

Dartmouth College

Dartmouth Digital Commons

Computer Science Technical Reports

Computer Science

10-1-1997

Computing Dense Clusters On-line for Information Organization

Javed Aslam

Dartmouth College

Katya Pelekhov

Dartmouth College

Daniela Rus

Dartmouth College

Follow this and additional works at: https://digitalcommons.dartmouth.edu/cs_tr



Part of the [Computer Sciences Commons](#)

Dartmouth Digital Commons Citation

Aslam, Javed; Pelekhov, Katya; and Rus, Daniela, "Computing Dense Clusters On-line for Information Organization" (1997). Computer Science Technical Report PCS-TR97-324.

https://digitalcommons.dartmouth.edu/cs_tr/157

This Technical Report is brought to you for free and open access by the Computer Science at Dartmouth Digital Commons. It has been accepted for inclusion in Computer Science Technical Reports by an authorized administrator of Dartmouth Digital Commons. For more information, please contact dartmouthdigitalcommons@groups.dartmouth.edu.

Computing Dense Clusters On-line for Information Organization

Javed Aslam Katya Pelekhov Daniela Rus

Department of Computer Science
Dartmouth College
Hanover, NH 03755

phone: (603) 646 1691

fax: (603) 646 1672

{jaa,katya,rus}@cs.dartmouth.edu

Abstract

We present and analyze the off-line star algorithm for clustering static information systems and the on-line star algorithm for clustering dynamic information systems. These algorithms partition a document collection into a number of clusters that is naturally induced by the collection. We show a lower bound on the accuracy of the clusters produced by these algorithms. We use the random graph model to show that both star algorithms produce correct clusters in time $\Theta(V + E)$. Finally, we provide data from extensive experiments.

1 Introduction

Modern information systems have vast amounts of unorganized data that changes dynamically. Consider, for example, the flow of information that arrives continuously on news wires, or is aggregated by a news organization such as CNN. Some stories are brand new. Other stories are follow-ups of previous stories. Yet another type of stories make previous reportings obsolete. The news focus changes regularly with this flow of information. In such dynamic systems users need to locate information fast. Users often don't know what they need until they need it. In dynamic, time-pressured situations such as emergency relief for weather disasters, access to the latest information needs to happen fast. Current information systems such as Inquiry[Tur90], Smart[Sal91], or Alta Vista provide some simple automation by computing ranked (sorted) lists of documents, but it is ineffective for users to scan a list of hundreds of document titles. To cull the critical information out of a large set of potentially useful dynamic sources we need methods for organizing and reorganizing dynamic information as accurate clusters, and ways of presenting

users with the topic summaries at various levels of detail.

There has been extensive research on clustering and applications to many domains [HS86, AB84]. For a good overview see [JD88]. For a good overview of using clustering in information retrieval see [Wil88].

The use of clustering in information retrieval was mostly driven by *the cluster hypothesis* [Rij79] which states that relevant documents tend to be more closely related to each other than to non-relevant documents. Efforts have been made to find whether the cluster hypothesis is valid. Voorhees [Voo85] discusses a way of evaluating whether the cluster hypothesis holds and shows negative results. Croft [Cro80] describes a method for bottom-up cluster search that could be shown to outperform a full ranking system for the Cranfield collection. The single link method [Cro77] does not provide any guarantees for the topic similarity within a cluster. In [JR71] Jardine and van Rijsbergen show some evidence that search results could be improved by clustering. Hearst and Pedersen [HP96] re-examine the cluster hypothesis by focusing on the Scatter/Gather system [CKP93] and conclude that it holds for browsing tasks.

Systems like Scatter/Gather [CKP93] provide a mechanism for user-driven organization of data in a fixed number of clusters, but the users need to be in the loop and the computed clusters do not have accuracy guarantees. Scatter/Gather uses fractionation to compute nearest-neighbor clusters. In a recent STOC paper, Charika et al. [CCFM97] consider a dynamic clustering algorithm to partition a collection of text documents into a *fixed number* of clusters. Since in dynamic information systems the the number of topics is not known *a priori*, a fixed number of clusters

can't generate a natural partition for the information.

Our work on clustering presented in this paper and in [APR97] provides positive evidence for the cluster hypothesis. We propose an off-line algorithm for clustering static information and an on-line version of this algorithm for clustering dynamic information. These two algorithms compute clusters induced by the natural topic structure of the space. Thus, this work is different than [CKP93, CCFM97] in that we do not impose the constraint to use a fixed number of clusters. As a result, we can guarantee a lower bound on the topic similarity between the documents in each cluster. The model for topic similarity is the standard vector space model used in the information retrieval community [Sal89] which is explained in more detail in this paper in Section 2.

To compute accurate clusters, we formalize clustering as covering graphs by cliques. Covering by cliques is NP-complete, and thus intractable for large document collections. Recent graph-theoretic results have shown that the problem cannot even be approximated in polynomial time [LY94, Zuc93]. Recent results for covering graphs by dense subgraphs [KP93] are encouraging. We used a cover by dense subgraphs that are star-shaped and that can be computed off-line for static data, or on-line for dynamic data. We show that the off-line and on-line algorithms produce correct clusters efficiently. Asymptotically, the running time of both algorithms is linear in the number of edges in the similarity graph that defines the information space (explained in detail in Section 2). We also show lower bounds on the topic similarity in the clusters. Finally, we provide experimental data.

Our algorithms for organizing information systems can be used in several ways. The off-line algorithm can be used as a pre-processing step in a static information system or as a post-processing step on the specific documents retrieved by a query. As a pre-processor, this system assists users with deciding how to browse a database of free text documents by highlighting relevant topics and irrelevant subtopics. Such clustered data is useful for narrowing down the database over which detailed queries can be formulated. As a post-processor, this system classifies the retrieved data into clusters that capture topic categories and subcategories. The on-line algorithm can be used as a basis for constructing self-organizing information systems. As the content of dynamic information systems changes, the on-line algorithm can efficiently automate the process of organization and re-organization to compute accurate topic summaries at various level of similarity.

2 Clustering static data with star-shaped subgraphs

In this section we motivate and present an off-line algorithm for organizing information systems. The algorithm is very simple and efficient, and computes high-density clusters.

We formulate our problem by representing an information system by its *similarity graph*. A similarity graph is an undirected, weighted graph $G = (V, E, w)$ where vertices in the graph correspond to documents and each weighted edge in the graph corresponds to a measure of similarity between two documents. We measure the similarity between two documents by using the cosine metric in the vector space model of the Smart information retrieval system [Sal91, Sal89].

The vector space model for textual information aggregates statistics on the occurrence of words in documents. The premise of the vector space model is that two documents are similar if they use the same words. A vector space can be created for a collection (or corpus) of documents by associating each important word in the corpus with one dimension in the space. The result is a high dimensional vector space. Documents are mapped as points in this space according to their word frequencies. Similar documents map to nearby points. In the vector space model, document similarity is measured by the angle between the corresponding document vectors. The standard in the information retrieval community is to map the angles to the interval $[0, 1]$ by taking the cosine of the vector angles.

G is a complete graph with edges of varying weight. An organization of the graph that produces reliable clusters of similarity σ (*i.e.*, clusters where documents have pairwise similarities of at least σ) can be obtained by (1) thresholding the graph for σ and then (2) performing a *minimum clique cover* with maximal cliques on the resulting graph G_σ . The *thresholded graph* G_σ is an undirected graph obtained from G by eliminating all the edges whose weight is lower than σ . The minimum clique cover has two features. First, by using cliques to cover the similarity graph, we are guaranteed that all the documents in a cluster have the desired degree of similarity. Second, minimal clique covers with maximal cliques allow vertices to belong to *several* clusters. In our information retrieval application this is a desirable feature as documents can have multiple subthemes.

Unfortunately, this approach is computationally intractable. For real corpora, similarity graphs can be very large. The clique cover problem is NP-complete, and it does not admit polynomial-time approximation algorithms [LY94, Zuc93]. While we cannot perform a

clique cover nor even approximate such a cover, we can instead cover our graph by *dense subgraphs*. What we lose in intra-cluster similarity guarantees, we gain in computational efficiency. In the sections that follow, we describe off-line and on-line covering algorithms and analyze their performance and efficiency.

2.1 Dense Star-Shaped Covers

We approximate a clique cover by covering the associated thresholded similarity graph with *star-shaped subgraphs*. A star-shaped subgraph on $m + 1$ vertices consists of a single *star center* and m *satellite vertices*, where there exist edges between the star center and each of the satellite vertices. While finding cliques in the thresholded similarity graph G_σ guarantees a pairwise similarity between documents of at least σ , it would appear at first glance that finding star-shaped subgraphs in G_σ would provide similarity guarantees between the star center and each of the satellite vertices, but no such similarity guarantees *between satellite vertices*. However, by investigating the geometry of our problem in the vector space model, we can derive a *lower bound* on the similarity between satellite vertices as well as provide a formula for the *expected* similarity between satellite vertices. The latter formula predicts that the pairwise similarity between satellite vertices in a star-shaped subgraph is high, and together with empirical evidence supporting this formula, we shall conclude that covering G_σ with star-shaped subgraphs is a reliable method for clustering a set of documents.

Consider three documents C , S_1 and S_2 which are vertices in a star-shaped subgraph of G_σ , where S_1 and S_2 are satellite vertices and C is the star center. By the definition of a star-shaped subgraph of G_σ , we must have that the similarity between C and S_1 is at least σ and that the similarity between C and S_2 is also at least σ . In the vector space model, these similarities are obtained by taking the cosine of the angle between the vectors associated with each document. Let α_1 be the angle between C and S_1 , and let α_2 be the angle between C and S_2 . We then have that $\cos \alpha_1 \geq \sigma$ and $\cos \alpha_2 \geq \sigma$. Note that the angle between S_1 and S_2 can be at most $\alpha_1 + \alpha_2$, and therefore we have proven the following lower bound on the similarity between satellite vertices in a star-shaped subgraph of G_σ .

Proposition 1 *Let G_σ be a similarity graph and let S_1 and S_2 be two satellites in the same star in G_σ . Then the similarity between S_1 and S_2 must be at least*

$$\cos(\alpha_1 + \alpha_2) = \cos \alpha_1 \cos \alpha_2 - \sin \alpha_1 \sin \alpha_2.$$

If $\sigma = 0.7$, $\cos \alpha_1 = 0.75$ and $\cos \alpha_2 = 0.85$, for instance, we can conclude that the similarity between

the two satellite vertices must be at least¹

$$(0.75) \cdot (0.85) - \sqrt{1 - (0.75)^2} \sqrt{1 - (0.85)^2} \approx 0.29.$$

Note that while this may not seem very encouraging, the above analysis is based on absolute worst-case assumptions, and in practice, the similarities between satellite vertices are much higher. We further undertook a study to determine the *expected* similarity between two satellite vertices.

2.2 The random graph model

The model we use for computing the expected similarity error in a star is the *random graph model* [Bol95]. A random graph $G_{n,p}$ is an undirected graph with n vertices, where each of its possible edges is inserted randomly and independently with probability p . Our problem fits the random graph model if we make the mathematical assumption that “similar” documents are essentially “random perturbations” of one another in the vector space model. This assumption is equivalent to viewing the similarity between two related documents as a random variable. By thresholding the edges of the similarity graph at a fixed value, for each edge of the graph there is a random chance (depending on whether the value of the corresponding random variable is above or below the threshold value) that the edge remains in the graph. This thresholded similarity graph is thus a random graph. We will use the random graph model for the analysis and the experimental verification of all the algorithms presented in this paper.

The following upper bound on the expected similarity between two satellite vertices holds in the random graph model:

Proposition 2 *The expected similarity between two satellite vertices S_1 and S_2 in the same star in a similarity graph G_σ is:*

$$\cos \alpha_1 \cos \alpha_2 + \frac{\sigma}{1 + \sigma} \sin \alpha_1 \sin \alpha_2.$$

Proof: The proof of this proposition relies on the random graph model and is omitted for space considerations. \square

Note that for the previous example, the above formula would predict a similarity between satellite vertices of approximately 0.78. We have tested this formula against real data, and the results of the test with the MEDLINE data set² are shown in Figure 1. In this

¹Note that $\sin \theta = \sqrt{1 - \cos^2 \theta}$.

²MEDLINE is a large collection of medical abstracts that is often used as benchmark in information retrieval experiments.

plot, the x - and y -axes are similarities between cluster centers and satellite vertices, and the z -axis is the actual mean squared prediction error of the above formula for the similarity between satellite vertices. Note that the absolute error (roughly the square root of the mean squared error) is quite small (approximately 0.13 in the worst case), and for reasonably high similarities, the error is negligible. From our tests with real data, we have concluded that the random graph model holds and that this formula is quite accurate. We can conclude that star-shaped subgraphs are reasonably “dense” in the sense that they imply relatively high pairwise similarities between documents.

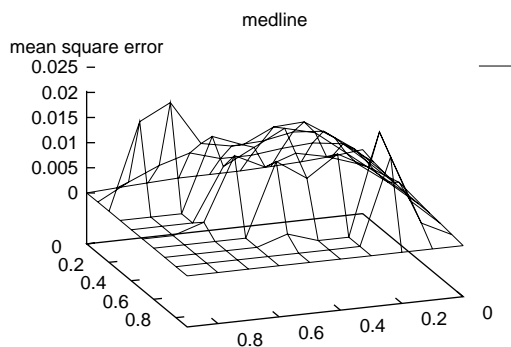


Figure 1: This figure shows the error for a 6000 abstract subset of MEDLINE.

3 The off-line star algorithm

Motivated by the discussion of the previous section, we now present the *star algorithm* which can be used to organize documents in an information system. The star algorithm is based on a greedy cover of the thresholded similarity graph by star-shaped subgraphs; the algorithm itself is summarized in Figure 2 below.

Theorem 1 *The running time of the off-line star algorithm on a similarity graph G_σ is $\Theta(V + E)$.*

Proof: The following implementation of this algorithm has a running time *linear* in the number of edges of the graph. Each vertex v has a data structure associate with it that contains $v.degree$, the degree of the vertex, $v.adj$, the list of adjacent vertices, $v.marked$, which is a bit denoting whether the vertex belongs to a star or not, and $v.center$, which is a bit denoting whether the vertex is a star center. The implementation starts by sorting the vertices in V by degree. The

For any threshold σ :

1. Let $G_\sigma = (V, E_\sigma)$ where $E_\sigma = \{e : w(e) \geq \sigma\}$.
2. Let each vertex in G_σ initially be *unmarked*.
3. Calculate the degree of each vertex $v \in V$.
4. Let the highest degree unmarked vertex be a star center, and construct a cluster from the star center and its associated satellite vertices. Mark each node in the newly constructed cluster.
5. Repeat step 4 until all nodes are marked.
6. Represent each cluster by the document corresponding to its associated star center.

Figure 2: The star algorithm

program then scans the sorted vertices from the highest degree to the lowest as a greedy search for star centers. Only vertices that do not belong to a star already (that is, they are unmarked) can become star centers. Upon selecting a new star center v , its $v.center$ and $v.marked$ bits are set and for all $w \in v.adj$, $w.marked$ is set. Only one scan of V is needed to determine all the star centers. Upon termination, the star centers and only the star centers have the *center* field set. We call the set of star centers the *star cover* of the graph. Each star is fully determined by the star center, as the satellites are contained in the adjacency list of the center vertex. \square

This algorithm has two interesting features. The first feature is that the star cover is not unique. A similarity graph may have several different star covers because when there are several vertices of the same highest degree, the algorithm arbitrarily chooses one of them as a star center (whichever shows up first in the sorted list of vertices). The second feature of this algorithm is that it provides a simple encoding of a star cover by assigning the types “center” and “satellite” (which is the same as “not center” in our implementation) to vertices. We define a *correct star cover* as a star cover that assigns the types “center” and “satellite” in such a way that (1) a star center is not adjacent to any other star center and (2) every satellite vertex is adjacent to at least one center vertex of higher degree. It immediately follows that:

Theorem 2 *The off-line star algorithm produces a correct star cover.*

We will use the two features of the off-line algorithm mentioned above in the analysis of the on-line version of the star algorithm, in the next section.

4 Clustering dynamic data with the star algorithm

In this section we consider algorithms for computing the organization of a dynamic information system. We derive an on-line version of the star algorithm for information organization that can incrementally compute clusters of similar documents. We continue assuming the vector space model and the cosine metric to capture the pairwise similarity between the documents of the corpus, and the random graph model for analyzing the expected behavior of the new algorithm.

We assume that documents are inserted or deleted from the collection one at a time. For simplicity, we will focus our discussion on adding documents to the collection. The delete algorithm is similar. The intuition behind the incremental computation of the star cover of a graph after a new vertex is inserted is depicted in Figure 3. The top figure denotes a similarity graph and a correct star cover for this graph. Suppose a new vertex is inserted in the graph, as in the middle figure. The original star cover is no longer correct for the new graph. The bottom figure shows the correct star cover for the new graph. How does the addition of this new vertex affect the correctness of the star cover? In general, the answer depends on the degree of the new vertex and on its adjacency list. If the adjacency list of the new vertex does not contain any star centers, the new vertex can be added in the star cover as a star center. If the adjacency list of the new vertex contains any center vertex c whose degree is higher, the new vertex becomes a satellite vertex of c . The difficult case that destroys the correctness of the star cover is when the new vertex is adjacent to a collection of star centers, each of whose degree is lower than that of the new vertex. In this situation, the star structure already in place has to be modified to assign the new vertex as a star center. The satellite vertices in the stars that are broken as a result have to be re-evaluated.

4.1 The on-line star algorithm

Motivated by the intuition in the previous section, we now describe an on-line algorithm for incrementally computing star covers of dynamic graphs. The algorithm is shown in Figure 4. This algorithm uses a special data structure to efficiently maintain the star covers of an undirected graph $G = (V, E)$. For each vertex $v \in V$, we maintain the following data.

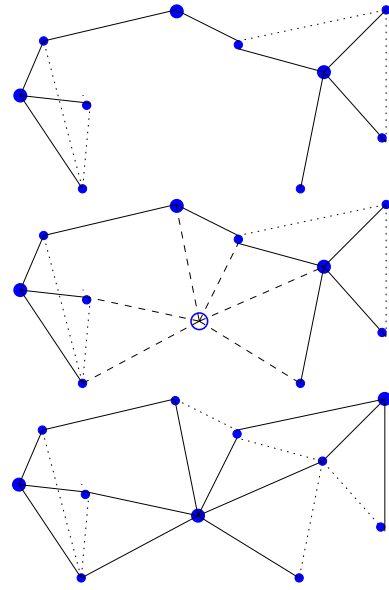


Figure 3: This figure shows the star cover change after the insertion of a new vertex. The larger-radius disks denote star centers, the other disks denote satellite vertices. The star edges are denoted by solid lines. The inter-satellite edges are denoted by dotted lines. The top figure shows an initial graph and its star cover. The middle figure shows the graph after the insertion of a new document. The bottom figure shows the star cover of the new graph.

$v.type$	satellite or center
$v.degree$	degree of v
$v.adj$	list of adjacent vertices
$v.centers$	list of adjacent centers
$v.inQ$	flag specifying if v being processed

Note that while $v.type$ can be inferred from $v.centers$ and $v.degree$ can be inferred from $v.adj$, it will be convenient to have all five pieces of data in the algorithm. Let α be a vertex to be added to G , and let L be the list of vertices in G which are adjacent to α . The algorithm in Figure 4 will appropriately update the star cover of G . The algorithm maintains a priority queue Q of vertices not yet correctly placed in the star cover. When a star is broken, its center and satellites are placed in Q .

The on-line star cover algorithm is significantly more complex than its off-line counterpart. We devote the rest of this section to proving that the algorithm is correct and to analyzing its expected running time.

4.2 Correctness

In this section we show that the on-line algorithm is correct by proving that it produces the same star cover as the off-line algorithm, when the off-line algorithm is run on the final graph considered by the on-line algorithm. Before we state the result, we note that the off-line star algorithm does not produce a unique cover. When there are several vertices of the same highest degree, the algorithm arbitrarily chooses one of them as the next star center. We will show that the cover produced by the on-line star algorithm is the same as one of the covers that can be produced by the off-line algorithm

Theorem 3 *The cover generated by the on-line star algorithm when $G = (V, E)$ is constructed incrementally (by inserting its vertices one at a time) is identical with one of the covers generated by the off-line star algorithm on G .*

Proof: We can view a star cover of G as a correct assignment of types (that is, “center” or “satellite”) to the vertices of G . The off-line star algorithm assigns correct types to the vertices of G . We will prove the correctness of the on-line star cover by induction. The induction invariant is that at all times, the types of vertices not in Q are correct, assuming that the types of vertices in Q are correct. This would imply that when Q is empty, all vertices are assigned a correct type and thus the star cover is correct.

The invariant is true initially: as the type of the new node α is unknown and α is in Q ; the type of all the neighbors of α are unknown and all the neighbors are in Q ; all the other vertices have correct types from the original cover. We now show that the induction invariant is maintained throughout the algorithm. Consider Figure 4. The first thing to note is that the type of all the vertices in Q is “satellite”. Statements 14 and 18 enqueue satellite vertices. Statement 38 enqueues both satellite and center vertices, but the types of all the center vertices enqueued here is set to “satellite” in statement 33. We now argue that every time a vertex ϕ of highest degree is pulled out of Q , it is assigned a correct type. When ϕ has no centers on its adjacency list, its type should be “center” (which is assigned correctly by statement 22). When ϕ is adjacent to star centers δ_i , each of which has a lower degree than ϕ , the correct type for ϕ is “center” (statement 28). This action has a side effect: all δ_i cease to be star centers and thus get enqueued for further evaluation (statements 32-39). Otherwise, the correct type for ϕ is the default “satellite”. Since ϕ was extracted

from Q and all vertices in Q are satellites, the type of ϕ is correct in this case as well.

To complete the argument, what remains to be shown is that eventually the queue Q becomes empty. The termination of the while loop at statement 19 in Figure 4 is guaranteed by the following result.

Lemma 1 *The degree of the stars broken by the on-line star algorithm is strictly monotonically decreasing.*

The lemma is equivalent to the following statement: node ϕ in Q has the potential of becoming a star center and has the capability of adding new nodes γ to Q that can become stars of degree strictly less than the degree of node ϕ .

Suppose ϕ becomes a new star center. We show that its satellite neighbors γ cannot become star centers. Two cases arise. (Case 1) γ_i is not a star center because its degree is smaller than the degree of the new star center that covers ϕ in the new cover. (Case 2) γ_i is not a star center because it is a satellite of a much larger star, so its degree is larger than the degree of the new star that covers ϕ . But this condition still holds after making the new star. This completes the proof sketch for the termination lemma and. It follows that the types assigned by the on-line algorithm are correct, in other words that there exists an off-line algorithm that produces the same cover. \square

4.3 Running Time Analysis and Experimental Results

In this section, we argue that the running time of the on-line star algorithm is excellent in practice, asymptotically matching the running time of the off-line star algorithm ($\Theta(V + E)$) to within lower order factors. We first note, however, that there exist worst-case thresholded similarity graphs G_σ and corresponding vertex insertion sequences which cause the on-line star algorithm to run in $\Theta(V^2)$ time.³ These graphs and insertion sequences rarely arise in practice though. An analysis more closely modeling practice is the random graph model in which G_σ is a random graph and the insertion sequence is random. In this model, the *expected* running time of the on-line star algorithm can be determined.

In the sections that follow, we first give intuition for the expected running time of the on-line star algorithm. In subsequent sections, we give experimental

³An example is a graph consisting of two connected vertices A and B of very high but identical degree (not both of which can be star centers) and an insertion sequence which causes the “local” center to repeatedly switch between A and B .

results showing that the on-line star algorithm is quite efficient with respect to two different types of random data and a large collection of real documents.

4.3.1 Intuition

We have implemented the on-line star algorithm using a heap for the priority queue and simple linked lists for the various lists required. The time required to insert a new vertex and associated edges into a thresholded similarity graph and to appropriately update the star cover is largely governed by the number of stars that are broken during the update, since breaking stars requires inserting new elements into the priority queue. In practice, very few stars are broken during any given update. This is due partly to the fact that relatively few stars exist at any given time (as compared to the number of vertices or edges in the thresholded similarity graph) and partly to the fact that the likelihood of breaking any individual star is also small. We begin with the former, noting that the number of stars expected to cover a random graph $G_{n,p}$ is only $\Theta(\log n)$.

Theorem 4 *The expected size of the star cover for $G_{n,p}$ is $\frac{\log n}{\log(\frac{1}{1-p})}$.*

Proof: The star cover algorithm is greedy: it iterates by selecting the highest degree vertex not yet covered as a star center and marking this node and all its adjacent vertices as covered. Each iteration creates a new star. We will argue that the number of iterations is $\frac{\log n}{\log(\frac{1}{1-p})}$. The argument relies on the random graph model described in Section 2.1.

Let $G_{n,p}$ correspond to a random graph on n vertices where each edge exists with probability p . The degree of each vertex of G is distributed binomially: $Pr[\text{deg} = k] = \text{bin}(k; n-1, p) = \binom{n-1}{k} p^k (1-p)^{n-1-k}$. The mean of this distribution is $\mu = (n-1)p$ and its variance is $\sigma = \sqrt{(n-1)p(1-p)}$. Note that while the degrees of the vertices do exhibit some dependence, for practical purposes they can be considered independent [Bol95]. This means that on average, each star covers $(n-1)p + 1 \geq np$ vertices⁴. Since the np vertices covered by each star are randomly chosen, there will be some overlap between the star covers. Each new star leaves uncovered a $(1-p)$ fraction of the previously uncovered vertices. In other words, after the first iteration, $(1-p)n$ vertices remain uncovered. After i iterations, $(1-p)^i n$ vertices remain uncovered. The algorithm terminates when all the vertices are covered, or $(1-p)^i n < 1$. By taking logs of both

⁴The star covers its center and $(n-1)p$ satellites.

sides of this inequality, it follows that $i > \frac{\log n}{\log(\frac{1}{1-p})}$ is sufficient. \square

Thus, the number of stars is expected to be relatively small. Furthermore, the probability any individual star will be broken is quite small as well. A star can only be broken if the star center has the same degree as one of its associated satellite vertices and if the vertex being added to the graph is connected to that satellite but not to the star center.⁵ In practice, the expected number of stars broken during an update is a small constant even for graphs containing thousands of vertices (though asymptotically it is certainly a slowly growing function of n). In Figures 5, 6, and 7, we give experimental results showing that the total number of stars broken during runs on three different types of data is roughly a linear function of the number of vertices; thus, the expected number of stars broken during any given update is roughly a constant (or more likely a very slowly growing function of n).

The time to break a star is roughly proportional to its size (the degree of its associated star center), and since the degrees of all vertices are expected to be similar in distribution ($\text{bin}(k; n-1, p)$), this is on the order of the number of edges being inserted into the graph. Since only a constant number of stars are expected to be broken, the expected time to perform an update will be roughly proportional to the number of edges inserted in the graph during the update. Thus, the total time to perform n updates should be roughly proportional to the total number of edges in the final graph. In the sections that follow, we give experimental results which confirm this fact.

4.3.2 Experimental results

We have experimented with the on-line clustering algorithm on three scenarios. The first type of data matches our random graph model and consists of random similarity graphs. While this type of data is useful as a benchmark for the running time of the algorithm, it does not satisfy the geometric constraints of the vector space model. To match reality better, we generated another set of random graphs where all the vertices in the graph are on an n -dimensional sphere. Finally, we used real data from the TREC collection⁶ as a third type of benchmark for the algorithm.

⁵Once a star is broken during an update, however, other stars can be broken in different ways via a cascading effect.

⁶TREC is the annual text retrieval conference. TREC is organized as a competition. Each participant is given on the order of 5 giga bytes of data and a standard set of queries on which to test their systems. The results and the system descriptions are presented as papers at the TREC conference.

We now detail our data generation procedure and the experimental running time of the on-line star algorithm on each data type.

Generating Random Data We run the on-line star cover algorithm on a random graph with 1000 nodes. The edges in this graph were inserted randomly with probability $p = 0.2$. The on-line algorithm was run 30 times. Each time, the vertices of the random graph were inserted in random order. The results were averaged over the 30 experiments. Figure 8 shows the data from these experiments. Note that the the running time is linear in the number of edges in the graph, and we can see the effects of lower order terms.

Generating Random Data on the Sphere

While the random data generation procedure described above is very useful in evaluating clustering algorithms, the data created will not necessarily meet the geometric constraints imposed by the vector space model on real data. In this section, we briefly describe a procedure for generating random clustering data which does meet the geometric constraints imposed by the vector space model.

In the vector space model, documents are represented by vectors in a high-dimensional space, and the similarity between pairs of documents is given by the cosine of the angle between the associated vectors. In the previous sections, we described a mechanism for generating the *similarity graph* associated with a collection. In this new data generation procedure, we instead randomly create the *vectors* in high-dimensional space which correspond to documents, and then construct the associated similarity graph from these vectors. In brief, well-spaced *cluster centers* are generated on a unit sphere of high-dimension, and the clusters of documents themselves are generated by randomly perturbing these cluster centers. By carefully varying the “spacing” of the cluster centers as well as the amount of perturbation allowed in generated the cluster documents, we can again allow for a specified overlap of clusters as well as a varying degree of faulty data.

We run the on-line star cover algorithm on a random sphere graph consisting of 1050 vertices generated as explained above. Figure 9 shows running-time data from this experiment. Note that the the running time is linear in the number of edges in the graph, and we can see the effects of lower order terms.

Experiments with real data We run the on-line star cover algorithm on a document collection that

consists of a slice of TREC documents augmented with our department’s technical reports. The resulting collection consists of 2224 documents. We ran 4 experiments. Each time we set a different threshold and added the similarity graph nodes in random order. The results of these experiments were averaged and the running time measurements appear to be linear in the number of edges of the similarity graph. Figure 10 shows the data from these experiments. Note that the the running time is linear in the number of edges in the graph, and we can see the effects of lower order terms.

5 Discussion

We presented and analyzed an off-line clustering algorithm for static information organization and an on-line clustering algorithm for dynamic information organization. We discussed the random graph model for analyzing these algorithms and showed that in this model, the algorithms have expected running time that is linear in the number of edges. The data we gathered from experimenting with these algorithms provides support for the validity of our model and analysis. The empirical tests show that both algorithms have an asymptotic linear time performance in the number of edges in the graph.

This work departs from previous clustering algorithms used in information retrieval that use a fixed number of clusters for partitioning the space. Since the number of clusters produced by our algorithms is given by the underlying topic structure in the information system, our clusters are dense and accurate. Our work extends previous results [HP96] that support using clustering for browsing applications and presents positive evidence for the cluster hypothesis. In [APR97], we argue that by using a clustering algorithm that guarantees the cluster quality through separation of dissimilar documents and aggregation of similar documents, clustering is beneficial for information retrieval tasks that require high precision and high recall. Precision-recall are the standard measurements for the performance of an information retrieval algorithm [Sal89].

Acknowledgements

This work is supported in part by the Navy under contract ONR N00014-95-1-1204. Special thanks to Philip Klein for stimulating discussions in the initial phase of this work.

References

- [AB84] M. Aldenderfer and R. Blashfield, *Cluster Analysis*, Sage, Beverly Hills, 1984.

- [APR97] J. Aslam, K. Pelehov, and D. Rus, Generating, visualizing, and evaluating high-accuracy clusters for information organization, Technical Report PCS-TR97-319, Department of Computer Science, Dartmouth, 1997.
- [Bol95] B. Bollobás, *Random Graphs*, Academic Press, London, 1995.
- [Can93] F. Can, Incremental clustering for dynamic information processing, in *ACM Transactions on Information Systems*, no. 11, pp143-164, 1993.
- [CCFM97] M. Charikar, C. Chekuri, T. Feder, and R. Motwani, Incremental clustering and dynamic information retrieval, in *Proceedings of the 29th Symposium on Theory of Computing*, 1997.
- [Cro80] W. B. Croft. A model of cluster searching based on classification. *Information Systems*, 5:189-195, 1980.
- [Cro77] W. B. Croft. Clustering large files of documents using the single-link method. *Journal of the American Society for Information Science*, pp189-195, November 1977.
- [CKP93] D. Cutting, D. Karger, and J. Pedersen. Constant interaction-time scatter/gather browsing of very large document collections. In *Proceedings of the 16th SIGIR*, 1993.
- [FG88] T. Feder and D. Greene, Optimal algorithms for approximate clustering, in *Proceedings of the 20th Symposium on Theory of Computing*, pp 434-444, 1988.
- [HP96] M. Hearst and J. Pedersen. Reexamining the cluster hypothesis: Scatter/Gather on Retrieval Results. In *Proceedings of the 19th SIGIR*, 1996.
- [HS86] D. Hochbaum and D. Shmoys, A unified approach to approximation algorithms for bottleneck problems, *Journal of the ACM*, no. 33, pp533-550, 1986.
- [JD88] A. Jain and R. Dubes. *Algorithms for Clustering Data*, Prentice Hall 1988.
- [JR71] N. Jardine and C.J. van Rijsbergen. The use of hierarchical clustering in information retrieval, 7:217-240, 1971.
- [KP93] G. Kortsarz and D. Peleg. On choosing a dense subgraph. In *Proceedings of the 34th Annual Symposium on Foundations of Computer Science (FOCS)*, 1993.
- [LY94] C. Lund and M. Yannakakis. On the hardness of approximating minimization problems. *Journal of the ACM* 41, 960-981, 1994.
- [Rij79] C.J. van Rijsbergen. *Information Retrieval*. Butterworths, London, 1979.
- [Sal89] G. Salton. *Automatic Text Processing: the transformation, analysis, and retrieval of information by computer*, Addison-Wesley, 1989.
- [Sal91] G. Salton. The Smart document retrieval project. In *Proceedings of the Fourteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval*, pages 356-358.
- [Tur90] H. Turtle. Inference networks for document retrieval. PhD thesis. University of Massachusetts, Amherst, 1990.
- [Voo85] E. Voorhees. The cluster hypothesis revisited. In *Proceedings of the 8th SIGIR*, pp 95-104, 1985.
- [Wil88] P. Willett. Recent trends in hierarchical document clustering: A critical review. *Information Processing and Management*, 24:(5):577-597, 1988.
- [Wor71] S. Worona. Query clustering in a large document space. In Ed. G. Salton, *The SMART Retrieval System*, pp 298-310. Prentice-Hall, 1971.
- [Zuc93] D. Zuckerman. NP-complete problems have a version that's hard to approximate. In *Proceedings of the Eight Annual Structure in Complexity Theory Conference*, IEEE Computer Society, 305-312, 1993.

```

UPDATE( $\alpha, L$ )
1  $\alpha.type \leftarrow satellite$ 
2  $\alpha.degree \leftarrow 0$ 
3  $\alpha.adj \leftarrow \emptyset$ 
4  $\alpha.centers \leftarrow \emptyset$ 
5 forall  $\beta$  in  $L$ 
6    $\alpha.degree \leftarrow \alpha.degree + 1$ 
7    $\beta.degree \leftarrow \beta.degree + 1$ 
8   INSERT( $\beta, \alpha.adj$ )
9   INSERT( $\alpha, \beta.adj$ )
10  if ( $\beta.type = center$ )
11    INSERT( $\beta, \alpha.centers$ )
12  else
13     $\beta.inQ \leftarrow true$ 
14    ENQUEUE( $\beta, Q$ )
15  endif
16 endfor
17  $\alpha.inQ \leftarrow true$ 
18 ENQUEUE( $\alpha, Q$ )
19 while ( $Q \neq \emptyset$ )
20    $\phi \leftarrow \text{EXTRACTMAX}(Q)$ 
21   if ( $\phi.centers = \emptyset$ )
22      $\phi.type \leftarrow center$ 
23     forall  $\beta$  in  $\phi.adj$ 
24       INSERT( $\phi, \beta.centers$ )
25     endfor
26   else
27     if ( $\forall \delta \in \phi.centers, \delta.degree < \phi.degree$ )
28        $\phi.type \leftarrow center$ 
29       forall  $\beta$  in  $\phi.adj$ 
30         INSERT( $\phi, \beta.centers$ )
31       endfor
32       forall  $\delta$  in  $\phi.centers$ 
33          $\delta.type \leftarrow satellite$ 
34         forall  $\mu$  in  $\delta.adj$ 
35           DELETE( $\delta, \mu.centers$ )
36           if ( $\mu.inQ = false$ )
37              $\mu.inQ \leftarrow true$ 
38             ENQUEUE( $\mu, Q$ )
39           endif
40         endfor
41       endfor
42     endif
43   endif
44    $\phi.inQ \leftarrow false$ 
45 endwhile

```

Figure 4: The on-line star algorithm for clustering.

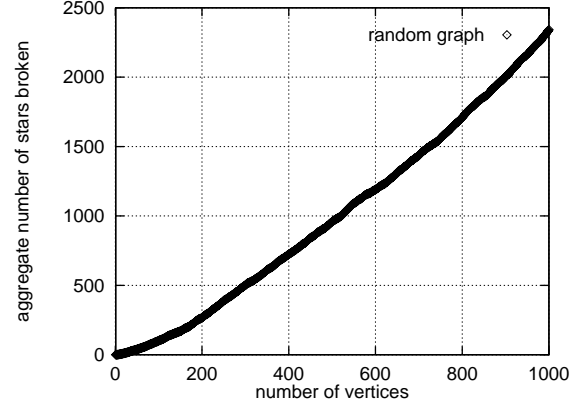


Figure 5: The dependence of the number of broken stars on the number of vertices in a random graph.

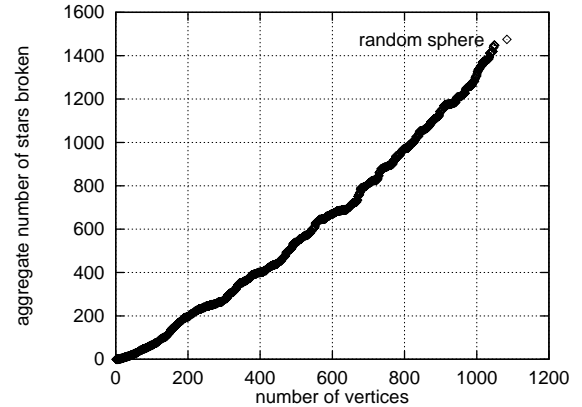


Figure 6: The number of broken stars on the number of vertices in a random sphere graph.

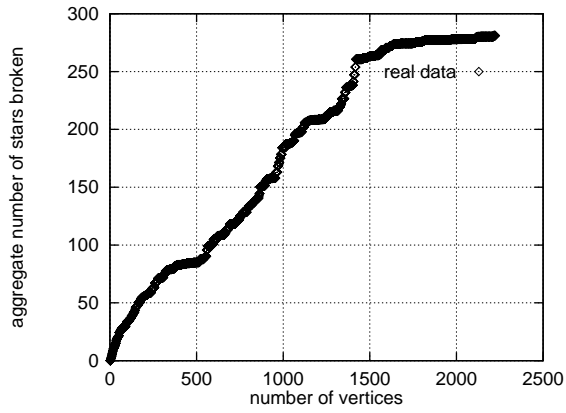


Figure 7: The number of broken stars on the number of vertices in for real data.

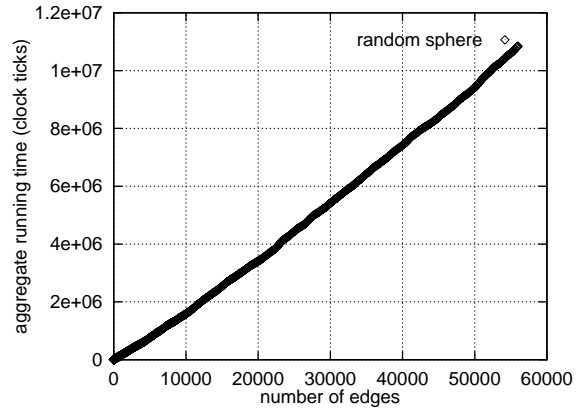


Figure 9: This figure shows the dependence of the running time of the on-line star algorithm on the number of edges on the sphere.

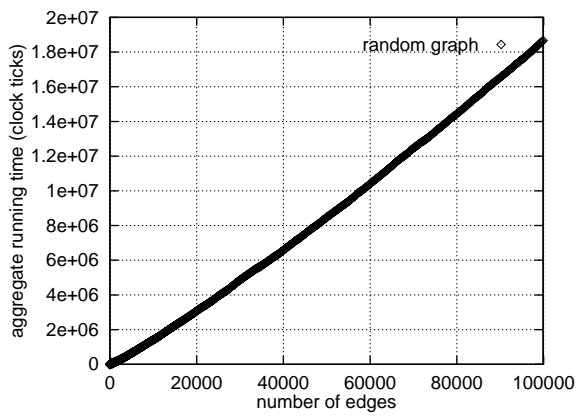


Figure 8: This figure shows the dependence of the running time of the on-line star algorithm on the number of edges in a random graph.

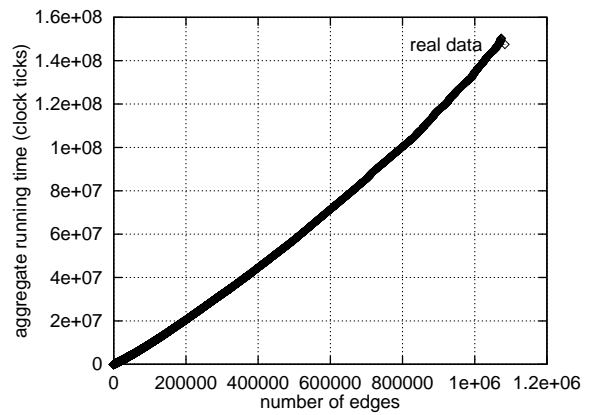


Figure 10: This figure shows the dependence of the running time of the on-line star algorithm on the number of edges in a real collection.