



ELSEVIER

Contents lists available at ScienceDirect

Journal of Informetrics

journal homepage: www.elsevier.com/locate/joi

Visualization of co-readership patterns from an online reference management system

Peter Kraker^{a,*}, Christian Schlögl^b, Kris Jack^c, Stefanie Lindstaedt^a^a Know-Center, Inffeldgasse 13, 8010 Graz, Austria^b University of Graz, Universitätsstraße 15, 8010 Graz, Austria^c Mendeley, 144a Clerkenwell Road, EC2 RF London, UK

ARTICLE INFO

Article history:

Received 5 August 2014

Received in revised form 3 December 2014

Accepted 8 December 2014

Available online 12 January 2015

Keywords:

Relational scientometrics

Topical distribution

Knowledge domain visualization

Mapping

Altmetrics

Readership statistics

ABSTRACT

In this paper, we analyze the adequacy and applicability of readership statistics recorded in social reference management systems for creating knowledge domain visualizations. First, we investigate the distribution of subject areas in user libraries of educational technology researchers on Mendeley. The results show that around 69% of the publications in an average user library can be attributed to a single subject area. Then, we use co-readership patterns to map the field of educational technology. The resulting visualization prototype, based on the most read publications in this field on Mendeley, reveals 13 topic areas of educational technology research. The visualization is a recent representation of the field: 80% of the publications included were published within ten years of data collection. The characteristics of the readers, however, introduce certain biases to the visualization. Knowledge domain visualizations based on readership statistics are therefore multifaceted and timely, but it is important that the characteristics of the underlying sample are made transparent.

© 2015 Elsevier Ltd. All rights reserved.

1. Introduction

In recent scientometric literature, usage data is being discussed as a valuable alternative to citations. With the advent of e-journals, digital libraries, and web-based archives, click and download data have been suggested as a potential alternative to citations (Kurtz et al., 2005; Rowlands & Nicholas, 2007). Compared to citation data, usage data has the advantage of being available earlier, shortly after a paper has been published. In many instances, usage statistics are also easier to obtain and collect (Bollen, Sompel, Smith, & Luce, 2005; Brody, Harnad, & Carr, 2006; Haustein & Siebenlist, 2011). Furthermore, usage statistics allow for an analysis of publications and research outputs that do not receive citations or for which citations are not tracked (Priem & Hemminger, 2010).

Another type of usage data besides clicks and downloads is created in social reference management systems like BibSonomy¹ and Mendeley.² These systems enable users to store their references in a personal library and share them

* Corresponding author. Tel.: +43 316 87330844.

E-mail addresses: pkraker@know-center.at (P. Kraker), christian.schloegl@uni-graz.at (C. Schlögl), kris.jack@mendeley.com (K. Jack), slind@know-center.at (S. Lindstaedt).¹ <http://bibsonomy.org>² <http://mendeley.com>

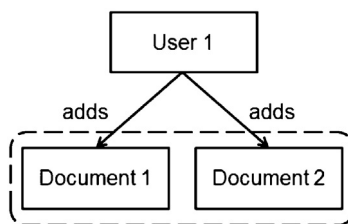


Fig. 1. Co-readership of two documents is established when at least one user has added the two documents to his or her user library.

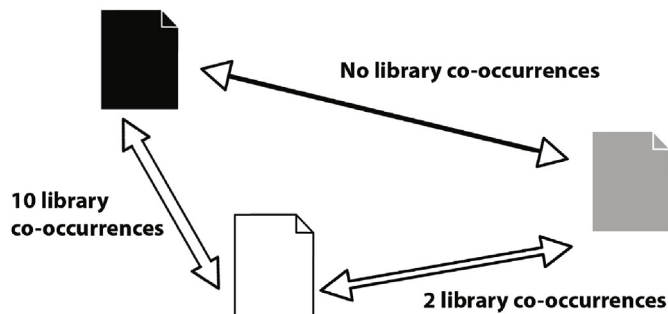


Fig. 2. Relationships between documents in a field based on co-readership. Co-occurrence in user libraries is employed as a measure of subject similarity.

with other people. The number of times an article has been added to user libraries is commonly referred to as the number of readers, or in short readership.³

Readership statistics have been of high scientometric interest in recent years. It has been shown that readership statistics provide a good coverage of top publications (Bar-Ilan et al., 2012), and that there is a medium correlation between readership data and citations (Schlögl et al., 2013) and a medium to high correlation between the impact factor and journal readership (Kraker, Körner, Jack, & Granitzer, 2012). Furthermore, Jiang, He, and Ni (2011) employ readership statistics from CiteULike to form clusters based on the occurrence and co-occurrence of articles in user libraries. They also correlate these clusters with ISI subject categories, and find them as effective as citation-based clusters when removing journals that cannot be found in CiteULike.

Therefore, we consider co-readership as a measure of subject similarity. Co-readership relation between two documents is established when at least one user has added the two documents to his or her user library (see Fig. 1). We assume that the more often the same two documents have been added to user libraries, the more likely they are of the same or a similar subject. The topical relationship established by co-readership can then be exploited for visualizations by clustering those papers that have high co-readership numbers (see Fig. 2). To the best of our knowledge, this measure has not been exploited before for knowledge domain visualization.

In this study, we first investigate the distribution of subject areas in user libraries of educational technology researchers on Mendeley. Then, we employ co-readership patterns for knowledge domain visualization to explore the field of educational technology. Educational technology is multi-disciplinary and highly dynamic in nature, as it is influenced by changes in pedagogical concepts and emerging technologies (Siemens & Tittenberger, 2009), as well as social change (Czerniewicz, 2010). Therefore, it seemed to be especially appropriate for this kind of analysis.

2. Related work

Traditionally, knowledge domain visualizations are based on citations. Small (1973) and Marshakova (1973) proposed co-citation as a measure of subject similarity and co-occurrence of ideas (see Fig. 3, left side, for a graphical representation of the relationship). This relationship can be employed to cluster documents, authors, or journals from a certain field and to map them in a two-dimensional space. Co-citation analysis has been used to map many fields, for instance information management (Schlögl, 2001, p. 48), hypertext (Chen & Carr, 1999), and educational technology (Chen & Lien, 2011) to name just a few. Furthermore, co-citation analysis has also been used to map out all of science (Boyack, Klavans, & Börner, 2005; Small, 1999).

³ Initially, the term readership might seem a bit misleading, because the addition of an article to a user library does not guarantee that the article has actually been read by said user. Nevertheless, researchers need to make a second decision after downloading an article before they add it to their user libraries. Furthermore, the term is already well established among researchers (see e.g. Bar-Ilan et al., 2012; Hausteijn & Larivière, 2014; Thelwall & Mafrahi, 2014; Zahedi, Costas, & Wouters, 2014); thus we use it in our research for reasons of consistency and to avoid neologisms.

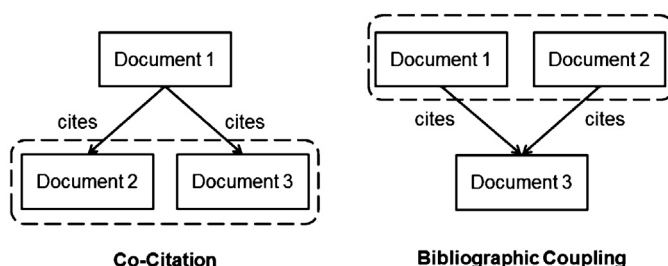


Fig. 3. Relationships between documents on the basis of citations.

Adapted from Schlögl (2001).

There is, however, a significant problem with citations: they take a long time to appear. It takes around two to six years after an article is published before the citation count peaks (Amin & Mabe, 2000). Therefore, visualizations based on co-citations – and indeed all analyses that are based on incoming citations – have to deal with this time lag. Bibliographic coupling (Kessler, 1963) presents an alternative to co-citation analysis; it is formed when two documents cite the same source document (see Fig. 3, right side). The more publications in the reference list the two documents have in common, the more related they are.

Bibliographic coupling is based on outgoing citations available at the time of publication and can therefore be used to map the research front. One difference between bibliographic coupling and co-citation analysis is that the former is a retrospective method (Garfield, 2001), which means that the relationship between two documents cannot change over time. For an overview of the properties and the accuracy of the two citation-based mapping techniques refer to Egghe and Rousseau (1990, chap. III.4) and Boyack and Klavans (2010).

In contrast to citations, usage statistics have been almost exclusively used in evaluative scientometrics (see e.g. Bollen, Rodriguez, & Van de Sompel, 2007; Darmoni, Roussel, Benichou, Thirion, & Pinhas, 2002; Schloegl & Gorraiz, 2010). There are only a handful of examples in relational scientometrics and knowledge domain visualization. One of the first are Polanco, Ivana, and Dominique (2006), who propose to use co-occurrences of document requests for clustering and mapping. Bollen and van de Sompel (2006) use consecutive accesses to journal articles as a measure of journal relationships. They derive clusters of journals which are statistically significantly related to ISI subject categories. In a later study, Bollen et al. (2009) create an overview map of all of science. The authors collect hundreds of millions of user interactions with digital libraries and bibliographic databases. Then, they re-create click-streams for each user, aggregated by journal, and apply network analysis to them. Among the challenges of the approach, the authors name that clickstreams need to be aggregated from various data sources. The varying user interfaces and the difference between reader and author population may introduce biases to the visualization (Bollen et al., 2008).

In social reference management systems we can address these challenges. First, we are able to use library co-occurrence from a single service as a basis for mapping the intellectual structure of a scientific domain. Second, being able to precisely attribute papers to individual readers allows for a better understanding of the results as the information found in the user profile adds further context. With the help of profile information, we can for example analyze the influence of different user groups. When using library co-occurrence, however, we are missing the temporal aspect represented in clickstreams, which may play a role when establishing subject similarity.

3. Data source

All data for this study was sourced from Mendeley on 10 August 2012. Mendeley provides users with software tools that support them in conducting research (Henning & Reichelt, 2008). One of the most popular of these tools is Mendeley Desktop, a cross-platform, freely downloadable PDF and reference management application. It allows users to organize their personal libraries into folders and apply tags to them for later retrieval. The articles, added by users around the world, are then crowd-sourced into a single collection called the Mendeley research catalog (Hammerton, Granitzer, Harvey, Hristakeva, & Jack, 2012). At the time of writing, this catalog contains more than one hundred million unique articles, crowd-sourced from over two and a half million users.

The users of Mendeley do not only help with building the catalog but also with structuring it. Users can identify themselves as belonging to a scientific discipline and optionally also to a sub-discipline. In August 2012, Mendeley offered 25 disciplines (see Table 1), and 473 sub-disciplines (see Table 2 for the sub-disciplines of “Education”). Each time, a user from a certain (sub-)discipline adds a document to his or her library, the document is automatically assigned to this (sub-)discipline in the catalog.⁴

⁴ As a result, a document can be assigned to more than one (sub-)disciplines.

Table 1

List of the 25 disciplines in the Mendeley catalog.

Arts and Literature	Astronomy/Astrophysics/Space Science
Biological Sciences	Business Administration
Chemistry	Computer and Information Science
Design	Earth Sciences
Economics	Education
Electrical and Electronic Engineering	Engineering
Environmental Sciences	Humanities
Law	Linguistics
Management Science/Operations Research	Materials Science
Mathematics	Medicine
Philosophy	Physics
Psychology	Social Sciences
Sports and Recreation	

Source: <http://www.mendeley.com/research-papers/>.**Table 2**

List of the 18 sub-disciplines of “Education”.

Business Education	Comparative Education
Counselling	Curriculum Studies
Education Research	Educational Administration
Educational Change	Educational Technology
Language Education	Mathematics Education
Medical Education	Miscellaneous
Physical Education	Science Education
Sociology of Education	Special Education
Teacher Education	Testing and Evaluation

Source: <http://www.mendeley.com/disciplines/education/>.

Furthermore, Mendeley Web enables users to create and maintain a user profile that includes their discipline, organization, location, career stage, research interests, biographical information, education, professional experience, contact details, and their own publications.

The following data sets have been sourced on 10 August 2012 and represent data for the sub-discipline educational technology that had been accumulated in the system up to that point:

- User profiles and user libraries: all user profiles and their accompanying user libraries in the sub-discipline of educational technology ($n = 2154$ users).⁵
- Documents: metadata of all documents in the field of educational technology ($n = 144,500$ documents).
- Co-occurrences: co-occurrences of these documents in all Mendeley user libraries ($n = 56,049,431$ co-occurrences).⁶

4. Distribution of subject areas in user libraries

Subject homogeneity is a necessary precondition that the results of co-readership analysis are valid; otherwise the assumption that co-occurrence of articles in user libraries implies subject similarity cannot be upheld. Therefore, we analyzed the subject distribution of articles included in Mendeley user libraries and compared it to the subject area distribution of reference lists of articles in Web of Science. The basis of this analysis is the user profiles and user libraries data set of researchers in educational technology ($n = 2154$ users). As already mentioned, categorization of users into sub-disciplines is determined by self-ascription of users.

In a first step, we analyzed the distribution of journal articles in user libraries. We used SCImago, which is a bibliometric service based on the bibliographic database Scopus, as an external validation source. SCImago categorizes each journal into one of 28 subject areas. The documents from the field of educational technology were matched to these subject areas through the journals they appear in. We used a semi-automated approach for matching journal names in Mendeley and SCImago.⁷

⁵ User profiles and user libraries were sourced at a later point (23 January 2013). Only users that signed up before 10 August 2012 were considered to ensure congruency with the rest of the data set. However, (minor) shifts in the user base cannot be excluded, since in the case users changed their sub-discipline, Mendeley provided only the most recent one chosen.

⁶ Co-occurrence calculation is a computationally intensive process. Therefore, the number of documents per user library was limited to 500. If a user library contained more documents, 500 documents were randomly selected. Then the co-occurrences were calculated.

⁷ Journal names from both sources were transliterated (if necessary) and converted to lowercase. White space at the beginning and at the end was stripped. Colons, commas, and dashes were removed as well as a potential starting definite article “The”. The resulting strings were compared, and all complete matches were taken. In a next step, the list of matches was searched for near-misses and other apparent mismatches, e.g. “User Model-ling and User-Adapted Interaction” as compared to “User Modelling and User-Adapted Interactions”

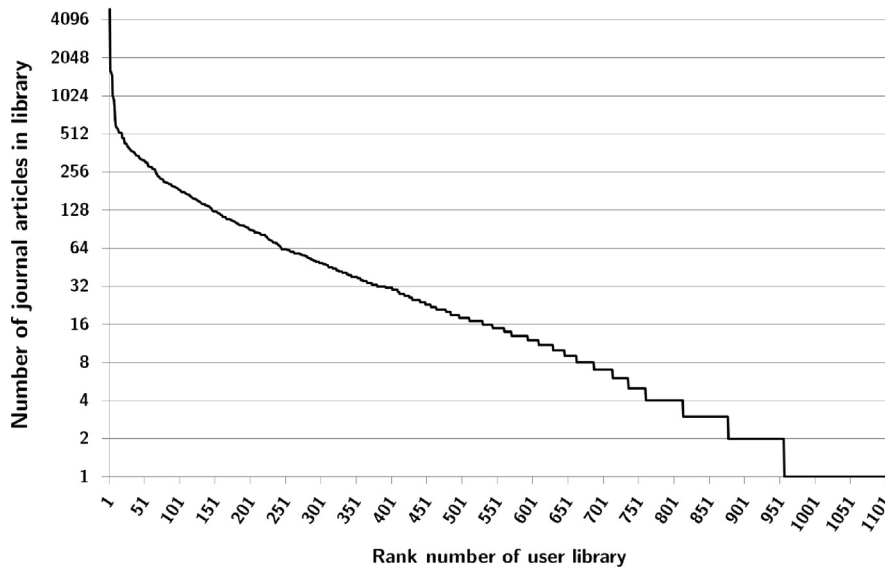


Fig. 4. Distribution of the size of user libraries (no. of journal articles) in educational technology on a logarithmic scale ($n = 72,721$ journal articles in 1107 user libraries).

Table 3

Subject area frequency distribution of articles in user libraries from educational technology ($n = 72,721$ journal articles in 1107 user libraries) and cited references in WoS articles ($n = 13,841$ cited references in 1394 documents). Ranks 11–25 (Mendeley) and 11–12 (WoS articles) were summed up.

Subject area	Mendeley		Web of science	
	Mean	SD	Mean	SD
1	69.19%	22.18%	75.91%	22.12%
2	14.65%	12.94%	15.59%	14.63%
3	6.23%	6.83%	5.24%	7.38%
4	3.59%	4.42%	1.93%	4.06%
5	2.14%	2.95%	0.80%	2.37%
6	1.41%	2.15%	0.35%	1.44%
7	0.97%	1.65%	0.11%	0.70%
8	0.61%	1.17%	0.03%	0.33%
9	0.41%	0.87%	0.02%	0.28%
10	0.29%	0.67%	0.01%	0.20%
>10	0.51%	1.77%	0.01%	0.17%

After this procedure, 1107 user libraries, which contained at least one article in a journal that is indexed by SCImago, were left.

A user library in educational technology has on average 155.7 documents ($SD = 460$, median = 17); slightly more than a third (56.7) of these documents are on average journal articles that appeared in journals indexed by SCImago ($SD = 202.2$, median = 15). As Fig. 4 shows, the distribution of the size of user libraries (number of journal articles) is highly skewed; 10% of all user libraries cover 62.4% of total journal articles.

We also created a data set of cited references from Web of Science. We searched for articles and reviews with the topic “educational technology” in the WoS Core Collection. This resulted in 1,394 documents. We retrieved the cited references for these documents; each document has on average 29.2 cited references ($SD = 23.8$, median = 25). We then applied the procedure outlined above to match references to subject areas via their journals. This resulted in 1221 reference lists which contained at least one document that is indexed by SCImago; 38% of these (11.1 documents) are on average journal articles that appeared in journals indexed by SCImago ($SD = 12.7$, median = 7).

Finally, we calculated the distribution of SCImago categories for each Mendeley user library from educational technology and each cited reference list for the article set retrieved from Web of Science. Afterwards, we ranked the results by subject area. For each library, the percentage of articles that are categorized into a common subject area was calculated. Then, the areas were ranked according to their frequency. The average subject area distribution for all educational technology user libraries can be seen in Fig. 5.

For Mendeley, on average, 69.2% of articles in a user library fall into the top subject area (see Table 3). 14.6% of articles in libraries were on average assigned to the second most frequent second area, while only 6.3% and 3.6% of articles are devoted to the third and fourth ranked subject area respectively. For WoS, 76.0% of cited references in journal articles fall into the top subject area. 15.6% of articles were assigned to the second most frequent area, and 5.2% into the third. In Mendeley, three

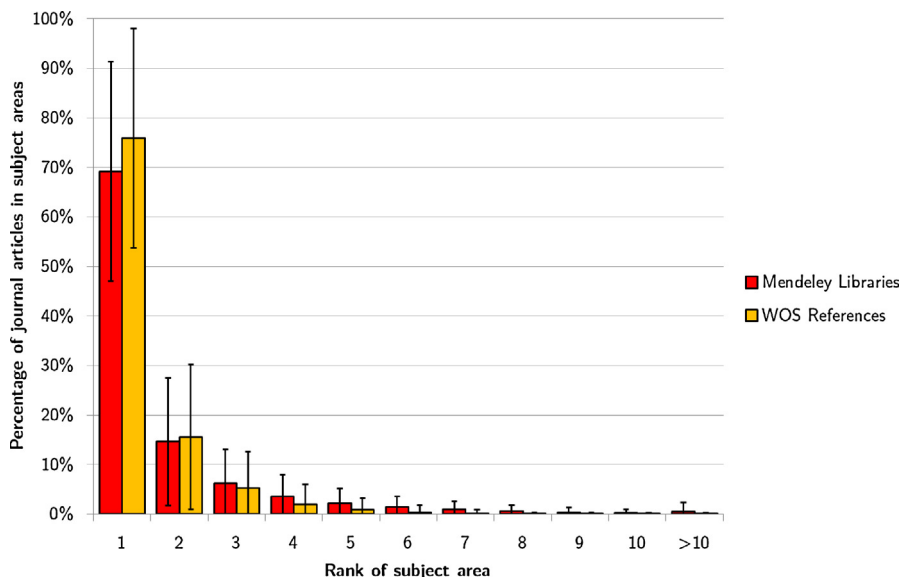


Fig. 5. Subject area frequency distribution of articles in user libraries from educational technology ($n = 72,721$ journal articles in 1107 user libraries) and cited references in WoS articles ($n = 13,841$ cited references in 1,394 documents). Ranks 11–25 (Mendeley) and 11–12 (WoS articles) were summed up.

subject categories account for more than 90% of all articles in an average user library, whereas in journal articles, the top two categories account for more than 90% of all cited references.

These results show that, as was expected, cited references in journal articles are very homogeneous with regards to their subject area distribution. Mendeley user libraries are less homogeneous, and they spread out over more subject areas. The top subject area, however, still accounts for 69.2% of articles in an average user libraries (compared to 76.0% in cited references), even though the number of journal articles in an average user library (56.7) is 5 times higher than the number of cited references in an average journal article (11.2). Therefore, although co-readership probably offers a weaker indication of subject similarity than co-citation, it can still be expected to serve as a useful indication of subject similarity.

5. Visualization of co-readership patterns

For the visualization of co-readership patterns, we followed the knowledge domain visualization process as proposed by Börner, Chen, and Boyack (2003). It consists of four steps: (1) selection of an appropriate data source, (2) determination of the unit of analysis, (3) analysis of the data using dimensionality reduction techniques, and (4) visualization and interaction design. Each of these steps is detailed below. The whole procedure can be seen in Fig. 6.

5.1. Data selection and pre-processing

The documents included in the analysis were taken from the Mendeley sub-discipline of educational technology.⁸ As mentioned before, a document is added to a sub-discipline, if it has at least one reader from this sub-discipline. At the point of data collection, there were approximately 2150 users that had indicated educational technology in their user profile.

To retrieve the most important documents, the document list was sorted by the number of library occurrences within the sub-discipline. A threshold of 16 occurrences was introduced as selection criterion. This means, a document needs to have been added to at least 16 libraries owned by users who identified themselves as being in the field of educational technology to be included in the analysis, leading to a total of 91 documents. The threshold was chosen upon manual inspection of the results. Among the evaluated solutions (thresholds between 11 and 25), the solution with 16 occurrences had the highest purity (0.91).⁹ Since sub-discipline is an optional field in Mendeley, only a minority of users have filled out this field. In order to include more users in Mendeley, the co-occurrence calculation was extended to all user libraries. The 91 documents appeared in 7414 user libraries with a total of 19,402 co-occurrences.

⁸ <http://www.mendeley.com/disciplines/education/educational-technology/>

⁹ Purity is an external cluster evaluation technique. It is defined as the number of correctly assigned documents divided by the number of all documents. A document is correctly assigned when it corresponds to the class that is most frequent within its cluster (Manning, Raghavan, & Schütze, 2009, p. 356f).

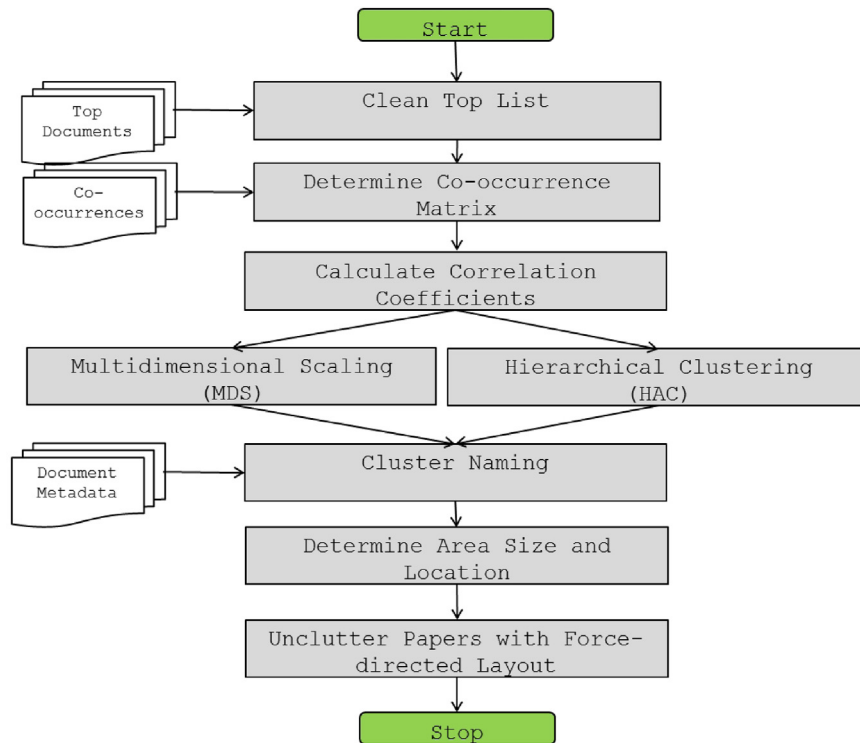


Fig. 6. Overview of the procedure used to create the co-readership visualization.

In a next step, a co-occurrence matrix was created. In line with McCain (1990), diagonal values were treated as missing values. In addition, document pairs with no combined readership were treated as missing values.¹⁰

Based on the co-occurrence matrix, we computed the Pearson correlation coefficient matrix with pairwise complete observations. These correlation coefficients were then used to calculate Euclidean distances between the documents.¹¹

5.2. Clustering and mapping

The matrix of correlation coefficients was the basis for multidimensional scaling (MDS) and hierarchical agglomerative clustering (HAC). Multidimensional scaling was used to project the documents into a two-dimensional space, clustering to find topic areas in the projection. For hierarchical agglomerative clustering, we employed Ward's method (minimum variance) using the R command *hclust*. Ward's method successively merges those two clusters that minimize the increase in the total within-cluster variance (Hair, Black, & Babin, 2010, p. 510). It is known to join smaller clusters and to produce clusters of approximately the same size (Tan, Steinbach, & Kumar, 2006, p. 523).

The number of clusters was determined by the elbow method using the R function *elbow.batch*. This function defines an elbow when the number of clusters k explains at least 80% of the variance in the model, and when the increment is lower than 1% for a bigger k . This criterion was reached at an explained variance of 84% and led to a total of 13 clusters.

In a second step, we plotted the results in a two-dimensional space with non-metric multidimensional scaling (NMDS). NMDS is often employed in scientific mapping efforts. Examples can be found in White and McCain (1998) and in Tsay, Xu, and Wu (2003). NMDS is an iterative approach: beginning with a random start configuration, it tries to minimize a given stress function in consecutive steps. Since NMDS is prone to reaching local minima, usually a number of random starts are used to find an optimum solution.

¹⁰ Usually, these cases are put down as zero co-occurrences. As mentioned above, however, we were limited to a maximum of 500 documents per library when calculating the co-occurrences due to computational constraints. Therefore, we cannot say for sure whether no co-occurrence was found and thus we put down these cases as missing values. We did not find much difference between the two variations, but in the case of missing values, topics were not spread over clusters and the solution was more stable in a bootstrapping analysis. One reason for this could be that the matrix in co-readership analysis is less sparse than in co-citation analysis. Treating document pairs with no combined readership as missing values might therefore serve as a better indicator of discrimination between documents. Therefore, the missing values approach was chosen. Nevertheless, it remains to be determined whether this will hold true for future data sets.

¹¹ Note that the Pearson correlation coefficient is disputed as a measure of subject similarity (Ahlgren, Jarneving, & Rousseau, 2003). For a discussion of alternate similarity measures see e.g. Egghe (2010).

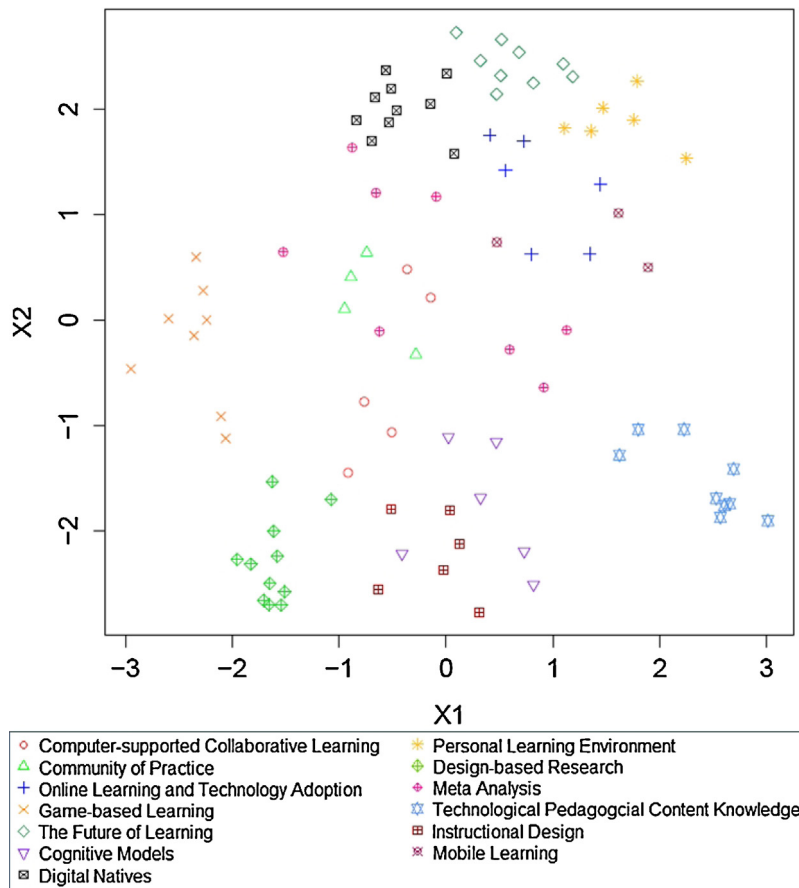


Fig. 7. Result of NMDS, HAC, and the naming algorithm. Each symbol represents a document, all of the documents with the same symbol constitute a topic area.

We selected the NMDS implementation provided by the R *ecodist* package (Goslee & Urban, 2007). In the NMDS, stress is reported as 0.2 and the R^2 is reported as 0.86. According to Hair et al. (2010), acceptable results for R^2 start at 0.60.

To create labels for the clusters, titles and abstracts of the documents in each cluster were submitted to the APIs of Zemanta¹² and OpenCalais.¹³ Both services crawl the semantic web and return a number of concepts that describe the content. The returned concepts were compared to word n -grams generated from titles and abstracts. The more words a concept has (and therefore, the more information it contains), and the more often it occurs within the text, the more likely it is to be the label of the cluster. The results of this procedure were manually checked and corrected if needed.

A plot of the results from the procedure described above can be seen in Fig. 7. Each symbol represents a document. The type of symbol signifies the topic area it belongs to. These 13 areas are listed in the legend below the graph.

5.3. Web visualization

In order to allow users to interact with this graph, we developed an interactive web visualization prototype. The visualization was created using D3.js.¹⁴ In the prototype, documents are represented as rectangles with dogears, a common metaphor, used in many icons and graphics. The size of the document signifies the number of readers it has. To avoid coding the documents with symbols (as in Fig. 7), topic areas are represented as bubbles. The center of each bubble is calculated as the mean of the coordinates of the publications based on the NMDS result. The size of the bubble is determined by the number of combined readers of the publications in the topic area.

¹² <http://zemanta.com>

¹³ <http://opencalais.com>

¹⁴ <http://d3js.org>

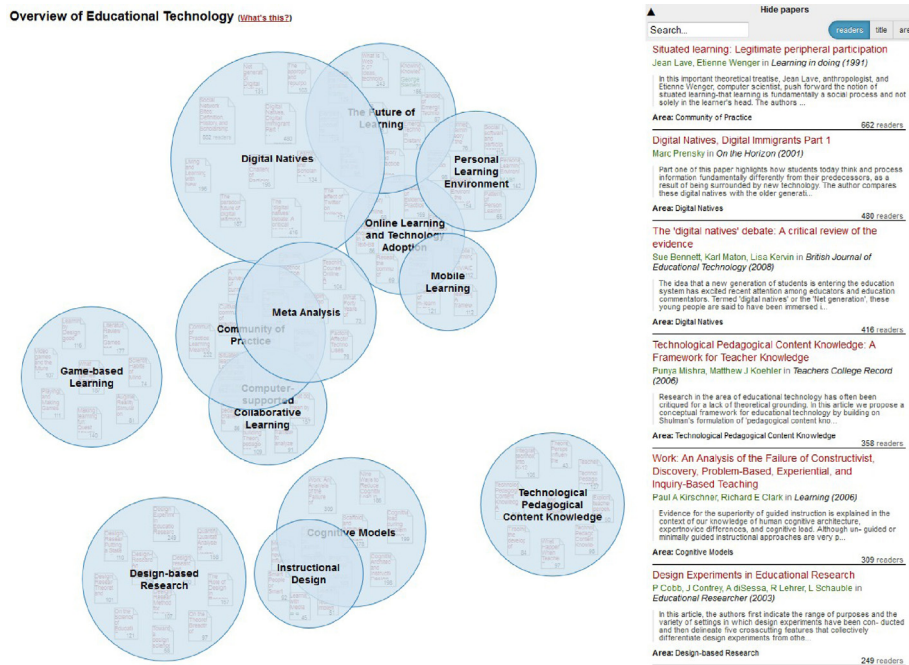


Fig. 8. Knowledge domain visualization of educational technology. The bubbles represent topic areas within the domain. The size of a bubble relates to the number of combined readers.

Additionally, force-directed placement (Fruchterman & Reingold, 1991) was employed on the documents to unclutter the visualization and move documents into their respective areas.¹⁵ To prevent overlapping documents, the collision detection algorithm by Mike Bostock¹⁶ was used.

It is important to note that – in contrast to the topic areas – the relative closeness of documents in the visualization does not necessarily imply relative subject similarity.¹⁷ To review the relationship between individual papers, one needs to go back to the original output of the MDS shown in Fig. 7.

5.4. Results

The resulting visualization prototype, which can be accessed on Mendeley Labs,¹⁸ is shown in Fig. 8. In the first few seconds of the visualization, the force-directed placement algorithm is executed. The papers are untangled and pulled into their respective areas, represented by the blue bubbles. After the force-directed algorithm has finished, users can interact with the visualization.

5.4.1. Topic area description and distribution

As can be seen in Fig. 8, there are 13 topic areas in the visualization with a combined readership of 13,630 at the time of data collection (10 August 2012). Table 4 gives an overview of the topic areas. It shows that they differ in terms of the number of documents and the number of readers. *Digital Natives* has the highest readership with over 20% of all readers. It has twice the readership of the second largest area: *Design-based Research* (DBR). DBR includes the most documents (11) of all areas. *Community of Practice* has only four documents, but still sports the fourth most readers. The area with the least readers and the least number of documents is *Mobile Learning* with just 3 documents and a combined readership of 345.

The topic areas can again be assigned to meta-areas. These meta-areas are formed by areas that are close to each other, as is assumed by multidimensional scaling. On the top of the map (see Fig. 8), social and technological developments are being discussed (in *Digital Natives* and *The Future of Learning*). Beneath, there is a large cluster of learning methods and technologies, spanning *Mobile Learning*, *Personal Learning Environment*, *Online Learning and Technology Adoption*, *Community of Practice*, and *Game-based Learning*. On the bottom of the visualization, there is a cluster of topic areas that form the psychological,

¹⁵ The area centers were denoted as gravitational centers. Documents not within the limits of the topic area were instructed to move toward the gravitational center. Edges and corresponding edge weights were not set; they are therefore initialized to default values by D3.

¹⁶ <http://bl.ocks.org/mbostock/3231298>

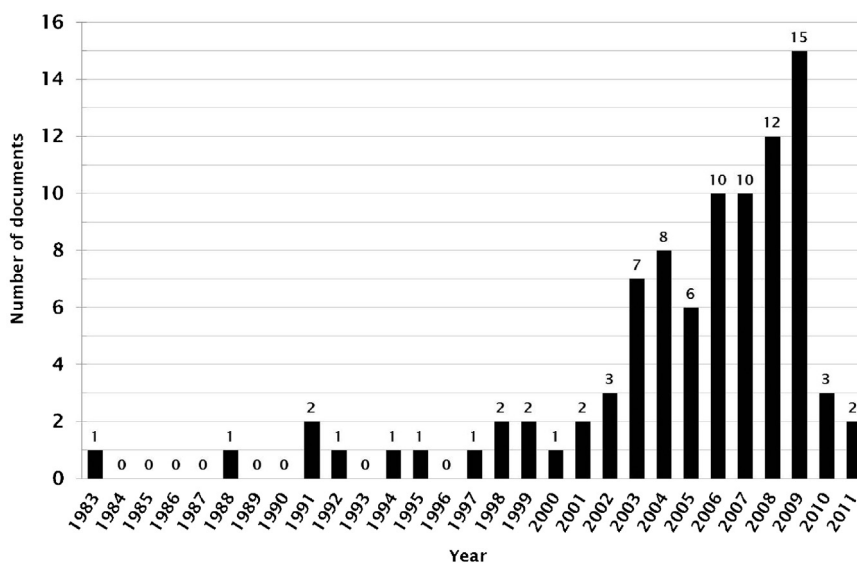
¹⁷ In the uncluttering effort using force-directed placement, the positions of documents are changed in a way that does not necessarily preserve the relative distances. Therefore, the distances between documents in the visualization do not represent the distances calculated with MDS anymore.

¹⁸ <http://labs.mendeley.com/headstart>. The source code can be obtained from <https://knowminer.at/svn/opensource/other-licenses/iglp.v3/headstart/>.

Table 4

Topic areas in the visualization.

Topic area	No. documents	No. readers	% readership
Digital Natives	10	2865	21.0%
Design-based Research	11	1477	10.8%
The Future of Learning	9	1183	8.7%
Community of Practice	4	1175	8.6%
Cognitive Models	6	1169	8.6%
Technological Pedagogical Content Knowledge	9	1049	7.7%
Game-based Learning	8	993	7.3%
Meta Analysis	8	991	7.3%
Personal Learning Environment	6	648	4.8%
Online Learning and Technology Adoption	6	637	4.7%
Computer-supported Collaborative Learning	5	615	4.5%
Instructional Design	6	483	3.5%
Mobile Learning	3	345	2.5%
Sum	91	13,630	100.0%

**Fig. 9.** Distribution of publication years of documents in the visualization ($n=91$).

pedagogical, and methodological foundations of the field. The areas *Computer-supported Collaborative Learning*, *Instructional Design* and *Cognition* relate to psychology, while *Technological Pedagogical Content Knowledge* relates to pedagogy. Research methods are represented by *Design-based Research*.

From what was mentioned above, it follows that pedagogical and psychological topics are covered very well in the visualization. However, topic areas that are largely influenced by computer science such as *Adaptive Hypermedia* or knowledge management (e.g. *Work-integrated Learning*) are missing from the overview. The reason for this is most likely the discipline taxonomy in Mendeley (see Section 6.2).

5.4.2. Publication types and age of publications

The 91 documents in the visualization can be divided into five different types of publications. The majority are journal articles (71 items, or 78%), followed by reports (7), books (6) and book chapters (5), and conference papers (2). The 71 journal articles were published in a variety of journals. The highest number of articles was published in “Computers & Education” (8), followed by “Educational Technology Research & Development” and “The Internet and Higher Education” (both 6) and “Review of Educational Research”, “Educational Researcher” and “Educational Psychologist” (all 5). These publication outlets are among the highest impact journals in the Journal Citation Reports (Thomson Reuters, 2013). All of the documents in the visualization are in English.

Fig. 9 shows the age distribution of the 91 publications covered in the visualization. 80% of publications were published from 2003 onwards, meaning that they were younger than ten years at the time of data collection (10 August 2012). Most

documents were published in 2009. The median age of publications is 6.0 years (mean = 7.3 years).¹⁹ The relative small amount of publications from 2010 and 2011 can be explained by the circumstance that it is more difficult for recent publications to reach the threshold value than for older ones.

6. Discussion

6.1. Recency

In the conducted co-readership analysis, the mean age of publications is 7.3 years with 80% of articles published within 10 years of data collection. While this constitutes a contemporary selection of publications, the relative low proportion of articles younger than two years indicates that not all of the latest developments might be represented in the visualization. However, in a comparable co-citation mapping effort in educational technology by [Cho, Park, Jo, and Suh \(2012\)](#), the mean age of papers was 14.1 years (median = 14 years) which is almost double the age of publications in the co-readership analysis. In addition, only 18% of the 28 papers included in the co-citation analysis were less than 10 years of age.²⁰

This suggests that the results of a co-readership analysis may be much more up-to-date than co-citation analysis. In contrast to bibliographic coupling, however, there is still a certain time lag after publication that needs to be taken into account. Therefore, a co-readership analysis may be most appropriate when a contemporary overview is desired but a dynamic method is preferred over a static one.

6.2. Biases in the visualization

An analysis of the results shows that the visualization is not free from biases. First, all of the papers are in English, even though educational technology is often researched by local communities that communicate in their native language ([Ely, 2008](#)). Second, the knowledge domain visualization represents an education-dominated view that lacks topic areas related to computer science.

Biases in usage statistics analyses were first mentioned by [Bollen and Sompel \(2008\)](#) in a study of downloads in an institutional repository. The authors found great differences in the correlation of usage impact factor and journal impact factor depending on the user base. The authors therefore concluded that these biases occur due to sample characteristics.

Consequently, we looked into the sample characteristics of users in educational technology that we investigated based on their user profiles ($n = 2153$ user profiles). At first, we analyzed the geographical distribution of users. One of the reasons for the fact that all of the papers are in English is surely that English is the *lingua franca* in science and research ([Tardy, 2004](#)). But most likely, this dominance of English also stems from the fact that there is a strong bias toward English-speaking countries on Mendeley.

This assumption is backed up by the results of the geographical analysis (see [Fig. 10](#)). Out of 2153 users, 927 (43.1%) have chosen to list a country in their user profile. In total, 70 countries have been specified, but the distribution is highly skewed. There is an emphasis on the US and the UK with a total number of 423 users (45.6%). In fact, when adding Canada and Australia, English-speaking countries have a share of over 54.3%. 56 countries with a low share of users have been summed up under "Other" (21.7%). This shows that Mendeley users come from a wide variety of countries, but that there is a strong focus on English speaking countries.

[Fig. 10](#) shows the comparison of this distribution to the geographic distribution of educational technology authors in the Web of Science Core Collection and the distribution of researchers according to [The World Bank \(2014\)](#). Of the 2965 unique authors with an assigned geographic location that have contributed to an article with the topic "educational technology" (out of 4602 in total), 1298 (43.8%) come from one of the four major English-speaking countries. Although this proportion is very high (for instance, the share of researchers from these 4 countries is 24.8% according to the World Bank), the dominance of Mendeley users from these 4 countries is even stronger. Two facts play an important role with regards to this imbalance: first, Mendeley originated in the UK and has an office in the USA. Second, the Mendeley software is only available in English for now.

The bias toward disciplines strongly related to education can be explained by Mendeley's discipline taxonomy which was used to determine the paper pre-selection in this study. Even though educational technology is an interdisciplinary field, it appears solely as a sub-discipline of education. The sign-up process in Mendeley requires a user to first select a discipline such as education, social science, or computer and information science. In a second step, a user can select a sub-discipline, such as educational technology. Therefore, a scholar in educational technology with a background in computer science will conclude after the first step that his or her sub-discipline is not represented in Mendeley and choose another one.

¹⁹ Calculation based on the exact date of data collection on 10/08/2012.

²⁰ All calculations based on the publication year of the most recent article in the analysis (2011).

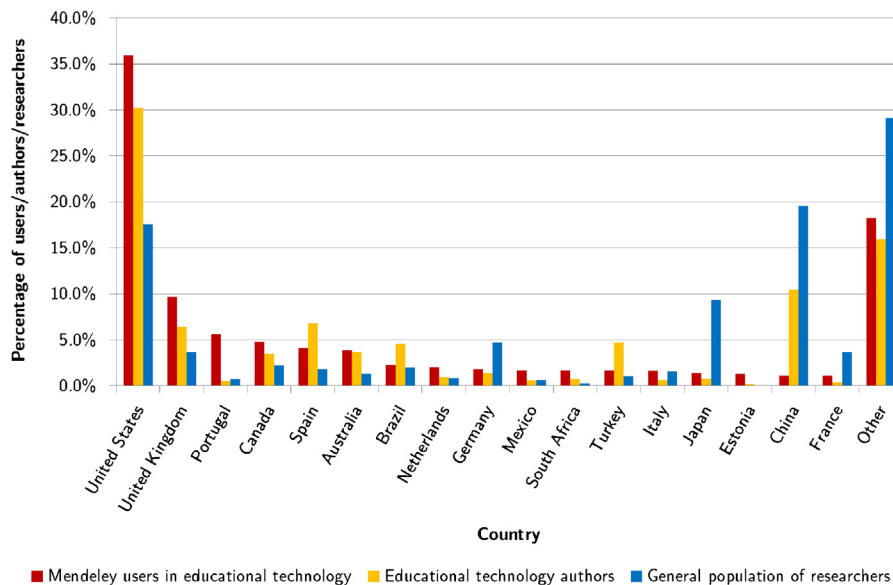


Fig. 10. Geographic distribution of users from educational technology in Mendeley ($n = 927$ users), educational technology authors in WoS ($n = 2965$ authors) and researchers in general ($n = 7,043,655$ researchers). Only countries that were present in all three data sets were taken into account when calculating the distribution. Data sources: Mendeley/Web of Science Core Collection/World Bank Development Indicators (researchers in R&D (per million people); population, total).

7. Conclusions and future work

In this paper, we analyze the adequacy and applicability of readership statistics recorded in social reference management systems for creating knowledge domain visualizations. We propose co-readership as a measure of subject similarity. An analysis of the distribution of subject areas in user libraries of educational technology researchers on Mendeley shows that 69.2% of the journal articles in an average user library can be attributed to a single subject area. This is in line with an earlier study by [Jiang et al. \(2011\)](#) which finds that clusters based on the occurrence and co-occurrence of articles in user libraries of CiteULike are as effective as citation-based clusters.

The prototypical visualization based on co-readership patterns of the field of educational technology comprises of 13 topic areas, which can be aggregated to meta-clusters, therefore strengthening the assumption that co-readership indicates subject similarity. The visualization is a recent representation of the field: 80% of the publications included are from within ten years of data collection. However, not all of the latest developments were represented in the visualization due to the fact that it is harder to reach threshold values for the most recent publications. Nevertheless, the papers included in the co-readership analysis are on average almost half as young as the papers included in a comparable co-citation analysis by [Cho et al. \(2012\)](#). This suggests that co-readership analysis may be able to represent more recent aspects than co-citation. In order to generalize this statement and to better understand the differences between co-citation analysis, bibliographic coupling, and co-readership analysis, however, comparison studies between the different similarity measures must be carried out.

The characteristics of the readers introduce certain biases to the visualization. All scientometric analyses are subject to bias; it is therefore important that the characteristics of the underlying sample are made transparent. In the co-readership analysis, information encoded in the user profiles can be used to explain these characteristics. In the present study, a majority of readers were self-ascribed to the field of education and they came from an English-speaking country. This resulted in a map that represents an education science-dominated view from mainly an Anglo-American perspective.

One of the limitations of this work is that the methodology has only been tested for a single field of research. Educational technology is a diverse field with many influences; but it would not be appropriate to generalize the results to all research fields. The question is whether the same analysis would work as well on a larger set of publications and for other fields and disciplines. Each discipline has its own theories, methods, accepted practices, in short: its own culture. Just like publication and citation practices are fundamentally different for the natural sciences and the humanities, cultural differences might also affect the usage of social reference management systems. In the future, this study must therefore be repeated in other fields of research. This could be especially interesting for those fields that are dynamic in nature, and those that have not been scientometrically analyzed before due to the lack of citation data.

When applied to larger collections of documents, the procedure used in this paper may be problematic. Both hierarchical clustering and multidimensional scaling have a high computational complexity. Therefore, it will be important to investigate algorithms that can deal with large data sets such as force-directed layout for ordination, and community detection for the establishment of topic areas. For a further discussion see [Fortunato \(2010\)](#) and [Gibson, Faith, and Vickers \(2012\)](#).

Finally, it seems promising to harness information encoded in the user profiles, such as location, discipline, and career stage, not only for a better understanding of the results (see above), but also for filtering the visualization. This would make it possible to compare visualizations, for instance between countries or career stages. Furthermore, with the availability of timestamps, it becomes possible to show the evolution of a research field over time at a granular level of detail.

Acknowledgements

We would like to thank the reviewers of this paper for their comprehensive comments and suggestions which have improved the paper considerably. The research presented in this work is in part funded by the European Commission as part of the FP7 Marie Curie IAPP project TEAM (Grant No. 251514). The Know-Center is funded within the Austrian COMET program – Competence Centers for Excellent Technologies – under the auspices of the Austrian Federal Ministry of Transport, Innovation and Technology, the Austrian Federal Ministry of Economy, Family and Youth, and the State of Styria. COMET is managed by the Austrian Research Promotion Agency FFG.

References

- Ahlgren, P., Jarneving, B., & Rousseau, R. (2003). Requirements for a cocitation similarity measure with special reference to Pearson's correlation coefficient. *Journal of the American Society for Information Science and Technology*, 54(April (6)), 550–560. <http://dx.doi.org/10.1002/asi.10242>
- Amin, M., & Mabe, M. (2000). Impact factors: Use and abuse. *Perspectives in Publishing*, 1(January (2000)), 1–6. <http://www.ncbi.nlm.nih.gov/pubmed/14518149>
- Bar-Ilan, J., Haustein, S., Peters, I., Priem, J., Shema, H., & Terliesner, J. (2012). Beyond citations: Scholars' visibility on the social Web. In *17th international conference on science and technology indicators* (pp. 1–14). arxiv:1205.5611.
- Bollen, J., Rodriguez, M., & Van de Sompel, H. (2007). MESUR: Usage-based metrics of scholarly impact. In *Proceedings of the 7th ACM/IEEE-CS joint conference on digital libraries* (p. 474). ACM. <http://dl.acm.org/citation.cfm?id=1255273>
- Bollen, J., & Sompel, H. V. D. (2008). Usage impact factor: The effects of sample characteristics on usage-based impact metrics. *Journal of the American Society for Information Science*, 59(1998), 136–149.
- Bollen, J., Sompel, H. V. D., Smith, J., & Luce, R. (2005). Toward alternative metrics of journal impact: A comparison of download and citation data. *Information Processing and Management*, 41(6), 1419–1440. pii:S0306457305000324.
- Bollen, J., & van de Sompel, H. (2006). Mapping the structure of science through usage. *Scientometrics*, 69(November (2)), 227–258. <http://dx.doi.org/10.1007/s11192-006-0151-8>
- Bollen, J., Van de Sompel, H., Hagberg, A., Bettencourt, L., Chute, R., Rodriguez, M. A., et al. (2009). Clickstream data yields high-resolution maps of science. *PLoS ONE*, 4(March (3)), e4803.
- Bollen, J., Van de Sompel, H., Hagberg, A., Bettencourt, L., Chute, R., Rodriguez, M. A., et al. (2008). A clickstream map of science. In K. Börner, & E. F. Hardy (Eds.), *Places & spaces: Mapping science* <http://scimaps.org>
- Börner, K., Chen, C., & Boyack, K. W. (2003). Visualizing knowledge domains. *Annual Review of Information Science and Technology*, 37(January (1)), 179–255.
- Boyack, K., Klavans, R., & Börner, K. (2005). Mapping the backbone of science. *Scientometrics*, 64(3), 351–374.
- Boyack, K. W., & Klavans, R. (2010). Co-citation analysis bibliographic coupling and direct citation: Which citation approach represents the research front most accurately? *Journal of the American Society for Information Science and Technology*, 61(December (12)), 2389–2404. <http://dx.doi.org/10.1002/asi.21419>
- Brody, T., Harnad, S., & Carr, L. (2006). Earlier web usage statistics as predictors of later citation impact. *Journal of the American Society for Information Science and Technology*, 57(8), 1060–1072. <http://dx.doi.org/10.1002/asi.20373/full>
- Chen, C., & Carr, L. (1999). Trailblazing the literature of hypertext: Author co-citation analysis (1989–1998). In *Proceedings of the tenth ACM conference on hypertext and hypermedia* (pp. 51–60). New York: ACM.
- Chen, L.-C., & Lien, Y.-H. (2011). Using author co-citation analysis to examine the intellectual structure of e-learning: A MIS perspective. *Scientometrics*, 89(July (3)), 867–886.
- Cho, Y., Park, S., Jo, S. J., & Suh, S. (2012, August). The landscape of educational technology viewed from the ETR&D journal. *British Journal of Educational Technology*.
- Czerniewicz, L. (2010 Dec). Educational technology – Mapping the terrain with Bernstein as cartographer. *Journal of Computer Assisted Learning*, 26(6), 523–534. <http://dx.doi.org/10.1111/j.1365-2729.2010.00359.x>
- Darmoni, S., Roussel, F., Benichou, J., Thirion, B., & Pinhas, N. (2002). Reading factor: A new bibliometric criterion for managing digital libraries. *Journal of the Medical Library Association*, 90(3), 323–327. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC116406/>
- Egghe, L. (2010). Good properties of similarity measures and their complementarity. *Journal of the American Society for Information Science and Technology*, 61(October (10)), 2151–2160. <http://dx.doi.org/10.1002/asi.21380>
- Egghe, L., & Rousseau, R. (1990). *Introduction to informetrics. Quantitative methods in library, documentation and information science*. Amsterdam: Elsevier.
- Ely, D. (2008). Frameworks of educational technology. *British Journal of Educational Technology*, 39(March (2)), 244–250. <http://dx.doi.org/10.1111/j.1467-8535.2008.00810.x>
- Fortunato, S. (2010). Community detection in graphs. *Physics Reports*, 486(February (3–5)), 75–174. <http://linkinghub.elsevier.com/retrieve/pii/S0370157309002841>
- Fruchterman, T., & Reingold, E. (1991). Graph drawing by force-directed placement. *Software: Practice and Experience*, 21(11), 1129–1164.
- Garfield, E. (2001). From bibliographic coupling to co-citation analysis via algorithmic historio-bibliography. Speech delivered at Drexel University, Philadelphia, PA, November, 27.
- Gibson, H., Faith, J., & Vickers, P. (2012). A survey of two-dimensional graph layout techniques for information visualisation. *Information Visualization*, 12(September (3–4)), 324–357. <http://dx.doi.org/10.1177/1473871612455749>
- Goslee, S., & Urban, D. (2007). The ecodist package for dissimilarity-based analysis of ecological data. *Journal of Statistical Software*, 22(7), 1–19. <http://www.jstatsoft.org/v22/i07/paper/>
- Hair, J., Black, W., & Babin, B. (2010). *Multivariate data analysis*. Upper Saddle River, NJ: Pearson Prentice Hall.
- Hammerton, J. A., Granitzer, M., Harvey, D., Hristakeva, M., & Jack, K. (2012). On generating large-scale ground truth datasets for the duplication of bibliographic records. In *Proceedings of the 2nd international conference on web intelligence, mining and semantics* <http://mics.fim.uni-passau.de/wp-content/papercite-data/pdf/hammerton/2012.pdf>
- Haustein, S., & Larivière, V. (2014). Mendeley as a source of readership by students and postdocs? Evaluating article usage by academic status. In *Proceedings of the 35th IATUL conference* <http://docs.lib.purdue.edu/iatul/2014/altmetrics/2/>
- Haustein, S., & Siebenlist, T. (2011). Applying social bookmarking data to evaluate journal usage. *Journal of Informetrics*, 5(3), 446–457.
- Henning, V., & Reichelt, J. (2008). Mendeley – A Last.fm for research? In *IEEE fourth international conference on eScience* (pp. 327–328).
- Jiang, J., He, D., & Ni, C. (2011). Social reference: Aggregating online usage of scientific literature in CiteULike for clustering academic resources. In *Proceeding of the 11th annual international ACM/IEEE joint conference on digital libraries* (pp. 401–402). ACM.

- Kessler, M. M. (1963). Bibliographic coupling between scientific papers. *American Documentation*, 14(1), 10–25. <http://dx.doi.org/10.1002/asi.5090140103/abstract>
- Kraker, P., Körner, C., Jack, K., & Granitzer, M. (2012). Harnessing user library statistics for research evaluation and knowledge domain visualization. In *Proceedings of the 21st international conference companion on world wide web* (pp. 1017–1024). Lyon: ACM.
- Kurtz, M. J., Eichhorn, G., Accomazzi, A., Grant, C., Demleitner, M., Murray, S. S., et al. (2005). The bibliometric properties of article readership information. *Journal of the American Society for Information Science and Technology*, 56(January (2)), 111–128.
- Manning, C. D., Raghavan, P., & Schütze, H. (2009). *An introduction to information retrieval*. No. c. Cambridge University Press.
- Marshakova, I. (1973). System of document connections based on references. *Scientific and Technical Information Serial of VINITI*, 6(2), 3–8.
- McCain, K. W. (1990). Mapping authors in intellectual space: A technical overview. *Journal of the American Society for Information Science*, 41(6), 433–443.
- Polanco, X., Ivana, R., & Dominique, B. (2006). User science indicators in the Web context and co-usage analysis. *Scientometrics*, 66(January (1)), 171–182.
- Priem, J., & Hemminger, B. M. B. (2010). Scientometrics 2.0: Toward new metrics of scholarly impact on the social Web. *First Monday*, 15(7) <http://firstmonday.org/ojs/index.php/fm/article/view/2874/2570>
- Rowlands, I., & Nicholas, D. (2007). The missing link: Journal usage metrics. *Aslib Proceedings*, 59(3), 222–228.
- Schloegl, C., & Gorraiz, J. (2010). Comparison of citation and usage indicators: The case of oncology journals. *Scientometrics*, 82(February (3)), 567–580. <http://dx.doi.org/10.1007/s11192-010-0172-1>
- Schlögl, C. (2001). *Bestandsaufnahme Informationsmanagement: Eine szientometrische, qualitative und empirische Analyse*. Wiesbaden: Gabler.
- Schlögl, C., Gorraiz, J., Gumpendorfer, C., Jack, K., Kraker, P., & Schloegl, C. (2013). Download vs. citation vs. readership data: The case of an information systems journal. In *14th international society of scientometrics and informetrics conference* (pp. 626–634).
- Siemens, G., & Tittenberger, P. (2009). *Handbook of emerging technologies for learning*. Accessed 06.10.13.
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24(4), 265–269.
- Small, H. (1999). Visualizing science by citation mapping. *Journal of the American Society for Information Science*, 50(9), 799–813.
- Tan, P.-N., Steinbach, M., & Kumar, V. (2006). *Introduction to data mining*. Boston: Addison-Wesley.
- Tardy, C. (2004). The role of English in scientific communication: Lingua franca or Tyrannosaurus rex? *Journal of English for Academic Purposes*, 3(3), 247–269.
- The World Bank. (2014). *World development indicators*. <http://databank.worldbank.org>
- Thelwall, M., & Maflahi, N. (2014, July). Are scholarly articles disproportionately read in their own country? An analysis of Mendeley readers. *Journal of the Association for Information Science and Technology*, <http://dx.doi.org/10.1002/asi.23252>
- Thomson Reuters. (2013). *Journal citation reports 2012*. Tech. rep.
- Tsay, M.-Y., Xu, H., & Wu, C.-W. (2003). Journal co-citation analysis of semiconductor literature. *Scientometrics*, 57(1), 7–25.
- White, H. D. H., & McCain, K. K. W. (1998). Visualizing a discipline: An author co-citation analysis of information science, 1972–1995. *Journal of the American Society for Information Science*, 49(4), 327–355.
- Zahedi, Z., Costas, R., & Wouters, P. (2014, March). How well developed are altmetrics? A cross-disciplinary analysis of the presence of 'alternative metrics' in scientific publications. *Scientometrics*, <http://dx.doi.org/10.1007/s11192-014-1264-0>