1

# How thematic maps can assist collection management: A qualitative assessment of Journals' thematic focus

Fidelia Ibekwe-SanJuan

*Department of Information-Communication, University of Lyon 3, 4, cours Albert Thomas, 69008 Lyon, France*

2

3

4

5

6

## Abstract

We present a method for mapping the content of a text collection. This method uses linguistic analysis to relate terms extracted from the texts and clusters them into thematic topics mapped onto a 2D space. While the graphic display of domain topics is useful for several information-driven tasks, the focus of the paper is more on the comparison of journal ranking by productivity (number of published papers in the collection) and by content representativity (ranking by number of terms and clusters). The results show that the two rankings are not identical, thus pointing to possible discrepancies between pure productivity and terminological density.
© 2005 Published by Elsevier Inc.

*Keywords:* Journal collection management; Content analysis; Data analysis; Thematic maps; Query refinement

7

8

9

10

11

12

13

14

15

16

17

## 1. Introduction

18

The issue of journal representativity vis-à-vis fields of knowledge is a crucial one for library collection management. Identifying the leading journals in a field and thus the journals to subscribe to has been a constant preoccupation for librarians and information scientists as a whole. This problem was addressed as early as 1934 by the world famous Bradford's law. Bradford found in essence that about 10% of the journals publishing in a field are responsible for producing 90% of the articles in that field. To recover the missing 10% of the articles,

19

20

21

22

23

24

about 90% of journals are needed. A lot of research has been carried out around modeling    25
Bradford' law to suit different situations. In the same vein, the *Journal Citation Report* (*JCR*)    26
computes impact factors of journals to measure the actual use by scientists of works published    27
by certain journals. Quoting the Institute for Scientific Information (ISI), Giles C.L writes    28
"the impact factor is a measure of the frequency with which the average article in a journal    29
has been cited for a particular year (actually averaged over 2 years) and is calculated dividing    30
the number of citations to articles published in two previous years by the total number of    31
articles published in those years. This produces a normalized parameter so that small and    32
large journals can be compared."    33

Among the target users of impact factors are librarians who have to "manage and maintain    34
journal collections and budget for subscriptions." The *JCR* covers more than 7500 most    35
highly cited, peer-reviewed journals in approximately 200 disciplines, 3300 editors across 60    36
countries.    37

Undoubtedly, tools like *JCR* and Bradford's law are important at the macrolevel for    38
selecting core journals in the disciplines covered by a library and thus for collection    39
management scheme based on journals representativity. However, for content-level analysis    40
of journal representativity per topic (specialties within disciplines), a microlevel and fine-    41
grained approach is needed. Such an approach can actually "enter into" the texts of articles    42
published by journals and map out the core topics. This can be utilized in specialized    43
collection management where identifying core journals is not the issue (they would already    44
have been identified using *JCR* or Bradford's law) but librarians or information scientists    45
actually need to understand from what angle and on what specific topics the subscribed    46
journals make publications on. This could be a further criteria for ranking journal relevance    47
for specific users needs (highly specialized libraries or libraries with different categories of    48
users, needing different levels of expertise).    49

We propose to this end, a thematic mapping system developed by Ibekwe-SanJuan and    50
SanJuan [10], which takes as input raw texts from a journal collection and returns topic maps    51
represented in a 2D space, which can be used to synthesize the contents of the text collections.    52

After a review of related works on automatic theme mapping in the Related work section,    53
an overview of the TermWatch system is given in the System overview section. The Mapping    54
domain topics from a collection of IR journals section shows the application of TermWatch to    55
a collection of bibliographic records in the information retrieval field. The Conclusion section    56
explores how the clusters obtained can be mapped onto the source journals of the texts in    57
order to gauge their representativity with regard to the specific topics identified through the    58
clusters. As parts of this research have been published elsewhere [10,11], this paper will focus    59
on a new dimension: possible application of thematic mapping to assist library journal    60
collection management.    61

## 2. Related work    62

Evaluating the state of the art of research in a scientific or technical field has been the    63
object of research since the early sixties. This has led to the emergence of bibliometrics in    64

ARTICLE IN PRESS

*F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx*                           3

1969 then to scientometrics in 1977 and later to informetrics. Today other objects of metrics have appeared: cybermetrics or webometrics. The two major methods used in these studies are the co-citation [17,20] and co-word analyses [3]. Co-citation analysis remains the most popular measure of author–journal contribution to a specific field [18]. After generating a matrix of co-occurrence of citations or of keywords, the underlying methods use a clustering algorithm to reduce the information space and obtain clusters of frequently co-occurring authors, journals, or keywords. These clusters are then mapped onto a 2D space in order to depict their layout and understand the scientific structure of the discipline surveyed. These methods have proved their utility at the macrolevel where entire disciplines are mapped out in order to perceive the social networks and leading actors of the field. They are, however, not targeted to fine-grained content analysis of sub-specialties, thus not very successful at the microlevel. One of the reasons is that clustering being based on occurrences, they need high occurrence thresholds in order to obtain meaningful results.

Clustering techniques are also used in the Information retrieval community (IR) and can be traced to Salton [16], Jardine and Van Rijsbergen [7], and Sparck Jones [19]. The underlying assumption is widely known as the "cluster hypothesis," which postulates that "closely associated documents tend to be relevant to the same requests." The basic approach to clustering in IR consists in partitioning a collection of documents into many small clusters or groups. The intent being later to map user queries to the most similar cluster. This is particularly useful in a context where users do not know a priori which search words to use or do not know the contents nor the indexing vocabularies of the database, as is the case with very large databases or the Internet. Clustering has also been used to address the specific issue of query expansion. Query expansion consists in formulating new query terms using the relevant set of documents. Thus, there is an underlying notion of cluster in this activity: it is hypothesized that the relevant "cluster" of documents "contains terms which can be used to describe a larger cluster of relevant documents [2]." Some techniques like the latent semantic indexing model have been introduced to this end [4]. Another domain in IR, which makes use of clustering, is the presentation of results of a query. Hearst [8] reviewed methods of text categorization or of clustering that enhance the presentation of retrieval results. The aim of these studies is not to explain the layout of research topics but to present groups of "similar" documents in answer to a user's request. To enhance this presentation, considerable interest is being given recently to the use of graphic display interfaces offering 2D or 3D facilities to enable users identify the situation of the relevant documents. Recently, clustering methods are being applied to information search on the Internet [22] and also to gene expression data in the bioinformatics field [21].

The aim of the TermWatch system [10] system is similar to that of co-citation analysis and co-word analyses. However, the thrust here is on text clustering through prior linguistic processing. Consequently, the system builds on recent advances in computational linguistics and particularly on computational terminology to enhance the input to the clustering scheme. Typically, the end user, a domain specialist wishes to know what are the major topics contained in a huge corpus, what topics are evolving and how each topic is related to one another. He or she needs a global view, a map of the domain research

ARTICLE IN PRESS

4                    F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx

topics embodied in the corpus. Additionally, he or she may want to see factual information (authors, laboratories, countries) on each of the topic nodes. The novelty of TermWatch over existing thematic mapping systems is that clustering is based not on co-occurrence of text units but on linguistic relations among them. It focuses on mapping the text content whereas dominant bibliometrics and scientometrics focus on factual data (author co-citation counts, country's or laboratory's publication counts, etc.). In these cases, there is no major difficulty in extracting the units to be counted.

## 3. System overview

The TermWatch system is a joint research program between two associate professors from two French universities, University of Lyon 3, and University of Metz (LITA). TermWatch has three major components: a term extractor, a linguistic relations miner, and a clustering module.

### 3.1. Term extraction module

This module extracts terminological units directly from the text collection to be analyzed. The terms extracted reflect the different topics addressed in each text, and thereby the different topics in the whole text collection. Terms should be taken here in their terminological sense (i.e., text units that refer to domain concepts or objects). Our term extraction rules rely on the recent research in the computational terminology field [1]. Most terms appear as noun phrases (NPs) although some verb and prepositional phrases can be terms. We currently extract only terminological NPs, which are multiword expressions that can appear as compounds (information retrieval system) or as syntagmatic NPs with prepositional attachments (special terminology of information science). Term extraction is performed in using the LTPOS tagger developed by the University of Edinburgh. LTPOS is a probabilistic part-of-speech tagger based on Hidden Markov Models. It has been trained on a large corpus and achieves an acceptable performance. It uses the Penn Treebank tag set, which ensures portability of the output with many other systems. Since LTCHUNK, a component of this system only identifies simplex NPs without prepositional attachments, we wrote contextual rules based on the output of the chunker to identify complex terminological NPs.

### 3.2. Linguistic relations miner

In order to cluster the extracted terms, this module searches for meaningful linguistic relations between them. The idea is that clustering can be performed based on other dimensions than the co-occurrence one. This dimension being linguistic will ensure the semantic coherence of the terms gathered into one cluster. To this end, we studied a variety of linguistic operations, which have come to be known in the terminology community as "variations." Systems aiming to extract domain terms need to address the variation issue in

**ARTICLE IN PRESS**

*F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx* 5

order to capture the actual state of a domain's terminology. This is particularly useful in several applications like acquisition of domain terminology and update, automatic indexing, question answering, information extraction, knowledge representation, and scientific and technological watch. Variations are local morphosyntactic and semantic operations affecting the form and structure of an existing term, thus yielding new terms, which are close to the initial one. Variations cover a wide spectrum of linguistic phenomena occurring at different linguistic levels, thus making their identification impossible without integrating computational linguistics techniques.

At the morphological level, we have spelling variants (specialization/specialization; centre/center; programme/program); inflection variants (academic library/academic libraries) including derivational morphological variants with prefix and suffix addition (tumor promoter/tumor promotion); abbreviations (www/World Wide Web); and compounding process (online Web access/on line Web access/on-line Web access).

Syntactic variants involve structural or formal changes in a term (information retrieval, retrieval of information, efficient retrieval of information), the addition of new modifier or head words in an existing term, that is, syntactic variants of "academic library" found in the corpus are *Canadian academic library privilege*, *changing culture in academic library*, *electronic communication in academic library*, *greater utilization of academic library service*, *Hellenic academic library link*, *service in Malaysian academic library*, *directors of academic library*, *future of academic library*. These relations can be distinguished according to the grammatical function of the word affected: head variation involves the addition or substitution of a new head word in a term as in "academic library" and "directors of academic library" whereas modifier variations implies that only modifier words are affected as in "academic library" and "Canadian academic library." Modifier and head roles are determined by the position of constituent words in a term.

Although morphological and syntactic variants also hold semantic relations, there are explicit semantic variants, which can be realized by surface linguistic markers. For instance, in the following sentences, the sequence "such as" signals a hypernym/hyponym (generic/specific) relation between the NP found on its left (nonlinear systems) and the following one (robotic manipulations). Likewise, the sequence "known as" creates a synonymy relation between "mathematical operation" and "convolution." These relational markers have been studied by Hearst [8], Morin and Jacquemin [15].

(1) The main motivation for this design was to control some known nonlinear systems, such as robotic manipulators, which violate the conventional assumption of the linear PID controller.
(2) This combination is performed by a mathematical operation known as convolution.

Given that all the semantic relations existing between domain terms may not be realized through surface linguistic markers, it is necessary to complete the semantic relation mining using an external resource such as WordNet [5]. WordNet is a general-purpose lexical taxonomy with synonymy, hypernym, and association (see also) relations between words. Synonymous words are gathered into the same "synsets" (classes of words used in the same

ARTICLE IN PRESS

6                    F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx

sense). In our system, WordNet is used to look up word–word relationship between terms 188
when the two terms share common lexical elements but differ by one word. For instance, 189
WordNet, enabled us to establish a relationship between "automatic categorization" and 190
"automatic classification" because "categorization" and "classification" were found in the 191
same WordNet synset. At the moment, morphosyntactic relations and WordNet semantic 192
relations have been implemented in this module. 193

All the relations mined between terms allow us to build a graph of term variants, which 194
serve as input to the clustering algorithm. 195

### 3.3. Clustering module                                                                      196

TermWatch implements a clustering approach, Classification by Preferential Clustered 197
Link (CPCL) presented in Ibekwe-SanJuan [13]. It works in two stages. A first level of 198
clustering consists in grouping together terms sharing the same headword and semantic 199
relations (either given by an external resource like WordNet or harvested through other 200
lexico-syntactic patterns). This results in connected components. For instance the following 201
terms were put into the same component "information department, information science 202
department, Sheffield University's information department." The result of the component 203
building stage is a monothematic organization, which is not the desired result. What we seek 204
to highlight is the transversal relation between these lone themes (i.e., what associations have 205
the authors been making between these themes?) To highlight these association, we now 206
cluster the connected components into classes using the second subset of variation relations, 207
those that involve a shift in the head noun, thus a shift in the topical focus of the noun phrase 208
(NP) as in "academic library" and "Canadian academic library privilege." Like in most 209
clustering methods, we need to compute a similarity index in order to build clusters. This 210
coefficient is defined as follows: 211

$$d(i,j) = \sum_{R \in \text{CLAS}} \frac{N_R(i,j)}{|R|}$$

where $N_R(i,j)$ denotes the number of $R$ variations between two connected components $i$ and $j$. 213
A notable difference with other clustering algorithms is that we do not compute this index on 214
the list of terms, but on the set of connected components. The user can set the number of 215
iterations at which the algorithm is stopped and the minimal similarity index to be considered 216
or let the algorithm converge and then choose the results of a given iteration. 217

The results of the clustering are mapped onto an integrated visualization tool, *Aisee* 218
(http://www.aisee.com). The system architecture is given in Fig. 1. 219

## 4. Mapping domain topics from a collection of IR journals                                   220

The text collection used in this experiment consists of titles and abstracts extracted from 221
16 scientific journals publishing articles in the IR and related fields (computer sciences). 222

# ARTICLE IN PRESS

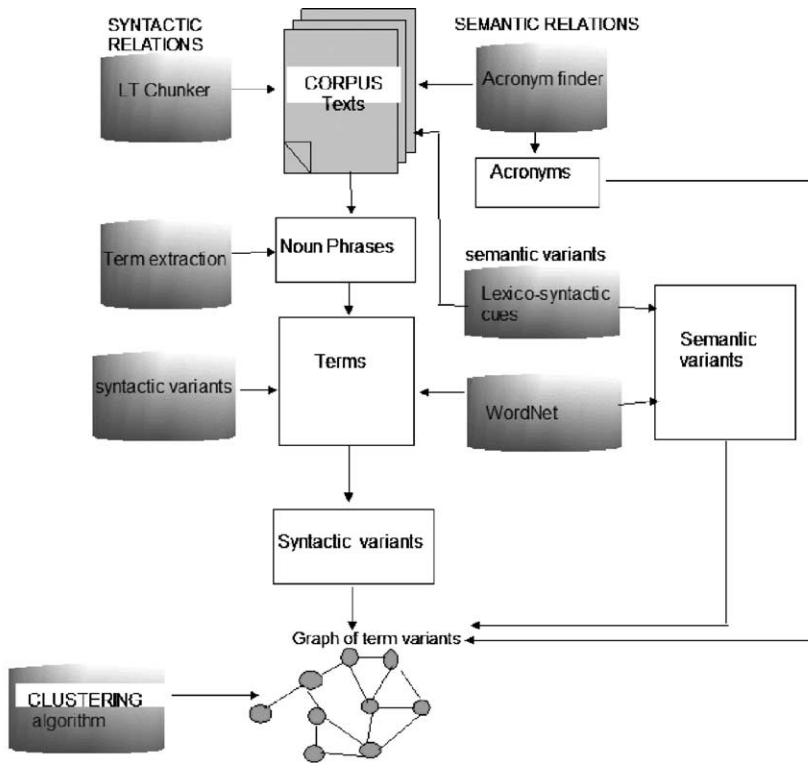*F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx* 7

Fig. 1. TermWatch's architecture.

The aim is to map out the research topics addressed in these abstracts over a period of 8 years (1997–2003). The 3,355 titles and abstracts records were thus extracted from the PASCAL multidisciplinary database maintained by the French Institute for Scientific Information (INIST) (http://www.inist.fr). These make up roughly 455,000 words. Although we worked on abstracts rather than on full texts, they were the authors' own texts and shorter texts like abstract are known to be more information dense than full texts. Thus, abstracts represent in our view, adequate surrogates of the full papers. The table below shows the ranking of the journals according to number of bibliographic records. Column 1 is the journal rank, column 2 gives the number of bibliographic records per journal, column 3 the proportion in the entire corpus, column 4 the cumulative, and the last column the journal name (Table 1).

As we can see, the journal that contributed most to the text collection is *Information Sciences*, followed closely by *JASIST*. This is the ranking obtained when using quantitative indicator (number of published papers) as the sole measure of journal representativity vis-à-vis a scientific field. We now look at the fine-grained content analysis of the journals contents as mapped out by TermWatch. We will map the clusters obtained onto the journals to see if the same ranking by productivity is maintained. Table 2 below gives some clustering details obtained from this collection. Because clustering is an iterative process, the user can choose the level of iteration at which to stop the process depending on cluster

**ARTICLE IN PRESS**

8 F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx

| | | | | | |
|---|---|---|---|---|---|
| t1.1 | Table 1 | | | | |
| t1.2 | Collection of 16 journals from the IR and related fields | | | | |
| t1.3 | 1 | 831 | 25% | 831 | 25% | *Information Sciences* |
| | 2 | 688 | 21% | 1519 | 45% | *Journal of the American Society for Information Science and* |
| t1.4 | | | | | | *Technology* |
| t1.5 | 3 | 283 | 8% | 1802 | 54% | *Information Processing and Management* |
| t1.6 | 4 | 272 | 8% | 2074 | 62% | *Journal of Information Science* |
| t1.7 | 5 | 267 | 8% | 2341 | 70% | *Information Systems Management* |
| t1.8 | 6 | 175 | 5% | 2516 | 70% | *Journal of Documentation* |
| t1.9 | 7 | 176 | 5% | 2692 | 80% | *Information Systems* |
| t1.10 | 8 | 116 | 3% | 2808 | 84% | *Information Systems Security* |
| t1.11 | 9 | 108 | 3% | 2916 | 87% | *Library and Information Science Research* |
| t1.12 | 10 | 108 | 3% | 3024 | 90% | *Online Information Review* |
| t1.13 | 11 | 87 | 3% | 3111 | 93% | *Journal of Internet Cataloging* |
| t1.14 | 12 | 70 | 2% | 3181 | 95% | *Information Retrieval and Library Automation* |
| t1.15 | 13 | 67 | 2% | 3248 | 97% | *Knowledge Organization* |
| t1.16 | 14 | 44 | 1% | 3292 | 98% | *Journal of Information Science and Engineering* |
| t1.17 | 15 | 34 | 1% | 3326 | 99% | *International Forum on Information and Documentation* |
| t1.18 | 16 | 29 | 1% | 3355 | 100% | *Information Retrieval* |
| t1.19 | | 3355 | 100% | | | |

granularity (size). In this experiment, we chose the results of the second iteration because 242
classes and their layout seemed meaningful. The 674 classes of variable sizes were thus 243
obtained containing a total of 5632 terms. 244

This clustering output is an improved version of the one already carried out on the 245
IRcorpus and published in Ibekwe-SanJuan and SanJuan [10,11]. In this experiment, we 246
refined the definitions of the variation relations (cf. Section 3.2) and changed their roles 247
during clustering. 248

### 4.1. Graphic display of collection thematics 249

We show here below the output of the system viewed through a graphic display package, 250
Aisee. Aisee interprets clusters built by TermWatch and aligns them according to their 251
centrality (number of outgoing links). Thus, in the cropped image below (we only show the 252

| | | |
|---|---|---|
| t2.1 | Table 2 | |
| t2.2 | Details of the clustering | |
| t2.3 | Clusters obtained from the IR corpus | |
| t2.4 | Number of iterations | 1 |
| t2.5 | Number of components | 1595 |
| t2.6 | Number of clusters | 674 |
| t2.7 | Size biggest cluster | 135 |
| t2.8 | Size smallest cluster | 4 |
| t2.9 | Total terms in clusters | 5632 |

ARTICLE IN PRESS

*F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx*                9

central part of the image), the most central topic is "information retrieval" this is not       253
surprising because the collection is built on this topic.       254

However, the selection of the texts was not based on keywords but on journal titles.       255
There was therefore no guarantee that "information retrieval" will be found as the most       256
active term with many linguistic relations (variants) in the corpus. It could have been       257
considered as a "meta term" by authors and as such, not used in their abstracts because       258
these journals were more or less about information retrieval. Surprisingly, this turned out       259
not to be the case. The fact that this term actually appeared with a lot of variants shows that       260
researchers actually use the macroterm together with more specific qualifiers to refer to       261
their works or to applications of their studies. Unfolding a cluster shows the most active       262
term variants. Unfolding the "information retrieval" cluster showed that it dealt with objects       263
and methods of information retrieval systems, hence the presence of variants like "content-       264
based image retrieval systems, NLP information retrieval systems, bibliographic retrieval       265
systems, modern text retrieval systems, natural language information retrieval systems,       266
online information retrieval systems…" Thus, the label is the most generic term while the       267
cluster contents point to more specific and current research concerns.       268

Surrounding this most central clusters are other clusters like "semantic similarity       269
measure," which deals with different similarity measures used in information retrieval like       270
Cosine, Jaccard, angle-based similarity measures, collocation-based similarity measure,       271
distance similarity measure, etc. The cluster labeled "vector space" refers to research on       272
vector space model of information retrieval. The cluster "wide web sites" deals with       273
different types of Web sites (academic, commercial, etc.). "Natural language" cluster       274
portrays research on natural language query processing. The cluster "online information       275
sources" concerns studies dealing with different online resources as shown by variants like       276
"electronic consumer health information, Web information sources, commercially produced       277
online information sources, distributed information sources, sources of bibliographic       278
information…" "Online catalog" contained variants like "Web on-line catalog, operational       279
online catalogs, commercially available Web browsers, needs of online catalog users, on-       280
line catalog searching, next generation of online catalogs, next generation of retrieval       281
systems" showing clearly the theme reflected by the cluster (Fig. 2).       282

The topographic layout of clusters offered by TermWatch is useful for grasping rapidly       283
the contents of a large text collection. This is particularly important for science and       284
technology watch, that is, understanding the interactions between domain topics and       285
following their evolution through time stamps [10] but also for query refinement.       286

## 4.2. Ranking journals by term and cluster representativity       287

The focus of this paper is to determine how the clusters of domain topics mapped       288
TermWatch can assist library collection management. Hence, we will seek to ascertain       289
how the 16 journals are distributed across the 674 clusters by assigning journals to       290
clusters in which they have the highest number of terms. Assuming that the most       291
productive journals as shown in Table 1 will also be the most productive in terms of       292
"terms" and "variation relations," we should obtain the same ranking with respect to a       293

## ARTICLE IN PRESS

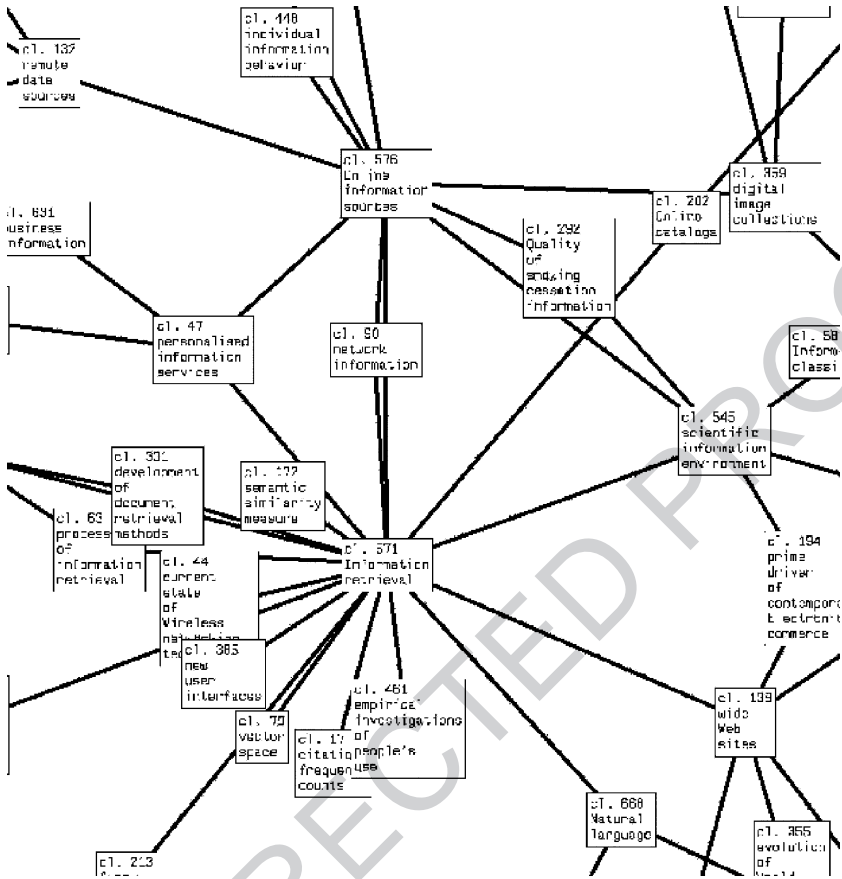10                    F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx

Fig. 2. Part of the thematic map produced by TermWatch on the IRcorpus.

given cluster. We ranked the journals by number of terms they contained, then by number 294
of clusters (Table 3). 295

As we can observe from the table below, journal representativity by number of terms is 296
roughly correlated with their representativity by number of clusters (except for two positions, 297
5th and 7th). However, comparison with the journal ranking by number of articles (Table 1) 298
shows some discrepancies. *JASIST* turns out to be the most productive in terms of domain terms 299
and variants whereas it was 2nd by number of articles. Conversely, *Information Sciences* now 300
comes 2nd by representativity in clusters. *Information Processing and Management*, *Journal of* 301
*Information Science*, *Journal of Documentation*, and *Journal of Information Science and* 302
*Engineering* maintain their respective positions in the two rankings. On the other hand, 303
"information systems, library and information science research, online information review," 304
and "knowledge organization" gain two places by arriving at the 5th, 7th, 8th, and 11th 305
positions, respectively, by number of term variants. *International Forum on Information and* 306
*Documentation* and *Information Retrieval* also gain three places by arriving at the 12th and 307
13th positions, respectively. *Information Systems Management*, *Information Systems Security*, 308

# ARTICLE IN PRESS

*F. Ibekwe-SanJuan / Libr. Coll. Acq. & Tech. Serv. xx (2005) xxx–xxx* 11

t3.1 Table 3
t3.2 Journal ranking by number of terms in clusters

| | Rank | Journal | Number of clusters | Number of terms |
|---|------|---------|--------------------|-----------------|
| t3.3 | | | | |
| t3.4 | 1 | *Journal of the American Society for Information Science and Technology* | 468 | 3616 |
| t3.5 | 2 | *Information Sciences* | 382 | 3115 |
| t3.6 | 3 | *Information Processing and Management* | 304 | 1582 |
| t3.7 | 4 | *Journal of Information Science* | 252 | 1067 |
| t3.8 | 5 | *Information Systems* | 219 | 997 |
| t3.9 | 6 | *Journal of Documentation* | 249 | 899 |
| t3.10 | 7 | *Library and Information Science Research* | 140 | 517 |
| t3.11 | 8 | *Online Information Review* | 153 | 488 |
| t3.12 | 9 | *Information Systems Management* | 121 | 438 |
| t3.13 | 10 | *Journal of Internet Cataloging* | 90 | 422 |
| t3.14 | 11 | *Knowledge Organization* | 85 | 227 |
| t3.15 | 12 | *International Forum on Information and Documentation* | 75 | 164 |
| t3.16 | 13 | *Information Retrieval* | 69 | 161 |
| t3.17 | 14 | *Journal of Information Science and Engineering* | 58 | 122 |
| t3.18 | 15 | *Information Systems Security* | 29 | 83 |
| t3.19 | 16 | *Information Retrieval and Library Automation* | 25 | 45 |

and *Information Retrieval and Library Automation* descend by four, five, and four places, 309
respectively. On the whole, seven out of the 16 journals showed consistency in the two rankings 310
while nine journals showed notable differences. 311

## 5. Conclusion 312

We have presented in this paper, an alternative to the journal collection management 313
problem. This could be through thematic mapping using linguistic and data analysis 314
techniques. The proposed approach, embodied in the TermWatch system enables a librarian to 315
grasp more readily the contents of a collection of journal through an in-depth analysis of their 316
texts. The resulting maps can be used for positioning research topics vis-à-vis one another and 317
contribute also to answering specific search needs of certain categories of users. The journal 318
ranking by thematic content also portrays differences with ranking by pure numerical factor 319
(i.e., journal productivity). This finding suggests that while some journals may publish a 320
considerable amount of papers in a given field, this number may not necessarily be correlated 321
with density of domain terms. 322

## 6. Uncited references 323

[6] 324
[9] 325
[12] 326
[14] 327

# References

328
329

[1] D. Bourigault, C. Jacquemin, & M. -C. L'Homme, (Eds.). (2001). *Recent advances in computational terminology, vol. 2*. Amsterdam: John Benjamins.

[2] Baeza-Yates, & Ribeiro-Neto, B. (1999). Query operations. *Modern information retrieval* (pp. 117–139). ACM Press.

[3] Callon, M., Law, J., & Rip, A. (1986). *Mapping the dynamics of science and technology.* Basingstoke: MacMillian Press.

[4] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, R., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science, 41*6), 391–407.

[5] Fellbaum, C. (1998). *WordNet, an electronic lexical database*. MIT Press.

[6] Giles, C. L. *Citation index: Journal impact factors*. Retrieved 23/02/2005 from http://www.neci.nj.nec.com/homepages/giles/html/cites.html

[7] Jardine, N., & Van Rijsbergen, C. J. (1971). The use of hierarchic clustering in information retrieval. *Information Storage and Retrieval, 7,* 217–240.

[8] Hearst, M. (1999). The use of categories and clusters for organizing retrieval results. In T. Strzalkowski (Ed.), *Natural language information retrieval. Text, Speech and Language Technology, vol. 7* (pp. 333–374). Kluwer Academic Press.

[9] Hearst, M. A. (1992). Automatic acquisition of hyponyms from large text corpora. *Proceedings of the International conference on Computational Linguistics (COLING'92), Nantes,* 539–545.

[10] Ibekwe-SanJuan, F., & SanJuan, E. (2004). Mining textual data through term variant clustering: The Termwatch system. *Proceedings of the International conference on "Recherche d'Information assistée par ordinateur. Avignon"* (pp. 487–503).

[11] Ibekwe-SanJuan, F., & SanJuan, E. (2004). Mining for knowledge chunks in a terminology network. *8th International ISKO conference, University College London, 13–16 July 2004* (pp. 41–47).

[12] Ibekwe-SanJuan, F., & SanJuan, E. (2002). *From term variants to research topics*. *International Journal on Knowledge Organization (ISKO), special issue on Human Language Technology, 29*(3/4), 181–197.

[13] Ibekwe-SanJuan, F. (1998). A linguistic and mathematical method for mapping thematic trends from texts. *Proceedings of the 13th European Conference on Artificial Intelligence, Brighton UK, 23–28 August* (pp. 170–174).

[14] Jacquemin, C., & Bourigault, D. (2003). *Term extraction and automatic indexing*. In R. Mitkov (Ed.), *Handbook of computational linguistics* (pp. 599–615). Oxford University Press.

[15] Morin, E., & Jacquemin, C. (2004). Automatic acquisition and expansion of hypernym links. *Computer and the Humanities, 38*(4), 343–362.

[16] Salton, G. (1968). *Automatic information organization and retrieval*. New York: McGraw-Hill (18 pp.).

[17] Small, H. (1999). Visualizing science by citation mapping. *Journal of the American Society for Information Science, 50*(9), 799–813.

[18] Schiffrin, R., & Börner, K. (2004). *Mapping knowledge domains*. *Publication of the National Academy of Science (PNAS), 101*(1), 5183–5185.

[19] Sparck Jones, K. (1971). *Automatic keyword classification for information retrieval*. London: Butterworths.

[20] White, H. D., & Mccain, K. W. (1989). *Bibliometrics*. In M. E. Williams (Ed.), *Annual review of information science and technology* (pp. 119–186). New York: Elsevier Science.

[21] Yeung, K. Y., Haynor, H., & Ruzzo, W. L. (2001). Validating clustering for gene expression data. *Bioinformatics, 17,* 309–318.

[22] Zamir, O., & Etzioni, O. (1998). Web document clustering, a feasibility demonstration. *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval* (pp. 46–54).

330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374