

Unifying Text and Link Analysis

Brian D. Davison

Department of Computer Science & Engineering
 Lehigh University, Bethlehem, PA 18015 USA
 davison@lehigh.edu

Abstract

This position paper presents a simple yet profound idea. By thinking about the relationships between and within terms and documents, we can generate a richer representation that encompasses aspects of Web link analysis as well as text analysis techniques from information retrieval. This paper shows one path to this unified representation, and demonstrates the use of eigenvector calculations from Web link analysis by stepping through a simple example. We further speculate that this general approach, which we term *relationship analysis*, can apply to other domains as well.

1 Introduction

The ubiquity of the World-Wide Web has placed information retrieval systems at the fingertips of millions of people, in the form of Web search engines. While those search engines initially used textual analysis to match documents with queries, the use of link analysis techniques have become more common, such that all major search engines now incorporate some kind of link analysis.

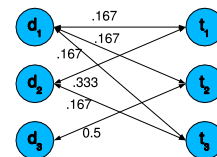
By link analysis, we refer to the study or use of algorithms operating over the Web’s link graph. This graph defines the relationships between pages, based on the hyperlinks from page to page. Such algorithms might help to find relevant documents for a query [Kleinberg, 1999; Page *et al.*, 1998], or find similar or duplicate documents. [Dean and Henzinger, 1999; Bharat and Broder, 1999]. Link analysis has similarities to, and has benefited from, social network analysis [Wasserman and Faust, 1994] and bibliographic citation analysis, in particular co-citation [Small, 1973] and bibliographic coupling [Kessler, 1963]. Thus, web link analysis is one form of the more generic problem of what we term *relationship analysis* which includes any algorithm operating over a network graph or matrix built from relationships between entities in the graph.

However, in the Web, most link analysis has only limited textual components. In systems based on Kleinberg’s HITS [Kleinberg, 1999], information retrieval based on text is used to select the initial core set of documents. PageRank [Page *et al.*, 1998], in contrast, doesn’t use text at all to determine a document’s authority score.

	d_1	d_2	d_3
t_1	1	2	0
t_2	1	0	3
t_3	1	1	0

$$M := \begin{bmatrix} .167 & .333 & 0.0 \\ .167 & 0.0 & 0.5 \\ .167 & .167 & 0.0 \end{bmatrix}$$

(a) Sample terms and documents, and their matrix view.



(b) Graph view.

Figure 1: A simple term-document matrix M and equivalent graph with three documents and three terms, using length-normalized TF-IDF weighting.

Likewise, in traditional information retrieval research, the use of eigenvectors — which underlies most Web link analysis — is commonly limited to the dimension-reduction approach found in LSI [Deerwester *et al.*, 1990].

In the rest of this position paper, we present a small example and step through its representations (Section 2) from typical term vectors to our unified representation. In Section 3, we describe the application of web link analysis algorithms to this unified representation. We conclude with Section 4 where we discuss future extensions and summarize our thoughts.

2 Representations

Typical information retrieval approaches represent documents as vectors of term weights. When placed together, these vectors form a term-document matrix, in which one axis enumerates each document, and the other enumerates each term found in the collection. Figure 1(a) displays this matrix, using a simple length-normalized form of TF-IDF [Salton and

	Terms	Docs
Terms	term-term	term-doc
Docs	doc-term	doc-doc

Figure 2: The generic augmented matrix with sub-matrices.

$$M_1 := \begin{bmatrix} 0.0 & 0.0 & 0.0 & .167 & .333 & 0.0 \\ 0.0 & 0.0 & 0.0 & .167 & 0.0 & 0.5 \\ 0.0 & 0.0 & 0.0 & .167 & .167 & 0.0 \\ .167 & .167 & .167 & 0.0 & 0.0 & 0.0 \\ .333 & 0.0 & .167 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.5 & 0.0 & 0.0 & 0.0 & 0.0 \end{bmatrix}$$

Figure 3: The expanded matrix M_1 with three documents and three terms.

McGill, 1983] weighting. This matrix corresponds to a bipartite graph in which terms and documents are nodes, with undirected links between them (Figure 1(b)) wherever there are non-zero entries in the matrix. In our thinking, these links are really relationships. That is, a document is linked to a term when that term is found in the document and vice versa. However, this model is limited to undirected links. To use directed links, we need to expand the matrix. In this work we propose the use of four submatrices (as shown in Figure 2). Continuing with our example, Figure 3 provides our expanded matrix, which incorporates the same term-document submatrix (M), but now also includes a document-term submatrix (M^T), plus new term-term and document-document submatrices. This particular matrix still represents the same graph, but now there are directed links that happen to have the same weights in both directions (shown in Figure 4). However, this matrix provides a richer representation. If desired, we can now have different weights for links in different directions, as we will demonstrate shortly. We can also have non-zero weights between nodes of the same type — that is, we can have links between terms, or between documents. The links between documents easily correspond to citations (whether hypertext or bibliographic), and the links between terms might be similarity values.

The doc-doc submatrix is exactly the matrix used in Web link analysis algorithms (typically a weighted adjacency matrix of some kind). By varying the mechanisms to determine weights, one can specify different algorithms. For example, the simplified form of the PageRank algorithm [Page *et al.*, 1998] doesn’t use the adjacency matrix directly — it uses the inverse of the number of outgoing links. In text analysis we often do something similar — one way to normalize the term frequencies is to divide by the document length. In our augmented matrix, we can use different weights for different link directions. Therefore, we can still use a normalized term frequency, just as in PageRank, for the links from docs to terms. We can even use the same principle (the inverse of the number of outgoing links) to weight the links from terms to docu-

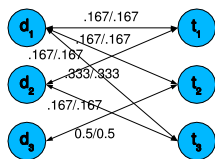


Figure 4: Graph view of the expanded matrix with three documents and three terms.

	EV_1	EV_2	EV_3	EV_4	EV_5	EV_6
$E.value$	0.541	-0.541	0.420	-0.420	0.061	-0.061
t_1	0.200	0.200	0.570	0.570	0.367	-0.367
t_2	0.662	0.662	-0.248	-0.248	0.024	-0.024
t_3	0.148	0.148	0.337	0.337	-0.604	0.604
d_1	0.311	-0.311	0.262	-0.262	-0.578	-0.578
d_2	0.169	-0.169	0.587	-0.587	0.357	0.357
d_3	0.612	-0.612	-0.295	0.295	0.196	0.196

Figure 5: Eigenvalues and eigenvectors of the running example.

ments. This, in a sense, corresponds to the inverse document frequency of TF-IDF. However, other variations are possible (potentially corresponding to other term weighting schemes investigated in the information retrieval literature).

3 Algorithms

Most Web link analysis algorithms revolve around the use of eigenvector calculations, and differ in their matrix weights and generation. For example, the random surfer model in PageRank, while described as a probabilistic jump from any document to any other, can be implemented as a simple operation over the existing doc-doc matrix (adding low-weight links from every page to every other page). This factor in PageRank is described as helping to prevent rank sinks.

Note that in our model, while this mechanism may be helpful, it is not strictly needed as the term-doc connectivity should be sufficient to prevent sinks.

Web link analysis generates eigenvectors primarily for two reasons. The first is to generate a total ordering of documents. Typically this comes from the principal eigenvector (e.g., in PageRank), but can also be from a combination of eigenvectors (e.g., as in DiscoWeb [Davison *et al.*, 1999]). This total ordering may then be combined with other factors, such as textual relevance, to generate a final ordering for presentation to the user. The second reason for generating eigenvectors is to investigate communities within some topic, as suggested by Kleinberg [Kleinberg, 1999] and used in DiscoWeb and elsewhere. These “communities” are effectively clusters of highly-interconnected pages on a topic. This is often successful because authors of pages on a particular topic tend to link to each other, or are linked from one or more common hub pages.

We propose using these approaches on text as well. We can directly calculate eigenvectors of our augmented matrix, and consider some fraction of documents with high absolute values in the principal eigenvector and at both the positive and negative ends of non-principal eigenvectors as clusters. A nicety of our model is that not only will it generate clusters, but those clusters will be (at least in part) self-describing because the clusters will include terms as well as documents.

Figure 5 shows the results of calculating the eigenvectors from our running example matrix M_1 . In the principal eigenvector, we see that term t_2 scores highly, followed closely by document d_3 . This matches the original distribution of terms to documents (since d_3 contained only instances of t_2). Likewise, in EV_3 , we see that document d_2 is ranked highest, followed closely by term t_1 .

Calculating the principal eigenvector of the standard augmented matrix doesn’t make much sense, as while it will gen-

erate an order for the terms and documents in the matrix, that order has no underlying meaning (unlike the calculation in the Web, which corresponds to authoritativeness, assuming links convey authority upon the target of the link). Additionally, it doesn't take the query into consideration, and so separate text analysis would be needed.

Our approach is based on the idea that PageRank can be personalized [Page *et al.*, 1998; Haveliwala, 2002]. Instead of using a uniform probability of jumping randomly from one page to any other, an emphasis can be made to a certain page or set of pages. This has the effect of ranking other pages in relation to the emphasized ones. We can do the same with our augmented matrix. A given query, either in the form of a page in the corpus or a set of terms can be emphasized by modifying the augmented matrix so that all objects have a link to the emphasized object(s). With enough emphasis, the query objects will rise to the top of the principal eigenvector, and remaining objects will be ordered in terms of their relevance to the initial query objects.

For example, we saw above that the principal eigenvector placed terms t_1 and t_3 well below the value assigned to t_2 . The documents were ordered d_3, d_1, d_2 . We can skew the eigenvector by emphasizing t_1 , as if it were the given query. To do so, we arbitrarily modify the network to incorporate links from all objects to t_1 with weight .2, and decrease all existing weights by 20%. This modified matrix M_2 is shown in Figure 6.

$$M_2 := \begin{bmatrix} 0.200 & 0.0 & 0.0 & 0.133 & 0.267 & 0.0 \\ 0.200 & 0.0 & 0.0 & 0.133 & 0.0 & 0.400 \\ 0.200 & 0.0 & 0.0 & 0.133 & 0.133 & 0.0 \\ 0.333 & 0.133 & 0.133 & 0.0 & 0.0 & 0.0 \\ 0.467 & 0.0 & 0.133 & 0.0 & 0.0 & 0.0 \\ 0.200 & 0.400 & 0.0 & 0.0 & 0.0 & 0.0 \end{bmatrix}$$

Figure 6: The skewed matrix M_2 emphasizing t_1 .

The principal eigenvector of M_2 contains the values [0.854, 0.100, 0.148, 0.246, 0.415, 0.067], resulting in a rank ordering of $t_1, d_2, d_1, t_3, t_2, d_3$, which is in fact the desired ordering.

4 Discussion

This position paper outlines work in progress, and so the ideas presented here are still exploratory in nature. Additional effort will be needed to experimentally verify (on a larger scale) the claims made here. Moreover, we make no claims of optimality in the details of our approach.

We envision many extensions to this work, the most obvious being the incorporation of term-term and doc-doc links, as mentioned in Section 2. However, some care may be needed to balance the relative influence of various kinds of relationships, and prevent domination by a single submatrix (such as, for example, the doc-doc link submatrix over contributions of the term-doc and doc-term submatrices). We also expect alternative weighting schemes, and efforts to streamline the computational overhead of these approaches.

The general approach we have taken to relationship analysis is not specific to textual content. It should also be applicable to other domains, such as collaborative filtering and recommender systems for various items, including music, movies, etc.

In summary, this position paper has promoted a unified representation for text and link data, and the operation of web link analysis algorithms for retrieval and clustering. It is hoped that the use of a simple example has helped provide an intuitive understanding of this approach to stimulate future research.

References

- [Bharat and Broder, 1999] Krishna Bharat and Andrei Broder. Mirror, mirror on the Web: a study of host pairs with replicated content. *The International Journal of Computer and Telecommunications Networking*, 31(11-16):1579–1590, May 1999. Proceedings of the 8th International Conference on the World Wide Web.
- [Davison *et al.*, 1999] Brian D. Davison, Apostolos Gerasoulis, Konstantinos Kleisouris, Yingfang Lu, Hyunju Seo, Wei Wang, and Baohua Wu. DiscoWeb: Applying link analysis to Web search. In *Poster proceedings of the Eighth International World Wide Web Conference*, pages 148–149, Toronto, Canada, May 1999.
- [Dean and Henzinger, 1999] Jeffrey Dean and Monika R. Henzinger. Finding related pages in the World Wide Web. In *Proceedings of the Eighth International World Wide Web Conference*, pages 389–401, Toronto, Canada, May 1999.
- [Deerwester *et al.*, 1990] Scott C. Deerwester, Susan T. Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
- [Haveliwala, 2002] Taher H. Haveliwala. Topic-sensitive PageRank. In *Proceedings of the Eleventh International World Wide Web Conference*, Honolulu, Hawaii, May 2002.
- [Kessler, 1963] M. M. Kessler. Bibliographic coupling between scientific papers. *American Documentation*, 14:10–25, 1963.
- [Kleinberg, 1999] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46(5):604–632, 1999. An earlier version of this paper appeared in the Proceedings of the ACM-SIAM Symposium on Discrete Algorithms, 1998.
- [Page *et al.*, 1998] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank citation ranking: Bringing order to the Web. Unpublished draft, 1998.
- [Salton and McGill, 1983] Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*. McGraw Hill, New York, 1983.
- [Small, 1973] Henry G. Small. Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24:265–269, 1973.
- [Wasserman and Faust, 1994] Stanley Wasserman and Katherine Faust. *Social Network Analysis*. Cambridge University Press, 1994.