

Definition of Cancer Research: Journals, Titles, Abstracts or Keywords?

Grant Lewison

University College London, Gower Street, London WC1E 6BT, UK
E-mail: g.lewison@ucl.ac.uk

ABSTRACT

Three versions of a 'filter' used to identify papers on cancer research, as defined by Cancer Research, UK, and interpreted by four experts, have been compared. The first filter was based only on specialist journals and had unacceptably low recall. The second filter was based both on journals and on title words, and had both precision and recall above 0.9. The third filter was based additionally on words in the abstract and/or keywords provided in the paper: it improved the recall to almost unity but the precision was severely degraded, with many false positives. The three filter versions were compared in terms of the outputs of 15 countries in the *Web of Science* in recent years, and in some instances, gave differing indicators of their performance (numbers of papers and citations) which could give conflicting messages for science policy.

Keywords: *Web of Science*, cancer research journals, filters, science policy

1. INTRODUCTION

The evaluation of research often involves making comparisons between the outputs from different actors—countries¹⁻³, regions⁴⁻⁸, universities⁹⁻¹², institutes and departments¹³⁻¹⁵, even persons¹⁶⁻¹⁸ on a number of criteria. These typically include the numbers of papers and some of their parameters, usually including numbers of citations. As an alternative, the output of one or more actors may be compared with the world average, particularly when countries are being compared in terms of their citation performance. It is well attested that the norms of production and of citation vary greatly between fields¹⁹⁻²², and between subject areas within them²³, so that such comparisons must respect these differences if they are to be valid.

Another common activity for bibliometricians is to examine a particular scientific field to determine its dynamics (how fast it is growing relative to all science, for example; e.g., Gupta & Dhawan²⁴), its structure (the relationships between sub-areas and how they are changing, often shown as maps²⁵⁻²⁷, and the principal actors²⁸⁻²⁹). Both of these tasks require the field or subject area to be defined, for details of the relevant papers to be extracted from a database by means of a 'filter', and for the filter to be calibrated in terms of its precision (or specificity) and recall (or sensitivity). Somewhat surprisingly, the first and last of these three jobs are often

omitted. But they are fundamental to a rigorous analysis of a subject area that will command confidence among the study's readership. Very often, the 'filter' simply consists of a set of journals allocated to pre-set subject areas by the database publisher (e.g., the *Web of Science–WoS*^{9,30,31}; *Scopus*³², or determined from cognitive relationships³³). However, now that several databases also contain searchable abstracts of many of the papers that they process, these have sometimes been used to generate additional papers for the analysis^{29,34}. Sets of keywords are also increasingly being added to the paper record—some given by the authors, some by the journal, or by the database provider (e.g., *MedLine*). More complicated filters have also been devised, based on citations either from or to papers to or from a 'core set'³⁵⁻³⁷. If the subject area of interest is not too large, then it may be possible to improve the precision of the filter by inspection of the individual papers with a view to the rejection of ones deemed irrelevant.

The lack of attention to how well the filter performs is surprising, as a poorly-designed filter can give spurious and misleading information about a subject area—how big it is, how well-cited it is, and its structure and the principal actors within it. Moreover, it is often difficult for others to check the stated results and see how sensitive they might be to small changes in the filter used to generate them.

This paper examines three filters that can be used to define the subject of 'cancer research', based first on oncology journals, second on journals and title words, and third on these plus terms in the abstract and keywords. The three successive filters will yield increasing numbers of papers. Which is best in terms of precision and recall? and how much difference does the choice of filter make to the dynamics of the subject area, its citation norms, and the relative ranking of some individual countries? The first task in any work of this type is to provide a simple and clear definition of the subject, (see Appendix in Webster³⁸) and as was done recently for nanotechnology by Maghrebi³⁹. Usually 50-100 words are enough, and these tell readers what is included and what is excluded, so that they know the definition used, even if they might have defined it differently. For cancer research, the definition provided by Cancer Research, UK, a leading charity, was used which reads as follows and has just 53 words: "The study and treatment of cancer or tumours. This incorporates academic oncology and clinical oncology. Academic oncology is aimed at identifying the causative agents or underlying genetic defects producing cancer and at developing these discoveries into effective drugs and other therapies. Clinical oncology is oriented towards the treatment, management and cure of cancer".

2. METHODOLOGY

The process of filter development is, or should be, a progressive process and it need to be tested at each stage to check that precision and recall are improving and approaching unity. The simplest way to start is to select some very obvious title words, or address words, that indicate the subject. In the present study, title words could be *cancer**, *carcinoma**, *leukemi**, *oncol**, *tumor** (where *denotes any character(s) or none) and address words or contractions could be *CANC*, *ONCOL*, *TUMOR*. These were used to search the database (*Web of Science* (WoS), *Science Citation Index Expanded*, which was limited to articles, proceedings papers and reviews) for 2005 and 2009 publication years.

The sources (i.e., journal, year, volume, issue, pages) were then downloaded to file and the names of all the journals were listed that had one or more papers. From this list, all those journals with appropriate strings in their titles, such as *CANCER*, *ONCOL*, *ONKOL*, *LEUKAEM*, *LEUKEM*, *TUMOR* were marked, plus a few others clearly relevant such as *CHEMOTHERAPY*.

The first filter was the list of all these journals, but for practical purposes, it was collapsed into a much shorter set of search strings by the use of asterisks, thus the six journals: *Advances in Cancer Research* or *American Journal of Clinical Oncology-Cancer Clinical Trials* or *Anti-Cancer Agents in Medicinal Chemistry* or *Anti-Cancer Drugs* or *Anticancer Research* or *Asian Pacific Journal of Cancer Prevention* were all be represented by the

contracted statement: *A*CANCER**. This procedure ensured that the journal list was up-to-date, but it could be repeated for an earlier year in order that the filter should capture papers in specialist journals that are no longer in existence, or no longer processed for the database, although the contracted statements mostly did this automatically.

The second filter used both specialist journals and title words. At this stage, it was necessary to engage the services of an expert in the subject area. The titles of all the papers in the specialist journals in the most recent year available were downloaded from the database to a file, and after some cleaning to remove punctuation marks, all the title words were listed in descending order of frequency of occurrence. Many of them were common words not relevant to the subject, but the experts were able to identify relevant ones and mark them. Some necessary to be qualified by the presence (or absence) of another word to ensure that they were used in the correct sense or context. Thus 'tumor' need not to be accompanied by 'necrosis factor' to be relevant to cancer, and 'irradiation' must be accompanied by 'fractionated'. The title words were conveniently sorted alphabetically and formed into a set of search statements, which could be combined with the search statements based on specialist journal names.

The third filter was similar to the second, but the list of words were applied not only to the titles of the papers but also to the abstracts and keywords.

The list of title words, and possibly also the list of specialist journals, needed to be tested to check that it did not generate too many false positives or false negatives. There are several ways to perform this calibration^{40,41} but the simplest is based on the assumption that research teams whose addresses contain one or more of the selected contractions (here, *CANC*, *ONCOL*, *TUMOR*) will publish similar papers to those without such address strings. Three sets of papers were then identified and the bibliographic details (title, source) of samples of them (perhaps a few hundred) were downloaded to file:

- Set A: Papers captured by the filter AND with the contractions in their address(es)
- Set B: Papers with the contractions in their address(es) but NOT captured by the filter
- Set C: Papers captured by the filter but WITHOUT the contractions in their address(es)

The expert was then invited to mark these papers as relevant (1) or not relevant (0); she/he might shade the mark with a decimal fraction for papers where the title did not give enough information for a firm decision to be made. For fairness, it was advisable to mix up papers from the

three sets so that the expert marked them without knowing to which set they belonged - of course, they had hidden codes or other markings so that they could subsequently be identified for analysis purposes.

To calibrate the filter, it was needed to determine the number of missing papers, i.e., ones not captured by the filter and not having the contractions in their addresses. If the number of papers in the database in a selected year or years in set A is a , and the precision of this set, based on the sample, is $p(a)$, then the true number of papers = $a^* = a \times p(a)$, and similarly for b^* and c^* .

The above assumption yields $d^* = c^* \times b^*/a^*$, and the true total of papers is $a^* + b^* + c^* + d^*$. The true number retrieved is $a^* + c^*$ and the actual number is $a + c$, so precision $p = (a^* + c^*)/(a + c)$ and recall $r = (a^* + c^*)/(a^* + b^* + c^* + d^*)$. Filter development proceeded in steps, and at each step it was necessary to check that p and r were increasing until a point was reached when gains in one were offset by losses in the other.

The calibration factor, $CF = p/r$, and might be either greater than or less than unity.

3. RESULTS

3.1 Precision and Recall of Filters

The latest version of the cancer research filter, labelled ONCOL, is actually the sixth, earlier versions having been supplemented with the names of new drugs and newly-discovered genes that code for an increased cancer risk. It is now quite complex, with 55 journal search strings, 8 title/abstract words with Boolean conditions, and 293 single title/abstract words or pairs (e.g., xeroderma pigmentosum). When applied to the WoS for 2009, the numbers of papers identified were as in Table 1.

For the first version of the filter, based only on specialist journals, the calculation of the numbers of cancer research papers and p and r are given in Table 2. For this filter, $p = (14991 + 9494)/25512 = 0.96$, and $r = 25512/30705 = 0.83$. This might be deemed fairly satisfactory, but there was need to investigate the other

Table 1. Numbers of papers retrieved from WoS by three versions of the ONCOL filter, publication year = 2009

Filter	Based on	Set A	Set B	Set C	Retrieved
1	Journals only	15147	28629	10365	25512
2	Journals & titles	28374	15402	40081	68455
3	Journals, titles, abstracts, keywords	35143	8633	89283	124426

two versions before relying on the assumption that researchers in eponymous (cancer) departments publish in a similar range of journals to those in non-eponymous departments.

The second version of the filter, which used title words as well as specialist journals, gave a much larger estimate of the size of the cancer research output as shown in Table 3.

For this version, $p = 0.93$ and $r = 0.93$, but the estimated true total is more than twice as large as with the first version. Evidently, there are many cancer research papers not published in specialist journals—in fact, the majority.

Table 2. Calculation of precision, p , and recall, r , of the first version of the filter

Set	n	Sample	OK	p	n^*
A	15147	213	210.8	0.990	14991
B	28629	524	69.7	0.133	3808
C	10365	80	73.3	0.916	9494
D					2412
Total	25512				30705

Table 3. Calculation of precision, p , and recall, r , of the second version of the filter

Set	n	Sample	OK	p	n^*
A	28374	509	492.4	0.967	27438
B	15402	524	69.7	0.133	2048
C	40081	470	422.4	0.899	36033
D					2690
Total	68455				68209

Table 4. Calculation of precision, p , and recall, r , of the third version of the filter

Set	n	Sample	OK	p	n^*
A	35143	1008	830.5	0.824	28947
B	8633	1000	6.5	0.0065	56
C	89283	967	456.4	0.472	42134
D					81
Total	124426				71218

Finally, the terms in the filter were also applied to abstracts and keywords and the results are presented in Table 4. The estimated true total is now somewhat larger, but only by 4.4 per cent. The addition of papers retrieved because of words in the abstract or keywords has apparently given almost complete retrieval ($r = 0.998$) but the precision is now severely degraded to $p = 0.571$. It is reasonable to conclude that the second version of the

Table 5. Numbers of cancer research papers (articles, proceedings papers and reviews) for 15 leading countries in the WoS, publication years 2005 through 2009, according to the three versions of the filter: Percentages of world total

ISO	Country	2005 (F1)	2005 (F2)	2005 (F3)	2009 (F1)	2009 (F2)	2009 (F3)
	World (papers)	21236	53072	96106	25512	68455	124426
US	United States	41.95	38.21	38.51	38.77	34.59	34.70
JP	Japan	10.50	11.12	10.64	9.67	9.19	8.70
DE	Germany	9.10	9.20	8.99	8.32	8.12	8.20
UK	United Kingdom	8.11	7.73	7.78	7.47	6.97	7.21
IT	Italy	7.60	6.55	5.85	7.22	6.63	6.07
FR	France	5.53	5.58	5.53	6.12	5.54	5.42
CA	Canada	4.82	4.24	4.34	4.93	4.28	4.44
CN	China (P. R.)	3.18	3.71	3.82	7.43	8.11	8.18
NL	Netherlands	4.14	3.37	3.02	4.09	3.20	2.89
ES	Spain	2.68	2.71	2.64	3.06	3.09	3.10
KR	South Korea	1.99	2.56	2.59	3.51	4.29	4.06
SE	Sweden	2.71	2.29	2.16	2.49	2.00	1.86
AU	Australia	2.35	2.24	2.31	2.91	2.57	2.58
CH	Switzerland	1.76	1.71	1.82	2.19	1.85	1.89
BE	Belgium	1.78	1.55	1.48	1.88	1.45	1.41

F1, F2, F3-three versions of filter

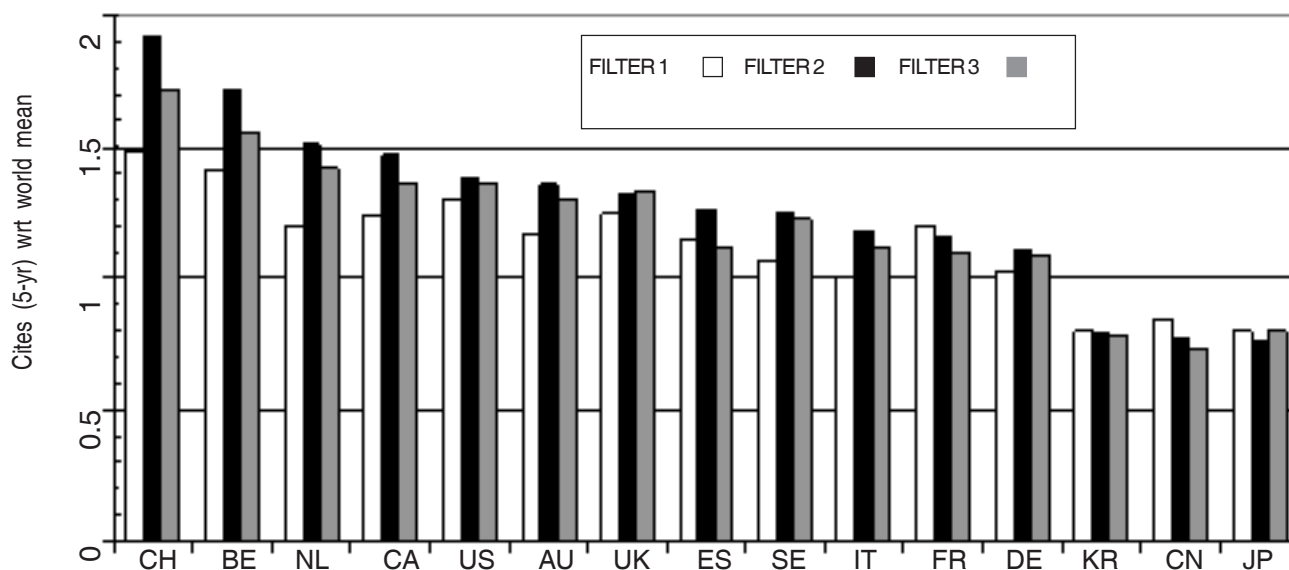


Figure 1. Five-year citation scores relative to the world mean values for cancer research papers from 15 leading countries (codes given in Table 5) published in 2005 and cited 2005 thru 2009, based on three versions of the cancer research filter.

filter is the most nearly correct one, as version 1 has a low recall (based on the results of versions 2 and 3) and version 3 a low precision.

3.2 Comparison of Outputs of Leading Countries

The three versions of the filter were applied to the WoS for publication years 2005 and 2009, and the number of papers world-wide and from 15 leading countries were determined. The results are shown in Table 5, and are given as percentages of the world totals using integer (whole) counting.

Between 2005 and 2009, world cancer research output increased by 20 per cent according to version 1 of the filter, but by 29 per cent according to both version 2 and version 3. Since these two versions gave very similar totals (about 70,000 papers per year), the latter growth rate can be accepted rather than the former. This means, incidentally, that an increasing percentage of cancer research papers are not published in specialist cancer journals but in general journals.

The results for the individual countries are somewhat varied, as would be expected. Some countries have similar percentage presence in the world on all three versions of the filter, such as France (in 2005), Germany, and Spain. A few have a higher presence according to the fuller versions of the filter ($F3 > F2 > F1$), such as China and South Korea; but most show the reverse, with a higher presence in the specialist journals and a lower one in the titles, abstracts and keywords, notably Italy, the Netherlands, Sweden, and Belgium. Between 2005 and 2009, despite the steady increase in international co-authorship, the four leading countries in Table 5 (the USA, Japan, Germany, and UK), and also the Netherlands and Sweden, all decreased their percentage presence in cancer research according to all three versions of the filter. Six countries (Canada, China, Spain, South Korea, Australia, and Switzerland) all increased their presence, again based on all three filter versions. But for the other three (France, Italy, and Belgium) the message was mixed, and the change could have been reported as either a gain or a loss of presence.

The rating of countries was determined based on the mean citation scores of their papers. In Fig. 1, these have all been compared with the world mean values in a five-year window, i.e., the numbers of citations in 2005 thru 2009 for the 2005 publications. These were respectively 18.35, 16.1, and 16.04 cites for filter versions 1, 2 and 3. It appeared that papers in the specialist cancer journals received more citations than ones in the general journals that were retrieved because of their titles or abstracts/keywords.

The countries have been ordered in Fig. 1 on the basis of their citation performance on the second version of the filter, and this ranking puts three small European countries ahead of Canada and the USA. Their performance, and that of Canada and several other European countries, is much better than that shown by the specialist journals, where the USA shows to advantage, but is still behind Switzerland and Belgium. The three East Asian nations all score relatively low on all three filter versions, as has been found elsewhere⁴².

4. CONCLUSIONS

This paper has examined one particular field, namely, cancer research, in some detail and has shown that the world output in the WoS was of the order of 70,000 research papers per year in 2009. The best filter in terms of both precision and recall was one based both on specialist journals and title words. The omission of title words meant that fewer than half the relevant papers were identified, and the addition of words in abstracts and/or keyword lists was not helpful as nearly all the additional papers identified were false positives.

The effects of using version 1 or version 3 of the filter instead of version 2 were rather variable, and some countries benefited in terms of their percentage presence or relative citation score, and some were disadvantaged. Few of the differences were large, but countries are often looking^{43,44} for evidence of small improvements to their relative position in order to claim that their science is being well managed and providing good value for money, as with the European agri-food research programmes⁴⁵, where different search strategies sometimes produced very different outcomes. It seems important, therefore, that any such claims should be based on the best approximation to the true set of papers in the selected field or subject area, even though it is really a 'fuzzy set' rather than one that can be precisely defined without argument.

ACKNOWLEDGEMENTS

The author is grateful to the four cancer specialists who kindly marked the sets of papers from the three filters: Dr Lynne Davies of Cancer Research, UK; Dr Katie Hyde of the National Cancer Research Institute; and Professor Arnie Purushotham and Professor Richard Sullivan of King's College, London. The author also thanks Dr Philip Roe who wrote the macros that enabled him to download and analyse the papers on the *Web of Science* and their citations.

REFERENCES

1. May, R.M. The scientific wealth of nations. *Science*, 1997, **275**, 793-96.

2. King, D.A. The scientific impact of nations. *Nature*, 2004, **430**, 311-16.
3. Prathap, G. An iCE map approach to evaluate performance and efficiency of scientific production of countries. *Scientometrics*, 2010, **85**, 185-91.
4. Lewison, G. The scientific output of the EC's less favoured regions. *Scientometrics*, 1991, **21**, 383-402.
5. Shapira, P.; Youtie, J. & Mohapatra, S. Linking research production and development outcomes at the regional level. *Research Evaluation*, 2003, **12**, 105-16.
6. Holbrook, J.A. & Clayman, B.P. Research funding by city: An indicator of regional technological competitiveness. *Research Evaluation*, 2006, **15**, 221-31.
7. Zhou, P.; Thijs, B. & Glänzel, W. Regional analysis on Chinese scientific output. *Scientometrics*, 2009, **81**, 839-57.
8. Levitt, J.M. & Thelwall, M. Does the higher citation of collaborative research differ from region to region? A case study of economics. *Scientometrics*, 2010, **85**, 171-83.
9. Adams, J. Benchmarking international research. *Nature*, 1998, **396**, 615-18.
10. Prathap, G. & Gupta, B.M. Ranking of Indian universities for their research output and quality using a new performance index. *Current Science*, 2009, **97**, 751-52.
11. Aguillo, I.F.; Bar-Ilan, J.; Levene, M. & Ortega, J.L. Comparing university rankings. *Scientometrics*, 2010, **85**, 243-56.
12. Docampo, D. On using the Shanghai ranking to assess the research performance of university systems. *Scientometrics*, 2011, **86**, 77-92.
13. Yi, C.G. & Kang, K.B. Developments of the evaluation system of government-supported research institutes in Korean science and technology. *Research Evaluation*, 2000, **9**, 158-79.
14. Hyvarinen, J. Evaluation of TEKES funding for research institutes and universities – the role of talent. *Research Evaluation*, 2009, **18**, 365-73.
15. Tores-Salinas, D.; Lopez-Cozar, E.D. & Jimenes-Contreras, E. Ranking of departments and researchers within a university using two different databases: *Web of Science* versus *Scopus*. *Scientometrics*, 2009, **80**, 761-74.
16. Qiu, J.P.; Ma, R.M. & Cheng, N. New exploratory work of evaluating a researcher's output. *Scientometrics*, 2008, **77**, 335-44.
17. Alonso, S.; Cabrerizo, F.J.; Herrera-Viedma, E. & Herrera, F. *hg-index*: A new index to characterise the scientific output of researchers based on the *h*- and *g*-indices. *Scientometrics*, 2010, **82**, 391-400.
18. Claro, J. & Costa, C.A.V. A made-to measure indicator for cross-disciplinary bibliometric ranking of researchers' performance. *Scientometrics*, 2011, **86**, 113-23.
19. Schubert, A. & Braun, T. Cross-field normalization of scientometric indicators. *Scientometrics*, 1996, **36**, 311-24.
20. Kostoff, R.N. Citation analysis cross-field normalisation: A new paradigm. *Scientometrics*, 1997, **39**, 225-30.
21. Nederhof, A.J. & Visser, M.S. Quantitative deconstruction of citation impact indicators—waxing field impact but waning journal impact. *Journal of Documentation*, 2004, **60**, 658-72.
22. Podlubny, I. Comparison of scientific impact expressed by the number of citations in different fields of science. *Scientometrics*, 2005, **64**, 95-99.
23. Zitt, M.; Ramanana-Rahary, S. & Bassecoulard, E. Relativity of citation performance and excellence measures: From cross-field to cross-scale effects of field-normalisation. *Scientometrics*, 2005, **63**, 373-401.
24. Gupta, B.M. & Dhawan, S.M. Condensed matter physics: An analysis of India's research output, 1993-2001. *Scientometrics*, 2008, **75**, 123-44.
25. Cambrosio, A.; Keating, P.; Mercier, S.; Lewison, G. & Mogoutov, A. Mapping the emergence and development of translational cancer research. *Euro J. Cancer*, 2006, **42**, 3140-148.
26. Calero-Medina, C. & Noyons, E.C.M. Combining mapping and citation network analysis for a better understanding of the scientific development: the case of the absorptive capacity field. *Journal of Informetrics*, 2008, **2**, 272-79.
27. Jeong, S. & Kim, H-G.. Intellectual structure of biomedical informatics reflected in scholarly events. *Scientometrics*, 2010, **85**, 541-51.
28. Karisiddappa, C.R., Gupta, B.M. & Kumar, S. Scientific productivity of authors in theoretical population genetics. *Scientometrics*, 2002, **53**, 73-93.

29. Kostoff, R.N. & Morse, S.A. Structure and infrastructure of infectious agent research literature: SARS. *Scientometrics*, 2011, **86**, 195-209.
30. Ugolini, D. & Mela, G.S. Oncological research overview in the European Union. A 5-year survey. *Euro. J. Cancer*, 2003, **39**, 1888-894.
31. Sooryamoorthy, R. Scientific publications of engineers in South Africa, 1975-2005. *Scientometrics*, 2011, **86**, 211-26.
32. Gupta, B.M. & Dhawan, S.M. Status of India in science and technology as reflected in its publication output in the Scopus international database, 1996-2006. *Scientometrics*, 2009, **80**, 473-90.
33. Leydesdorff, L. The delineation of nanoscience and nanotechnology in terms of journals and patents: A most recent update. *Scientometrics*, 2008, **76**, 159-67.
34. Meyer, M.; Debackere, K. & Glänzel, W. Can applied science be "good science"? Exploring the relationship between patent citations and citation impact in nanoscience. *Scientometrics*, 2010, **85**, 527-39.
35. Schwechheimer, H. & Winterhager, M. Mapping interdisciplinary research fronts in neuroscience: A bibliometric view to retrograde amnesia. *Scientometrics*, 2001, **51**, 311-18.
36. Glänzel, W.; Janssens, F. & Thijs, B. A comparative analysis of publication activity and citation impact based on the core literature in bioinformatics. *Scientometrics*, 2009, **79**, 109-29.
37. Bolaños-Pizarro, M.; Thijs, B. & Glänzel, W. Cardiovascular research in Spain: A comparative scientometric study. *Scientometrics*, 2010, **85**, 509-26.
38. Webster, B.M. International presence and impact of the UK biomedical research, 1989-2000. *Aslib Proceedings*, 2005, **57**, 22-47.
39. Maghrebi, M.; Abbasi, A.; Amiri, S.; Monsefi, R. & Harati, A. A collective and abridged lexical query for delineation of nanotechnology publications. *Scientometrics*, 2011, **86**, 15-25.
40. Lewison, G. The definition of biomedical research subfields with title keywords and application to the analysis of research outputs. *Research Evaluation*, 1996, **6**, 25-36.
41. Lewison, G. The definition and calibration of biomedical subfields. *Scientometrics*, 1999, **46**, 529-37.
42. López-Illescas, C.; de Moya-Anegón, F. & Moed, H.F. The actual citation impact of European oncological research. *Euro. J. Cancer*, 2008, **44**, 228-36.
43. Grant, J. & Lewison, G. Government funding of research and development. *Science*, 1997, **278**, 878-79.
44. May, R.M. Government funding of research and development—response. *Science*, 1997, **278**, 879-80.
45. Borsi, B. & Schubert, A. Agrifood research in Europe: A global perspective. *Scientometrics*, 2011, **86**, 133-54.