

# The Minimum Transfer Cost Principle for Model-Order Selection

*Mario Frank*

*Morteza Haghir Chehreghani*

*Joachim M. Buhmann*

Department of Computer Science, ETH Zurich



# Outline

- The minimum Transfer Cost Principle
- Model order selection for truncated SVD
- Model order selection for correlation clustering
- Transfer costs for  $k$ -means clustering
- Conclusion

# Motivation

- Given:
  - A set of  $N$  objects with the measurements  $\mathbf{X}$ .
  - Two sets of objects  $\mathbf{O}^{(m)}, m \in \{1,2\}$  with corresponding measurements  $\mathbf{X}^{(m)}$  generated from the same source.
  - A data **model** characterized by a cost function  $R(\mathbf{s}, \mathbf{X}, k)$ , where the solution  $\mathbf{s}$  incorporates all relevant parameters.
- Question:
  - **What is the appropriate model order  $k$ ?**

# Transfer Costs: Intuition

- Cross-validation:
  - Good choice of the model-order based on a given dataset should also yield low costs on a second dataset.
- How to transfer a solution from  $\{\mathbf{O}^{(1)}, \mathbf{X}^{(1)}\}$  to  $\{\mathbf{O}^{(2)}, \mathbf{X}^{(2)}\}$ ?
- Classification:
  - Class labels make mapping obsolete.
- Unsupervised learning:
  - No labels are available.

# Transfer costs for factorial models

- The cost function is: 
$$R(\mathbf{s}, \mathbf{X}, k) = \sum_{i=1}^N R_i(\mathbf{s}(i), \mathbf{x}_i, k)$$

- We define an object-wise mapping function  $\psi$ :

$$\psi : \mathcal{O}^{(2)} \times \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{O}^{(1)}$$

$$\left( i', \mathbf{X}^{(1)}, \mathbf{X}^{(2)} \right) \mapsto \psi(i', \mathbf{X}^{(1)}, \mathbf{X}^{(2)})$$

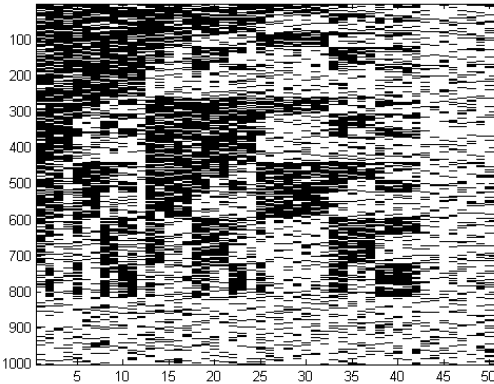
- $\psi$  aligns each object in  $\mathcal{O}^{(2)}$  with its nearest neighbor in  $\mathcal{O}^{(1)}$ .
- Transfer costs:
 
$$R^T(\mathbf{s}^{(1)}, \mathbf{X}^{(2)}, k) := \frac{1}{N_2} \sum_{i'=1}^{N_2} \sum_{i=1}^{N_1} R_{i'}(\mathbf{s}^{(1)}(i), \mathbf{x}_{i'}^{(2)}, k) \mathbb{I}_{\{\psi(i', \mathbf{X}^{(1)}, \mathbf{X}^{(2)})=i\}}$$

# The Minimum Transfer Cost Principle

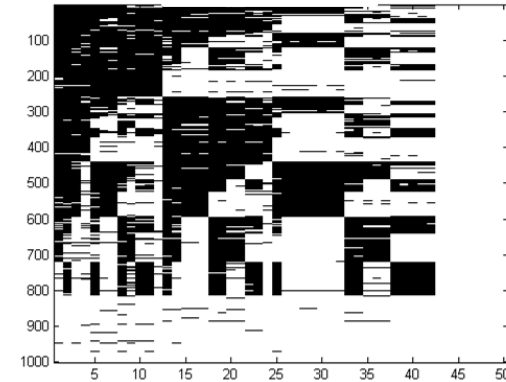
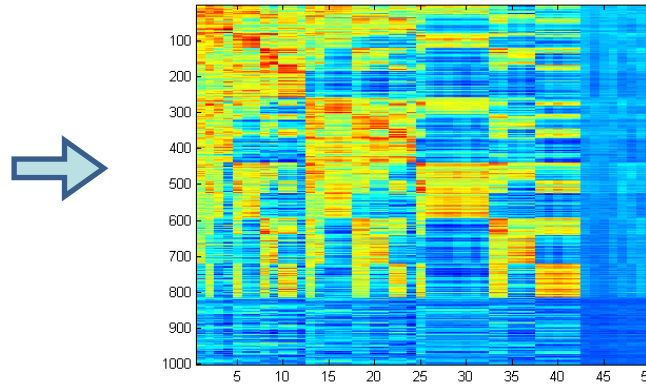
- The **minimum transfer cost** principle (MTC) selects the model order  $k$  with lowest transfer costs.
  - Too simple models underfit and achieve high costs on both datasets.
  - Too complex models overfit to the fluctuations of  $\mathbf{X}^{(1)}$  which results in high costs on  $\mathbf{X}^{(2)}$  where the fluctuations are different.

# Denoising Matrices via truncated SVD

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}$$

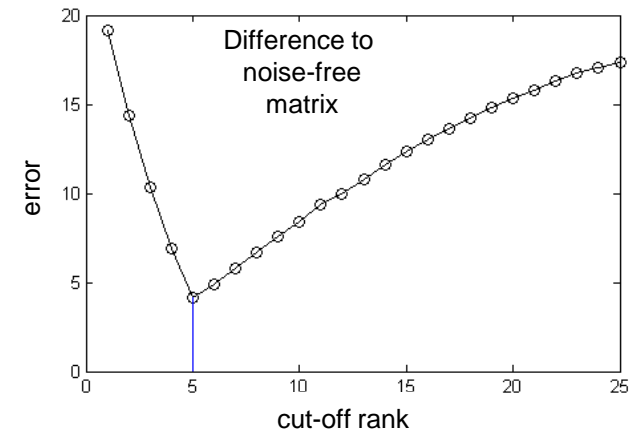


$$\mathbf{X}_5 = \mathbf{U}_5 \mathbf{S}_5 \mathbf{V}_5$$

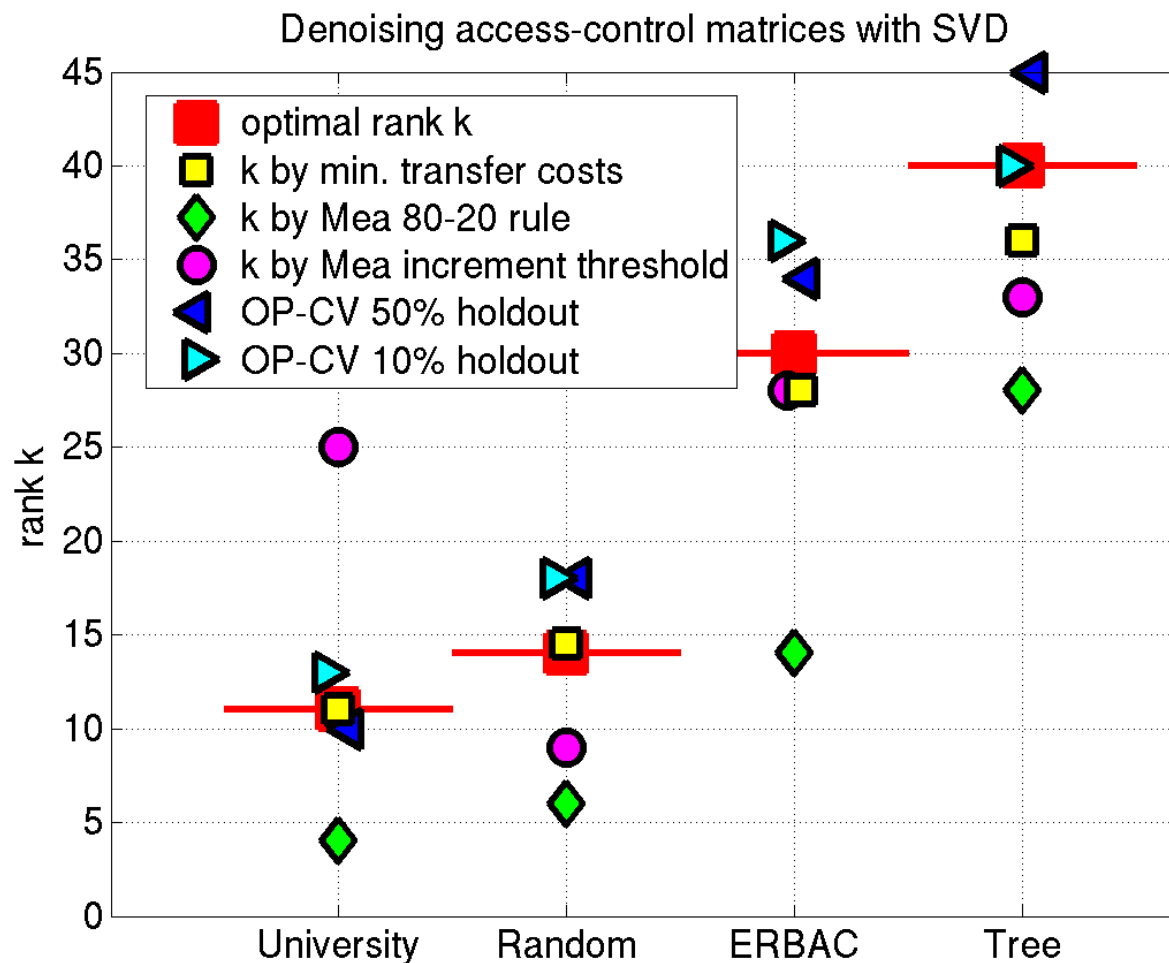


- Where to cut-off the spectrum?
- Error depends heavily on  $k$ .
- Transfer costs with nearest- $n$ . mapping:

$$R^T(\mathbf{s}, \mathbf{X}, k) = \frac{1}{N_2} \left\| \psi_{NN}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) \circ \mathbf{X}^{(2)} - \left( \mathbf{U}_k^{(1)} \mathbf{S}_k^{(1)} \mathbf{V}_k^{(1)\top} \right) \right\|_2^2$$



# Selecting the rank for truncated SVD





# Correlation Clustering

## ■ Formulation

- Given: graph  $G(\mathbf{O}, \mathbf{X})$  with similarity matrix  $\mathbf{X} := \{X_{ij}\} \in \{\pm 1\}^{\binom{N}{2}}$ .
- For the tuple  $(\mathbf{s}, \mathbf{X})$ , the costs are

$$R(\mathbf{s}, \mathbf{X}, k) = -\frac{1}{2} \sum_{1 \leq u \leq k} \sum_{(i,j) \in E_{u,u}} (X_{ij} - 1) + \frac{1}{2} \sum_{1 \leq u \leq k} \sum_{1 \leq v < u} \sum_{(i,j) \in E_{u,v}} (X_{ij} + 1)$$

where,  $E_{u,v} = \{(i,j) \in E : \mathbf{s}(i) = u \wedge \mathbf{s}(j) = v\}$ .

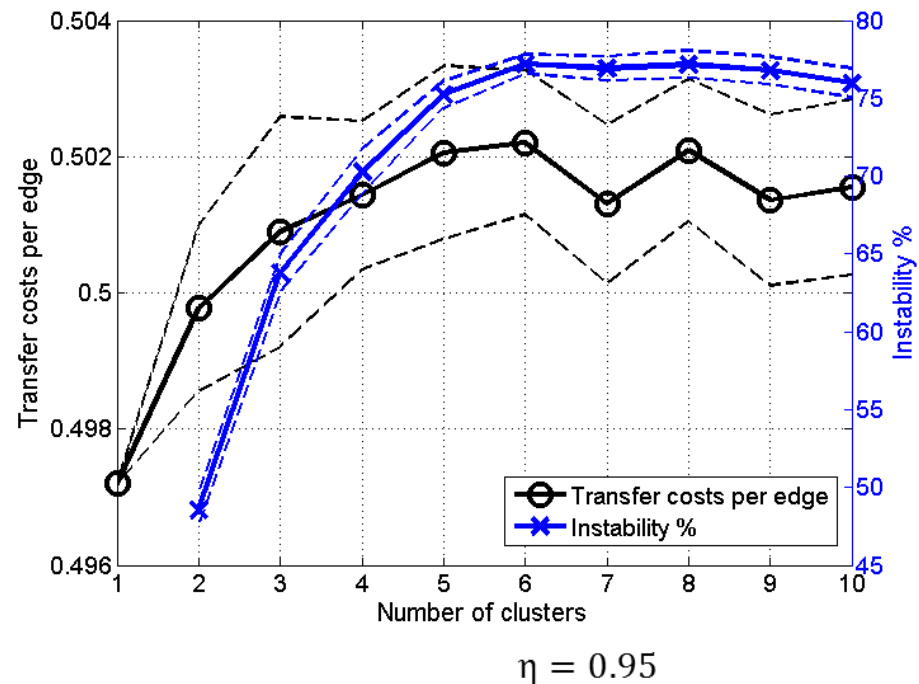
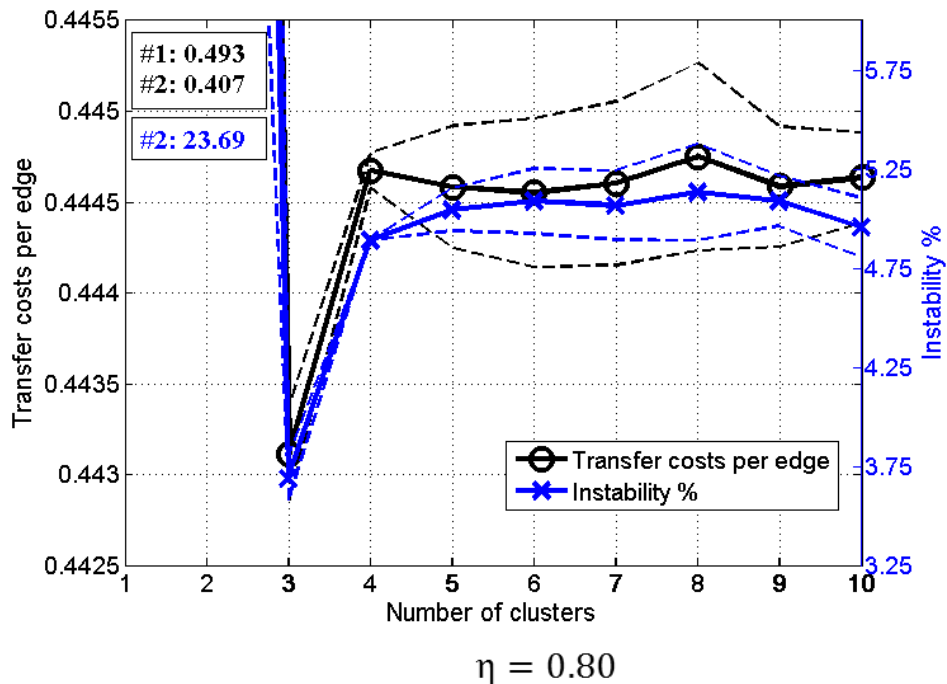
- The cluster index of object  $i'$  from  $\mathbf{O}^{(2)}$  is determined by

$$\mathbf{s}^{(1)}(i') = \arg \min_{1 \leq v \leq k} H(i', \mathbf{s}_v^{(1)}), \text{ with}$$

$$H(i', \mathbf{s}_v^{(1)}) = -\frac{1}{2} \sum_{j \in \mathbf{s}_v} (X_{ij} - 1) + \frac{1}{2} \sum_{1 \leq u \leq k, u \neq v} \sum_{j \in \mathbf{s}_u} (X_{ij} + 1)$$

# Finding the number of clusters

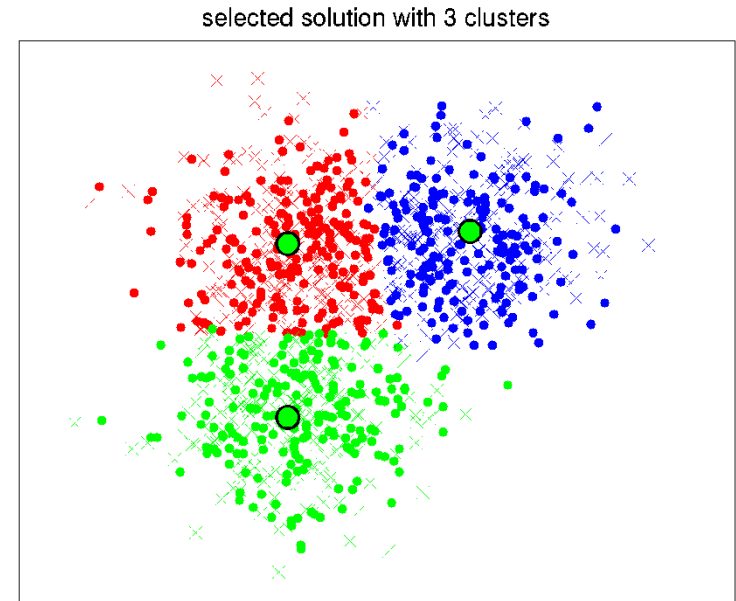
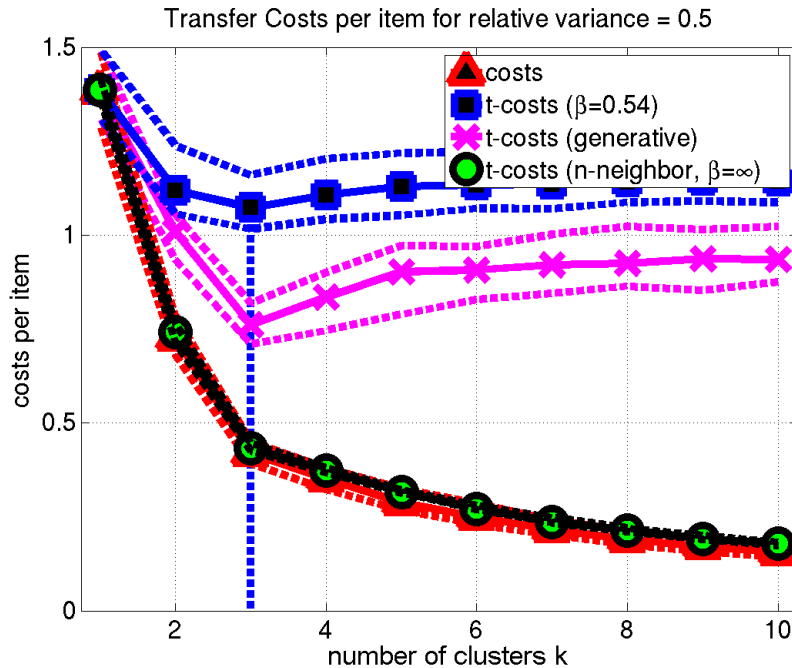
- Generate correlation data with varying noise level  $\eta$



# Transfer costs for k-means

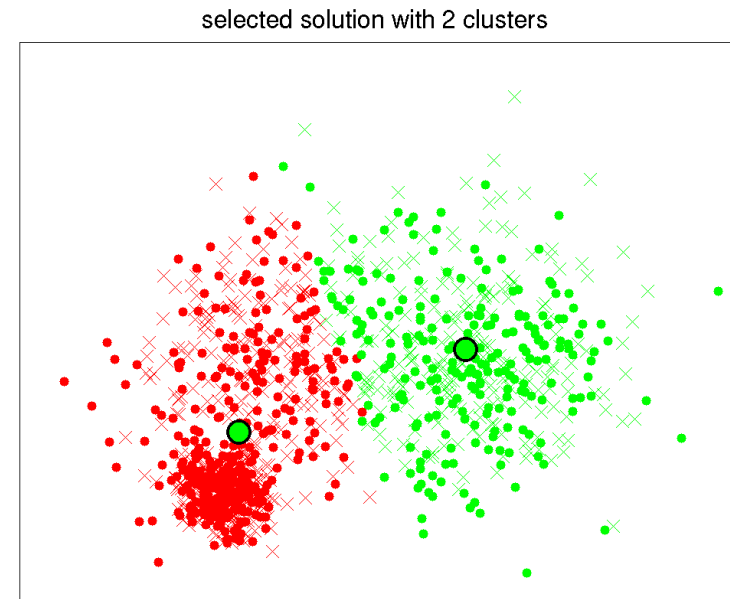
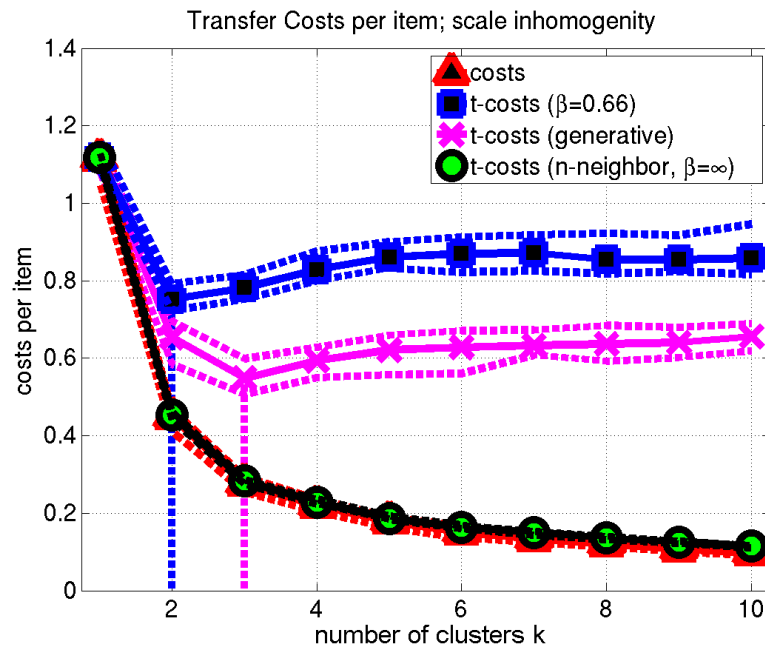
- Goal: challenge the criterion with a **model mismatch**.
- Experiment:
  - Generate two datasets from a **mixture of Gaussians**.
  - Cluster with **k-means**.
  - Select number of clusters with minimum transfer costs.
- Model mismatch: **no variance for k-means** → vector quantization preferred.

# k-means: easy case



- Discrete mapping leads to **monotonically decreasing t-costs!**
- **Soft mapping** selects true number of clusters.

# k-means: scale inhomogeneity



- Soft mapping temperature serves as a **'global' variance**.
- Gaussian mixture models are required if variance changes locally.

# Conclusion

- MTC: an easily applicable method for model-order selection in unsupervised scenarios.
- Demonstration of model-order selection in
  - Gaussian Mixture Models
  - Truncated SVD: denoising of images and Boolean matrices
  - Boolean matrix factorization (role mining)
  - Correlation clustering
- Soft mapping for model mismatch
  - *k*-means case study