

COS 435: Notes on Clustering Algorithms for 3/30/06

On 3/28/06 I presented on the board single-pass agglomerative clustering and agglomerative clustering for “single link” cluster distance using minimum spanning trees. These are covered in the assigned reading in essentially the same way. On 3/30/06 we covered general hierarchical agglomerative clustering, the k-means algorithm, and divisive clustering. Martin Makowiecki’s slides are posted for his presentation on the k-means algorithm. In this note, I am summarizing my development on the board of general hierarchical agglomerative clustering and divisive clustering. These may vary somewhat from the presentations in the book. – there are many variations that use the same principles. I am also adding a *small* amount of detail that I skipped due to time constraints in class.

General hierarchical agglomerative algorithm

Uses any computable cluster similarity measure $\text{sim}(C_i, C_j)$:

1. For n objects v_1, \dots, v_n , assign each to a singleton cluster $C_i = \{v_i\}$.
2. repeat {
 - a. identify two most similar clusters C_j and C_k (could be ties – chose one pair)
 - b. delete C_j and C_k and add $(C_j \cup C_k)$ to the set of clusters} until only one cluster

Dendograms diagram the sequence of cluster merges.

Running time if must compare all clusters each time to do 2a:

$$\sum_{i=1}^{n-1} \left(\text{cost of } 2b \text{ for } i^{\text{th}} \text{ merge} + \text{cost to compute/look up } \text{sim}(C_j, C_k) \text{ for all clusters } C_j, C_k \text{ in the } i^{\text{th}} \text{ repetition of the loop} \right)$$

If one keeps a matrix of $\text{sim}(C_j, C_k)$ for all current pairs, then the running time is:

$$\sum_{i=1}^{n-1} \left(\text{cost of } 2b \text{ for } i^{\text{th}} \text{ merge and update of } \text{sim}(C_i, C_j) \text{ matrix} + \frac{(n-i+1)(n-i)}{2} \right)$$

If the update of the $\text{sim}(C_j, C_k)$ matrix can be done in $O(n^2)$ steps, the whole algorithm is $O(n^3)$.

For the single-link distance measure, recall that one algorithm builds a minimum spanning tree by sorting all the edges between objects and then processing the edges in sorted order, smallest first, and keeping for the MST any edge that does not form a cycle. This algorithm is doing what the general hierarchical agglomerative algorithm does – the partial MSTs are the clusters. By pre-sorting all the edges between objects, the combined cost of step 2a over all repetitions is $O(n^2)$ – assuming we can detect in constant time that an edge would cause a cycle. This is easy with a random access table in main memory giving the cluster for each object and a challenge if the amount of data forces some of the

tables onto disk. The cost of step 2b to update such a table is $O(n)$ in main memory. This gives a total running time of $O(n^2 \log n)$ assuming random access tables can be used.

Space: for a general similarity matrix at least space $\frac{1}{2}n^2$ is required, assuming $\text{sim}(v_i, v_j) = \text{sim}(v_j, v_i)$. If one can calculate $\text{sim}(v_i, v_j)$, e.g. using a cosine metric for vector objects, then one might need only $O(n)$ space. However, this might cost in running time because one might have to recalculate all similarities between objects or clusters each time one needed a most similar pair. $O(n)$ space does not provide enough space to store the similarities between all pairs of clusters until the number of clusters has been reduced to about \sqrt{n} .

For “complete linkage” and “all pairs average distance” similarity measures for clusters, we can update the similarity between clusters after a merge in constant time by taking advantage of the following equations (in terms of distance):

- For “complete linkage”,

$$\text{dist}(C_j \cup C_k, C_p) = \max(\text{dist}(C_j, C_p), \text{dist}(C_k, C_p))$$
- For “all pairs average distance”,

$$\text{dist}(C_j \cup C_k, C_p) = 1 / (|C_j| + |C_k|) * (|C_j| * \text{dist}(C_j, C_p) + |C_k| * \text{dist}(C_k, C_p))$$

Then we can maintain a matrix of $\text{sim}(C_j, C_k)$ and update it in $O(n^2)$ steps per merge.

Divisive algorithms and cut-based measures of cluster similarity

Cut-based measures of similarity measure the weakness of the connection between objects in different clusters rather the strength of the connection of objects within a cluster.

Divisive algorithms start with one cluster of all objects and recursively split up the cluster and subsequent clusters. Since this is a cutting operation, cut-based measures seem to be a natural choice. However, it is not necessary to use a cut-based measure with a divisive algorithm.

There are many cut-based similarity measures for clusters. We looked at one:

Let U be the set of all objects, i.e. $U = \{v_1, \dots, v_n\}$. For any clustering C_1, C_2, \dots, C_k of the objects, $U = \bigcup_{i=1, \dots, k} C_i$. Define:

$$\text{cutcost}(C_p, U - C_p) = \sum_{(v_i \text{ in } C_p, v_j \text{ in } U - C_p)} \text{sim}(v_i, v_j).$$

$$\text{intracost}(C_p) = \sum_{(v_i, v_j \text{ in } C_p)} \text{sim}(v_i, v_j).$$

Then the cost of a clustering C_1, \dots, C_k is defined to be:

$$\text{cost}(C_1, \dots, C_k) = \sum_{p=1, \dots, k} (\text{cutcost}(C_p, U-C_p) / \text{intracost}(C_p))$$

For this measure, the cost of a cluster measures the total similarity to objects outside the cluster (its cutcost) relative to the total similarity among items in the cluster. Finding a clustering that minimizes cost (C_1, \dots, C_k) is called *mi-max cut* clustering.

An iterative improvement algorithm to heuristically find a clustering with low min-max cut is as follows:

```

Choose initial partition  $C_1, \dots, C_k$ 
repeat {
  unlock all vertices
  repeat {
    choose some  $C_i$  at random
    choose an unlocked vertex  $v_j$  in  $C_i$ 
    if moving  $v_j$  to any other cluster improves min-max cut cost
      move vertex  $v_j$  to cluster giving best improvement
    "lock" vertex  $v_j$ 
  } until all vertices locked
}until converge

```

Convergence is usually determined by improvement being smaller than some chosen threshold between executions of the inner "repeat" loop. Vertex "locking" just insures that all vertices are examined before examining any vertex twice. There are many variations on this iterative improvement algorithm.

This iterative improvement algorithm shares several properties with the k-means algorithm: the number of clusters, k , is chosen in advance; an initial clustering is chosen (possibly at random); iterative improvement is used to modify a clustering to a better clustering. The most important difference is that the min-max cut cost is minimized rather than the sum of the squares of the distances from vertices to their cluster centroids.