

Advanced Data Analysis: Spectral Clustering

Masashi Sugiyama (Computer Science)

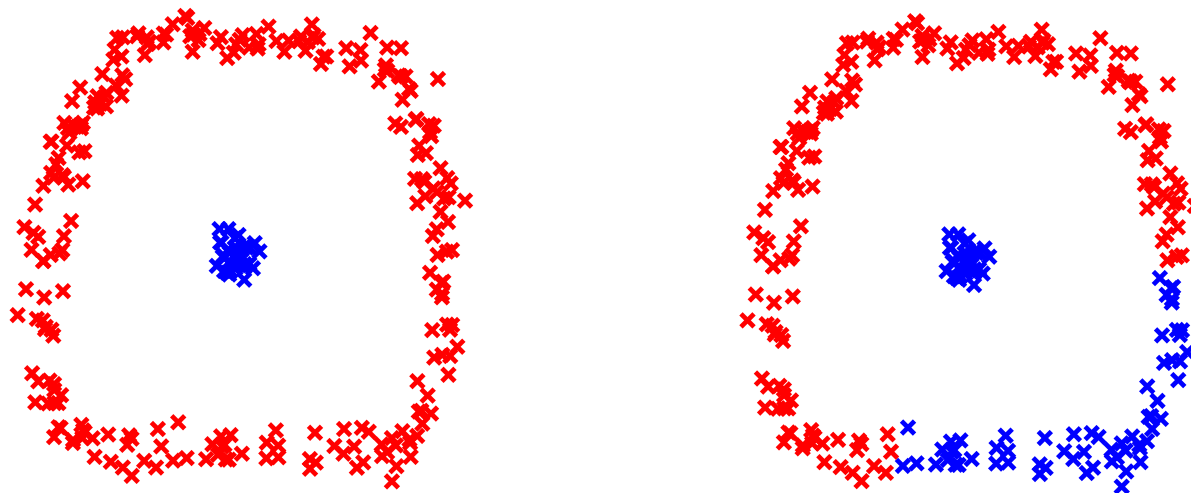
W8E-505, sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi>

Kernel K-Means

153

- Ordinary k-means clustering does not work well if the data crowds have non-convex shapes.
- Kernel k-means is more flexible.
- However, solution depends **crucially** on the initial cluster assignments since clustering is carried out in a high-dimensional feature space.



Similarity-Based Clustering 154

- Similarity matrix W : $W_{i,j}$ is large if x_i and x_j are similar.
- Assumptions on W :
 - Symmetric: $W_{i,j} = W_{j,i}$
 - Positive entries: $W_{i,j} \geq 0$
 - Invertible: $\exists W^{-1}$
 - Positive semi-definite: $\forall \mathbf{y}, \langle W\mathbf{y}, \mathbf{y} \rangle \geq 0$

Examples of Similarity Matrix¹⁵⁵

$$W_{i,j} = W(\mathbf{x}_i, \mathbf{x}_j)$$

■ Distance-based:

$$W(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / \gamma^2) \quad \gamma > 0$$

■ Nearest-neighbor-based:

$W(\mathbf{x}_i, \mathbf{x}_j) = 1$ if \mathbf{x}_i is a k' -nearest neighbor of \mathbf{x}_j or \mathbf{x}_j is a k' -nearest neighbor of \mathbf{x}_i .
Otherwise $W(\mathbf{x}_i, \mathbf{x}_j) = 0$.

■ Combination of two is also possible.

$$W(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / \gamma^2) \\ 0 \end{cases}$$

Local Scaling Heuristic

156

- γ_i : scaling around the sample \mathbf{x}_i

$$\gamma_i = \|\mathbf{x}_i - \mathbf{x}_i^{(k)}\|$$

$\mathbf{x}_i^{(k)}$: k-th nearest neighbor sample of \mathbf{x}_i

- Local scaling based similarity matrix:

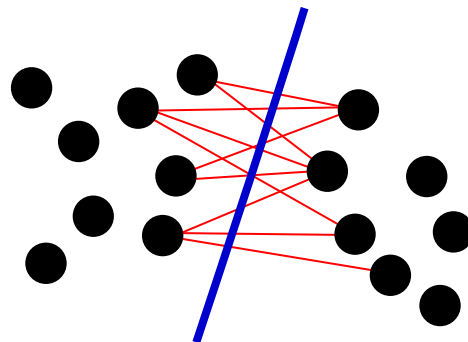
$$\mathbf{W}_{i,j} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (\gamma_i \gamma_j))$$

- A heuristic choice is $k = 7$.

- **Idea:** Minimize sum of similarities between samples inside and outside the cluster
- In two-cluster cases:

$$\min_{\mathcal{C}_1, \mathcal{C}_2} \left[\sum_{x \in \mathcal{C}_1} \sum_{x' \in \mathcal{C}_2} W(x, x') + \sum_{x \in \mathcal{C}_2} \sum_{x' \in \mathcal{C}_1} W(x, x') \right]$$

- From a graph-theoretic viewpoint, this corresponds to finding **minimum cut**.

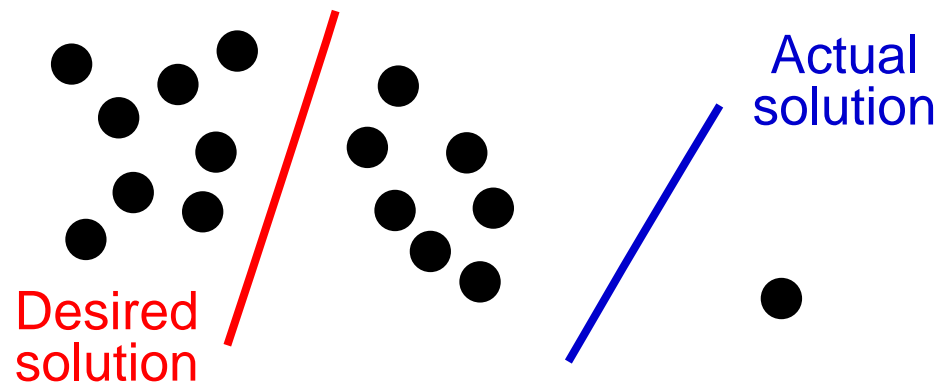


Cut Criterion (cont.)

158

$$\min_{\mathcal{C}_1, \mathcal{C}_2} \left[\sum_{x \in \mathcal{C}_1} \sum_{x' \in \mathcal{C}_2} W(x, x') + \sum_{x \in \mathcal{C}_2} \sum_{x' \in \mathcal{C}_1} W(x, x') \right]$$

- Mincut method tends to give a cluster with a **very small number of samples**.



Normalized Cut Criterion

159

- **Idea:** Penalize small clusters
- In two-cluster cases:

$$\min_{\mathcal{C}_1, \mathcal{C}_2} \left[\frac{\sum_{\mathbf{x} \in \mathcal{C}_1} \sum_{\mathbf{x}' \in \mathcal{C}_2} W(\mathbf{x}, \mathbf{x}')}{n} + \frac{\sum_{\mathbf{x} \in \mathcal{C}_2} \sum_{\mathbf{x}' \in \mathcal{C}_1} W(\mathbf{x}, \mathbf{x}')}{n} \right]$$
$$\left[\frac{\sum_{\mathbf{x}'' \in \mathcal{C}_1} \sum_{j=1} W(\mathbf{x}'', \mathbf{x}_j)}{\sum_{\mathbf{x}'' \in \mathcal{C}_1} \sum_{j=1} W(\mathbf{x}'', \mathbf{x}_j)} + \frac{\sum_{\mathbf{x}'' \in \mathcal{C}_2} \sum_{j=1} W(\mathbf{x}'', \mathbf{x}_j)}{\sum_{\mathbf{x}'' \in \mathcal{C}_2} \sum_{j=1} W(\mathbf{x}'', \mathbf{x}_j)} \right]$$

- Denominator is a normalization factor, which is the sum of similarities between samples inside the class and all samples.

Normalized Cut Criterion (cont.)¹⁶⁰

- In k -cluster cases, normalized cut is defined as

$$\operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{Ncut}]$$

$$J_{Ncut} = \sum_{i=1}^k \left[\frac{\sum_{\mathbf{x} \in \mathcal{C}_i} \sum_{\mathbf{x}' \notin \mathcal{C}_i} W(\mathbf{x}, \mathbf{x}')}{\sum_{\mathbf{x}'' \in \mathcal{C}_i} \sum_{j=1}^n W(\mathbf{x}'', \mathbf{x}_j)} \right]$$

Normalized Cut As Weighted¹⁶¹ Kernel K-Means (Homework)

■ Weighted kernel k-means criterion with

- **Weight:** $d(\mathbf{x}) = \sum_{i=1}^n W(\mathbf{x}, \mathbf{x}_i)$

- **Kernel:** $K(\mathbf{x}_i, \mathbf{x}_j) = W(\mathbf{x}_i, \mathbf{x}_j) / (d(\mathbf{x}_i)d(\mathbf{x}_j))$

shares the same optimal solution as
the normalized cut criterion:

$$\operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{Ncut}] = \operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{WS}]$$

$$J_{WS} = \sum_{i=1}^k \sum_{\mathbf{x} \in \mathcal{C}_i} d(\mathbf{x}) \|\phi(\mathbf{x}) - \boldsymbol{\mu}_i\|^2$$

$$\boldsymbol{\mu}_i = \frac{1}{s_i} \sum_{\mathbf{x}' \in \mathcal{C}_i} d(\mathbf{x}') \phi(\mathbf{x}')$$
$$s_i = \sum_{\mathbf{x} \in \mathcal{C}_i} d(\mathbf{x})$$

Algorithm 1

162

- Clustering based on the normalized cut criterion can be obtained by **weighted kernel k-means algorithm** with

$$d(\mathbf{x}) = \sum_{i=1}^n W(\mathbf{x}, \mathbf{x}_i)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = [\mathbf{D}^{-1} \mathbf{W} \mathbf{D}^{-1}]_{i,j}$$

1. Randomly initialize partition: $\{\mathcal{C}_i\}_{i=1}^k$
2. Update cluster assignments until convergence:

$$\mathbf{x}_j \rightarrow \mathcal{C}_t$$

$$t = \operatorname{argmin}_i \left[-\frac{2}{s_i} \sum_{\mathbf{x}' \in \mathcal{C}_i} d(\mathbf{x}') K(\mathbf{x}_j, \mathbf{x}') + \frac{1}{s_i^2} \sum_{\mathbf{x}', \mathbf{x}'' \in \mathcal{C}_i} d(\mathbf{x}') d(\mathbf{x}'') K(\mathbf{x}', \mathbf{x}'') \right]$$

Normalized Cut As Weighted¹⁶³ Kernel K-Means (cont.)

- Normalized-cut clustering looks reasonable.
- But it is solved by (weighted) kernel k-means in the end.
- Thus the drawback (strong dependency on initial cluster assignment) of kernel k-means still remains.

Dual Formulation

164

$$\operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{Ncut}]$$

- Instead of optimizing $\{\mathcal{C}_i\}_{i=1}^k$, we optimize **cluster indicator** A :

$$A_{i,j} = \begin{cases} 1 & \text{if } x_j \in \mathcal{C}_i \\ 0 & \text{o.w.} \end{cases}$$

- An optimizer of J_{Ncut} is given by

$$\operatorname{argmin}_{A \in \mathcal{B}^{k \times n}} [\operatorname{tr}(A L A^\top)]$$

(Homework)

$$\text{subject to } A D A^\top = I_k$$

$\mathcal{B}^{k \times n}$: Set of all $k \times n$ matrices such that one of the elements in each column takes one and others are all zero

Relation to Laplacian Eigenmap¹⁶⁵

- Let us allow A to take any real values.
- Then relaxed problem is given as

$$\min_{A \in \mathbb{R}^{k \times n}} \left[\text{tr}(A L A^\top) \right]$$

$$\text{subject to } A D A^\top = I_k$$

$$L = D - W \quad D = \text{diag}\left(\sum_{j=1}^n W_{i,j}\right)$$

- This is equivalent to **Laplacian eigenmap!**
- **Implication:** Laplacian eigenmap embedding “softly” clusters the data samples!

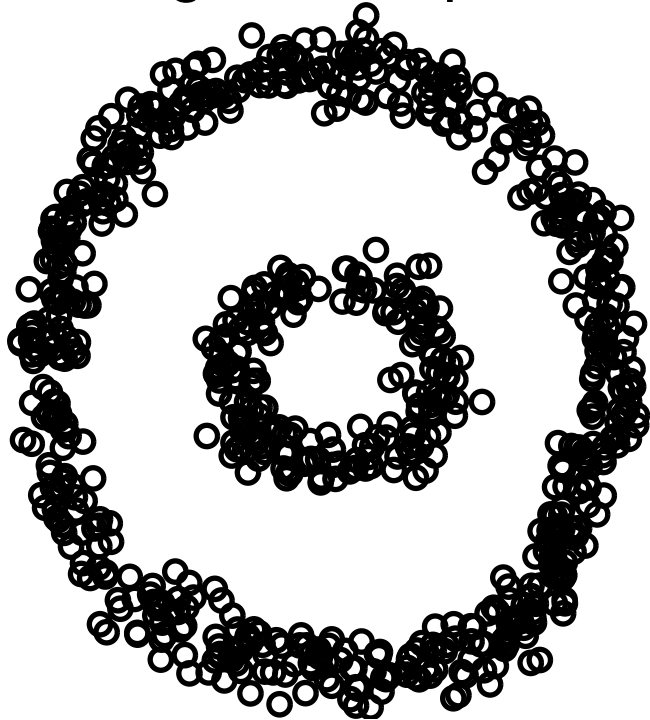
Algorithm 2 (Spectral Clustering)¹⁶⁶

1. Embed $\{x_i\}_{i=1}^n$ into $(k - 1)$ - dimensional space by **Laplacian eigenmap embedding**.
 2. Cluster the embedded samples by **(non-kernelized) k-means clustering algorithm**.
- Kernel k-means had a drawback that the clustering results crucially depend on the **initial cluster assignment**.
 - Since Laplacian eigenmap has **soft clustering property**, the above algorithm is less dependent on initialization.

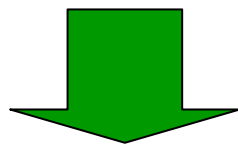
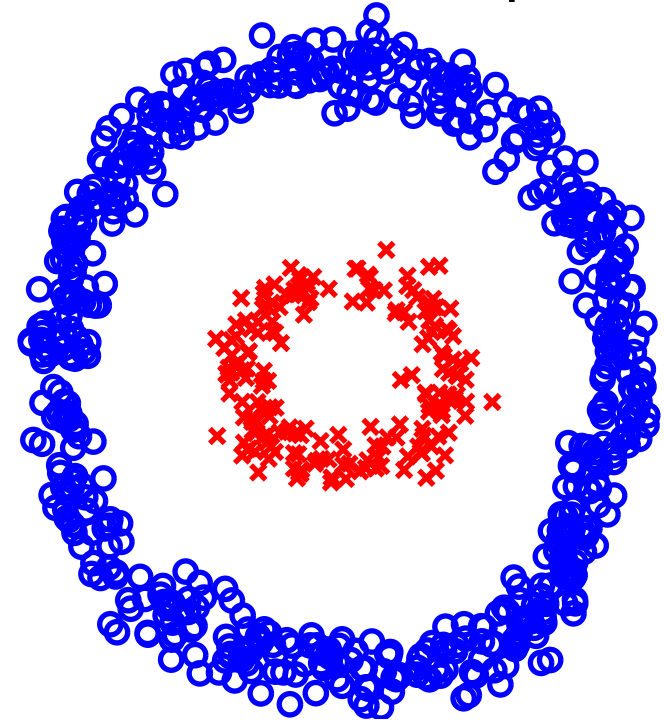
Examples

167

Original samples



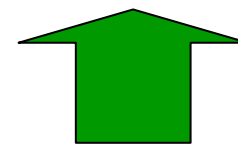
Clustered samples



Laplacian
eigenmap



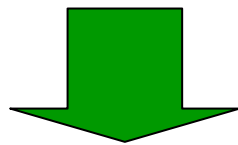
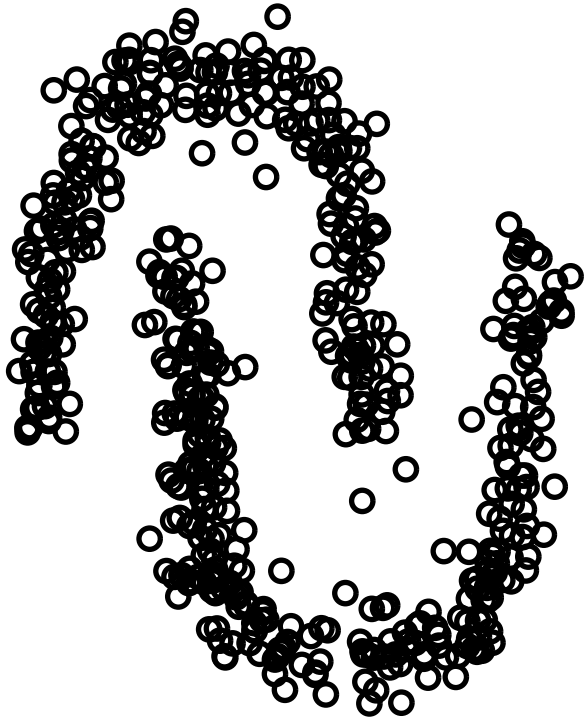
Ordinary
k-means



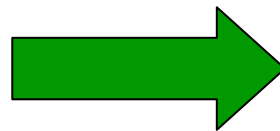
Examples (cont.)

168

Original samples

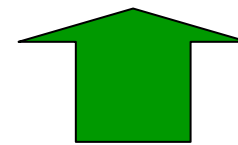
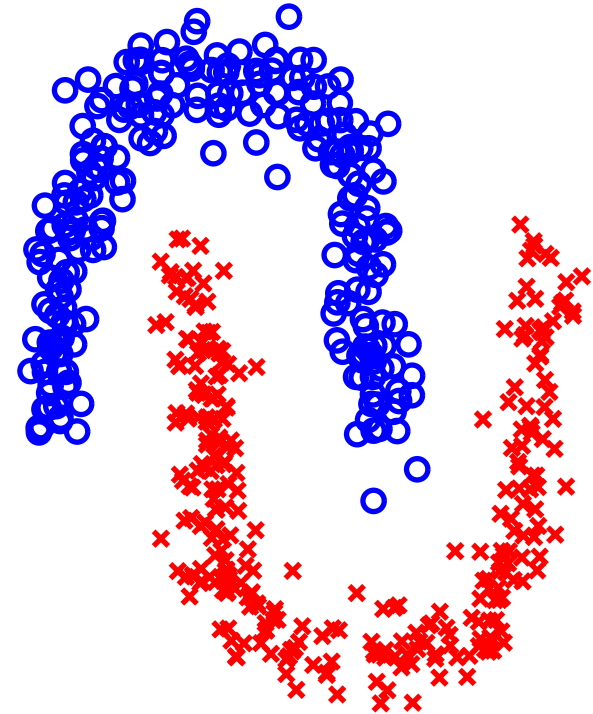


Laplacian
eigenmap



Ordinary
k-means

Clustered samples



Summary of Clustering Methods¹⁶⁹

- Three different families result in the same criterion!!

- K-means
- Kernel k-means
- Weighted kernel k-means

- Min-cut
- Normalized min-cut

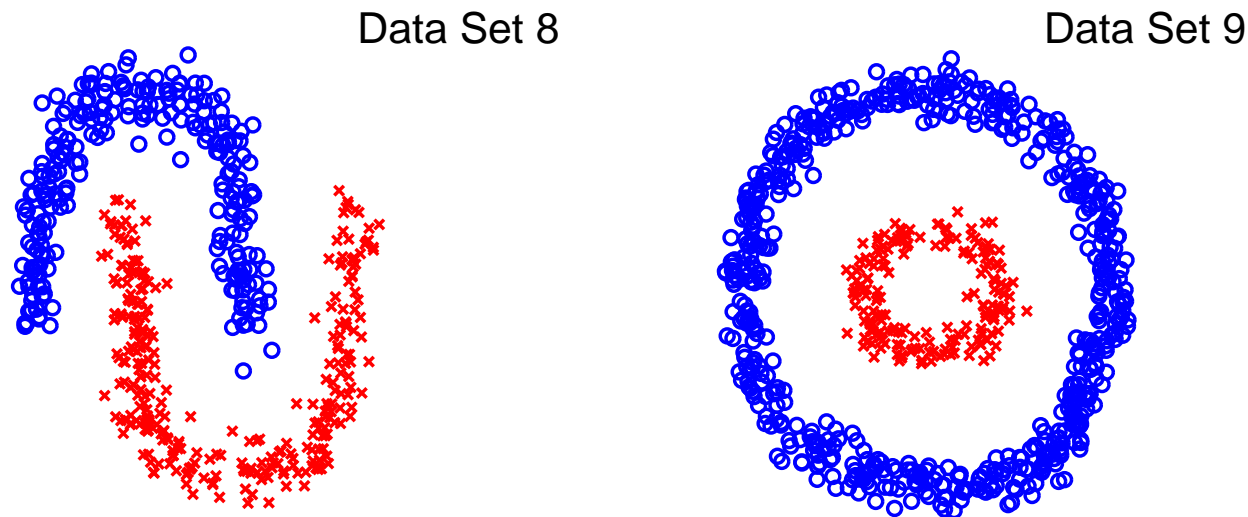
- Locality preserving projection
- Laplacian eigenmap
- “Hard” Laplacian eigenmap

Homework

170

1. Implement Algorithm 2 (spectral clustering) and reproduce the 2-dimensional examples shown in the class.

<http://sugiyama-www.cs.titech.ac.jp/~sugi/data/DataAnalysis>



Test the algorithm with your own (artificial or real) data and analyze their characteristics.

Homework (cont.)

171

2. Prove that weighted kernel k-means criterion with

- **Weight:** $d(\mathbf{x}) = \sum_{i=1}^n W(\mathbf{x}, \mathbf{x}_i)$
- **Kernel:** $K(\mathbf{x}_i, \mathbf{x}_j) = W(\mathbf{x}_i, \mathbf{x}_j) / (d(\mathbf{x}_i)d(\mathbf{x}_j))$

shares the same optimal solution as the normalized cut criterion:

$$\operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{Ncut}] = \operatorname{argmin}_{\{\mathcal{C}_i\}_{i=1}^k} [J_{WS}]$$

Homework (cont.)

172

2. Hint:

Express all elements in J_{WS} in terms of the affinity $W(\mathbf{x}, \mathbf{x}')$, e.g.,

$$s_i = \sum_{\mathbf{x}'' \in \mathcal{C}_i} \sum_{j=1}^n W(\mathbf{x}'', \mathbf{x}_j)$$

$$J_{WS} = \sum_{i=1}^k \sum_{\mathbf{x} \in \mathcal{C}_i} d(\mathbf{x}) \|\phi(\mathbf{x}) - \boldsymbol{\mu}_i\|^2$$

$$\boldsymbol{\mu}_i = \frac{1}{s_i} \sum_{\mathbf{x}' \in \mathcal{C}_i} d(\mathbf{x}') \phi(\mathbf{x}')$$

$$s_i = \sum_{\mathbf{x} \in \mathcal{C}_i} d(\mathbf{x})$$

$$J_{Ncut} = \sum_{i=1}^k \left[\frac{\sum_{\mathbf{x} \in \mathcal{C}_i} \sum_{\mathbf{x}' \notin \mathcal{C}_i} W(\mathbf{x}, \mathbf{x}')}{\sum_{\mathbf{x}'' \in \mathcal{C}_i} \sum_{j=1}^n W(\mathbf{x}'', \mathbf{x}_j)} \right]$$

Homework (cont.)

173

3. Prove that an optimizer of J_{Ncut} is given by

$$\operatorname{argmin}_{\mathbf{A} \in \mathcal{B}^{k \times n}} \left[\operatorname{tr}(\mathbf{A} \mathbf{L} \mathbf{A}^\top) \right]$$

$$\text{subject to } \mathbf{A} \mathbf{D} \mathbf{A}^\top = \mathbf{I}_k$$

$\mathcal{B}^{k \times n}$: Set of all $k \times n$ matrices such that one of the elements in each column takes one and others are all zero

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

$$\mathbf{D} = \operatorname{diag}(\sum_{j=1}^n \mathbf{W}_{i,j})$$

$$\mathbf{A}_{i,j} = \begin{cases} 1 & \text{if } \mathbf{x}_j \in \mathcal{C}_i \\ 0 & \text{o.w.} \end{cases}$$

Homework (cont.)

174

3. Hint:

Let $A = (a_1 | a_2 | \cdots | a_k)^\top$ and express all elements in J_{Ncut} in terms of $\{a_i\}_{i=1}^k$, e.g.,

$$\sum_{x'' \in \mathcal{C}_i} \sum_{j=1}^n W(x'', x_j) = \langle W a_i, \mathbf{1}_n \rangle = \langle D a_i, a_i \rangle$$

$$J_{Ncut} = \sum_{i=1}^k \left[\frac{\sum_{x \in \mathcal{C}_i} \sum_{x' \notin \mathcal{C}_i} W(x, x')}{\sum_{x'' \in \mathcal{C}_i} \sum_{j=1}^n W(x'', x_j)} \right]$$

Notification of Final Assignment

- **Data Analysis:** Apply dimensionality reduction or clustering techniques to your own data set and “mine” something interesting!

Mini-Conference on Data Analysis

- At the end of the semester, we have a **mini-conference on data analysis**.
- Some of the students may present their data analysis results.
- Those who give a talk at the conference will have **very good grades!**

Schedule

177

- June 9th: Regular lecture (spectral clustering)
- June 16th: Preparation for the mini-conference
(no lecture)
- June 23rd: Regular lecture (projection pursuit 1)
- June 30th: Regular lecture (projection pursuit 2)
- July 7th: Preparation for the mini-conference
(no lecture)
- July 14th: Mini-conference (day 1)
- July 21st: Mini-conference (day 2 if necessary)

Mini-Conference on Data Analysis

- Application procedure: On **June 23rd**, just say to me “**I want to give a talk!**”.
- Presentation: **approx. 10 min (?)**
 - Description of your data
 - Methods to be used
 - Outcome
- Slides should be in English.
- Better to speak in English, but Japanese may also be allowed (perhaps your friends will provide simultaneous translation!).