

Recommenders in a personalized, collaborative digital library environment

Henri Avancini · Leonardo Candela ·
Umberto Straccia

Received: 22 December 2004 / Revised: 6 October 2005 /
Accepted: 12 October 2005
© Springer Science + Business Media, LLC 2006

Abstract We envisage an information source not only as an information resource where users may submit queries to satisfy their daily information need, but also as a collaborative working and meeting space of people sharing common interests. Indeed, we will present a highly personalized environment where not only users may organize (and search into) the information space according to their individual taste and use, but which provides advanced features of collaborative work among the users. It is up to the system to discover interesting properties about the users' interests, relationships between users and user communities and to make recommendations based on preference patterns of the users, which is the main topic of this paper.

Keywords Digital library · Recommenders · Collaborative work

1 Introduction

It is widely recognized that *Digital Libraries* (DLs) (Fox & Marchionini, 2001), will play an important role in the next future not merely in terms of the 'controlled' digital information they allow access to, but in terms of the *services* they provide to the information society at large. Informally, DLs can be defined as consisting of collections of information (usually, heterogeneous in content and format), which have associated services delivered to users and user communities using a variety of

H. Avancini · L. Candela (✉) · U. Straccia
Istituto di Scienza e Tecnologie dell'Informazione – C.N.R. – Via G. Moruzzi,
1 I-56124 Pisa, Italy
e-mail: leonardo.candela@isti.cnr.it

H. Avancini
e-mail: avancini@isti.cnr.it

U. Straccia
e-mail: straccia@isti.cnr.it

technologies. The services offered on such information can be various, ranging from content operations to rights management, and can be offered to *individuals* as well as to *user communities*.

Even though DLs have evolved rapidly over the past decade, typically, they still are limited to provide a search facility to the digital society at large. Indeed, they are oriented towards a generic user, as they answer queries crudely rather than learn the long-term requirements of a specific user. In practice, users use the same information resource over and over and would benefit from automatization: the time consuming effort that the user put in searching documents and possibly downloading them from the DL is often forgotten and lost. This requires a repetition of the manual labour in searching and browsing to find the documents just like the first time.

As DLs will become more commonplace, there is a need for them to move from being *passive* with little adaptation to their users, to being more *proactive* and *personalized* in offering and tailoring information for individual users. Nowadays, in several DLs some personalization functionalities are provided. Mainly they fall into the category of personalized *alerting services* (Bollacker, Lawrence, & Giles, 1999; Faensen, Faulstich, Schweppe, Hinze, & Steidinger, 2001; Fernandez, Sanchez, & Garcia, 2000; Moukas, 1996; Rocha, 1999). These services notify a user (usually, by sending an e-mail), with a list of references to newly available documents deemed as relevant to some of the (manually) user specified topics of interests. Typically, such alerting services are provided by scientific journal publishers. Some other DLs, in addition, support the users in being able to organize the information space they are accessing, according to *their own subjective perspective* (Fernandez et al., 2000). This is important as not necessarily all the information provided by a DL may be of interest to a user, but just some ‘slices’ of it. Users and communities of users might well profit from being able to organize the information space in a personalized fashion, both in terms of restricting the information space in which to search into, as well as in terms of organizing it, not necessarily in the way the DL manager thought would be well suited for anyone.

Seldom, (Di Giacomo, Mahoney, Bollen, Monroy-Hernandez, & Meraz, 2001) is an exception, DLs also can be considered as a *collaborative meeting place* of people sharing common interests. Indeed, our vision is that DLs may be viewed as *common working places* where users may become aware of each other, open communication channels, and exchange information and knowledge with each other or with experts. In fact, usually users and/or communities access a DL in search of some information. This means that it is quite possible that users may have overlapping interests if the information available in a DL matches their expectations, backgrounds, or motivations. Such users might well profit from each other’s knowledge and expertise by sharing opinions or experiences or offering advice. Some users might enter into long-term relationships and eventually evolve into a community if only they were to become aware of each other. Hence, we are moving from services supporting an individual user towards services supporting *groups* of users as well.

A major service in a personalized and collaborative DL is the *recommendation* of items to a user based not only on preference patterns of the user itself, but also on those of other users. The use of opinions and knowledge of other users to predict the relevance value of items to be recommended to each user in a community is known as *Collaborative* or *Social Filtering* (Billsus & Pazzani, 1998; Breese, Heckerman, & Kadie, 1998; Goldberg, Nichols, Oki, & Terry, 1992; Good et al., 1999; Herlocker,

Konstan, Borchers, & Riedl, 1999; Konstan et al., 1997; Resnick, Iacovou, Suchak, Bergstorm, & Riedl, 1994). These methods are built on the assumption that a good way to find interesting content is to find other users who have similar interests, and then recommend items that those similar users like.

In this paper we first present the CYCLADES system and stress its ‘personalization’ and alerting features.¹ A major distinction of CYCLADES is that it is indeed a DL environment supporting collaboration and personalization at various levels. Users and communities may search, share and organize their information space according to their own view. The system is able to generate recommendations of various types based on user and community profiles. We then extensively present the recommendation algorithms used in CYCLADES. Worth noting is that the system not only recommend to users *documents* deemed as relevant, like in alerting services, but also recommends *users* and *communities* with similar interests, and *collections* where to search.

The outline of the paper is as follows. In the next section we present the main features of CYCLADES; Section 3 introduces the recommendation algorithms adopted within the final version of CYCLADES; Section 4 reports our experimental results discussing the effectiveness of the proposed algorithms; Section 5 concludes.

This work concludes a series of works on the subject (Avancini & Straccia, 2004a,b; Candela & Straccia, 2004; Renda & Straccia, 2002). This paper summaries these works, but all tests and the collection recommendation algorithm are novel.

2 A brief overview of CYCLADES

The objective of CYCLADES is to provide an integrated environment for users and groups of users (communities) that want to use, in a highly personalized and flexible way, ‘open archives,’ i.e., electronic archives of documents compliant with the Open Archives Initiative² (OAI) standard. Informally, the OAI is an agreement among several digital archives providers in order to provide some minimal level of interoperability between them. In particular, the OAI defines an easy-to-implement gathering protocol over HTTP, which gives *data providers* (the individual archives) the possibility to make the metadata of the documents stored in their archives externally available. Indeed, the agreement specifies that each document of an archive should possess a *metadata* record describing the document properties and content. In particular, the default format of the metadata records should be DublinCore³ and any other additional metadata format is supported also. The metadata record consists of several attributes describing author, title, abstract, etc. of the document. The protocol allows then to gather these metadata records (in place of the real documents). A link to the ‘real’ document also is present if the document is accessible. A metadata record may be understood as a statement of existence and short description of a

¹CYCLADES has been partially funded by the European Community. See, <http://www.ercim.org/cyclades>.

²Open Archive Initiative web site. <http://www.openarchives.org>.

³DublinCore web site <http://dublincore.org>.

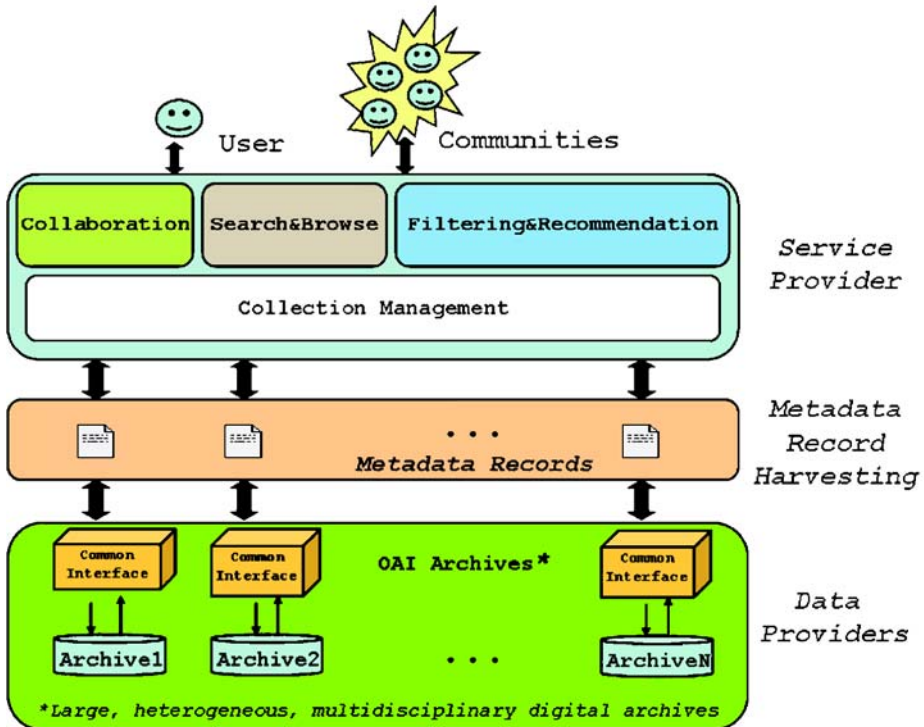


Fig. 1 Logical view of CYCLADES functionality

document. The document may be then accessible to a user according to the access policies of the archive, which owns the document. To date, there is a wide range of archives available (more than one hundred registered archives) in terms of their content, forming a quite *heterogeneous and multidisciplinary information space* (e.g., covering the fields of biology, mathematics, computer science, physics, etc.).

The availability of the metadata records from the OAI compliant archives makes it possible for *service providers* to build higher-level functionality. In this sense, CYCLADES allows the access to the metadata provided by these archives, as it gathers these records and through them provides access to the referenced complete documents (if they exist and their access is allowed). On top of them, CYCLADES acts as an OAI service provider (see Fig. 1).

The architecture of the CYCLADES system is depicted in Fig. 2. Each box is a service accessible via the Web. The exchange of data among the services is based on XML-RPC.⁴ XML-RPC is a specification and a set of implementations that allow software running on disparate operating systems, running in different environments to make procedure calls over the Internet. It is remote procedure calling using HTTP as the transport and XML as the encoding. XML-RPC is designed to be as simple as possible, while allowing complex data structures to be transmitted, processed and returned.

⁴XML-RPC is a popular protocol that uses XML over HTTP to implement remote procedure calls. See <http://www.xmlrpc.com>.

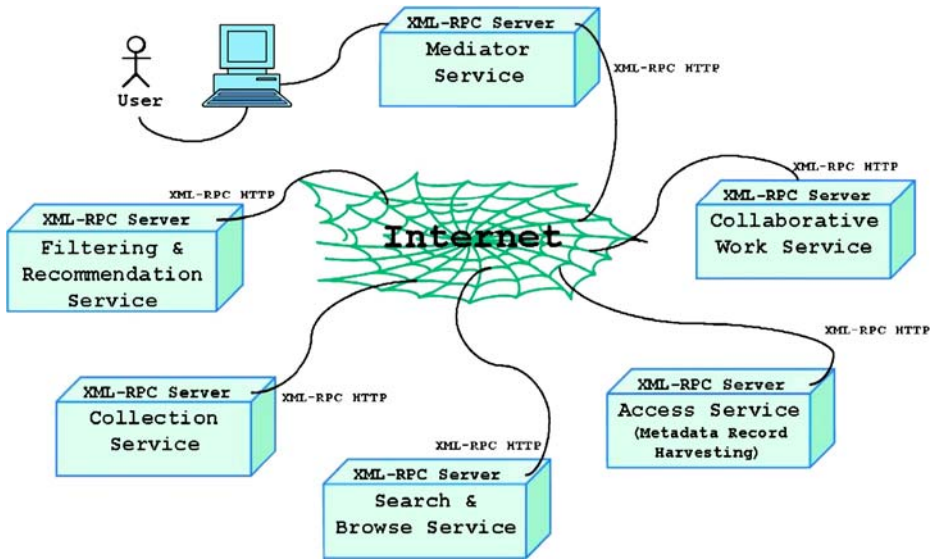


Fig. 2 CYCLADES architecture

The CYCLADES system, accessible through Web browsers, provides the user with different environments, according to the actions the user wants to perform. Roughly, the Mediator Service, the main entry point to the CYCLADES system, acts as a registry for the other services, checks if a user is entitled to use the system, and ensures that the other services are only called after proper authentication. The Access Service is in charge of interfacing with the underlying metadata archives. In CYCLADES only archives adhering to the OAI specification will be accounted for; however, the system is extensible to other kinds of archives by modifying the Access Service only. A user may also ask CYCLADES to include newly OAI compliant archives as well (see Fig. 3).

The Access Service gathers the OAI records from archives and populates the system with them (currently, CYCLADES gathers metadata from more than 100 archives). The Collaborative Work Service provides an environment for managing metadata records and collaborative activities. The Search and Browse Service supports the activity of searching records from the various collections, while the Collection Service provides collection management functionality (i.e., their definition, creation, and update). Finally, the Filtering and Recommendation Service provides filtering and recommendation functionality. In summary, CYCLADES provides functionality for:

1. Advanced search in *large, heterogeneous, multidisciplinary digital archives*;
2. Collaborative work;
3. Information filtering;
4. Various forms of information recommendation; and
5. Management of records grouped into so-called *collections* (see later on).⁵

⁵In CYCLADES, there is one collection of metadata records, for each OAI archive. But, additionally, users may define their own sets.

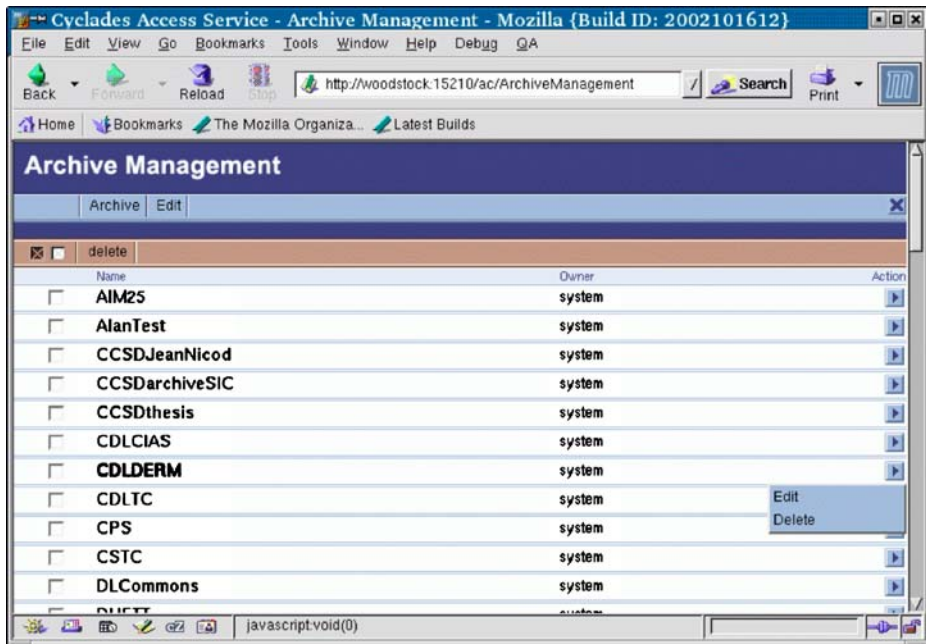


Fig. 3 User interface: OAI archive list within CYCLADES

We now describe the functionality in more detail. The main principle underlying CYCLADES is the *folder paradigm* (see Fig. 4). That is, users and communities of users may organize the information space into their own folder hierarchy, as e.g., may be done with directories in operating systems, bookmark folders in Web browsers and folders in e-mail programs. Each folder typically corresponds to one user related subject (or discipline, or field), so that it may be viewed as a thematic and usually semantically related repository of data items.

The CYCLADES functionality are available in several system environments. Each functionality is supported by a web service, which we describe next.

2.1 Collaborative work environment

This is an extension of the BSCW environment (Basic Support for Collaborative Work) (Bentley et al., 1997) and provides the folder-based environment. The collaborative work environment allows the management of various data items, such as metadata records, queries, collections, external documents (i.e., documents that can be uploaded by the user to the folder), ratings and annotations.

There are two types of folders: (a) *private folders*, i.e., a folder owned by one user only. This kind of folder can only be accessed and manipulated by its owner. They are invisible to other users; and (b) *community folders*, which can be accessed and manipulated by all members of the community that owns the folder. Community folders are used to share data items with other users and to build up a common folder hierarchy. Community folders may also contain *discussion forums* where notes may be exchanged in threaded discussions (similar to news groups). For instance,

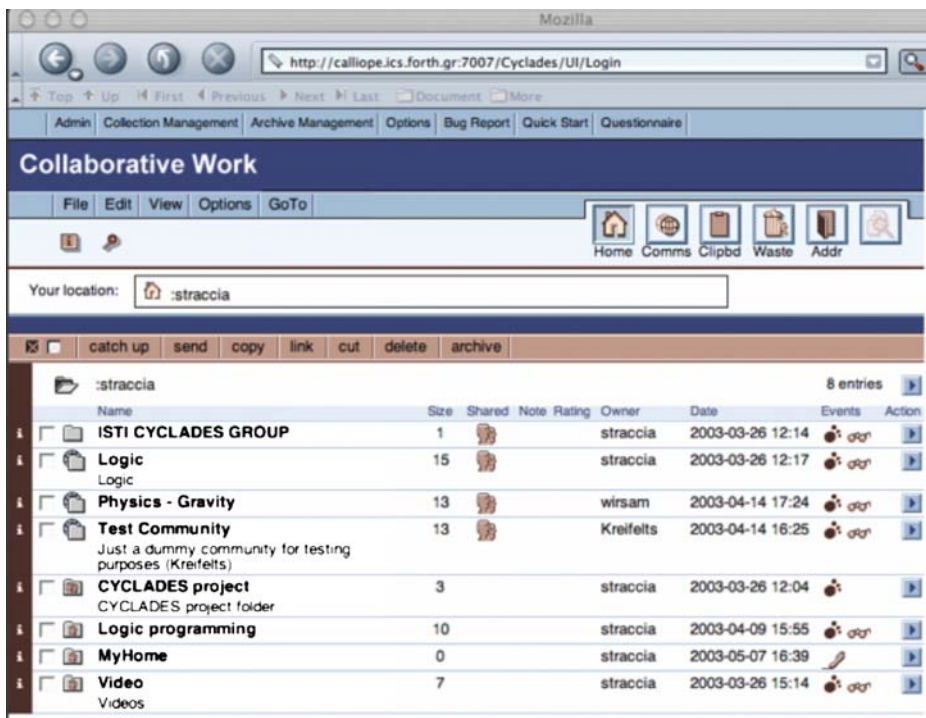


Fig. 4 User interface: a user home folder

Fig. 4 shows the home (top level) folder of a user. It contains several sub-folders. Among them, there are some (shared) folders belonging to communities (created by someone) to which the user joined, like the ‘Physics-Gravity’ folder (community), while others are private folders and have been created directly by the user, e.g., the ‘Logic Programming’ folder. These folders contain community or user collected OAI records relevant to some topics (e.g., gravity and logic programming, respectively). Figure 5 shows the content of a specific user folder, in our case the ‘Physics-Gravity’ folder of the community of physicists. In order not to lose shared activity in the collaborative DL environment, mutual awareness can be supported through event icons displayed in the environment. Activity reports that are daily received by email are also possible. Users can view the list of all existing communities and can join a community directly if the community allows this policy, or contact the community administrator in order to be invited to the community. In the collaborative work environment, the access policies can be set-up, as well as the notification (alerting) modalities.

2.2 Search and browse environment

It supports the activity of searching records in the various metadata record collections accessible from within CYCLADES as well as to search into the shared folders or private folders a user owns. Users can issue a query and are allowed to store selected records

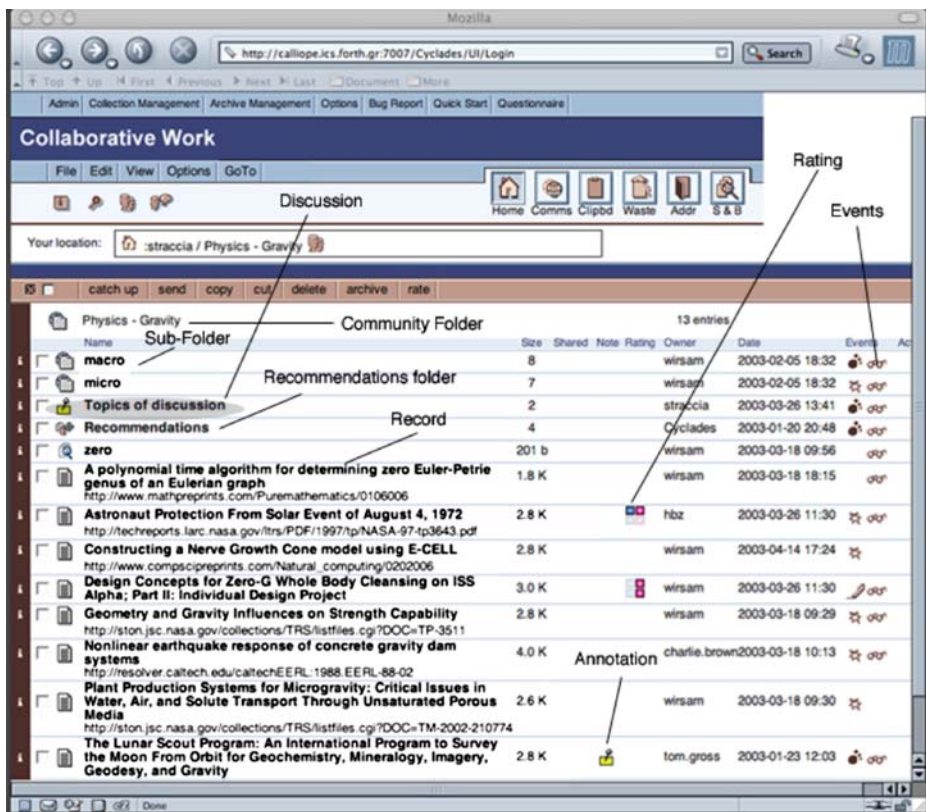


Fig. 5 User interface: folder content

within their folders and community folders they have access to (see Fig. 6 for a query and Fig. 7 for a query result list).

Essentially, three types of search are supported:

1. *Ad-hoc search*, where a user specifies a query and the system looks for relevant records within a specified collection;
2. *Filtered search* is like the usual ad-hoc search, except that the user specifies, additionally to a query (e.g., 'zero'), also a target folder (e.g., 'Physics-Gravity'). The goal of the system consists then to find documents not only relevant to the query, but also relevant to the topic of the target folder (in our example, the request is something like 'find records about zero gravity'); and
3. *What's new, on-demand*, where the user specifies a target folder, without specifying a query, and the goal of the system consists of finding all records, relevant to the target folder, which have become available to the system since the last time the user asked for this request. This corresponds roughly to the functionality provided by alerting services, except that the profile is build implicitly from the folder content, and that records are delivered to the user on-demand.

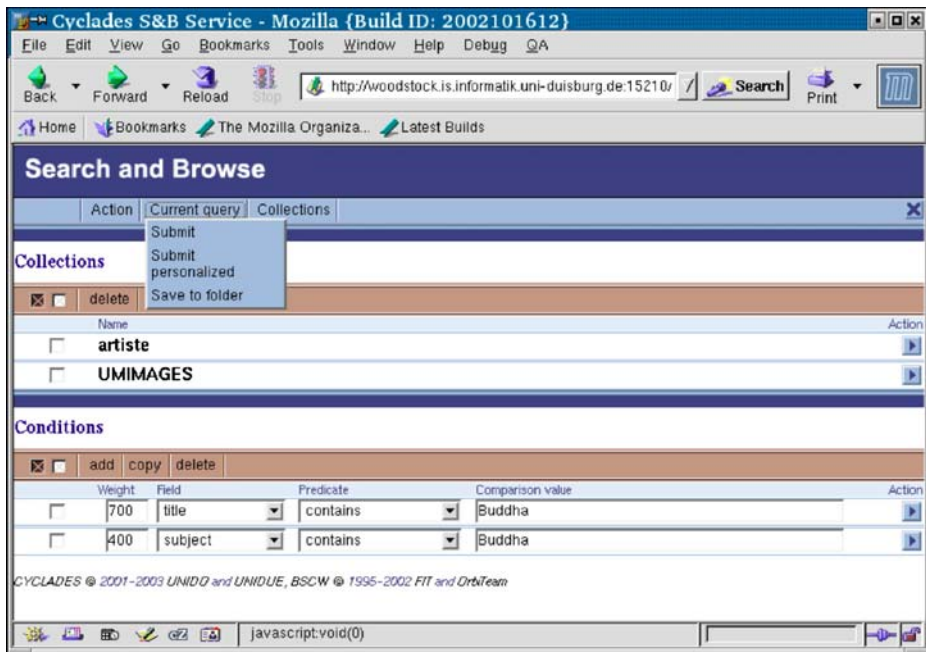


Fig. 6 User interface: search metadata records

2.3 Filtering and recommendation environment

It supports the personalized search (i.e., “filtered search” and “what’s new, on-demand”) and provides the recommendations functionalities. All recommendations are specific to a given user folder (topic of interest), i.e., they have always to be understood in the context not of the general interests of the user, but of the specific interests (topic) of the user represented by a folder. A user may get recommendations of

- *Metadata records*, suggesting to the user to access relevant documents;
- *Collections*, suggesting to the user to search within a relevant information space;
- *Users*, suggesting to the user to enter into relationship with a user or give a look to the publicly available documents of the recommended user; and
- *Communities*, suggesting to the user to join the community.

Recommendations are issued to users based on user and/or community profiles. For instance, Fig. 8 shows the recommendations related to the ‘Physics-Gravity’ folder, deemed by the system as relevant to this folder.

2.4 Collection management

It manages collections (i.e., their definition, creation, and update) of metadata records. Its aim is to allow a dynamic partitioning of the information space according to the users’ interests. Usually, a collection is meant to reflect a topic of interest of

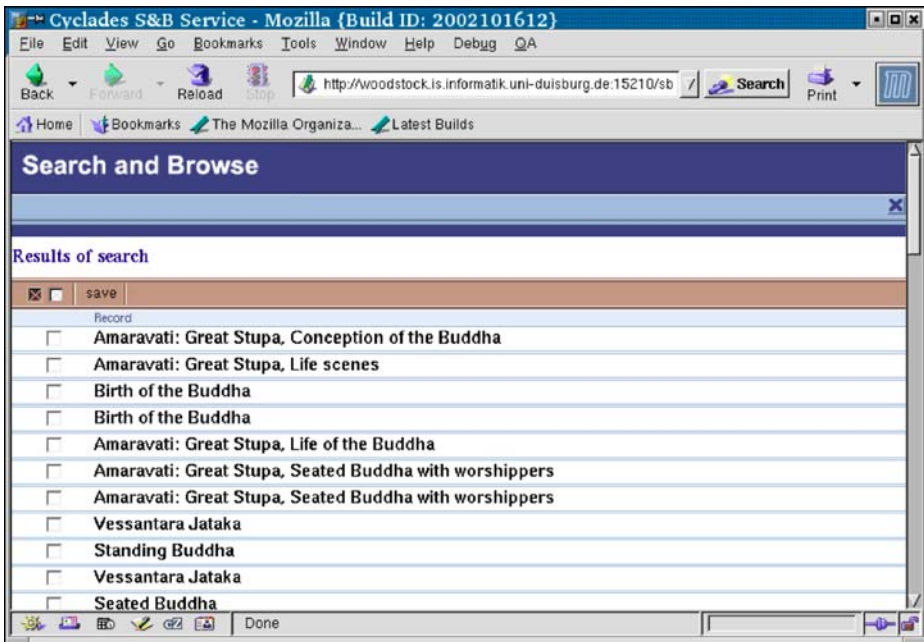


Fig. 7 User interface: result list of metadata records

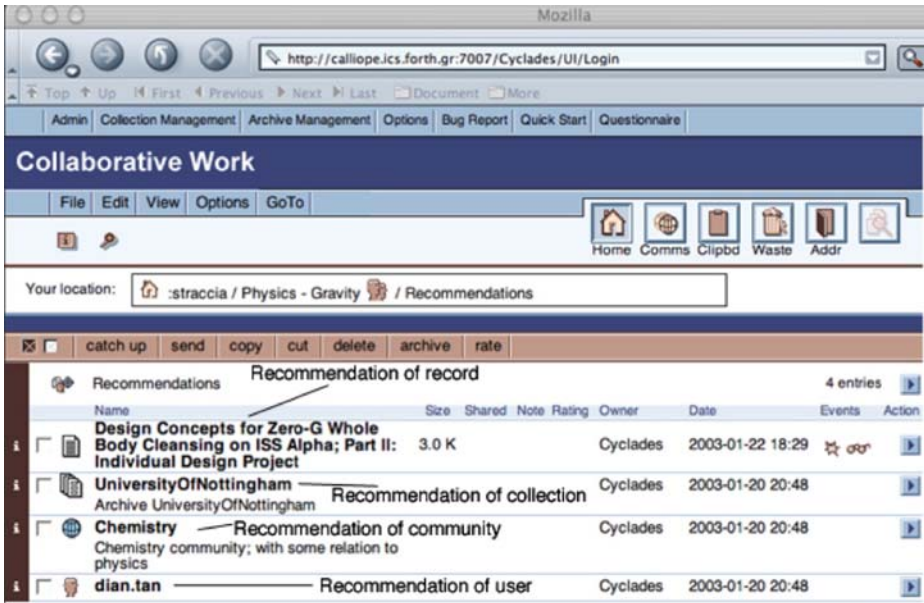


Fig. 8 User interface: folder content and recommendations

a user or a community, e.g., the collection of records about ‘Information Retrieval.’ Informally, a collection specification is the definition of a non-materialized view over the information space. Thus, the collection specification is expressed via a query expression. It is up to the system to automatically determine the ‘physical’ collections in which to search for relevant records (this is accomplished by means of a technique called *automated source selection* see, e.g., Callan, 2000; Fuhr, 1999; Meng et al., 1998; Rasolofo, Abbaci, & Savoy, 2001).

3 Recommendation algorithms

As pointed out, filtering and recommendation play an important role in making CYCLADES a personalized and collaborative environment. In this section, we detail all the algorithms used: namely, for (a) metadata record recommendation; (b) user recommendation; and (c) collection recommendation. We also report experimental results of their effectiveness. We point out that *the recommendation of communities, i.e., the suggestion to join a community is realized as for the case of collection recommendation*. Indeed, from a topic point of view, we identify a community with its community folder, i.e., a set of metadata records. That is, the records belonging to the community folder describe the topics the community is interested in. As such, a community can be seen as a collection of records and, thus, the recommendation of a community can be implemented like the recommendation of a collection. Therefore, we do not address it further.

Before we present the recommendation algorithms, we need to introduce some notation. Then, for each recommendation type (record, user, collection) we provide the objective and the algorithm.

3.1 Preliminaries

In the following, we denote with u_k a user, with F_i a folder (of the users) and with d_j a metadata record.

We consider a metadata record as an unstructured piece of plain text (of course, more sophisticated algorithms can be devised by taking into account the metadata structure). Metadata records belong to folders and each user may rate a document within a folder he has access to. A record may belong to multiple folders. With r_{ijk} we indicate the rating value given by a user u_k to record d_j , which is stored in folder F_i . We further assume that whenever a data item d_j belongs to a folder F_i of a user u_k , an *implicit* default rating r_{ijk} is assigned. Indeed, a record belonging to a folder of a user is an implicit indicator of being the record relevant to the user folder.

Given a folder F_i and a record d_j belonging F_i , we compute an average rating from all the ratings given to record d_j in folder F_i , i.e., we average the ratings r_{ijk} relative to the same record-folder pair (i, j) and indicate it as $r_{ij} = \text{mean}_{k \geq 1} \{r_{ijk}\}$ (see matrix c in Table 1). The rating r_{ij} indicates the average rating given by the users to record d_j in folder F_i . This allows us to re-conduct our three dimensional recommendation setting (user-folder-record) to an easier to manage two dimensional setting (folder-record) to which we apply our recommendation algorithms. Of course, in place of the mean, other measures may be applied to compute r_{ij} as well.

Table 1 (a) The records matrix, (b) the folder profile matrix, (c) the folder-record rating matrix, and (d) the user profile matrix

(a)					
	t_1	...	t_k	...	t_m
d_1	w_{11}	...	w_{1k}	...	w_{1m}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
d_j	w_{j1}	...	w_{jk}	...	w_{jm}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
d_n	w_{n1}	...	w_{nk}	...	w_{nm}

(b)					
	t_1	...	t_k	...	t_m
f_1	w_{11}	...	w_{1k}	...	w_{1m}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
f_i	w_{i1}	...	w_{ik}	...	w_{im}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
f_v	w_{v1}	...	w_{vk}	...	w_{vm}

(c)					
	d_1	...	d_j	...	d_n
F_1	r_{11}	...	r_{1j}	...	r_{1n}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
F_i	r_{i1}	...	r_{ij}	...	r_{in}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
F_v	r_{v1}	...	r_{vj}	...	r_{vn}

(d)					
	t_1	...	t_k	...	t_m
p_1	w_{11}	...	w_{1k}	...	w_{1m}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
p_u	w_{u1}	...	w_{uk}	...	w_{um}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
p_h	w_{h1}	...	w_{hk}	...	w_{hm}

All records in the folders are indexed according to the well-known vector space model (Salton & McGill, 1983). Therefore, with $d_j = \langle w_{j1}, \dots, w_{jm} \rangle$ we indicate the indexed representation of a record, where $0 \leq w_{jk} \leq 1$ is the ‘weight’ of term (keyword) t_k in the record d_j (see matrix a in Table 1).

For each folder, we compute its *folder profile*, which is a machine representation of what a folder is *about*. For a given folder F_i , the folder profile f_i is computed as the *centroid*, or average, of the records belonging to F_i , i.e., $f_i = (1/|F_i|) \sum_{d_j \in F_i} d_j$, and, thus, it is represented as a vector of weighted terms as well, i.e., $f_i = \langle w_{i1}, \dots, w_{im} \rangle$ (see matrix b in Table 1). Note that the folder profile does not take into account the ratings associated with the records, but only the records text content.

Similarly, the *user profile* of a user u (denoted p_u) is built as the centroid of the folder profiles the user has access to and is an indicator of the interests of a user, i.e. if \mathcal{F}_u is the set of folders the user u has access to, then $p_u = (1/|\mathcal{F}_u|) \sum_{F_i \in \mathcal{F}_u} f_i$. Like folder profiles, the user profile is represented as a vector of weighted terms as well, i.e. $p_u = \langle w_{u1}, \dots, w_{um} \rangle$.

The data we represent is summarized in Table 1. Matrix a represents the records, matrix b represents the folder profiles, matrix c represents the average ratings associated with the records belonging to folders, while matrix d represents the user profiles. Worth noting is that w_{ij} represents the weight term t_j assumes (i) in the record d_i in matrix a, (ii) in the folder f_i in matrix b, and (iii) in the user profile p_i in matrix d, respectively.

By relying on matrix a of Table 1, the correlation (similarity) between two rows establishes a similarity between records. Similarly, in matrix b, the correlation between two rows establishes a correlation between folder profiles. Of course, the same applies to matrix d so that a correlation among users may be computed. Furthermore, as all three matrixes (a), (b) and (d) are term based, we also compute

the mixed correlations among them. For instance, we may compute the similarity among folders profiles and user profiles to determine whether the user is interested in the implicit topic described by the folder. All these measures are based on *content* only (using the weight of the terms, but no ratings are taken into account). The measure used for content correlation between two vectors v_1, v_2 taken from matrixes a, b or d (denoted $CSim(v_1, v_2)$) is the well-know *cosine*, i.e., the scalar product between two row vectors (we assume that the rows are already normalized):

$$CSim(v_1, v_2) = \sum_k w_{1k} \cdot w_{2k}. \quad (1)$$

By relying on matrix c, a correlation among folders can be determined by taking into account the ratings issued by users. This similarity is called *rating similarity* of two folders F_1 and F_2 , denoted $RSim(F_1, F_2)$, and is determined using the well-known *Pearson correlation coefficient* (Breese et al., 1998), i.e.,

$$RSim(F_1, F_2) = \frac{\sum_j (r_{1j} - \bar{r}_1) \cdot (r_{2j} - \bar{r}_2)}{\sigma_1 \cdot \sigma_2}, \quad (2)$$

where \bar{r}_i is the mean of the ratings $r_{i1} \dots r_{in}$, and σ_i is their standard deviation.

The combined similarity between two folders is then obtained by taking into account the content similarity ($CSim$) and the rating similarity ($RSim$). In what follows, the *combined similarity* or simply *similarity*, denoted $Sim(F_1, F_2)$, between two folders F_1 and F_2 will be determined as a linear combination between their content similarity and their rating similarity, i.e.,

$$Sim(F_1, F_2) = \alpha \cdot CSim(f_1, f_2) + (1 - \alpha) \cdot RSim(F_1, F_2), \quad (3)$$

where $0 \leq \alpha \leq 1$.

3.2 Record recommendation algorithm

Objective. The objective of the record recommendation algorithm is, given a user u and a folder F_t belonging to u (indicating a topic the user is interested in), called the *target folder*, to recommend to F_t (and, thus, to the user) records relevant to the topic represented by F_t .

Algorithm. Our recommendation algorithm follows a standard four-step schema, see for instance Breese et al. (1998):

1. Select a set $MS(F_t)$ of k -most similar folders to the target folder F_t , according to the similarity measures $xSim$ (we can use $CSim$, $RSim$ or Sim with $\alpha = 0.5$);
2. From this set of similar folders, determine a pool P_D of candidate records, i.e., set of records belonging to the similar folders $F_i \in MS(F_t)$;
3. For each of the records $d_j \in P_D$ compute a recommendation score of d_j to F_t , according to content (cosine between target folder and record) and rating, i.e.,

$$s^R(F_t, d_j) = \bar{r} + \frac{\sum_{F_i \in MS(F_t)} (r_{ij} - \bar{r}_i) \cdot RSim(F_t, F_i)}{\sum_{F_i \in MS(F_t)} RSim(F_t, F_i)}, \quad (4)$$

where $\bar{r}(\bar{r}_i)$ is the mean of the ratings in the target folder F_t , i.e., the mean of the F_t row values in matrix c (mean of $F_i \in MS(F_t)$ rating row). This formula is derived directly from Breese et al. (1998).

4. Recommend to folder F_t records having a positive score.

3.3 User recommendation algorithm

Objective. The goal of the user recommendation algorithm is, given a target folder F_t of user u , to recommend to F_t (and, thus, to user u) those users, which by the system are thought to have overlapping interests with the topic addressed by the folder F_t (and, thus, may be related to user u).

Algorithm. We have analyzed three different algorithms, with increasing level of effectiveness. The first user recommendation algorithm follows a four-step schema, like for the record recommendation case:

1. Select the set $MS(F_t)$ of s -most similar folders to F_t , according to similarity measures. We can use either $CSim$, $RSim$ or the combination of both (Sim), but in Avancini and Straccia (2004a,b) it has already been observed that $Sim(\alpha = 0.5)$ has better effectiveness, so we use it here as well;
2. From this set of similar folders, determine a pool $P_U(F_t)$ of candidate users to be recommended, i.e., let $P_U(F_t)$ be the set of users being owners of the folders in $MS(F_t)$;
3. Compute the *recommendation score* for each possible recommendable user, i.e., for each user $u_k \in P_U(F_t)$ determine the *user hits factor* (where $F_i \in u_k$ means that folder F_i is accessible by user u_k) $h(u_k) = |\{F_i : F_i \in MS(F_t), F_i \in u_k\}|$, i.e., the number of folders F_i judged as similar to the target folder F_t belonging to the same user u_k . For each user $u_k \in P_U(F_t)$ the recommendation score $s(F_t, u_k)$ is computed as follows: $s(F_t, u_k) = h(u_k) \cdot \sum_{F_i \in MS(F_t), F_i \in u_k} Sim(F_t, F_i)$;
4. Recommend to folder F_t , the top- n ranked users, ranked according to the recommendation score.

The intuition behind Step 3 is that the more a user appears among the owners of the top- s similar folders, the more he is considered as relevant to the target folder.

The second algorithm is a variation of the first one in which Step 3 is replaced with:

3. For each user $u_k \in P_U(F_t)$, consider the profile of u_k , p_{u_k} , and compute the recommendation score as the similarity between the user profile p_{u_k} and the profile f of the target folder F_t , i.e., $s(F_t, u_k) = CSim(f, p_{u_k})$;

The intuition here is to use the user profile of recommendable users $u_k \in P_U(F_t)$ directly in place of the folder profiles of similar folders.

Finally, the third algorithm does not consider the set of similar folders, but just compares the profile of the target folder against all user profiles, using $CSim$. Note that in this way, no ratings are taken into account. Therefore, we remove Step 1, and in Step 2, the pool of candidate users, $P_U(F_t)$, is given by all users known to the system, and Step 3 is as in algorithm 2.

3.4 Collection recommendation algorithm

Objective. The objective of the collection recommendation algorithm is, given a target folder F_t of user u , to recommend to F_t those collections that potentially contain records relevant for the topic addressed by the folder itself.

Algorithm. The recommendation algorithm is based on a simple two-step schema:

1. The computation of an approximation of the content of each collection; and
2. The selection of the top- n collections deemed as most relevant to the target folder, relying on the approximations of the collections content.

The first step is done only once for all collections and consists of the computation of a representation of what a collection is about, i.e., the so-called *collection topic* or *language model* of the collection. This data is then used in the second step to compute the *recommendation score* for each collection, i.e., a measure of similarity between the folder topic and the collection topic and, thus, establishes the relatedness of the collection to the target folder.

The language model of a collection consists of a list of terms with their term weight information. We rely on the so-called *query-based sampling method* (Callan & Connell, 2001), which has been proposed for automatically acquiring statistical information about the content of a collection. A major feature is that it requires only that a collection provides a query facility and access to the records, that are in the result of a query. Informally, the method is an iteration of the following steps 2 and 3:

1. Issue a random query to the collection (as start-up);
2. Add the top- k records to the sample;
3. Select randomly a record from the current sample. Select randomly n terms from this record. Build a new query using these terms and issue the query to the collection.

The iteration continues until a stop criterion is satisfied (more details about this technique in our context can be found also in Candela & Straccia, 2004).

As a result, a *sample* set of records for each collection has been gathered. This set is our resource description, or approximation of a collection.

Once we have gathered a sample (approximation) of each collection, we are ready to perform collection recommendation. So, let F_t be a target folder and let f_t be its profile (vector of terms t_{tk} with relative weight w_{tk}). The collection recommendation score $G(F_t, C_l)$ for a collection C_l with respect to the folder F_t is defined as follows:

$$G(F_t, C_l) = \frac{\sum_k p(t_{tk}|C_l)}{|f_t|}, \quad (5)$$

where $|f_t|$ is the number of non-zero weighted terms in the profile f_t . The *belief* $p(t_{tk}|C_l)$ in C_l , for term t_{tk} appearing in F_t profile f_t is computed using a variant of the well-known CORI algorithm (Callan, Lu, & Croft, 1995):

$$p(t_{tk}|C_l) = T_{l,tk} \cdot I_{tk} \cdot w_{tk} \quad (6)$$

$$T_{l,tk} = \frac{df_{l,tk}}{df_{l,tk} + 50 + 150 \cdot \frac{cw_l}{\overline{cw}}} \quad (7)$$

$$I_{tk} = \frac{\log \left(\frac{|C|+0.5}{cf_{tk}} \right)}{\log (|C| + 1.0)} \quad (8)$$

where:

- $df_{l,tk}$ is the number of records in the approximation of C_l containing term t_{tk} ;
- cw_l is the number of terms in the approximation of C_l ;
- \overline{cw} is the mean value of cw over the approximation of C_l ;
- cf_{tk} is the number of approximated collections containing t_{tk} ;
- $|C|$ is the number of the collections.

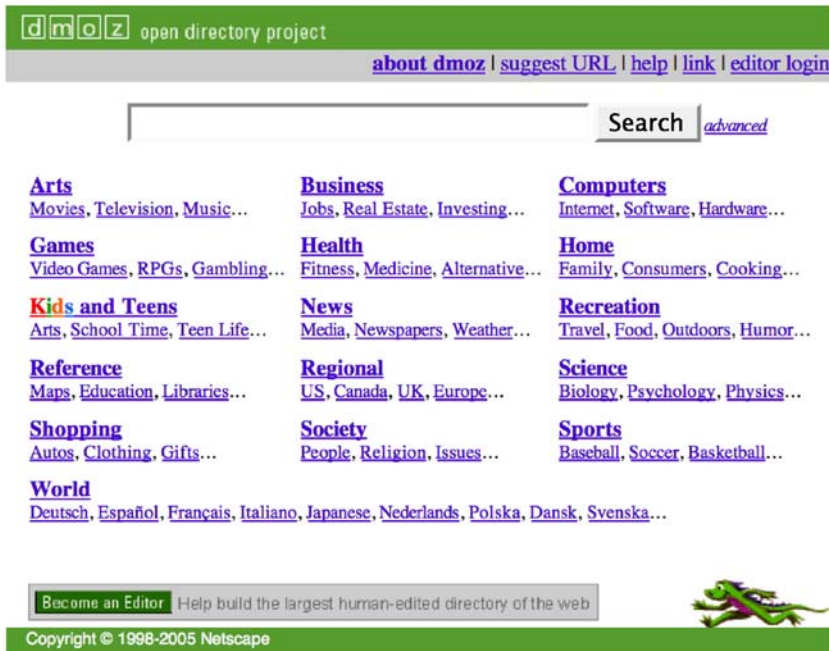
Finally, given a target folder F_t , all collections C_l of a collection environment are ranked according to their collection recommendation score $G(F_t, C_l)$ and the top- n are recommended to F_t .

4 Experimental evaluation

Before we proceed with the description and effectiveness tests of the various algorithms we describe the test corpus we adopted for our experimentation. Then,

ODP - Open Directory Project

<http://dmoz.org/>



5,156,479 sites - 69,719 editors - over 590,000 categories

Fig. 9 DMOZ top categories

for each recommendation type (record, user, collection) we provide the test set, the evaluation method, and the result analysis.

4.1 Corpus

As neither the corpus of data at the end of the CYCLADES project was sufficiently big nor (to the best of our knowledge) there exists an available corpus from the literature, which fits to our setting, we built a suitable corpus by taking the data from the Internet.

The corpus was selected from the *Open Directory Project* hierarchy (ODP or DMOZ).⁶ ODP is among the largest human-edited directories of the Web. The ODP data includes over 5.1 million sites, about 69,000 editors and over 590,000 categories. ODP powers the core directory services for the Web largest search engines and portals, e.g., Google.⁷ Figures 9, 10, 11 and 12 are screenshots of DMOZ categories.

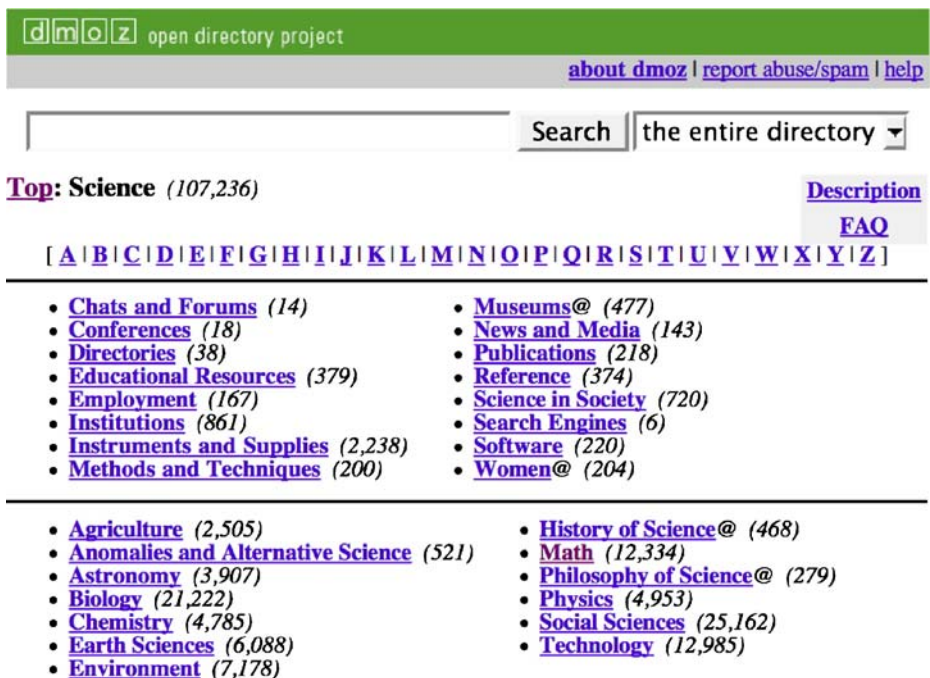



Fig. 10 DMOZ category “Science”

⁶Open Directory Project website <http://dmoz.org>.

⁷Google web site <http://www.google.com>.

 open directory project
 [about dmoz](#) | [suggest URL](#) | [update listing](#) | [report abuse/spam](#) | [help](#)

Top: Science: Math: Number Theory (1,466)
Description

- [Chats and Forums](#) (8)
- [Education](#) (17)
- [Events](#) (254)
- [History](#) (7)
- [Open Problems](#) (48)
- [People](#) (497)
- [Publications](#) (79)
- [Research Groups](#) (72)
- [Software](#) (64)
- [Tables](#) (31)

- [Algebraic](#) (13)
- [Analytic](#) (35)
- [Computational](#) (12)
- [Diophantine Equations](#) (51)
- [Elliptic Curves and Modular Forms](#) (72)
- [Factoring](#) (60)
- [Prime Numbers](#) (106)

See also:

- [Science: Math: Geometry: Algebraic Geometry](#) (135)

This category in other languages:

[French](#) (34)
 [German](#) (18)
 [Russian](#) (8)
 [Spanish](#) (16)
 [Document](#)

- [Number Theory Web](#) ★ - Things of interest to number theorists collected by Keith Matthews.
- [About Two Months of Computing](#) - Work on the 196 palindrome problem by Tim Irvin.
- [Active Elementary Number Theory](#) - Interactive (Javascript) expression of a number as a sum of two squares.
- [Application of Digital Root](#) - A test for the terms of a sequence being expressed by a polynomial, by Roupam Ghosh.
- [The Arithmetic Properties of Binomial Coefficients](#) - Activated text by Andrew Granville.
- [The Canadian Number Theory Association](#) - A loosely defined group of individuals, mainly in Canada, active in research in Number Theory.
- [Covers, Sumsets and Zero-sums](#) - A unified approach to covering systems, restricted sumsets and zero-sum problems by Zhi-Wei Sun.
- [Fermat's Little Theorem](#) - With notes on Carmichael numbers and the life of R.D. Carmichael.
- [Frequently Asked Questions](#) - Number Theory section of the sci.math FAQ list.
- [Hakmem Continued Fractions](#) - Some notes from the MIT collection. Includes Gosper's algorithms for CF arithmetic.
- [Infinite Series Theorem](#) - Addresses the question whether a series of rational functions converges to a rational number.
- [Integer Relations](#) - To determine linear integer dependence among numerical constants and to determine the minimal polynomial of an approximate algebraic number. Interactive or via email.
- [Introduction to Bernoulli Numbers](#) - A web article with a brief history and account of their relationship with the Riemann zeta function and Fermat's Last Theorem (HTML/PS).

Fig. 11 A DMOZ category “Number Theory”

COVERAGE

	#records 10	#records 100	#records 200
<i>Csim</i>	15	30	30
<i>Rsim</i>	17	17	17
<i>Csim+Rsim</i>	19	33	33

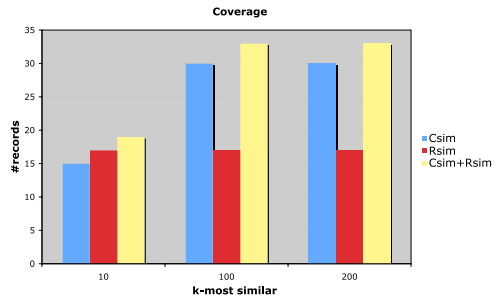


Fig. 12 Record recommendation: coverage

Each category in ODP contains a set of Web documents, which have been evaluated by one or more editors for their relevance to the category. We quote

“Editors select, evaluate, describe and organize Web sites. They are responsible for reviewing submissions to their categories, and list sites according to prescribed editorial guidelines. Editors join the ODP by applying to edit a category that corresponds to their interests. Generally, new editors apply to edit in small categories at first, and then apply to edit additional areas after they have accumulated a number of edits.”

Furthermore, to each document within a category, Google assigns a score, using the so-called PageRank algorithm (Brin & Page, 1998). The PageRank determines a relevance score of the document to the category.

We construct a suitable corpus, which may simulate roughly a CYCLADES scenario, as follows.

1. The set of users is the set of ODP editors.
2. The set of records is the set of documents in ODP.
3. The set of folders is the set of categories in ODP.
4. We considered a subset of the categories under “Science,” together with the involved records and users (1,415 folders, 61 users, 18,091 records). All the profiles (of the folders and users) have been restricted to the top weighted 100 terms.
5. To each record d_j in folder F_i , evaluated by user u_k , we simulate the user’s rating r_{ijk} randomly such that the average rating over all users rating record d_j in folder F_i , is equal to the score s_{ij} assigned to record d_j w.r.t. folder F_i .
6. In order to test the effectiveness of the collection recommendation algorithm, we consider different environments of 23, 85 and 100 collections of records, respectively. Each set is characterized by a different degree of heterogeneity of the collections (see description below). The idea is that the more heterogeneous is a collection the more difficult is the task.

The three different collection environments we consider in our tests have been built as follows from the ODP records. In Table 2 we report the characteristics of these three environments of collections.

Table 2 The experimental environments

Environment	Number of collections	Collection size		
		Max	Min	Avg
ScienceI	23	15,722	11	3,027.47
Quasi-random	85	5,660	1	809.67
Random	100	768	621	696.32

- ScienceI:* Each of the 23 collections contains the records that are in the corresponding first-level sub-category of Science into the DMOZ hierarchy. Therefore, the records of a collection share the *same* topic;
- Quasi-random:* Each of the 85 collections contains a set of records, which have been selected randomly, from a subset of the first-level sub-category of the category Science. In this case, the collections are *more heterogeneous* than the previous case as the records of a collection belong to a different, but limited, set of categories. That is, a collection is multi-topic;
- Random:* Each of the 100 collections contains a set of records randomly selected from the whole corpus. Therefore, the collections are *highly heterogeneous*.

4.2 Record recommendation algorithm evaluation

- Test set.* To create the test set, we considered the set \bar{D} of all records of the corpus, which belong to at least two folders. For each of these records $d_j \in \bar{D}$ (300 records), we randomly choose a folder F_t in which d_j occurs. The set of chosen folders $\bar{F} = \{F_t\}$ (250 folders), the records d_j and the relative average rating r_{ij} forms the test set, i.e., the test set in this case is a set of triples (F_t, d_j, r_{ij}) , where we check whether our recommendation algorithm recommends record d_j to the *target* folder F_t and how well the computed recommendation score agrees with the rating r_{ij} .
- Evaluation method.* For each target folder $F_t \in \bar{F}$, we compute the set of recommended records $\{d_j\}$ and consider their recommendation score $s^R(F_t, d_j)$. We measure the *coverage*, which is the percentage of records correctly recommended to the target folders, i.e., the system recommends a record to a folder and the record belongs to the test set.
- Result analysis.* The results of our experiments are summarized below. We considered two different parameters:
1. The strategy to select the similar folders w.r.t. a target folder; and
 2. The impact of varying k , the k – *most* similar folders threshold.

In case 1, we have three different options: select similar folders by relying on rating similarity only (using *RSim*), or using content similarity only (using *CSim*), or to use both (*Sim*). In case 2, k took one of the values 10, 100 and 200.

Figure 12 reports the coverage in using *CSim* (left bar), *RSim* (center bar) and *Sim* (right bar) for various k . The combined approach *Sim* is the best performing one.

Note that all cases reach a stability at $k \geq 100$. For $k = 100$, the combined approach performs 49% better than relying on *RSim* only and 10% better than relying on *CSim* only.

In Fig. 13 we report the fraction of records, which have been recommended by using both *RSim* and *CSim*. For instance, consider the “CSim Coverage” figure. The whole bar reports the coverage of recommended records using *CSim*. The part labeled “*CSim & RSim*” indicates the amount of records that have been recommended both using *CSim* and *RSim*, while the label “Only CSim” indicates the amount of records that have been recommended using *CSim* but *can not* be recommended using *RSim* and, thus, indicates the added value of using *CSim*. The explanation for the “*RSim Coverage*” is similar. Interestingly, we can observe that a large part of records, that have been recommended using *CSim* only, can not be recommended using *RSim* and vice-versa. This confirms that the combination of both strategies can improve the coverage significantly.

4.3 User recommendation algorithm evaluation

Test set.

To create the test set, we considered the set \bar{U} of all users of the corpus, which have access to at least two folders. For each

CSIM COVERAGE

	#records 10	#records 100	#records 200
<i>Csim & Rsim</i>	7	10	10
<i>Only Csim</i>	8	20	20

RSIM COVERAGE

	#records 10	#records 100	#records 200
<i>Csim & Rsim</i>	7	10	10
<i>Only Csim</i>	10	7	7

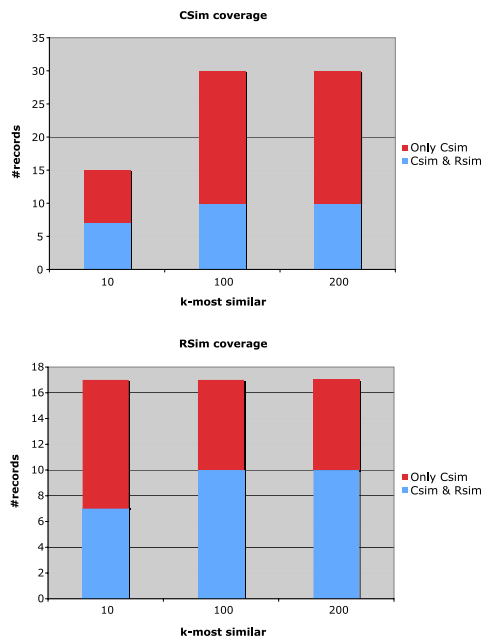


Fig. 13 Record recommendation: *RSim* vs. *CSim*

of these users $u_k \in \bar{U}$ (61 users), we randomly choose a folder F_t user u_k has access to. The set of chosen folders $\bar{F} = \{F_t\}$ (49 folders) and the users u_k forms the user test set, i.e., the test set in this case is a set of pairs (F_t, u_k) , where we check whether our recommendation algorithm recommends user u_k to the target folder F_t .

Evaluation method. For each target folder $F_t \in \bar{F}$, we compute the set of recommended users $\{u_k\}$ and ranked them according to their recommendation score $s(F_t, u_k)$. In this rank, we highlight the rank position of user u_k . If the recommendation score is 0, user u_k is ranked 0. We compute *precision* and *recall*, by recommending the top- n users for each target folder, where $n \in \{1, 2, 5, 10\}$.

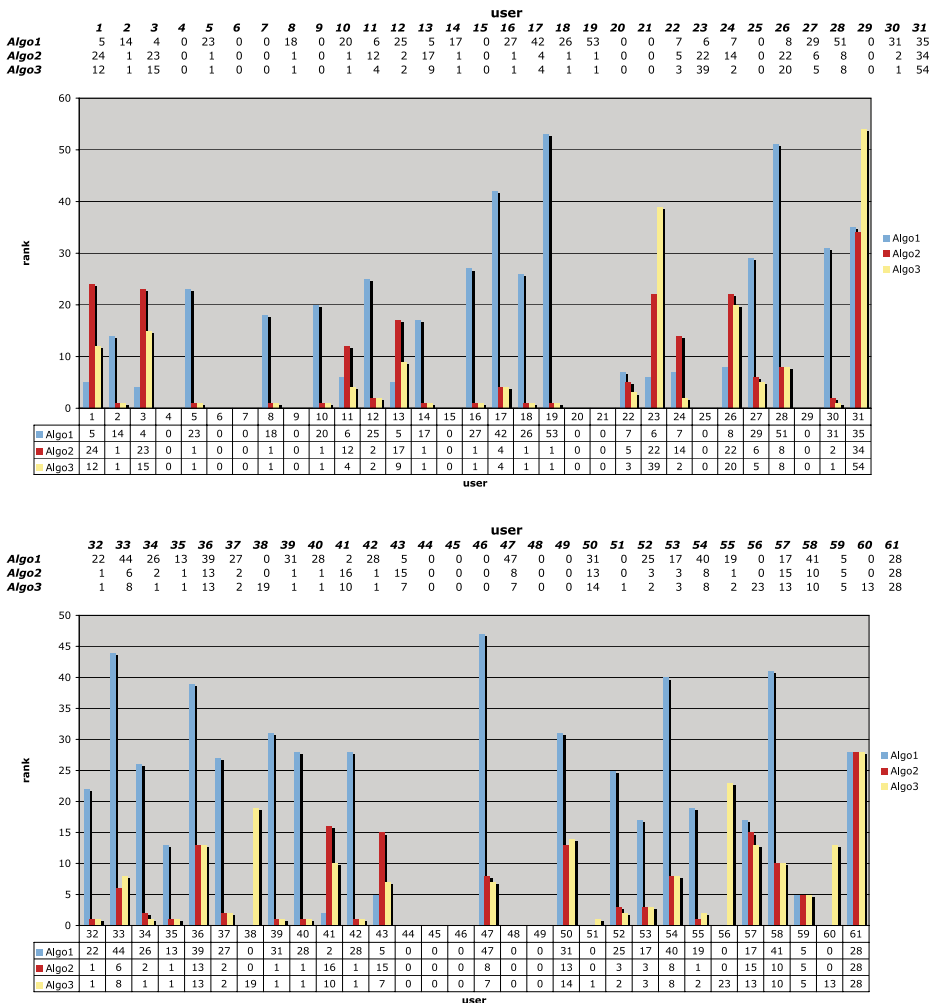


Fig. 14 Rank position of correctly recommended users

Precision is the fraction between correctly recommended users and the total number of recommended users, i.e.,

$$Precision_n = \frac{|CorRec_n|}{49 \cdot n}, \quad (9)$$

Where $CorRec_n$ is the set of correctly recommended users among all 49 tests and n is the number of recommended users for each test. Recall is the ratio between correctly recommended users and the number of test users, i.e.,

$$Recall_n = \frac{|CorRec_n|}{61}. \quad (10)$$

Of course, the more users we recommend (i.e., n increases), the more correctly recommended users we have (recall improves), but the less precise we are (precision decreases). We also report the standard measure

$$F1_n = \frac{2 \cdot Precision_n \cdot Recall_n}{Precision_n + Recall_n}, \quad (11)$$

which gives us an estimate of the combination between precision and recall, as well as their average rank.

Result analysis.

The results of our experiments are summarized in Fig. 14 and Table 3. In Fig. 14, for each test pair (F_i, u_k) we report the rank position of user u_k , according to our recommendation algorithms 1, 2 and 3, respectively. It turns out that algorithm 3 performs best, algorithm 2 is second, while the least effective one is algorithm 1. As we can see, for algorithms 2 and 3, in

Table 3 Effectiveness measures

	Algorithm	Top-1	Top-2	Top-5	Top-10
Precision	1	0.00	0.01	0.02	0.02
	2	0.29	0.18	0.09	0.06
	3	0.33	0.21	0.11	0.07
Recall	1	0.00	0.02	0.10	0.18
	2	0.23	0.30	0.38	0.48
	3	0.26	0.34	0.44	0.57
F1	1	undef.	0.01	0.04	0.04
	2	0.25	0.23	0.15	0.11
	3	0.29	0.26	0.18	0.13

Algorithm	Average rank
1	23.18
2	8.28
3	8.09

many cases the correctly recommended user is highly ranked. It turns out that algorithm 3, which does not rely on the ratings provided by the users, is the most effective one. Interestingly, recommending just the top ranked user is the most satisfactory compromise between precision and recall. An exception is algorithm 1, where taking the top ranked user only failed to be the right one in all test cases. Therefore, precision and recall are 0 and, thus, $F1$ is undefined.

4.4 Collection recommendation algorithm evaluation

- Test set.* To create the test set, we considered the 300 folders belonging to three different categories of Science, i.e., Agriculture, Anomalies and Alternative Science, and Astronomy. These folders are our target folders F_t , where we check how good our recommendation algorithm is in recommending relevant collections to the folder F_t ;
- Evaluation method.* The evaluation phase is rather involved and consists of two parts:
1. The evaluation of the effectiveness of a collection approximation, i.e., comparing the learned resource description of a collection with the real resource description for that collection; and
 2. The evaluation of the effectiveness of the collection recommendation.

Concerning the evaluation of the effectiveness of a collection approximation, in accordance with Callan and Connell (2001), we have used two metrics to evaluate the quality of the approximation: the *ctf ratio* (CTF) to measure the portion of term occurrences in the collection vocabulary (V) that are covered by the terms in the learned vocabulary (V') and the *Spearman Rank Correlation Coefficient* (SRCC) to measure the correspondence between the learned and the actual document frequency information. These metrics are calculated using Eqs. 12 and 13 below, where ctf_i is the number of times term t_i occurs in the collection C , $\delta(i)$ is the rank difference of a common term $t_i \in V' \cap V$. The two term rankings are based on the learned and the actual document frequency df_i . n is the total number of common terms, i.e., $n = |V' \cap V|$.

$$CTF = \frac{\sum_{t_i \in V'} ctf_i}{\sum_{t_i \in V} ctf_i} \quad (12)$$

$$SRCC = 1 - \frac{6}{n^3 - n} \sum_{t_i \in V' \cap V} \delta(i)^2. \quad (13)$$

Concerning the second evaluation point (i.e., the effectiveness of collection recommendation) we require some additional steps. Indeed, we neither know a priori, which are the correct collections to be recommended to a target folder, nor we are

willing to evaluate it manually. A way to automated effectiveness evaluation of a collection recommendation algorithm is as follows. For each of the three collections environments, we have built a collection containing all the records belonging to the environment. This collection is called *baseline collection*. The baseline collection has been indexed. This index is called *baseline index*. For indexing (and retrieval, see later on) we have used the Jakarta Lucene search engine.⁸ So, we have three baseline collections and indexes, one for each environment. We also indexed all collections of all environments, which we call *collection index*.

For each target folder F_t , we consider its profile f_t . Each folder profile represents a query that is matched against a baseline collection of a collection environment. As a result, the retrieval engine give us back a ranked list $rank(F_t)$ (called *baseline rank*) of records belonging to the environment, deemed by the system as relevant to the query, i.e., relevant to the folder F_t .

Now, we proceed as follows. For each collection environment and for each target folder F_t , we use the recommendation algorithm to compute a *recommendation score* $G(F_t, C_i)$ for each collection C_i of the environment. We select the top- n collections according to their score. Then we match the query f_t against each collection C_i of this list and get a ranked list of records $ret_i(F_t)$. Now we merge these ranked lists into a unique ranked list $ret(F_t)$. The fusion of multiple ranked lists is done by computing, for each record, a normalized global retrieval score, $RSV_g(d_j)$, obtained from its normalized retrieval score, $RSV'_i(d_j)$, within the ranked lists, i.e.,

$$RSV_g(d_j) = \frac{RSV'_i(d_j) + 0.4 \cdot RSV'_i(d_j) \cdot G(F_t, C_i)}{1.4} \quad (14)$$

$$RSV'_i(d_j) = \frac{RSV_i(d_j) - \min_k(RSV_i(d_k))}{\max_k(RSV_i(d_k)) - \min_k(RSV_i(d_k))}, \quad (15)$$

where $RSV_i(d)$ is the score of record d in the ranked list $ret_i(F_t)$ (note that a record does appear at most in one ranked list, as there is no record overlapping among the collections). Therefore, for target folder F_t , on one hand we have the baseline rank $rank(F_t)$ which has been computed over the baseline collection (without selecting collections) and on the other hand we have the rank $ret(F_t)$ obtained from the fusion of rank lists *with* collection selection.

Table 4 Collections and their samples characteristics

Environment	Avg record			Avg terms		
	Source	Sample	Sample(%)	Source	Sample	Sample(%)
Science I	3,027.47	279.17	9.22	102,375	17,707.2	17.29
Quasi-random	809.67	238.3	29.43	46,197.7	18,083.17	39.14
Random	696.32	295.74	42.47	48,571.31	27,121.01	55.83

⁸Jakarta Lucene search engine <http://jakarta.apache.org/lucene>.

Table 5 Sample statistics

Environment	CTF			SRCC		
	Max(%)	Min(%)	Avg(%)	Max(%)	Min(%)	Avg(%)
ScienceI	98	74	87	97	53	80
Quasi-random	100	71	98	100	63	85
Random	92	78	87	90	75	85

The effectiveness of the recommendation algorithm is computed as the similarity among the two ranked lists $rank(F_t)$ and $ret(F_t)$ and is based on set-based measures like Precision and Recall defined as follows:

$$Precision = \frac{|rank(F_t) \cap ret(F)|}{|ret(F_t)|} \quad (16)$$

$$Recall = \frac{|rank(F_t) \cap ret(F_t)|}{|rank(F_t)|} \quad (17)$$

We have also considered the harmonic mean ($F1$ -score) of these two values computed as $F1 = 2 \cdot Precision \cdot Recall / (Precision + Recall)$.

As the measure of Precision and Recall does not consider the ordering of the two ranked lists to compare, we have also used the $SRCC$ to measure the correspondence between the two ranked lists, i.e., $SRCC = 1 - \frac{6}{n^3 - n} \sum \delta(i)^2$, where this time $\delta(i)$ is the rank difference of a record i appearing in both ranks and n is the total number of records common to both ranks.

Result analysis. At first, we report the evaluation of the effectiveness of the sampling method.

In Table 4 we report the characteristics of the collections and their approximations in terms of number of records gathered and number of terms in the samples, respectively.

In Table 5 we report the results of CTF and $SRCC$ effectiveness metrics. For example, observing the ScienceI case, we can note that acquiring just the 9% of the records of the source we are able to have a very close representation of the content of the collections as we obtain a CTF of about 90% and an $SRCC$ of 80%. Moreover,

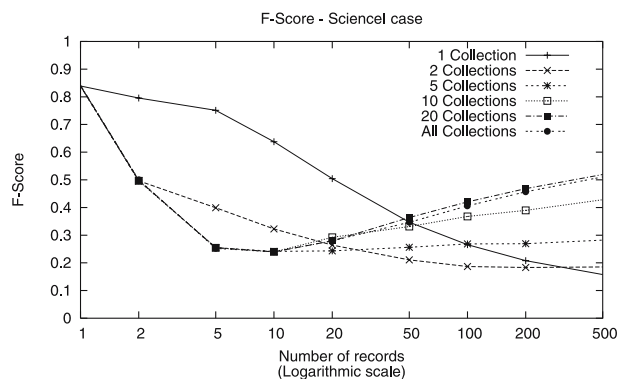
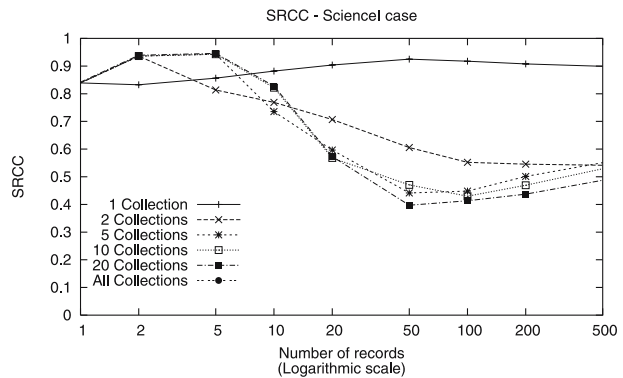
Fig. 15 ScienceI case: F1 values

Fig. 16 ScienceI case: SRCC values

note that the effectiveness of the approximations are quite independent from the content homogeneity of the collections in the collection environments. But, for instance, by observing the Random environment case, we can note that the number of records acquired by the sample process is much greater in percentage than that for the ScienceI case in order to get similar *CTF* and *SRCC* values. Essentially, and quite intuitively, the more heterogeneous is a collection with respect to its content, the more records have to be gathered in its sample to reasonably approximate the collection content.

At second, we report the evaluation of the effectiveness of the recommendation algorithm. For each target folder, we have evaluated the effect of varying the number of collections to be recommended and the number of records (from 1 to 500) to be considered in each of the two ranked lists $rank(F_t)$ and $ret(F_t)$.

In Figs. 15 and 16, we report the results for the ScienceI case. As previously observed, in this case a collection is homogeneous, i.e., its records belong to the same topic or set of topics. The recommendation of just the top-1 collection for each target folder produces high F1 value, about 80%, and a high SRCC value. This means that our algorithm is able to find the most appropriate collection w.r.t. the target folder in case the collection has a homogeneous topic. Recommending more than one collection produces a deterioration of the results as the selected collections contain documents less pertinent to the target folder topic.

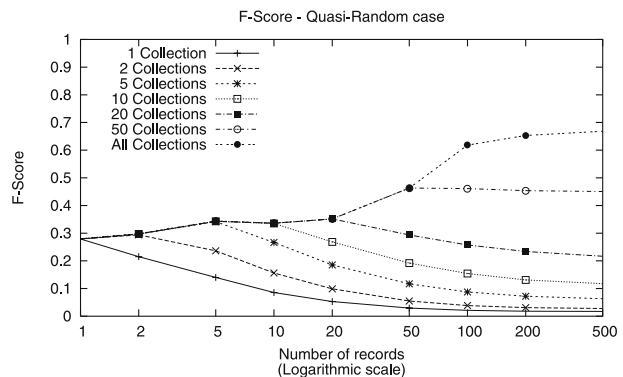
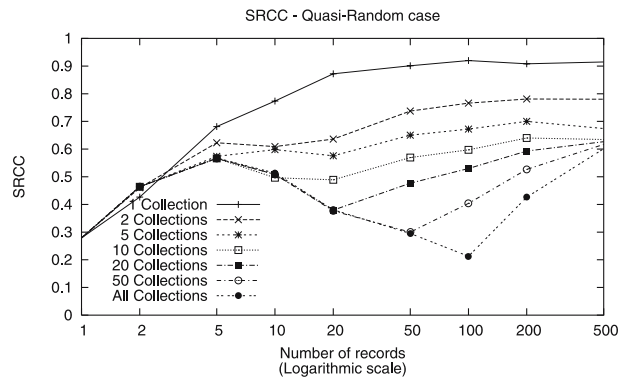
Fig. 17 Quasi-random case: F1 values

Fig. 18 Quasi-random case: SRCC values



In Figs. 17 and 18 we report the results for the Quasi-random case. As previously observed these collections are more heterogeneous than the previous ones, i.e., the records of a collection may belong to a different but limited set of topics, and records about a topic are distributed among a limited number of collections. In this case recommending just the top-1 collection to each target folder produces a lower F1 than for the previous case. Moreover, the performance decreases if the number of records to be recommended increases. Recommending more than one collection produces an improvement of the results as the selected collections contain records relevant to the target folder. Concerning the *SRCC* curve we note that increasing the number of records in the rank list, each curve has a decreasing phase and finally it increases. The end of the decreasing phase coincides with the point where the F1 value starts to decrease. This indicates that the number of common records, between the baseline rank and the test rank, decreases and, thus this improves the *SRCC*. On the other hand, the numerator of both precision and recall decreases and, thus, F1 decreases.

In Figs. 19 and 20 we report the results for the Random case. In this environment, the collections are highly heterogeneous, i.e., the records of a collection can belong to many different topics, and the records of a category are distributed, potentially, among all collections. This is the worst case. We can note that the performance decreases if the number of records in the rank list increases, while it increases if the

Fig. 19 Random case: *F*-score values

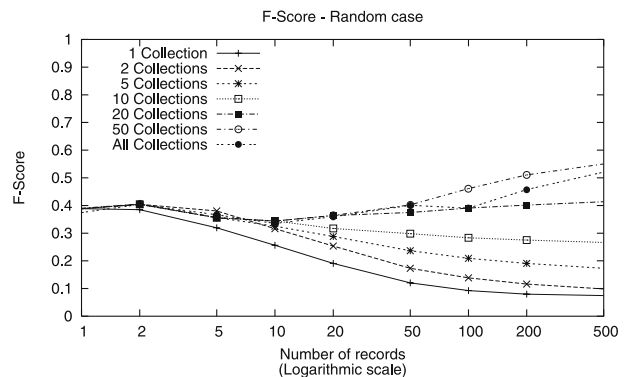
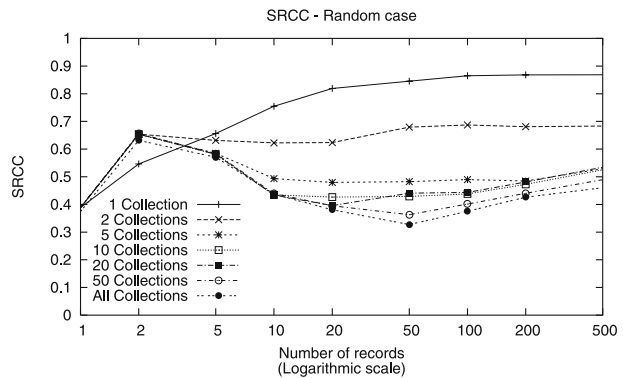


Fig. 20 Random case: *SRCC* values

number of collections to recommend increases. However, our algorithm is still able to find the most appropriate collection w.r.t. the target folder.

5 Conclusion

Since the Web, and consequently the information contained in it, is growing rapidly, every day a huge amount of “new” information is electronically published and new Digital Libraries are available to satisfy the user information needs. We described here a Digital Library environment that is not only an information resource where users may submit queries to get what they are searching for, but also a collaborative working and meeting space. Indeed, users looking within an information resource for relevant data might have overlapping interests, which may turn out to be of reciprocal interest for the users: users might well profit from each other’s knowledge by sharing opinions and experiences. As such, we have presented CYCLADES, a personalized collaborative Digital Library environment, in which the user functionality may be organized into four categories: users may (i) search for information; (ii) organize the information space (according to the “folder paradigm”); (iii) collaborate with other users sharing similar interests; and (iv) get recommendations. CYCLADES not only provides recommendation of records, as it usually happens in personalization system dealing with documents, but by taking advantage of the highly collaborative environment, it may recommend also communities, collections and users. Particular attention has been paid to the recommendation part and to the experiments showing the effectiveness of the adopted algorithms.

Acknowledgement This work has been partially funded by the European Community in the context of the CYCLADES project IST-2000-25456, under the Information Societies Technology programme (see, <http://www.ercim.org/cyclades/>).

References

- Avancini, H., & Straccia, U. (2004a). Personalization, collaboration, and recommendation in the digital library environment CYCLADES. In *Proceedings of the IADIS International Conference Applied Computing (AC-04)* (pp. 589–596). Lisbon, Portugal: IADIS.

- Avancini, H., & Straccia, U. (2004b). User awareness in CYCLADES: a collaborative and personalized environment for open digital archives. In *Proceedings of the IADIS International Conference on Web Based Communities (WBC-04)* (pp. 67–74). Lisbon, Portugal: IADIS.
- Bentley, R., Appelt, W., Busbach, U., Hinrichs, E., Kerr, D., Sikkell, K., et al. (1997). Basic support for cooperative work on the World Wide Web. *International Journal of Human Computer Studies*, 46, 827–846.
- Billsus, D., & Pazzani, M. J. (1998). Learning collaborative information filters. In *Proceedings of the 15th International Conference on Machine Learning* (pp. 46–54). San Francisco, California: Morgan Kaufmann.
- Bollacker, K., Lawrence, S., & Giles, C. L. (1999). A system for automatic personalized tracking of scientific literature on the web. In *Digital Libraries 99 - The Fourth ACM Conference on Digital Libraries* (pp. 105–113). New York: ACM.
- Breese, J. S., Heckerman, D., & Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the 14th Annual Conference on Uncertainty in Artificial Intelligence (UAI-98)* (pp. 43–52). Madison, Wisconsin, USA.
- Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual Web search engine. *Computer Networks ISDN Systems*, 30(1–7), 107–117.
- Callan, J. (2000). Distributed information retrieval. In W. Croft (Ed.), *Advances in Information Retrieval* (pp. 127–150). Hingham, MA, USA: Kluwer Academic.
- Callan, J., & Connell, M. (2001). Query-based sampling of text databases. *ACM Transactions on Information Systems*, 19(2), 97–130.
- Callan, J., Lu, Z., & Croft, B. W. (1995). Searching distributed collections with inference networks. In *Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (ACM SIGIR-95)* (pp. 21–28). New York: ACM.
- Candela, L., & Straccia, U. (2004). The personalized, collaborative digital library environment CYCLADES and its collections management. In J. Callan, F. Crestani, and M. Sanderson (Eds.), *Multimedia Distributed Information Retrieval, Lecture Notes in Computer Science*, 2924 (pp. 156–172). Berlin Heidelberg New York: Springer.
- Di Giacomo, M., Mahoney, D., Bollen, J., Monroy-Hernandez, A., & Meraz, C. M. R. (2001). MyLibrary, a personalization service for digital library environments. In *Proceedings of the Second DELOS Network of Excellence Workshop on Personalisation and Recommender Systems in Digital Libraries*. Ireland: Dublin City University.
- Faensen, D., Faulstich, L., Schweppe, H., Hinze, A., & Steidinger, A. (2001). Hermes: a notification service for digital libraries. In *ACM/IEEE Joint Conference on Digital Libraries* (pp. 373–380). New York: ACM.
- Fernandez, L., Sanchez, J. A., & Garcia, A. (2000). MiBiblio: personal spaces in a digital library universe. In *ACM DL* (pp. 232–233). New York: ACM.
- Fox, E. A., & Marchionini, G. (2001). Digital libraries: Introduction. *Communications of the ACM*, 44(5), 30–32.
- Fuhr, N. (1999). A decision-theoretic approach to database selection in networked IR. *ACM Transactions on Information Systems*, 3(17), 229–249.
- Goldberg, D. J., Nichols, D., Oki, B. M., & Terry, D. (1992). Using collaborative filtering to weave information tapestry. *Communications of the ACM*, 35(12), 61–70.
- Good, N., Schafer, J., Konstan, J., Borchers, A., Sarwar, B., Herlocker, J., et al. (1999). Combining collaborative filtering with personal agents for better recommendations. In *Proceedings 1999 Conf. American Association of Artificial Intelligence* (pp. 439–446). Menlo Park, California: American Association for Artificial Intelligence.
- Herlocker, J. L., Konstan, J. A., Borchers, A., & Riedl, J. (1999). An algorithmic framework for performing collaborative filtering. In *Proceedings of the 22th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (ACM SIGIR-99)* Berkeley, California, USA (pp. 230–237). New York: ACM.
- Konstan, J., Miller, B., Maltz, D., Herlocker, J., Gordon, L., & Riedl, J. (1997). Applying collaborative filtering to Usenet news. *Communications of the ACM*, 40(3), 77–87.
- Meng, W., Liu, K., Yu, C., Wang, X., Chang, Y., & Rishe, N. (1998). Determining text databases to search in the internet. In *Proceedings of the 24th International Conference on Very Large Data Bases (VLDB-98)* (pp. 14–25). New York, USA.
- Moukas, A. (1996). *Amalthaea*: Information discovery and filtering using a multiagent evolving ecosystem. In *Proceedings Practical Applications of Agents and Multiagent Technology* (pp. 437–457). New York: Taylor and Francis.

- Rasolofo, Y., Abbaci, F., & Savoy, J. (2001). Approaches to collection selection and results merging for distributed information retrieval. In *Proceedings of the 10th International Conference on Information and Knowledge Management (CIKM-01)* (pp. 191–198). New York: ACM.
- Renda, M. E., & Straccia, U. (2002). A personalized collaborative digital library environment. In *5th International Conference on Asian Digital Libraries (ICADL-02), Lecture Notes in Computer Science 2555*, Singapore, Republic of Singapore (pp. 262–274). Berlin Heidelberg New York: Springer.
- Resnick, P., Iacovou, N., Suchak, M., Bergstorm, P., & Riedl, J. (1994). GroupLens: An open architecture for collaborative filtering of netnews. In *Proceedings of ACM 1994 Conference on Computer Supported Cooperative Work*, Chapel Hill, North Carolina (pp. 175–186). New York: ACM.
- Rocha, L. M. (1999). Talkmine and the adaptive recommendation project. In *ACM DL* (pp. 242–243). New York: ACM.
- Salton, G., & McGill, J. M. (1983). *Introduction to modern information retrieval*. Reading, Massachusetts: Addison Wesley.