Advanced Data Analysis: K-Means Clustering

Masashi Sugiyama (Computer Science)

W8E-505, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

## **Data Clustering**

129

We want to divide data samples {x<sub>i</sub>}<sup>n</sup><sub>i=1</sub> into k (1 ≤ k ≤ n) disjoint clusters so that samples in the same cluster are similar.
We assume that k is prefixed.

0



0

# Within-Cluster Scatter Criterion<sup>30</sup>

Idea: Cluster the samples so that withincluster scatter is minimized

 $\mathcal{C}_i$ : Set of samples in cluster *i* 

$$igcup_{i=1}^k \mathcal{C}_i = \{oldsymbol{x}_j\}_{j=1}^n$$

$$\min_{\{\mathcal{C}_i\}_{i=1}^k} \left[ \sum_{i=1}^k \sum_{\boldsymbol{x} \in \mathcal{C}_i} \|\boldsymbol{x} - \boldsymbol{\mu}_i\|^2 \right]$$

$$oldsymbol{\mu}_i = rac{1}{|\mathcal{C}_i|} \sum_{oldsymbol{x}' \in \mathcal{C}_i} oldsymbol{x}'$$

 $\mathcal{C}_i \cap \mathcal{C}_j = \phi$ 

# Within-Cluster Scatter Minimization

$$\min_{\{\mathcal{C}_i\}_{i=1}^k} \left[ \sum_{i=1}^k \sum_{oldsymbol{x} \in \mathcal{C}_i} \|oldsymbol{x} - oldsymbol{\mu}_i\|^2 
ight]$$

- When all possible cluster assignment is tested in a greedy manner, computation time is proportional to  $k^n$ .
- Actually, the above optimization problem is NP-hard, i.e., we do not yet have a polynomial-time algorithm.

# K-Means Clustering Algorithm<sup>132</sup>

- Randomly initialize partition:  $\{C_i\}_{i=1}^k$
- Repeat the following until convergence:

• Update cluster assignment:  $j = 1, 2, \dots, n$ 

$$\boldsymbol{x}_j 
ightarrow \mathcal{C}_{t_j} \qquad t_j = \operatorname*{argmin}_i \| \boldsymbol{x}_j - \boldsymbol{\mu}_i \|^2$$

• Update cluster centroids:  $i = 1, 2, \dots, k$ 

$$oldsymbol{\mu}_i = rac{1}{|\mathcal{C}_i|} \sum_{oldsymbol{x}' \in \mathcal{C}_i} oldsymbol{x}'$$

Note: Only local optimality is guaranteed

Examples





K-means method can successfully separate the two data crowds from each other.

#### Examples (cont.)

134



# However, it does not work well if the data crowds have non-convex shapes.

# Non-Linearizing K-Means <sup>135</sup>

Map the original data to a feature space by a non-linear transformation:

$$\phi: \boldsymbol{x} \to \boldsymbol{f} \qquad \{\boldsymbol{f}_i \mid \boldsymbol{f}_i = \phi(\boldsymbol{x}_i)\}_{i=1}^n$$

Run the k-means algorithm in the feature space.

$$\min_{\{\mathcal{C}_i\}_{i=1}^k} \left[ \sum_{i=1}^k \sum_{\boldsymbol{x} \in \mathcal{C}_i} \|\phi(\boldsymbol{x}) - \boldsymbol{\mu}_i\|^2 \right]$$

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

### Kernel K-Means Algorithm <sup>136</sup>

Randomly initialize partition:  $\{C_j\}_{j=1}^k$ 

Update cluster assignments until convergence:

 $\boldsymbol{x}_j \to \mathcal{C}_{t_j}$   $j = 1, 2, \dots, n$ 

$$t_j = \underset{i}{\operatorname{argmin}} \left[ -\frac{2}{|\mathcal{C}_i|} \sum_{\boldsymbol{x}' \in \mathcal{C}_i} K(\boldsymbol{x}_j, \boldsymbol{x}') + \frac{1}{|\mathcal{C}_i|^2} \sum_{\boldsymbol{x}', \boldsymbol{x}'' \in \mathcal{C}_i} K(\boldsymbol{x}', \boldsymbol{x}'') \right]$$

$$\begin{split} \|\phi(\boldsymbol{x}) - \boldsymbol{\mu}_i\|^2 &= \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}) \rangle - 2 \langle \phi(\boldsymbol{x}), \boldsymbol{\mu}_i \rangle + \langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_i \rangle \\ &= K(\boldsymbol{x}, \boldsymbol{x}) - \frac{2}{|\mathcal{C}_i|} \sum_{\boldsymbol{x}' \in \mathcal{C}_i} K(\boldsymbol{x}, \boldsymbol{x}') + \frac{1}{|\mathcal{C}_i|^2} \sum_{\boldsymbol{x}', \boldsymbol{x}'' \in \mathcal{C}_i} K(\boldsymbol{x}', \boldsymbol{x}'') \\ & \quad \text{constant} \end{split}$$



Kernel k-means method can separate the two data crowds successfully.

# Examples of Kernel K-Means (coht.)

$$K(x, x') = \exp(-||x - x'||^2/c^2)$$



#### It also works well for data with nonconvex shapes.

# Examples of Kernel K-Means (coh?)

$$K(x, x') = \exp(-||x - x'||^2/c^2)$$



- Choice of kernels (type and parameter) depends on the result.
- Appropriately choosing kernels is not easy in practice.

## Examples of Kernel K-Means (colff.)

$$K(x, x') = \exp(-||x - x'||^2/c^2)$$



Solution depends crucially on the initial cluster assignments since clustering is carried out in a high-dimensional feature space.

# Weighted Scatter Criterion <sup>141</sup>

We assign a positive weight d(x) for each sample x:

 $\min_{\{\mathcal{C}_i\}_{i=1}^k} \left[ J_{WS} \right]$ 

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

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

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



#### Prove that

$$\underset{i}{\operatorname{argmin}} \left[ d(\boldsymbol{x}) \| \phi(\boldsymbol{x}) - \boldsymbol{\mu}_i \|^2 \right]$$

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

is equivalent to

$$\underset{i}{\operatorname{argmin}} \left[ -\frac{2}{s_i} \sum_{\boldsymbol{x}' \in \mathcal{C}_i} d(\boldsymbol{x}') K(\boldsymbol{x}_j, \boldsymbol{x}') \right]$$

$$+\frac{1}{s_i^2}\sum_{\boldsymbol{x}',\boldsymbol{x}''\in\mathcal{C}_i}d(\boldsymbol{x}')d(\boldsymbol{x}'')K(\boldsymbol{x}',\boldsymbol{x}'')\Big]$$

#### Proof



independent of i

# Weighted Kernel K-Means <sup>144</sup>

- Randomly initialize partition:  $\{C_i\}_{i=1}^k$
- Update cluster assignments until convergence:

 $x_j o \mathcal{C}_t$ 

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

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

# **Hierarchical Clustering**

Hierarchical cluster structure can be obtained recursively clustering the data.

Perhaps we may fix k=2.



145

# Homework

146

Implement linear/kernel k-means algorithms and reproduce the 2-dimensional examples shown in the class.

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



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

# 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

148

- 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

 June 9<sup>th</sup>: Regular lecture (spectral clustering)
 June 16<sup>th</sup>: Preparation for the mini-conference (no lecture)

June 23<sup>rd</sup>: Regular lecture (projection pursuit 1)

June 30<sup>th</sup>: Regular lecture (projection pursuit 2)

- July 7<sup>th</sup>: Preparation for the mini-conference (no lecture)
- July 14<sup>th</sup>: Mini-conference (day 1)

July 21<sup>st</sup>: Mini-conference (day 2 if necessary)

# Mini-Conference on Data Analysis

- Application procedure: On June 23<sup>rd</sup>, 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!).