



ELSEVIER

Computational Statistics & Data Analysis 41 (2002) 19–45

COMPUTATIONAL
STATISTICS
& DATA ANALYSIS

www.elsevier.com/locate/csda

Web document clustering using hyperlink structures

Xiaofeng He^{a,b,*,1}, Hongyuan Zha^{a,1}, Chris H.Q. Ding^{b,2},
Horst D. Simon^{b,2}

^a*Department of Computer Science and Engineering, The Pennsylvania State University, University Park, PA 16802, USA*

^b*NERSC Division, Lawrence Berkeley National Laboratory, University of California, Berkeley, CA 94720, USA*

Abstract

With the exponential growth of information on the World Wide Web, there is great demand for developing efficient methods for effectively organizing the large amount of retrieved information. Document clustering plays an important role in information retrieval and taxonomy management for the Web. In this paper we examine three clustering methods: K-means, multi-level METIS, and the recently developed *normalized-cut* method using a new approach of combining textual information, hyperlink structure and co-citation relations into a single similarity metric. We found the normalized-cut method with the new similarity metric is particularly effective, as demonstrated on three datasets of web query results. We also explore some theoretical connections between the normalized-cut method and the K-means method. © 2002 Elsevier Science B.V. All rights reserved.

Keywords: World Wide Web; Graph partitioning; Cheeger constant; Clustering method; K-means method; Normalized cut method; Eigenvalue decomposition; Link structure; Similarity metric

1. Introduction

Currently the World Wide Web contains a huge number of documents and it is still growing rapidly. Finding the relevant documents to satisfy a user's information need is

* Corresponding author.

E-mail addresses: xhe@cse.psu.edu (X. He), zha@cse.psu.edu (H. Zha), chqding@lbl.gov (Chris H.Q. Ding), hdsimon@lbl.gov (H. D. Simon).

¹ This work was supported in part by NSF grant CCR-9901986.

² Supported by Department of Energy through an LBL LDRD fund.

0167-9473/02/\$ - see front matter © 2002 Elsevier Science B.V. All rights reserved.

PII: S0167-9473(02)00070-1

a very important and challenging task. Many commercial search engines have been developed and used by millions of people all over the world. However, the relevancy of documents returned in search engine result sets is still lacking, and further research and development is needed to really make search engines a ubiquitous information-seeking tool. The World Wide Web has a rich structure: it contains both textual web documents and the hyperlinks that connect them. The web documents and hyperlinks between them form a *directed* graph in which the web documents can be viewed as vertices and the hyperlinks as directed edges. Algorithms have been developed utilizing this directed graph to extract information contained in a collection of hyperlinked web documents. Kleinberg proposed the HITS algorithm based purely on hyperlink information to retrieve the most relevant information: *authority* and *hub* documents for a user query Kleinberg (1998). However, if the hypertext collection consists of several topics, authority and hub documents may only cover the most popular topics and leave out the less popular ones. One way to remedy this situation is to first partition the hypertext collection into topical groups, and present the search results as a list of topics to the user. This leads to the need to cluster web documents based on both the textual and hyperlink information.

There exists a large amount of literature on clustering methods and algorithms Everitt (1993), Gordon (1981). Generally speaking, the purpose of *cluster analysis* is to organize the data into meaningful groups: the data objects in the same group are highly similar and those in different groups are dissimilar. Judging the effectiveness of a clustering algorithm is difficult and usually application-dependent. In this paper, we apply a similarity-based clustering method to the problem of clustering web documents. It utilizes a graph-theoretic criterion called *normalized cut* which has its root in the study of graph isoperimetric problems Cheeger (1970), Chung (1997). This method was proposed by Shi and Malik and has been successfully used in image segmentation Shi and Malik (1997).

Document clustering has been studied by many people Willett (1988). Clustering retrieval results have been examined by Hearst and Paderson (1996) based on textual information only, with emphasis on summarization. More recently this approach is taken in the Grouper web interface Zamir and Etzioni (1999). In the context of clustering web documents, in addition to the textual contents of documents, many other sources of information can be effectively used to enhance clustering effectiveness, for example, hyperlinks between documents, and co-citation (co-reference) patterns among documents. In our web document clustering approach, we incorporate information from hyperlink structure, co-citation patterns and textual contents of documents to construct a new similarity metric for measuring the topical homogeneity of web documents. Specifically, the hyperlink structure is used as the dominant factor in the similarity metric, and the textual contents are used to modulate the strength of each hyperlink. The similarity metric is used to construct a weighted graph which is then fed to the clustering method based on the normalized cut. This combination gives rise to a powerful method which effectively organizes the retrieved information against a user query and presents the organized information in a more accessible form for the user.

The rest of the paper is organized as follows: In Section 2, we give some background materials about search engines and explain why it is important to cluster the search

engine result sets effectively. In Section 3, we explain how the textual contents, the hyperlink structure and co-citation patterns can be combined to construct a weighted graph for partitioning. In Section 4, we discuss the *normalized cut* in spectral graph partitioning. We also outline a spectral clustering algorithm based on the normalized cut. The relationship between the normalized cut and the Cheeger constant is discussed in Section 5. In Section 6, we explore the connection between the normalized-cut based clustering method and the K-means method in the context of similarity-based clustering. Section 7 presents our experiments with the normalized-cut clustering algorithm. We give details of our data preparation process, and analyze the structures of clusters obtained for several broad-topic queries.

2. Organizing query result sets for search engines

The World Wide Web is a rich source of information. It has grown to more than one billion unique web documents and continues to grow roughly at a rate of one million documents per day Bharat and Broder (1998). While the users are enjoying the large amount of information available on the World Wide Web, it poses a problem to both the users and the search engines. The problem lies in the sheer quantity of the web documents. When users submit a query, the search engine returns the retrieved web documents as the search result set. The query which users tend to submit to the search engine typically has 1–3 words with little query formulation Anick (1994), Croft et al. . Sometimes the users are uncertain about their information need when submitting the query Efthimiadis (1993), and for ordinary users, it is difficult to come up with query terms that accurately specify their information needs. These situations often result in broad-topic queries. Since many search engines retrieve the web documents based on the text similarity and link structures, after a broad-topic query is submitted, it is likely for the search engine to form a search result set with thousands, even up to millions of web documents. A particular issue arises: How to effectively organize such a large number of retrieved web documents for a broad-topic user query on a search engine according to their relevance, and provide the users with the information they are looking for? It is often the case that for some broad-topic queries, the documents for the popular topic tend to dominate the result set. Under this circumstance, web document ranking algorithms such as HITS are unable to solve this problem. The authoritative documents returned to users by HITS may represent only the most popular topic. Documents related to under-represented topics have less chance to be returned to the users while perhaps they are what the users are expecting. To overcome this problem, before ranking the web documents and presenting them to users, it is necessary to group the retrieved web documents into distinct topic areas, then return the ranked documents for each group according to their relevance.

The hyperlink structure of the World Wide Web provides us with rich information on web communities. It contains a lot of latent human annotation of the web society. This motivates us to cluster the web documents by partitioning the web link graph. If two web documents have very small text similarity, it is less likely that they belong to the same topic, even though they are connected by a hyperlink. Therefore, to improve the

quality of the graph representation, the text information can be incorporated into the link graph as a factor of edge weight. The co-citation is another relevance measure between two web documents based on how many other web documents create hyperlinks to both of them. For further improving the quality of the clustering results, combining the co-citation information into the link graph is also desirable.

For a clustering algorithm to be useful for the Web, it needs to be both effective and efficient. However, it is not necessary that the result sets be processed and the clusters be generated in real-time. Here, we are dealing with popular broad-topic queries, search engines can cache the result sets for those popular queries so that the processing can be done offline if necessary. We are currently investigating fast methods for spectral decomposition along the lines proposed in Frieze et al. (2000) to make our approach more suitable for real-time processing.

3. Similarity metric

One popular way for clustering data objects into subgroups is based on a similarity metric between objects, with the goal that objects within a subgroup are very similar, and objects between different subgroups are less similar. In clustering a graph, similarity between nodes is represented by edge weight.

In the web documents clustering problem, we take a new approach in defining the similarity between two web documents. We incorporate link structure, textual information, and co-citation information into a similarity metric which gives rise to the weight matrix W . The link structure is the dominant factor, and the textual similarity is used to modulate the strength of each hyperlink. To further enhance the link structure, co-citation is also incorporated.

3.1. Hyperlink structure

The link information is obtained directly from the link graph. The method to form the link graph is introduced in Section 7.1. Given a link graph $G = (V, E)$, which is directed, we define matrix A to be:

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \text{ or } (j, i) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

A is the adjacency matrix of the link graph where directionality of the hyperlinks is ignored. Link structure alone provides us with rich information on the topics of the document collection. By exploring the link structure, we are able to extract useful information from the web Chakrabarti et al. (1998), Gibson et al. (1998), Chakrabarti et al. (1999), [Flake et al. (2000), Kumar et al. (1999), Pirolli et al. (1996), Larson (1996)]. One of the most popular algorithms to extract document ranking information from the link structure is the HITS algorithm developed by Jon Kleinberg which will be briefly discussed in Appendix A.

3.2. Textual information

The textual information can be included to better cluster the web documents. Moreover, compared to printed literature, web documents reference each other more randomly. This is another reason that the text information is incorporated in order to regulate the influence of the document. One approach to incorporating the textual information is to measure the similarity between a user query and the anchor text (text between $\langle A \text{ HREF} = \dots \rangle$ and $\langle /A \rangle$) as in Li (1998), Chakrabarti et al. (1999). We experimented with this approach and found it did not work as effectively in our data sets as we had expected.

Here we use a new approach that (a) utilizes the entire text of a web document, not just the anchor text; (b) measures the textual similarity S_{ij} between two web documents i, j , instead of between the user query and the web document; (c) uses S_{ij} as the strength of the hyperlink between web documents i, j . The key observation here is that if two web documents have very little text similarity, it is unlikely that they belong to the same topic, even though they are connected by a hyperlink. Therefore S_{ij} properly gauges the extent or the importance of an individual hyperlink.

We represent each web document as a vector in the vector space model of IR (Information Retrieval) RijsBergen (1979) then compute the *similarity* between them. The higher the similarity, the more likely the two documents deal with the same topic. For each element of the vector we use the standard tf.idf weighting: $tf(i, j) * idf(i)$. $tf(i, j)$ is the Term Frequency of word i in document j , representing the number of occurrences of word i in document j . idf is the Inverse Document Frequency corresponding to word i , defined as

$$idf(i) = \log \left(\frac{\text{no. of total docs}}{\text{no. of docs containing word } i} \right).$$

Some words appear too frequently in many documents. We assume these words are not very useful to identify the documents. Inverse Document Frequency can effectively decrease the influence of these words.

Since the term vector lengths of the documents vary, we use cosine normalization in computing similarity. That is, if x and y are vectors of two documents d_1 and d_2 , then the similarity between d_1 and d_2 is:

$$S(d_1, d_2) = S(d_2, d_1) = \frac{\sum_i x_i y_i}{\|x\|_2 \|y\|_2}$$

where $\|x\|_2 = \sqrt{\sum_i x_i^2}$.

The similarities between documents form the similarity matrix S .

3.3. Co-citation patterns

Co-citation is another metric to measure the relevance of two web documents. If there are many documents pointing to both of them, then these two documents are likely to address a similar issue. The co-citation pattern is used by H. Small and others

to trace and map scientific literatures Small (1973). The co-citation C_{ij} of documents i and j is the number of web documents pointing to both i and j .

Incorporating the above information into the similarity metric, we form the weight matrix of the graph:

$$W = \alpha \frac{A \otimes S}{\|A \otimes S\|_2} + (1 - \alpha) \frac{C}{\|C\|_2} \quad (1)$$

where A is the adjacency matrix of the link graph. S is the similarity matrix. C is the co-citation matrix and $(A \otimes S)_{ij} = A_{ij}S_{ij}$, $\|C\|_2 = \sqrt{\sum_{i,j} C_{ij}^2}$. The meaning of notation \otimes applies to the rest of the paper. α is a real value between 0 and 1. The algorithms we introduce will be applied on matrix W .

The link and text information is retrieved *a priori*. The adjacency matrix A is formed by sweeping all the link graph edges once, requiring $O(|E|)$ time. Since each web document has an average out-degree (the number of links originated from it) around 8 Kleinberg et al. (1999), Kumar et al. (1999), on average there are 28 pairs of documents pointed to by a web document. It takes about $28|V|$ operations to compute the co-citation matrix C , requiring $O(|V|)$ time. Note that we only need the similarity of two documents which have a link between them. Suppose the average number of words in each document is N_w , then the running time to compute the similarity matrix S is $O(N_w|E|)$. Thus the total running time for computing the weight matrix W is $O(N_w|E| + |V|)$.

4. Partitioning method based on normalized cut

After forming the weight matrix W representing the web graph as in Section 3, we can cluster the web documents into distinct topic areas by applying some existing graph partitioning methods on W . The criteria used to measure the goodness of the partition will largely determine the final quality of the partition results. In graph theory, a set of edges that separates the graph into two disconnected parts is called an edge-cut (edge separator) of the graph. The standard minimum cut algorithms minimize the size of edge-cut alone. This often causes an unbalanced partition: it may cut a portion of a graph with a small number of vertices. In the context of graph clustering, this is in general not desirable. To avoid partitioning out a small part of a graph G by using edge-cut alone, many criteria utilize various normalized forms of edge-cut which are generally obtained by dividing the edge-cut by some measure of the size of the partitions. The *normalized cut* is one such measure.

4.1. Normalized cut

Given an undirected graph $G = (V, E)$ with edge weight matrix W , let S be a subset of V and $\bar{S} = V - S$, the complement of S in V . The normalized cut is defined as

$$\text{Ncut}(S, \bar{S}) = \frac{\text{cut}(S, \bar{S})}{W(S, V)} + \frac{\text{cut}(S, \bar{S})}{W(\bar{S}, V)}, \quad (2)$$

where $\text{cut}(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} W_{ij}$, $W(S, V) = \sum_{i \in S, j \in V} W_{ij}$. The partitioning problem here is to find the partitions that minimize the normalized cut $\text{Ncut}(S, \bar{S})$. If we let D be the diagonal matrix with $d_{ii} = \sum_{k=1}^n W_{ik}$, $i = 1, \dots, n$, the minimization problem associated with normalized cut has been shown to be equivalent to the following discrete minimization problem Shi and Malik (1997):

$$\frac{x^T(D - W)x}{x^T Dx} \quad (3)$$

subject to the constraint that $x_i \in \{1, -b\}$ and $x^T D e = 0$ where e is a vector with all elements equal to 1, and b is positive. Expression (3) is in the form of a Rayleigh quotient. Relaxing the condition $x_i \in \{1, -b\}$ and allowing the elements of x to take any real values, we can easily see that (3) can be minimized by the second smallest eigenvector of the generalized eigensystem:

$$(D - W)x = \lambda Dx. \quad (4)$$

With the introduction of the normalized-cut criterion, we can easily apply it to measure web graph partitioning problems. It effectively avoids the problem of cutting a small part of the graph. Notice that letting $y = D^{1/2}x$, Eq. (4) can be transformed to

$$D^{-1/2} W D^{-1/2} y = (1 - \lambda)y. \quad (5)$$

Therefore, we just need to compute the second largest eigenvector of $D^{-1/2} W D^{-1/2}$.

4.2. Spectral method

The normalized cut method is a further development in the class of spectral graph partitioning methods. Spectral methods are based on the early work of Donath and Hoffman (1972), Fiedler (1973, 1975) and many others (see a review by Mohar (1992)). It was popularized by the work of Pothen et al. (1990) and many other follow-up studies Hendrickson and Leland (1995), Spielman and Teng (1996) on graph bisection ($|S| = |\bar{S}|$) cases.

Spectral partitioning methods use eigenvectors of the Laplacian matrix $L = D - W$ of a graph to partition the graph. In particular, the eigenvector corresponding to the second smallest eigenvalue is most useful to partition the graph. This eigenvector is called the *Fiedler Vector* in recognition of Fiedler's work.

In Section 4.1, we use the second smallest eigenvector of the generalized eigensystem (4) to partition the graph. Clearly the normalized cut method is a spectral graph partitioning method with the removal of the constraint $|S| = |\bar{S}|$. Here we call the second smallest eigenvector of (4) the scaled Fiedler vector. One interesting property of the scaled Fiedler vector is its relation to algebraic connectivity:³ *theorem* given connected graph $G(V, E)$. Let f be the eigenvector corresponding to the second smallest eigenvalue λ_2 of the generalized eigensystem (4). Define $V_1 = \{v \in V: f_v \leq 0\}$, then the subgraph induced by V_1 is connected. Similarly, define $V_2 = \{v \in V: f_v \geq 0\}$, then the subgraph induced by V_2 is also connected.

³ The proof of the theorem follows similar arguments as in Theorem 2.1 in Pothen et al. (1990).

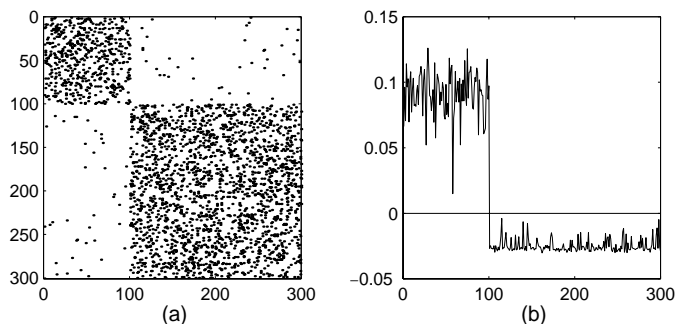


Fig. 1. (a) The weight matrix of a graph. (b) The scaled Fiedler vector values on each node.

Remark. The above theorem implies that if the original graph is connected, after partitioning using the scaled Fiedler vector, the subgraphs obtained are also connected if the scaled Fiedler vector has no zero entries. But in general this conclusion is not true for the K-means method which we will discuss later.

Now we give an illustrative example of using the scaled Fiedler vector to partition the graph. Fig. 1(a) is a weight (adjacency) matrix generated by matlab. The first diagonal block is 100×100 , and the second diagonal block is 200×200 . These two blocks represent two highly connected subgraphs. The off-diagonal blocks represent sparse connectivity between the two subgraphs. The adjacency matrix is symmetric with diagonal elements set to zero. Fig. 1(b) shows the plot of the scaled Fiedler vector of the generalized eigensystem corresponding to this matrix. From the plot, we can easily see that the scaled Fiedler vector cuts the nodes of the matrix into two distinct parts, if we choose zero as the cutting point.

4.3. Hierarchical divisive clustering

In this algorithm, we partition the weighted graph with the scaled Fiedler vector and normalized cut. After the scaled Fiedler vector f of the generalized eigensystem (4) is available, we sort the vertices of G based on their corresponding values in f and check all possible cutting points for the best partitioning which we define as the one with the smallest normalized cut. We then decide whether to accept or reject the partitioning based on the smallest normalized cut and predefined threshold. If the smallest normalized cut computed is below the threshold, we accept the partitioning. If the smallest normalized cut is above the threshold, it means there are more connections between these two subgraphs. We reject the partitioning and consider them as one cluster. This algorithm is shown below:

Given a weighted graph $G = (V, E)$, its weight matrix is W .

- (1) $d = We$, where $e = (1, 1, \dots, 1)$. $D = \text{Diag}(d)$.
- (2) Solve the generalized eigensystem (4) for the scaled Fiedler vector f .

- (3) Sort the vertices of G based on their corresponding values in f .
- (4) Check every possible cutting point in the order of the sorted vertices, find the one with the smallest normalized cut.
- (5) If the normalized cut is below a certain threshold, accept the partition and recursively partition the sub-graphs. Otherwise, stop.

The running time for this algorithm depends on the number of iterations. Here we analyze the running time for one iteration only. To simplify the notation, we still use E and V to represent the edge set and node set of the graph to be partitioned in each recursion step. Because the running time is dominated by the eigensolver and the sorting procedure, here we consider the running time of steps 2 and 3 only. Using Lanczos method to compute the scaled Fiedler vector f of (4) has complexity roughly proportional to $\text{nnz}(D - W)$, the number of nonzero elements of $D - W$, which is $O(|E|)$. If the number of Lanczos iteration steps is N_r , then the running time for step 2 is $O(N_r|E|)$. In step 3, the sorting of f takes $O(|V|\log(|V|))$. Thus the running time for one recursion is $O(N_r|E| + |V|\log(|V|))$. Note that $|V|$ and $|E|$ become smaller in each subgraph, so the running time is smaller for each recursion as the recursion level goes deeper.

The number of recursion is bounded by the number of clusters we finally obtained, which is very small comparing with $|E|$ and $|V|$, hence the total time for the partition is $O(N_r|E| + |V|\log(|V|))$.

5. Connection with the Cheeger constant

The definition of the normalized cut (2) can be rewritten as

$$\begin{aligned} \text{Ncut}(S, \bar{S}) &= \frac{\text{cut}(S, \bar{S})}{\min(W(S, V), W(\bar{S}, V))} + \frac{\text{cut}(S, \bar{S})}{\max(W(S, V), W(\bar{S}, V))} \\ &= h_G(S) + \frac{\text{cut}(S, \bar{S})}{\max(W(S, V), W(\bar{S}, V))} \end{aligned}$$

where

$$h_G(S) = \frac{\text{cut}(S, \bar{S})}{\min(W(S, V), W(\bar{S}, V))}.$$

The Cheeger constant h_G of graph G is defined as

$$h_G = \min_S h_G(S).$$

The Cheeger constant is one of the earliest normalized forms of edge-cut measuring the quality of the graph partition. It was studied by Cheeger in early 1970s Cheeger (1970). The normalized cut can be regarded as a slight variation of the Cheeger constant.

Defining a constant h_{ncut} for the normalized cut cost function (2) as $h_{\text{ncut}} = \min_S \text{Ncut}(S, \bar{S})$, we show that the upper and lower bounds for the Cheeger

constant h_G can be expressed in terms of h_{ncut} :

$$h_{\text{ncut}} \geq h_G \geq \frac{h_{\text{ncut}}}{2}. \quad (6)$$

Clearly we have $h_{\text{ncut}} \geq h_G$. On the other hand,

$$h_G(S) \geq \frac{1}{2} \left(\frac{\text{cut}(S, \bar{S})}{W(S, V)} + \frac{\text{cut}(S, \bar{S})}{W(\bar{S}, V)} \right) = \text{Ncut}(S, \bar{S})/2$$

which implies $h_G \geq h_{\text{ncut}}/2$. The inequality (6) shows the close relationship between the Cheeger constant and the normalized cut.

6. Connection with the K-means method

In this section, we explore the connection of the normalized-cut based clustering method and the popular K-means clustering method. For ease of discussion, we will only consider the case where the weight matrix W satisfies $W_{ij} \in \{0, 1\}$. Similar conclusions can be drawn for the more general cases.

6.1. K-means method

The general idea of the K-means method is the following: (1) partition the nodes into two arbitrary clusters; (2) for each node re-assign it to the cluster that the node is in some sense closest to; repeat the above two steps until convergence, i.e., until when the cluster membership of no node will change by running through the two steps. One natural and reasonable heuristics for measuring the closeness of a node to a cluster is the following: for a node k , we examine all the nodes that are adjacent to k , if there are more neighbors of k belonging to cluster one, we then re-assign i to cluster one, otherwise we re-assign k to cluster two. It is easy to see that this approach can be generalized to the multiple cluster case: we re-assign i to the cluster in which i has the biggest number of neighbors. If we denote the cluster assignment of k as $x_k \in \{-1, 1\}$, then

$$\sum_{j=1}^n W_{kj} x_j$$

is the difference in the numbers of neighbors of k belonging to the two clusters. Therefore, the re-assignment of k can be computed as

$$x_k^{(i)} = \text{sign} \left(\sum_{j=1}^n W_{kj} x_j^{(i-1)} \right),$$

where i denotes the iteration step. Rewriting in matrix-vector format, we have

$$x^{(i)} = \text{sign}(Wx^{(i-1)}).$$

As will be discussed in more detail later, the above iteration can be considered as a variation of the power method for computing the largest eigenpair of a symmet-

ric matrix Golub and Loan (1989), the difference being that a different normalization process is used here. Restricting $x_k \in \{-1, 1\}$ corresponds to hard-clustering: a node either belongs to cluster one or cluster two. For soft-clustering, i.e., we can consider $x_k \in [0, 1]$ as the probability of node k belonging to cluster one, and use normalization

$$y^{(i)} = Wx^{(i-1)}, \quad x^{(i)} = y^{(i)} / \|y^{(i)}\|_2.$$

This corresponds exactly to the power method for W . In the above discussion we only considered the neighbors of a node k . We can also turn the heuristic argument around: for a node k , we examine all the nodes that are *not* adjacent to k , if there are more of these belonging to cluster one, we then re-assign i to cluster two, otherwise we re-assign k to cluster one. These two arguments can certainly be combined, and we obtain the following re-assignment rule

$$x_k^{(i)} = \text{sign} \left(\sum_{j \sim k} x_j^{(i-1)} - \alpha \sum_{j! \sim k} x_j^{(i-1)} \right), \quad (7)$$

where $\alpha > 0$ is a balancing factor. Here we also use $j \sim k$ to denote that node j is adjacent to node k and $j! \sim k$ to denote that node j is not adjacent to node k .

6.2. Connection between normalized cut and K-means method

The scaled Fiedler vector f of (4) can be obtained by first computing the second largest eigenvector y_2 of (5). We know that $D^{1/2}e$ is the eigenvector of $D^{-1/2}WD^{-1/2}$ corresponding to its largest eigenvalue $\lambda_1 = 1$.⁴ The second largest eigenvector y_2 can be obtained by computing the largest eigenpair of

$$\hat{W} = D^{-1/2}WD^{-1/2} - \frac{D^{1/2}ee^T D^{1/2}}{e^T D e}. \quad (8)$$

One approach for computing the largest eigenpair of a symmetric matrix \hat{W} is the power method Golub and Loan (1989):

Start with a unit vector $x^{(0)}$. For $i = 1, 2, \dots$, until the result converges

$$y^{(i)} = \hat{W}x^{(i-1)}$$

$$x^{(i)} = y^{(i)} / \|y^{(i)}\|_2$$

At the end of convergence, an approximate eigenvalue can be computed as $\hat{\lambda} = (x^{(i)})^T \hat{W}x^{(i)}$.

Applying the power method on \hat{W} , we get the largest eigenvector y_2 of \hat{W} which is the second largest eigenvector of $D^{-1/2}WD^{-1/2}$. Then the scaled Fiedler vector f of (4) is obtained by $f = D^{-1/2}y_2$.

⁴ It can be readily proved that all the eigenvalues of $D^{-1/2}WD^{-1/2}$ are at most 1.

Apply the power method to the matrix in (8), and let the k th entry of $x^{(i)}$ be $x_k^{(i)}$, we have

$$y_k^{(i)} = \sum_{j=1}^n \frac{W_{kj}x_j^{(i)}}{\sqrt{d_k d_j}} - \frac{\sqrt{d_k d_j}x_j^{(i)}}{d_s}$$

where $d_s = e^T d$, the sum of the diagonal elements of D . We can rewrite the above equation as

$$y_k^{(i)} = \sum_{j \sim k} \left(\frac{x_j^{(i)}}{\sqrt{d_k d_j}} - \frac{\sqrt{d_k d_j}x_j^{(i)}}{d_s} \right) - \sum_{j! \sim k} \frac{\sqrt{d_k d_j}x_j^{(i)}}{d_s}.$$

Now comparing the above with the K-means iteration in (7), we notice the following: for the neighbors of node k , the normalized-cut method uses the modified weight

$$\frac{1}{\sqrt{d_k d_j}} - \frac{\sqrt{d_k d_j}}{d_s}$$

and for those nodes not adjacent to k , normalized-cut uses

$$\frac{\sqrt{d_k d_j}}{d_s}$$

as the weight. Therefore, the normalized-cut method can be considered as a variation of the K-means method.

7. Experiments

We tested the algorithm on the link graphs of queries *amazon*, *star* and *apple*. We choose these three query terms because they all have multiple distinct meanings. *Amazon* has at least three meanings. One is related to amazon.com, one of the largest on-line shopping web sites. Another is the famous rain forest in South America. The third is the name of ancient female warriors from Alecto, a female ruled monarchy. For the query *star*, we can think about a natural luminous body visible in the sky, a movie star, a famous athlete and the movie *Star Wars*. Mentioned *apple*, we will associate it with a kind of fruit or the apple computer.

We then apply our clustering algorithm on the data sets. Since each cluster often has a large number of web documents, we choose only the most important web documents among a cluster. The most important or authoritative web documents in each cluster are determined using the HITS algorithm introduced in Appendix A.

We also compare the results obtained from our algorithm to the results by recursively applying simple K-means based algorithm. As our results will show, the normalized-cut based technique performs better than the K-means based algorithm. The stop criteria for the K-means algorithm can be either the maximal iteration number reached, the result is stable, or the cluster size is below a certain threshold. In our implementation, we choose a threshold on maximal iteration number to be our stop criteria. The K-means

method is applied recursively to each subgraph obtained until the size of the subgraph is below certain threshold.

7.1. Data preparation

To obtain the link graphs, we first provide a text-based search engine *hotbot* with query terms. *Hotbot* returns as the query result a list of URLs with highest ranked web documents. We limit the number of returned URLs to be 120 which form the root set. By limiting the number of returned URLs to be 120, we can keep the overall data set within a reasonable size. Then we expand the root set by including all web documents that point to a web document in the root set and those pointed to by a web document in the root set. This is a level one expansion. The link graph can be viewed as a directed graph. It is easy to convert it to the adjacency matrix of an undirected graph.

After the full list of URLs is available, we run a web crawler to get the text information of these web documents. The text of a document is obtained using a web crawler that we write in *Perl*. To accommodate the vast differences in web document lengths, we limit the length of each document to be 500 words. The rest of the document is discarded if it has more than 500 words. Stopwords (such as *I*, *is*, *a*, etc.) are discarded using a standard list. Words are stemmed using Porter stemming Porter (1980), so the words *linking* and *linked* are both stemmed to the same root *link*. Once the text information of all documents is available, we compute the document similarity as in Section 3.2.

At this point, we have the necessary data for the experiment. In our experiment, we set α in expression (1) to be 0.5.

7.2. Clustering results

We apply our algorithm on the data sets of three query terms with the threshold of normalized cut set to 0.06, then run the HITS algorithm on each cluster obtained. The top authorities of each significant clusters are listed. By significant we mean the size of the cluster is not too small.

The order of the clusters has no special meanings. It is merely the sequence by which we process the clusters.

7.2.1. Query term: *amazon*

There are a total of 2294 URLs in this data set. Applying our algorithm, we obtain the following significant clusters. The clusters with small size are not counted. Some small clusters exist as the result of the algorithm because they form the connected components of the link graph. It shows that our algorithm has the capability to find the cluster of small size but tightly connected.

Cluster 1.

- www.amazon.com/
- www.amazon.co.uk/
- www.amazon.de

These three web documents are the home pages of amazon.com, one of the most famous shopping companies in the world. The first website is located in the USA, the second in the UK and third in Germany. It is perfect to cluster them together. We list only three authorities here because the rest of the documents ranked among the top 10 have very very low authority weights compared with the first three, so we do not list them here.

Cluster 2.

- www.amazoncity.com/
- www.amazoncityradio.com/
- www.amazoncity.com/spiderwoman/
- radio.amazoncity.com/
- www.wired.com/news/news/culture/story/6751.html

This cluster is about female issue.

Cluster 3.

- www.amazon.org/
- www.amazonfembks.com/
- www.igc.apc.org/women/bookstores/
- www.teleport.com/~rockyl/queer.shtml
- www.advocate.com/html/gaylinks/resources.html

The topic of this cluster is on female related issues: bi-sexuality and female books.

Cluster 4.

- [sothebys.amazon.com/execlvarzealtg/special-sales/...](http://sothebys.amazon.com/execlvarzealtg/special-sales/)
- sothebys.amazon.com/execlvarzealsubstlhome/sothebys.html...
- [sothebys.amazon.com/execlvarzealtg/special-sales/...](http://sothebys.amazon.com/execlvarzealtg/special-sales/)
- s1.amazon.com/execlvarzealsubstlhome/home.html...
- sothebys.amazon.com/execlvarzealsubstlhome/sothebys.html...

All five authorities listed here are web documents of a large on-line auction company formed by Sothebys and amazon.com. It is not clustered into Cluster 1 because there are relatively few links between them.

Cluster 5.

- www.swalliance.com/
- timeline.ehostation.com
- www.ehostation.com:8080/~1
- downtime.ehostation.com
- rpg.ehostation.com/

The topic of this cluster is about *Star Wars*, surprisingly. But there are a total of 68 documents on this topic, most created by *Star Wars* fans. Some are on-line shopping companies selling goods related to the movie *Star Wars*.

There are two clusters, but for each of them, the web documents are from the same site: www.langenberg.com and www.latingrocer.com, respectively. These two clusters

are:

Cluster 6.

- misc.langenberg.com/
- cooking.langenberg.com/
- shipping.langenberg.com/
- money.langenberg.com/
- weather.langenberg.com/

and

Cluster 7.

- www.latingrocer2.com
- www.latingrocer.com
- www.latingrocer.com/Pages/customer.html
- www.latingrocer.com/Pages/privacy.html
- www.latingrocer.com/Pages/contact.html

Cluster 8.

- www.internext.com.br/lariau

This cluster has only one authority. Other web documents in this cluster all point to it. It is a website of a hotel in the Amazon River valley.

After applying the K-means based algorithm on the data, we also get several clusters. Cluster 1 is separated into three clusters, that is, the websites in the three different countries form three clusters under this algorithm. Second, there are more clusters with small size than the result from our algorithm. Many of them should belong to larger clusters, but were partitioned incorrectly.

By adjusting the threshold, the algorithm can group or separate the following web documents:

- www.amazon.com/
- www.amazon.co.uk/ and
- www.amazon.de

while the K-means based algorithm can only partition them into different clusters. This verifies that our algorithm is not only more flexible than the other, but also clusters more reasonably.

From the clusters we obtained, we found that, unlike what we had expected, no cluster has a focused topic on the rain forest while Cluster 8 does have something to do with the rain forest in Brazil. Checking the entire data set, we only found a couple of web documents that mention the rain forest. They do not form a cluster with significant size. If we return to *hotbot* and enter the query *amazon*, the web sites presented are dominantly related to amazon.com. The documents about the rain forest are not among the highest-ranked list returned by *hotbot*. It means the number of web documents on amazon the rain forest is small, which makes the web documents about this topic have low rank weight. That is the reason that we cannot have a cluster on this topic.

As for the third meaning of amazon, although female warrior dose not directly appear as a distinct topic in any cluster, Cluster 3 focused on female topics, or even the bi-sexual issue. By examining the content of these documents, we are sure that these issues are the extent of the original meaning of amazon as female warriors.

There are clusters with small size that apparently have nothing to do with *amazon*, such as the clusters about websites *langenberg.com* and *latingrocer.com*. They are explored as separate clusters because the web documents on the same site point to each other, raising the importance themselves. To avoid such situations, before applying any clustering algorithm, we can coarsen the link graph first, so that the web documents from the same site collapse to one node in the graph.

7.2.2. Query term: *star*

We have 3504 URLs for this query. Setting the threshold to be 0.06, we run our algorithm on the data set of query term *star*. The authorities of each cluster are listed below:

Cluster 1.

- www.starwars.com/
- www.lucasarts.com/
- www.sirstevesguide.com/
- www.jediknight.net/
- www.surfthe.net/swmal/

This cluster is focused on *star wars*.

Cluster 2.

- www.kcstar.com/
- www.dailystarnews.com/
- www.kansascity.com/
- www.starbulletin.com/
- www.trib.com/

This cluster includes the web documents of some news media with the word *star* as part of their names.

Cluster 3.

- www.weatherpoint.com/starnews
- www.starnews.com/digest/sports.html
- www.starnews.com/digest/citystate.html
- www.indy.com
- speednet.starnews.com/

The top authorities are all web documents of *starnews.com* in this cluster.

Cluster 4.

- www.state.mn.us/mainmenu.html
- www.mda.state.mn.us/
- www.doli.state.mn.us/

- www.legalethics.com/palstates/state/mn.htm
- www.exploreminnesota.com

This cluster's topic is the state of Minnesota. The reason that Minnesota is a topic of the cluster under query *star* is because the official State of Minnesota web site is called *North Star*, which is named second among all government web sites.

Cluster 5.

- www.star-telegram.com/
- www.dfw.com/
- www.virtualtexan.com/
- marketplace.dfw.com
- www.star-telegram.com/advertiselvshops/

A cluster of star-telegram.com located in Texas.

Cluster 6.

- www.starpages.net/

There is only one authority in this cluster. This authority is a web site introducing preeminent persons in various fields, such as movie star, sport stars, etc. This cluster is what we have expected to be obtained after partitioning.

Cluster 7.

- www.aavso.org/
- www.astro.wisc.edu/dolan/constellations/
- ourworld.compuserve.com/homepages/rahwide_home_page
- adc.gsfc.nasa.gov/adc/adc_amateurs.html
- heasarc.gsfc.nasa.gov/docs/www_info/webstars.html

This cluster talks about space and astronomy. It is also what we have expected. Now all three meanings of *star* are found to be topics of separate clusters.

Clusters 2, 3, 5 are all web documents of news media. They are partitioned to be different clusters since there is no link among the three clusters.

Using the K-means based algorithm, we obtained many more clusters than obtained by our algorithm. The same as the result of the query *amazon* under the K-means based algorithm, many of them are not really clusters with focused topics. On the other hand, Clusters 2 and 4 are grouped together, although there is no link between them. This is because the K-means based algorithm is not global. It is easy to be trapped to a local minimum while the first algorithm avoids this situation by using the scaled Fiedler vector to partition.

In Fig. 2, graph (a) is the original weighted link graph of data for the query *star*, (b) is the re-ordered graph after clustering using our algorithm, (c) is the re-ordered graph after clustering using the K-means based algorithm. In graph (b) and (c), each diagonal block corresponds to a resulting cluster. Comparing graphs (b) and (c), we can see that, in graph (b), the off-diagonal blocks are much sparser than off-diagonal blocks in graph (c). This means our algorithm creates clusters with much fewer connections between them than the K-means based algorithm does.

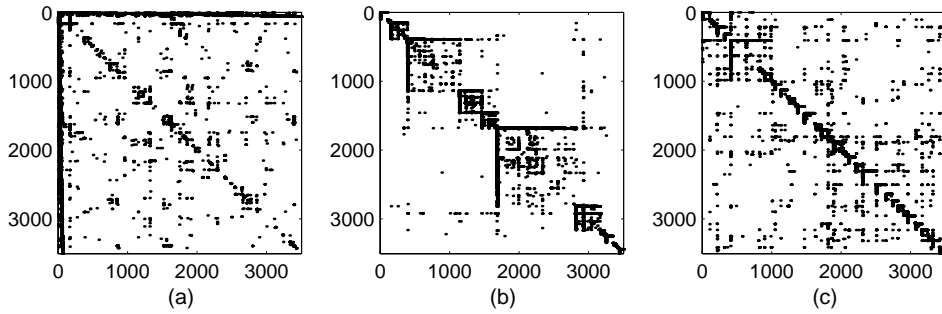


Fig. 2. Clustering result of query *star*. (a) Weighted link graph of *star*. (b) Clustering result using our algorithm. The threshold is 0.06. (c) Clustering result using K-means based algorithm.

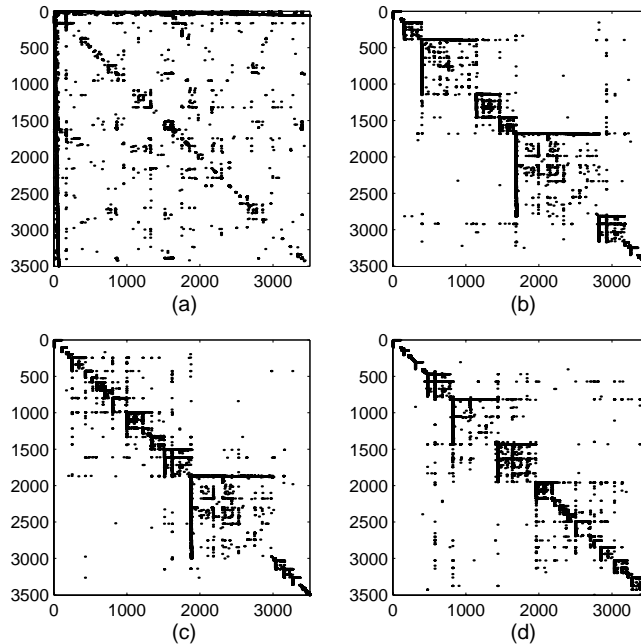


Fig. 3. Clustering result of query *star* using our algorithm. (a) Weighted link graph of *star*. (b) Clustering result with threshold=0.06. (c) Clustering result with threshold=0.1. (d) Clustering result with threshold=0.2.

In Fig. 3, graph (a) is the original weighted link graph of data for the query *star*, (b) is the re-ordered graph after clustering using our algorithm with the threshold equal to 0.06, (c) is the re-ordered graph after clustering using our algorithm with the threshold equal to 0.1, (d) is the re-ordered graph after clustering using our algorithm with the threshold equal to 0.2. The graph clearly shows that by increasing the threshold, we can obtain more clusters with the size of each cluster smaller. In this way, we can

control the clustering result by changing the value of the threshold in the algorithm. It makes this algorithm more flexible than the K-means based one.

7.2.3. Query term: *apple*

In this data set, there are 2757 URLs returned by the search engine. The threshold is still 0.06. The authorities of each cluster are listed below after running the algorithm on the data of the query term *apple*:

Cluster 1.

- www.apple.com/
- www.apple.com/support/
- www.apple.com/education/
- www.apple.com/quicktime/
- www.apple.com/hotnews/

Here all top authorities are from the same web site: www.apple.com. This cluster is dominant in this query. Most URLs belong to it. After running the HITS algorithm with the site information considered, that is, the URLs from the same web site are collapsed to one node, we obtain different top authorities:

- www.apple.com/powermacserver/
- www.claris.com/
- www.apple.rulhardware/displays
- www.cs.brandeis.edu/~xray/oldmac.html
- www.next.com/

The second URL in this cluster is the website of a computer software company. It produces software used for the Macintosh. The other four URLs are all about the apple computer. This list of authorities provides more useful information than the previous one does.

Cluster 2.

- www.yabloko.rul
- www.cityline.rulpolitikal
- www.russ.rul
- www.forum.msk.rul
- www.novayagazeta.ru

All web documents in this cluster are written in Russian.

Cluster 3.

- www.michiganapples.com/

This cluster has only one authority. It is related to apple, the fruit.

The following cluster is formed in this query simply because its name happens to contain the query term *apple*:

Cluster 4.

- www.ci.apple-valley.mn.us/

which is the website of the city of Apple Valley, MN.

Cluster 5.

- www.valleyweb.com/

This is the website of Annapolis Valley in Canada where THERE is the Apple Blossom Festival in the spring celebrating the traditions and agricultural heritage.

There are not many clusters formed by the algorithm, nor do we find very interesting topics which are beyond our expectation.

The performance of the K-means algorithm is mixed. On one hand it does not form a cluster as large as Cluster 1. For example, it forms a new cluster around www.france.euro.apple.com/ which belongs to Cluster 1 obtained with our algorithm. On the other hand, it groups totally different web documents together. For instance, the top authorities of one cluster are:

- www.jokewallpaper.com/
- the-tech.mit.edu/Macmadel/
- www.geocities.com/SiliconValley/Vista/7184/guitool.html

But www.jokewallpaper.com/ and the-tech.mit.edu/Macmadel/ mention totally different things.

8. Further discussion*8.1. Importance of text information*

Intuitively, the links in the webgraph can not be regarded as equally important, specifically many web links are created much less carefully than the references in scientific literature. The text similarity between two web documents provide us with a useful metric to address this issue. Without the similarity as the measure of the link strength, that is, if we form the weight matrix as

$$W = \alpha \frac{A}{\|A\|_2} + (1 - \alpha) \frac{C}{\|C\|_2},$$

then we unduly raise the strength of some links which connect two web documents with little in common. When applied to the data set *amazon*, our algorithm groups Clusters 1, 3 and 4 into one single cluster, using the same threshold as the stopping criterion, namely 0.06 in our previous experiments. Cluster 3 addresses the female issues and Clusters 1 and 4 are related to *amazon.com*. Apparently this clustering result is not a good one. This result justifies our choice of incorporating the text information into the weight metric.

Recall that our weight matrix is formed as

$$W = \alpha \frac{A \otimes S}{\|A \otimes S\|_2} + (1 - \alpha) \frac{C}{\|C\|_2}.$$

Table 1
 (α, T) pairs that separate three authoritative web documents in Cluster 1 of query *amazon*

| | | | | | | | | | | |
|----------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| α | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
| T | 1.1 | 1.0 | 0.8 | 0.7 | 0.6 | 0.5 | 0.5 | 0.4 | 0.3 | 0.3 |

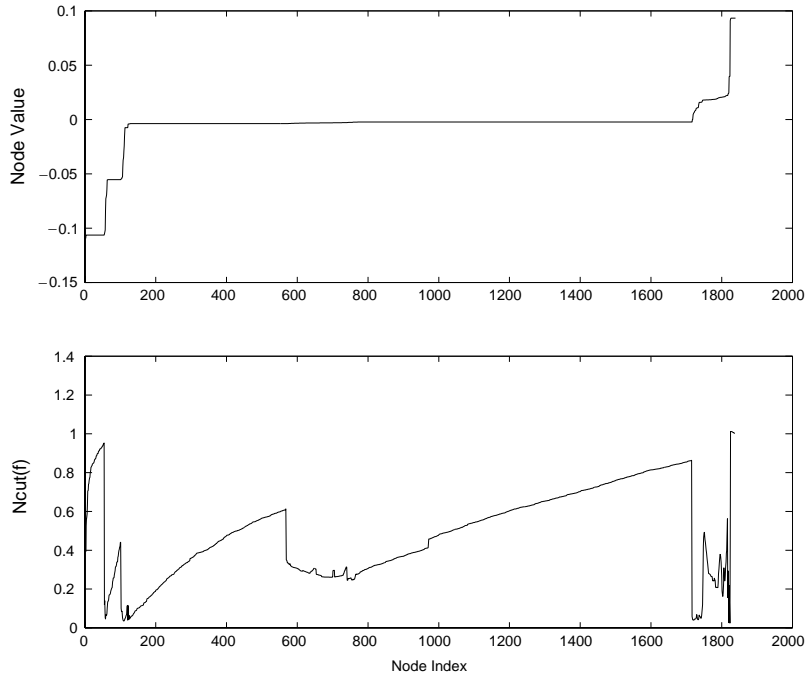


Fig. 4. The sorted scaled Fiedler vector value(top) and corresponding normalized cut(bottom).

If the value of α changes, the weight matrix will change accordingly. Then we can set different thresholds in the algorithm to have three authoritative web documents in Cluster 1 of the data set *amazon* separated into different clusters. Table 1 lists the α value and the corresponding threshold T such that these three web documents are separated by the algorithm. From the (α, T) pairs in the table, we see clearly that when the value of α increases, to partition these three web documents into different clusters, the threshold decreases monotonically.

8.2. Scaled Fiedler vector

Now we use a concrete example to show that partitioning a graph with the scaled Fiedler vector and normalized cut is feasible. Here we use the subgraph which includes URLs in Clusters 1, 3 and 4 of the data set *amazon*. There are a total of 1839 web documents in this subgraph. From the discussion above, we know that this subgraph contains multiple topics and can be produced by using the weight matrix without the text information. Fig. 4 (top) is the sorted scaled Fiedler vector of the subgraph. The

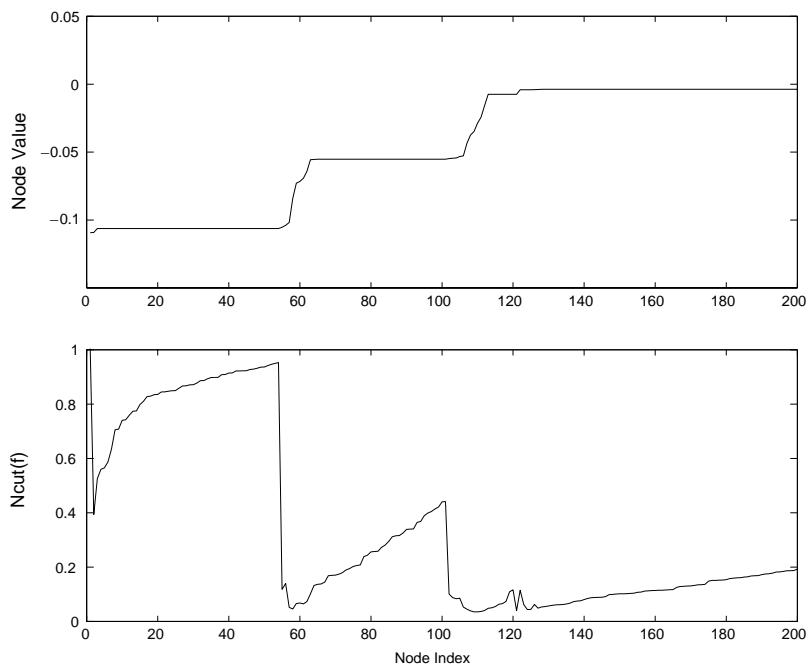


Fig. 5. The sorted scaled Fiedler vector value(top) and corresponding normalized cut(bottom) for the first 200 nodes.

horizontal axis is the node index and the vertical axis is the corresponding entry of the scaled Fiedler vector. Fig. 4 (bottom) is the node index and corresponding normalized cut value. We see that the sorted scaled Fiedler vector behaves like step function. Each segment takes nearly constant value. The jumping points can be used as the cutting points because they correspond to small normalized cut values. See Fig. 4 (bottom) for their values. After we have determined the threshold, we can obtain the cutting points, and partition the graph into clusters. If we set the threshold to be 0.06, we obtain the following cutting points:

(57 58)
 (106 107 108 109 110 111 112 113 121 124 125 127 128 129 130 131)
 (1716 1717 1718 1719 1720 1721 1722 1723 1724 1725 1726 1727 1728 1731
 1732 1733 1734 1735 1739 1740 1741 1742)
 (1820 1823)

The numbers in the same pair of () are almost consecutive, we can choose any one of them as the cutting point. This will not affect the resulting clusters significantly since these nodes correspond to URLs which are not important, i.e., not good authorities nor good hubs.

If we enlarge Fig. 4 to include only the first 200 nodes, we can see the result more clearly in Fig. 5.

For convenience, we choose the cutting points which correspond to the smallest normalized cut in each group: 58 110 1721 1820. This gives us 5 groups:

- (1) 1–57: with top authority *amazon.org*
- (2) 58–109: with top authority *fembks.com*
- (3) 110–1720: corresponds to Cluster 1 of *amazon*
- (4) 1721–1819: corresponds to Cluster 4 of *amazon*
- (5) 1820–1839: with top authority *ethnobotany.org*

Groups 1 and 2 correspond to Cluster 3. They are separated here because we compute the normalized cut value at point 58 with respect to the whole subgraph of 1839 nodes. The value is below the threshold 0.06. While in our algorithm we first partition this subgraph and obtain a smaller subgraph with nodes from 1 to 109, then continue to partition this subgraph. At this point, the smallest normalized cut value is still obtained at node 58, but with the value of 0.0715 which is above the threshold. The partition is rejected. So groups 1 and 2 form one cluster in our algorithm. The result here indicates that our algorithm may be improved by checking the points with the normalized cut values below the threshold and partition the graph into several subgraphs at the same time.

Our algorithm does find the cluster corresponding to group 5. Since its size is too small, we did not list it in our results.

8.3. Robustness

Fig. 6 is the plot of the α value in Expression (1) and the corresponding average normalized cut value. The graph shows that for the query term *amazon*(top), the smallest average normalized cut value corresponds to $\alpha = 0.7$ while for the query term *star*(bottom), $\alpha = 0.1$ leads to the smallest average normalized cut value. From this figure, we know that we cannot choose an α value or a range of α values that can minimize the average normalized cut in all cases.

For the query *amazon*, $\alpha = 0.7$ gives the best average normalized cut. We choose this α value to test the average normalized cut under different threshold. In case of the query *star*, it is $\alpha = 0.1$ that gives the best average normalized cut. But we cannot make α too small, since doing so will lower the importance of text similarity. So we choose $\alpha = 0.4$ to test because it is a local minimum. From Fig. 7 we see that, as the threshold increases, the average normalized cut value increases, too. The relationship between them are almost linear.

8.4. Multi-level partition method

We also applied multi-level K-way partition method METIS Karypis and Kumar. The result is not quite as good. (In fact, we use normalized cut method after working with METIS for a while.) METIS tends to produce subgraphs with nearly equal sizes, while sizes of different clusters are generally different in clustering, especially for web graph. More fundamentally, it appears that the node contraction procedure in METIS is not suitable for highly random graphs like webgraphs.

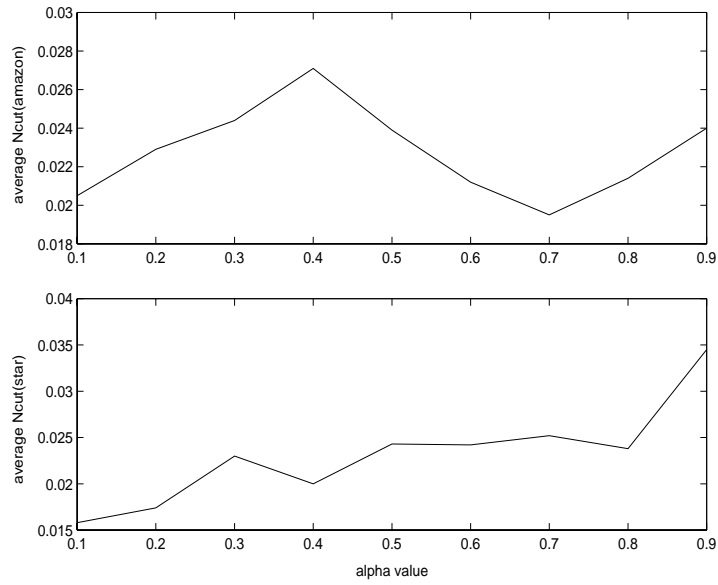


Fig. 6. α value and corresponding average normalized cut (top):*amazon*, (bottom):*star*.

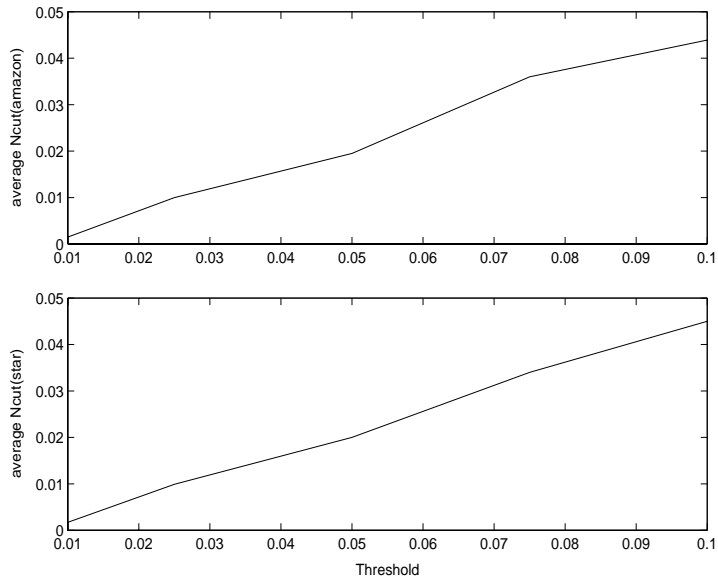


Fig. 7. Threshold and corresponding average normalized cut (top):*amazon* $\alpha = 0.7$, (bottom):*star*, $\alpha = 0.4$.

9. Concluding remarks

In this paper, we present an algorithm to solve the web document clustering problem. We treat the problem as graph partitioning, measure the partitioning result using the *normalized cut* criterion which was first proposed in the field of image segmentation. Combining normalized cut and the scaled Fiedler vector together, this approach forms a global, unbiased algorithm which can effectively extract different topics contained in the webgraph. Compared with the K-means based algorithm, it avoids creating clusters with small size by controlling the threshold on the value of normalized cut. The clusters obtained have high similarity within clusters and dissimilarity between clusters. In our experiment, after choosing suitable threshold, we obtain the clusters, each with distinct topics.

Acknowledgements

We thank Daniel Denton of the CITG group, Lawrence Berkeley National Laboratory for the proofreading of this paper.

Appendix A

Here we briefly introduce Kleinberg's HITS algorithm to find the authorities and hubs of each cluster we obtained. Kleinberg, 1998 defines the *authorities* as the most relevant documents for the topic. The *Hubs* are defined as the web documents which link to many related authorities. They implicitly represent an "endorsement" of the authorities they point to. The authority and hub information can be retrieved based entirely on the link structure information. Since a good *authority* is pointed to by many good *hubs* and a good *hub* points to many good authorities, such mutually reinforcing relationship can be represented as:

$$x_p = \sum_{q:(q,p) \in E} y_q, \quad (\text{A.1})$$

$$y_p = \sum_{q:(p,q) \in E} x_q, \quad (\text{A.2})$$

where x_p is the authority weight of web document x and y_p is the hub weight. E is the set of links(edges). Iteratively update the authority and hub weights of every web document, using Eqs. (A.1) and (A.2), and sort the web documents in decreasing order according to their authority and hub weights, respectively, we can obtain the authorities and hubs of the topic. Many other issues need to be taken into consideration, such as the hyperlinks from the same web site, etc.

Let A be the adjacency matrix of the link graph, x be the vector of authority weights and y be the vector of hub weights, then (A.1) and (A.2) can be transformed to

$$x = A^T y, \quad y = Ax$$

which leads to

$$x = A^T A x, \quad y = A A^T y$$

Clearly, x and y are the principal eigenvectors of $A^T A$ and $A A^T$, respectively.

References

- Anick, P.G., 1994. Adapting a full-text information retrieval system to compute the troubleshooting domain. *Proceedings of ACM SIGIR '94*, pp. 349–358.
- Bharat, K., Broder, A., 1998. A technique for measuring the relative size and overlap of public web search engines. *Proceedings of the Seventh World-Wide Web Conference (WWW7)*, pp. 379–388.
- Chakrabarti, S., Dom, B.E., Raghavan, P., Rajagopalan, S., Iyengar, D., Kleinberg, J., 1998. Automatic resource compilation by analyzing hyperlink structure and associated text. *Comput. Networks ISDN Systems* 30 (1–7), 65–74.
- Chakrabarti, S., Dom, B.E., Gibson, D., Kleinberg, J.M., Kumar, R., Raghavan, P., Rajagopalan, S., Tomkins, A., 1999. Mining the link structure of the world wide web. *32(8):60–67*.
- Cheeger, J., 1970. A Lower Bound for the Smallest Eigenvalue of the Laplacian, *Problems in Analysis*. Princeton University Press, Princeton, NJ.
- Chung, F.R.K., 1997. *Spectral Graph Theory*. American Mathematical Society, Providence, RI.
- Croft, W.B., Cook, R., Wilder, D. Providing government information on the internet: Experience with 'thomas'. Technical Report, University of Massachusetts, pp. 95–45.
- Donath, W.E., Hoffman, A.J., 1972. Algorithms for partitioning of graphs and computer logic based on connected matrices. *IBM Tech. Disclosure Bull.* 15, 938–944.
- Efthimiadis, E.N., 1993. A user-centered evaluation of ranking algorithms for interactive query expansion. *Proceedings of ACM SIGIR '93*, pp. 146–159.
- Everitt, B., 1993. *Cluster Analysis*. Edward Arnold, London.
- Fiedler, M., 1973. Algebraic connectivity of graphs. *Czechoslovak Math. J.* 23, 298–305.
- Fiedler, M., 1975. A property of eigenvectors of non-negative symmetric matrices and its application to graph theory. *Czechoslovak Math. J.* 25, 619–633.
- Flake, G.W., Lawrence, S., Giles, C.L., 2000. Efficient identification of web communities, *SIGKDD*. pp. 150–160.
- Frieze, A., Kannan, R., Vempala, S., 2000. Fast Monte-Carlo methods for finding low-rank approximations. <http://www.cs.yale.edu/homes/kannan/Papers/pubs.html>.
- Gibson, D., Kleinberg, J., Raghavan, P., 1998. Inferring web communities from link topology. *Proceedings of the Ninth ACM Conference on Hypertext and Hypermedia (HYPER-98)*, June 20–24, ACM Press, New York, pp. 225–234.
- Golub, G., Loan, C.V., 1989. *Matrix Computations*, 2nd Edition. Johns Hopkins, Baltimore, MD.
- Gordon, A.D., 1981. *Classification*. Chapman & Hall, London.
- Hearst, M.A., Paderson, J.O., 1996. Re-examining the cluster hypothesis: Scatter/gather on retrieval results. *Proceedings of the SIGIR'96*, pp. 246–255.
- Hendrickson, B., Leland, R., 1995. An improved spectral graph partitioning algorithm for mapping parallel computations. *SIAM J. Sci. Comput.* 16 (2), 452–469.
- Karypis, G., Kumar, V. Metis* a software package for partitioning unstructured graphs, partitioning meshes, and computing fill-reducing orderings of sparse matrices. <http://www.cs.umn.edu/~karypis>.
- Kleinberg, J.M., Kumar, R., Raghavan, P., Rajagopalan, S., Tomkins, A.S., 1999. The web as a graph: measurements, models, and methods. *Proceedings of the Fifth Annual International Computing and combinatorics Conference*, 1999. pp. 26–28.
- Kleinberg, J.M., 1998. Authoritative sources in a hyperlinked environment. *Proceedings of Ninth ACM-SIAM Symposium on Discrete Algorithm*, 25–27 January, pp. 668–677.
- Kumar, R., Raghavan, P., Rajagopalan, S., Tomkins, A., 1999. Extracting large-scale knowledge bases from the web. *Proceedings of the 25th VLDB Conference*. pp. 639–650.

- Larson, R.R., 1996. Bibliometrics of the world wide web: an exploratory analysis of the intellectual structures of cyberspace. *Proceedings of the SIGIR'96*, pp. 71–78.
- Li, Y., 1998. Towards a qualitative search engine. *IEEE Internet Comput.*, 2(4):24–29.
- Mohar, B., 1992. Laplace eigenvalues of graphs—a survey. *Discrete Math.* 109, 171–183.
- Pirolli, P., Pitkow, J., Rao, R., 1996. Silk from a sow's ear: extracting usable structures from the web. *Proceedings of the SIGCHI'96*, pp. 118–125.
- Porter, M.F., 1980. An algorithm for suffix stripping. *Program* 14, 130–137.
- Pothen, A., Simon, H.D., Liou, K.P., 1990. Partitioning sparse matrices with eigenvectors of graph. *SIAM J. Matrix Anal. Appl.* 11, 430–452.
- Rijsbergen, C.J.V., 1979. *Information Retrieval*. Butterworths, London.
- Shi, J., Malik, J., 1997. Normalized cuts and image segmentation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June. pp. 731–737.
- Small, H., 1973. Co-citation in the scientific literature: A new measure of the relationship between two documents. *J. Amer. Soc. Inform. Sci.* 24, 265–269.
- Spielman, D.A., Teng, S., 1996. Spectral partitioning works: planar graphs and finite element meshes. *IEEE Symposium on Foundations of Computer Science*, pp. 96–105.
- Willett, P., 1988. Recent trends in hierarchical document clustering. *Inform. Process. Manage.* 24, 577–597.
- Zamir, O., Etzioni, O., 1999. Grouper: a dynamic clustering interface to web search results, WWW8.