

# Discovering shared interests through co-outlinking in a municipal web space

## Abstract

In webometric research interlinking and co-inlinking has been studied extensively, but co-outlinking has been used as a data collection and research method only in a couple of studies before. Yet it may have some potential as a tool for mapping shared interests of organizations creating the outgoing links. This research studied this potential of co-outlinking and mapped the shared interests of 54 municipalities in Finland. The research showed that co-outlinking can in fact be used to map shared interests and perhaps even cooperation but the method required a lot of manual cleansing and classification of the data. Although the results were promising co-outlinking cannot be recommended unless some improvements are developed for classification of the link data.

## Introduction

Webometric research has mostly focused on measuring link counts and mapping interlinking patterns between universities (e.g. Kim, 2000; Bar-Ilan, 2004b; Thelwall, 2004; 2006; Stuart, Thelwall & Harries, 2007; Ortega et al., 2008). Some studies have been made on co-inlinking (e.g. Dean & Henzinger, 1999; Vaughan, Kipp & Gao, 2007; Ortega et al., 2008; Holmberg, 2010), but co-outlinking has not been extensively studied before. Interlinking has been shown to reflect existing cooperation and other relationships between the researched organizations (e.g. Wilkinson, et al., 2003) and co-inlinking has been shown to e.g. be a useful tool to map competitive positions between companies (e.g. Vaughan & You, 2005) or to map cooperation between municipalities (Holmberg, 2010). Co-outlinking may reveal another aspect or a relationship between the investigated websites and organizations. Because linking can be seen as referencing or showing awareness or interest towards the target website and organization, co-outlinking could be assumed to show how organizations share various interests towards the organizations they link to. Shared interests would be interesting and valuable to map as they may potentially lead to cooperation and other stronger relationships in the future. An example of such a scenario could be two municipalities sharing an interest for recycling and waste management. If these shared interests become realized the municipalities could start cooperating in the matter. Co-outlinking however remains a relatively unexplored webometric measure, something that this research aims to change.

## Literature review

Kessler (1963) suggested *bibliographic coupling* as a similarity measure by. The assumption with bibliographic coupling is that if two articles both cite the same third one, they are similar because they have a shared interest towards it. The webometric measure co-outlinking is based on bibliographic coupling. Co-outlinking is based on the assumption that two documents on the web are similar or that they have something in common if they both link to a third document. There is something in the target document or website that both of the source pages, or source organizations, have found interesting or valuable enough to link to. They have a shared interest towards the target website and the content of the target website may reveal something about what that shared interest is about.

There are only a few earlier webometric studies on co-outlinking. Björneborn (2001) studied the possibilities of using a combination or a sequence of co-inlinks and co-outlinks to create so called co-linkage chains that could possibly function as short cuts between different topical communities on the web, and hence, shorten the path length between any two pages on the web. Thelwall and Wilkinson (2004) studied interlinking, co-inlinking and co-outlinking as methods to locate academic websites with similar content. The study showed that combining the three different measures increased the probability for identifying similar sites, but a

combination of the three was only a small improvement against using interlinking alone. García-Santiago and de Moya-Anegón (2009) used co-outlinks to map triple-helix networks in the Spanish society. Using co-outlinking links gave them the opportunity to study relationships between the organizations even though the organizations did not have direct connections between them on the web. Co-outlinking provides an indirect method to analyze connections and relationships between organizations on the web.

*Classifying linking motivations* is important to gain knowledge about the reasons behind link creation. Links can be created for many different reasons and unless these reasons are studied nothing can be said about any trends or patterns that may be discovered from linking networks. There is no uniform classification scheme that could be used for all webometric research and all type of content on the web. So far, researchers have created their own classification schemes that have best suited the goals of their current research. There are almost as many different classification schemes as there are studies about the subject (e.g. Haas & Grams, 1998; Kim, 2000; Wilkinson, et al., 2003; Bar-Ilan, 2004b; 2005; Stuart, Thelwall & Harries, 2007).

Different studies have looked at slightly different aspects of linking and link creation, but earlier studies have been interested in who links to whom and why. To find answers to these questions various aspects of the links and linking have been researched. Haas and Grams (1998) classified link types based on the context the link was created in, the link itself, and reasons why the link had been created. Bar-Ilan (2004a) categorized links according to the intention of the link by looking at the surroundings of the link on the source page. Thelwall, Harries and Wilkinson (2003) classified the type of the links based on the content of the page where the link was created and the context of the link. Smith (2004) classified links using the source and target pages, and any apparent reasons for linking, while Thelwall and Harries (2003) classified only target pages. Stuart, Thelwall and Harries (2005) classified links based on the source pages and reasons for link placement. Vaughan, Gao and Kipp (2006) studied link creation motivations by classifying the links according to the content of the source page and the context in which the link had been created. Bar-Ilan (2004b) proposed a classification of link types, in which she included many different aspects of both source and target pages and the context in which the links were created. This categorization scheme was later further refined and developed by Bar-Ilan (2005). These earlier studies have all mainly concentrated on two aspects of linking: 1) content and context of the source and the target pages and 2) the reasons for creating the link.

### **Research questions**

The goal of this research is to study if co-outlinking from municipal websites could be used to map shared interests of the municipalities. This research will also try to map what these shared interests are. The goal of this research is summarized in the following two research questions:

1. Can co-outlinking be used to map shared interests of the municipalities in the region of Finland Proper?
2. What are these shared interests?

By answering the research questions above this study will give some new information about the relationships between the municipalities in the region of Finland Proper, and also help develop new research methods for webometric research. If successful, co-outlinking could

prove to be a method to discover some other aspects of relationships between organizations on the web than those possible to discover from interlinking or co-inlinking.

## **Methods**

The research studied co-outlinking from the 54 municipalities in the region of Finland Proper. The municipalities varied greatly in size, with Turku being the largest city in the region with over 170,000 residents and with municipalities like Iniö and Velkua with less than 250 residents. At the time of data collection Finland was in the middle of municipal merges. The goal with the merges was to create larger municipalities and with that to remove excessive operations and create more efficient local administration.

To answer the research questions the link data was converted into a binary matrix and visualized as a network graph. The graphs were visually analyzed and statistically compared with another matrix that was created based on existing cooperation between the municipalities. Schneider and Borlund (2007) state that the choice of an appropriate measure to compare and analyze matrices should always be done based on the material and the research goals, and that there may in fact not be a single measure that would in every case be “the best”. They continue by suggesting that in every case the results should be confirmed with some complementary techniques, such as content analysis or link classifications in webometrics. Only by studying the motivations for creating the studied links can we say something about the networks they create and the relationships that are discovered.

Precision and recall are basic operations for evaluating information retrieval systems like web search engines (e.g. Harter & Hert, 1997; Clarke & Willett, 1997), but the concept could also be used to compare how similar two binary squared matrices are to each other. When comparing binary matrices precision would be a measure for how many links in the first matrix would match with links in the second matrix. To get a high precision most of the links in the first matrix should match with links in the second matrix, but precision can be high even if there are several links in the second matrix that do not match with links from the first matrix. Recall is a measure of comprehensiveness of the match. Recall is a measure for how many of the links in the second matrix actually matched links in the first one. This can be written as a matrix operation where precision equals  $M_{11}/(M_{11}+M_{10})$  and recall equals  $M_{11}/(M_{11}+M_{01})$ . If for example a Matrix A would have 10 links and a Matrix B would have 100 links and all the links in Matrix A would match links in Matrix B, then the precision would be 100%. This would seem like a good match, however, 90 of the links in Matrix B were not matched by links in Matrix A. The recall for this match would be 10%, clearly showing that the match between the two matrices was not as good as it first seemed. From these a precision/recall ratio can be calculated, but because a high precision and low recall would give the same ratio as medium precision and medium recall, the difference between precision and recall also has to be taken into consideration. For a good match both precision and recall should be high and the difference between them should be low.

## **Data collection and analysis**

SocSciBot web crawler was used for data collection (Thelwall, 2001). The crawler has been developed for accurate and complete crawl of specified large sites and it has been used in several other webometric studies before (e.g. Thelwall, 2002; Smith & Thelwall, 2002; Björneborn, 2004; Thelwall & Wilkinson, 2004; Li et al., 2005a; Stuart, Thelwall & Harries, 2007). The design and use of the crawler are also well documented (Thelwall, 2004) and hence it was also chosen for this research. All municipal websites were crawled in July 2006 and all the links and all the pages found on the municipal websites were collected.

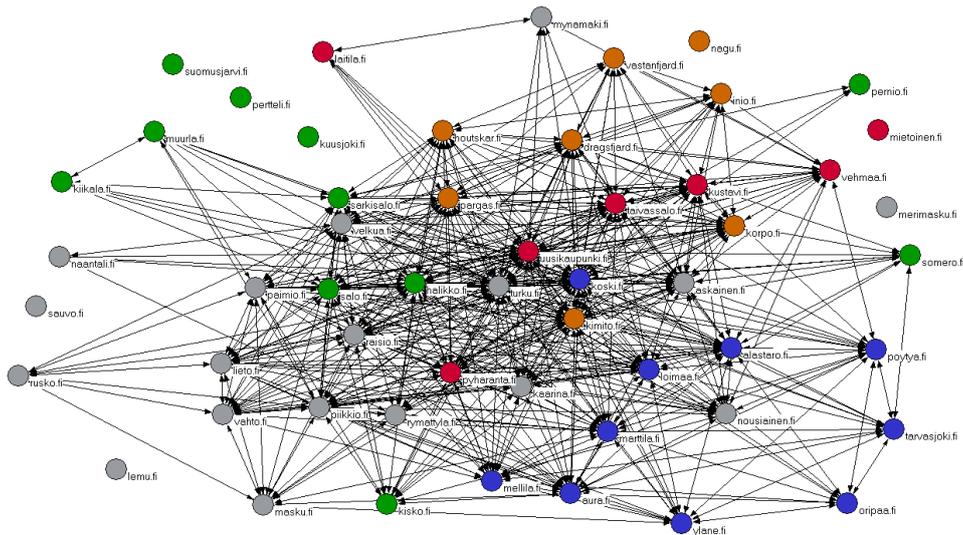
All outlinks that targeted some page outside the municipal web space were chosen for closer analysis. All the co-outlinks were also outlinks on the municipal websites, but not all outlinks were co-outlinks. All links that appeared only on a single municipality's website were removed and the co-outlinks were manually extracted from the data. This meant manually searching for links that two municipalities had created to target the same website or web page outside the municipal web space. All duplicate links at the site level were removed (Thelwall & Wilkinson, 2003) and the link data was converted into a binary matrix with the BibExcel software (Persson, 2006) for matrix comparisons and visualizations.

The municipal websites are official representations of the municipalities and the content on the websites represent the services provided by the municipalities and the duties and responsibilities they have. The links published on the municipal websites also indicated the official nature of the municipal websites and reflected different aspects of the local government authorities. Due to the official nature of linking from municipal websites simple classification schemes of the different linking types were considered as appropriate. The motivations to create links were classified based on 1) the type of link (i.e. in what context the link had been created in) and 2) the purpose of the link (what purpose the link fulfilled). For each linking type the classification was built inductively during the classification and new categories were created as new links that did not fit in any existing categories appeared in the classification process. The purpose of the link was classified by visiting the source and target pages and by judging what purpose the link fulfilled by combining the two web pages and by directing traffic from the source page to the target page or by showing a connection between the two websites or pages.

## **Results**

The co-outlinks were outlinks on the municipal websites targeting some other sites or pages outside the municipal web space of the region of Finland Proper. The connection between the municipalities formed by co-outlinks was created when two or more municipalities had outlinks to the same page or the same site outside the municipalities. There were a total of 7535 co-outlinks that targeted 1950 websites.

The binary connectivity matrix based on co-outlinking had a total of 986 links. The density of the matrix was 34.5%, which meant that about one third of all possible links were present in the network. The co-outlinks were visualized with Pajek and the Kamada-Kawai algorithm (Kamada & Kawai, 1989) in Figure 1, but no visible patterns emerged from this visualization as there were plenty of links connecting municipalities from different areas of the graph. Neighboring municipalities and municipalities from the same functional regions (a higher more unofficial level of local administration) were scattered all over the map without any clear patterns.



**Figure 1. Co-outlinking from the municipal websites**

It was possible that there were many different overlapping categories of link types representing different duties and responsibilities that the municipalities have and that could be the reason why there were no clear patterns in the network graph. If these different link types and motivations could be extracted and analyzed separately some useful information could be discovered. Categorizing the co-outlinking could show the shared interests of the municipalities and the categorization could tell something about what those shared interests were.

#### *Co-outlinking motivations*

It was hypothesized that the co-outlinking data included many different types of links and purposes for creating them, and that analyzing these types or genres separately could reveal some new information about the relationships between the municipalities. Categorizing the co-outlinking from the municipalities according to areas of municipalities responsibilities could also reveal something about the shared interests related to that particular type of co-outlinking or that particular type of municipal activities. If this approach leads to a successful mapping of the shared interests of the municipalities, then these categories would be the answer to the second research question.

When the co-outlinking motivations were investigated it was discovered that the purpose of almost all of the co-outlinks was the same, to link to some external information sources or to extend the information provided on the municipalities' websites with some additional resources. In that sense the co-outlinking links indicated some official action of the municipalities and hence classification of the links did not leave much room for interpretations. Because the purpose of the links was the same for almost all of the links, the links were categorized according to the type of information associated with the links and the context they were created in.

All 7535 co-outlinks were manually categorized in 21 different categories according to the content on the source and target pages. The categorization was done by creating new categories when new types of content that did not fit in existing categories were found. The link types and the categories formed reflected the services provided by the municipalities and their responsibilities. The category *Libraries* had 392 co-outlinks to small local libraries, bigger national libraries, online library catalogues and public archives. Among other

categories that were created were e.g. *Education* with 540 co-outlinks, *Health care* with 370 co-outlinks, *Recycling* with 99 co-outlinks, *Recreation* with 723 co-outlinks and *Business* with 679 co-outlinks. All the categories and their respective link types are listed in table 1 below.

**Table 1. Categories of co-outlinking**

<b>Category</b>	<b>Co-outlinks</b>
Miscellaneous	1135
Newspapers	998
Recreation	723
Business	679
Education	540
Central government administration	425
Libraries	392
Health care	370
Sports	298
Agriculture	277
Traffic	253
Nature	248
Search engines	221
Tourism	203
Software and hardware	184
Living and building	154
Banking and insurance	106
Other regions	103
Recycling	99
Kids and youth	99
Religion	28
<b>Total</b>	<b>7535</b>

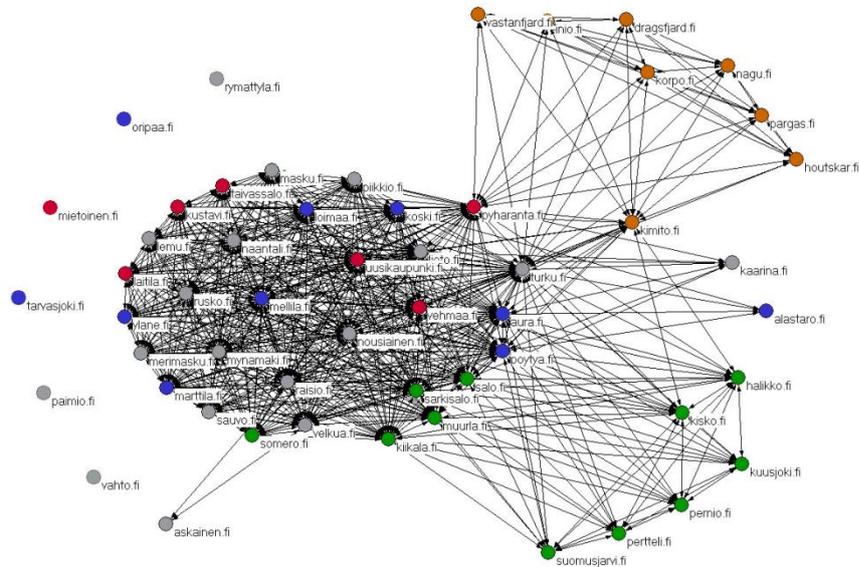
Cooperation between the libraries in the region of Finland Proper is very active. There were six library networks with online databases in the region and most of the libraries belonged to one of them. For instance the online database called *Blanka* (<http://www.blanka.fi/>) was a shared resource of the municipalities in the functional region of Turunmaa, namely Dragsfjärd, Houtskär, Iniö, Kimito, Korpo, Nagu, Pargas and Västanfjärd. The other databases were called *Loisto* (<http://loisto.kirjas.to/>), *Loki* (<http://www.lokikirjastot.fi/>), *Salomo* (<http://kirjastot.salonseutu.fi/>), *Peimari* (<http://kaarina.kirjas.to/>) and *Kirvakka* (<http://www.pyharanta.fi/index.php?id=142>). Only Lieto, Somero and Turku were not part of any library cooperation or there was no information found about such cooperation. Because the cooperation between the libraries in the region was so well defined through these library networks, the category with co-outlinking to libraries and library related websites and pages was chosen for closer analysis. Other co-outlinking categories could also have been used but they did not have as clearly defined cooperation as the libraries did and hence lacked a clearly defined offline network that the linking network could have been compared with.

The existing cooperation between libraries in the region was used to create a binary connectivity matrix based on the member municipalities sharing an online database. This

binary connectivity matrix was then compared with the co-outlinking network based on library related linking.

### *Total co-outlinking to libraries*

As a result of dividing the total co-outlinks to categories according to the content of the source and target pages some patterns started to emerge in the graph. In Figure 2 below the co-outlinks to libraries or library related websites or pages were visualized. The density of the matrix and the graph was 42%.



**Figure 2. Co-outlinking to libraries and library related websites**

The graph showed some clustering. The library cooperation in the functional regions of Turunmaa and Salo could be separated from the graph, but there was still a large, very tightly connected large cluster with no obvious patterns in it. In this tightly connected central cluster were municipalities from different functional regions mixed together. Six municipalities were disconnected from the graph.

### *Co-outlinking to libraries with manually cleaned linking data*

The library related co-outlinking data was also manually cleaned by removing all co-outlinks that lead to websites or pages located outside the studied region of municipalities (Figure 3). The removed links included links to other libraries, archives, and information sources located outside the region of Finland Proper. The density of this manually cleaned matrix was 12%, which meant that 12% of all possible links were present in the graph.



municipalities, but unless some parts of the categorizing or the data cleaning can be automated, using co-outlinking data cannot be recommended for this purpose.

## Conclusions

This research studied if co-outlinking from the municipalities could indicate some shared interests of the municipalities. In the first analysis it seemed that in the total co-outlinking data there were some different overlapping types or categories of interests indicated by the vast variety of target sites. By classifying the co-outlinking to categories before analyzing the data it was possible to get focused data sets of different co-outlinking categories, each representing shared interests of the municipalities. The category with libraries and library related content was chosen for closer analysis because there was clearly defined library cooperation that had manifested itself in six online library databases and library cooperation in the region. Based on this cooperation a matrix was built that was used to test for similarities with the co-outlinking matrix.

The manually cleaned data produced the best match with the actual situation of library cooperation in the region. This is supported by the fact that the co-outlinking links were usually created to extend the information provided on municipal websites or to direct visitors to other resources. With that it can be concluded that co-outlinking has some potential to be used to map shared interests and in the case of libraries, it can also be used to map cooperation, but not all categories or types of co-outlinking have as clearly defined cooperation as libraries do and they may therefore not produce as reliable results.

Manually cleaning the data required that every website that was co-outlinked to had to be visited and it had to be determined whether the organization or service behind the website was in the region of Finland Proper or not. Manually cleaning the data gave the best possible data about the shared interests of the municipalities, but the method required a lot of work and it cannot be recommended unless some automatic detection method could be developed.

Developing methods to automatically extract and filter co-outlinking data could help in finding useful applications for the co-outlinking data. Future research and development may find methods to automate some of the manual work that was done in this research to filter the best possible co-outlinking data. One possibility that could be researched further would be to filter the links according to domain, i.e. include all co-outlinks to .fi country code top-level domains.

## References

- Bar-Ilan, J. (2004a). Self-linking and self-linked rates of academic institutions on the Web. *Scientometrics*, vol. 59, no. 1, pp. 29-41.
- Bar-Ilan, J. (2004b). A microscopic link analysis of academic institutions within a country – the case of Israel. *Scientometrics*, vol. 59, no. 3, pp. 391-403.
- Bar-Ilan, J. (2005). What do we know about links and linking? A framework for studying links in academic environments. *Information Processing and Management*, vol. 41, pp. 973-986.
- Björneborn, L. (2001). Small-world linkage and co-linkage. In *Proceedings of the 12th ACM conference on Hypertext and Hypermedia*, pp. 133-137. Retrieved May 3, 2009, from [http://vip.db.dk/lb/papers/ht01\\_bjorneborn.pdf](http://vip.db.dk/lb/papers/ht01_bjorneborn.pdf).
- Björneborn, L. (2004). *Small-World Link Structures across an Academic Web Space - a Library and Information Science Approach*. [PhD thesis] Copenhagen: Royal School of Library and Information Science. Retrieved April 16, 2009, from <http://vip.db.dk/lb/phd/phd-thesis.pdf>.
- Clarke, S.J. & Willett, P. (1997). Estimating the recall performance of Web search engines. *Aslib Proceedings*, vol. 49, no. 7, pp. 184-189.

- Dean, J. & Henzinger, M. (1999). Finding related pages in the World Wide Web. *Computer Networks and ISDN Systems*, vol. 31, pp. 389-401.
- García-Santiago, L. & Moya-Anegón, F. de (2009). Using co-outlinks to mine heterogeneous networks. *Scientometrics*, vol. 79, no. 3, pp. 681-702.
- Haas, S.W. & Grams, E. (1998). Page and link classifications: connecting diverse resources. *ACM digital libraries 1998*, pp. 99-107.
- Harter, S.P. & Hert, C.A. (1997). Evaluation of information retrieval systems: approaches, issues and methods. *Annual Review of Information Science and Technology*, vol. 32, pp. 3-79.
- Holmberg, K. (2010). Co-inlinking to a municipal Web space: A webometric and content analysis. *Scientometrics*, vol. 83, pp. 851-862.
- Kamada, T. & Kawai, S. (1989). An algorithm for drawing general undirected graphs. *Information Processing Letters*, vol. 31, no. 1, pp. 7-15.
- Kessler, M.M. (1963). Bibliographic coupling between scientific papers. *American Documentation*, vol. 14, no. 1, pp. 10-25.
- Kim, H.J. (2000). Motivations for hyperlinking in scholarly electronic articles: a qualitative study. *Journal of the American Society for Information Science and Technology*, vol. 51, no. 10, pp. 887-899.
- Li, X., Thelwall, M., Wilkinson, D. & Musgrove, P. (2005). National and international university departmental Web site interlinking. Part 1: Validation of departmental link analysis. *Scientometrics*, vol. 64, no. 2, pp. 151-185.
- Persson, O. (2006) *Bibexcel - a toolbox for bibliometricians* (Version 2008-04-11). [Computer software]. Retrieved May 8, 2008, from <http://www.umu.se/inforsk/Bibexcel/>.
- Ortega, J.L., Aguillo, I., Cothey, V. & Scharnhorst, A. (2008). Maps of the academic web in the European Higher Education Area – an exploration of visual web indicators. *Scientometrics*, vol. 74, no. 2, pp. 295-308.
- Schneider, J. & Borlund, P. (2007). Matrix comparison, Part 2: Measuring the resemblance between proximity measures or ordination results by use of the Mantel and Procrustes statistics. *Journal of the American Society for Information Science and Technology*, vol. 58, no. 11, pp. 1596-1609.
- Smith, A. & Thelwall, M. (2002). Web impact factors for Australasian Universities. *Scientometrics*, vol. 54, no. 3, pp. 363-380.
- Smith, A.G. (2004). Web links as analogues of citations. *Information Research*, vol. 9, no. 4, paper 188. Retrieved April 26, 2009, from <http://informationr.net/ir/9-4/paper188.html>.
- Stuart, D. & Thelwall, M. (2005). What can university-to-government web links reveal about university-government collaborations? In Ingwersen, P. & Larsen, B. (eds.), *Proceedings of the 10th International Conference of the International Society of Scientometrics and Informetrics*. Stockholm, Sweden: Karolinska University press, vol. 1, pp. 188-192.
- Stuart, D., Thelwall, M. & Harries, G. (2007). UK academic web links and collaboration – an exploratory study. *Journal of Information Science*, vol. 33, no. 2, pp. 231-246.
- Thelwall, M. (2001). A web crawler design for data mining. *Journal of Information Science*, vol. 27, no. 5, pp. 319-325.
- Thelwall, M. (2002). Evidence for the existence of geographic trends in university Web site interlinking. *Journal of Documentation*, vol. 58, no. 5, pp. 563-574.
- Thelwall, M. (2003). What is this link doing here? Beginning a fine-grained process of identifying reasons for academic hyperlink creation. *Information Research*, vol. 8, no. 3. Retrieved May 25, 2009, <http://informationr.net/ir/8-3/paper151.html?text=1>.
- Thelwall, M. (2004). *Link Analysis: An Information Science Approach*. Elsevier Academic Press.
- Thelwall, M. (2006). Interpreting social science link analysis research: A theoretical framework. *Journal of the American Society for Information Science and Technology*, vol. 57, no. 1, pp. 60-68.
- Thelwall, M. & Harries, G. (2003). The connection between the research of a university and counts of links to its Web pages: An investigation based upon a classification of the relationships of pages to the research of the host university. *Journal of the American Society for Information Science and Technology*, vol. 54, no. 7, pp. 594-602.
- Thelwall, M., Harries, G. & Wilkinson, D. (2003). Why do web sites from different academic subjects interlink? *Journal of Information Science*, vol. 29, no. 6, pp. 453-471.

- Thelwall, M. & Wilkinson, D. (2003). Three target document range metrics for University Web sites. *Journal of the American Society for Information Science and Technology*, vol. 54, no. 6, pp. 490-497.
- Thelwall, M. & Wilkinson, D. (2004). Finding similar academic Web sites with links, bibliometric couplings and colinks. *Information Processing & Management*, vol. 40, pp. 515-526.
- Vasileiadou, E. & Besselaar, P. van den (2006). Linking shallow, linking deep. How scientific intermediaries use the web for their network of collaborators. *Cybermetrics*, vol. 10, no. 1, paper 4. Retrieved March 18, 2011, from <http://www.cindoc.csic.es/cybermetrics/articles/v10i1c4.html>.
- Vaughan, L., Gao, Y. & Kipp, M. (2006). Why are hyperlinks to business Websites created? A content analysis. *Scientometrics*, vol. 67, no. 2, pp. 291-300.
- Vaughan, L., Kipp, M. & Gao, Y. (2007). Why are Websites co-linked? The case of Canadian universities. *Scientometrics*, vol. 72, no. 1, pp. 81-92.
- Vaughan, L. & You, J. (2005), Mapping business competitive positions using Web co-link analysis. In: Ingwersen, P. & Larsen B. (Eds), *Proceedings of ISSI 2005 – the 10th International Conference of the International Society for Scientometrics and Informetrics*, pp. 534–543, Stockholm, Sweden, July 24–28, 2005.
- Wilkinson, D., Harries, G., Thelwall, M. & Price, L. (2003). Motivations for academic web site interlinking: evidence for the Web as a novel source of information on scholarly communication. *Journal of Information Science*, vol. 29, no. 1, pp. 49-56.