

A Divide-and-Merge Methodology for Clustering

David Cheng*
Massachusetts Institute of Technology
Cambridge, MA 02139
drcheng@mit.edu

Santosh Vempala*
Massachusetts Institute of Technology
Cambridge, MA 02139
vempala@mit.edu

Ravi Kannan†
Yale University
New Haven, CT 06520
kannan@cs.yale.edu

Grant Wang*
Massachusetts Institute of Technology
Cambridge, MA 02139
gjw@mit.edu

ABSTRACT

We present a divide-and-merge methodology for clustering a set of objects that combines a top-down “divide” phase with a bottom-up “merge” phase. In contrast, previous algorithms either use top-down or bottom-up methods to construct a hierarchical clustering or produce a flat clustering using local search (e.g., k -means). Our divide phase produces a tree whose leaves are the elements of the set. For this phase, we use an efficient spectral algorithm. The merge phase quickly finds an optimal tree-respecting partition for many natural objective functions, e.g., k -means, min-diameter, min-sum, correlation clustering, etc.. We present a meta-search engine that uses this methodology to cluster results from web searches. We also give empirical results on text-based data where the algorithm performs better than or competitively with existing clustering algorithms.

1. INTRODUCTION

The rapidly increasing volume of readily accessible data presents a challenge for computer scientists: find methods that can locate relevant information and organize it in an intelligible way. This is different from the classical database problem in at least two ways. First, there may neither be the time nor (in the long term) the computer memory to store and structure all the data (e.g., the world-wide web or a portion of it) in a central location. Second, one would like to find interesting patterns in the data without knowing what to look for in advance.

Clustering refers to the process of classifying a set of data objects into groups so that each group consists of similar

objects and objects from different groups are dissimilar. The classification could either be flat (a partition of the data set usually found by a local search algorithm such as k -means [14]) or hierarchical [16]. Clustering has been proposed as a method to aid information retrieval in many contexts (e.g., [9, 31, 28, 21, 12]). Document clustering can help generate a hierarchical taxonomy efficiently (e.g., [6, 35]) as well as organize the results of a web search (e.g., [33, 32]). It has also been used to learn (or fit) mixture models to data sets [15] and for image segmentation [30].

Most hierarchical clustering algorithms can be described as either divisive methods (i.e., top-down) or agglomerative methods (i.e., bottom-up) [2, 16, 17]. Both methods create trees, but do not provide a flat clustering. A divisive algorithm begins with the entire set and recursively partitions it into two (or more) pieces, forming a tree. An agglomerative algorithm starts with each object in its own cluster and iteratively merges clusters. We combine top-down and bottom-up techniques to create both a hierarchy and a flat clustering. In the divide phase, in principle we could apply any divisive algorithm to form a tree T whose leaves are the objects. This is followed by the merge phase in which we start with each leaf of T in its own cluster and merge clusters going up the tree. The final clusters form a partition of the data set and are tree-respecting, i.e., each cluster is the complete subtree rooted at some node of T . For a large class of natural objective functions, the merge phase can be executed optimally, producing the best tree-respecting clustering. Figure 1 shows a depiction of the methodology.

For the divide phase we use the spectral algorithm studied in [18]. There, the authors use a quantity called *conductance* to define a measure of a good clustering based on the graph of pairwise similarities. They prove that the tree constructed by recursive spectral partitioning contains a partition that has reasonable worst-case guarantees with respect to conductance. However, the running time for a data set with n objects could be $O(n^4)$. We describe an efficient implementation of this algorithm when the data is presented in a document-term matrix and the similarity function is the inner product. For a document-term matrix for n objects with M nonzeros, our implementation runs in $O(Mn \log n)$ in the worst case and seems to perform much better in practice (see Figure 2(a)). The data need not be text; all that is needed is for the similarity of two objects to be the inner

*Supported in part by NSF ITR Award CCR-0312339

†Supported in part by NSF Award CCR-0312354

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS 2005 Baltimore, Maryland USA

Copyright 2005 ACM 1-59593-062-0/05/06 ... \$5.00.

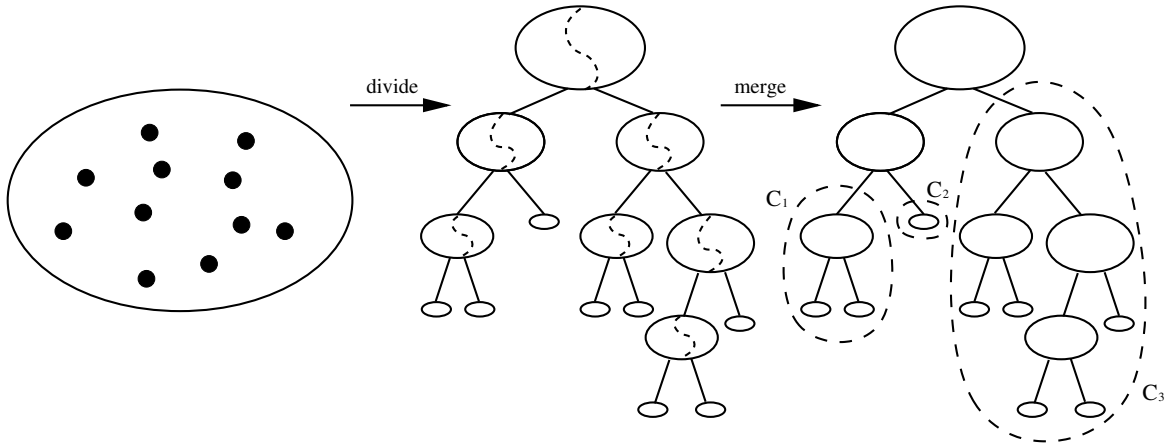


Figure 1: The Divide-and-Merge methodology

product between the two vectors representing the objects.

The class of functions for which the merge phase can find an optimal tree-respecting clustering include standard objectives such as k -means [14], min-diameter [8], and min-sum [24]. It also includes correlation clustering, a formulation of clustering that has seen recent interest [3, 7, 11, 13, 29]. Each of the corresponding optimization problems is NP-hard to solve for general graphs. Although approximation algorithms exist for these problems, many of them have impractical running times. Our methodology can be seen as an efficient alternative.

We show promising empirical results for the methodology. The first application is a meta-search engine, EigenCluster [1], that clusters the results of a query to a standard web search engine. EigenCluster consistently finds the natural clustering for queries that exhibit polysemy, e.g., for the query *monte carlo*, EigenCluster finds clusters pertaining to the car model, the city in Monaco, and the simulation technique. We describe EigenCluster and show results of example queries in Section 3. We also apply the methodology to clustering text-based data whose correct classification is already known. In Section 4, we describe the results of a suite of experiments that show that a good clustering exists in the tree built by the spectral algorithm.

2. DIVIDE-AND-MERGE METHODOLOGY

As mentioned in the introduction, there are two phases in our approach. The divide phase produces a hierarchy and can be implemented using any algorithm that partitions a set into two disjoint subsets. The input to this phase is a set of objects whose pairwise similarities or distances are given (or can be easily computed from the objects themselves). The algorithm recursively partitions a cluster into two smaller sets until it arrives at singletons. The output of this phase is a tree whose leaves are the objects themselves; each internal node represents a subset of the objects, namely the leaves in the subtree below it. Divisive algorithms that can be applied in the divide phase are known for a variety of data representations such as graphs [12] and high-dimensional vectors [6]. In Section 2.1, we use a spectral algorithm for the divide phase when the objects are represented as a document-term matrix and the similarity

between the objects is the inner product between the corresponding vectors.

The merge phase is applied to the tree T produced by the divide phase. The output of the merge phase is a partition C_1, \dots, C_k of the set of objects and each C_i is a node of T . The merge phase uses a dynamic program to find the optimal tree-respecting clustering for a given objective function g . The optimal solutions are computed bottom-up on T ; to compute the optimal solution for any interior node C , we *merge* the optimal solutions for C_l and C_r , the children of C . The optimal solution for any node need not be just a clustering; an optimal solution can be parameterized in a number of ways. Indeed, we can view computing the optimal solution for an interior node as computing a Pareto curve; a value on the curve at a particular point is the optimal solution with the parameters described by the point. A specific objective function g can be efficiently optimized on T if the Pareto curve for a cluster can be efficiently computed from the Pareto curves of its children. In Section 2.2, we describe dynamic programs to compute optimal tree-respecting clusterings for several well-known objective functions: k -means, min-diameter, min-sum, and correlation clustering.

2.1 Divide phase

The spectral algorithm given here deals with the common case in which the objects are given as a sparse document-term matrix A . The rows are the objects and the columns are the features. We denote the i th object, a row vector in A , by $A_{(i)}$. The similarity of two objects is defined as the inner product of their term vectors: $A_{(i)} \cdot A_{(j)}$. The algorithm can easily be applied to the case when the pairwise similarities are given explicitly in the form of a similarity matrix. However, when the similarity function is the inner product, computation of the similarity matrix can be avoided and the sparsity of A can be exploited.

The algorithm constructs a hierarchical clustering of the objects by recursively dividing a cluster C into two pieces through a cut $(S, C \setminus S)$. To find the cut, we compute v , an approximation of the second eigenvector of the similarity matrix AA^T normalized so that all row sums are 1. The ordering of the objects in v gives a set of cuts, and we take the “best” one. The algorithm then recurses on the subparts. To compute the approximation of the second eigenvector, we

Input: An $n \times m$ matrix A .

Output: A tree with the rows of A as leaves.

1. Let $\rho \in \mathbb{R}^n$ be a vector of the row sums of AA^T , and $\pi = \frac{1}{(\sum_i \rho_i)} \rho$.
2. Let $R = \text{diag}(\rho)$, and $D = \text{diag}(\sqrt{\pi})$ be diagonal matrices.
3. Compute the second largest eigenvector v' of $Q = DR^{-1}AA^TD^{-1}$.
4. Let $v = D^{-1}v'$, and sort v so that $v_i \leq v_{i+1}$.
5. Find t such that the cut

$$(S, T) = (\{1, \dots, t\}, \{t+1, \dots, n\})$$

minimizes the conductance:

$$\phi(S, T) = \frac{c(S, T)}{\min(c(S), c(T))}$$

where $c(S, T) = \sum_{i \in S, j \in T} A_{(i)} \cdot A_{(j)}$, and $c(S) = C(S, \{1 \dots, n\})$.

6. Let \hat{A}_S, \hat{A}_T be the submatrices of A . Recurse (Steps 1-5) on \hat{A}_S and \hat{A}_T .

Table 1: Divide phase

use the power method, a technique for which it is not necessary to explicitly compute the normalized similarity matrix AA^T . We discuss this in Section 2.1.1. The algorithm is given in Table 1.

In Step 5, we consider $n - 1$ different cuts and use the cut with the smallest conductance. This itself is an approximation to the minimum conductance of the entire subgraph at this point [18]. Why should the "best" way to divide the set of objects be a cut of small conductance? Imagine a graph where the nodes are the objects and the edges are weighted according to the similarity between two objects. One might think that if the weight crossing a cut is small, the partition induced by the cut is a good candidate. However, just looking at the weight across the cut hides too much information. Indeed, consider two cuts, C_1 and C_2 , with the same weight crossing them. Suppose C_1 partitions the set into two subsets of equal size, both of which hold high weight while C_2 partitions the set into a singleton and the rest of the set. Then C_1 seems to be a more effective cut than C_2 ; the measure of conductance formalizes this idea by normalizing a cut by the smaller weight of the two parts it induces. More intuition for why conductance is a good measure for clustering can be found in [18].

The cut (S, T) we find using the second eigenvector in Step 5 is not the cut of minimum conductance; finding such a cut is NP-hard. However, the conductance of (S, T) is not much worse than the minimum conductance cut.

For a document-term matrix with n objects and M nonzeros, Steps 1-5 take $O(M \log n)$ time. Theoretically, the worst-case time to compute a complete hierarchical clustering of the rows of A is $O(Mn \log n)$. Empirical experiments, however, show that the algorithm usually performs much better (see Section 2.1.2) and seems to be almost linear in M .

2.1.1 Details

Any vector or matrix that the algorithm uses is stored using standard data structures for sparse representation. The main difficulty is to ensure that the similarity matrix AA^T is not explicitly computed; if it is, we lose sparsity and our running time could grow to m^2 , where m is the number of terms. We briefly describe how we avoid this in Steps 1 and 3.

Step 1: Computing row sums. Observe that

$$\rho_i = \sum_{j=1}^n A_{(i)} \cdot A_{(j)} = \sum_{j=1}^n \sum_{k=1}^m A_{ik} A_{jk} = \sum_{k=1}^m A_{ik} \left(\sum_{j=1}^n A_{jk} \right).$$

Because $\sum_{j=1}^n A_{jk}$ does not depend on i , we can compute $u = \sum_{i=1}^n A_{(i)}$ so we have that $\rho_i = A_{(i)} \cdot u$. The total running time is $\theta(M)$ and the space required is $\theta(n + m)$.

Step 3: Computing the eigenvector. The algorithm described in [18] uses the second largest eigenvector of $B = R^{-1}AA^T$, the normalized similarity matrix, to compute a good cut. To compute this vector efficiently, we compute the second largest eigenvector v of the matrix $Q = DBD^{-1}$. The eigenvectors and eigenvalues of Q and B are related; if $Bv = \lambda v$, then $Q(Dv) = \lambda Dv$.

The matrix Q is symmetric; it is easy to see this from $D^2B = B^TD^2$. Therefore, we can compute the second largest eigenvector of Q using the power method, an iterative algorithm whose main computation is a matrix-vector multiplication.

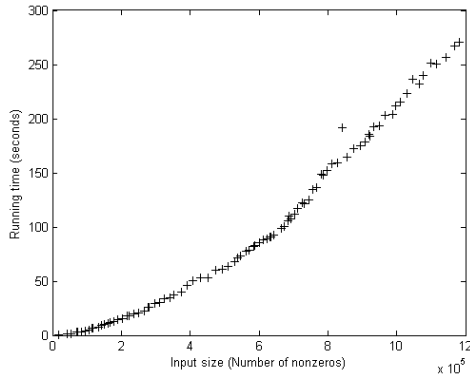
Power Method

1. Let $v \in \mathbb{R}^n$ be a random vector orthogonal to $\pi^T D^{-1}$.
2. Repeat
 - (a) Normalize v , i.e. set $v = v/|v|$.
 - (b) Set $v = Qv$.

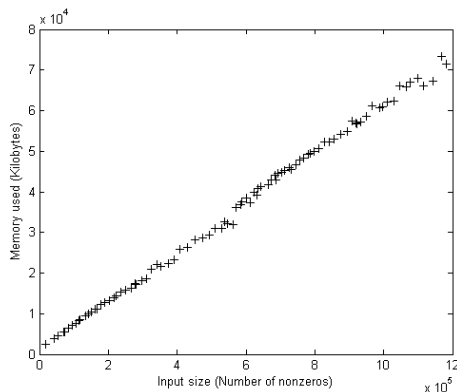
Step 1 ensures that the vector we compute is the second largest eigenvector. Note that $\pi^T D^{-1}Q = \pi^T D^{-1}$ so $\pi^T D^{-1}$ is a left eigenvector with eigenvalue 1. To evaluate $Qv = v$ in Step 3, we only need to do four sparse matrix-vector multiplications, since $Q = (DR^{-1}AA^TD^{-1})$, and each of these matrices is sparse. Note that we do not form Q explicitly. The following lemma shows that the power method takes $\Theta(\log n)$ iterations to convert to the top eigenvector. Although stated for the top eigenvector, the lemma and theorem still hold when the starting vector is chosen uniformly over vectors orthogonal to the top eigenvector $\pi^T D^{-1}$; in this case, the power method will converge to the second largest eigenvector. The proof is given in the Appendix.

LEMMA 1. *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix, and let $v \in \mathbb{R}^n$ be chosen uniformly at random from the unit n -dimensional sphere. Then for any positive integer k , the following holds with probability at least $1 - \delta$:*

$$\frac{\|A^{k+1}v\|}{\|A^k v\|} \geq \left(n \ln \frac{1}{\delta} \right)^{-\frac{1}{2k}} \|A\|_2.$$



(a) Time vs. input size



(b) Space vs. input size

Figure 2: Performance of spectral algorithm in experiments

The next theorem follows directly from the lemma and quantifies the number of steps to run the power method to find a good approximation.

THEOREM 1. *If $k \geq \frac{1}{2\epsilon} \ln(n \ln(\frac{1}{\delta}))$, then with probability at least $1 - \delta$, we have:*

$$\frac{\|A^{k+1}v\|}{\|A^k v\|} \geq (1 - \epsilon)\lambda_1.$$

2.1.2 Time and space requirements

In practice, our algorithm seems to be quite efficient. Figures 2(a) and 2(b) show the results of a performance experiment. In this experiment, we computed a complete hierarchical clustering for N newsgroup articles, where N ranged from 200 to 18,000, in the 20 newsgroups data set [20] and measured the running time and memory used. When we clustered 18,000 documents (for a total of 1.2 million nonzeros in the document-term matrix), we were able to compute a complete hierarchical clustering in 4.5 minutes on commodity hardware (a 3.2 Ghz Pentium IV with 1 gigabyte of RAM).

2.2 Merge phase

The merge phase finds the optimal clustering in the tree T produced by the divide phase. In this section, we give dynamic programs to compute the optimal clustering in the tree T for many standard objective functions. The running time of the merge phase depends on both the number of times we compute the objective function and its evaluation time. Suppose at each interior node we compute a Pareto curve at t points from the Pareto curves of the node’s children. Let c be the cost of evaluating the objective function. Then the total running time is $O(nt^2 + ntc)$, linear in n and c with a small polynomial dependence on t .

k -means: The k -means objective function seeks to find a k -clustering such that the sum of the squared distances of the points in each cluster to the centroid p_i of the cluster is minimized:

$$g(\{C_1, \dots, C_k\}) = \sum_i \sum_{u \in C_i} d(u, p_i)^2.$$

The centroid of a cluster is just the average of the points in the cluster. This problem is NP-hard; several heuristics (such as the k -means *algorithm*) and approximation algorithms exist (e.g. [14, 19]).

Let $\text{OPT}(C, i)$ be the optimal clustering for C using i clusters. Let C_l and C_r be the left and right children of C in T . Then we have the following recurrence: when $i = 1$,

$$\text{OPT}(C, i) = C$$

since we are constrained to only use 1 cluster. When $i > 1$, we have:

$$\text{OPT}(C, i) = \text{argmin}_{1 \leq j < i} g(\text{OPT}(C_l, j) \cup \text{OPT}(C_r, i - j)).$$

By computing the optimal clustering for the leaf nodes first, we can determine the optimal clustering efficiently for any interior node. Then $\text{OPT}(\text{root}, k)$ gives the optimal clustering. Note that in the process of finding the optimal clustering, the dynamic program finds the Pareto curve $\text{OPT}(\text{root}, \cdot)$ that describes the tradeoff between the number of clusters used and the “error” incurred.

Min-diameter: We wish to find a k -clustering for which the cluster with maximum diameter is minimized:

$$g(\{C_1, \dots, C_k\}) = \max_i \text{diam}(C_i).$$

The diameter of any cluster is the maximum distance between any pair of objects in the cluster. A similar dynamic program to that above can find the optimal tree-respecting clustering. This objective function has been studied in [8].

Min-sum: Another objective that has been considered is minimizing the sum of pairwise distances within each cluster:

$$g(\{C_1, \dots, C_k\}) = \sum_{i=1}^k \sum_{u, v \in C_i} d(u, v).$$

We can compute an optimal answer in the tree T by a similar dynamic program to the one above. Although approximation algorithms are known for this problem (as well as the one above), their running times seem too large to be useful in practice [10].

data set	Spectral	p-QR	p-Kmeans	K-means
alt.atheism/comp.graphics	93.6 ± 2.6	89.3 ± 7.5	89.6 ± 6.9	76.3 ± 13.1
comp.graphics/comp.os.ms-windows.misc	81.9 ± 6.3	62.4 ± 8.4	63.8 ± 8.7	61.6 ± 8.0
rec.autos/rec.motorcycles	80.3 ± 8.4	75.9 ± 8.9	77.6 ± 9.0	65.7 ± 9.3
rec.sport.baseball/rec.sport.hockey	70.1 ± 8.9	73.3 ± 9.1	74.9 ± 8.9	62.0 ± 8.6
alt.atheism/sci.space	94.3 ± 4.6	73.7 ± 9.1	74.9 ± 8.9	62.0 ± 8.6
talk.politics.mideast/talk.politics.misc	69.3 ± 11.8	63.9 ± 6.1	64.0 ± 7.2	64.9 ± 8.5

Table 2: 20 newsgroups data set (Accuracy)

Correlation clustering: Suppose we are given a graph where each pair of vertices is either deemed similar (red) or dissimilar (blue). Let R and B be the sets of red and blue edges, respectively. Correlation clustering seeks to find a partition that minimizes the number of blue edges within clusters plus the number of red edges between clusters:

$$g(\{C_1 \dots C_k\}) = \sum_i |\{(u, v) \in B \cap C_i\}| + \frac{1}{2} |\{(u, v) \in R : u \in C_i, v \in U \setminus C_i\}|.$$

Let C be a cluster in the tree T , and let C_l and C_r be its two children. The dynamic programming recurrence for $\text{OPT}(C)$ is:

$$\text{OPT}(C) = \text{argmin} \{g(C), g(\text{OPT}(C_l) \cup \text{OPT}(C_r))\}.$$

If, instead, we are given pairwise similarities in $[0, 1]$, where 0 means dissimilar and 1 means similar, we can define two thresholds t_1 and t_2 . Edges with similarity greater than t_1 are colored red and edges with similarity less than t_2 are colored blue. The same objective function can be applied to these new sets of edges $R_{(t_1)}$ and $B_{(t_2)}$. Approximation algorithms are known for this problem, although the techniques used (linear and semidefinite programming) incur large computational overhead [3, 7, 11, 13, 29].

3. APPLICATION TO WEB SEARCHING: EIGENCLUSTER

In a standard web search engine such as Google or Yahoo, the results for a given query are ranked in a linear order. Although suitable for some queries, the linear order fails to show the inherent structure of the results for queries with multiple meanings or contexts. For instance, consider the query **mickey**. The query can refer to multiple people (Mickey Rooney and Mickey Mantle) or even a fictional character (Mickey Mouse).

We have implemented our methodology in a meta-search engine that discovers the clustered structure for queries and identifies each cluster by its three most significant terms. The website can be found at <http://eigencluster.csail.mit.edu>. The user inputs a query which is then used to find 400 results from Google, a standard search engine. Each result contains the title of the webpage, its location, and a small snippet. We construct a document-term matrix representation of the results; each result is a document and the words in its title and snippet make up its terms. Standard text pre-processing such as TF/IDF and removal of too frequent/infrequent terms is applied. The similarity between two results is the inner product between their two term vectors.

The divide phase was implemented using our spectral algorithm. For the merge phase, we used the correlation clustering objective function with a threshold. A number of other natural objective functions seem to do comparably well. For instance, we have seen similar performance for minimizing the following objective function (for appropriate choice of α, β):

$$\sum_i \alpha \left(\sum_{u, v \in C_i} 1 - A_{(u)} \cdot A_{(v)} \right) + \beta \left(\sum_{u \in C_i, v \notin C_i} A_{(u)} \cdot A_{(v)} \right).$$

One advantage of using these objective functions is that they do not depend on a predefined number of clusters k . This is appropriate for our application, since the number of meanings or contexts of a query is not known beforehand.

Sample queries can be seen in Figure 3; in each example, EigenCluster identifies the multiple meanings of the query as well as keywords corresponding to those meanings. Furthermore, many results are correctly labeled as singletons. In Figure 3, the pictures on the left are screenshots of EigenCluster. The pictures on the right are before and after depictions of the similarity matrix. In the before picture, the results are arranged in the order received from Google. In the after picture, the results are arranged according to the cuts made by the spectral algorithm. Here, the cluster structure is apparent. EigenCluster takes roughly .7 seconds to fetch and cluster results on a Pentium III 700 megahertz with 512 megabytes of RAM.

4. EXPERIMENTS ON TEXT-BASED DATA

The appropriate objective function for an application will naturally depend on the specific application. To show the applicability of our methodology, we show experimental evidence that a *good* clustering exists in the hierarchical clustering constructed by the spectral algorithm. Finding this clustering in the merge phase amounts to determining the right objective function to use. We used our spectral algorithm to create a hierarchical clustering for different data sets of text-based data. In each of the data sets, there was a pre-defined correct classification. We found the partition in the hierarchy that “agrees” the most with the correct classification. The amount of agreement was evaluated using three standard measures: F -measure, entropy, and accuracy. Descriptions of the measures can be found in the Appendix.

We performed experiments on the Reuters, SMART and 20 newsgroups data sets as well as data sets that were used in experiments for other clustering algorithms [6]. We compare the performance of the spectral algorithm in these experiments with known results of other algorithms on the data sets. In all of the experiments, we perform better or compet-

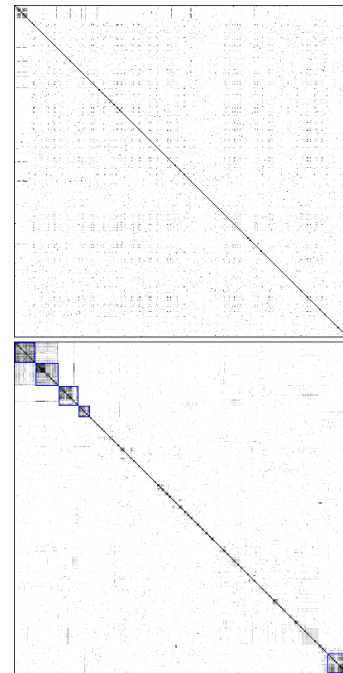
EigenCluster

5 clusters and 294 additional results found in 1.290 seconds.
Explore a cluster or click on a **keyword** to refine your search.

<p>coffee roast senseo (28 pages)</p>	<p>Make Your Own Coffee Pods (INeedCoffee.com) [http://www.ineedcoffee.com/04/...] Make Your Own Coffee Pods. sure everything fits together. The better they fit, Coffee pods are about to invade the US [...] [http://www.springwise.com/...] Coffee pods are about to invade the US. So who's next? New niche companies</p>
<p>espresso coffee illy (25 pages)</p>	<p>Easier Espresso with Pods [http://coffeetea.about.com/library/...] Espresso Pods. Pods are the newest thing in espresso making. Check them out! What Are Pods? [http://coffeetea.about.com/cs/...] Espresso Pod. Email to a friendPrint this page. Related Resources. Espresso</p>
<p>seeds magnolia poppy (24 pages)</p>	<p>Magnolia Seed Pods [http://home.att.net/~SpanishMoss/...] Magnolia Seed Pods are Magnificent! Magnolia Seed Pods can be used to decorate How to Make Christmas Ornaments From Magnolia... [http://www.ehow.com/...] They drop their handsome seed pods just in time to make unusual Christmas</p>
<p>sigmod conference acm (24 pages)</p>	<p>PODS [http://www.informatik.uni-trier.de/...] Symposium on Principles of Database Systems (PODS). ACM Digital Library: PODS 21. PODS 2002: Madison, Wisconsin USA [http://www.informatik.uni-trier.de/...] 21. PODS 2002: Madison, Wisconsin, USA. Lucian Popa (Ed.): Proceedings of the</p>
<p>pea recipe stir (14 pages)</p>	<p>Cook's Thesaurus: Edible Pods [http://www.foodsubs.com/Pods.html] home legumes nuts edible pods. Edible Pods. Chinese pea pod. Chinese pea. Pods - A poem by Carl Sandburg - American... [http://www.americanpoems.com/poets/...] Carl Sandburg - Pods. PEA pods cling to stems. Neponset, the village, Clings to</p>

[PODS](#) [http://www.podsusa.com/]
PODS portable moving and storage, on-site storage containers, mini-storage and

(a) Query: pods



(b) Before/after: pods

EigenCluster

4 clusters and 233 additional results found in 2.120 seconds.
Explore a cluster or click on a **keyword** to refine your search.

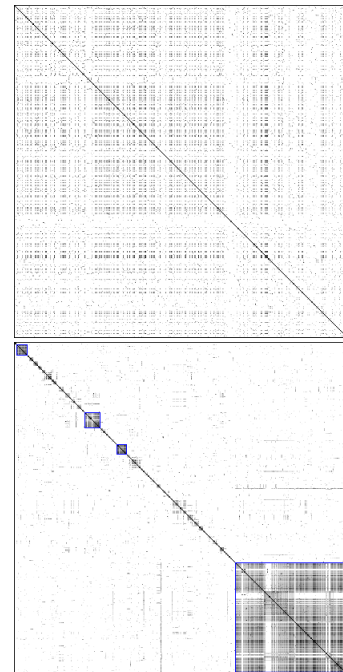
<p>mouse disney walt (138 pages)</p>	<p>The Main Mouse Is In The House [http://www.mickey-mouse.com/...] MICKEY MOUSE, Walt Disney's most famous character, made his screen debut on Hidden Mickeys of Disney [http://www.hiddenmickeys.org/] Hidden Mickeys of Disney is your guide for what's new, Hidden Mickey sightings,</p>
<p>mantle foundation division (20 pages)</p>	<p>IMDb name search [http://www.imdb.com/...] Search Web. Mickey Rooney Characters Plots Biographies Quotes more Tribute to Mickey Mantle [http://www.theswearingens.com/mick/] With this web page, I hope to allow today's youth the opportunity to sneak a</p>
<p>hart discography dead (14 pages)</p>	<p>welcome [http://www.mickeyhart.net/] www.mickeyhart.net/ - 3k - Cached - Similarpages MICKEY HART NET Drummerworld: Mickey Hart [http://www.drummerworld.com/...] Mickey Hart Mickey Hart is best known for his nearly three decades as</p>
<p>rooney star show (13 pages)</p>	<p>The Official Web Site of Mickey Rooney [http://www.mickeyrooney.com/] walk! Now you can bring home any star you wish, including Mickeys, thanks The Mickey Rooney Experience [http://pages.prodigy.net/mshimkus/...] Mickey Rooney Resources: Complete Filmography Mickey Rooney biography from</p>

[CNN Kicks Out the Jams! - Plus--Why the left..](#) [http://www.kausfiles.com/]
CNN Kicks Out the Jams! Plus--Why the left could love Bush's ownership society.

[kaus files dot com](#) [http://www.kausfiles.com/...]
Join the kausfiles.com mailing list! Enter your email address below, then click

[Mickey Rourke](#) [http://www.imdb.com/name/nm000620/]

(c) Query: mickey



(d) Before/after: mickey

Figure 3: EigenCluster search examples

itively with known results. The rest of this section describes the data sets and results.

4.0.1 20 newsgroups

The 20 newsgroups resource [20] is a corpus of roughly 20,000 articles that come from 20 specific Usenet newsgroups. We performed a subset of the experiments in [34]. Each experiment involved choosing 50 random newsgroup articles each from two newsgroups.¹ The results can be seen in Table 2. Note that we perform better than p-QR, the algorithm proposed in [34] on all but one of the experiments. We also outperform k -means and a variation of the k -means algorithm, p-Kmeans. In each of these experiments, the measure of performance was accuracy. Since the experiment involved choosing 50 random newsgroup articles, the experiment was run 100 times and the mean and standard deviation of the results were recorded.

4.0.2 Reuters

The Reuters data set [22] is a corpus of 8,654 news articles that have been classified into 135 distinct news topics. We performed same two experiments on this data set as were conducted in [5, 21, 23]. The first experiment, performed by [5, 21], constructed a complete hierarchical tree for a document-term matrix that includes all 8,654 news articles. In the second experiment, a complete hierarchical tree was produced for a document-term matrix containing only 6,575 news articles from 10 of the 135 largest news topics. This experiment was conducted by [23]. Our algorithm outperformed the results of prior experiments under the F -measure (see Table 3).

data set	Spectral	BEX02	LA99	NJM01
8,654 articles	.713	.57	.63	N/A
6,575 articles	.733	N/A	N/A	.665

Table 3: Reuters data set (F-measure)

4.0.3 Web pages

Boley [6] performs a series of experiments on clustering 185 webpages that fall into 10 distinct categories. In each of the 11 experiments (J1-J11), the term vector for each webpage was constructed in a slightly different way (the exact details can be found in [6]). A comparison of results under the entropy measure can be found in Table 4(a). In 7 of the 11 experiments, our algorithm performs better.

4.0.4 SMART data set

The SMART data set is a set of abstracts originating from Cornell University [25] that have been used extensively in information retrieval experiments. The makeup of the abstracts is as follows: 1,033 medical abstracts (Medline), 1,400 aeronautical systems abstracts (Cranfield), and 1,460 information retrieval abstracts (Cisi). We performed the same four experiments as those found in [12]. In the first three experiments, the data sets were the mixture of abstracts from two classes. In the fourth experiment, the data

¹We used the BOW toolkit for processing the newsgroup data. More information on the BOW toolkit can be found on <http://www-2.cs.cmu.edu/~mccallum/bow>.

set was the set of all abstracts. We perform competitively in the entropy measure (see Table 4(b)).

data set	Spectral	B97
J1	.77	.69
J2	.81	1.12
J3	.54	.85
J4	1.12	1.10
J5	.81	.74
J6	.81	.83
J7	.63	.90
J8	.84	.96
J9	.65	1.07
J10	1.77	1.17
J11	.90	1.05

(a) Webpage data set (Entropy)

data set	Spectral	Dhillon 2001
MedCran	.032	.026
MedCisi	.092	.152
CisiCran	.045	.046
Classic3	.090	.089

(b) SMART data set (Entropy)

Table 4: SMART and Webpage data sets

5. CONCLUSION

We have presented a divide-and-merge methodology for clustering, and shown an efficient and effective spectral algorithm for the divide phase. For the merge phase, we have described dynamic programming formulations that compute the optimal tree-respecting clustering for standard objective functions. Some questions for future work include: are there algorithms for the divide phase such that the tree-respecting clusterings found in the end are provably good approximations to the optimal clusterings? Does the tree produced by the spectral algorithm contain a provably-good clustering for some standard objective functions? Questions also arise from our experimental work. In Section 4, the experiments suggest that for text data, a good clustering exists in the tree constructed by the spectral algorithm. Is there a general objective function that can be used to get the right clustering for this data in the merge phase? Formulating a dynamic program for this objective function would guarantee that the merge phase finds the desired clustering.

6. REFERENCES

- [1] Eigencluster. <http://eigencluster.csail.mit.edu>.
- [2] M. Anderberg. *Cluster Analysis for Applications*. Academic Press, 1973.
- [3] N. Bansal, A. Blum, and S. Chawla. Correlation clustering. In *Proceedings of the 43rd IEEE Symposium on Foundations of Computer Science*, pages 238–247, 2002.

- [4] D. Barbara, Y. Li, and J. Couto. Coolcat: an entropy-based algorithm for categorical clustering. In *Proceedings of the eleventh international conference on Information and knowledge management*, pages 582–589, 2002.
- [5] F. Beil, M. Ester, and X. Xu. Frequent term-based text clustering. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 436–442, 2002.
- [6] D. Boley. Principal direction divisive partitioning. *Data Mining and Knowledge Discovery*, 2(4):325–344, 1998.
- [7] M. Charikar, V. Guruswami, and A. Wirth. Clustering with qualitative information. In *Proceedings of the 44th Annual IEEE Symposium on Foundations of Computer Science*, pages 524–533, 2003.
- [8] C. Chekuri, T. Feder, and R. Motwani. Incremental clustering and dynamic information retrieval. In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing*, pages 626–635, 1997.
- [9] D. R. Cutting, D. R. Karger, J. O. Pedersen, and J. W. Tukey. Scatter/gather: a cluster-based approach to browsing large document collections. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 318–329, 1992.
- [10] W.F. de la Vega, M. Karpinski, C. Kenyon, and Y. Rabani. Approximation schemes for clustering problems. In *Proceedings of the 35th Annual ACM Symposium on Theory of Computing*, pages 50–58, 2003.
- [11] E.D. Demaine and N. Immerlica. Correlation clustering with partial information. In *Proceedings of the 6th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems*, pages 1–13, 2003.
- [12] I. S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 269–274, 2001.
- [13] D. Emanuel and A. Fiat. Correlation clustering—minimizing disagreements on arbitrary weighted graphs. In *Proceedings of the 11th European Symposium on Algorithms*, pages 208–220, 2003.
- [14] J.A. Hartigan and M.A. Wong. A k-means clustering algorithm. In *Applied Statistics*, pages 100–108, 1979.
- [15] T. Hofmann. The cluster-abstraction model: Unsupervised learning of topic hierarchies from text data. In *International Joint Conference on Artificial Intelligence*, pages 682–687, 1999.
- [16] A. Jain and R. Dubes. *Algorithms for Clustering Data*. Prentice Hall, 1988.
- [17] A. Jain, M. Murty, and P. Flynn. Data clustering: A review. *ACM Computing Surveys*, 31, 1999.
- [18] R. Kannan, S. Vempala, and A. Vetta. On clusterings: Good, bad, and spectral. *Journal of the ACM*, 51(3):497–515, 2004.
- [19] A. Kumar, S. Sen, and Y. Sabharwal. A simple linear time $(1+\epsilon)$ -approximation algorithm for k -means clustering in any dimensions. In *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science*, pages 454–462, 2004.
- [20] K. Lang. 20 newsgroups data set. <http://www.ai.mit.edu/people/jrennie/20Newsgroups/>.
- [21] B. Larsen and C. Aone. Fast and effective text mining using linear-time document clustering. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 16–22, 1999.
- [22] D. Lewis. Reuters data set. <http://www.daviddlewis.com/resources/testcollections/reuters21578/>.
- [23] A. Nickerson, N. Japkowicz, and E. Milios. Using unsupervised learning to guide re-sampling in imbalanced data sets. In *Proceedings of the Eighth International Workshop on AI and Statistics*, pages 261–265, 2001.
- [24] S. Sahní and T. Gonzalez. P-complete approximation problems. *Journal of the ACM*, 23(3):555–566, 1976.
- [25] G. Salton. SMART Data Set. <ftp://ftp.cs.cornell.edu/pub/smart>.
- [26] C. E. Shannon. A mathematical theory of communication. *Bell Systems Technical Journal*, 27:379–423, 1948.
- [27] N. Slonim and N. Tishby. Document clustering using word clusters via the information bottleneck method. In *Proceedings of the 23rd Annual International ACM Conference on Research and Development in Information Retrieval*, pages 208–215, 2000.
- [28] M. Steinbach, G. Karypis, and V. Kumar. A comparison of document clustering techniques. In *KDD Workshop on Text Mining*, 2000.
- [29] C. Swamy. Correlation clustering: Maximizing agreements via semidefinite programming. In *Proceedings of ACM-SIAM Symposium on Discrete Algorithms*, pages 519–520, 2004.
- [30] J. Theiler and G. Gisler. A contiguity-enhanced k-means clustering algorithm for unsupervised multispectral image segmentation. In *Proceedings of the Society of Optical Engineering*, pages 108–111, 1997.
- [31] C. J. Van Rijsbergen. *Information Retrieval, 2nd edition*. Dept. of Computer Science, University of Glasgow, 1979.
- [32] W. Wong and A. Fu. Incremental document clustering for web page classification. In *IEEE International Conference on Information Society in the 21st Century: Emerging Technologies and New Challenges*, 2000.
- [33] O. Zamir, O. Etzioni, O. Madani, and R. M. Karp. Fast and intuitive clustering of web documents. In *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining*, pages 287–290, 1997.
- [34] H. Zha, C. Ding, M. Gu, X. He, and H. Simon. Spectral relaxation for k-means clustering. In *Neural Information Processing Systems*, pages 1057–1064, 2001.
- [35] Y. Zhao and G. Karypis. Evaluation of hierarchical clustering algorithms for document datasets. In *Proceedings of the Eleventh International Conference on Information and Knowledge Management*, pages 515–524, 2002.

7. APPENDIX

7.1 F -measure, Entropy, and Accuracy

For a data set, let the correct classification be $C_1 \dots C_k$. We refer to each C_i as a *class*. Let the nodes of a hierarchical clustering be $\hat{C}_1 \dots \hat{C}_l$. We refer to each \hat{C}_i as a *cluster* – the subset of nodes in the tree below it.

F -measure: For each class C_i , the F -measure of that class is:

$$F(i) = \max_{j=1}^l \frac{2P_j R_j}{P_j + R_j}$$

where:

$$P_j = \frac{|C_i \cap \hat{C}_j|}{|\hat{C}_j|}, R_j = \frac{|C_i \cap \hat{C}_j|}{|C_i|}$$

The F -measure of the clustering is defined as:

$$\sum_{i=1}^k F(i) \cdot \frac{|C_i|}{|C|}$$

The F -measure score is in the range $[0, 1]$ and a **higher** F -measure score implies a better clustering. For a more in-depth introduction and justification to the F -measure, see e.g. [31, 21, 5, 23].

Entropy: For each cluster \hat{C}_j , we define the entropy of \hat{C}_j as:

$$E(\hat{C}_j) = \sum_{i=1}^k - \left(\frac{|C_i \cap \hat{C}_j|}{|\hat{C}_j|} \right) \log \left(\frac{|C_i \cap \hat{C}_j|}{|\hat{C}_j|} \right)$$

The entropy of a cluster is a measure of the disorder within the cluster. As such, a **lower** entropy score implies that a clustering is better; the best possible entropy score is 0. Entropy was first introduced in [26] and has been used as a measure of clustering quality in [6, 12, 4].

The entropy of a k -clustering $\hat{C}_1 \dots \hat{C}_k$ is the weighted sum of the entropies of the clusters. The entropy of a hierarchical clustering $\{\hat{C}_1 \dots \hat{C}_l\}$ is the minimum entropy of any choice of k nodes that partition C .

Accuracy: The accuracy of a cluster \hat{C}_j is:

$$A(\hat{C}_j) = \max_{i=1}^k \frac{|C_i \cap \hat{C}_j|}{|\hat{C}_j|}$$

As before, the accuracy of a k -clustering $C_1 \dots C_k$ is the weighted sum of accuracies. The accuracy of a hierarchical clustering is the maximum accuracy of any choice of k nodes that partition C . Note that the range of an accuracy score is between 0 and 1; the **higher** the accuracy score, the better.

Accuracy, which has been used as a measure of performance in supervised learning, has also been used in clustering (see [27]).

7.2 Convergence Proof

PROOF (OF LEMMA 1). Since A is symmetric, we can write

$$A = \sum_{i=1}^n \lambda_i u_i u_i^T,$$

where the λ_i 's are the eigenvalues of A arranged in the order $|\lambda_1| \geq |\lambda_2| \dots |\lambda_n|$ and the u_i are the corresponding eigenvectors. Express v in this basis as $v = \sum_i \alpha_i u_i$, where $\sum_i \alpha_i^2 = 1$. Since, v is random, we have that with probability at least $1 - \delta$, $\alpha_1^2 \geq 1/(n \ln(1/\delta))$. Then, using Hölder's inequality (which says that for any $p, q > 0$ satisfying $(1/p) + (1/q) = 1$ and any $a, b \in \mathbb{R}^n$, we have $\sum_i a_i b_i \leq (\sum_i |a_i|^p)^{1/p} (\sum_i |b_i|^q)^{1/q}$), we have

$$\|A^k v\|^2 = \sum_i \alpha_i^2 \lambda_i^{2k} \leq \left(\sum_i \alpha_i^2 \lambda_i^{2k+2} \right)^{k/(k+1)}$$

where the last inequality holds using Hölder with $p = 1 + (1/k)$ $q = k + 1$ $a_i = \alpha_i^{2k/(k+1)} \lambda_i^{2k}$ $b_i = \alpha_i^{2/(k+1)}$. Note that:

$$\left(\sum_i \alpha_i^2 \lambda_i^{2k+2} \right)^{k/(k+1)} \leq \left(\sum_i \alpha_i^2 \lambda_i^{2k+2} \right) / \lambda_1^2 \alpha_1^{2/(k+1)}$$

from which the lemma follows. \square

8. EIGENCLUSTER EXAMPLE SEARCHES

We give a few more EigenCluster example searches on the next page.

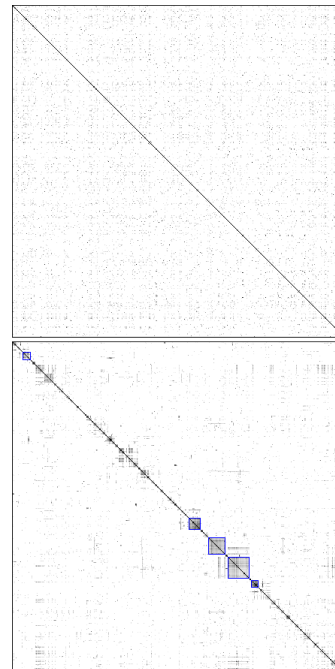
EigenCluster

5 clusters and 353 additional results found in 1.930 seconds.
Explore a cluster or click on a **keyword** to refine your search.

<p>plant database preserving (29 pages)</p> <p>shrubs plant information (23 pages)</p> <p>christmas national associated (16 pages)</p> <p>binary problems recipes (11 pages)</p> <p>decision overview tools (10 pages)</p>	<p>Trees - The National Arbor Day Foundation [http://www.arborday.org/trees/] Planting and caring for trees, identifying trees, buying trees, conferences</p> <p>American Forests [http://www.americanforests.org/] Plant Trees Now Planting trees in our Global ReLeaf Projects</p> <hr/> <p>Home Page: FOR 141/241: Tree and Shrub ID /... [http://oregonstate.edu/instruct/...] They each include information on ecology, principal uses, and natural history</p> <p>Trees and Shrubs [http://gardenline.usask.ca/trees/] GardenLine. Alkaline Tolerant Trees and Shrubs; Birch Tree Guide; Birch Dieback</p> <hr/> <p>National Christmas Tree Association [http://www.realchristmastrees.org/] The National Christmas Tree Association represents the growers of trees and</p> <p>Massachusetts Christmas Tree Association [http://christmas-trees.org/] If you are a Real Christmas Tree consumer you will find valuable information in</p> <hr/> <p>Ivan Galkin [http://ulcar.uml.edu/~iag/CS/...] Enumeration of the Binary Trees (Catalan numbers). For subtrees. The total</p> <p>Red/Black Tree Demo [http://www.ecs.uic.edu/~franco/...] Red/Black Trees: These are binary trees with the following properties. Due to</p> <hr/> <p>Mind Tools - Decision tree analysis [http://www.mindtools.com/pages/...] Outcomes. How to use tool: Decision Trees are excellent tools for helping you</p> <p>TREES Project Overview [http://www.treepeople.org/trees/] TREES logo. Overview. Welcome to the TREES Project. In these pages you will</p>
---	--

[British Trees Website Home Page - native...](http://www.british-trees.com/) [<http://www.british-trees.com/>]
Welcome to the British Trees Website! This site contains a wealth of reference

(a) Query: trees



(b) Before/after: trees

EigenCluster

5 clusters and 245 additional results found in 1.970 seconds.
Explore a cluster or click on a **keyword** to refine your search.

<p>teddy care artists (69 pages)</p> <p>polar arctic 2004 (36 pages)</p> <p>black ursus americanus (32 pages)</p> <p>chicago tickets team (26 pages)</p> <p>athletics university official (15 pages)</p>	<p>Teddy Bears on the NET ---- Teddy Bears and... [http://www.tbonnet.com/] Web Site designed and maintained by Terry Bauman of Teddy Bears on the NET</p> <p>Care-Bears.com : Official Website of the Care... [http://www.care-bears.com/] Care Bears, Find Hundreds of Care Bears items at great prices! Click Here! The</p> <hr/> <p>Polar Bears International, working to... [http://www.polarbearsalive.org/] Tons of educational information and research on polar bears, gorgeous polar bear</p> <p>The Bear Facts on Polar Bears [http://www.polarbearsalive.org/...] Compiled by the web's foremost source on polar bear information, Polar Bears</p> <hr/> <p>Bears [http://edtech.kennesaw.edu/web/...] Bears. Research/Informational Sites. All Asian Black Bear; Bear Den -</p> <p>Black Bear Home Page [http://www.bear.org/Black/...] Black bear facts, sounds, signs, populations, and pictures.</p> <hr/> <p>Chicago Bears [http://www.chicagobears.com/] Hutchinson passes inspired Bears to victory The carefree smile that creased</p> <p>Chicago Sun-Times - Sports - Bears [http://www.suntimes.com/index/...] Bears 24, Vikings 14 Chad Hutchinson, who was surfing in California a little</p> <hr/> <p>California Golden Bears - Official Athletic... [http://calbears.ccsn.com/] (more). 12/5/04 - W Swimming Bears Finish Second at Georgia Invitational The</p> <p>California Golden Bears - Official Athletic... [http://www.calbears.com/] 4 California Visits Southern Miss ESPN to Televis Cal-USM Game at 4:30 pm 16</p>
---	---

[Bears.Org, Bear Information and Resources](http://www.bears.org/) [<http://www.bears.org/>]
Dedicated to preserving information about bears and presenting it in a

(c) Query: bears



(d) Before/after: bears

Figure 4: EigenCluster search examples