

DEVELOPMENT OF A METHOD FOR DETECTION AND TREND ANALYSIS OF RESEARCH FRONTS BUILT BY LEXICAL OR COCITATION ANALYSIS⁺

M. ZITT*, ** E. BASSECOULARD*

* LERECO, INRA, BP 527 - F-44026 Nantes Cedex 03 (France)

** Observatoire des Sciences et des Techniques (OST), 93 rue de Vaugirard, F-75006 Paris (France)

(Received January 18, 1994)

Detecting homogeneous areas in research networks is a very common feature of bibliometric analysis, either for academic or policy purposes. The method presented here combines structural analysis and trend detection, by operating on a "thick-slice" of time, starting from co-citation or co-word analysis (applications of either type have already been carried on). Significance of "trend" of clusters is partially addressed, through an analysis of publication delays. Examples are given on a co-citation analysis in the field of astrophysics (1986-1989).

Introduction

Detecting homogeneous areas in research networks is a very common feature of bibliometric analysis, either for academic or policy purposes. Various multidimensional methods, either continuous or discontinuous, have been developed to represent scientific universes. Discontinuous methods such as classifications, are widely used to scale down the level of observation throughout the fractal structure of science or technology networks. The most common structuring relationships they rely upon, are cocitations (I.S.I. research fronts), word-proximity (e.g. C.S.I., LeximappeTM) or co-classification. Cocitation and co-word methods have often been compared and/or combined (*Braam et al.*¹).

Combining structural and temporal aspects is also an appealing domain for these methods (*Braam et al.*²). For this purpose we developed in our laboratory a sequence of programmes [1] applicable to either structuring relation with appropriate setting, and embodying direct "dynamic" approach in the very weak sense of trend characterization. The principle is very simple: clustering is carried on a "thick slice" of

⁺Paper presented at the Fourth International Conference on Bibliometrics, Informetrics and Scientometrics in Berlin (Germany), September 11-15, 1993.

time, and allows temporal coding of clusters, which gives a prospective view, again in a weak meaning of extrapolation. This information on trend is especially valuable when combined to other features such as cluster structure or immediacy.

Applications of both types (cocitation, cword) were conducted. The first one was a lexical analysis (*Bassecouard, Zitt*³) of a single journal on Food Science, "Science des Aliments" which is not an I.S.I. source journal, on a "long period" basis (11 years). Examples will be given here on the second one, which is a cocitation analysis on two periods (I-1986 through 1989 and II-1989 through 1992) of a set of two prominent international journals in the field of astrophysics, the *Astrophysical Journal* and *Astronomy and Astrophysics*. In the present account, we only show provisional results for period I. Validation by expert advice of findings in the field of Astrophysics is in process, especially "predictive" aspects through comparison of periods I and II.

Methods

Trend representation

Without speaking about more ambitious causal or formal dynamic models of science, several types of trend assessments, very limited in scope, are helpful to describe changes in science networks :

- structural analysis copes with the problem in an indirect way, assuming that structural and dynamic properties of a given system are linked. For instance, in a given cognitive network, dense areas are supposed to be more resistant to change than others. Leximappe™ strategic diagrams is a well known tool for such an approach. The position of an element in these synchronic diagrams may be interpreted as reflecting its position in a life-cycle (*Courtial*⁴).
- explicitly diachronic studies often use "comparative statics" and juxtapose successive situations, for instance successive pictures of synchronic relations (cocitation, co-word...) after factorial, MDS or clustering methods. For classification, a good example is I.S.I. chronological sequences of clusters.

Our approach is somewhat different. We consider a thick slice of time (for instance 4 years) to identify "heavy phenomena", namely clusters structure, and we evaluate trends inside this structure using a diachronic information, namely the dates of documents assigned to clusters. The heavy structure may be more reliable through averaging; in counterpart the slice must not be too thick, since the average structure

may obscure network changes; so various horizons may be adopted, depending on the immediacy characteristics of the discipline. Thick slice is more demanding in computer resources, which can be a serious drawback.

"Thick slice" analysis

The tool is typically aimed at "thick slice" study. Here, each cluster (theme) is given the weighted average date of the assigned documents, which depends on the variation of the amount of assigned publications within the period.

This temporal coding is maintained throughout the analysis and in most graphic outputs (classification trees, maps). Moreover, inside a several-year file, as documents are dated, virtually all derived information elements may be assigned a date: dating of elements is straightforward (average date of occurrence of words, authors, institutions, etc.). In graphic outputs presented below, temporal coding is expressed from light grey (beginning of period) to dark grey (end of period). Dating of links between elements may also be useful. For long periods Zitt⁵ proposed to characterize word co-occurrence links after their chronological profile. For a few-years period, where links are measured after co-occurrence proximities, a simpler way is to assign to the link the average date of the co-occurrence. Coauthorship, for instance, can be treated this way (link dating is not provided in the following examples).

The temporal coding of clusters calls for several caveats. For example if a thick slice view gives more reliability to cluster identification, the advantage does not hold for date determination.

Cluster average date significance

The average date of a cluster is computed after the publication date of documents, which is obviously an important landmark, but not the only one. A scientific research undergoes long processes which have been analysed from many points of view: transformation of scientific statements (sociology of "translation"), efficiency of scientific communication, etc... Ziman⁶, for instance, addressed the question of how scientific information becomes public knowledge. He confirmed the observations of Garvey and Griffith⁷ in their study of scientific communication in psychology, who stated that the various stages from hypothesis to publication take months or years to complete. Roland and Kirkpatrick⁸ scrutinized the various components of the "idea-to-research-to-paper-to-publication process" and evaluated

a total time of at least four years for medicine. Landmarks of research processes are only available through such special inquiries. This is also true for the "real date" of a work achievement (say T_0), obviously not a clearcut notion. More precise milestones are met in the publication process itself: T_1 , first date of submission to the journal, and sometimes T_2 , date of final acceptance, are documented in scientific journals, and T_3 (publication date) is known. Studies of publication time in various fields have been published, for instance in medical biology (*Carson and Wyatt*⁹), agriculture and allied sciences (*Jain and Goyal*¹⁰) or chemical engineering (*Pings*¹¹).

For our purpose, we consider that the submission date could be an indicator of front dynamism, at least as good as the publication date itself. We have then to accept that the measure of "date of papers" based on available publication dates embodies the random nature of the delay T_1 - T_3 . Though a long span T_1 - T_2 could mean that the work needed important corrections and completion by authors, through referees' advice (*Lay*¹²), this first delay also accounts for various phenomena such as selection of referees, referees' available time, mailing, etc... The time from acceptance to publication T_2 - T_3 combines purely technical aspects and journal periodicity. Empirically measured – examples are discussed below – the total delay T_1 - T_3 appears as a random phenomenon.

What we are primarily interested in, is the dispersion of the publication lag rather than its average value. A strong dispersion of this variable will clearly jeopardize the interpretation of clusters dynamism, evaluated after a weighted average of publication dates of papers: using submission instead of publication date would imply to go back for each study to individual papers, since this piece of information is not available in usual databases. So we will rely on estimates. n being the number of papers dated d_i assigned to the cluster with a weight p_i , the variance of the average date D under independence condition is $\text{var}(D) = \sum_{i(1,n)} [p_i^2 \text{var}(d_i)]$. n is known for each cluster; considering that $\text{var}(d_i)$ is only dependent on the journal, it can be obtained by sampling. Assuming that the distribution of the average delay is normal, the degree of significancy of the average date D for a cluster, with a zero hypothesis "no dynamism" (average date = central date of the period), can be stated from there with corrections for the class effect (e.g., when only the *year* of publication is documented in the database). Since the studies described here were both conducted at the journal level, the delay T_1 - T_3 was easy to document.

In the case of astrophysics, the dispersion of publication lags has been estimated on a sample of N papers on both journals, on the period of cocitation analysis: 1986-1992. Average lag is 288 ± 111 for *Astrophysical Journal*, 321 ± 113 for *Astronomy and*

Astrophysics ($N=56$ in both cases). Distributions of lag are almost normal and largely similar, so that global distribution is almost normal [2]. This, along with the sample size, allows to estimate the population standard deviation after the sample's, namely 113 days. Somewhat arbitrarily, we will use a more prudent value of 6 months for significance tests; higher deviations could be advocated invoking T_0 or of course original ideas dates. Fig. 1 displays the respective part of T_1-T_2 and T_2-T_3 , by classes of total delay for each journal: we observe that as could be perhaps expected, the T_2-T_3 is fluctuating around its mean, and that the instruction delay is steadily growing for *Astronomy and Astrophysics*. For *Astrophysical Journal*, T_2-T_3 also accounts for lengthening of total delay. *Abt*¹³ used *Astrophysical Journal* as representing Astrophysics in his study of publication practices in various disciplines: on a much bigger sample for the year 1990, he found an average of 314 ± 166 , with a coefficient of variation ($\approx 1/2$) higher than in other disciplines ($\approx 1/3$). On a longer time basis (our sample), the journal seems to approach the last value.

Cocitation study options

Cocitation analysis proposed in 1973 by *Small*¹⁴ and *Marshakova*¹⁵ is widely used as a powerful tool to represent science structure, with a few variants (*Small, Sweeney*¹⁶). Caveats and limits have often been discussed (see for example on a case study *Hicks*¹⁷). A clear drawback is the fact that cocitation clustering can only classify a fraction of the citing literature. Various solutions have been proposed to overcome this difficulty, e.g. complementary assigning using words (*Braam et al.*¹). In this respect author cocitation (*McCain*¹⁸, *Penan*¹⁹) may also improve recall rate. At this stage, this problem has not been addressed in the following experiments, that used a fairly classical document cocitation scheme, with a cosine form index (weighted Ochiai index [3]) and citing documents assignment. A drastic selection of structuring papers has been conducted [4].

Classification options

Hierarchical clustering, closely connected to the embedded or "fractal" structure of cognitive networks, often gives more legible results on large files than continuous techniques that can be advocated for statistical reasons. A very large number of clustering methods, and options inside a given method, are available from statistical packages; from the numerous comparisons reported, either on theoretical or

empirical grounds, one can hardly conclude to a single "best way". In the following examples average link was chosen; though not bias-free, this method is considered as limiting ultrametric distortion.

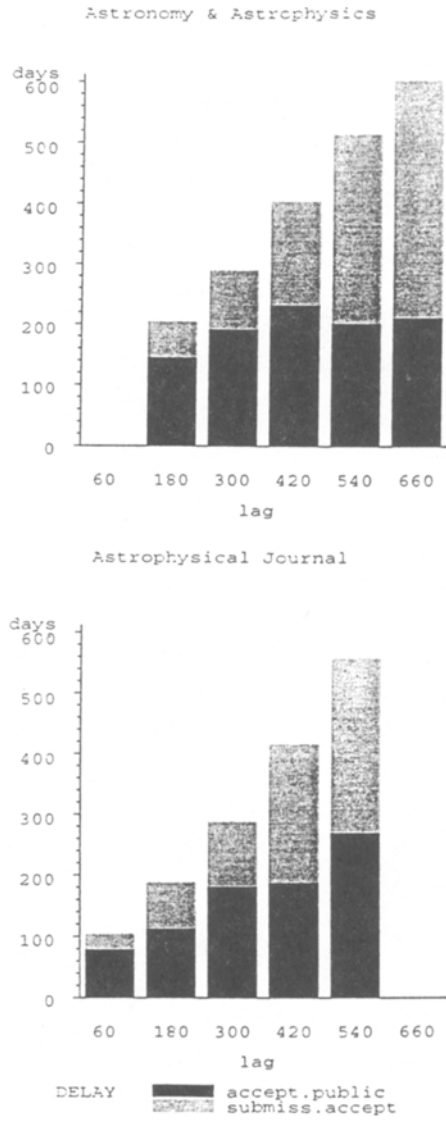


Fig. 1. Lag components

In hierarchical clustering, the search for natural frontiers between clusters is sometimes conducted with "modal analysis". The principle is to stop the merging of clusters if they are both above a chosen critical size. Though orthodox use of modal analysis is "density" classification (eg nth neighbour), following examples partially rely on modal analysis: instead of a unique cutting level, a range is determined for the clustering index, and modal cutting is activated inside the range [5]. In the same view, we introduce an empirical coefficient of "internality", defined after the path of an element in the ultrametric tree, reported in a normalized square diagram: cluster population (ordinates) vs clustering index or rank function. The area beneath the path, normalized by the embedding rectangle area, measures the "internality" of an element. For a given axes scale, the internality rank is preserved in the aggregation process: if the element *i* is less internal than *j* in the cluster *K*, it will be less internal than *j* in any cluster *L* containing *K*. This allows for instance to qualify the respective position of "core" documents inside a cluster and related superclusters, position which may be accounted for when assigning citing documents. Internality also suggests a way to measure the stability of the clusters [6].

Example: the astrophysics field

The maps displayed come from a cocitation analysis in the field of astrophysics. Among prestigious journals contributing to the subject, either specialized or multidisciplinary, the couple *Astrophysical Journal/Astronomy and Astrophysics* (the first one is published in the U.S.A., the second one in Europe) appears as a good subset. The "research fronts" built from this subset of citing documents do not only inform on the evolution of the contents of both journals but also, to a certain extent, on the whole field (on the position of the *Astrophysical Journal*, see *Davoust et al.*²⁰; recent citation analyses of astronomical research have been completed by *Trimble*²¹ for American astronomers and *Jaschek*²² for European astronomers).

Typical output

The sequence of programmes follows usual stages of co-word or co-citation analyses: reformatting, preparation, similarity matrices, classification, cartography for the whole subject and for individual clusters. A reformatting module was developed for I.S.I. sources exploitation. A critical task, as usual, is dictionaries normalizing.

We will not detail the description of clusters, which is classical: fronts are documented on the one hand by their core (ranked cocited documents), on the other hand by characteristics of assigned documents building the "front" or "theme" (citing documents). A given document may be partially assigned to several fronts; low fractions (typically <20%) are discarded. Usual statistics on actors, with probability index (cluster/whole subject), are available. Pseudo-titles are given to clusters using a sequence of top specific words from assigned documents titles [7].

The "average date" of the front (from documents dates, weighted by assignment fractions) is added to theses usual features; the statistical significance of this date, and consequently a possible non flat trend, is addressed after the hypotheses mentioned above. In the example of the cluster "planetary nebulae", the trend is significantly upward (threshold 0.025). Tables allowing to compare trends for the very center and the periphery of the cluster can be used to describe more deeply the "growth regime" of the clusters.

Mapping a cluster: examples on cluster 329 "planetary nebulae".

The Fig. 2 exhibits the core documents (resp. words for a co-word study) in a given cluster (inner circle on the map) and its neighbourhood (outer circle) and their cocitation links. A specificity of the representation is the "trend" characterization of the objects: the color or pattern of each element expresses its average citing date. In addition, for co-citation analysis, an auxiliary word similarity map (title words) can be drawn (Fig. 3.), for illustration sake. The proximity to the centre depends on the specificity of the word to the cluster. In the example shown, similarity is local (Ochiai index). If computer constraints allow it, profile (contextual) similarity could be a better approach for title words. As for the whole field (see below) coauthorship network inside the cluster could be drawn (not displayed).

Mapping the whole subject

The first characterisation is the general collaboration network. This map is of classical type. Co-authorship maps can be obtained for various levels of actors (typically country, institution, authors), or for several types of co-authorship index (all, fractional, semi-fractional count; gross or normalized). Normalized index maps show mutual affinities, gross figures reflect rather the dominance or centrality of actors. For Astrophysics for instance, the general map (not included) shows a "club"

of dominant interconnected countries (USA, Europe), with a few bilateral relations with other countries. Separate maps for each journal (Fig. 4a, 4b), with direct gross counting (thresholds are proportional to each journal number of source documents), show that the USA are completely central in the Astrophysical Journal [8], while for Astronomy and Astrophysics, published in Europe, a "club" effect (multilateral coauthorships) for a group of active European countries (+ USA, and Chile) is at work.

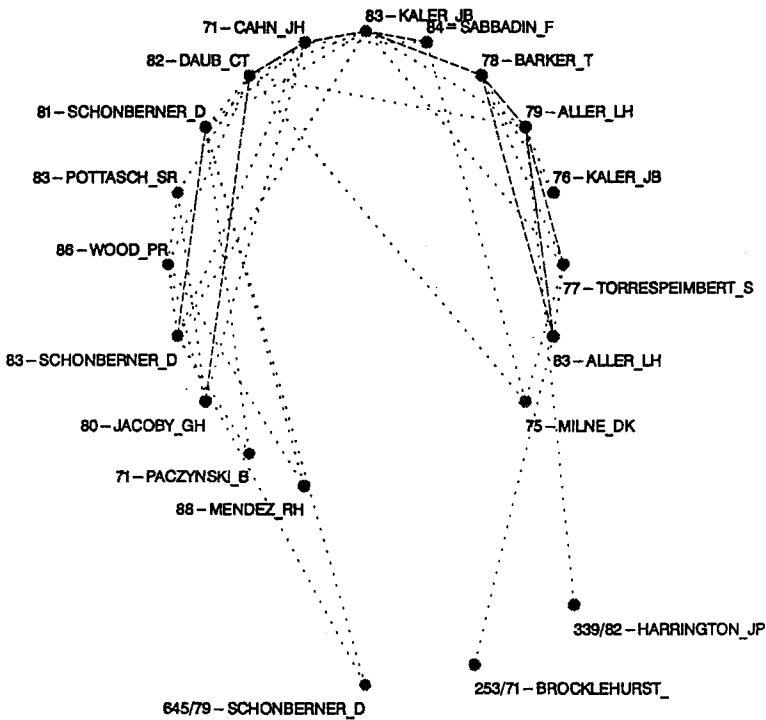


Fig. 2. Cocitation links on theme (theme planet - nebula (329))

All maps of the field only include biggest clusters. The next diagram (Fig. 5) is a summarized (25 levels) classification tree, limited here to 76 big clusters, that gradually combine into "superclusters". In order to detect trends for superclusters, trend coding (colors) is displayed for all levels. This may be useful to distinguish local from more extended trends.

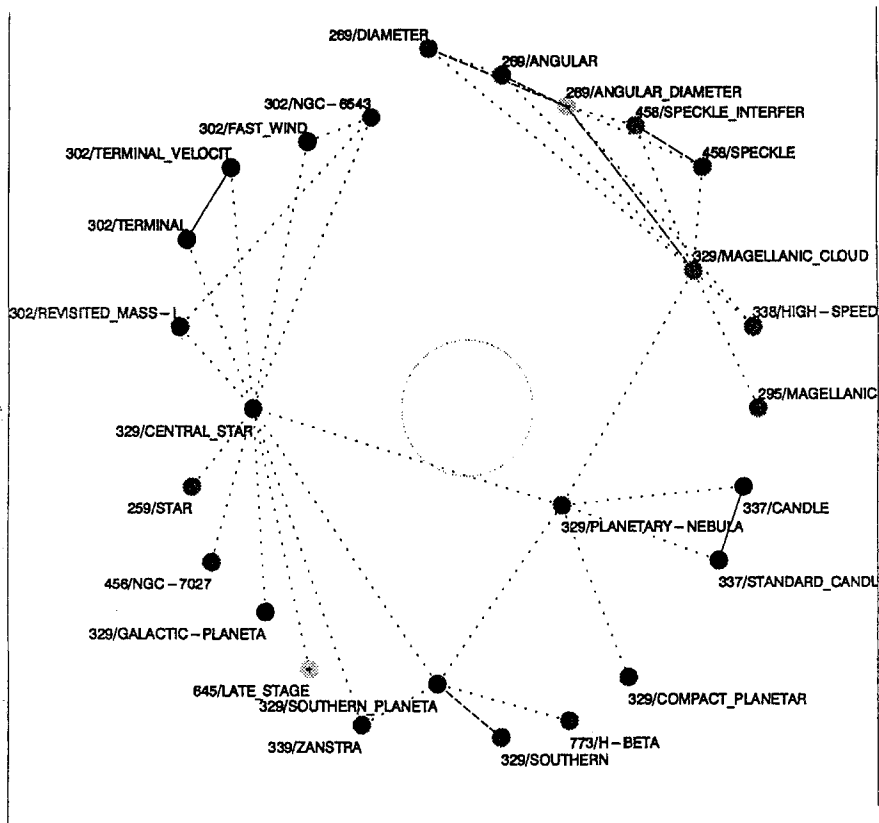


Fig. 3. Words map by theme (theme planet - nebul (329) theme words incl threshold = 0.15)

More links are documented on the classical connection map of Fig. 6 (links, more than visual proximity, show the connections). The population of a front, expressed in number of equivalent-documents assigned to the cluster, is coded by areas of circle symbol. The colors of symbols code the average date. This map is completed by a "size vs connections" diagram on Fig. 7: although connections are present, this diagram is less ambitious than Leximappe™ "strategic diagram" plotting centrality vs density. Here the diagram only looks towards "superclusters", and cluster density is not taken into account; however density, measured by average link, cannot be inferior to the upper cutting value. The interpretation is straightforward along the first diagonal direction, for upper-right (big clusters with strong connections) and lower left quadrant (small isolated clusters). In other quadrants, an interpretation in terms of indifference curves is sometimes possible, when changing the level of observation

induces interconnected small clusters to merge into large ones with fewer links. The front "planetary nebulae", for instance, is very large, but not closely connected with others. The trend coding adds a dynamic aspect to this structural characterization: for example, strongly connected and "hot" themes suggest areas of interest for the near future.

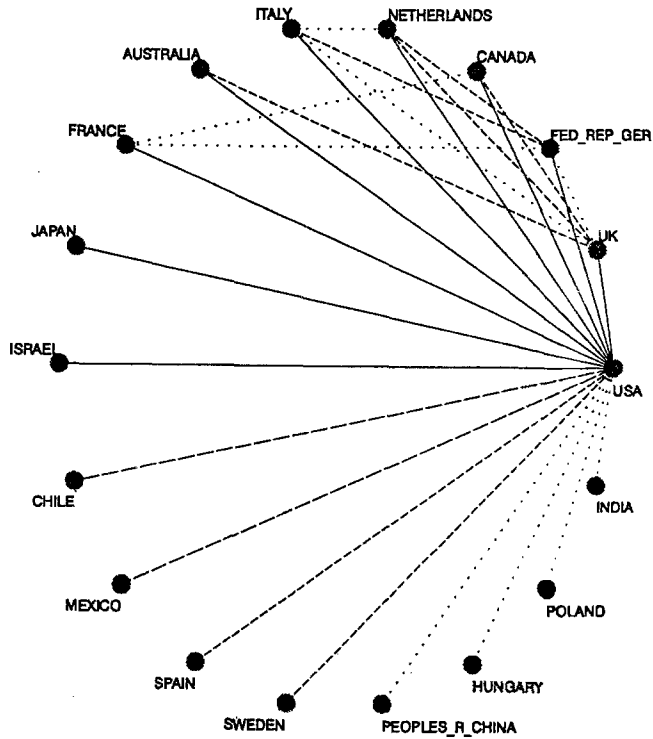


Fig. 4a. Coauthorship all themes (Astroph - Jnal)

Mutatis mutandis, all maps and diagrams above hold for co-word analysis with appropriate settings and interpretations. But the last diagram on Fig. 8 is specific of cocitation studies; it plots trend (average date value of the cluster, ie citing documents) versus immediacy (lag between average date of citing documents and median date of cited "core" documents). Interpretation for strategic purposes is fairly clear on the first diagonal: progressive themes with young science basis are placed in

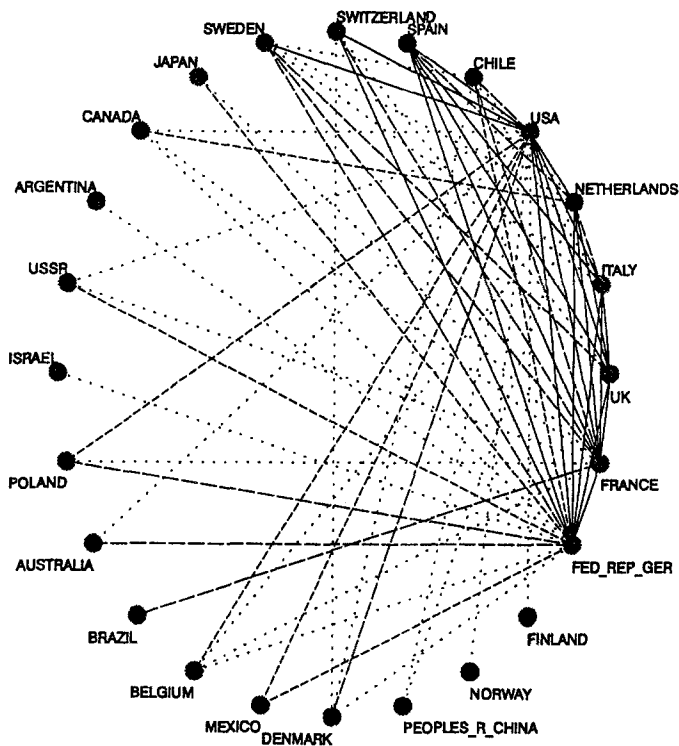


Fig. 4b. Coauthorship all themes (A. & A)

the upper right quadrant, and regressive themes with older basis in the lower left; for legibility, the figure displays the ordinal diagram. In lower right quadrant we can expect classical topics still growing ("planetary nebulae"), or fronts based on ideas revival; regressive themes based on "new science" (upper left) could at first sight signal flashes in the pan. But of course "trend" of clusters must be carefully interpreted: a few research fronts may be intrinsically linked to special events. Astrophysics give several examples of such fronts, e.g. dealing with Halley's Comet (upper left), or with a 1987 supernova (upper right). All these results need to be validated against scale changes, tests of predictive value [9], and experts advice.

SPNOVA - 1987A*656
 NEUTRIN - SNOV*347
 GALCHEM - EYOL*317
 SPECTHOT - CRA*322
 BLUE - COMCAL*322
 HYDROCARBON*450
 DUSTGRAIN - IR*340
 DR 21 - W51 - W3*354
 MAGELL - CLOUD*295
 T - TAURI*280
 HERBIG - AE*248
 BIPOL - OUTFLO*274
 HERBIG - HARC*338
 MOL - CLOUD*285
 INTER - CLOUD*272
 ORION - KL*384
 GAL CENTER*482
 RADIO - CYGNUS*258
 ECLIPS - BINAR*378
 LITHIUM - ABUN*428
 ELEM - ABUNDAN*258
 GLOB - CLUSTER*353
 WOLF - RAYET*259
 NONLTE - HOTST*250
 RAD - DRI - WIND*302
 BE - STAR*318
 GLUESE - FLARE*463
 COOL - STAR*321
 HR 5999PERIO*787
 EARLYTYP - GAL*384
 COOLING - FLOW*248
 ELLIPT - GALAX*284
 GLOBCLUS - DYN*441
 COSM - STRING*267
 MATTER - DARK*294
 CORREL - CLUST*489
 VOID - LS - STRU*422
 TULLY - FISHER*275
 COSM - RAY - ACC*398
 ASTROPHY - JET*319
 SUPERLUMINA*352
 LYMAN - QSO*283
 SEYFERT - GALA*258
 RADIO - CONTIN*257
 BROAD - LINE*418
 QUAS - NUCLEUS*288
 CHROMOM - COSM*280
 PLANET - NEBUL*329
 WHITE - DWARF*332
 NUCLEOS - AL26*372
 CARBON - STAR*289
 SNOVA - TYPE I*553
 STEL - ENVELOP*307
 RL - STAR - MASE*328
 HE 3 RICH*423
 SOL - FLOW - TUB*297
 SO - CHROMOSP*587
 ACCRET - DISK*239
 CATAcly - NOVA*270
 X - RAY - PULSAR*387
 QOSCILLATIO*363
 EXDSAT - X - RAY*290
 PULSAR - RADIO*655
 PULS - MAGNETO*577
 SOLAR - FLARE*182
 MAGN - ENERGY*500
 SWAVE - ALPVEN*320
 GRAVIT - LENS*450
 POUWER - CEPH*370
 HALLEY - WATER*759
 HALL - COMET P*670
 HALLEY - CRAFT*573
 COMET - HALLEY*525

THEMES TREND



Fig. 5. Themes classification and trends

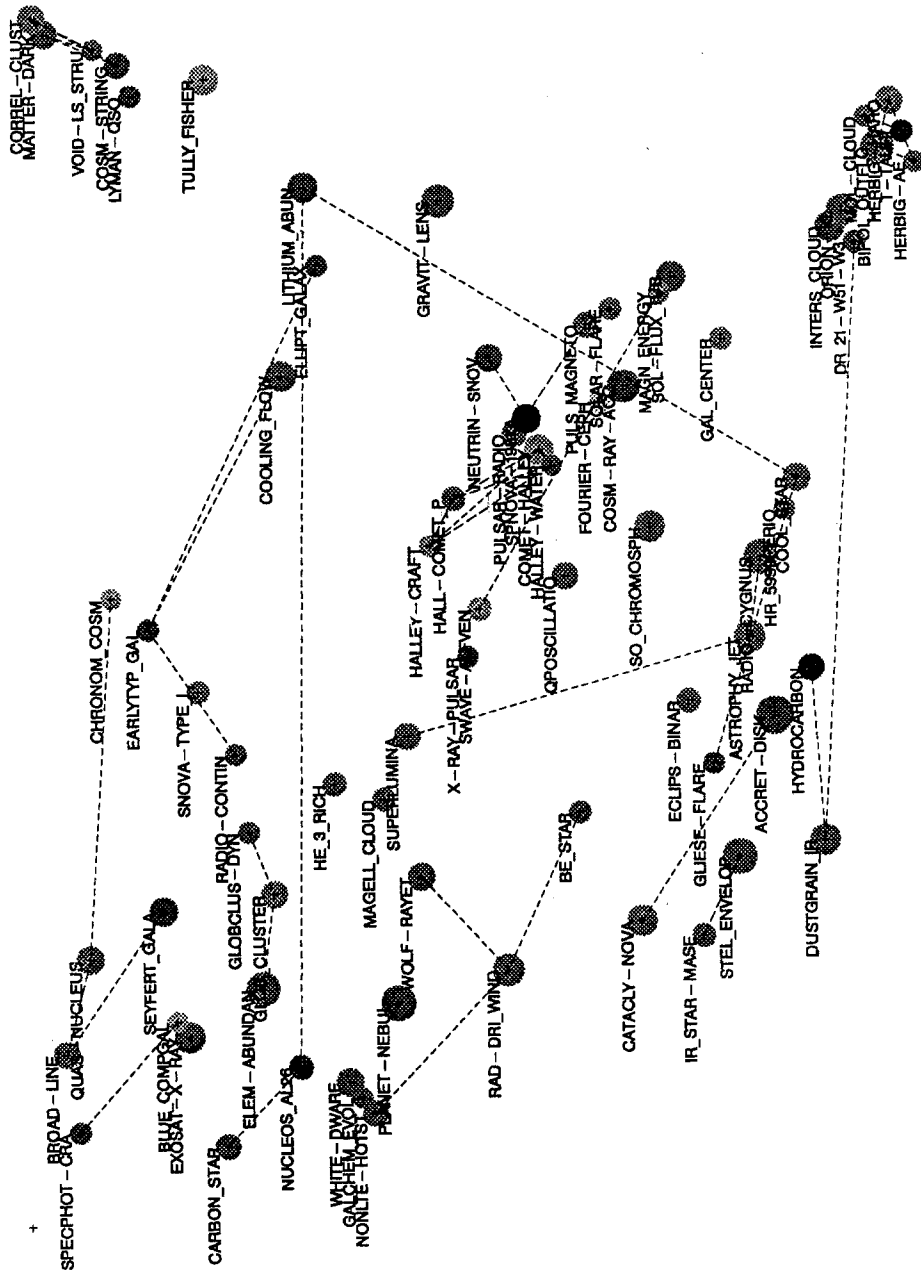


Fig. 6. Themes size, trend and linkage

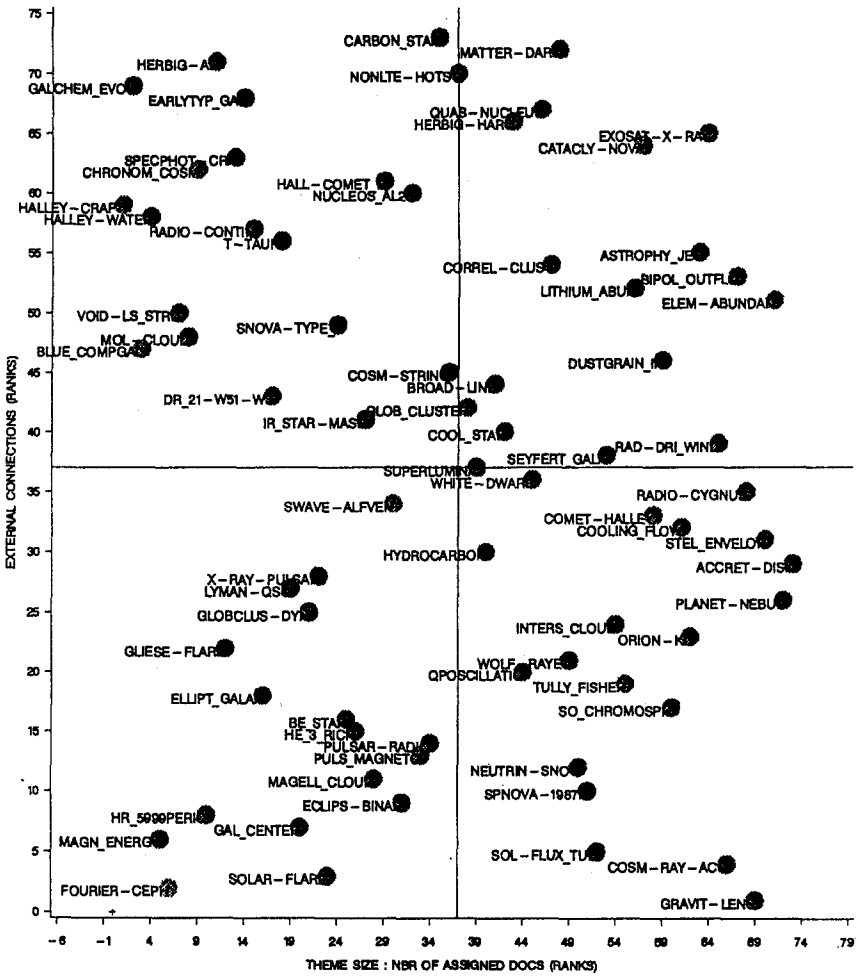


Fig. 7. Size - connections of themes (ordinal) (connect. sum of links, quadrants separators = medians of values)

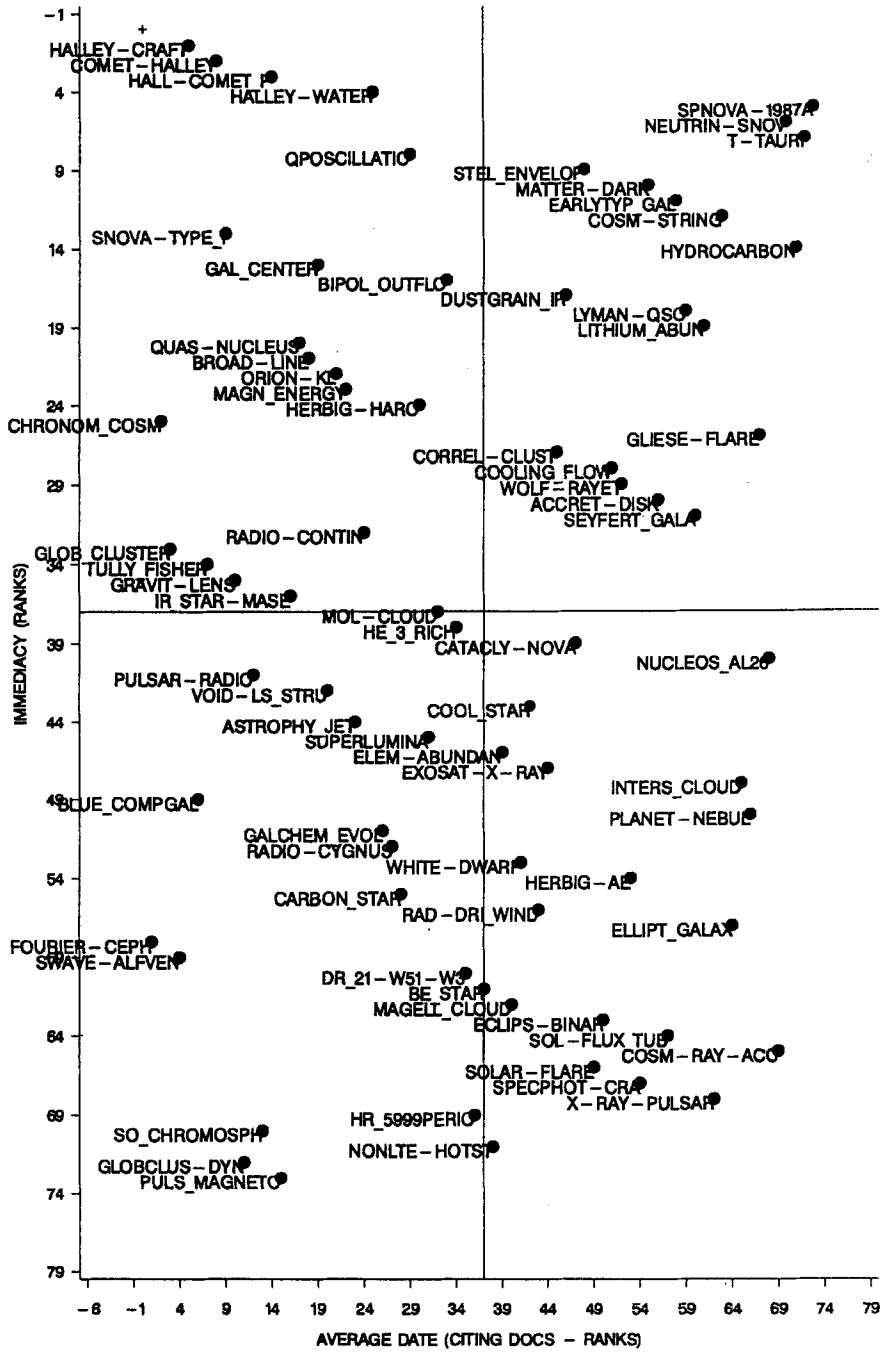


Fig. 8. Trend - immediacy of themes (ordinal), (quadrants separators central values)

Conclusion

The trend analysis of scientific areas, based on usual publication studies, relies heavily on assumptions concerning the meaning of publication date, present in all databases. Depending on which landmark is considered as the best in the history of papers, estimation of fields dynamism through publication dates is subject to specific corrections. Developments of electronic communication could also alter the problem outlines. Keeping this caveats in mind, combining structural and temporal aspects of network analysis (copublications, cocitations, word associations) seems to be very profitable in a bibliometric study of a field; furthermore, in cocitation analysis empirical characterization of trend versus immediacy may raise stimulating questions in term of life-cycle, and back more formal approaches.

*

The authors thank M. *Crance* (CNRS-UNIPS, France) and M. *Qannari* (ENITIAA, France) for their collaboration.

References

1. R. R. BRAAM, H. F. MOED, A. F. J. VAN RAAN, Mapping of Science by combined Cocitation and Word Analysis, I Structural Aspects, *Journal of the American Society of Information Science*, 42 (1991), 233-251.
2. R. R. BRAAM, H. F. MOED, A. F. J. VAN RAAN, Mapping of Science by combined Cocitation and Word Analysis, II Dynamical Aspects, *Journal of the American Society of Information Science*, 42 (1991), 252-266.
3. E. BASSECOULARD, M. ZITT, Une chaîne d'analyse scientométrique dynamique: application à l'étude en longue période d'une revue scientifique, Communication SFBA, L'Ile-Rousse, Corse, juin 1993.
4. J. P. COURTIAL, Qualitative Models, Quantitative Tools and Network Analysis, *Scientometrics*, 15, 5-6 (1989), 527-534.
5. M. ZITT, A simple method for dynamic scientometrics using lexical analysis, *Scientometrics*, 22, 1 (1991), 229-252.
6. J. M. ZIMAN, Information, Communication, Knowledge, *Nature*, 224 (1969), 318-324.
7. W. D. GARVEY, B. C. GRIFFITH, Scientific information exchange in psychology, *Science*, 142 (1964) 1655-1659.
8. C. G. ROLAND, R. A. KIRKPATRICK, Time lapse between hypothesis and publication in the medical sciences, *The New England Journal of Medicine* 292 (1975), 1273-1276.
9. J. CARSON, H. V. WYATT, Delays in the Literature of Medical Microbiology – before and after publication, *Journal of Documentation*, 39, 3 (1983), 155-165.
10. T. C. JAIN, S. P. GOYAL, A Study of the Time-lag in the Publication of Research Papers in some selected Periodicals in Agriculture and allied Sciences, *Annals of Library Science and Documentation*, 16, 1 (1969), 11-14.
11. C. J. PINGS, Publication Delays in the Chemical Engineering Literature, *Journal of Chemical Documentation*, 7,3 (1967) 179-181.

12. C. H. LAY, Publication delay: An analysis of journal days, reviewer days, and author days to revision, *Canadian Journal of Behavioral Science*, 19, 3 (1987) 324-331.
13. H. A. ABT, Publication Practices in Various Sciences, *Scientometrics*, 24, 3 (1992), 441-447.
14. H. G. SMALL, Co-citation in the scientific literature, *Journal of the American Society for Information Science*, 24 (1973) 265-269.
15. I. V. MARSHAKOVA, Document coupling system based on references taken from Science Citation Index (in Russian), *Nauchno - Tekhnicheskaya Informatsiya, Ser. 2* (1973) No. 6,3.
16. H. G. SMALL, E. SWEENEY, Clustering the Science Citation Index using Co-citations, *Scientometrics*, 7 (1985) 391-409.
17. D. HICKS, Limitations of Co-citation Analysis as a Tool for Science Policy, *Social Studies of Science*, 17,(1987) 295-316.
18. K. W. MCCAIN, The Author Cocitation Structure of Macroeconomics, *Scientometrics*, 5 (1983), 277-289.
19. H. PENAN, Analyse des Citations, Principes, Applications à la Théorie Microéconomique, in DESVALS, DOU, *La Veille Technologique*, Dunod 1992, 277-330.
20. E. DAVOUST, H. BERCEGOL, M. CALLON, L'Astronomie, cartographie d'une discipline, in La Scientométrie en action, *Les Cahiers de l'ADEST*, juillet 1993, 44-49.
21. C. JASCHEK, The "Visibility" of West European Astronomical Research, *Scientometrics* 23, 3 (1992) 377-393.
22. V. TRIMBLE, Patterns in Citations of Papers by American Astronomers, *Quarterly Journal of the Royal Astronomical Society*, 34 (1993) 235-250.

Notes

- 1 "Strategic Information through Dynamic Bibliometric Analysis of Areas Development" (SINDBAD), sequence of programmes based on SAS[®] statistical software.
- 2 N=112; Mean=304.7; StdDev=112.8; Skew=1.3; Kurt=2.63; CV=37%; W:Normal=0.908.
- 3 Weighted Ochiai index. Without weighting, a citing document with many references plays an excessive role through generating a large number of citation links. Adversely, a complete compensation of this effect (profile approach where each citing document is given the same weight), can be seen as overvaluing each link generated by documents with few references. Here a midway has been chosen, with a weighting in n, number of cited references. Assignment of citing documents takes the citation frequency and intra cluster position into account.
- 4 Types of document: articles, letters, notes; quasi-reviews articles (>=100 references) were discarded. nr of citing papers: 9152 / nr of cited papers: 66913 / citation threshold for cited papers: 21 (four years) / nr of selected top cited papers: 991 / nr of assigned citing papers: 7023
- 5 Ordinary cluster growth is permitted up to the lower bound, and forbidden above the upper bound of the range; inside the range, growth is restricted by a proper setting of modal rule, so that a cluster is allowed to extend to its "natural frontier" with another one.
- 6 Without detailing this point, let us mention a few remarks. Monotonous changes of scale may affect the internality measure. An interesting type of scale is double log on population vs rank function: the images of all paths from a regular dichotomic tree tend to confound along the diagonal, with a trivial constant internality. Adversely, irregularities in the tree result in scattered paths, and therefore different internalities. The conservative property of internality highly contrasts with metric approach of cluster means. It may be noticed that if one wants the assignment process of citing documents to take into account the respective position of the "core" papers they refer to, the above property leads to a fairly stable process of assignment when changing the level, from clusters to superclusters. Secondly, the internality of the more internal elements (a couple) in a cluster also appears as a possible empirical indicator of the "internal stability" of this cluster (resistance to level

variation). Symmetrically, an external stability can be defined (area above the path, starting from the cluster through the root). It must be stressed that such empirical measures assume the limits of the ultrametric transformation of the data.

- 7 For the cluster detailed in following examples, the complete pseudo-title is PLANETARY-NEBULA/CENTRAL_STAR/SOUTHERN_PLANETARY_NEBULA/MAGELLANIC_CLOUD_PLANETARY...)/NGC_6543/IONIZATION_STRUCTURE/NGC_6826. Contracted form "PLANET_NEBUL" is used in maps. Some more information on other contracted forms :
RAD-DRI-WIND = radiation driven wind/ SO_CHROMOSPH = solar chromosphere/
CORREL-CLUST = superclustering, correlation function, scale-invariant/ QUAS-NUCLEUS = quasar, active galactic nucleus/ QOSCILLATION= quasi-periodic oscillation/
SUPERLUMINA= superluminal motion/ NONLTE-HOTST= non LTE, hot star/ PULS-MAGNETO = pulsar magnetosphere/ NICLEOS_AL26 = nucleosynthesis, AL26/ SWAVE-ALFVEN = surface wave, Alfven/ IR_STAR_MASE = IR star, Maser/ RADIO-CONTIN= radio-continuum emission/ SPECTPHOT_CRA = spectrophotometry, Crab nebula/ HR5999PERIO = HR5999, orbital period/ CHRONOM_COSM = chronometric cosmology/ VOID-LS_STRU = void, large-scale structure/ BLUE_COMPGAL = blue compact galaxy/ GALCHEM_EVOL = galactic chemical evolution
- 8 "Fan" effect, for journals with a strong national preference; this kind of network was also obtained in Food Science for the journal "Science des Aliments"
- 9 Prospective aspects have been tested on other subjects with encouraging results for early variants of the SINDBAD chain.