surgeri; tumour; chemotherapi; mutat; inflammatori; dna; bone; recurr; surgic; cytokin; syndrom; hepat; apoptosi; biopsi; lymphoma; lymphocyt; neurostat; and provide antiparts and provide antiparts.
17	pancreat; gastric; resect; therapeut; rat; histolog; symptom; assai; prolifer; invas; inhibitor; asthma; median; coronari; arteri; cardiac; ventricular; myocardi; hypertens; cardiovascular; aortic; atrial; infarct; diabet; therapi; valv; vascular; stent; endotheli; surgeri; pulmonari; mortal; cholesterol; systol; bypass; syndrom; lv; graft; diastol; renal; vein; rat; echocardiographi; ischemia; hospit; chronic; dysfunct; receptor; angiotensin; anouryme: after fibril: athorosclaraci; mitral: yenou; perfus: ischem; corum;
18	reperfus; implant; inhibitor; ldl; vessel; fuzzi; schedul; robot; logic; machin; graph; nonlinear; traffic; web; asymptot; circuit; neural; nois; finit; stochast; fault; queri; custom; wireless; video; node; semant; heurist; antenna; motion; markov; polynomi; bayesian; iter; processor; sensor; covari; busi; execut; bandwidth; server; vehicl; internet; packet; wavelet; voltag; queue; alloc; algebra; intellig; bit; bardwar; theorem; ltd; parametr;
19	forest; habitat; fish; genu; egg; predat; season; sp; nest; sea; ecolog; larva; reproduct; lake; bird; prei; island; taxa; seed; river; veget; soil; breed; southern; ecosystem; nov; mate; genera; diet; biomass; insect; phylogenet; parasit; northern; marin; juvenil; forag; sediment; landscap; larval; winter; coastal; ocean; eastern: nutrient; summer; leaf; land; fisheri; assemblag;
20	film; alloi; laser; quantum; crystal; ion; steel; thermal; atom; beam; coat; glass; si; grain; microstructur; corros; silicon; dope; spin; ceram; powder; nm; scatter; fabric; neutron; diffract; dielectr; photon; cu; electrod; excit; ni; fe; emiss; sinter; deform; microscopi; fiber; voltag; hydrogen; sensor; anneal; spectra; lattic; spectroscopi; weld; semiconductor; nonlinear; machin; discharg;

109

Data Mining and Knowledge Discovery in Real Life Applications

Cluster	50 best terms
21	psycholog; adolesc; mental; emot; cognit; symptom; child; anxieti; psychiatr;
	abus; student; school; alcohol; schizophrenia; sexual; mother; gender; attitud;
	suicid; interview; cope; violenc; therapi; questionnair; youth; disabl; offend; men;
	belief; item; psychotherapi; aggress; mood; ptsd; client; satisfact; victim; peer;
	profession; distress; development; infant; interperson; crime; adhd; style;
	therapist; esteem; skill; childhood;
22	polit; urban; parti; gender; reform; economi; capit; democrat; democraci; employ;
	sector; war; land; sociolog; geographi; union; labour; rural; elect; welfar; ethnic;
	labor; discours; immigr; privat; actor; trade; civil; poverti; firm; citizen; busi;
	china; worker; incom; feminist; vote; liber; eu; household; elector; contemporari;
	agenc; job; inequ; domest; foreign; ideolog; agricultur; organiz;

Table 6. The 50 best TF-IDF terms describing the 22 hybrid citation-lexical clusters

In order to gain a better understanding of the cluster structure, we have ranked the journals of each of the 22 clusters according according to a modified version of Google's PageRank algorithm (Brin & Page, 1998) in which the number of citations is taken into account, normalised by the number of published papers. The following equation was used,

$$PR_{i} = \frac{(1-\alpha)}{n} + \alpha \sum_{j} PR_{j} \frac{a_{ji}/P_{i}}{\sum_{k} \frac{a_{jk}}{P_{k}}}$$
(1),

where PR_i is the PageRank of journal *i*, α is a scalar between 0 and 1 (α =0.9 in our implementation), *n* is the number of journals in the cluster, a_{ji} the number of citations from journal *i* to journal *i*, and P_i is the number of papers published by journal *i*, all in the period under study. Both sums iterate over the journals in the same cluster that contains journal *i*. Journal self-citations were removed prior to application of the algorithm. The five journals with highest PageRank are presented in Table 7. The PageRank of a journal can be understood here as the probability that a random reader will be reading that journal, when he randomly, continuously, and with equal probability looks up cited references to other journals (different from the current one), but once in a while randomly picks another journal from the library (cluster). Journals from arts & humanities (according to the ISI Subject Categories) were removed prior to application of the PageRank algorithm because of the low reliability of citation indicators in these disciplines. Zhang and Glänzel (2008) have shown that high entropy of journal cross-citations, relatively low impact and high share of journal self-citation makes it difficult to build reliable citation indicators for the humanities. This has to do with the subject-specific peculiarities in scholarly communication.

In general, the journals ranking highest represent their cluster in an adequate manner (cf. Table 7). Results of the PageRanking thus provide a realistic and representative picture of the hybrid clustering.

4.6 Comparison of subject and cluster structure

In this subsection we compare the structure resulting from the hybrid clustering with the ESI subject classification. This comparison is based on the *centroids* of the clusters and fields.

110

The centroid of a cluster or field is defined as the linear combination of all documents in it and is thus a vector in the same vector space. For each cluster and for each field, the centroid was calculated and the MDS of pairwise distances between all centroids is shown in Figure 7.

Cluster 1	Cluster 2	Cluster 3	Cluster 4
 YALE LAW J HARVARD LAW REV STANFORD LAW REV AM HIST REV COLUMBIA LAW REV 	1. ANNU REV ASTRON ASTR 2. ASTROPHYS J SUPPL S 3. EARTH-SCI REV 4. REV MINERAL GEOCHEM 5. ANNU REV EARTH PL SC	1. ANNU REV BIOCHEM 2. CELL 3. NAT REV MOL CELL BIO 4. ANNU REV CELL DEV BI 5. ANNU REV GENET	1. CHEM REV 2. PROG POLYM SCI 3. ACCOUNTS CHEM RES 4. ANNU REV PHYS CHEM 5. ADV DRUG DELIVER REV
Cluster 5	Cluster 6	Cluster 7	Cluster 8
1. PROG ENERG COMBUST 2. ANNU REV FLUID MECH 3. PROG AEROSP SCI 4. P COMBUST INST 5. COMBUST FLAME	1. Q J ECON 2. J FINANC 3. J ECON LIT 4. J POLIT ECON 5. J FINANC ECON	 ANNU REV NEUROSCI NAT REV NEUROSCI NEURON NAT NEUROSCI PROG NEUROBIOL 	 J AM MATH SOC ANN MATH ACTA MATH- DJURSHOLM INVENT MATH COMMUN PUR APPL MATH
Cluster 9	Cluster 10	Cluster 11	Cluster 12
1. PSYCHOL REV 2. BEHAV BRAIN SCI 3. COGNITIVE PSYCHOL 4. J EXP PSYCHOL GEN 5. COGNITION	1. MILBANK Q 2. ANNU REV PUBL HEALTH 3. JAMA-J AM MED ASSOC 4. HEALTH SERV RES 5. J HEALTH SOC BEHAV	1. REV EDUC RES 2. AM EDUC RES J 3. EDUC EVAL POLICY AN 4. EDUC PSYCHOL-US 5. J LEARN SCI	 CRIT REV ORAL BIOL M PROG RETIN EYE RES AM J SPORT MED PERIODONTOL 2000 SPORTS MED
Cluster 13	Cluster 14	Cluster 15	Cluster 16
Cluster 13 1. J ACQ IMMUN DEF SYND 2. AIDS 3. CLIN MICROBIOL REV 4. J INFECT DIS 5. CLIN DIAGN VIROL	Cluster 14 1. ADMIN SCI QUART 2. ACAD MANAGE J 3. ORGAN SCI 4. ACAD MANAGE REV 5. MIS QUART	Cluster 15 1. ANNU REV PLANT BIOL 2. PLANT CELL 3. CURR OPIN PLANT BIOL 4. ANNU REV PHYTOPATHOL 5. MICROBIOL MOL BIOL R	Cluster 16 1. ANNU REV IMMUNOL 2. NAT REV IMMUNOL 3. NAT IMMUNOL 4. CA-CANCER J CLIN 5. IMMUNITY
Cluster 13 1. J ACQ IMMUN DEF SYND 2. AIDS 3. CLIN MICROBIOL REV 4. J INFECT DIS 5. CLIN DIAGN VIROL Cluster 17	Cluster 14 1. ADMIN SCI QUART 2. ACAD MANAGE J 3. ORGAN SCI 4. ACAD MANAGE REV 5. MIS QUART Cluster 18	Cluster 15 1. ANNU REV PLANT BIOL 2. PLANT CELL 3. CURR OPIN PLANT BIOL 4. ANNU REV PHYTOPATHOL 5. MICROBIOL MOL BIOL R Cluster 19	Cluster 16 1. ANNU REV IMMUNOL 2. NAT REV IMMUNOL 3. NAT IMMUNOL 4. CA-CANCER J CLIN 5. IMMUNITY Cluster 20
Cluster 13 1. J ACQ IMMUN DEF SYND 2. AIDS 3. CLIN MICROBIOL REV 4. J INFECT DIS 5. CLIN DIAGN VIROL Cluster 17 1. CIRCULATION 2. CIRC RES 3. J AM COLL CARDIOL 4. ARTERIOSCL THROM VAS 5. CARDIOVASC RES	Cluster 14 1. ADMIN SCI QUART 2. ACAD MANAGE J 3. ORGAN SCI 4. ACAD MANAGE REV 5. MIS QUART Cluster 18 1. ACM COMPUT SURV 2. J ACM 3. J R STAT SOC B 4. VLDB J 5. IEEE T ROBOTIC AUTOM	Cluster 15 1. ANNU REV PLANT BIOL 2. PLANT CELL 3. CURR OPIN PLANT BIOL 4. ANNU REV PHYTOPATHOL 5. MICROBIOL MOL BIOL R Cluster 19 1. ANNU REV ECOL EVOL S 2. SYSTEMATIC BIOL 3. ANNU REV ENTOMOL 4. OCEANOGR MAR BIOL 5. TRENDS ECOL EVOL	Cluster 16 1. ANNU REV IMMUNOL 2. NAT REV IMMUNOL 3. NAT IMMUNOL 4. CA-CANCER J CLIN 5. IMMUNITY Cluster 20 1. REV MOD PHYS 2. MAT SCI ENG R 3 ANNU REV NUCL PART S 4. PHYS REP 5. PROG MATER SCI
Cluster 13 1. J ACQ IMMUN DEF SYND 2. AIDS 3. CLIN MICROBIOL REV 4. J INFECT DIS 5. CLIN DIAGN VIROL Cluster 17 1. CIRCULATION 2. CIRC RES 3. J AM COLL CARDIOL 4. ARTERIOSCL THROM VAS 5. CARDIOVASC RES Cluster 21	Cluster 14 1. ADMIN SCI QUART 2. ACAD MANAGE J 3. ORGAN SCI 4. ACAD MANAGE REV 5. MIS QUART Cluster 18 1. ACM COMPUT SURV 2. J ACM 3. J R STAT SOC B 4. VLDB J 5. IEEE T ROBOTIC AUTOM Cluster 22	Cluster 15 1. ANNU REV PLANT BIOL 2. PLANT CELL 3. CURR OPIN PLANT BIOL 4. ANNU REV PHYTOPATHOL 5. MICROBIOL MOL BIOL R Cluster 19 1. ANNU REV ECOL EVOL S 2. SYSTEMATIC BIOL 3. ANNU REV ENTOMOL 4. OCEANOGR MAR BIOL 5. TRENDS ECOL EVOL	Cluster 16 1. ANNU REV IMMUNOL 2. NAT REV IMMUNOL 3. NAT IMMUNOL 4. CA-CANCER J CLIN 5. IMMUNITY Cluster 20 1. REV MOD PHYS 2. MAT SCI ENG R 3 ANNU REV NUCL PART S 4. PHYS REP 5. PROG MATER SCI

Table 7. The five most important journals of each cluster according to a modified version of Google's PageRank algorithm (see Equation 1).



Fig. 7. Three-dimensional MDS map visualising distances between the centres (centroids) of the 22 ESI fields and the 22 clusters containing 8305 WoS journals.

In Figure 8, we use the Jaccard index to determine the concordance between our clustering solution and the ESI Scheme by comparing each cluster with every field, in order to detect the best matching fields for each cluster. The darker a cell in the matrix, the higher the Jaccard index, and hence the more pronounced the overlap between the corresponding cluster and ESI field. For example, cluster #4 (Chemistry) definitely corresponds to ESI field



Fig. 8. Concordance between our clustering solution and the ESI Scheme visualised by coloured cells representing the Jaccard index for each cluster and field pair. The darkest cells represent the best matching pairs of fields and clusters. In the upper figure, the Jaccard index is computed from the number of journals a cluster and a field have in common, while the lower figure takes the size of each journal into account by counting the numbers of overlapping papers.

#3 (Chemistry). The same applies to field and cluster #6 (Economics and business). Clearly, ESI field #21 has the least concordance as this field is spread over seven clusters. It is defined as one single field in social sciences. It is not a surprise that the strongest match is found with our somewhat 'fuzzy' multidisciplinary social cluster. On the other hand, clusters #13 and #14 are quite similarly spread over four ESI fields each.

4.7 Migration of journals among subject fields and clusters

If clustering algorithms are adjusted or changed, one can observe the following phenomenon. Some units of analysis are leaving clusters they formerly belonged to and end up in different clusters. This phenomenon is called 'migration'. We can distinguish between 'good migration' and 'bad migration'. 'Good migration' is observed if the goodness of the unit's classification improves, otherwise we speak about 'bad migration'. We can also apply this notion of migration to the comparison of clustering results with any reference classification. In the following we will use the ESI scheme as reference classification.

In the previous section we visualised the concordance between the clustering and the ESI classification. To determine for each ESI field the cluster that best matches the field, we used the Jaccard index on basis of the number of overlapping journals (cf. upper part of Figure 8). Out of 8305 journals under study, there were more than one third, namely, 3204 journals that were not assigned to the cluster which best matches their ESI field. As already mentioned above, we call these journals 'migrated journals'. The largest 'exodus' comprising 226 migrating journals occurred from the ESI "Engineering" field to cluster #18 (Computer science), whereas the best matching cluster for the Engineering field is actually Cluster #5 (Engineering). The top 10 strongest patterns of migration are listed in Table 8, which indicate possible improvements of journal assignments.

Migration pattern	Number of migrated journals
From ESI field 7 (Engineering) to Cluster 18	226
From ESI field 14 (Molecular Biology & Genetics) to Cluster 3	159
From ESI field 21 (Social Sciences, general) to Cluster 10	145
From ESI field 11 (Materials Science) to Cluster 20	139
From ESI field 4 (Clinical Medicine) to Cluster 7	132
From ESI field 19 (Plant & Animal Science) to Cluster 15	108
From ESI field 21 (Social Sciences, general) to Cluster 21	98
From ESI field 7 (Engineering) to Cluster 20	95
From ESI field 4 (Clinical Medicine) to Cluster 3	86
From ESI field 8 (Environment/Ecology) to Cluster 15	86

Table 8. Top 10 strongest migration patterns

To measure the quality of migrations, we calculated the differences in Silhouette values before and after migration (based on textual and citation distances), for each migrated journal. Most migrated journals improved their Silhouette values. In the following, we will give some examples of good migrations and bad migrations.

'Good migrations' are observed if journals improved their Silhouette values after migration. Based on their titles and scopes (not shown), apparently they should indeed be assigned to the cluster to which they have moved. We observed numerous good migrations and the following cases will serve just as examples.

The Journal of Analytical Chemistry and Chemia Analityczna migrated from ESI field #7 (Engineering) to cluster 4 (chemistry). The best matching ESI cluster were field #3 (Chemistry) in this case (cf. Figure 8). Similarly, Land Economics, Developing Economies and Economic Development and Cultural Change migrated from field #21 (Social Sciences, general) to the more specific cluster 6 (economics and business). Here, the corresponding ESI field were #6 (Economics & business). In the life sciences, we found the following good migration. The journals Neuropathology, Revista de Neurologia, Current Opinion in Neurology, Revue Neurologique, Lancet Neurology, European Journal of Neurology, Neurologist, Nervenheilkunde, Visual Neuroscience, Seminars in Neurology, Epilepsy & Behavior and Journal of Neuroimaging migrated from field #4 (Clinical Medicine) to cluster #7 (neuroscience and behaviour) which rather corresponds to ESI field #16 (Neuroscience and behavior). Finally, we mention a migration between engineering and mathematics. The journals Quarterly of Applied Mathematics, Bit Numerical Mathematics, Siam Journal on Discrete Mathematics and Discrete Applied Mathematics, which were assigned to the ESI field Engineering (field #7), were found in our 'Mathematics' cluster (#8) which in turn corresponds to WSI field #12 (Mathematics). In the case of bad migration, the Silhouette values decreased after migration, that is, their Silhouette values in the ESI scheme were better than in the hybrid clustering. The reasons for this phenomenon are not always clear. According to their titles and scopes this migration is not always convincing. For instance, Journal of Astrophysics and Astronomy, New Astronomy, Astrophysical Journal and Astronomy & Astrophysics migrated from the ESI field 22 (Space Science) to Cluster 2 (geosciences) corresponding to ESI field #9, where we have to admit that journals in astronomy and astrophysics are in general spread over the geosciences and physics clusters. Viral Immunology migrated from field #10 (Immunology) to cluster #13 (microbiology and veterinary science) and Canadian Journal of Microbiology migrated from field #13 (Microbiology) to cluster #15 (agricultural and environmental sciences). Both clusters are rather spread over several ESI fields each (see Figure 8).

The distinction between good and bad makes a target-oriented adjustment of the existing classification scheme possible. Good migration can be used to reassign journals within the old scheme on the basis of the concordance with the results of clustering.

5. Conclusions

The hybrid clustering using textual information and cross-citations provided good results and proved superior to its two components when applied separately. The goodness of the resulting classification was even better than that of the "intellectual" reference scheme, the ESI subject scheme. Both classification systems form partitions of the Web of Science so that the direct comparison of clusters and fields was possible. According to our expectations, not all clusters have a unique counterpart in the ESI scheme and *vice versa* although the number of clusters coincided with the number of ESI fields. Although the Silhouette and modularity values substantiate a more coherent structure of the hybrid clustering as compared with the

ESI subject scheme, not all clusters are of high quality. Problems have been found, for instance, in clusters #1 and #12 where interdisciplinarity and strong links with other clusters distort the intra-cluster coherence. However, intellectual classification schemes usually do have a category "multidisciplinary sciences" as well. Although the result of a hard clustering algorithm often does contain a cluster with objects (journals) not strongly related to any other cluster, forming a "multidisciplinary sciences" cluster is not an inherent goal of the algorithm, and actually is not really meaningful either in the light of our outset goal to improve the classification of the sciences. Consequently, real multidisciplinary journals are scattered around different clusters.

Based on the external validation of clustering results by expert knowledge present in ISI subject categories, seven clusters seem to yield best results. Although there is no adequate subject classification scheme with 7 categories to be used as reference system, a more detailed analysis of this solution will be part of future research. Additional ideas for future research are a further improvement of the hybrid clustering algorithm by iterative cleaning of clusters as a post-processing step; allowing multiple assignments by fuzzy clustering; evaluating other algorithms like spectral clustering; and, finally, dynamic analysis by dynamic hybrid clustering.

The continuous rise of computing power might one day allow a large-scale mapping of the scientific universe explorable at various levels of detail. What's more, application of advanced natural language processing and machine summarization at the scale of large bibliographic corpora might offer some insight into semantics beyond mere statistical processing.

6. References

- Bader, B.W. & Kolda, T.G. (2006). Algorithm 862: MATLAB tensor classes for fast algorithm prototyping. ACM Transactions on Mathematical Software, 32, 4, 635-653, ISSN: 0098-3500
- Baeza-Yates, R. & Ribeiro-Neto, B. (1999). *Modern Information Retrieval*, Addison-Wesley, ISBN-10: 020139829X, Cambridge
- Batagelj, V. & Mrvar, A. (2002). Pajek analysis and visualization of large networks. *Graph Drawing*, 2265, 477–478, ISSN: 0302-9743
- Berry, M.; Dumais, S.T. & O'brien, G.W. (1995). Using linear algebra for intelligent information retrieval. *SIAM Review*, 37, 4, 573–595, ISSN: 0036-1445
- Brin, S. & Page, L. (1998). The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30, 1-7, 107–117, ISSN: 0169-7552
- Boyack, K.W.; Börner, K. & Klavans, R. (2008). Mapping the structure and evolution of chemistry research. *Scientometrics*, forthcoming
- Braam, R.R.; Moed, H.F. & Van Raan, A.F.J. (1991a). Mapping of science by combined cocitation and word analysis, Part 1: Structural aspects. JASIS, 42, 4, 233-251, ISSN: 0002-8231
- Braam, R.R.; Moed, H.F. & Van Raan, A.F.J. (1991b). Mapping of science by combined co-citation and word analysis, Part II: Dynamical aspects. JASIS, 42, 4, 252-266, ISSN: 0002-8231

- Callon, M.; Courtial, J.P.; Turner, W.A. & Bauin, S. (1983). From translations to problematic networks – An introduction to co-word analysis. *Social Science Information*, 22, 2, 191–235, ISSN: 0539-0184
- Callon, M.; Law, J. & Rip, A. (Eds.). (1986). *Mapping the Dynamics of Science and Technology:* Sociology of Science in the Real World, Macmillan Press, ISBN-10: 0333372239, London
- Deerwester, S.; Dumais, S.T.; Furnas, G.W.; Landauer, T.K. & Harshman, R. (1990). Indexing by latent semantic analysis. *JASIS*, 41, 6, 391–407, ISSN: 0002-8231
- ESI, *Essential Science Indicators* (accessible via: http://www.esi-topics.com/fields/ index.html)
- Fano, R.M. (1956). Information Theory and the Retrieval of Recorded Information, In: Documentation in Action, J. H. Shera, A. Kent, J.W. Perry (Ed.), 238–244, Reinhold Publ. Co., New York
- Garfield, E. (1975). ISIS Atlas of Science may help students in choice of career in science. *Current Contents*, 29, 5–8
- Garfield, E. (1988). The encyclopedic ISI Atlas of Science launches three new sections biochemistry, immunology, and animal & plant sciences. *Current Contents*, 7, 3–8
- Glänzel, W. & Czerwon, H.J. (1996). A new methodological approach to bibliographic coupling and its application to the national, Regional and institutional level. *Scientometrics*, 37, 2, 195–221, ISSN: 0138-9130
- Glenisson, P.; Glänzel, W. & Persson, O. (2005a). Combining full-text analysis and bibliometric indicators. A pilot study. *Scientometrics*, 63, 1, 163–180, ISSN: 0138-9130
- Glenisson, P.; Glänzel, W; Janssens, F & De Moor B. (2005b). Combining full text and bibliometric information in mapping scientific disciplines. *Information Processing & Management*, 41, 6, 1548–1572, ISSN: 0306-4573
- Hatcher, E. & Gospodnetic, O. (2004). *Lucene in Action*, Manning Publications Co, ISBN-10: 1932394281, New York
- Hicks, D. (1987). Limitations of co-citation analysis as a tool for science policy. *Social Studies* of Science, 17, 2, 295–316, ISSN: 0306-3127
- Jain, A. & Dubes, R. (1988). *Algorithms for Clustering Data*. Prentice Hall, ISBN-10: 013022278X, New Jersey
- Janssens, F. (2007a). Clustering of Scientific Fields by Integrating Text Mining and Bibliometrics, Ph.D. Thesis, Faculty of Engineering, Katholieke Universiteit Leuven, Belgium, http://hdl.handle.net/1979/847
- Janssens, F.; Glänzel, W. & De Moor, B. (2007b). Dynamic hybrid clustering of bioinformatics by incorporating text mining and citation analysis. Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 360–369, ISBN 978-1-59593-609-7, San Jose, CA, USA, August 2007, ACM, New York
- Janssens, F.; Glanzel, W. & De Moor B. (2008). A hybrid mapping of information science. *Scientometrics*, 75, 3, 607–631
- Jarneving, B. (2005). The Combined Application of Bibliographic Coupling and the Complete Link Cluster Method in Bibliometric Science Mapping. PhD Thesis, University College of Borås/Göteborg University, Sweden

Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes

- Jarneving, B. (2007). Bibliographic coupling and its application to research-front and other core documents. *Journal of Informetrics*, 1, 4, 287–307, 1751-1577, ISSN: 1751-1577
- Kessler, M.M. (1963). Bibliographic coupling between scientific papers. *American Documentation*, 14, 1, 10–25, ISSN: 0096-946X
- Kostoff, R.N.; Toothman, D.R.; Eberhart, H.J. & Humenik, J.A. (2001). Text mining using database tomography and bibliometrics: A review. *Technological Forecasting and Social Change*, 68, 3, 223–253, ISSN: 0040-1625
- Kostoff, R.N.; Buchtel, H.A.; Andrews, J. & Pfeil, K.M. (2005). The hidden structure of neuropsychology: Text mining of the journal Cortex, 1991–2001. *Cortex*, 41, 2, 103–115, ISSN: 0010-9452
- Leydesdorff, L. (2006). Can scientific journals be classified in terms of aggregated journaljournal citation relations using the Journal Citation Reports? *JASIST*, 57, 5, 601–613, ISSN: 1532-2882
- Leydesdorff, L. & Rafols, I. (2008). A Global Map of Science Based on the ISI Subject Categories. *JASIST*, to be published, ISSN: 1532-2882
- Mardia, K.V.; Kent, J.T. & Bibby, J.M. (1979). *Multivariate Analysis*, Harcourt Brace & Co, Academic Press, ISBN-10: 0124712525, London, UK.
- Marshakova, I.V. (1973). System of connections between documents based on references (as the Science Citation Index), *Nauchno-Tekhnicheskaya Informatsiya, Seriya* 2, 6, 3–8
- Moya-Anegon, F. de; Vargas-Quesada, B.; Chinchilla Rodriguez, Z. ; Corera-Alvarez, E.; Munoz Fernandez, F.J. & Herrero-Solana, V. (2007). Visualizing the Marrow Science. JASIST, 58, 14, 2167-2179.
- Narin, F. (1976). Evaluative Bibliometrics: The Use of Publication and Citation Analysis in the Evaluation of Scientific Activity. Computer Horizons, Inc., Washington, D.C
- Newman, M.E.J. & GIRVAN, M. (2004). Finding and evaluating community structure in networks. *Physical Review E*, 69, 2, ISSN: 1063-651X
- Newman, M.E.J. (2006). Modularity and community structure in networks. *PNAS US*, 103, 23, ISSN: 0027-8424
- Noyons, E.C.M. (1999). *Bibliometric Mapping as a Science and Research Management Tool,* DSWO Press, Leiden University, ISBN 9090132503, Leiden, The Netherlands
- Pinski, G. & Narin, F. (1976). Citation influence for journal aggregates of scientific publications. *Information Processing and Management*, 12, 5, 297–312, ISSN: 0306-4573
- Rousseeuw, P. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics, 20,1,53–65.
- Salton, G. & Mcgill, M.J. (1986). Introduction to Modern Information Retrieval, McGraw-Hill, Inc, ISBN: 0070544840, New York.
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *JASIS*, 24, 4, 265–269, ISSN: 0002-8231
- Small, H. (1978). Cited documents as concept symbols. *Social Studies of Science*, 8, 327–240, ISSN: 0306-3127
- Small, H. (1998). A general framework for general large-scale maps of science in two or three dimensions: The SciViz system. *Scientometrics*, 41, 1–2, 125–133, ISSN: 0138-9130

Zhang, L. & Glänzel, W. (2008). Journal cross-citation matrices reconsidered. Tracing the role of individual journals in the communication network. In: *Proceedings of WIS 2008*, H. Kretschmer and F. Havemann (Ed.), Berlin. (accessible via: http://www.collnet.de/Berlin-2008/ZhangLinWIS2008jcm.pdf)







Data Mining and Knowledge Discovery in Real Life Applications Edited by Julio Ponce and Adem Karahoca

ISBN 978-3-902613-53-0 Hard cover, 436 pages **Publisher** I-Tech Education and Publishing **Published online** 01, January, 2009 **Published in print edition** January, 2009

This book presents four different ways of theoretical and practical advances and applications of data mining in different promising areas like Industrialist, Biological, and Social. Twenty six chapters cover different special topics with proposed novel ideas. Each chapter gives an overview of the subjects and some of the chapters have cases with offered data mining solutions. We hope that this book will be a useful aid in showing a right way for the students, researchers and practitioners in their studies.

How to reference

In order to correctly reference this scholarly work, feel free to copy and paste the following:

Frizo Janssens, Lin Zhang and Wolfgang Glänzel (2009). Hybrid Clustering for Validation and Improvement of Subject-Classification Schemes, Data Mining and Knowledge Discovery in Real Life Applications, Julio Ponce and Adem Karahoca (Ed.), ISBN: 978-3-902613-53-0, InTech, Available from:

http://www.intechopen.com/books/data_mining_and_knowledge_discovery_in_real_life_applications/hybrid_clu stering_for_validation_and_improvement_of_subject-classification_schemes



InTech Europe

University Campus STeP Ri Slavka Krautzeka 83/A 51000 Rijeka, Croatia Phone: +385 (51) 770 447 Fax: +385 (51) 686 166 www.intechopen.com

InTech China

Unit 405, Office Block, Hotel Equatorial Shanghai No.65, Yan An Road (West), Shanghai, 200040, China 中国上海市延安西路65号上海国际贵都大饭店办公楼405单元 Phone: +86-21-62489820 Fax: +86-21-62489821 © 2009 The Author(s). Licensee IntechOpen. This chapter is distributed under the terms of the <u>Creative Commons Attribution-NonCommercial-ShareAlike-3.0 License</u>, which permits use, distribution and reproduction for non-commercial purposes, provided the original is properly cited and derivative works building on this content are distributed under the same license.



