AIN311

Fundamentals of Machine Learning

Lecture 22:
K-Means Example Applications
Spectral clustering
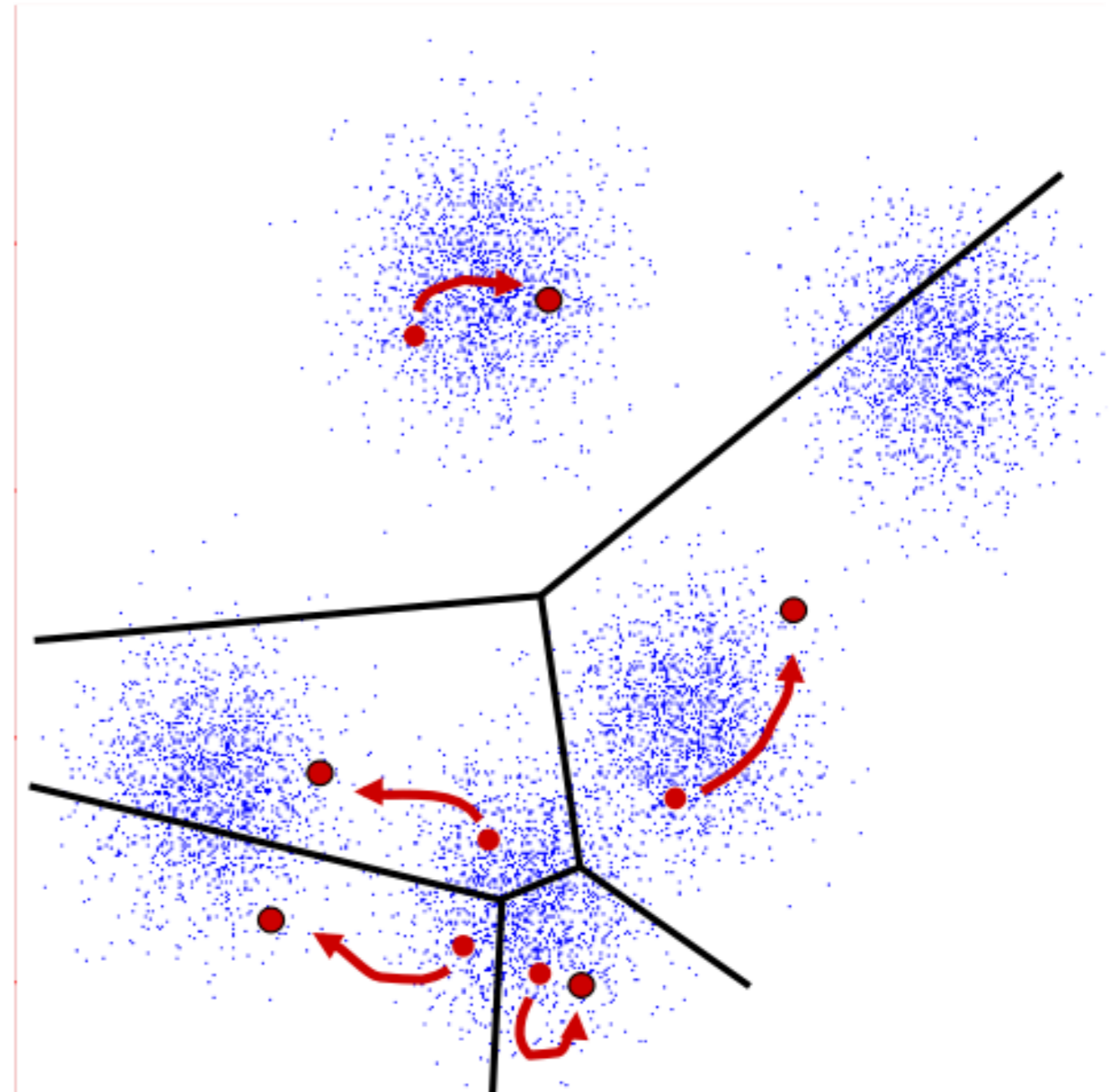Hierarchical clustering
What is a good clustering?

Photo by Unsplash user @foodiesfeed

Erkut Erdem // Hacettepe University // Fall 2023

# Last time... **K-Means**

- An iterative clustering algorithm
  - Initialize: Pick $K$ random points as cluster centers (means)
  - Alternate:
    - Assign data instances to closest mean
    - Assign each mean to the average of its assigned points
  - Stop when no points' assignments change

2

# Today

- K-Means Example Applications
- Spectral clustering
- Hierarchical clustering
- What is a good clustering?

# K-Means
# Example Applications

# Example: K-Means for Segmentation

K=2

$K = 2$ $K = 2$



**Goal of Segmentation is to partition an image into regions each of which has reasonably homogenous visual appearance.**

Original

5

# Example: K-Means for Segmentation
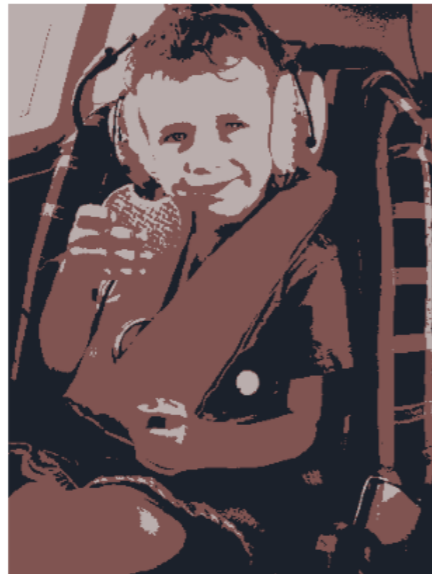
6

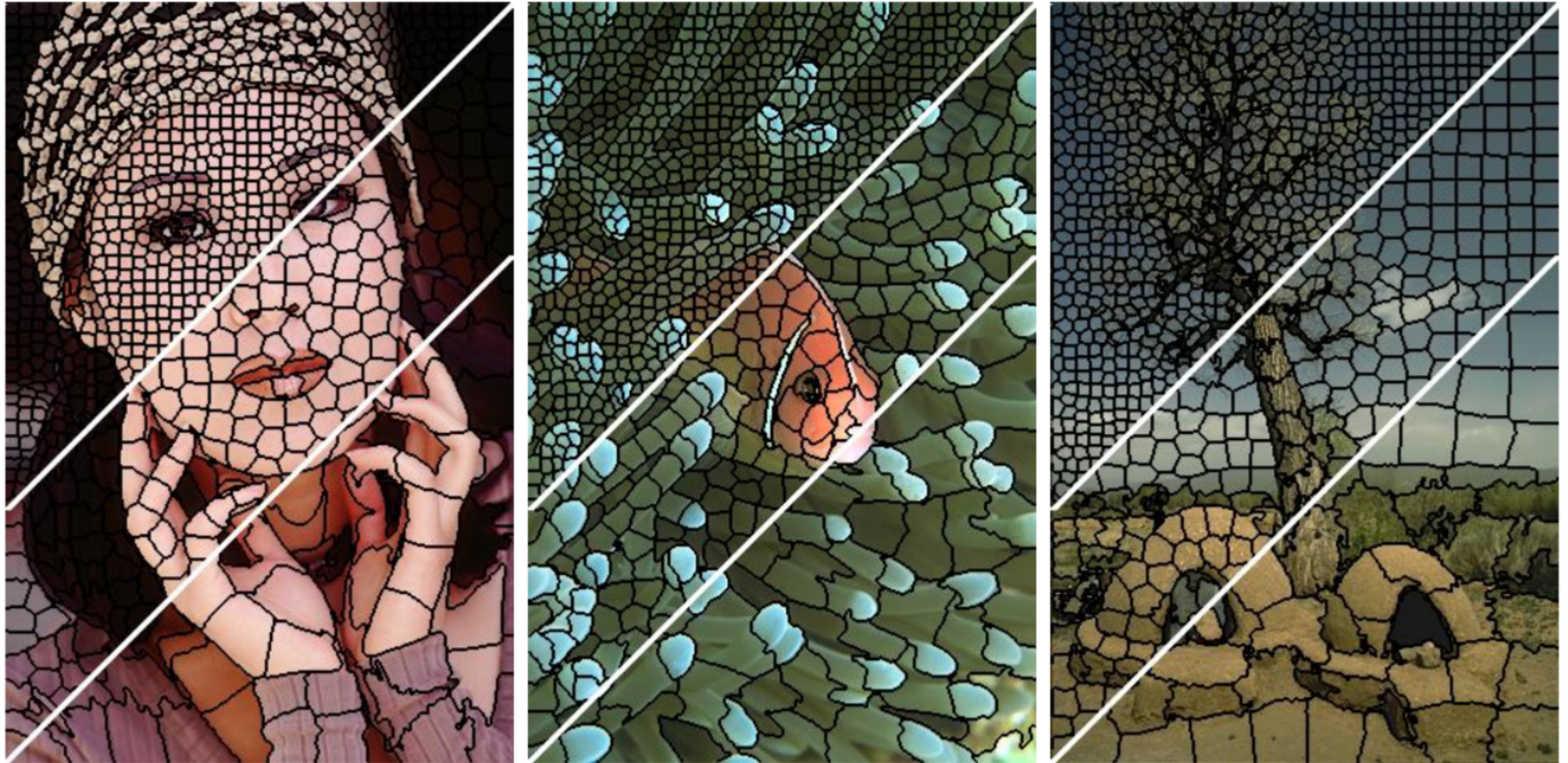# Example: K-Means for Segmentation



K=2

K=3

K=10

Original

# Example: Vector quantization



**FIGURE 14.9.** *Sir Ronald A. Fisher (1890 − 1962) was one of the founders of modern day statistics, to whom we owe maximum-likelihood, sufficiency, and many other fundamental concepts. The image on the left is a 1024×1024 grayscale image at 8 bits per pixel. The center image is the result of 2 × 2 block VQ, using 200 code vectors, with a compression rate of 1.9 bits/pixel. The right image uses only four code vectors, with a compression rate of 0.50 bits/pixel*

[Figure from Hastie *et al.* book]

# Example: Simple Linear Iterative Clustering (SLIC) superpixels



$$\Psi(x,y) = \begin{bmatrix} \lambda x \\ \lambda y \\ I(x,y) \end{bmatrix}$$

$\lambda$: spatial regularization parameter

R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Susstrunk SLIC Superpixels Compared to State-of-the-art Superpixel Methods, IEEE T-PAMI, 2012

# Bag of Words model

11

# Object → Bag of 'words'

slide by Fei Fei Li
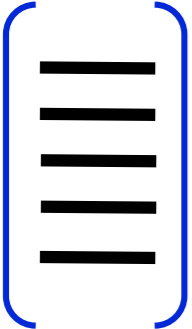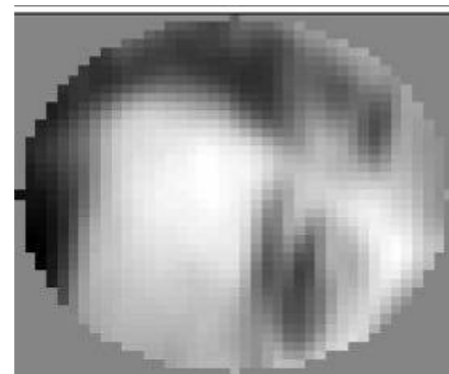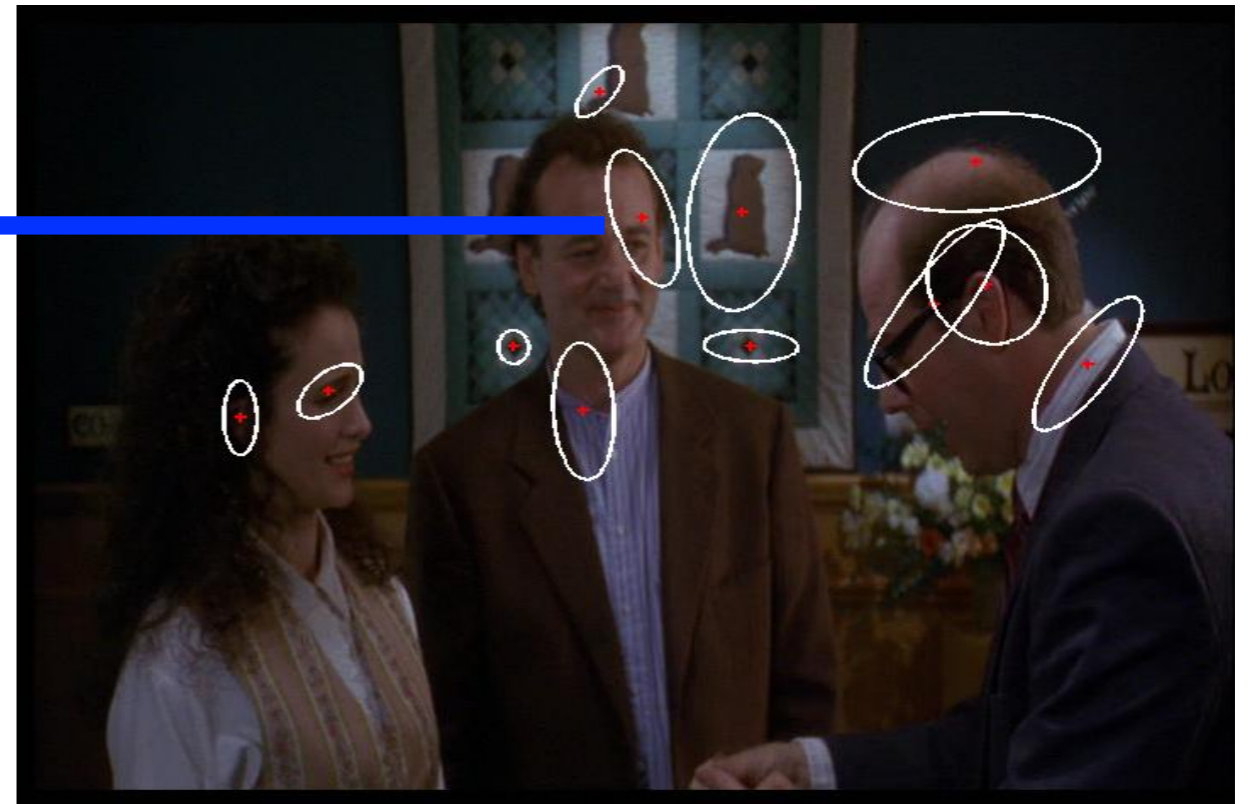
# Interest Point Features



**Compute SIFT descriptor**

[Lowe'99]

**Normalize patch**
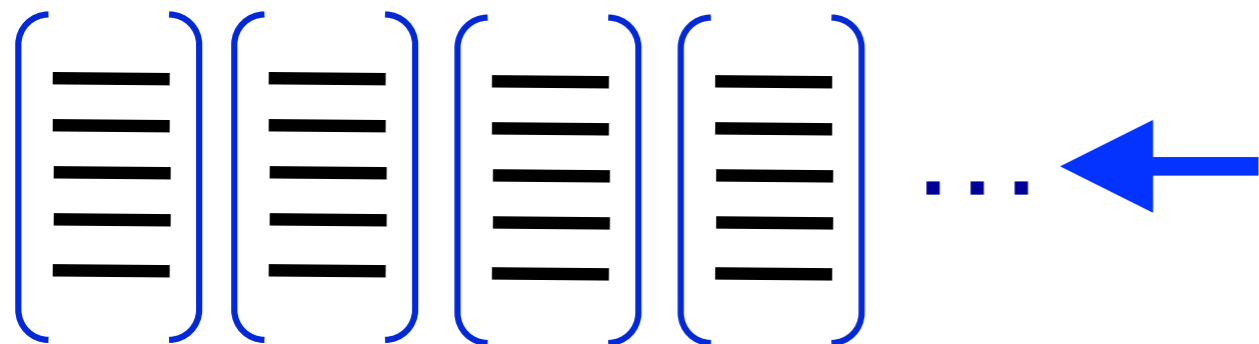
Detect patches

[Mikojaczyk and Schmid '02]

[Matas et al. '02]
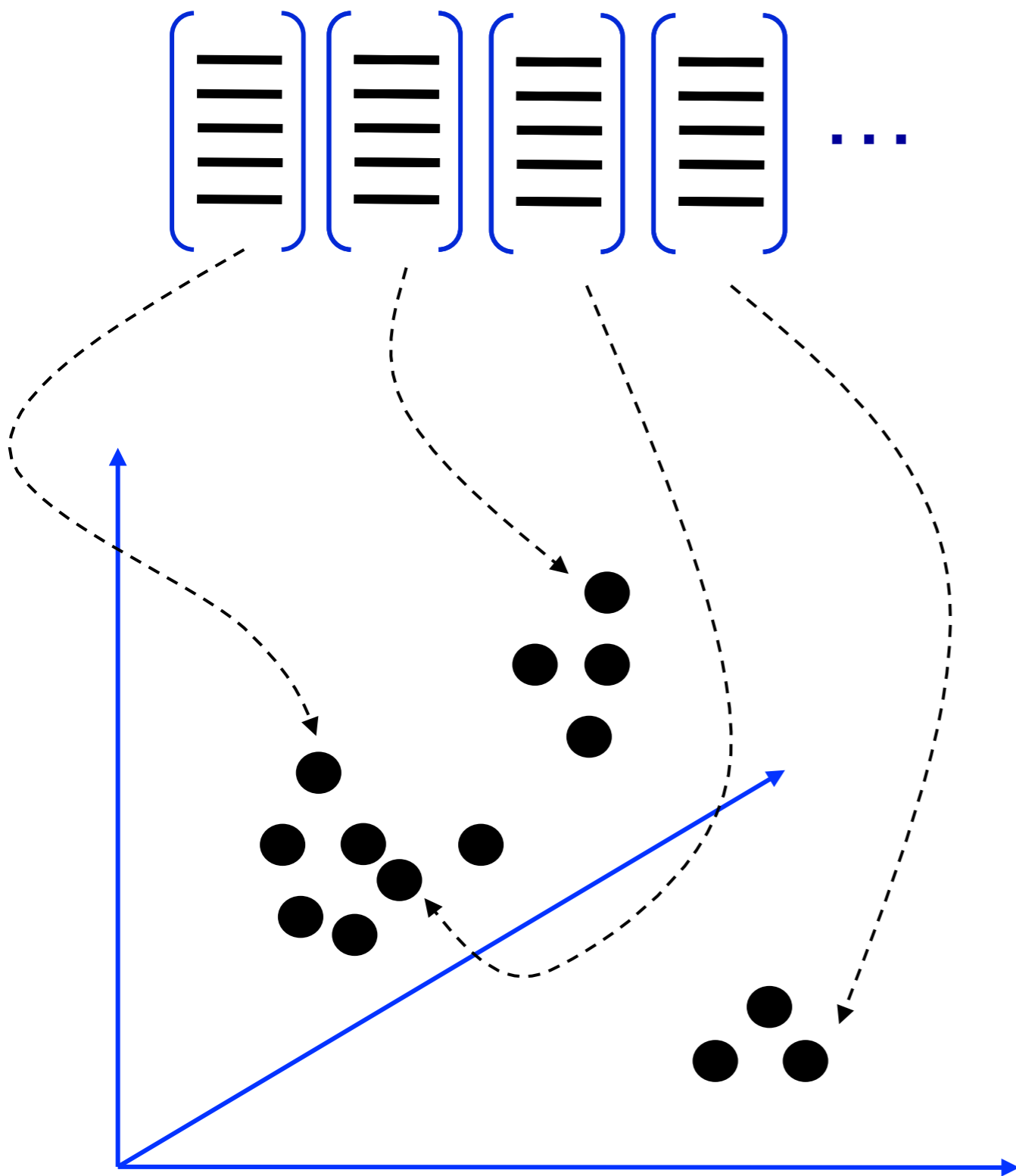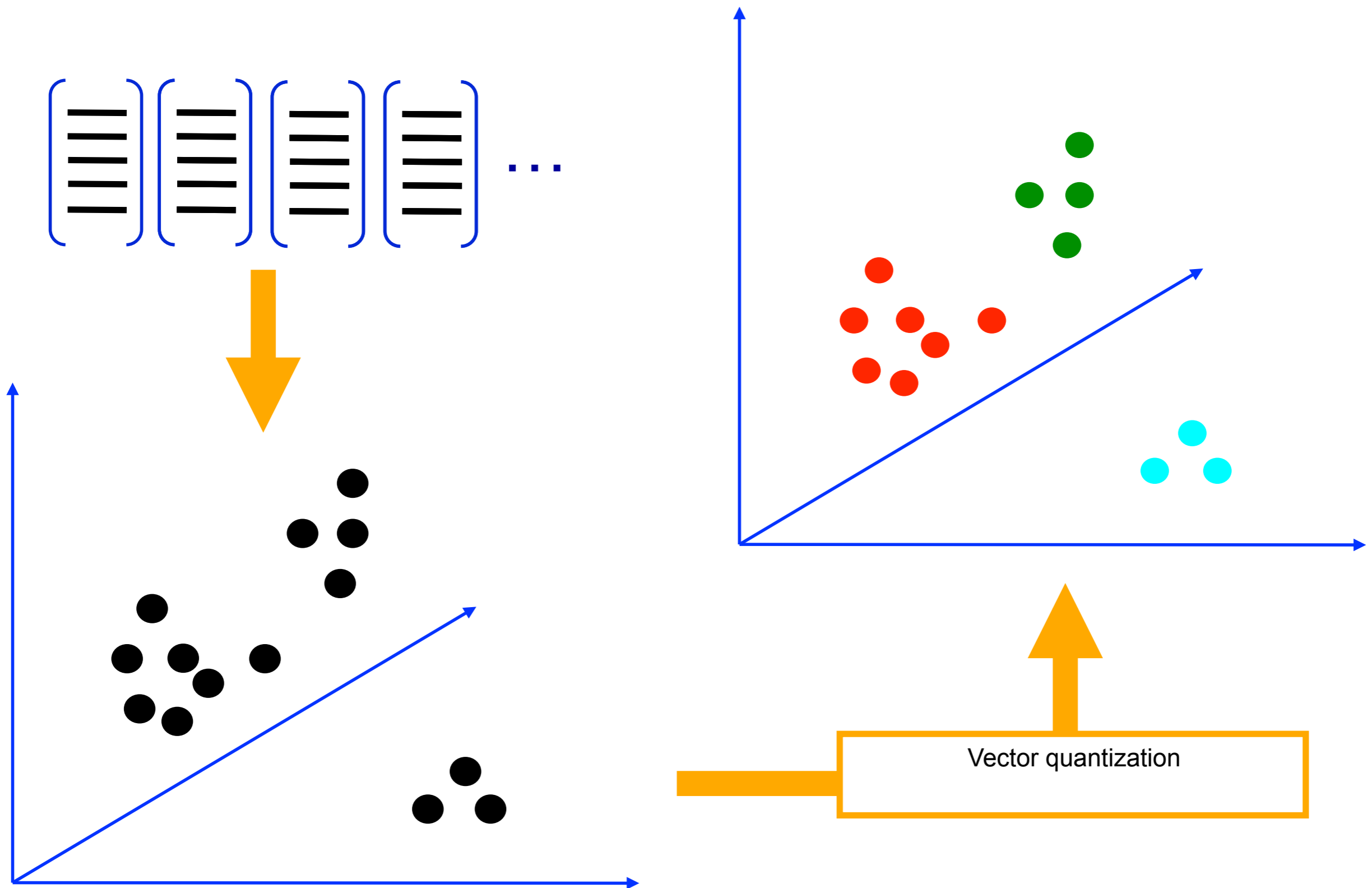
[Sivic et al. '03]

slide by Josef Sivic

13

# Patch Features

14

# Dictionary Formation

slide by Josef Sivic

# Clustering (usually K-means)



Vector quantization
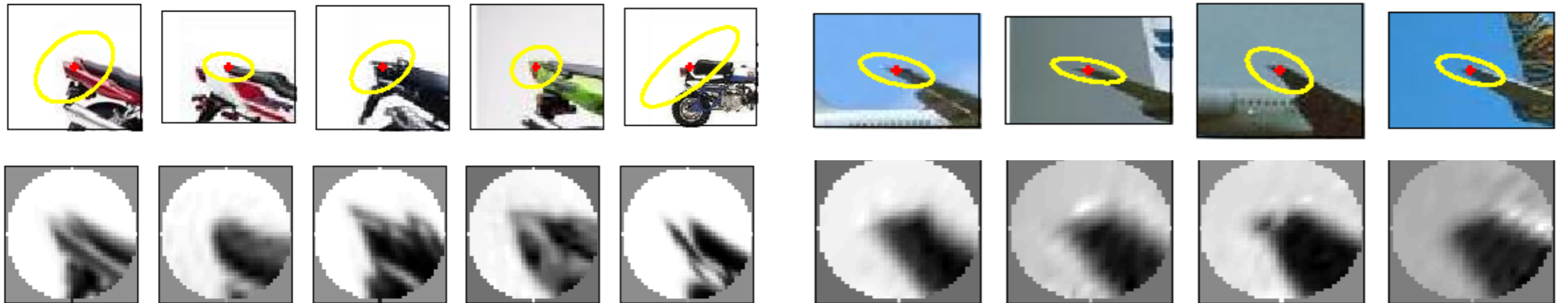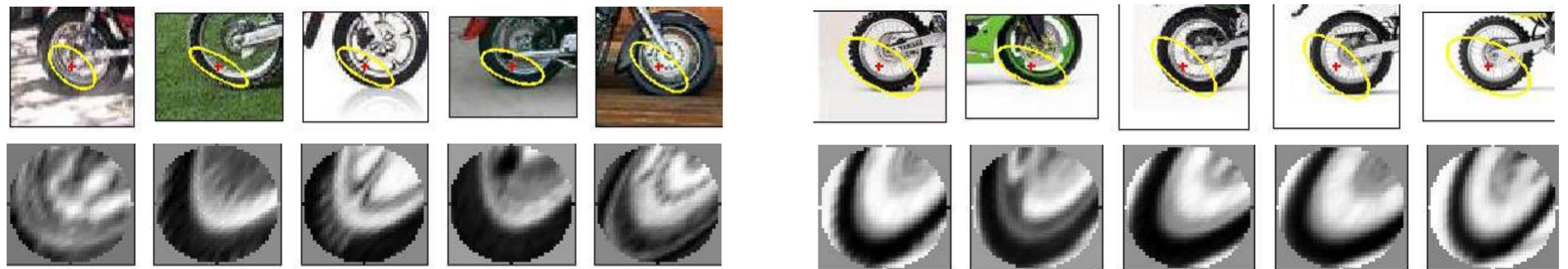
16

# Clustered Image Patches

# Visual synonyms and polysemy



Visual Polysemy. Single visual word occurring on different (but locally similar) parts on different object categories.



Visual Synonyms. Two different visual words representing a similar part of an object (wheel of a motorbike).

18

# Image Representation



frequency

codewords

# Spectral clustering

# Graph-Theoretic Clustering

**Goal:** Given data points $X_1, \ldots, X_n$ and similarities $W(X_i, X_j)$,

partition ... up
are si... ar.

Simila...

Data      Similarities

Data      Similarities      Similarity graph

Part
edg

# Graphs Representations



$$
\begin{array}{c c c c c c}
 & a & b & c & d & e \\
a & \begin{bmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 \end{bmatrix}
\end{array}
$$

Adjacency Matrix

# A Weighted Graph and its Representation

Affinity Matrix

$$W = \begin{bmatrix} 1 & .1 & .3 & 0 & 0 \\ .1 & 1 & .4 & 0 & .2 \\ .3 & .4 & 1 & .6 & .7 \\ 0 & 0 & .6 & 1 & 1 \\ 0 & .2 & .7 & 1 & 1 \end{bmatrix}$$

$W_{ij}$ : probability that i & j belong to the same cluster

# Similarity graph construction

- Similarity Graphs: Model local neighborhood relations between data points

- E.g. epsilon-NN

$$W_{ij} = \begin{cases} 1 & \|x_i - x_j\| \le \epsilon \\ 0 & \text{otherwise} \end{cases}$$

Controls size of neighborhood

or mutual k-NN graph ($W_{ij} = 1$ if $x_i$ or $x_j$ is k nearest neighbor of the other)



$G = \{V, E\}$

# Similarity graph construction

- Similarity Graphs: Model local neighborhood relations between data points

- E.g. Gaussian kernel similarity function

$$W_{ij} = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

$\epsilon\epsilon$

→ Controls size of neighborhood



$G = \{V, E\}$

# Scale affects affinity

- Small $\sigma$: group only nearby points
- Large $\sigma$: group far-away points

# Feature grouping by "relocalisation" of eigenvectors of the proximity matrix

Guy L. Scott
Robotics Research Group
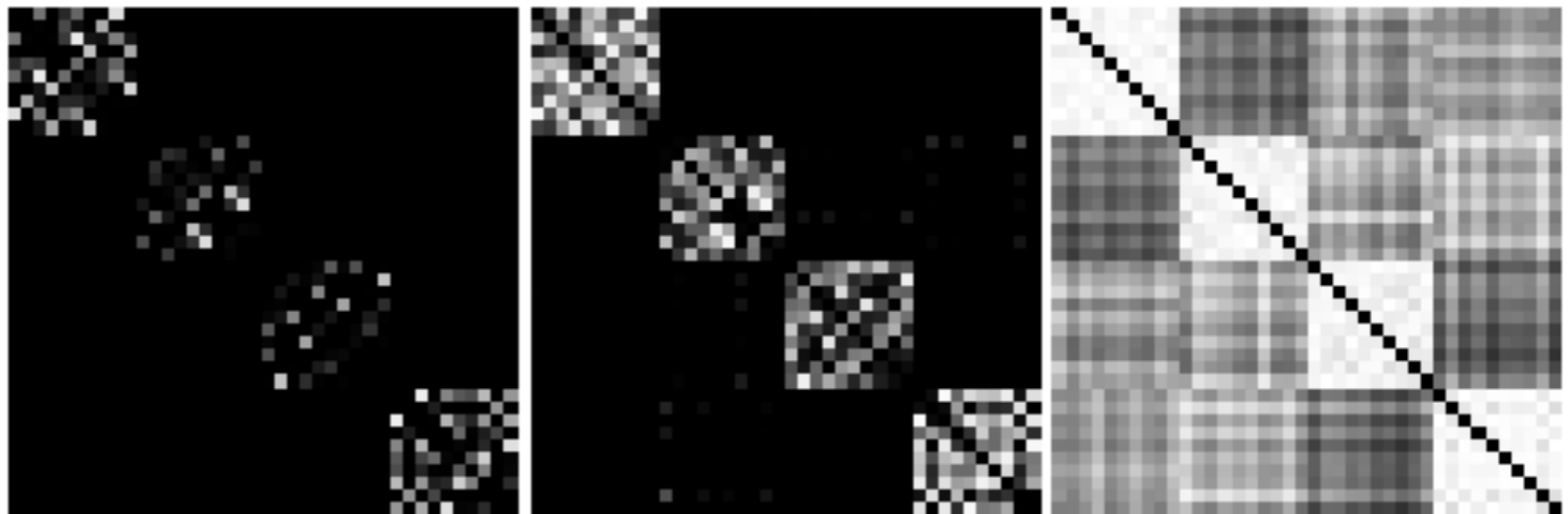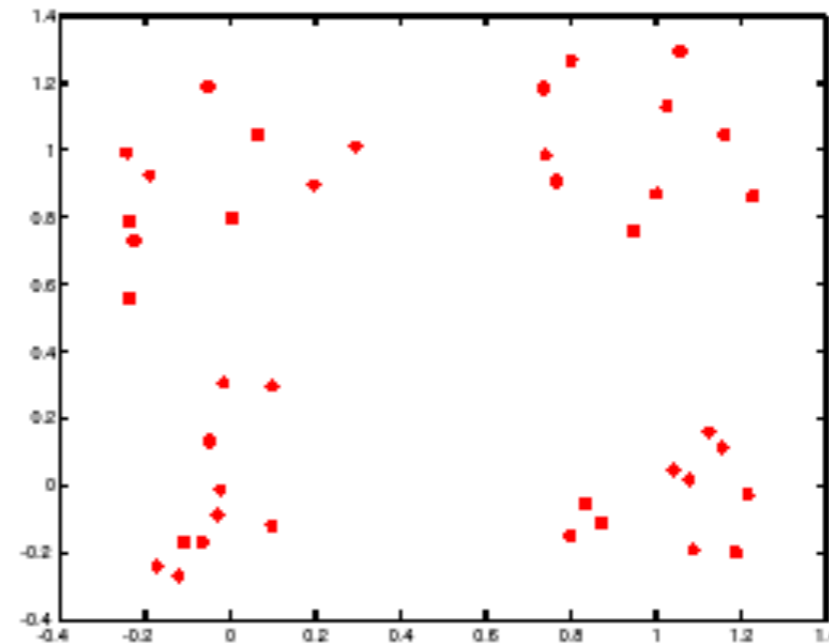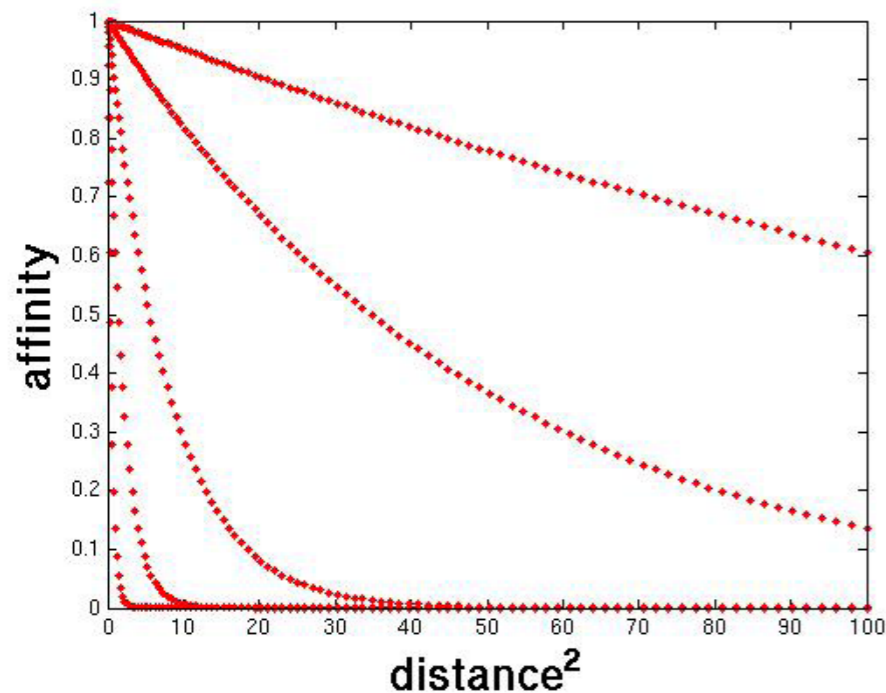Department of Engineering Science
University of Oxford

H. Christopher Longuet-Higgins
University of Sussex
Falmer
Brighton

$$W_{ij} = \exp(-\| z_i - z_j \|^2 / s^2)$$

With an appropriate s

W=

|   | A | B | C |
|---|---|---|---|
| A | 1.00 | 0.63 | 0.03 |
| B | 0.63 | 1.00 | 0.0 |
| C | 0.03 | 0.0 | 1.00 |

The eigenvectors of W are:

|             | $E_1$ | $E_2$ | $E_3$ |
|-------------|-------|-------|-------|
| Eigenvalues | 1.63  | 1.00  | 0.37  |
| A           | -0.71 | -0.01 | 0.71  |
| B           | -0.71 | -0.05 | -0.71 |
| C           | -0.04 | 1.00  | -0.03 |

The first 2 eigenvectors group the points as desired…

Three points in feature space

# Example eigenvector

points

eigenvector

Affinity matrix

28

# Example eigenvector



points

eigenvector

Affinity matrix

29

# Graph cut



- Set of edges whose removal makes a graph disconnected
- Cost of a cut: sum of weights of cut edges
- A graph cut gives us a partition (clustering)
  - What is a "good" graph cut and how do we find one?

# Minimum cut

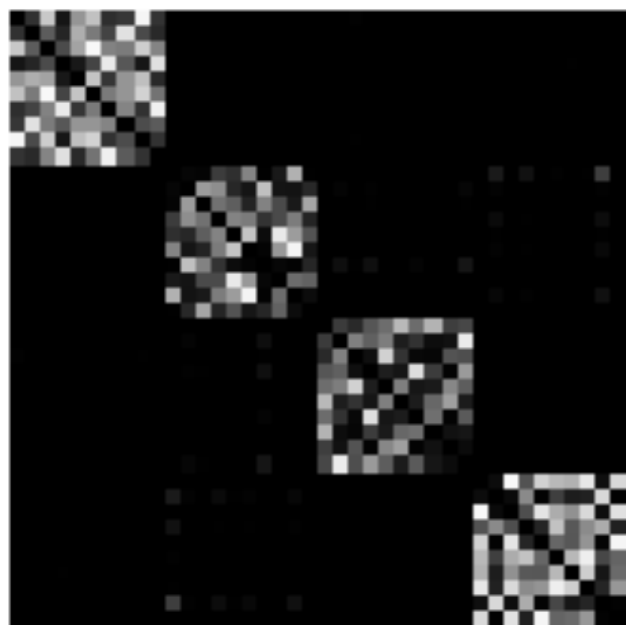- A cut of a graph *G* is the set of edges *S* such that removal of *S* from *G* disconnects *G*.

**Cut**: sum of the weight of the cut edges:



$$cut(\mathrm{A},\mathrm{B}) = \sum_{u\in\mathrm{A}, v\in\mathrm{B}} \mathrm{W}(u,v),$$

with $\mathrm{A} \cap \mathrm{B} = \varnothing$

# Minimum cut

- We can do clustering by finding the minimum cut in a graph
  - Efficient algorithms exist for doing this

Minimum cut example

# Minimum cut

- We can do segmentation by finding the <span style="color:green">minimum cut</span> in a graph
  - Efficient algorithms exist for doing this

Minimum cut example

33

# Drawbacks of Minimum cut

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



Ideal Cut

Cuts with lesser weight than the ideal cut

# Normalized cuts

Write graph as V, one cluster as A and the other as B



$$\text{Ncut}(A,B) = \frac{\text{cut}(A,B)}{\text{assoc}(A,V)} + \frac{\text{cut}(A,B)}{\text{assoc}(B,V)}$$

cut(A,B) is sum of weights with one end in A and one end in B

$$cut(\mathrm{A,B}) = \sum_{u \in \mathrm{A}, v \in \mathrm{B}} \mathrm{W}(u,v),$$
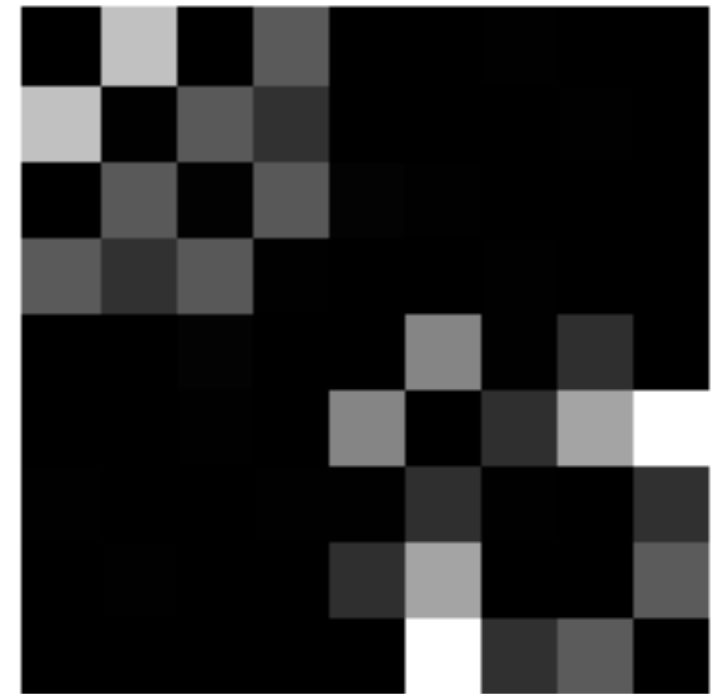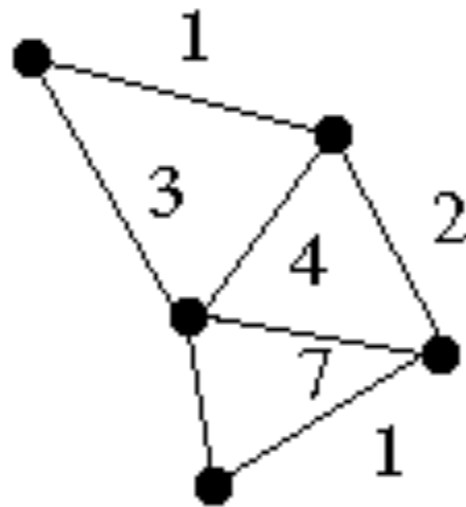
with $\mathrm{A} \cap \mathrm{B} = \varnothing$

assoc(A,V) is sum of all edges with one end in A.

$$assoc(\mathrm{A,B}) = \sum_{u \in \mathrm{A}, v \in \mathrm{B}} \mathrm{W}(u,v)$$

A and B not necessarily disjoint

J. Shi and J. Malik. Normalized cuts and image segmentation. PAMI 2000

# Normalized cut

- Let *W* be the adjacency matrix of the graph
- Let *D* be the diagonal matrix with diagonal entries
  $D(i, i) = \Sigma_j W(i, j)$
- Then the normalized cut cost can be written as

$$\frac{y^T (D - W) y}{y^T D y}$$

*D-W*: Graph Laplacian

where *y* is an indicator vector whose value should be 1 in the *i*-th position if the *i*-th feature point belongs to A and a negative constant otherwise
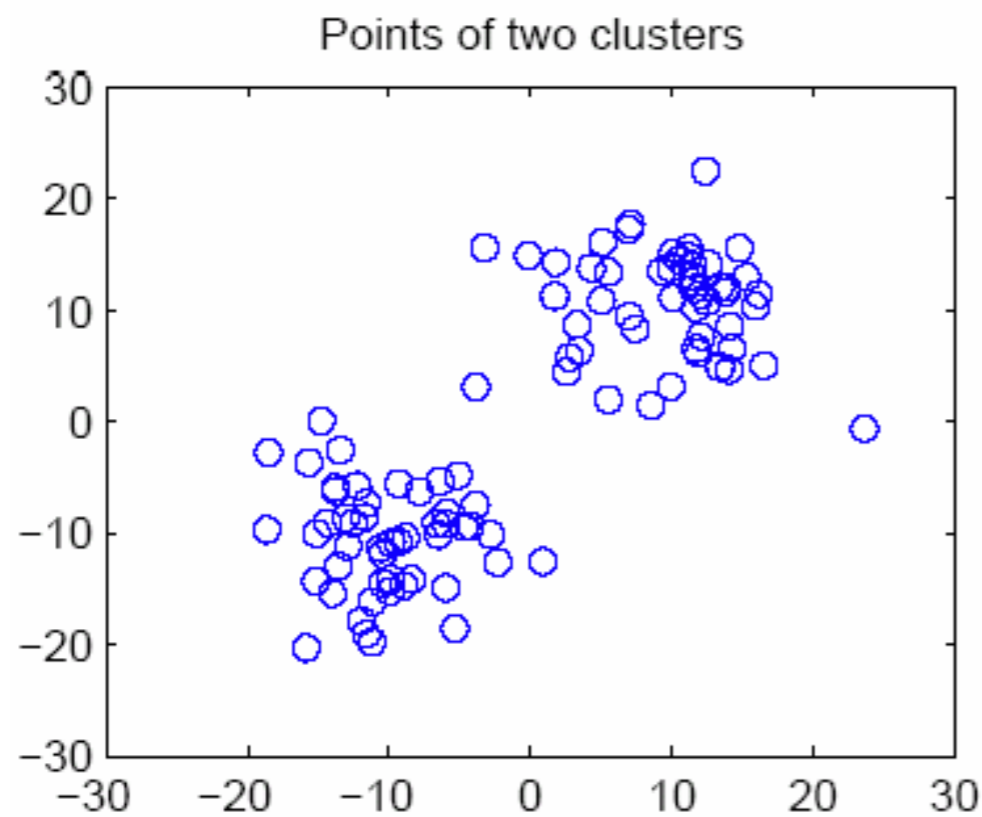
# Normalized cut

- Finding the exact minimum of the normalized cut cost is NP-complete, but if we *relax y* to take on arbitrary values, then we can minimize the relaxed cost by solving the **generalized eigenvalue problem** $(D − W)y = \lambda Dy$

- The solution *y* is given by the generalized eigenvector corresponding to the second smallest eigenvalue, aka the <u>Fiedler vector</u>

- Intuitively, the *i*-th entry of *y* can be viewed as a "soft" indication of the component membership of the *i*-th feature

  - Can use 0 or median value of the entries as the splitting point (threshold), or find threshold that minimizes the Ncut cost

J. Shi and J. Malik. <u>Normalized cuts and image segmentation.</u> PAMI 2000
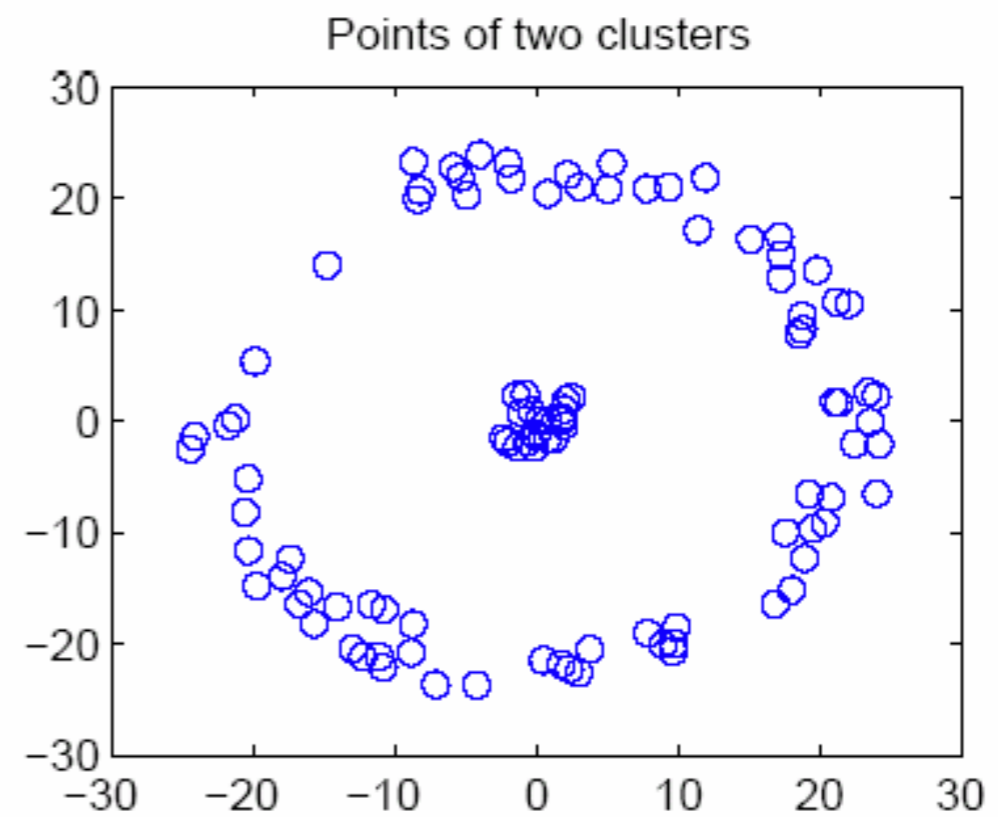
# Normalized cut algorithm

1. Given an image or image sequence, set up a weighted graph $\mathbf{G} = (\mathbf{V}, \mathbf{E})$, and set the weight on the edge connecting two nodes being a measure of the similarity between the two nodes.

2. Solve $(\mathbf{D} - \mathbf{W})\boldsymbol{x} = \lambda \mathbf{D} \boldsymbol{x}$ for eigenvectors with the smallest eigenvalues.

3. Use the eigenvector with second smallest eigenvalue to bipartition the graph.

4. Decide if the current partition should be sub-divided, and recursively repartition the segmented parts if necessary.

slide by Bill Freeman and Antonio Torralba

J. Shi and J. Malik. Normalized cuts and image segmentation. PAMI 2000    38

# K-Means vs. Spectral Clustering

- Applying k-means to Laplacian eigenvectors allows us to find cluster with non-convex boundaries.
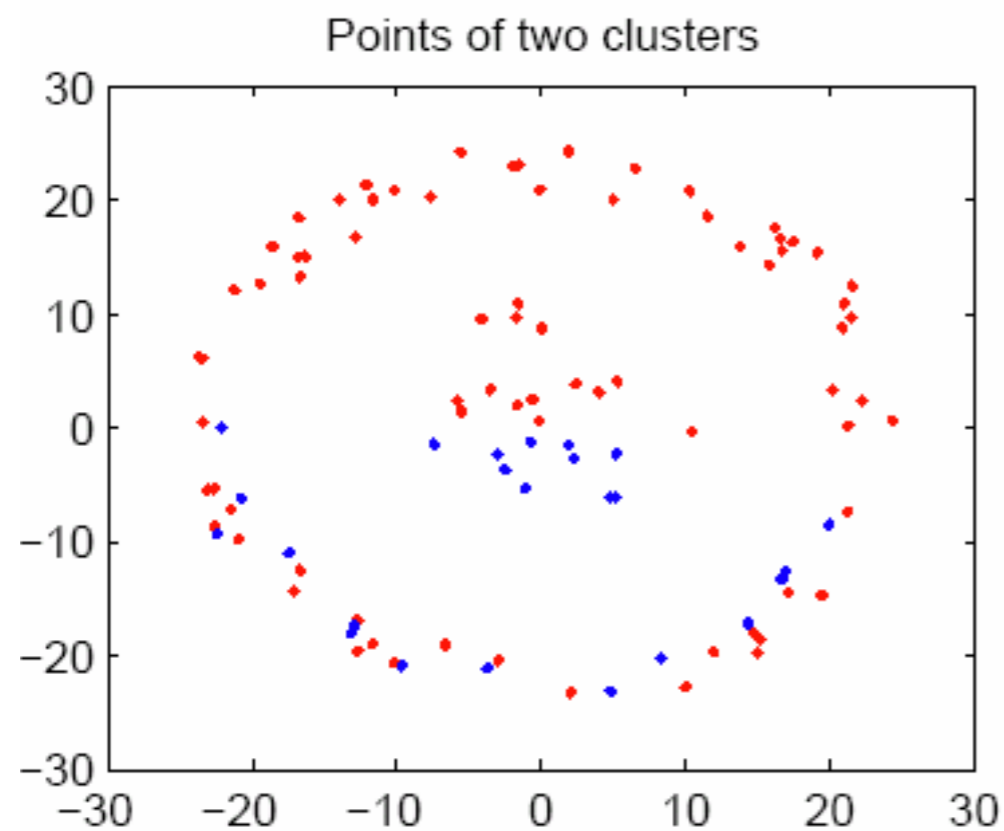


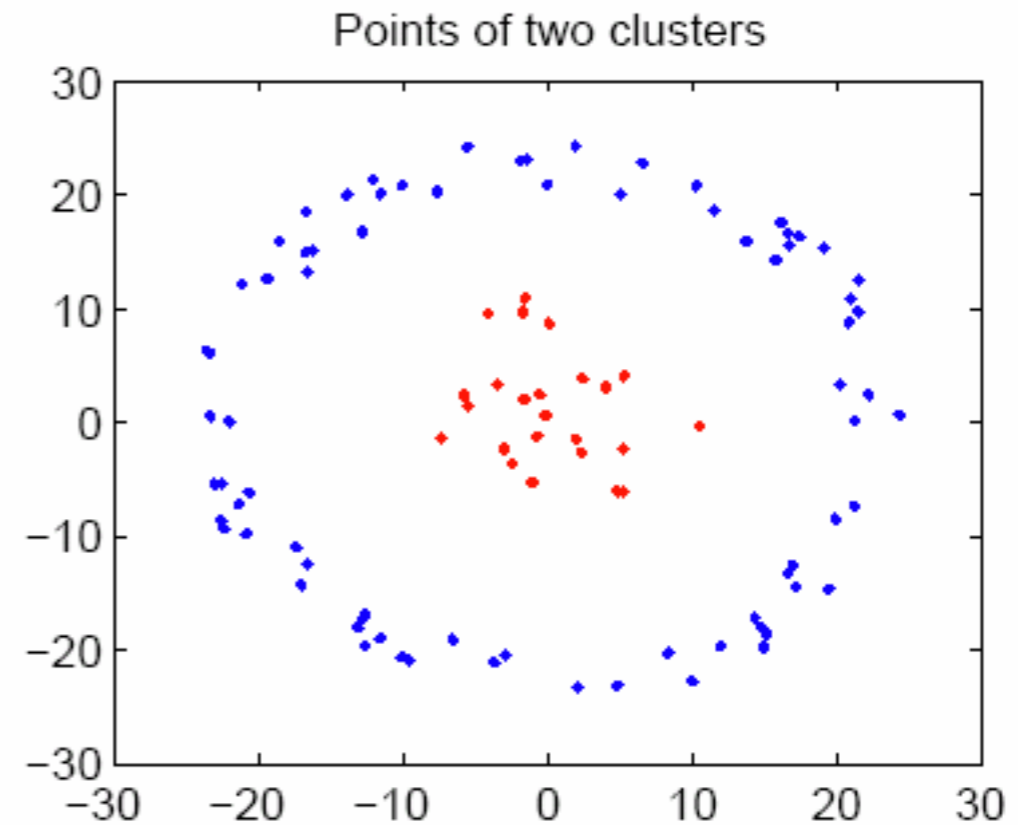Both perform same                    Spectral clustering is superior

# K-Means vs. Spectral Clustering

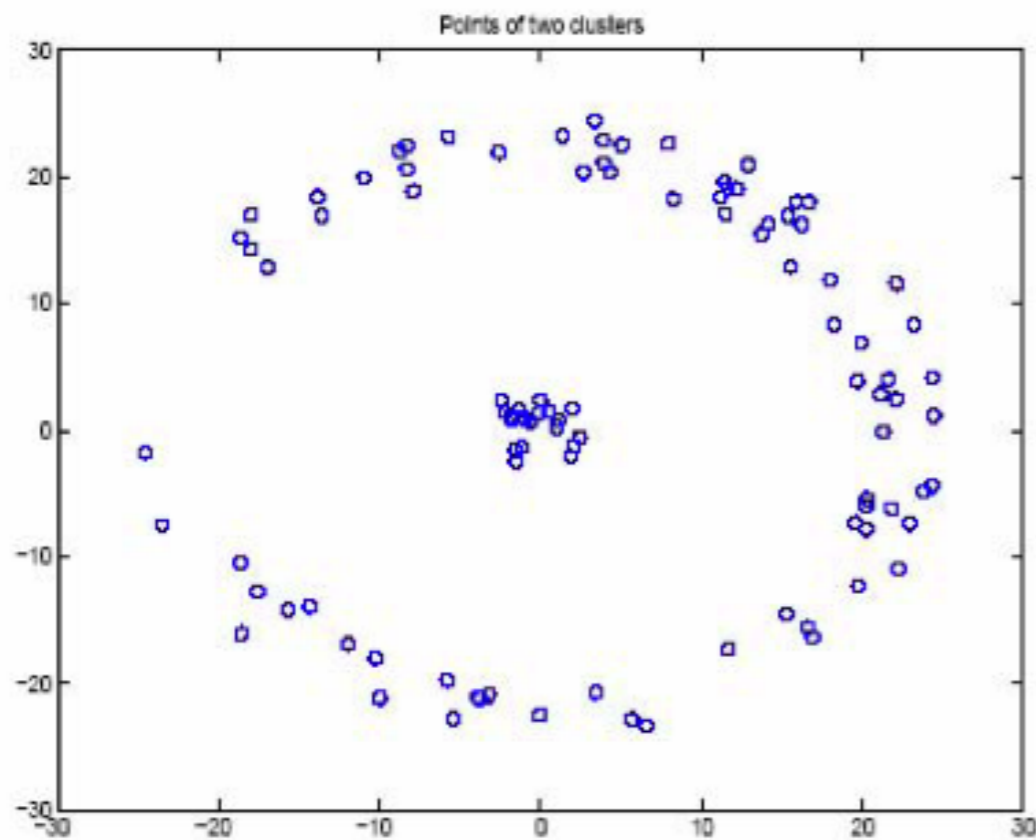- Applying k-means to Laplacian eigenvectors allows us to find cluster with non-convex boundaries.
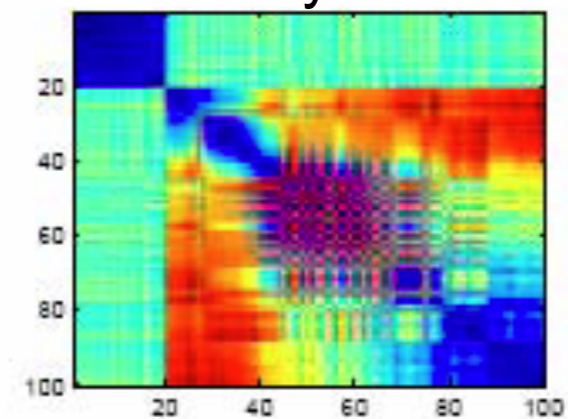


k-means output

Spectral clustering output
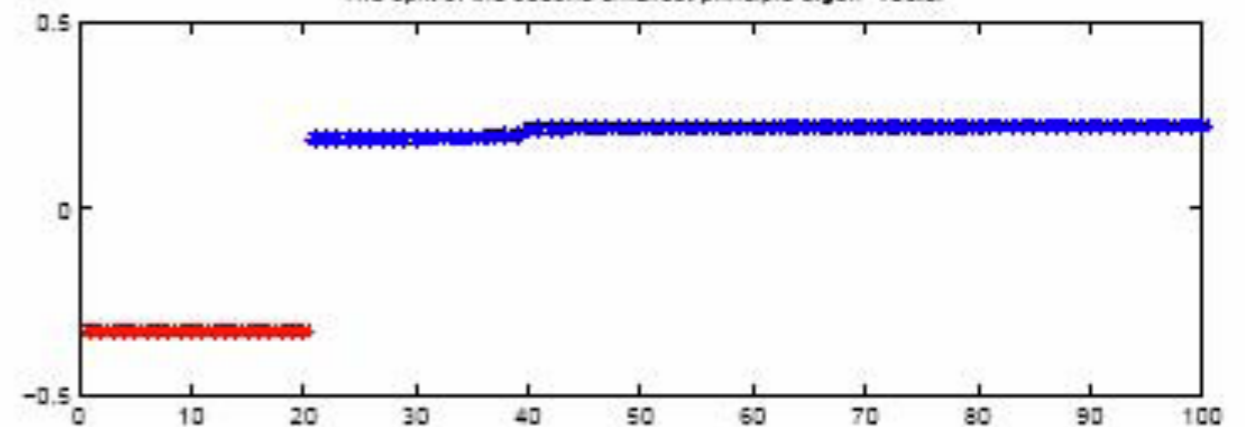
40

# K-Means vs. Spectral Clustering

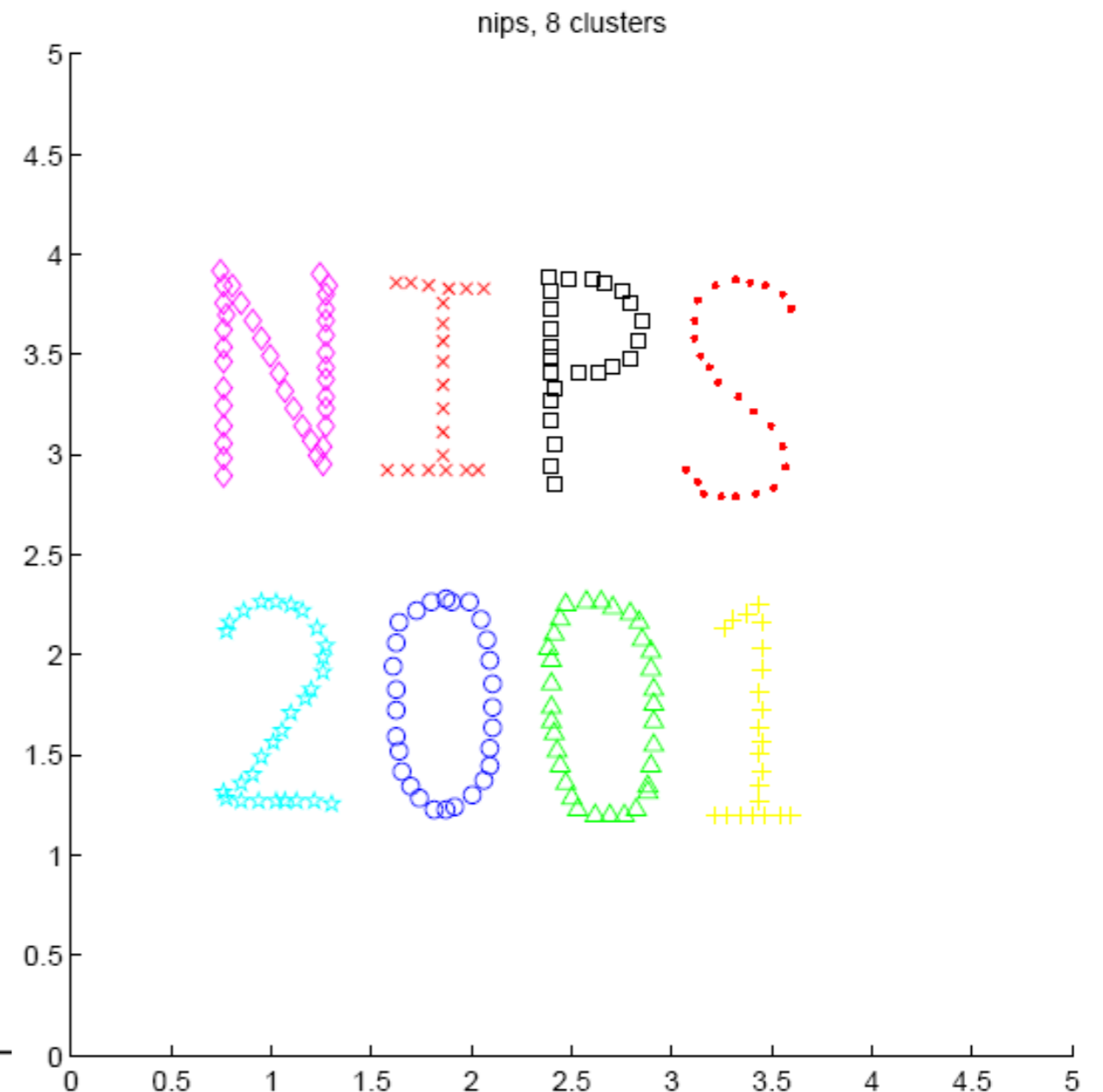- Applying k-means to Laplacian eigenvectors allows us to find cluster with non-convex boundaries.
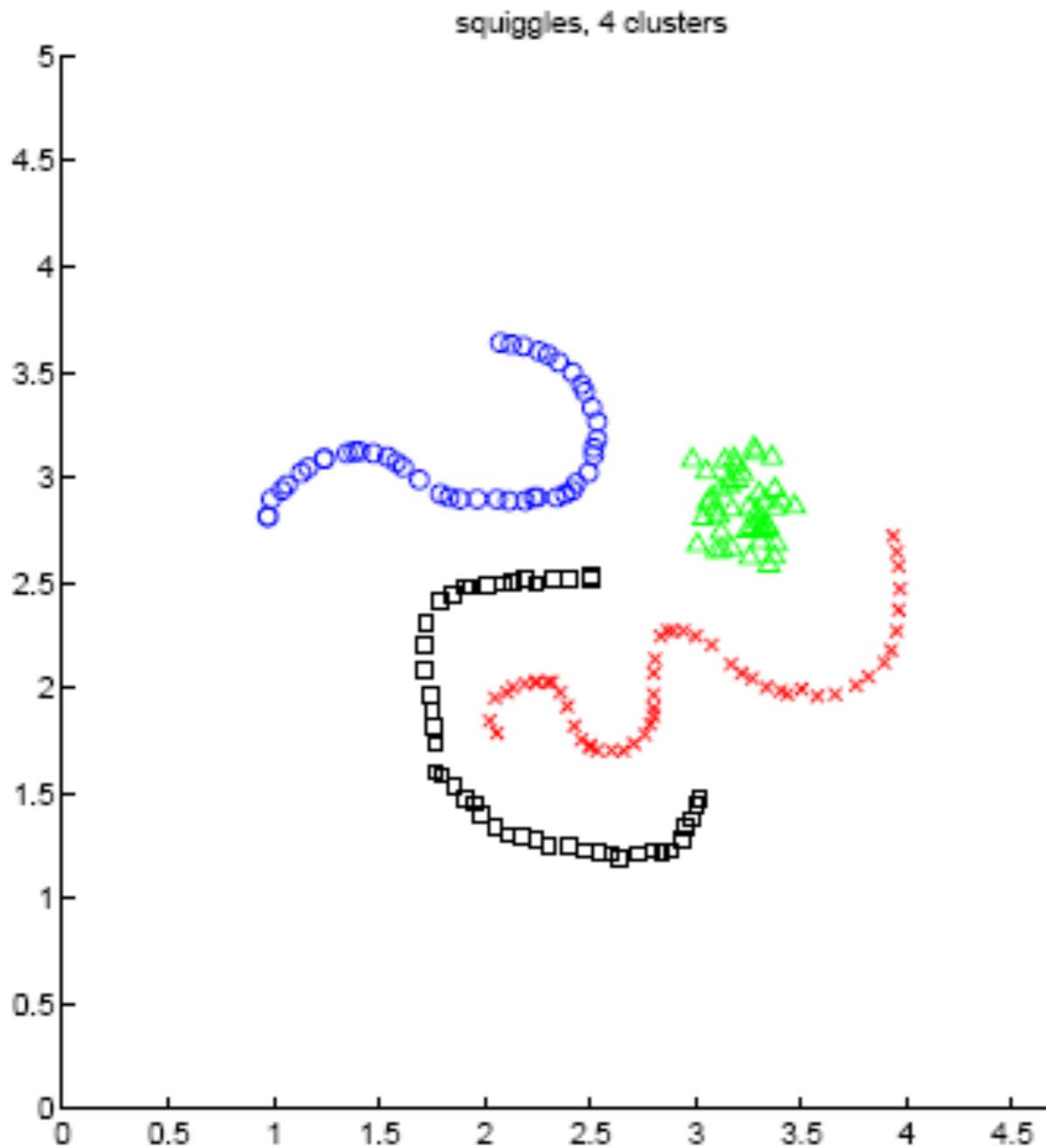


Points of two clusters

Second eigenvector of graph Laplacian

The split of the second smallest principle eigen−vector

41

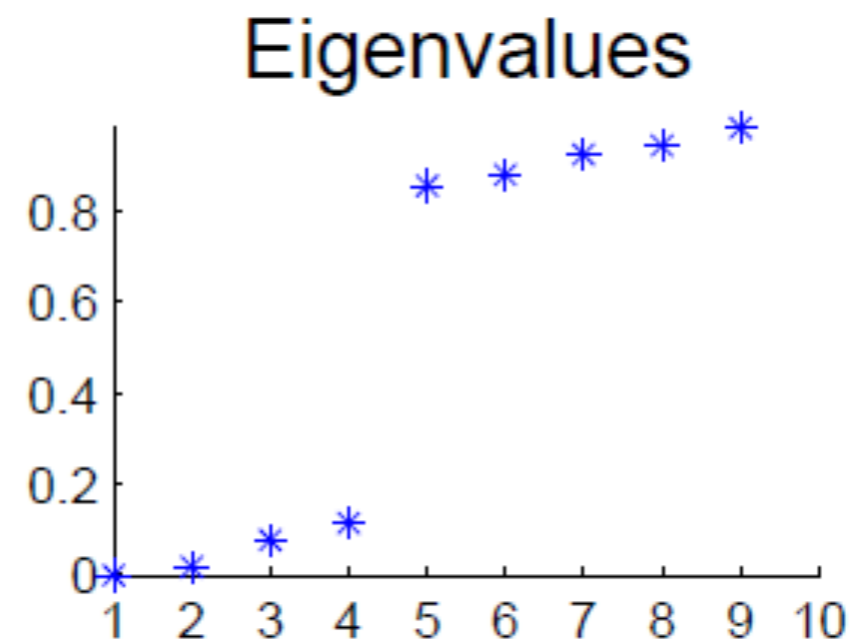# Examples



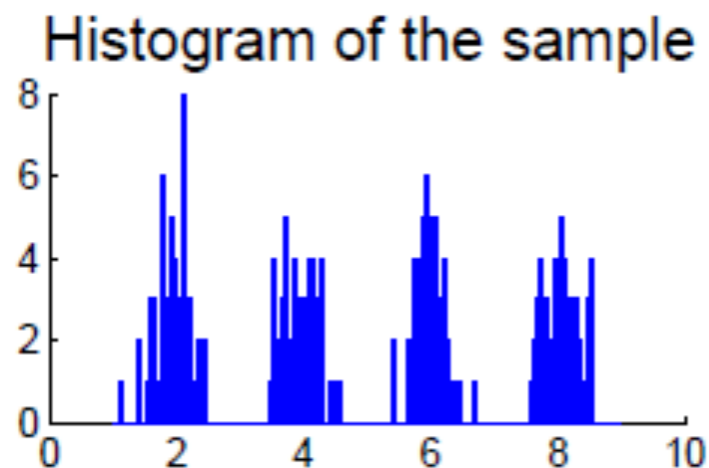squiggles, 4 clusters

nips, 8 clusters

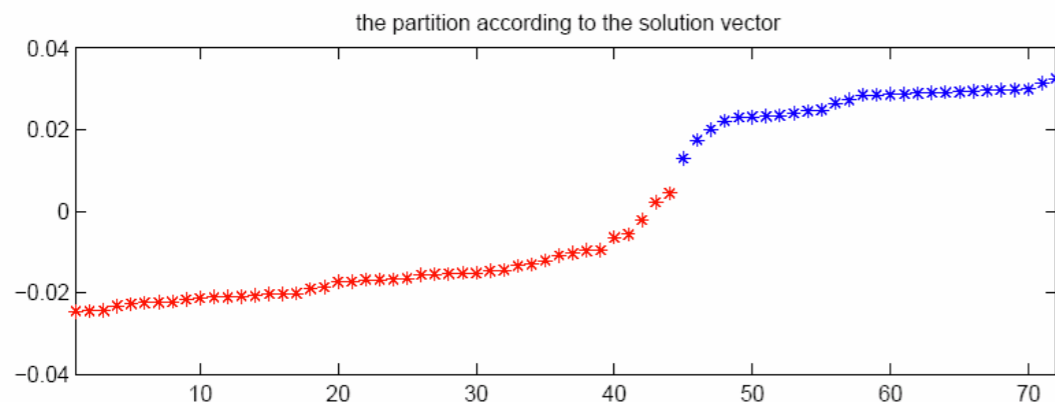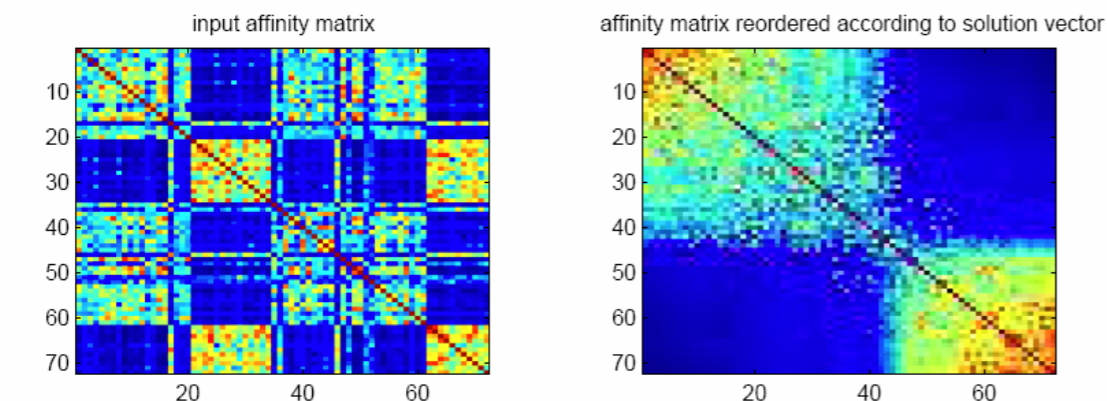[Ng et al., 2001]

42

# Some Issues

- Choice of number of clusters k
  - Most stable clustering is usually given by the value of k that maximizes the eigengap (difference between consecutive eigenvalues)

$$\Delta_k = \left\| \lambda_k - \lambda_{k-1} \right\|$$
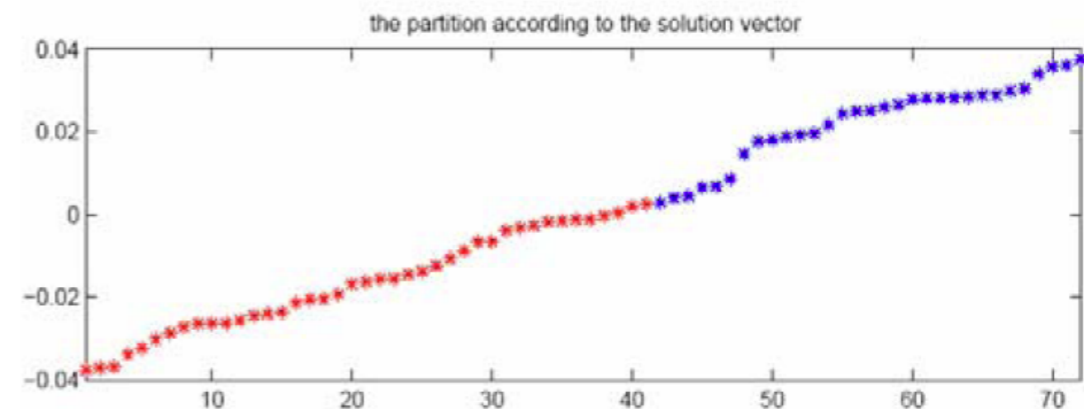


Histogram of the sample

Eigenvalues

43

# Some Issues

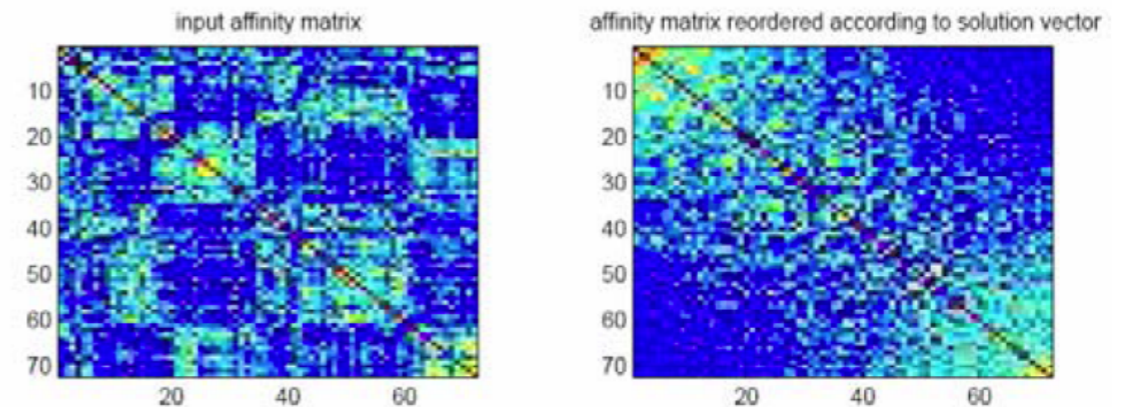- Choice of number of clusters k
- Choice of similarity
  - Choice of kernel
    for Gaussian kernels, choice of $\sigma$



Good similarity measure

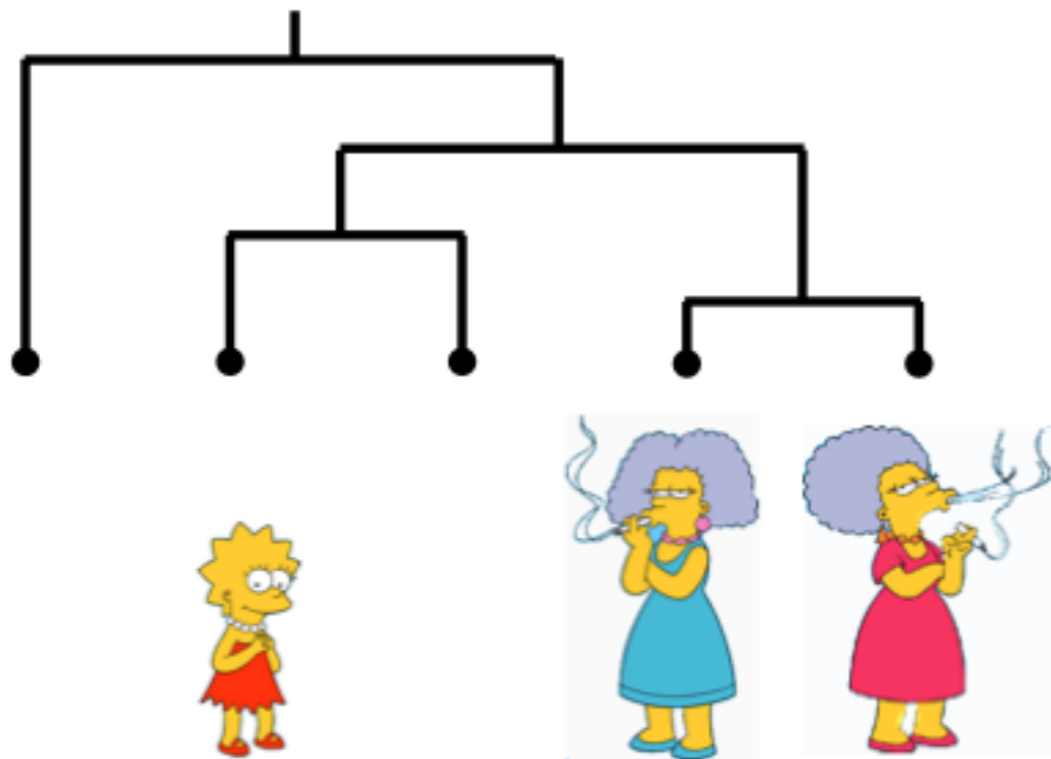Poor similarity measure

# Some Issues

- Choice of number of clusters k

- Choice of similarity
  - Choice of kernel
    for Gaussian kernels, choice of $\sigma$

- Choice of clustering method
  - k-way vs. recursive 2-way

# Hierarchical clustering

# Hierarchical Clustering

• Bottom-Up (agglomerative): Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.



• The number of dendrograms with

$n$ leafs = (2n -3)!/[(2(n -2)) (n -2)!]

| Number of leafs | Number of possible Dendrongrams |
|---|---|
| 2 | 1 |
| 3 | 3 |
| 4 | 15 |
| 5 | 105 |
| … | … |
| 10 | 34,459,425 |

We begin with a distance matrix which contains the distances between every pair of objects in our dataset

$$D( \text{ } , \text{ } ) = 8$$

$$D( \text{ } , \text{ } ) = 1$$

| 0 | 8 | 8 | 7 | 7 |
|---|---|---|---|---|
|   | 0 | 2 | 4 | 4 |
|   |   | 0 | 3 | 3 |
|   |   |   | 0 | 1 |
|   |   |   |   | 0 |

48

# Bottom-Up (agglomerative):

Start with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Consider all possible merges…

… Choose the best

49

# Bottom-Up (agglomerative):

Start with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.



Consider all possible merges… ... Choose the best

Consider all possible merges… ... Choose the best

50

# Bottom-Up (agglomerative):

Start with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.



Consider all possible merges...    Choose the best

Consider all possible merges...    Choose the best

Consider all possible merges...    Choose the best

51

# Bottom-Up (agglomerative):

Start with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Consider all possible merges…

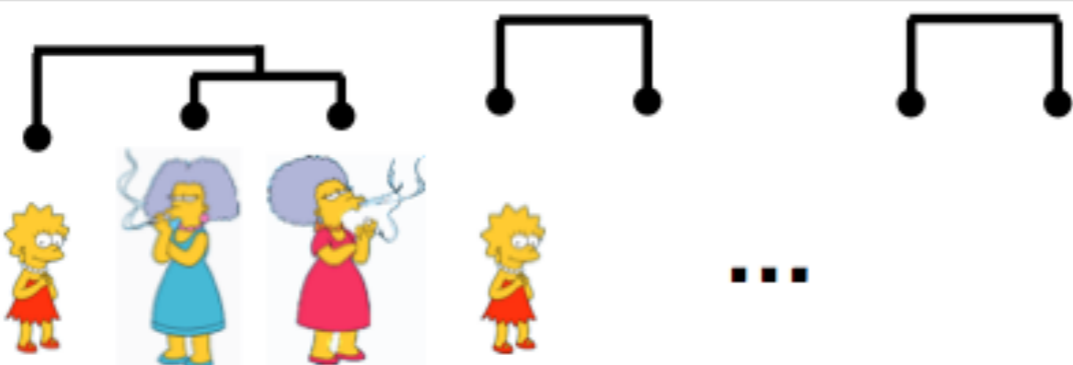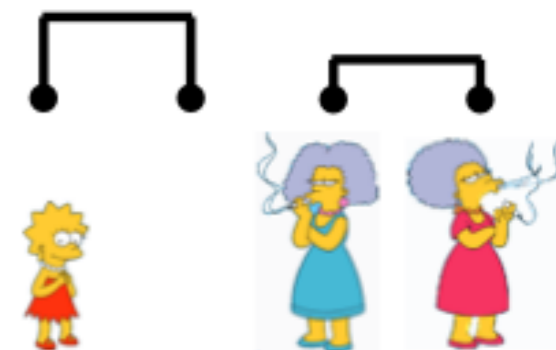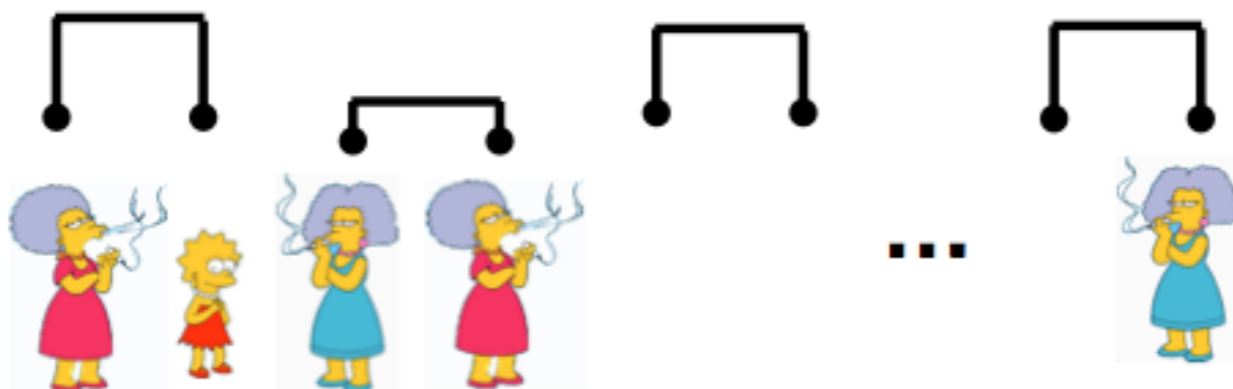Choose the best

Consider all possible merges…

…

But how do we compute distances between clusters rather than objects?

Consider all possible merges…

…

Choose the best

52

# Computing distance between clusters:
# **Single Link**

- Cluster distance = distance of two closest members in each class



- Potentially long and skinny clusters

53

# Computing distance between clusters:
# **Complete Link**

- Cluster distance = distance of two farthest members in each class

- Tight clusters

# Computing distance between clusters:
# **Average Link**

- Cluster distance = average distance of all pairs

- The most widely used measure
- Robust against noise

# Agglomerative Clustering

Good

- Simple to implement, widespread application
- Clusters have adaptive shapes
- Provides a hierarchy of clusters

Bad

- May have imbalanced clusters
- Still have to choose number of clusters or threshold
  - **— silhouette coefficient**
- Need to use an "ultrametric" to get a meaningful hierarchy

# What is a good clustering?

# What is a good clustering?

- Internal criterion: A good clustering will produce high quality clusters in which:
  - the intra-class (that is, intra-cluster) similarity is high
  - the inter-class similarity is low
  - The measured quality of a clustering depends on both the obj. representation and the similarity measure used

- External criteria for clustering quality
  - Quality measured by its ability to discover some or all of the hidden patterns or latent classes in gold standard data
  - Assesses a clustering with respect to ground truth
  - Example:
    - Purity
    - Entropy of classes in clusters (or Mutual Information between classes and clusters)

# External Evaluation of Cluster Quality

- Simple measure: purity, the ratio between the dominant class in the cluster and the size of cluster

  - Assume documents with C gold standard classes, while our clustering algorithms produce K clusters, $\omega_1, \omega_2, ..., \omega_K$ with $n_i$ members.

$$\text{purity}(\Omega, C) = \frac{1}{N} \sum_k \max_j |\omega_k \cap c_j|$$

  - Example:



Cluster I    Cluster II    Cluster III

purity = 1/17* (max(5, 1, 0)+max(1, 4, 1)+max(2, 0, 3))
= 1/17*(5+4+3) ≈ 0.71

Cluster II: Purity = 1/6 (max(1  4  1)) = 4/6

# External Evaluation of Cluster Quality

- Let:

$$TC = TC_1 \cup TC_2 \cup...\cup TC_n$$

$$CC = CC_1 \cup CC_2 \cup...\cup CC_m$$

be the target and computed clusterings, respectively.

- $TC = CC$ = original set of data

- Define the following:
  - $a$: number of pairs of items that belong to the same cluster in both $CC$ and $TC$

  - $b$: number of pairs of items that belong to different clusters in both $CC$ and $TC$

  - $c$: number of pairs of items that belong to the same cluster in $CC$ but different clusters in $TC$

  - $d$: number of pairs of items that belong to the same cluster in $TC$ but different clusters in $CC$

60

# External Evaluation of Cluster Quality

$$P = \frac{a}{a + c}$$

$$R = \frac{a}{a + d}$$

$$\frac{a+b}{a+b+c+d}$$

**Rand Index**

Measure of clustering agreement: how similar are these two ways of partitioning the data?

$$F = \frac{2 \times P \times R}{P + R}$$

**F-measure**

# External Evaluation of Cluster Quality

$$\frac{a+b}{a+b+c+d}$$

**Rand Index**

$$\frac{2(ab-cd)}{(a+c)(c+b)+(a+d)(d+b)}$$

**Adjusted Rand Index**

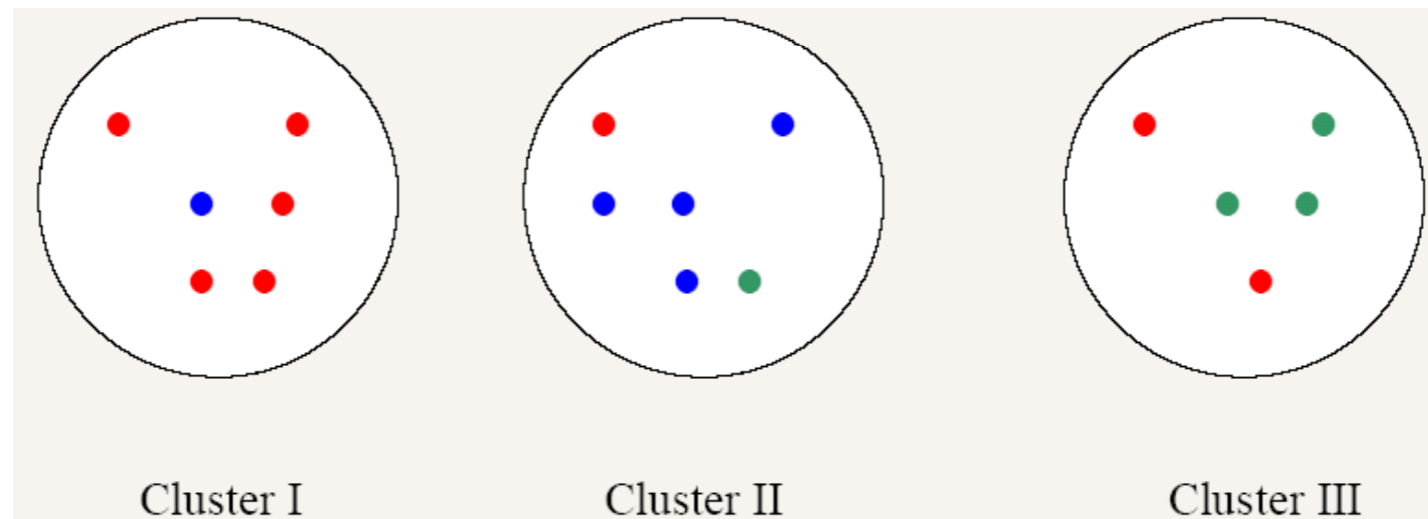Extension of the Rand index that attempts to account for items that may have been clustered by chance

slide by Christophe Giraud-Carrier

# External Evaluation of Cluster Quality

p        p   y

$$Entropy(CC_i) = \sum_{TC_j \in TC} -p(TC_j \,|\, CC_i) \log p(TC_j \,|\, CC_i)$$

$$AvgEntropy(CC) = \sum_{i=1}^{m} \frac{|CC_i|}{|CC|} Entropy(CC_i)$$

Measure of purity wrt the target clustering

- Example:



Cluster I          Cluster II          Cluster III

Entropy(CC$_1$) = (5/6)log(5/6) + (1/6)log(1/6) + (0/6)log(0/6) = -.650
Entropy(CC$_2$) = (1/6)log(1/6) + (4/6)log(4/6) + (1/6)log(1/6) = -1.252
Entropy(CC$_3$) = (2/5)log(2/5) + (0/5)log(0/5) + (3/5)log(3/5) = -.971

Cluster Purity = 1/6 (max (1 (4 1)) = 4/6

AvgEntropy(CC) = (-.650 * 6/17) + (-1.252 * 6/17) + (-.971 * 5/17)
AvgEntropy(CC) = -.956

63

# **Next Lecture:**
# Dimensionality Reduction