# Cluster Analysis (b)

# Outline

- **Grid-Based and Density-Based Algorithms**

- Graph-Based Algorithms

- Non-negative Matrix Factorization

- Cluster Validation

- Summary

# Density-Based Algorithms

- ☐ **One Motivation**
  - ■ Find clusters with arbitrary shape
- ☐ **The Key Idea**
  - ■ Identify fine-grained dense regions
  - ■ Merge regions into clusters
- ☐ **Representative Algorithms**
  - ■ Grid-Based Methods
  - ■ DBSCAN
  - ■ DENCLUE

# Grid-Based Methods

## ☐ The Algorithm

**Algorithm** *GenericGrid*(Data: $\mathcal{D}$, Ranges: $p$, Density: $\tau$ )
**begin**
  Discretize each dimension of data $\mathcal{D}$ into $p$ ranges;
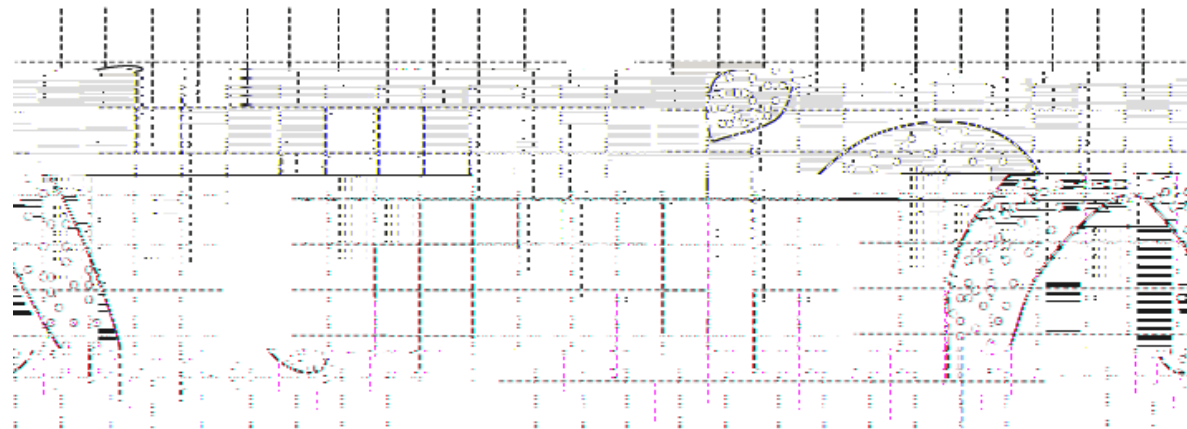  Determine dense grid cells at density level $\tau$;
  Create graph in which dense grids are connected if they are adjacent;
  Determine connected components of graph;
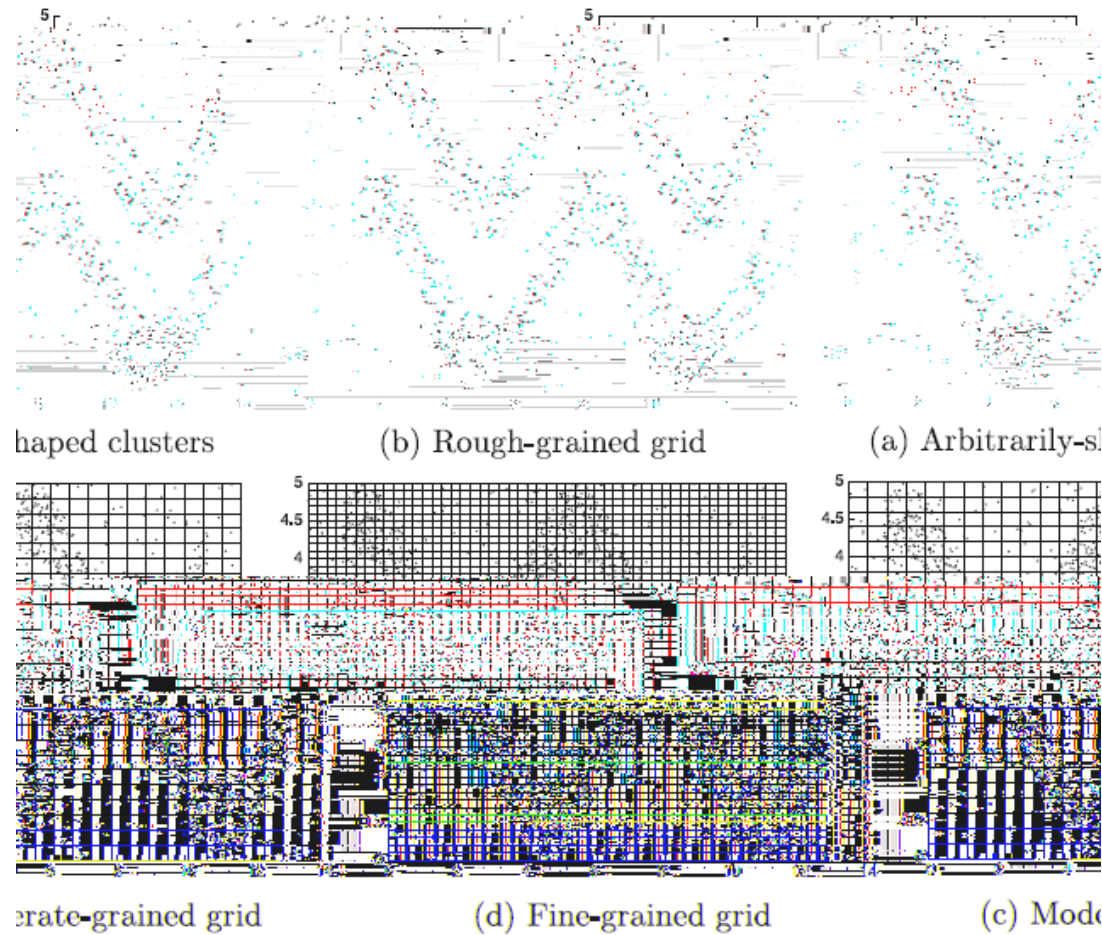  **return** points in each connected component as a cluster;
**end**

) Data points and grid          (b) Agglomerating adjacent grids          (a
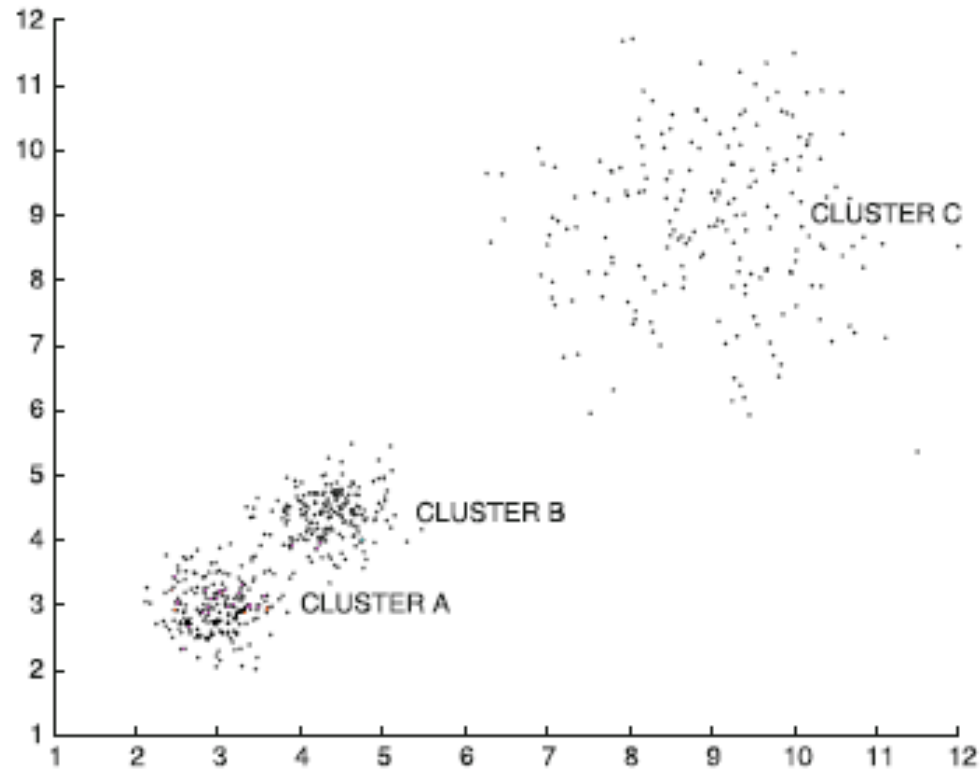
# Limitations-2 Parameters (1)

☐ The number of Grids



(b) Rough-grained grid    (a) Arbitrarily-s...haped clusters

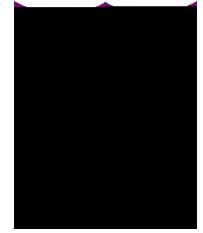(d) Fine-grained grid    (c) Mod...erate-grained grid

# Limitations-2 Parameters (2)
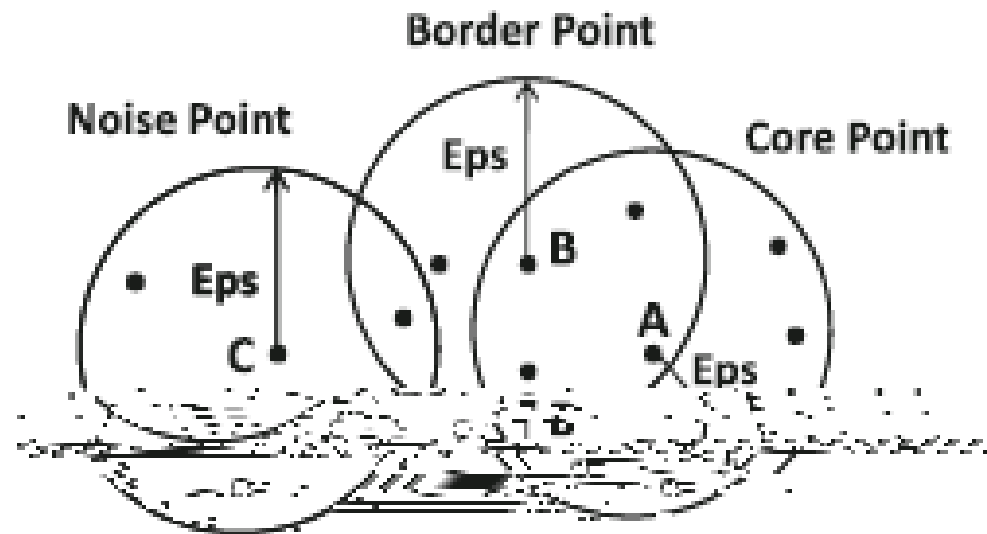
□ The Level of Density

# DBSCAN (1)

1. Classify data points into
   - **Core point**: A data point is defined as a core point, if it contains at least $\tau$ data points within a radius $Eps$.
   - **Border point**: A data point is defined as a border point, if it contains less than $\tau$ points, but it also contains at least one core point within a radius $Eps$.
   - **Noise point**: A data point that is neither a core point nor a border point is defined as a noise point.

# DBSCAN (2)

1. Classify data points into Core point, Border point, and Noise points.

# DBSCAN (3)

1. Classify data points into Core point, Border point, and Noise points.
2. A connectivity graph is constructed with respect to the core points

   - Core points are connected if they are within $Eps$ of one another

3. Determine connected components
4. Assign each border point to connected component

   - with which it is best connected

# Limitations of DBSCAN
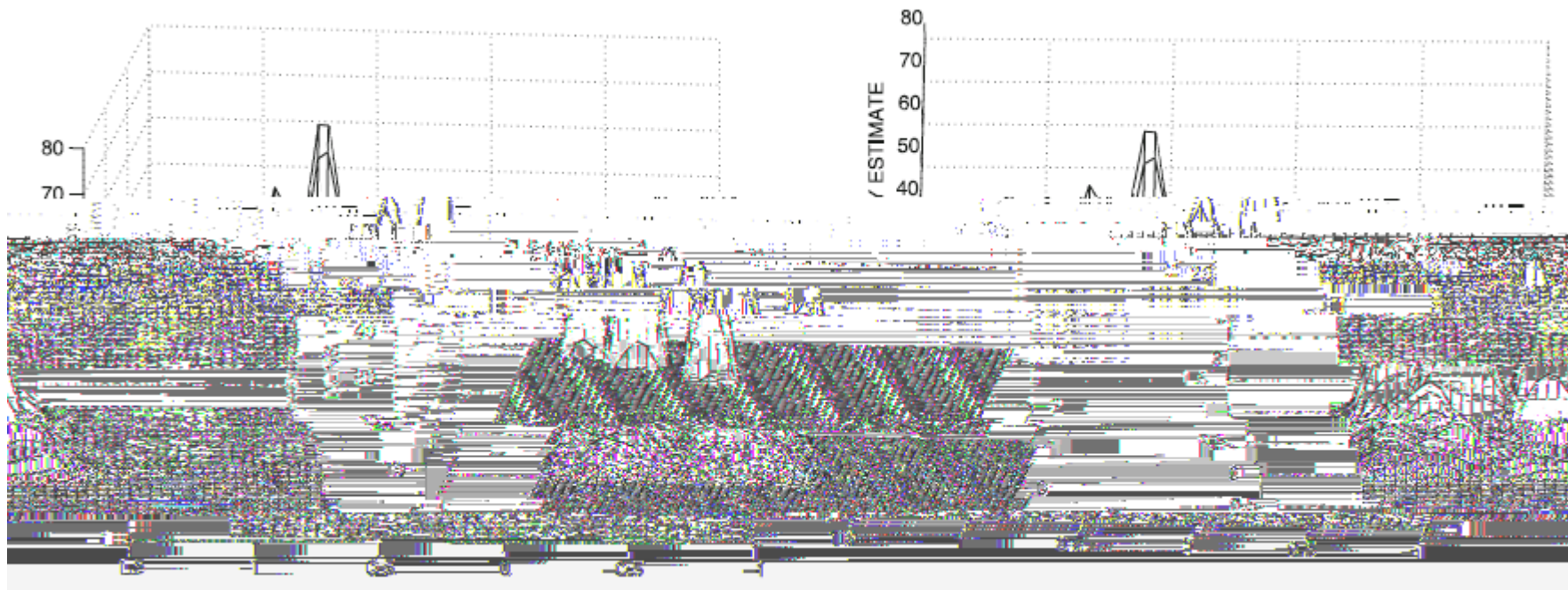
☐ Two Parameters
  ■ Radius

# DENCLUE—Preliminary

- Kernel-density Estimation
  - Given $n$ data points $X$

# DENCLUE—The Key Idea

☐ Determine clusters by using a density threshold $\tau$


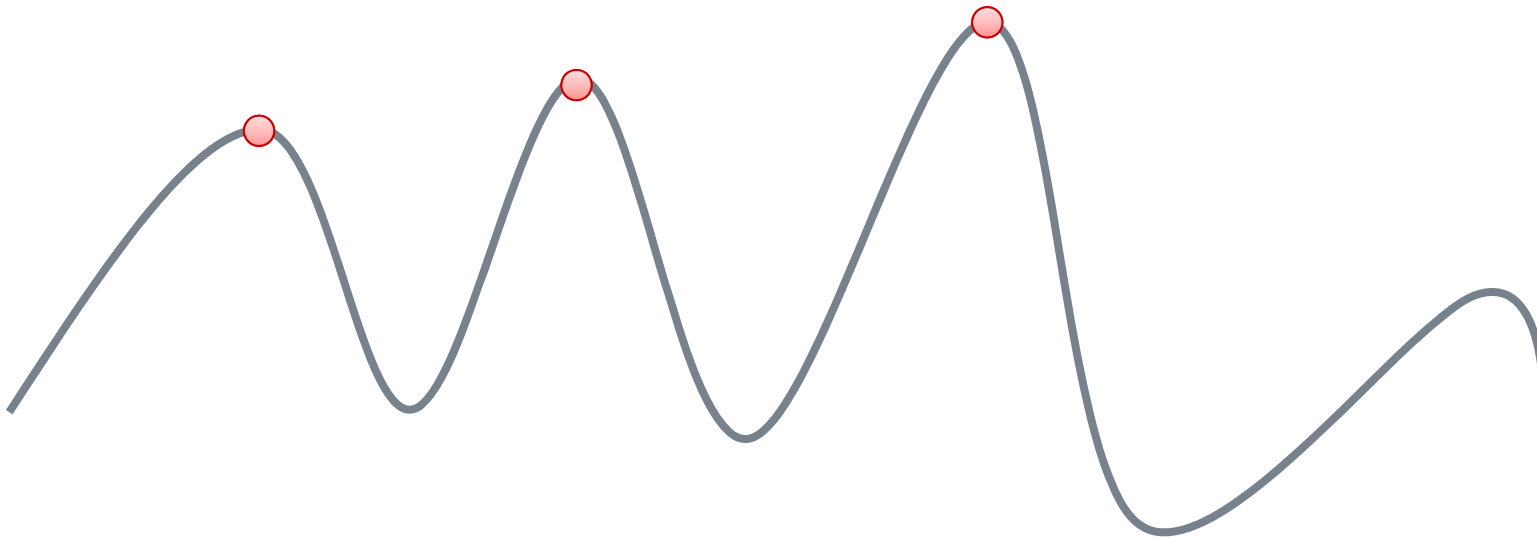
2 clusters                    3 clusters

# DENCLUE—Procedure
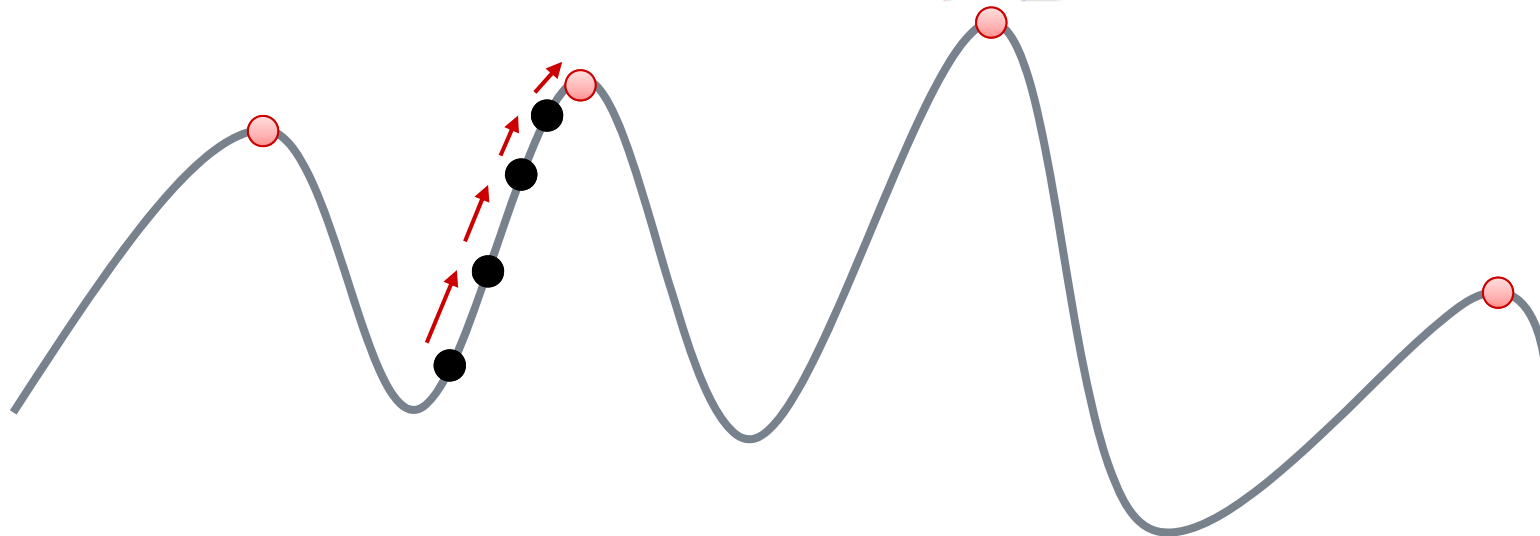
☐ **Density Attractors**
   ◼ Local Maximum/Peak

# DENCLUE—Procedure

- ☐ Density Attractors
    - ◼ Local Maximum/Peak
- ☐ Identify a Peak for Each Data Point
    - ◼ An iterative gradient ascent

$$X^{(t+1)} = X^{(t)} + a\nabla f(X^{(t)})$$

# DENCLUE—Procedure

- ☐ **Density Attractors**
  - ■ Local Maximum/Peak
- ☐ **Identify a Peak for Each Data Point**
  - ■ An iterative gradient ascent

$$X^{(t+1)} = X^{(t)} + a\nabla f(X^{(t)})$$

- ☐ **Post-Processing**
  - ■ Attractors whose density is smaller than $\tau$ are excluded
  - ■ Density attractors are connected to each other by a path of density at least $\tau$ will be merged

# DENCLUE—Implementation

☐ **Gradient Ascent**

  ■ Gradient

$$\nabla f(\overline{X}) = \frac{1}{n} \sum_{i=1}^{n} \nabla K(\overline{X} - \overline{X_i}).$$

  ■ Gaussian Kernel

$$\nabla K(\overline{X} - \overline{X_i}) \propto (\overline{X_i} - \overline{X}) K(\overline{X} - \overline{X_i})$$

☐ **Mean-shift Method**

$$\overline{X^{(t+1)}} = \frac{\sum_{i=1}^{n} \overline{X_i} K(\overline{X^{(t)}} - \overline{X_i})}{\sum_{i=1}^{n} K(\overline{X^{(t)}} - \overline{X_i})}$$
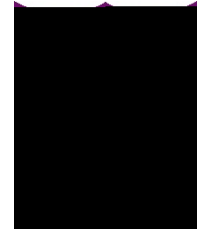
  ■ Converges much faster

# Outline

☐ Grid-Based and Density-Based Algorithms

☐ **Graph-Based Algorithms**

☐ Non-negative Matrix Factorization

☐ Cluster Validation

☐ Summary

# Graph Construction for a Set of $n$ Points $\mathcal{O} = \{O_1, \dots, O_n\}$

☐ A node is defined for each

# Spectral Clustering

☐ **Dimensionality Reduction**

  ■ Find a low-dimensional representation for each node in the graph



[   ]

  ■ Laplacian Eigenmap [Belkin and Niyogi, 2002]

☐ $k$-means

  ■ Apply $k$-means to new representations of the data

# Laplacian Eigenmap (1)

- ☐ The Objective Function ($k = 1$)
  - ■ $y \in \mathbb{R}$ is a 1-dimensional representation of $O$
  - ■ $w$ is the similarity between $O$ and $O$

$$( \quad )$$

  - ■ Similar points will be mapped closer
    - ✓ Similar points have larger weights

## The Objective Function ($k = 1$)

- Vector Form

- $\mathbf{y} = [y\,,\ldots,y\,]$

# Laplacian Eigenmap (3)

□ The Optimization Problem ($k = 1$)

$$\min \quad \mathbf{y}^\top L\mathbf{y}$$

$$\text{s.t.} \quad \mathbf{y}^\top D\mathbf{y} = 1$$

■ Add a Constraint to Remove Scaling Factor

✓ is introduced for normalization [Luxburg, 2007]

□ The Solution

$$L\mathbf{y} = \lambda D\mathbf{y}$$

■ Generalized Eigenproblem [Luxburg 2007]

■ The smallest eigenvector is

✓ Useless since

# Laplacian Eigenmap (3)

- ☐ The Optimization Problem $(k =$

# Laplacian Eigenmap (4)

- ☐ The Objective Function ($k > 1$)
  - ■ Vector Form

  $$\| \quad \|^2 \quad \text{trace}( \quad )$$

  - ■ $Y = [\mathbf{y}, \dots, \mathbf{y}] \in \mathbb{R}$
  - ■ $L = D - W \in \mathbb{R}$    is the graph Laplacian
  - ■ $W = [w] \in \mathbb{R}$    is the similarity matrix
  - ■ $D \in \mathbb{R}$    is a diagonal matrix with $D =$
    $\sum w$

# Laplacian Eigenmap (4)

- ☐ The Optimization Problem $(k > 1)$

$$\min \quad \text{trace}(Y^\top LY)$$

$$\text{s.t.} \quad Y^\top DY = I$$

- ☐ The Solution

$$L\mathbf{y} = \lambda D\mathbf{y}$$

  - ■ Generalized Eigenproblem [Luxburg 2007]
  - ■ Use [ ,..., ] as the optimal solution
    - ✓ is the -th generalized eigenvector
    - ✓ The new representation for is the -th row of
  - ■ Don't forget the normalization

# Properties of Spectral Clustering

□ Varying Cluster Shape and Density



CLUSTER C (SPARSE)

CLUSTER A
(ARBITRARY SHAPE)

THE TWO DENSELY
CONNECTED
COMMUNITIES OF
THE K-NEAREST

THE THREE DENSELY
CONNECTED
COMMUNITIES OF

NEIGHBOR GRAPH

NEIGHBOR-GRAPH

CLUSTER B

CLUSTER D (DENSE)

CLUSTER E (DENSE)

(b) Varying cluster density

(a) Varying cluster shap

■ Due to the nearest neighbor graph

□ High Computational Cost

# Outline

☐ Grid-Based and Density-Based
Algorithms

☐ Graph-Based Algorithms

☐ **Non-negative Matrix Factorization**

☐ Cluster Validation

☐ Summary

# Non-negative Matrix Factorization (NMF)

- ☐ Let $X = [\mathbf{x}, \ldots, \mathbf{x}] = \mathbb{R}$ be a non-negative data matrix
- ☐ NMF aims to factor $X$ as $U \times V$
  - ■ $U \in \mathbb{R}$ and $V \in \mathbb{R}$ are non-negative
- ☐ The Optimization Problem

$$\min_{,} \quad \|X - UV\|$$
$$\text{s.t.} \quad U \geq 0, V \geq 0$$

  - ■ Non-convex

# Interpretation of NMF (1)

- ☐ **Matrix Appromation**

$$X \approx UV$$

- ☐ **Element-wise**
  - ■ $X = [\mathbf{x}, ..., \mathbf{x}] \in \mathbb{R}$ , where $\mathbf{x} \in \mathbb{R}$
  - ■ $U = [\mathbf{u}, ..., \mathbf{u}] \in \mathbb{R}$ , where $\mathbf{u} \in \mathbb{R}$
  - ■ $V = [\mathbf{v}, ..., \mathbf{v}] \in \mathbb{R}$ , where $\mathbf{v} \in \mathbb{R}$
    - ✓ is the -th column of
    - ✓ is the -th row of
  - ■ Then,

    - ✓ is the -th element of vector

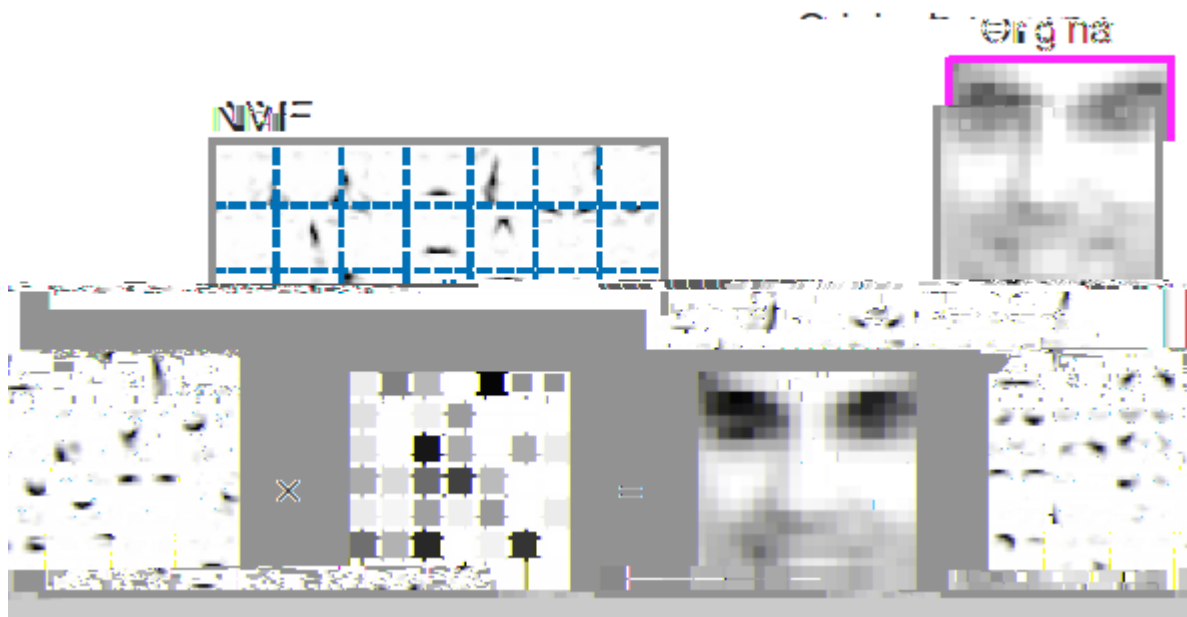# Interpretation of NMF (2)

# Parts-Based Representations

☐ When each **x** is a face image



■ [Lee and Seung, 1999]

# Clustering by NMF

☐ Vector Approximation

- ▪ $\mathbf{u}$ can be treated as an representative of the $j$-th cluster
- ▪ $v$ can be treated as the association between $\mathbf{x}$ and $\mathbf{u}$

☐ The cluster label $l$ for $\mathbf{x}$

$$\text{argmax}$$

- ▪ [Xu et al., 2003]

# An Example

☐ Discover both Row and Column Clusters

$$
\begin{array}{|c|c|c|c|c|c|}
\hline
2 & 2 & 1 & 2 & 0 & 0 \\
\hline
2 & 3 & 3 & 3 & 0 & 0 \\
\hline
1 & 1 & 1 & 1 & 0 & 0 \\
\hline
2 & 2 & 2 & 3 & 1 & 1 \\
\hline
0 & 0 & 0 & 1 & 1 & 1 \\
\hline
0 & 0 & 0 & 2 & 1 & 2 \\
\hline
\end{array}
\approx
\begin{array}{|c|c|}
\hline
2 & 0 \\
\hline
3 & 0 \\
\hline
1 & 0 \\
\hline
2 & 1 \\
\hline
0 & 1 \\
\hline
0 & 2 \\
\hline
\end{array}
\times
\begin{array}{|c|c|c|c|c|c|}
\hline
1 & 1 & 1 & 1 & 0 & 0 \\
\hline
0 & 0 & 0 & 1 & 1 & 1 \\
\hline
\end{array}
$$

# Optimization in NMF

□ Alternating between $U$ and $V$

$$u_{ij} \leftarrow u_{ij} \frac{(\mathbf{XV})_{ij}}{(\mathbf{UV}^T\mathbf{V})_{ij}}$$

$$v_{ij} \leftarrow v_{ij} \frac{(\mathbf{X}^T\mathbf{U})_{ij}}{(\mathbf{VU}^T\mathbf{U})_{ij}}$$

■ Local Optimal Solutions
  ✓ Run multiple times and choose the best one

□ Other Optimization Algorithms are also Possible

# Outline

- Grid-Based and Density-Based Algorithms

- Graph-Based Algorithms

- Non-negative Matrix Factorization

- **Cluster Validation**

- Summary

# Concepts

- Cluster validation
  - Evaluate the quality of a clustering

- Internal Validation Criteria
  - Do not need additional information
  - Biased toward one algorithm or the other

- External Validation Criteria
  - Ground-truth clusters are known
  - Ground-truth may not reflect the natural clusters in the data

# Internal Validation Criteria

☐ Sum of square distances to centroids

$$\| \qquad \|$$

☐ Intracluster to intercluster distance ratio

$$Intra = \sum_{(\overline{X_i}, \overline{X_j}) \in P} dist(\overline{X_i}, \overline{X_j}) / |P|$$

$$Inter = \sum_{(\overline{X_i}, \overline{X_j}) \in Q} dist(\overline{X_i}, \overline{X_j}) / |Q|.$$

☐ Silhouette coefficient

☐ Probabilistic measure

# External Validation Criteria

- ☐ **Class Labels**
  - ■ The Ground-truth
- ☐ **Confusion Matrix**
  - ■ Each row $i$ corresponds to the class label $j$
  - ■ Each column $j$ corresponds to the algorithm-determined cluster $j$

| Cluster Indices | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 97 | 0 | 2 | 1 |
| 2 | 5 | 191 | 1 | 3 |
| 3 | 4 | 3 | 87 | 6 |
| 4 | 0 | 0 | 5 | 195 |

| Cluster Indices | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 33 | 30 | 17 | 20 |
| 2 | 51 | 101 | 24 | 24 |
| 3 | 24 | 23 | 31 | 22 |
| 4 | 46 | 40 | 44 | 70 |

  - ■ Ideal clustering ⇒ a diagonal matrix after permutation

# Notations

☐ $m$ : number of data points from class (*ground-truth*) cluster $i$ that are mapped to (*algorithm-determined*) cluster $j$

☐ $N$ : number of data points in *true cluster i*

$$N_i = \sum_{j=1}^{k_d} m_{ij} \qquad\qquad \forall i = 1 \ldots k_t$$

☐ $M$ : number of data points in *algorithm-determined* cluster $j$

$$M_j = \sum_{i=1}^{k_t} m_{ij} \qquad\qquad \forall i = 1 \ldots k_t$$

# Purity

☐ For a given algorithm-determined cluster

# Gini index

- ☐ Limitation of Purity
  - ■ Only accounts for the dominant label in the cluster and ignores the distribution of the remaining points
- ☐ Gini index $G$ for column (algorithm-determined cluster) $j$

$$G_j = 1 - \sum_{i=1}^{k_t} \left( \frac{m_{ij}}{M_j} \right)^2$$

- ☐ The average Gini coefficient
  - ■ Low values

$$G_{average} = \frac{\sum_{j=1}^{k_d} G_j \cdot M_j}{\sum_{j=1}^{k_d} M_j}$$

# Outline

- ☐ Grid-Based and Density-Based Algorithms

- ☐ Graph-Based Algorithms

- ☐ Non-negative Matrix Factorization

- ☐ Cluster Validation

- ☐ **Summary**

# Summary

- **Grid-Based and Density-Based Algorithms**
  - Grid-Based Methods
  - DBSCAN, DENCLUE
- **Graph-Based Algorithms**
  - Laplacian Eigenmap
- **Non-negative Matrix Factorization**
- **Cluster Validation**
  - Purity, Gini index

# Reference

- [Belkin and Niyogi, 2002] Belkin, M. and Niyogi, P. (2002). Laplacian eigenmaps and spectral techniques for embedding and clustering. In NIPS 14, pages 585–591.

- [Luxburg, 2007] Luxburg, U. (2007). A tutorial on spectral clustering. Statistics and Computing, 17(4):395–416.

- [Lee and Seung, 1999] Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. Nature, 401(6755):788–791.

- [Xu et al., 2003] Xu, W., Liu, X., and Gong, Y. (2003). Document clustering based on non-negative matrix factorization. In SIGIR, pages 267–273.

- [Hinneburg and Keim, 1998] Hinneburg, A. and Keim, D. A. (1998). An efficient approach to clustering in large multimedia databases with noise. In KDD, pages 58–65.