

CSCI 5417
Information Retrieval Systems

Jim Martin

Lecture 15
10/13/2011

Today 10/13

- More Clustering
 - Finish flat clustering
 - Hierarchical clustering

K-Means

- Assumes documents are real-valued vectors.
- Clusters based on *centroids* (aka the *center of gravity* or mean) of points in a cluster, c :

$$\mu(c) = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$$

- Iterative reassignment of instances to clusters is based on distance to the current cluster centroids.
 - (Or one can equivalently phrase it in terms of similarities)

10/17/11

CSCI 5417 - IR

3

K-Means Algorithm

Select K random docs $\{s_1, s_2, \dots, s_K\}$ as seeds.
Until stopping criterion:

For each doc d_i :

Assign d_i to the cluster c_j
such that $dist(d_i, s_j)$ is minimal.

For each cluster c_j

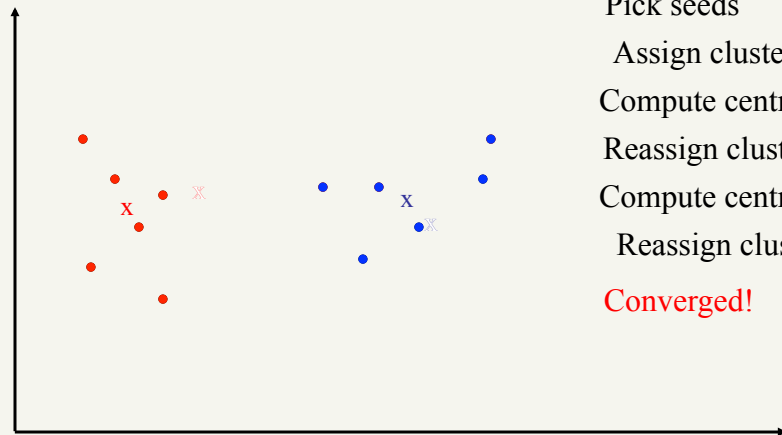
$$s_j = m(c_j)$$

10/17/11

CSCI 5417 - IR

4

K Means Example ($K=2$)



Pick seeds

Assign clusters

Compute centroids

Reassign clusters

Compute centroids

Reassign clusters

Converged!

10/17/11

CSCI 5417 - IR

5

Termination conditions

- Several possibilities
 - A fixed number of iterations
 - Doc partition unchanged
 - Centroid positions don't change

10/17/11

CSCI 5417 - IR

6

Convergence

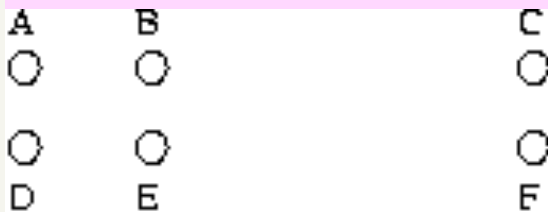
- Why should the K -means algorithm ever reach a *fixed point*?
 - A state in which clusters don't change.
- K -means is a special case of a general procedure known as the *Expectation Maximization (EM) algorithm*.
 - EM is known to converge.
 - Number of iterations could be large.

But in practice usually isn't

Seed Choice

- Results can vary based on random seed selection.
- Some seeds can result in poor convergence rate, or convergence to sub-optimal clusterings.
 - Select good seeds using a heuristic (e.g., doc least similar to any existing mean)
 - Try out multiple starting points
 - Initialize with the results of another method.

Do this with $K=2$



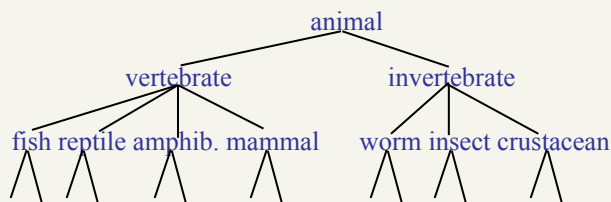
10/17/11

CSCI 5417 - IR

9

Hierarchical Clustering

- Build a tree-based hierarchical taxonomy (*dendrogram*) from a set of unlabeled examples.



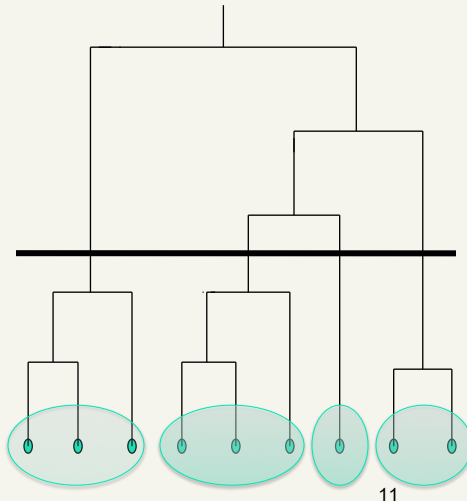
10/17/11

CSCI 5417 - IR

10

Dendrogram: Hierarchical Clustering

- Traditional clustering partition is obtained by cutting the dendrogram at a desired level: each connected component forms a cluster.



Break

- Past HW
 - Best score on part 2 is .437
 - Best approaches
 - Multifield indexing of title/keywords/abstract
 - Snowball (English), Porter
 - Tuning the stop list
 - Ensemble (voting)
 - Mixed results
 - Boosts
 - Relevance feedback

Descriptions

- For the most part, your approaches were pretty weak (or your descriptions were)
 - Failed to report R-Precision
 - Use of some kind of systematic approach
 - X didn't work
 - Interactions between approaches
 - Lack of details
 - Use relevance feedback and it gave me Z
 - I changed the stop list
 - Boosted the title field
 - Etc.

10/17/11

CSCI 5417 - IR

13

Next HW

- Due 10/25
- I have a new untainted test set
 - So don't worry about checking for the test document; it won't be there

10/17/11

CSCI 5417 - IR

14

Hierarchical Clustering algorithms

- **Agglomerative (bottom-up):**
 - Start with each document being a single cluster.
 - Eventually all documents belong to the same cluster.
- **Divisive (top-down):**
 - Start with all documents belong to the same cluster.
 - Eventually each node forms a cluster on its own.
- Does not require the number of clusters k to be known in advance
 - But it does need a cutoff or threshold parameter condition

10/17/11

CSCI 5417 - IR

15

Hierarchical -> Partition

- Run the algorithm to completion
 - Take a slice across the tree at some level
 - Produces a partition
- Or insert an early stopping condition into either top-down or bottom-up

10/17/11

CSCI 5417 - IR

16

Hierarchical Agglomerative Clustering (HAC)

- Assumes a similarity function for determining the similarity of two instances and two clusters.
- Starts with all instances in separate clusters and then repeatedly joins the two clusters that are most similar until there is only one cluster.
- The history of merging forms a binary tree or hierarchy.

10/17/11

CSCI 5417 - IR

17

Hierarchical Clustering

- Key problem: as you build clusters, how do you **represent each cluster**, to tell which pair of clusters is closest?

10/17/11

CSCI 5417 - IR

18

"Closest pair" in Clustering

- Many variants to defining closest pair of clusters
 - Single-link
 - Similarity of the most cosine-similar
 - Complete-link
 - Similarity of the "furthest" points, the least cosine-similar
 - "Center of gravity"
 - Clusters whose centroids (centers of gravity) are the most cosine-similar
 - Average-link
 - Average cosine between all pairs of elements

10/17/11

CSCI 5417 - IR

19

Single Link Agglomerative Clustering

- Use maximum similarity of pairs:

$$sim(c_i, c_j) = \max_{x \in c_i, y \in c_j} sim(x, y)$$

- Can result in "straggly" (long and thin) clusters due to chaining effect.
- After merging c_i and c_j , the similarity of the resulting cluster to another cluster, c_k , is:

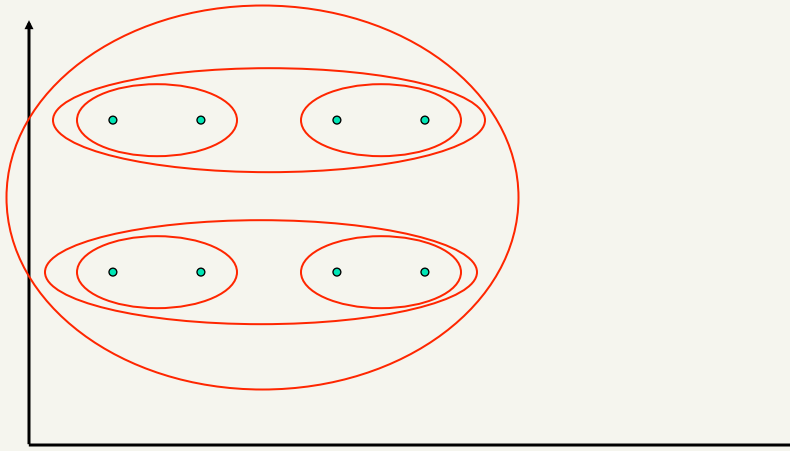
$$sim((c_i \cup c_j), c_k) = \max(sim(c_i, c_k), sim(c_j, c_k))$$

10/17/11

CSCI 5417 - IR

20

Single Link Example



10/17/11

CSCI 5417 - IR

21

Complete Link Agglomerative Clustering

- Use minimum similarity of pairs:

$$sim(c_i, c_j) = \min_{x \in c_i, y \in c_j} sim(x, y)$$

- Makes "tighter," spherical clusters that are typically preferable.
- After merging c_i and c_j , the similarity of the resulting cluster to another cluster, c_k , is:

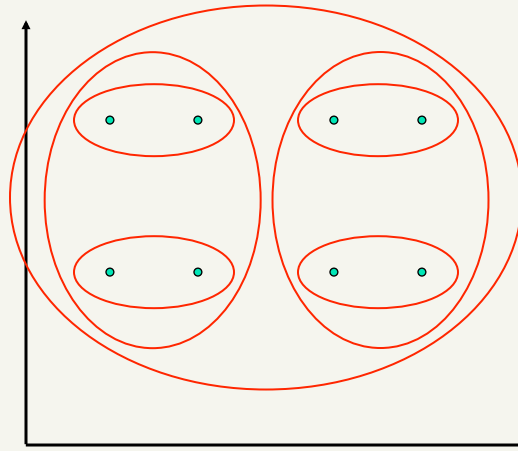
$$sim((c_i \cup c_j), c_k) = \min(sim(c_i, c_k), sim(c_j, c_k))$$

10/17/11

CSCI 5417 - IR

22

Complete Link Example



10/17/11

CSCI 5417 - IR

23

Misc. Clustering Topics

- Clustering terms
- Clustering people
- Feature selection
- Labeling clusters

10/17/11

CSCI 5417 - IR

24

Term vs. document space

- So far, we clustered docs based on their similarities in term space
- For some applications, e.g., topic analysis for inducing navigation structures, you can “dualize”:
 - Use docs as axes
 - Represent (some) terms as vectors
 - Cluster terms, *not* docs

10/17/11

CSCI 5417 - IR

25

Clustering people

- Take documents (pages) containing mentions of ambiguous names and partition the documents into bins with identical referents.
 - SemEval competition
 - Web People Search Task: Given a name as a query to google, cluster the top 100 results so that each cluster corresponds to a real individual out in the world

10/17/11

CSCI 5417 - IR

26

Labeling clusters

- After clustering algorithm finds clusters - how can they be useful to the end user?
- Need pithy label for each cluster
 - In search results, say “Animal” or “Car” in the *jaguar* example.

10/17/11

CSCI 5417 - IR

27

How to Label Clusters

- Show titles of typical documents
 - Titles are easy to scan
 - Authors create them for quick scanning
 - But you can only show a few titles which may not fully represent cluster
- Show words/phrases prominent in cluster
 - More likely to fully represent cluster
 - Use distinguishing words/phrases
 - Differential labeling
 - But harder to scan

10/17/11

CSCI 5417 - IR

28

Labeling

- Common heuristics - list 5-10 most frequent terms in the centroid vector.
 - Drop stop-words; stem.
- Differential labeling by frequent terms
 - Within a collection “Computers”, clusters all have the word **computer** as frequent term.
 - Discriminant analysis of centroids.
- Perhaps better: distinctive noun phrases
 - Requires NP chunking

10/17/11

CSCI 5417 - IR

29

Summary

- In clustering, clusters are inferred from the data without human input (unsupervised learning)
- In practice, it’s a bit less clear. There are many ways of influencing the outcome of clustering: number of clusters, similarity measure, representation of documents, . . .