### Clustering

1

### Informal goal Given set of objects and measure of similarity between them, group similar objects together What mean by "similar"? What is good grouping? Computation time / quality tradeoff



### Applications: Many – biology – astronomy – computer aided design of circuits – information organization – marketing – ...

### Clustering in information search and analysis

- · Group information objects
  - $\Rightarrow$  discover topics
  - ? other groupings desirable
- Clustering versus classifying
  - classifying: have pre-determined classes with example members
  - clustering:
    - get groups of similar objects
    - added problem of labeling clusters by topic
      - -e.g. common terms within cluster of docs.

# Example applications in search Query evaluation: cluster pruning (§7.1.6) cluster all documents choose representative for each cluster evaluate query w.r.t. cluster reps. evaluate query for docs in cluster(s) having most similar cluster rep.(s) Results presentation: labeled clusters cluster only query results e.g. Yippy.com (metasearch)

## Issues What, if any, attributes represent items for clustering purposes? What is measure of similarity between items? General objects and matrix of pairwise similarities Objects with specific properties that allow other specifications of measure Most common: Objects are d-dimensional vectors Euclidean distance cosine similarity What is measure of similarity between clusters?

### Issues continued

- Cluster goals?
  - Number of clusters?
  - flat or hierarchical clustering?
  - cohesiveness of clusters?
- How evaluate cluster results?
  - relates to measure of closeness between clusters
- Efficiency of clustering algorithms

   large data sets => external storage
- Maintain clusters in dynamic setting?
- Clustering methods? MANY!



- In applications, quality of clustering depends on how well solves problem at hand
- Algorithm uses measure of quality that can be optimized, but that may or may not do a good job of capturing application needs.
- Underlying graph-theoretic problems usually NP-complete
  - e.g. graph partitioning
- Usually algorithm not finding optimal clustering

### General types of clustering methods

- constructive: decide in what cluster each object belongs and don't change
   often faster
- iterative improvement: start with a clustering and move objects around to see if can improve clustering

10

often slower but better

Vector model: K- means algorithm

- Well known, well used
- Flat clustering
- Number of clusters picked ahead of time

11

- · Iterative improvement
- Uses notion of centroid
- Typically uses Euclidean distance

K-means overviews
Choose k points among set to be clustered
(al them k centroids)
not required to be in set to be clustered
to required to be in set to be clustered
assignments give initial clustering
Assignments give initial clustering
(ansignments give initial clustering)
Recompute centroids of clusters:
centroid of set of vectors {v<sub>i</sub> | si≤n} = sin s n<sub>i=1</sub> v<sub>i</sub>
New centroids may not be points of original set
Beassign all points to closest centroid
Updates clusters





































## Similarity between clusters, cont. Possible definitions: III. average of pairwise similarity between all pairs of objects, one from each cluster "centroid" similarity IV. average of pairwise similarity between all pairs of distinct objects, including w/in same cluster "group average" similarity Generally no representative point for a cluster; compare K-means If using Euclidean distance as metric centroid bounding box







