

Dissimilarities Approximation Issues

D. Fortin INRIA-Rocquencourt

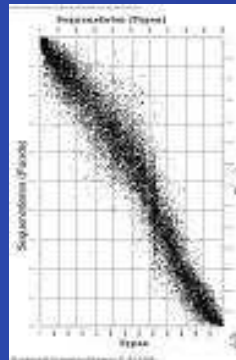
P. Pr ea LIF-Marseille

Plan

- Dissimilarities $D : d_{ij} \geq 0, d_{ii} = 0$
- Problem Taxonomy
 - Seriation
 - Phylogeny
 - Clustering
- Bottleneck Properties
- Approximations

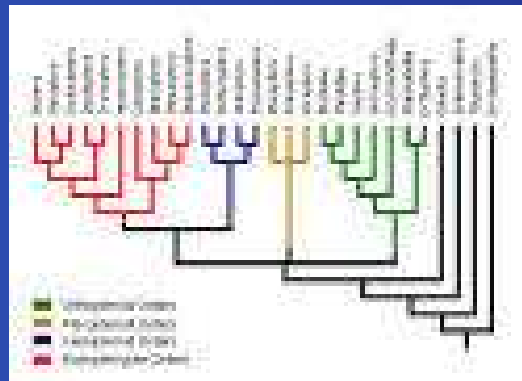
Seriation

- $D \in \mathbb{R}^{n \times n} \mapsto$ hamiltonian path
- assumption *innovation* propagates along shortest path



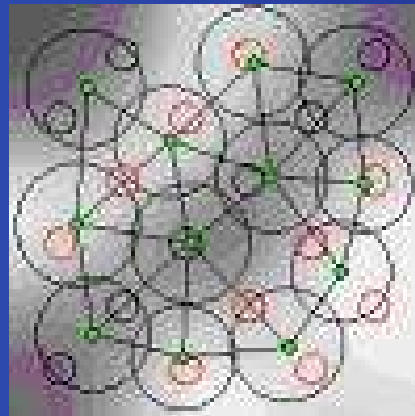
Phylogeny

- $D \in \mathbb{R}^{n \times n} \mapsto$ tree metrics
- hereditary relationship

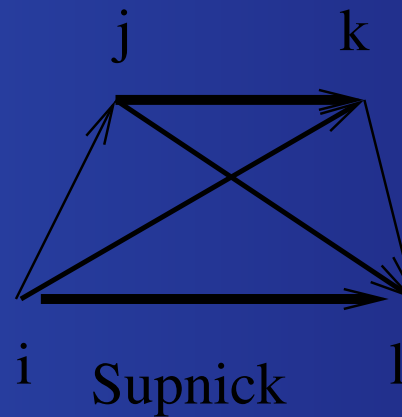
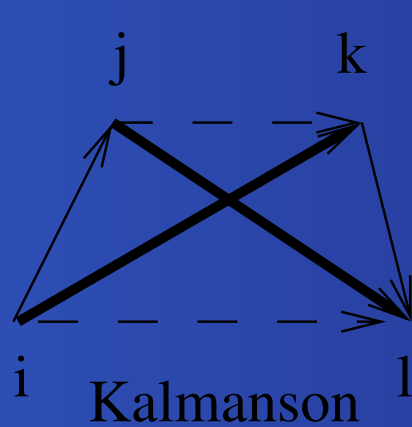
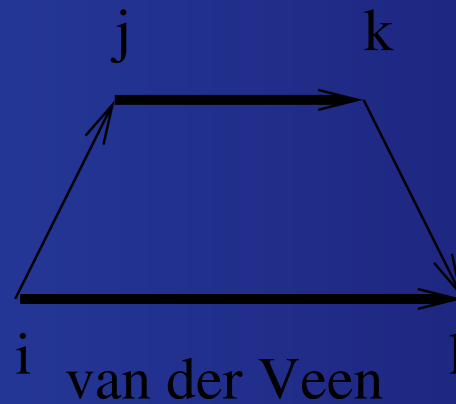
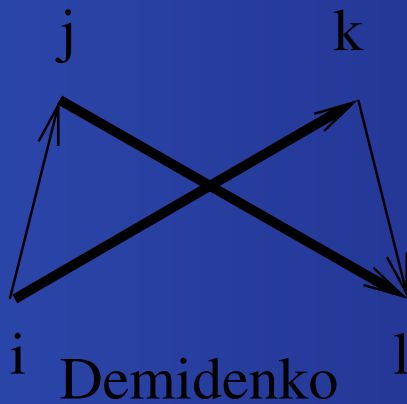
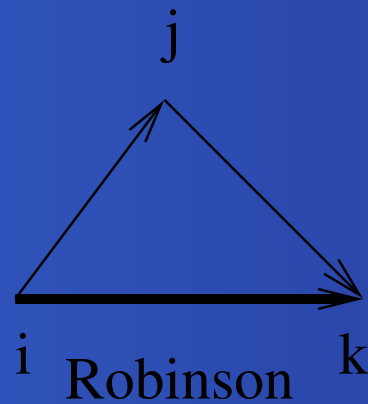


Clustering

- $D \in \mathbb{R}^{n \times n} \mapsto$ graph
- classes are dense (cliques)
- classes are separable (stables)



Bottleneck Properties



Easy Recognition

- ultrametrics $d_{ik} \leq \max(d_{ij}, d_{jk}) \quad o(n^2 \log n)$
- treemetrics (phylogeny)
 $d_{ij} + d_{kl} \leq d_{ik} + d_{jl} = d_{il} + d_{jk} \quad o(n^2 \log n)$
- Robinson (seriation, clustering)
 $d_{ik} \geq \max(d_{ij}, d_{jk}), i < j < k \quad o(n^2 \log n)$

Characterization $i < j < k < l$

- Kalmanson (sum case: master tour)

$$d_{ik} \oplus d_{jl} \geq \begin{cases} d_{ij} \oplus d_{kl} \\ d_{il} \oplus d_{jk} \end{cases}$$

- Supnick (sum case: pyramidal tour)

$$d_{ij} \oplus d_{kl} \leq d_{ik} \oplus d_{jl} \leq d_{il} \oplus d_{jk}$$

$(1, 3, 5 \dots n, n - 1, \dots 6, 4, 2)$



Conic Approximations

- definitions

$\text{vec}(\cdot) : \mathbf{R}^{n \times n} \mapsto \mathbf{R}^{n^2}$ (matrix vectorization)

$\text{Mat}(\cdot) = \text{vec}(\cdot)^{-1}$ (vector matrixization)

$\langle X, Y \rangle = \langle \text{vec}(X), \text{vec}(Y) \rangle$ (inner product)

$\mathcal{K}_{\hat{X}}^0 = \{Y \mid \langle Y, X - \hat{X} \rangle \leq 0 \text{ for all } X \in \mathcal{K}\}$ (normal cone)

$\mathcal{K}^\circ = \{Y, y \in \mathbf{R} \mid \langle X, Y \rangle \leq y \text{ for all } X \in \mathcal{K}\}$ (polar cone)

$\mathcal{K} \text{ cone} \Rightarrow \mathcal{K}^\circ = \mathcal{K}_0^0$

Dissimilarity Modeling

- $\min \max \mapsto \| \cdot \|_\infty$
- $\min \sum \mapsto \| \cdot \|_1$
- subdominant $\min(\cdot, \cdot, \dots, \cdot)$ multicriterion
 - $\min \max \text{vec}(X_1), \min \langle E, X_1 \rangle, \min \text{vec}(X_1)$
 - s.t. $X_1 + X_2 = D$
 - $X_1 \in \bar{\mathcal{N}}, X_2 \in \bar{\mathcal{K}} = \bar{\mathcal{N}} \cap \mathcal{K}$
- $\text{vec}([X_1, X_2]) = \begin{bmatrix} \text{vec}(X_1) \\ \text{vec}(X_2) \end{bmatrix}$

Dissimilarity Approximations

- primal

$$\begin{aligned} \min \quad & \langle [E, O], X \rangle \\ \text{s.t.} \quad & [I_{n^2}, I_{n^2}] \text{vec}(X) = D \\ & X \in \overline{\mathcal{N}} \times \overline{\mathcal{K}} \end{aligned}$$

- dual

$$\begin{aligned} \max \quad & \langle D, Y_1 \rangle \\ \text{s.t.} \quad & Y_1 + Y_2 = E \\ & Y_2 \in (\overline{\mathcal{N}} \times \overline{\mathcal{K}})^0 \end{aligned}$$

- strong duality (primal=dual) iff $\text{int}(\overline{\mathcal{N}} \times \overline{\mathcal{K}})^0 \neq \emptyset$

Clustering Modeling

- $x_{ic} = 1$ if $i \in c$, 0 otherwise
- intracluster $\min \sum_{i,j,c} d_{ij} x_{ic} x_{jc}$
- intercluster $\max \sum_{i,j,c} d_{ij} x_{ic} (1 - x_{jc})$
- overlapping $1 \leq \sum_c x_{ic} \leq k$
 - partition $k = 1$
 - hierarchy $k = 2 \quad \supset (D \text{ is Robinson})$

Clustering Approximation

- $(\min \langle DX, X \rangle, \max \langle DX, E - X \rangle)$
- $e \leq Xe \leq ke, X$ 0-1 matrix

$$\left\{ \begin{array}{l} \min \langle DX, Y \rangle \\ e \leq Ye \leq ke \\ Y = X \end{array} \right. \quad \left\{ \begin{array}{l} \min \langle -DE, X \rangle \\ e \leq Xe \leq ke \end{array} \right.$$

- lagrangian decomposition \simeq bimatrix game

Intracluster Approximation schemes

- k -Clustering $\min_c \sum_{ij \in c} d_{ij}$
- k -Median $\min_c \sum_{i \in c} d_{ir_c}$ to representative r_c
- k -Center $\min_c \max_{i \in c} d_{ir_c}$ to representative r_c
- metrics (facility location \rightarrow Guha)
- \mathbf{R}^d : $d(x, y) = \|x - y\|^2$ (\rightarrow Kenyon)

