15-853: Algorithms in the Real World

Clustering: Lectures 1 and 2

15-853 Page 1 Given a set of objects and a similarity (or distance) measure among the objects, cluster into groups of similar (close) objects

Also called:

- · Unsupervised learning
- Classification
- Typology
- · Numerical taxonomy



Page 2

Applications

Biology: Multiple allignments, evolutionary trees

Business: Market research, risk analysis

Liberal arts: Classifying painters, writers, musicians Sociology: personality types, classifying criminals,

classifying survey results.

Computer Science: compression, information retrieval, text mining, image segmentation, recommender systems, anomaly detection

> 15-853 Page 3

Types of clustering

15-853

Hard vs. soft Hierarchical vs flat Distance vs similarity based

Distance metrics:

- Euclidean, Minkowski, Hamming, Edit, ... Similarity measures:

- Cosine, Kernel functions, or S_{ii} = (1 + d_{ii})

15-853

Main Approaches

Centroid based: K-means

Distribution/Model-based (EM)

Mixture of gaussians

Spectral

Agglomerate

Density based

Neural nets

15-853 Page 5

K-means

<u>K-clustering</u>: given a metric space (X,d), and a point set $S \subset X$ partition S into k sets C_1, C_2, \dots, C_k

<u>Cost</u>: $\phi(C) = \sum_{i=1}^{k} \min_{c \in X} \sum_{x \in C_i} d(c, x)^2$

Goal: minimize the cost

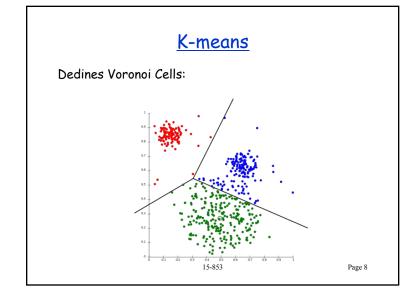
Problem is NP-hard. Can find approximations.

Typically looking for an approximation.

Related to expectation-maximization, mixture of gaussians, and k-median.

15-853

Page 7



Lloyds algorithm for K-means

A greedy local search algorithm:

Start with a set of centers: c1, c2, ..., ck in X Repeat until "convergence":

- Assign each x in S to nearest center
- Update location of each center to minimize sum of distances to points assigned to it

Will converge but perhaps slowly and perhaps to a local minimum

Often tried with many starting sets

Often dimensionality reduction is applied first

age 9

K-means++

Picks the starting set more intelligently:

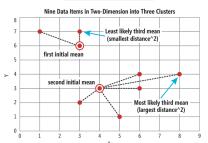
Pick a center uniformly at random from X and add to centers Y (initially empty)

For k-1 steps:

- For each x in X calculate min distance d(x) from points in Y
- Pick an x in X with probability proportional to $d(x)^2$ and add to Y

15-853 Page 10

K-means++



Gives an $8(\ln k + 2)$ approximation even without Lloyds Using Lloyds will only improve the result.

15-853 Page 11

Expectation Maximization (EM)

K-means is a special case.

Start with an arbitrary set of clusters defined by parameters: p1, p2, ..., pk

Repeat until "convergence":

- Expectation: Assign each x in S to cluster that best matches parameters.
 - This can be a probabilistic (soft) assignment
- Maximization: Update parameters to best fit the assignment.

Coverges to a local maximal likelihood estimator. $_{\tiny Page \, 12}$

age 12

EM: mixture of gaussians Here the parameters are (anisotropic) Gaussians. The parameters form a matrix Can deal with elongated structures. More generally can be any parameterized model of the data Useful if you know what form the clusters will have. 15-853 Page 13

Spectral Clustering

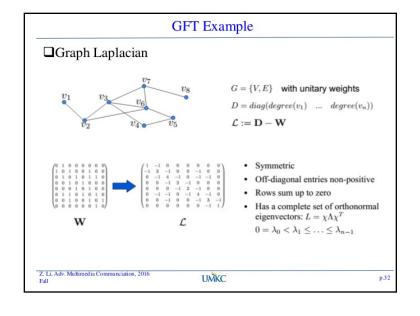
Input is a similarity graph (often sparse)

- Cosine measure (x dot y)/(||x|| ||y||)
- K-nearest neighbor graph

Uses Eigenvectors of the Graph Laplacian

- If W is the weight matrix
- and D is a diagonal matrix summing each row
- L = D W (often normalized)

15-853 Page 14



Spectral Clustering

Used in two possible ways:

- Divisive hierarchical clustering
 Use second eigenvector to split in two
 Recurse on the parts
- K-clustering
 Use first I eigenvectors as a reduced dimensional space
 Use K-means on the result

15-853 Page 16

The Eigenvectors

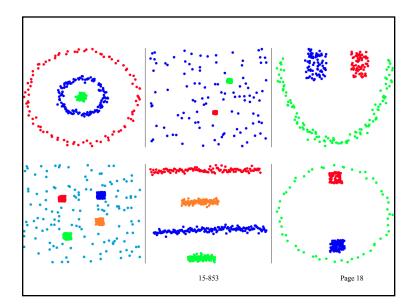
First is trivial (all 1s) with eigenvalue 0 (if normalized) The second gives information about how well the graph can be separated.

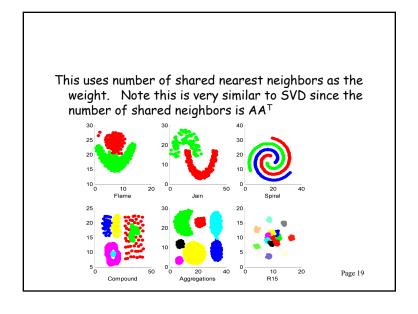
Cheeger constant:

 $E(A,V\setminus A)/|A|$ for any A, |A| < |V|/2A measure of how well graph separates Related to expander graphs (do not separate)

Related to second eigenvalue.

15-853





Agglomerate Clustering

Hierarchical: bottom up Assuming a distance measure

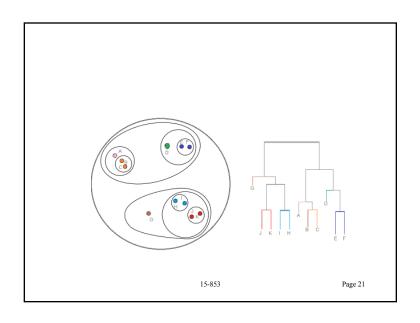
Initially one group per object: G = P

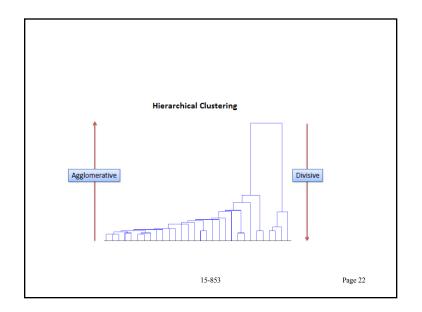
While |G| > 1find "closest" pair in G and join pair into group

Algorithms vary depending on distance between groups.

15-853

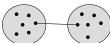
Page 20





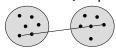
Agglomerate: Distances

Choice 1: min distance (single linkage)

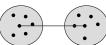


Min Spanning Tree

Choice 2: max distance (complete linkage)



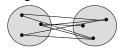
Choice 3: centroid distance



Page 23

Agglomerate: Distances

Choice 4: average distance



Choice 5: Min sum of squares (Ward's method)



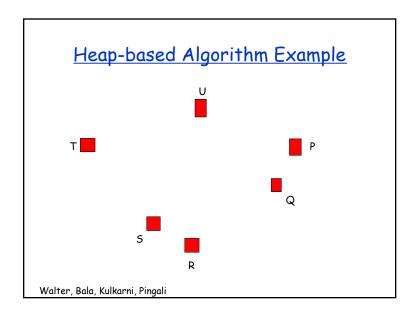


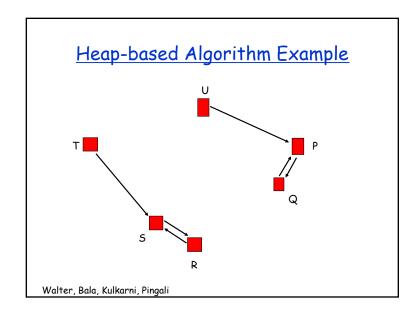
$$SS(G) = \frac{1}{2|G|} \sum_{x \in G, y \in G} d(x, y)^{2}$$

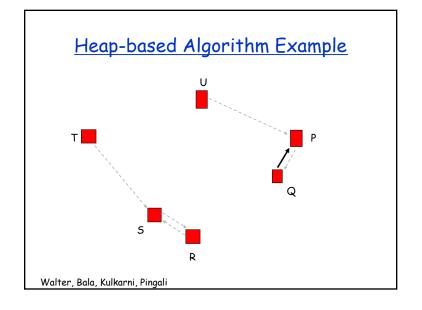
$$SS(G_{12}) - SS(G_{1}) - SS(G_{2})$$

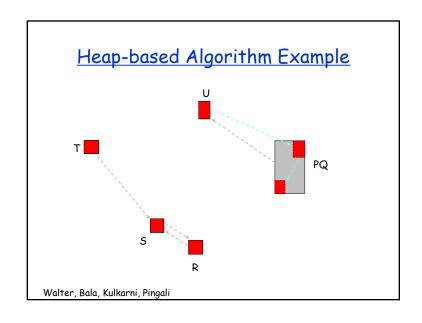
Page 24

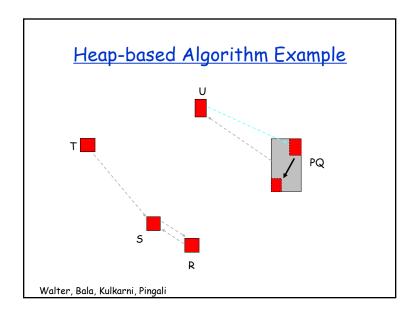
Heap-based Algorithm Initialize KD-Tree with elements Initialize heap with best match for each element Repeat { Remove best pair <A,B> from heap If A and B are active clusters { Create new cluster C = A+B Update KD-Tree, removing A and B and inserting C Use KD-Tree to find best match for C and insert into heap } else if A is active cluster { Use KD-Tree to find best match for A and insert into heap} } until only one active cluster left Walter, Bala, Kulkarni, Pingali

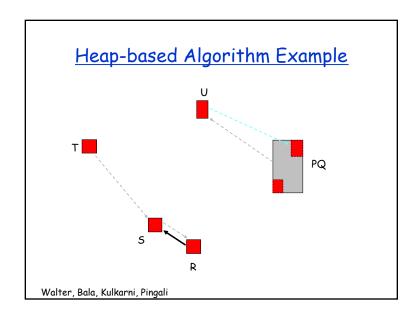


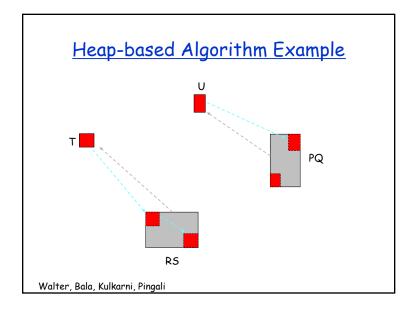












Main Approaches

Centroid based: K-means

Distribution/Model-based (EM)

Mixture of gaussians

Spectral

Agglomerate

Density based

Neural nets

15-853 Page 33

