

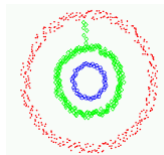
Machine Learning

10-701/15-781, Spring 2008

Spectral Clustering

Eric Xing

Lecture 23, April 14, 2008



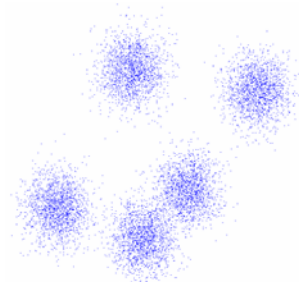
Eric Xing

Reading:

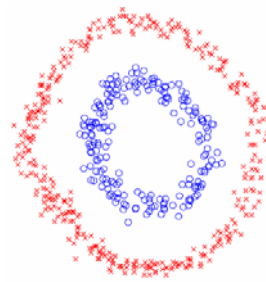
1

Data Clustering

- Two different criteria
 - Compactness, e.g., k-means, mixture models
 - Connectivity, e.g., spectral clustering



Compactness

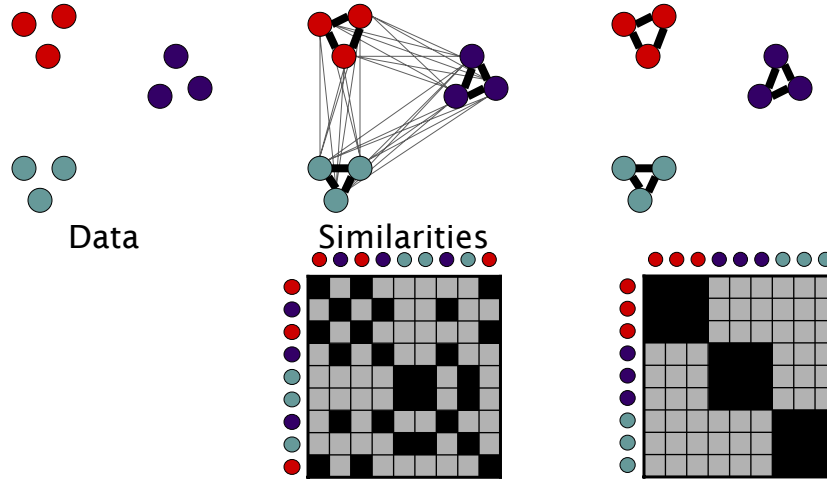


Connectivity

Eric Xing

2

Spectral Clustering



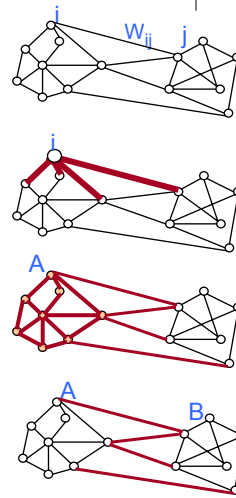
Eric Xing

3

Weighted Graph Partitioning



- Some graph terminology
 - Objects (e.g., pixels, data points)
 $i \in I = \text{vertices of graph } G$
 - Edges $(ij) = \text{pixel pairs with } W_{ij} > 0$
 - Similarity matrix $W = [W_{ij}]$
 - Degree
 $d_i = \sum_{j \in G} S_{ij}$
 $d_A = \sum_{i \in A} d_i$ degree of $A \subseteq G$
 - $\text{Assoc}(A,B) = \sum_{i \in A} \sum_{j \in B} W_{ij}$



Eric Xing

4

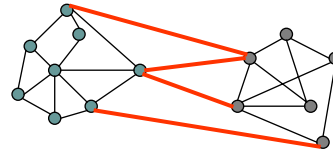
Cuts in a Graph



- (edge) cut = set of edges whose removal makes a graph disconnected

- weight of a cut:

$$\text{cut}(A, B) = \sum_{i \in A} \sum_{j \in B} W_{ij} = \text{Assoc}(A, B)$$



- Normalized Cut criteria: minimum $\text{cut}(A, \bar{A})$

$$\text{Ncut}(A, B) = \frac{\text{cut}(A, B)}{d_A} + \frac{\text{cut}(A, \bar{A})}{d_B}$$

More generally:

$$\text{Ncut}(A_1, A_2 \dots A_k) = \sum_{r=1}^k \left(\frac{\sum_{i \in A_r, j \in V \setminus A_r} W_{ij}}{\sum_{i \in A_r, j \in V} W_{ij}} \right) = \sum_{r=1}^k \left(\frac{\text{cut}(A_r, \bar{A}_r)}{d_{A_r}} \right)$$

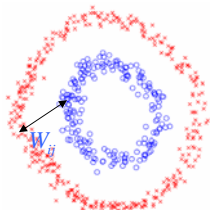
Eric Xing

5

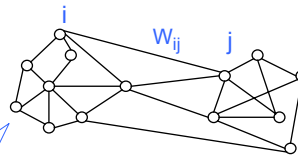
Graph-based Clustering



- Data Grouping



$$W_{ij} = f(d(x_i, x_j))$$



$$G = \{V, E\}$$

- Image segmentation



- Affinity matrix: $W = [w_{i,j}]$
- Degree matrix: $D = \text{diag}(d_i)$
- Laplacian matrix: $L = D - W$
- (bipartite) partition vector:

$$x = [x_1, \dots, x_N]$$

$$= [1, 1, \dots, 1, -1, -1, \dots, -1]$$

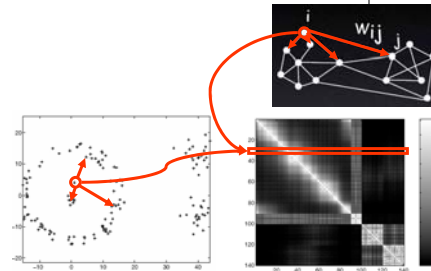
Eric Xing

6

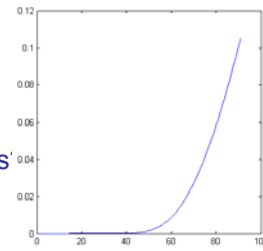
Affinity Function



$$W_{i,j} = e^{-\frac{\|X_i - X_j\|_2^2}{\sigma^2}}$$



- Affinities grow as σ grows \rightarrow
- How the choice of σ value affects the results
- What would be the optimal choice for σ ?



Eric Xing

7

Clustering via Optimizing Normalized Cut



- The normalized cut:

$$Ncut(A, B) = \frac{cut(A, B)}{d_A} + \frac{cut(A, B)}{d_B}$$
- Computing an optimal normalized cut over all possible y (i.e., partition) is NP hard
- Transform Ncut equation to a matrix form (Shi & Malik 2000):

$$\min_x Ncut(x) = \min_y \frac{y^T (D - W) y}{y^T D y} \quad \text{Rayleigh quotient}$$

Subject to: $y \in \{1, -b\}^n$
 $y^T D 1 = 0$

- Still an NP hard problem

$$Ncut(A, B) = \frac{cut(A, B)}{\deg(A)} + \frac{cut(A, B)}{\deg(B)}$$

$$= \frac{(1+x)^T (D-S)(1+x)}{k^T D 1} + \frac{(1-x)^T (D-S)(1-x)}{(1-k)^T D 1}; \quad k = \frac{\sum_{i>0} D(i, i)}{\sum_i D(i, i)}$$

$$= \dots$$

Eric Xing

Relaxation



$$\min_x Ncut(x) = \min_y \frac{y^T (D-W)y}{y^T Dy}$$

Rayleigh quotient

Subject to: $y \in \{1, -b\}^n$
 $y^T D \mathbf{1} = 0$

- Instead, relax into the continuous domain by solving generalized eigenvalue system:

$$\min_y y^T (D-W)y, \quad \text{s.t. } y^T Dy = 1$$

- Which gives: $(D-W)y = \lambda Dy$ *Rayleigh quotient theorem*
- Note that $(D-W)\mathbf{1} = \mathbf{0}$ so, the first eigenvector is $y_0=1$ with eigenvalue 0.
- The second smallest eigenvector is the real valued solution to this problem!!

Eric Xing

9

Algorithm



1. Define a similarity function between 2 nodes. i.e.:

$$w_{i,j} = e^{-\frac{\|x_{(i)} - x_{(j)}\|_2^2}{\sigma_x^2}}$$

2. Compute affinity matrix (W) and degree matrix (D).

3. Solve $(D-W)y = \lambda Dy$

- Do singular value decomposition (SVD) of the graph Laplacian $L = D - W$

$$L = V^T \Lambda V \Rightarrow y^*$$

4. Use the eigenvector with the second smallest eigenvalue, y^* , to bipartition the graph.

- For each threshold k , $A_k = \{i \mid y_i \text{ among } k \text{ largest element of } y^*\}$
 $B_k = \{i \mid y_i \text{ among } n-k \text{ smallest element of } y^*\}$
- Compute $Ncut(A_k, B_k)$
- Output $k^* = \arg \max Ncut(A_k, B_k)$ and A_{k^*}, B_{k^*}

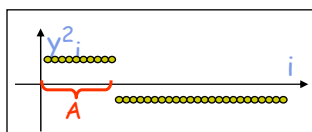
Eric Xing

10

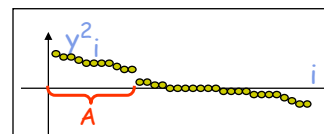
Ideally ...



$$Ncut(A,B) = \frac{y^T (D-S)y}{y^T Dy}, \text{ with } y_i \in \{1,-b\}, y^T D\mathbf{1} = 0.$$



$$(D-S)y = \lambda Dy$$



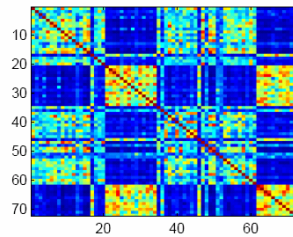
Eric Xing

11

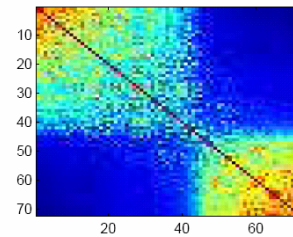
Example



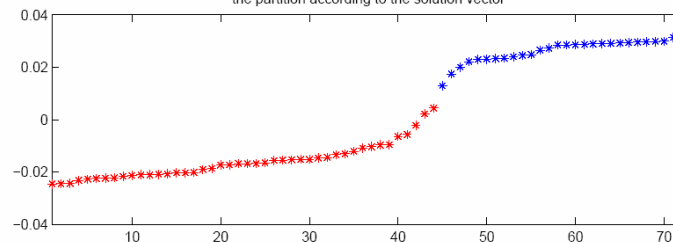
input affinity matrix



affinity matrix reordered according to solution vector



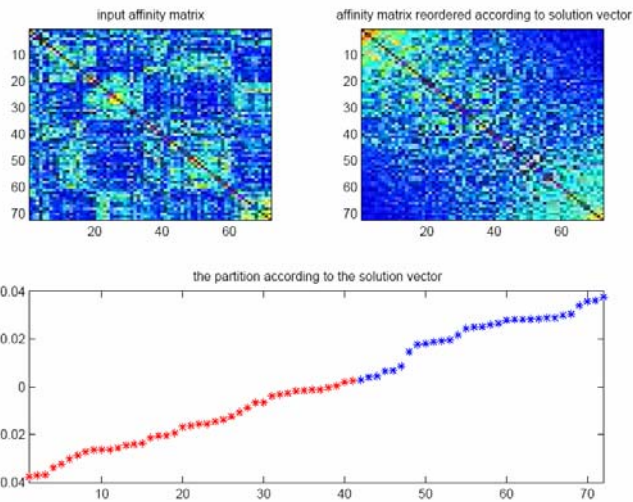
the partition according to the solution vector



Eric Xing

12

Poor features can lead to poor outcome (xing et al 2002)



Eric Xing

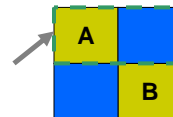
13

Cluster vs. block matrix

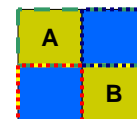


$$Ncut(A, B) = \frac{cut(A, B)}{d_A} + \frac{cut(A, B)}{d_B}$$

$$Degree(A) = \sum_{i \in A, j \in V} W_{i,j}$$



$$Ncut(A, B) = \frac{cut(A, B)}{d_A} + \frac{cut(A, B)}{d_B}$$



Eric Xing

14

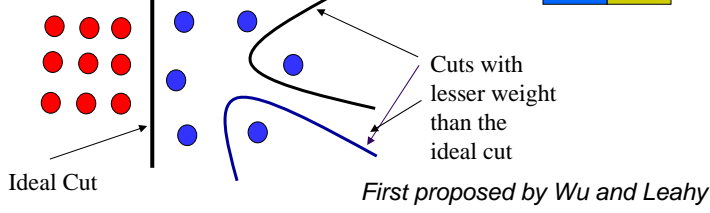
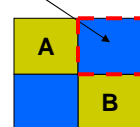
Compare to Minimum cut



- Criterion for partition:

$$\min_{A,B} \text{cut}(A, B) = \min_{A,B} \sum_{i \in A, j \in B} W_{i,j}$$

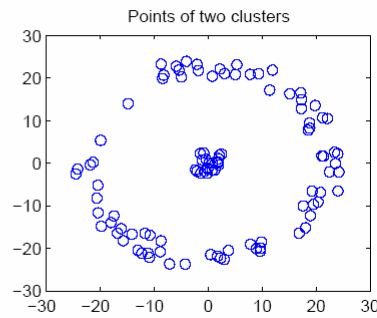
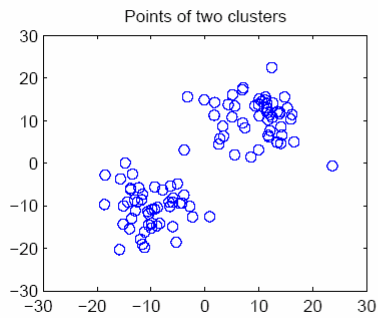
Problem!
Weight of cut is directly proportional to the number of edges in the cut.



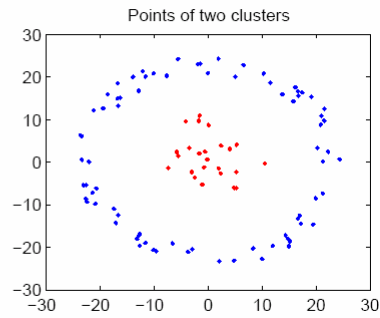
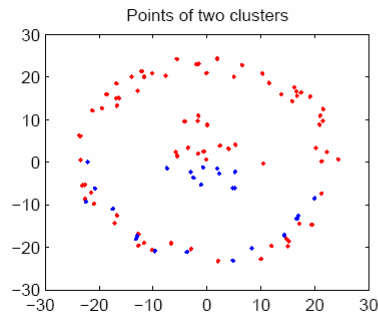
Superior performance?



- K-means and Gaussian mixture methods are biased toward convex clusters



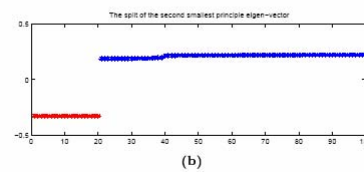
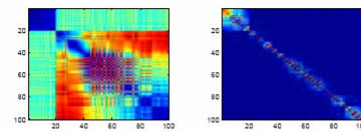
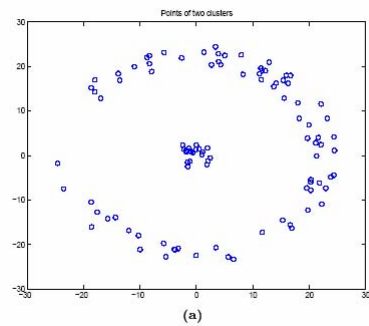
Ncut is superior in certain cases



Eric Xing

17

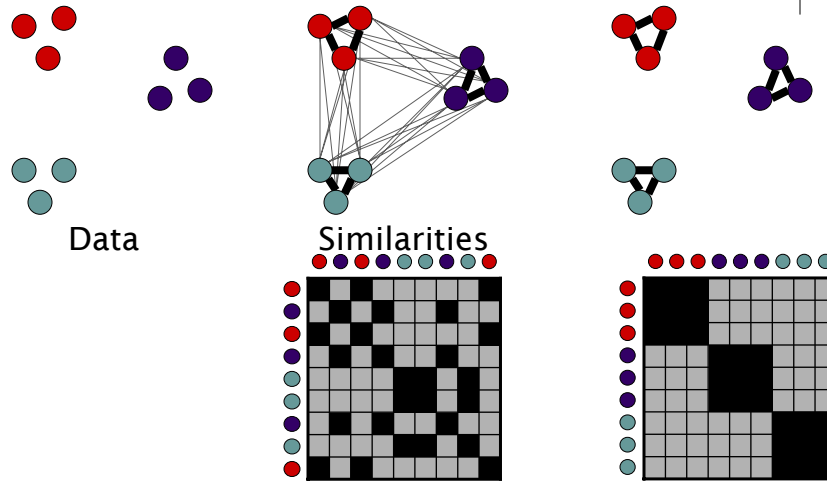
Why?



Eric Xing

18

General Spectral Clustering



Eric Xing

19

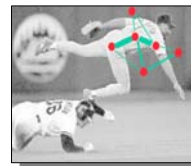
Representation



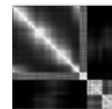
- Partition matrix X :

$$X = [X_1, \dots, X_K]$$

$$X = \begin{matrix} & \begin{matrix} \text{segments} \\ \text{pixels} \end{matrix} \\ \begin{matrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix} & \begin{matrix} \\ \\ \\ \\ \\ \end{matrix} \end{matrix}$$



- Pair-wise similarity matrix W : $W(i, j) = \text{aff}(i, j)$



- Degree matrix D : $D(i, i) = \sum_j w_{i, j}$
- Laplacian matrix L : $L = D - W$

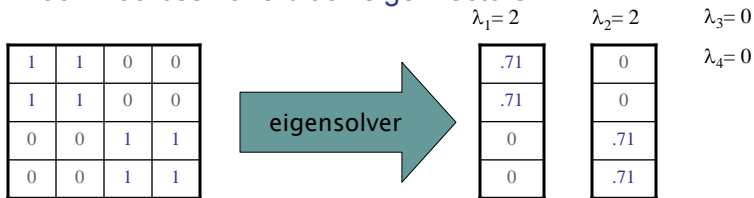
Eric Xing

20

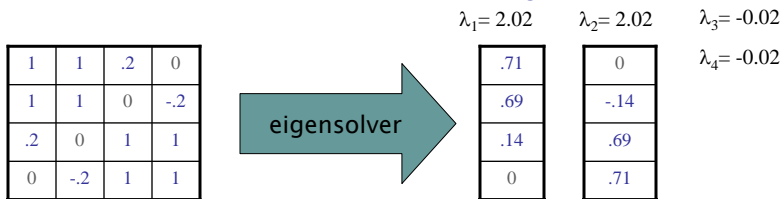
Eigenvectors and blocks



- Block matrices have block eigenvectors:



- Near-block matrices have near-block eigenvectors:



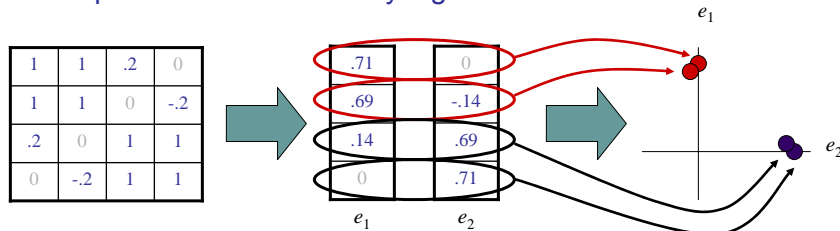
Eric Xing

21

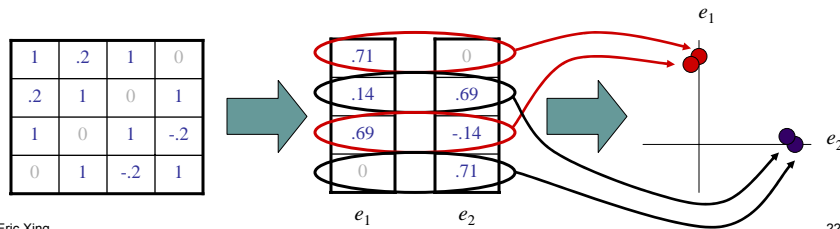
Spectral Space



- Can put items into blocks by eigenvectors:



- Clusters clear regardless of row ordering:



Eric Xing

22

Spectral Clustering



- Algorithms that cluster points using eigenvectors of matrices derived from the data
- Obtain data representation in the low-dimensional space that can be easily clustered
- Variety of methods that use the eigenvectors differently (we have seen an example)
- Empirically very successful
- Authors disagree:
 - Which eigenvectors to use
 - How to derive clusters from these eigenvectors
- Two general methods

Eric Xing

23

Method #1



- Partition using only one eigenvector at a time
- Use procedure recursively
- Example: Image Segmentation
 - Uses 2nd (smallest) eigenvector to define optimal cut
 - Recursively generates two clusters with each cut

Eric Xing

24

Method #2



- Use k eigenvectors (k chosen by user)
- Directly compute k -way partitioning
- Experimentally has been seen to be “better”

Spectral Clustering Algorithm

Ng, Jordan, and Weiss 2003

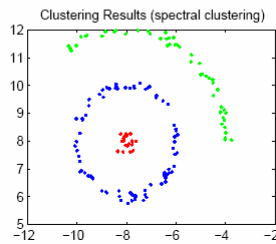
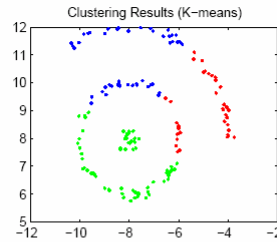
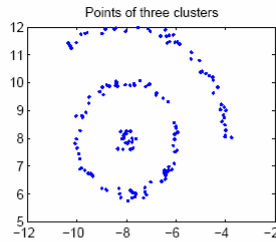


- Given a set of points $S = \{s_1, \dots, s_n\}$
- Form the affinity matrix $w_{i,j} = e^{-\frac{\|s_i - s_j\|_2^2}{\sigma^2}}$, $\forall i \neq j$, $w_{i,i} = 0$
- Define diagonal matrix $D_{ii} = \sum_k a_{ik}$
- Form the matrix $L = D^{-1/2} W D^{-1/2}$
- Stack the k largest eigenvectors of L to form the columns of the new matrix X :

$$X = \begin{bmatrix} | & | & \cdots & | \\ x_1 & x_2 & \cdots & x_k \\ | & | & & | \end{bmatrix}$$

- Renormalize each of X 's rows to have unit length and get new matrix Y . Cluster rows of Y as points in R^k

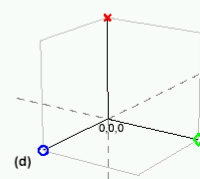
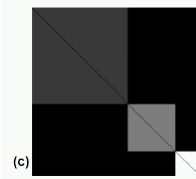
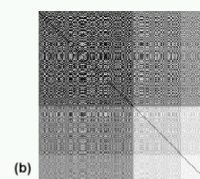
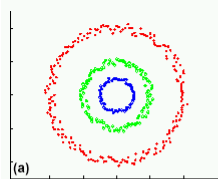
SC vs Kmeans



Eric Xing

27

Why it works?



- K-means in the spectrum space !

Eric Xing

28

More formally ...



- Recall generalized Ncut

$$\text{Ncut}(A_1, A_2 \dots A_k) = \sum_{r=1}^k \left(\frac{\sum_{i \in A_r, j \in V \setminus A_r} W_{ij}}{\sum_{i \in A_r, j \in V} W_{ij}} \right) = \sum_{r=1}^k \left(\frac{\text{cut}(A_r, \bar{A}_r)}{d_{A_r}} \right)$$

- Minimizing this is equivalent to spectral clustering

$$\min \text{Ncut}(A_1, A_2 \dots A_k) = \sum_{r=1}^k \left(\frac{\text{cut}(A_r, \bar{A}_r)}{d_{A_r}} \right)$$

$$\min Y^T D^{-1/2} W D^{-1/2} Y$$

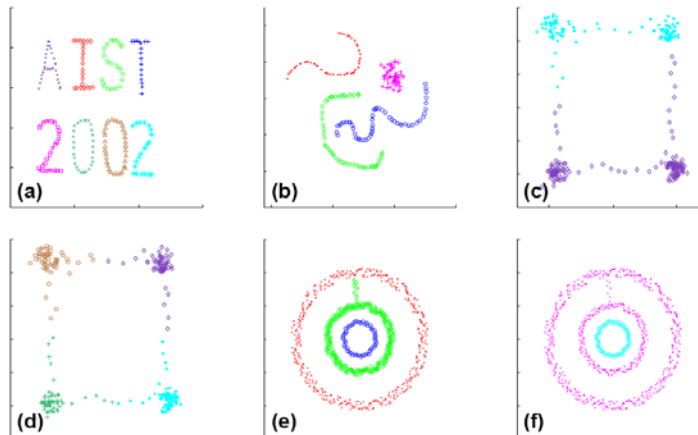
$$\text{s.t. } Y^T Y = I$$

$$Y = \begin{matrix} \text{segments} & \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ & \text{pixels} \end{matrix}$$

Eric Xing

29

Toy examples



Images from Matthew Brand (TR-2002-42)

Eric Xing

30

User's Prerogative



- Choice of k , the number of clusters
- Choice of scaling factor
 - Realistically, search over σ^2 and pick value that gives the tightest clusters
- Choice of clustering method: k -way or recursive bipartite
- Kernel affinity matrix

$$w_{i,j} = K(S_i, S_j)$$

Eric Xing

31

Conclusions



- Good news:
 - Simple and powerful methods to segment images.
 - Flexible and easy to apply to other clustering problems.
- Bad news:
 - High memory requirements (use sparse matrices).
 - Very dependant on the scale factor for a specific problem.

$$W(i, j) = e^{-\frac{\|X_{(i)} - X_{(j)}\|_2^2}{\sigma_x^2}}$$

Eric Xing

32