

Estimation of low-rank tensors via convex optimization

Ryota Tomioka¹, Kohei Hayashi², Hisashi Kashima¹

¹The University of Tokyo

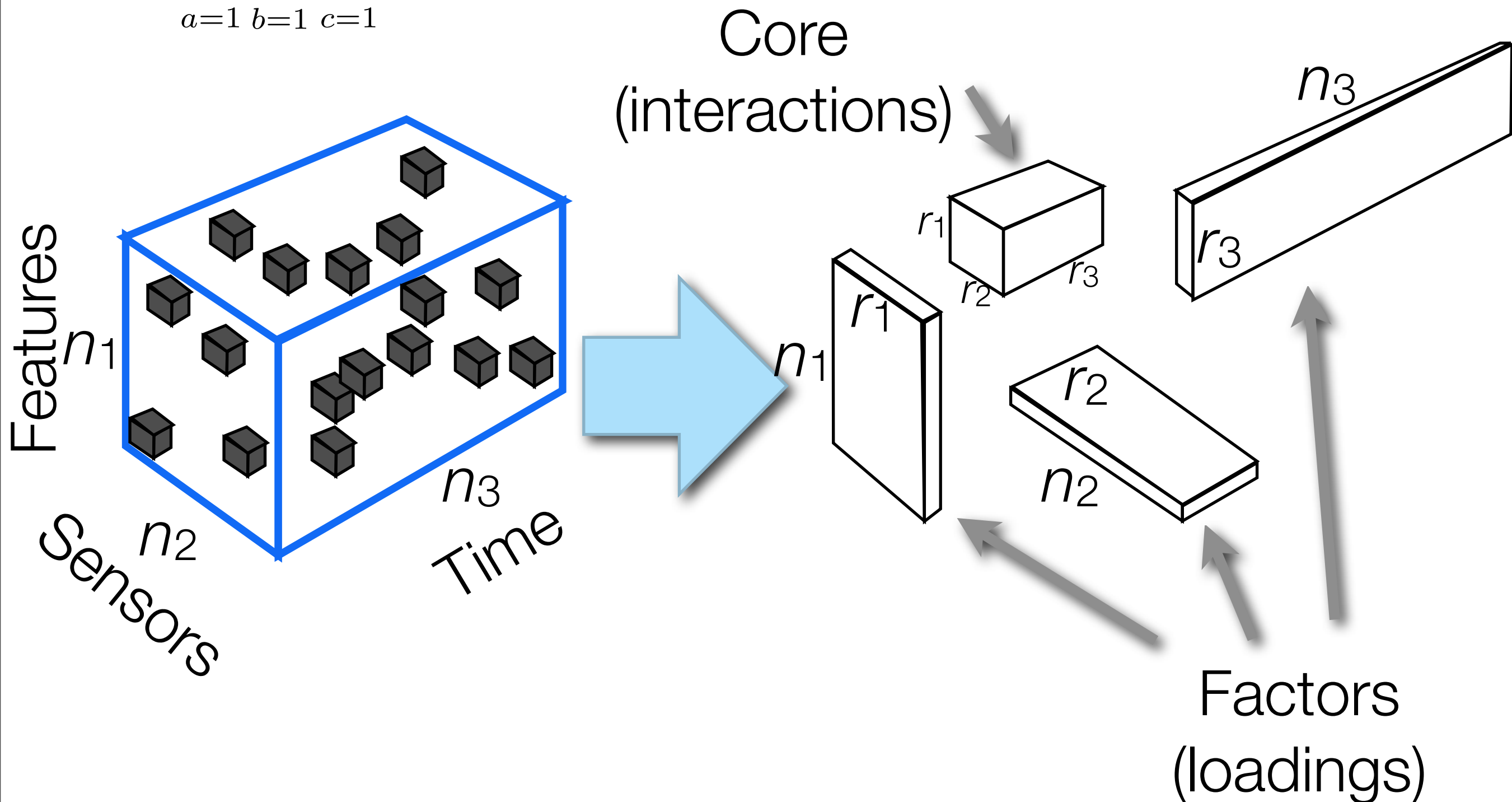
²Nara Institute of Science and Technology

2011/3/23 @ TU Berlin

Convex low-rank *tensor* completion

$$X_{ijk} = \sum_{a=1}^{r_1} \sum_{b=1}^{r_2} \sum_{c=1}^{r_3} C_{abc} U_{ia}^{(1)} U_{jb}^{(2)} U_{kc}^{(c)}$$

Tucker decomposition



Conventional formulation (nonconvex)

$$\underset{\mathcal{C}, \mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3}{\text{minimize}} \quad \|\Omega \circ (\mathcal{Y} - \mathcal{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3)\|_F^2 + \text{regularization}.$$

/ \

observation mode-k product

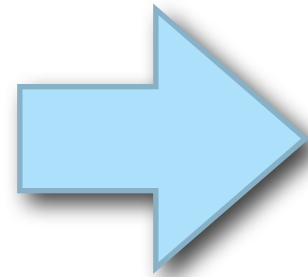
$$\underset{\mathcal{X}}{\text{minimize}} \quad \|\Omega \circ (\mathcal{Y} - \mathcal{X})\|_F^2 \quad \text{s.t.} \quad \text{rank}(\mathcal{X}) \leq (r_1, r_2, r_3).$$

- Alternate minimization
- Have to fix the rank beforehand

Our approach

Matrix

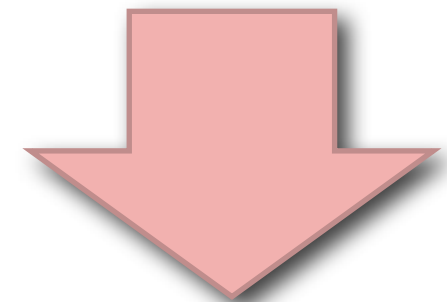
Estimation of
low-rank matrix
(hard)



Trace norm
minimization
(tractable)

[Fazel, Hindi, Boyd 01]

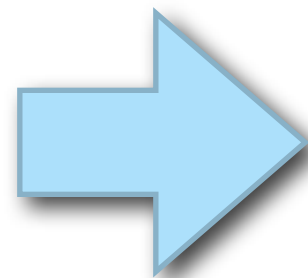
Generalization



Tensor

Estimation of
low-rank tensor
(hard)

Rank defined in the sense of
Tucker decomposition



Extended
trace norm
minimization
(tractable)

Trace norm (nuclear norm) regularization

$$\mathbf{X} \in \mathbb{R}^{R \times C}, \quad m = \min(R, C)$$

$$\|\mathbf{X}\|_* = \sum_{j=1}^m \sigma_j(\mathbf{X})$$

Linear sum of singular-values

- Roughly speaking, L1 regularization on the singular-values.
- Stronger regularization --> more **zero singular-values** --> low rank.
- Not obvious for tensors (no singular-values for tensors)

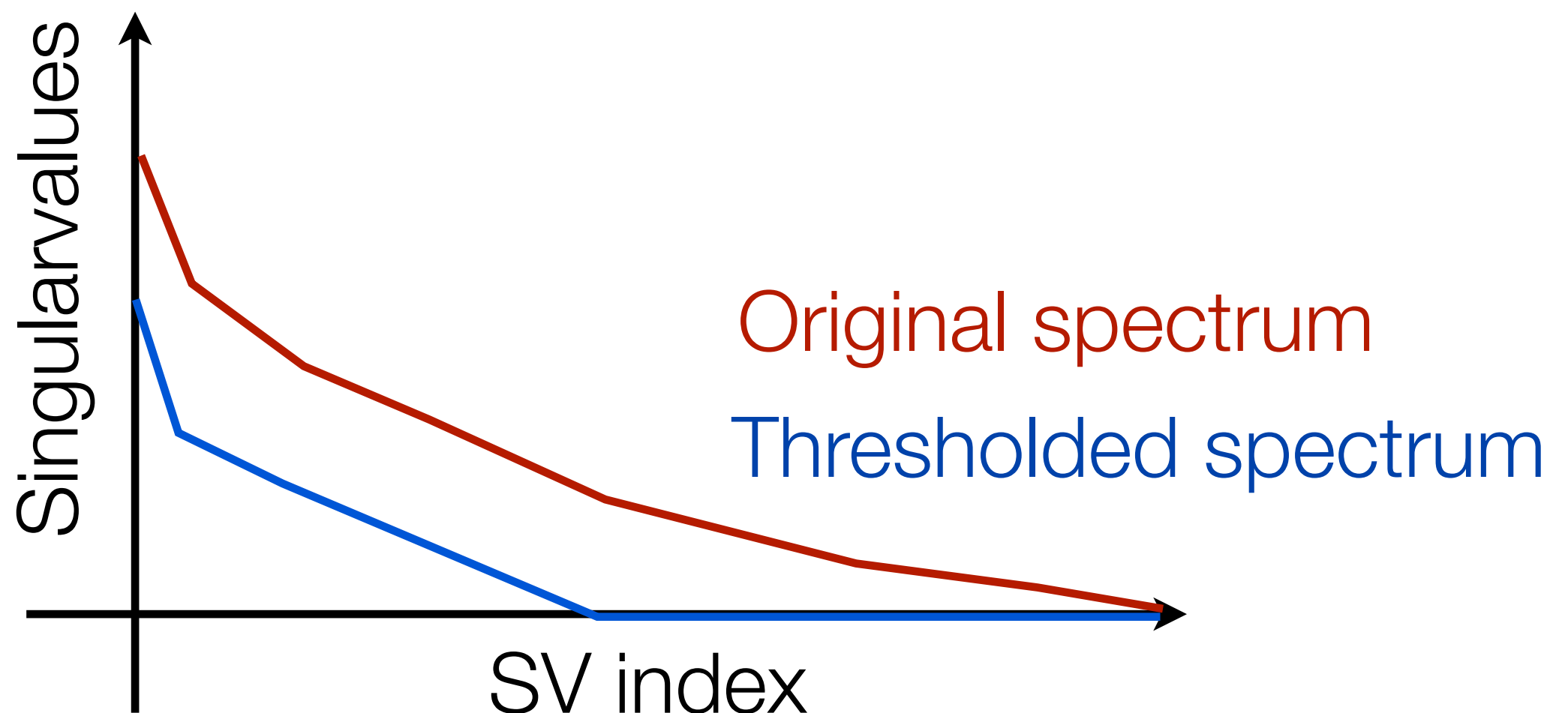
Spectral soft-threshold operation

all observed and matrix --> analytic solution

$$\text{softth}(\mathbf{X}) = \underset{\mathbf{Z} \in \mathbb{R}^{R \times C}}{\text{argmin}} \left(\frac{1}{2} \|\mathbf{Z} - \mathbf{X}\|_F^2 + \lambda \|\mathbf{Z}\|_* \right)$$

$$= \mathbf{U} \max(\mathbf{S} - \lambda, 0) \mathbf{V}^\top$$

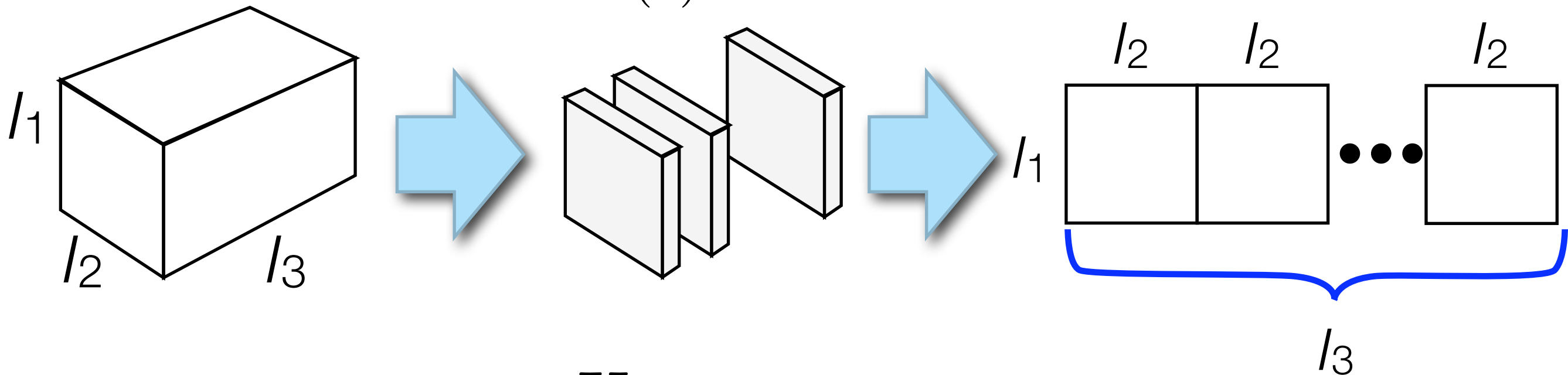
where $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$



Mode-k unfolding (matricization)

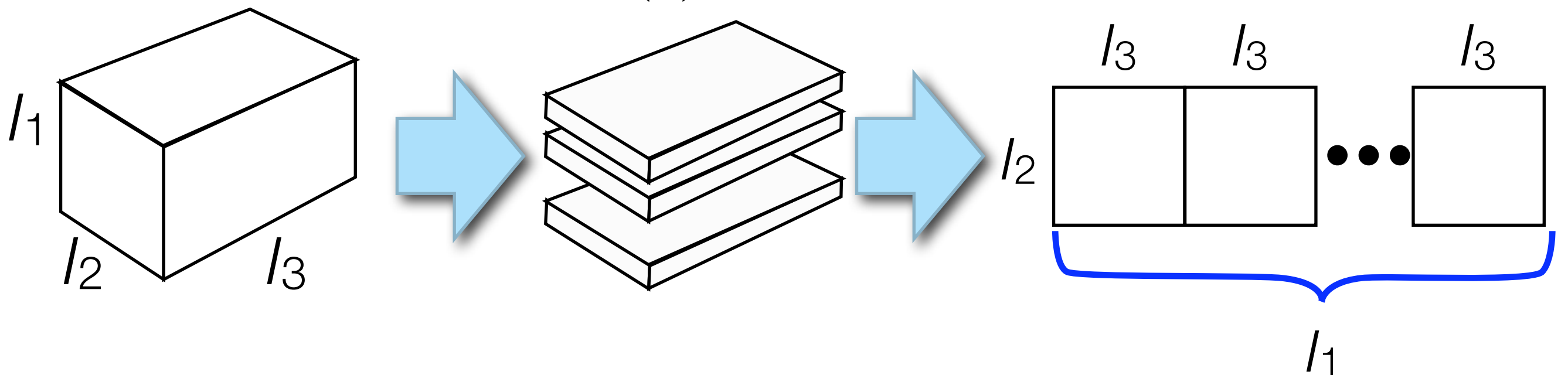
Mode-1 unfolding

$\mathbf{X}_{(1)}$



Mode-2 unfolding

$\mathbf{X}_{(2)}$



Low-rank tensor is a low-rank matrix

$$\mathcal{X} = \mathcal{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3$$

Mode-1 unfolding

$$\mathbf{X}_{(1)} = \mathbf{U}_1 \mathbf{C}_{(1)} (\mathbf{U}_3 \otimes \mathbf{U}_2)^\top$$

rank $\leq r_1$

Mode-2 unfolding

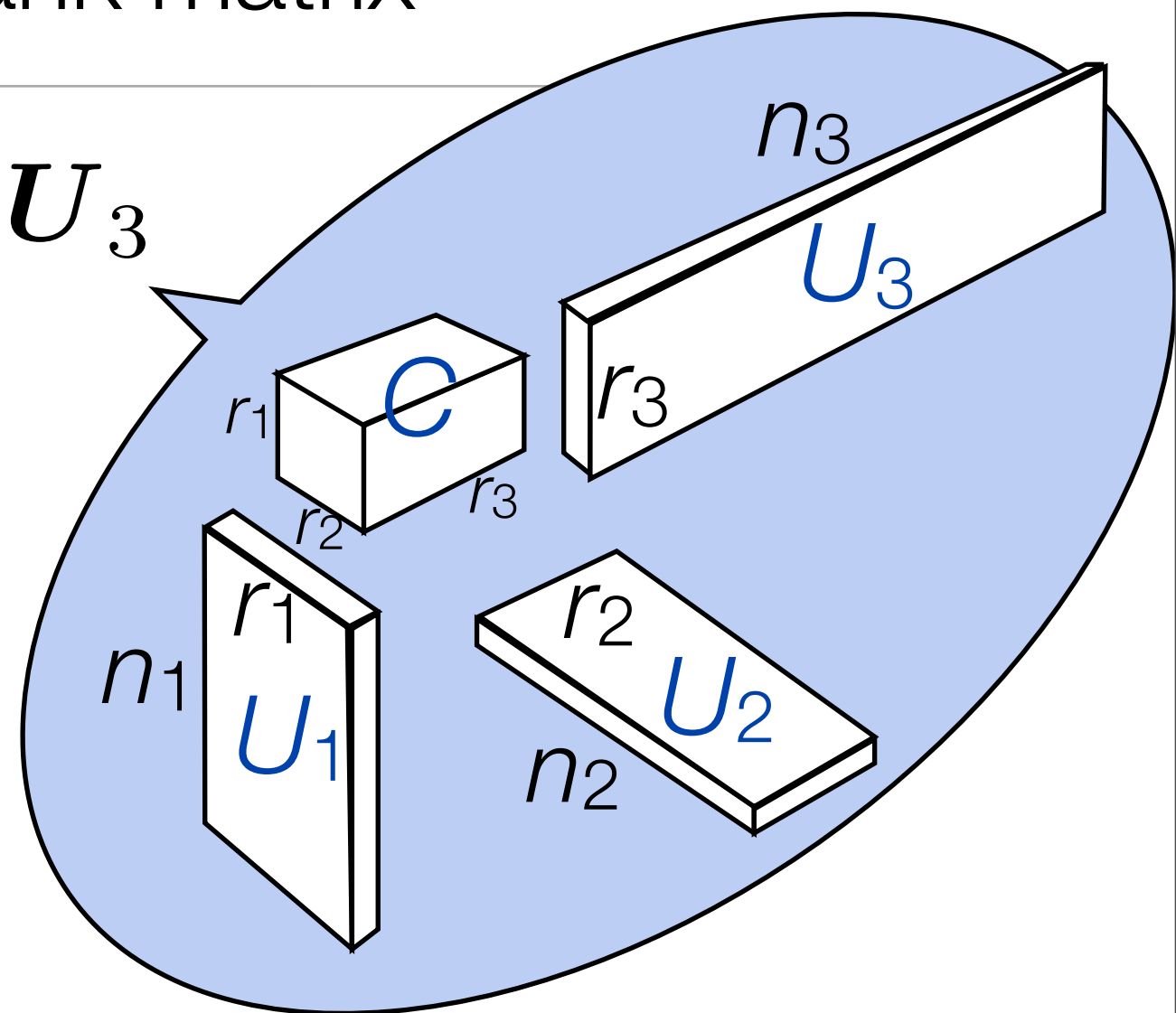
$$\mathbf{X}_{(2)} = \mathbf{U}_2 \mathbf{C}_{(2)} (\mathbf{U}_1 \otimes \mathbf{U}_3)^\top$$

rank $\leq r_2$

Mode-3 unfolding

$$\mathbf{X}_{(3)} = \mathbf{U}_3 \mathbf{C}_{(3)} (\mathbf{U}_2 \otimes \mathbf{U}_1)^\top$$

rank $\leq r_3$



The rank of $\mathbf{X}_{(k)}$ is no more than the rank of $\mathbf{C}_{(k)}$

Low-rank matrix is a low-rank tensor

- Given $X=USV^T$ (low-rank)
- Define

$$\mathcal{C} = SV^T$$

$$U_1 = U$$

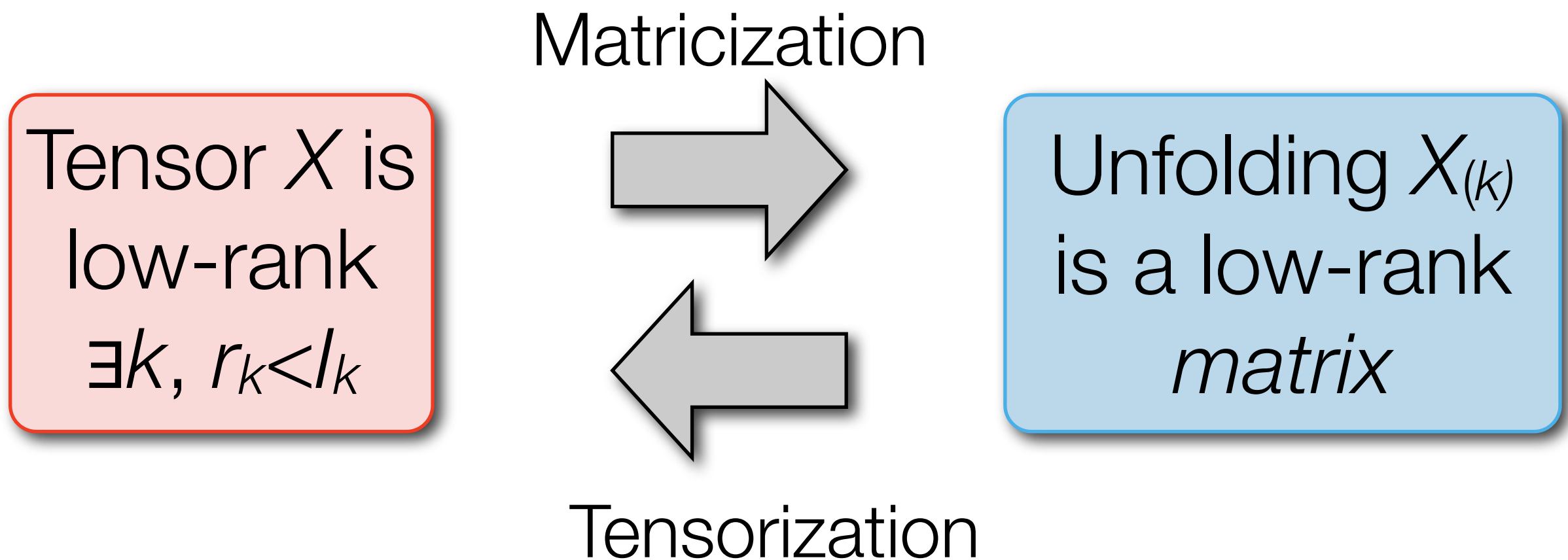
$$U_2 = I_{n_2}$$

$$U_3 = I_{n_3}$$

$\mathcal{X} = \mathcal{C} \times_1 U_1 \times_2 U_2 \times_3 U_3$ is low-rank
(at least for mode 1)

What it means

- We can use the trace norm of an unfolding of a tensor X to learn low-rank X .



Approach 1: As a matrix

- Pick a mode k , and hope that the tensor to be learned is low rank in mode k .

$$\underset{\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega \circ (\mathcal{Y} - \mathcal{X})\|_F^2 + \|\mathbf{X}_{(k)}\|_*,$$

Pro: Basically a matrix problem

→ Theoretical guarantee (Candes & Recht 09)

Con: Have to be lucky to pick the right mode.

Approach 2: Constrained optimization

- Constrain so that each unfolding of \mathbf{X} is simultaneously low rank.

$$\underset{\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega \circ (\mathcal{Y} - \mathcal{X})\|_F^2 + \sum_{k=1}^K \gamma_k \|\mathbf{X}_{(k)}\|_*.$$

Pro: Jointly regularize every mode

Con: Strong constraint

γ_k : tuning parameter usually set to 1.

(See also Signoretto et al., 10; Gandy et al. 11)

Approach 3: Mixture of low-rank tensors

- Each mixture component Z_k is regularized to be low-rank **only in mode- k** .

$$\underset{\mathcal{Z}_1, \dots, \mathcal{Z}_K}{\text{minimize}} \quad \frac{1}{2\lambda} \left\| \Omega \circ \left(\mathcal{Y} - \sum_{k=1}^K \mathcal{Z}_k \right) \right\|_F^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_{k(k)}\|_*,$$

Pro: Each Z_k takes care of each mode

Con: Sum is not low-rank

Optimization via Alternating Direction Method of Multipliers (ADMM)

(Gabay & Mercier 76)

- Useful when we have linear operation inside
sparsity penalty

$$\underset{\mathcal{X} \in \mathbb{R}^{n_1 \times \cdots \times n_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|_F^2 + \sum_{k=1}^K \gamma_k \|\mathbf{X}_{(k)}\|_*.$$



Permutation

Optimization via Alternating Direction Method of Multipliers (ADMM)

(Gabay & Mercier 76)

- Useful when we have linear operation inside sparsity penalty

$$\underset{\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|_F^2 + \sum_{k=1}^K \gamma_k \|\mathbf{X}_{(k)}\|_*.$$

↑
Permutation

- Split Bregman Iteration (Goldstein & Osher) is also an ADMM

Total-variation image reconstruction:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega(\mathbf{x}) - \mathbf{y}\|^2 + \sum_{j=1}^n \|D_j \mathbf{x}\|$$

↑
2D derivative at j th pixel

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

 Linear operation

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

 Linear operation

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

- Step 1: Split & Augment

$$\begin{aligned} &\underset{x,z}{\text{minimize}} \quad f(x) + g(z) + \frac{\eta}{2} \|Ax - z\|^2 \\ &\text{subject to} \quad z = Ax \end{aligned}$$

 Linear operation

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

- Step 1: Split & Augment

$$\begin{aligned} &\underset{x,z}{\text{minimize}} \quad f(x) + g(z) + \frac{\eta}{2} \|Ax - z\|^2 \\ &\text{subject to} \quad z = Ax \end{aligned}$$

 Linear operation

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

- Step 1: Split & Augment

$$\begin{aligned} &\underset{x,z}{\text{minimize}} \quad f(x) + g(z) + \frac{\eta}{2} \|Ax - z\|^2 \\ &\text{subject to} \quad z = Ax \end{aligned}$$

 Linear operation

ADMM preliminaries

- Problem

$$\underset{x}{\text{minimize}} \quad f(x) + g(Ax)$$

- Step 1: Split & Augment

$$\begin{aligned} &\underset{x,z}{\text{minimize}} \quad f(x) + g(z) + \frac{\eta}{2} \|Ax - z\|^2 \\ &\text{subject to} \quad z = Ax \end{aligned}$$

Linear operation



- Step 2: Augmented Lagrangian function

$$L_{\eta}(x, z, \alpha) = f(x) + g(z) + \alpha^{\top} (Ax - z) + \frac{\eta}{2} \|Ax - z\|^2$$

Ordinary Lagrangian



Augmented



ADMM algorithm (Gabay & Mercier 76)

- Minimize the AL function wrt X

$$x^{t+1} = \operatorname{argmin}_{x \in \mathbb{R}^n} L_{\eta}(x, z^t, \alpha^t),$$

- Minimize the AL function wrt Z

$$z^{t+1} = \operatorname{argmin}_{z \in \mathbb{R}^m} L_{\eta}(x^{t+1}, z, \alpha^t),$$

- Update the multiplier vector

$$\alpha^{t+1} = \alpha^t + \eta(Ax^{t+1} - z^{t+1}).$$

ADMM algorithm (Gabay & Mercier 76)

- Minimize the AL function wrt X

$$x^{t+1} = \operatorname{argmin}_{x \in \mathbb{R}^n} L_{\eta}(x, z^t, \alpha^t),$$

- Minimize the AL function wrt Z

$$z^{t+1} = \operatorname{argmin}_{z \in \mathbb{R}^m} L_{\eta}(x^{t+1}, z, \alpha^t),$$

- Update the multiplier vector

$$\alpha^{t+1} = \alpha^t + \eta(Ax^{t+1} - z^{t+1}).$$

Every limit point of ADMM is a minimizer of the original problem. [Eckstein & Bertsekas 92]

For approach “Constraint”

- Move the permutation out of the regularizer

$$\begin{aligned} & \underset{\mathcal{X}, \mathbf{Z}_1, \dots, \mathbf{Z}_K}{\text{minimize}} && \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_*, \\ & \text{subject to} && \mathbf{X}_{(k)} = \mathbf{Z}_k \quad (k = 1, \dots, K), \end{aligned}$$

- Augmented Lagrangian:

$$\begin{aligned} L_\eta(\mathcal{X}, \{\mathbf{Z}_k\}_{k=1}^K, \{\mathbf{A}_k\}_{k=1}^K) = & \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_* \\ & + \eta \sum_{k=1}^K \left(\langle \mathbf{A}_k, \mathbf{X}_{(k)} - \mathbf{Z}_k \rangle + \frac{1}{2} \|\mathbf{X}_{(k)} - \mathbf{Z}_k\|_F^2 \right). \end{aligned}$$

ADMM for “Constraint” ($\lambda \rightarrow 0$)

- Minimize the AL function wrt X

$$\begin{cases} \Omega(\mathcal{X}^{t+1}) = y & \text{(observed elem.)} \\ \bar{\Omega}(\mathcal{X}^{t+1}) = \bar{\Omega} \left(\frac{1}{K} \sum_{k=1}^K \text{tensor}_k(\mathbf{Z}_k^t - \mathbf{A}_k^t) \right) & \text{(unobserved elem.)} \end{cases}$$

- Minimize the AL function wrt Z

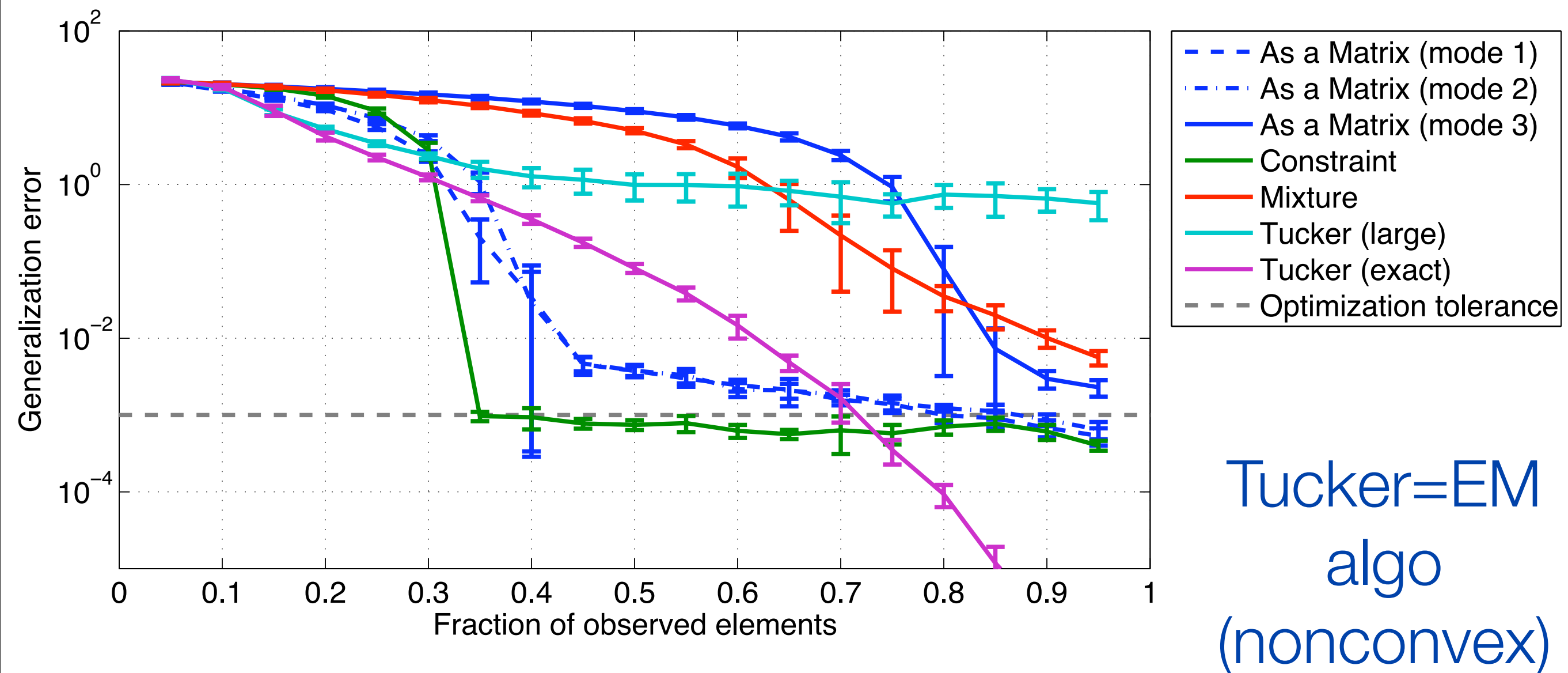
$$\mathbf{Z}_k^{t+1} = \text{softth}_{\gamma_k/\eta} \left(\mathbf{X}_{(k)}^{t+1} + \mathbf{A}_k^t \right) \quad (k = 1, \dots, K)$$

- Update multipliers

$$\mathbf{A}_k^{t+1} = \mathbf{A}_k^t + \left(\mathbf{X}_{(k)}^{t+1} - \mathbf{Z}_k^{t+1} \right) \quad (k = 1, \dots, K)$$

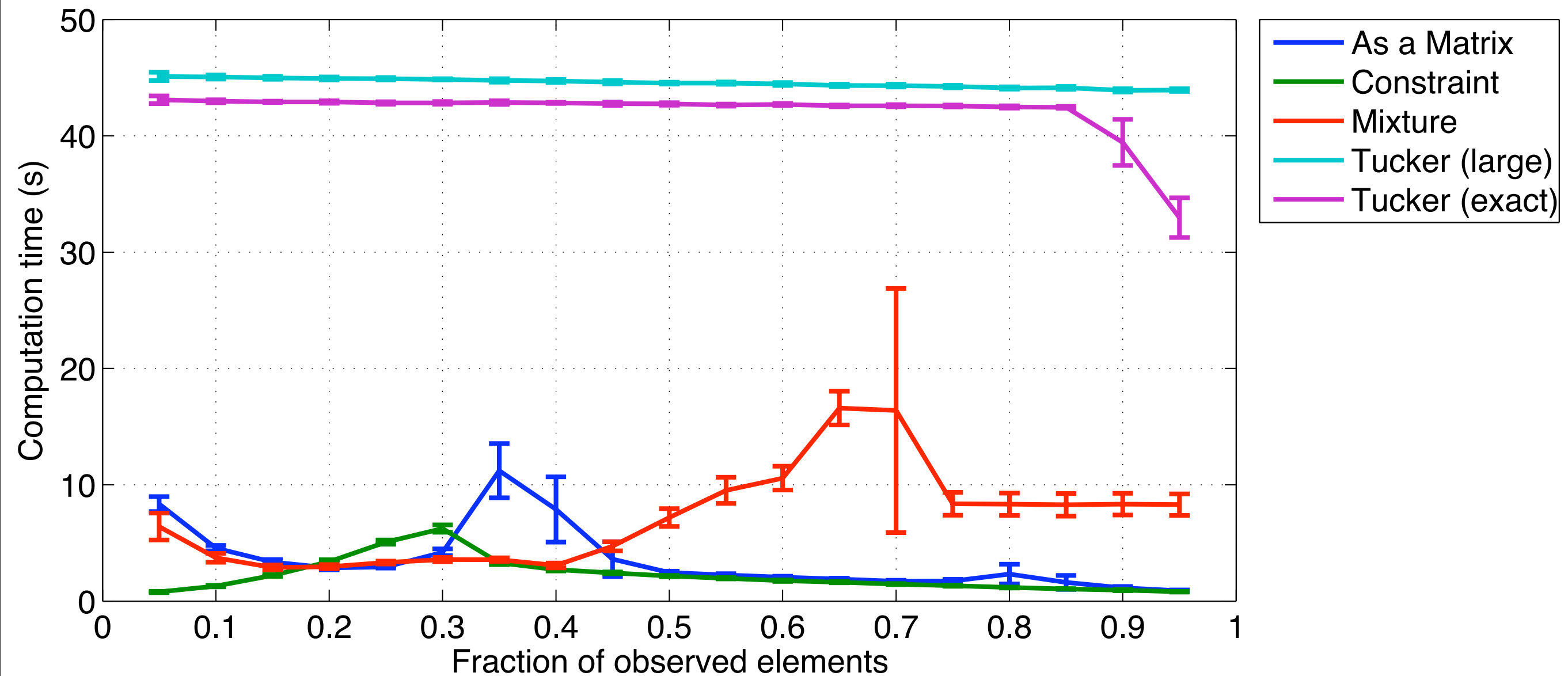
Numerical experiment

- True tensor: Size 50x50x20, rank 7x8x9. No noise ($\lambda=0$).
- Random train/test split.



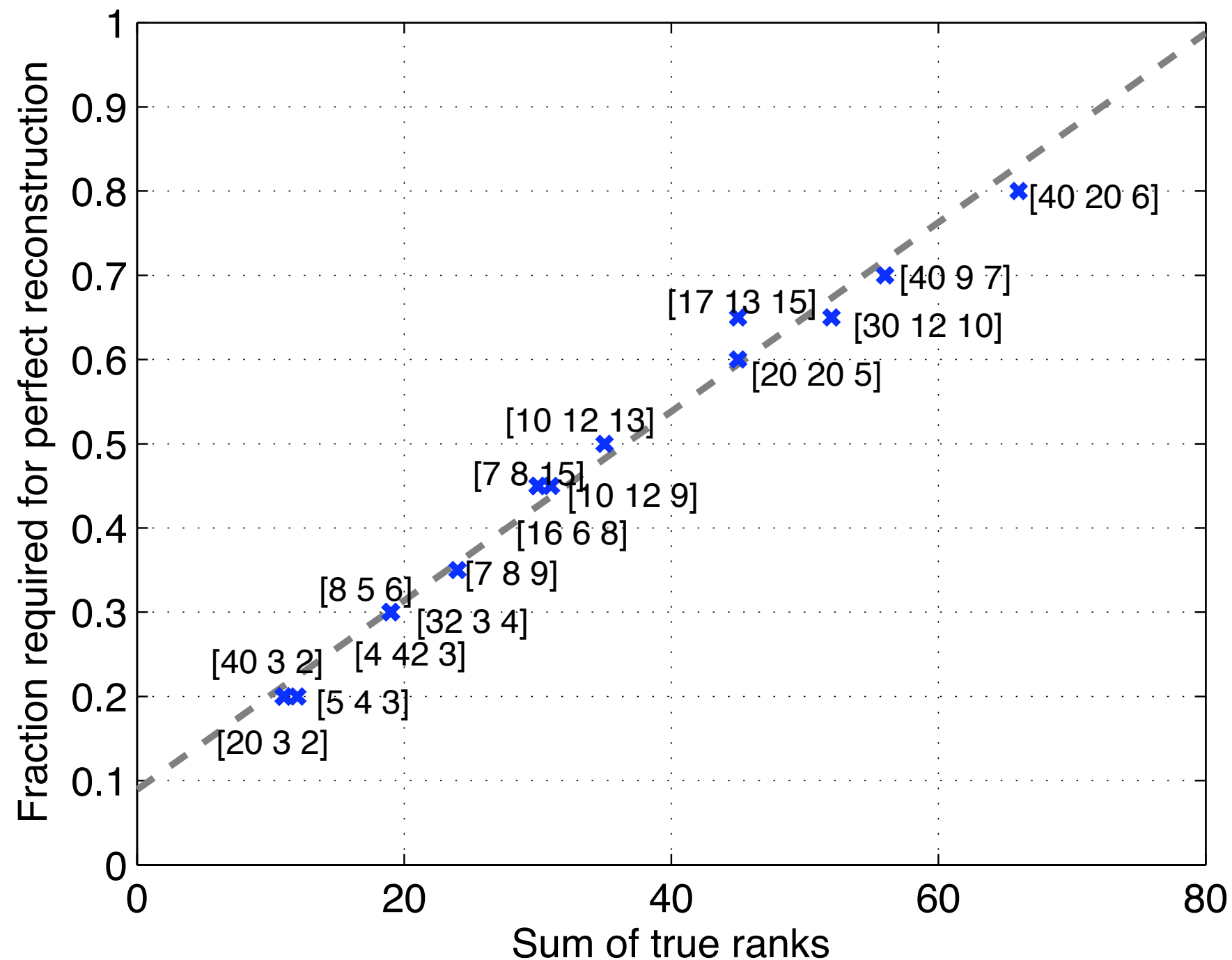
Computation time

- Convex formulation is also fast

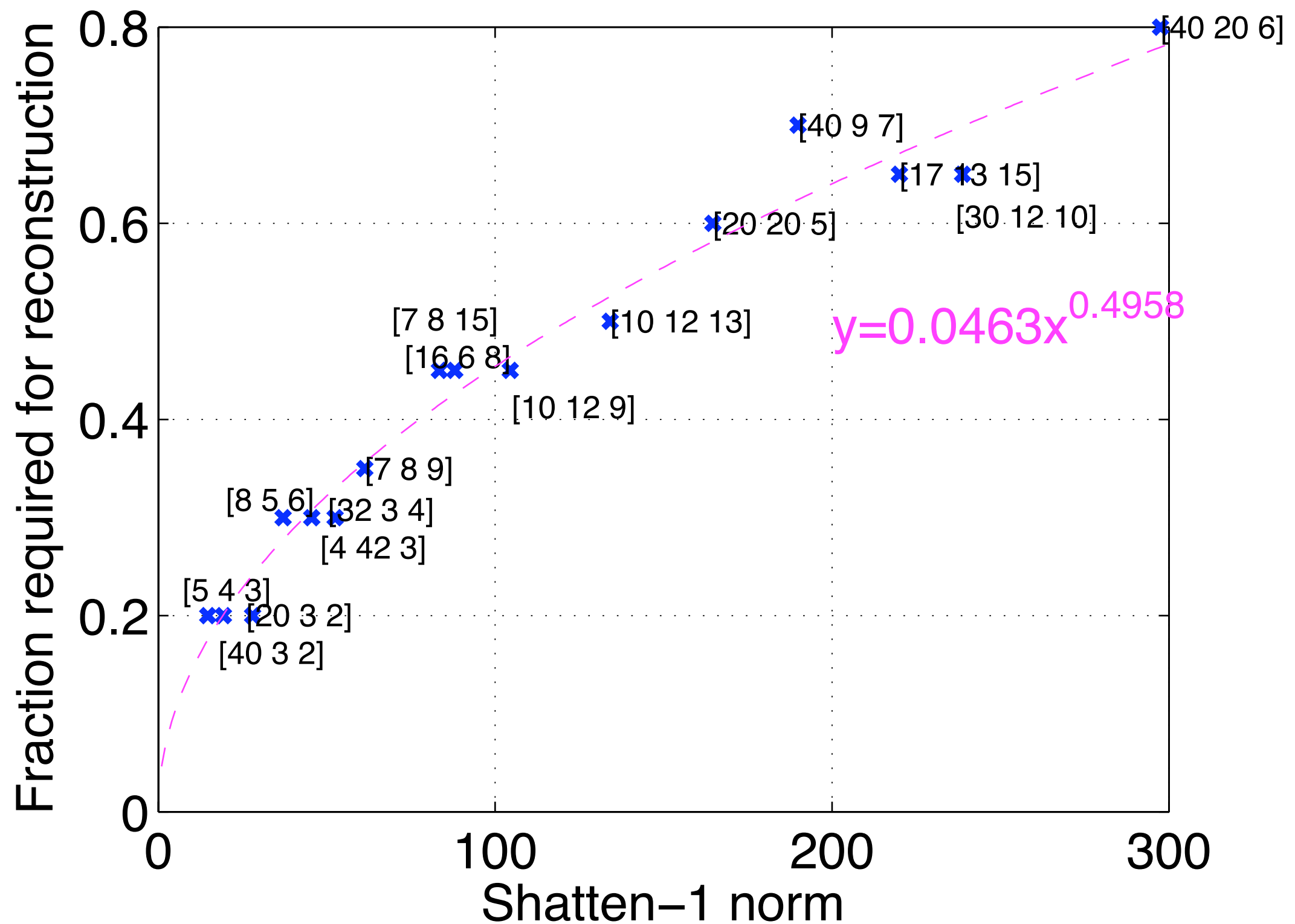


Phase transition behaviour

- Sum of true ranks = $\min(r_1, r_2, r_3) + \min(r_2, r_3, r_1) + \min(r_3, r_1, r_2)$

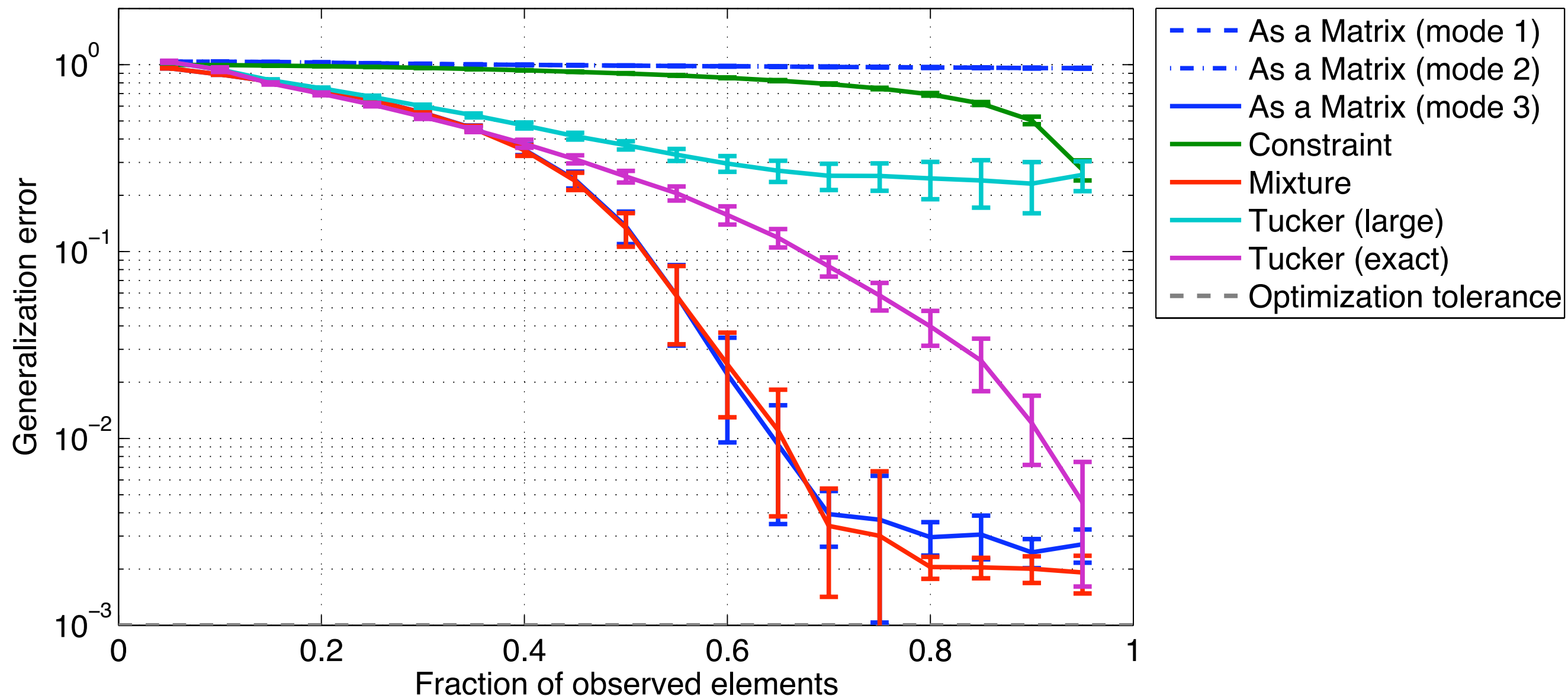


Phase transition (vs Shatten-1 norm)



“Mixture” is sometimes better

- True tensor: Size 50x50x20, rank 50x50x5. No noise ($\lambda=0$).

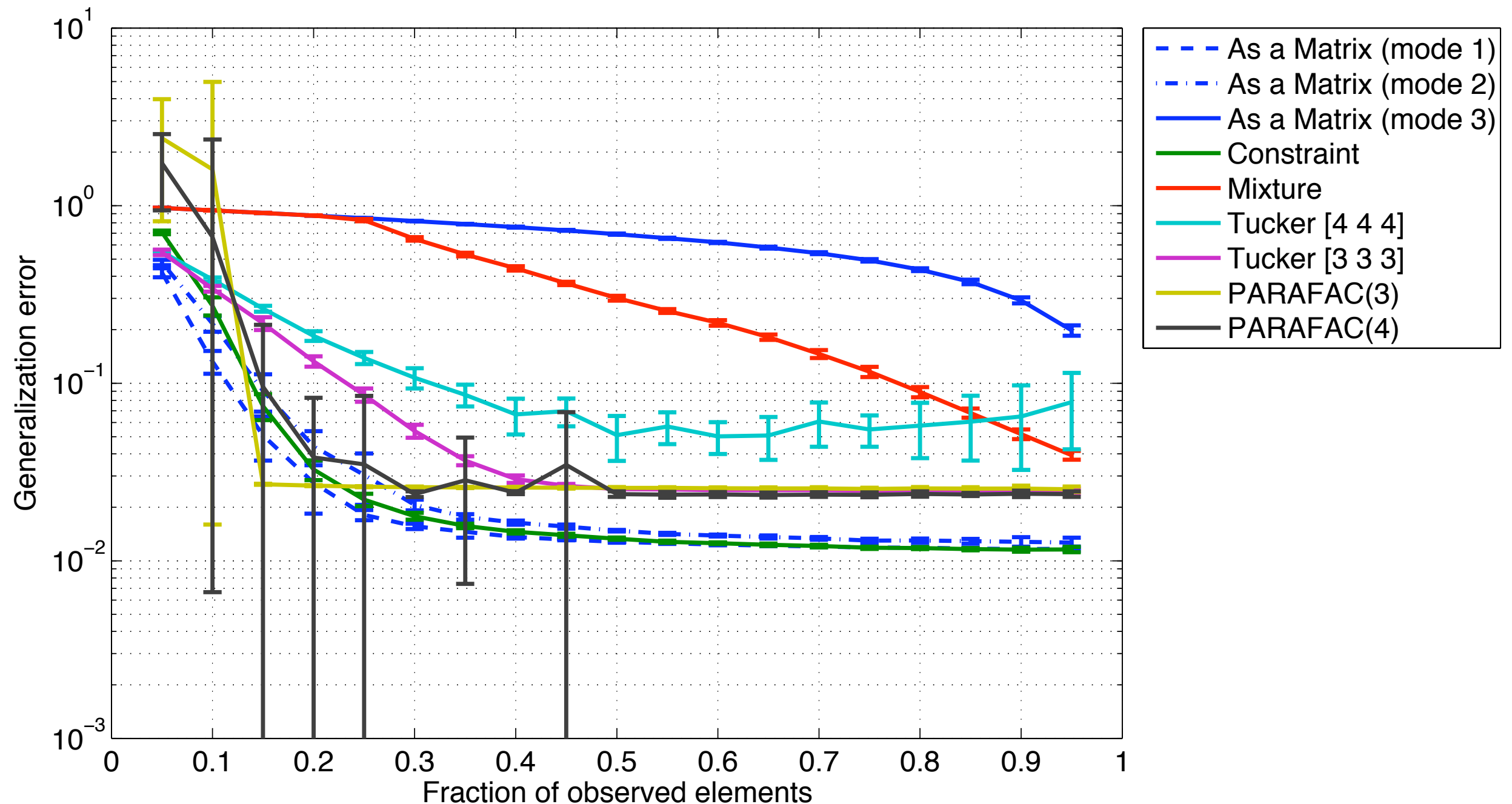


Amino acid fluorescence data [Bro & Andersson]

- Size 201x61x5.
- Five solutions with different amount of three amino acids (tyrosine, tryptophan, phenylalanine)
- Rank=3 PARAFAC is correct.
- Interested in
 - Generalization performance
 - Number of components
 - Interpretation

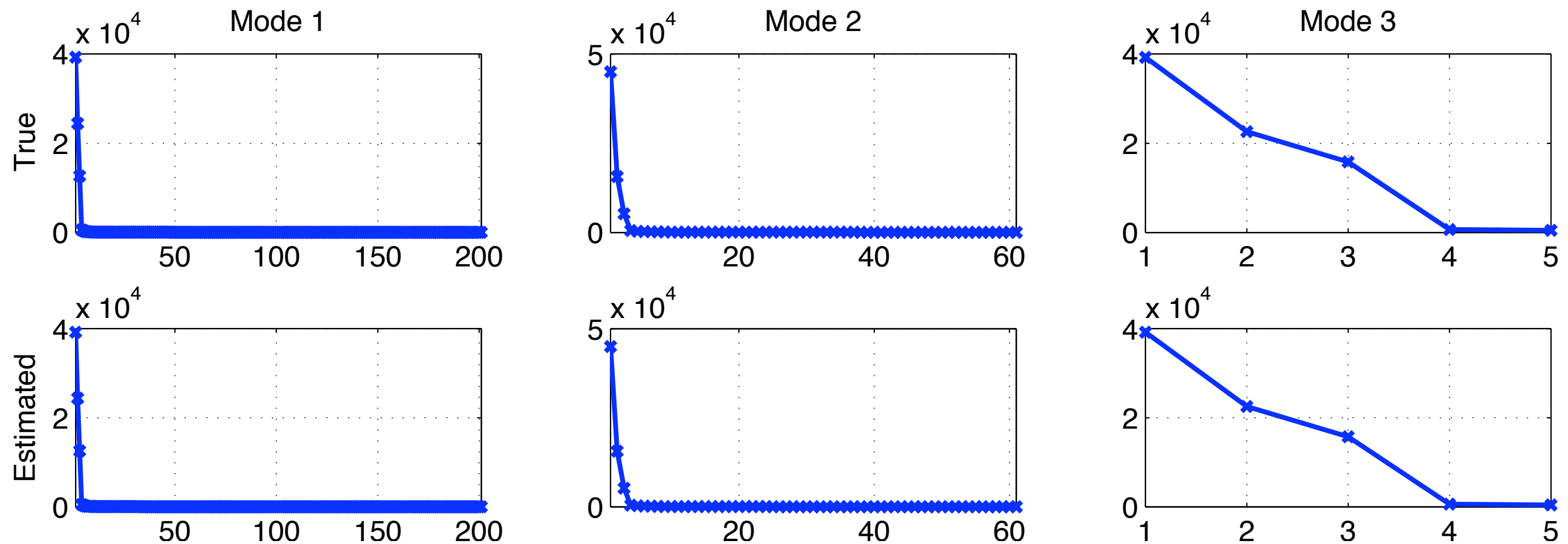
Amino acid: Generalization performance

- “Constraint” performs comparable to PARAFAC with the correct rank.



Amino acid: Singular-value spectra

Estimated spectra from half of the entries are almost identical to the truth.



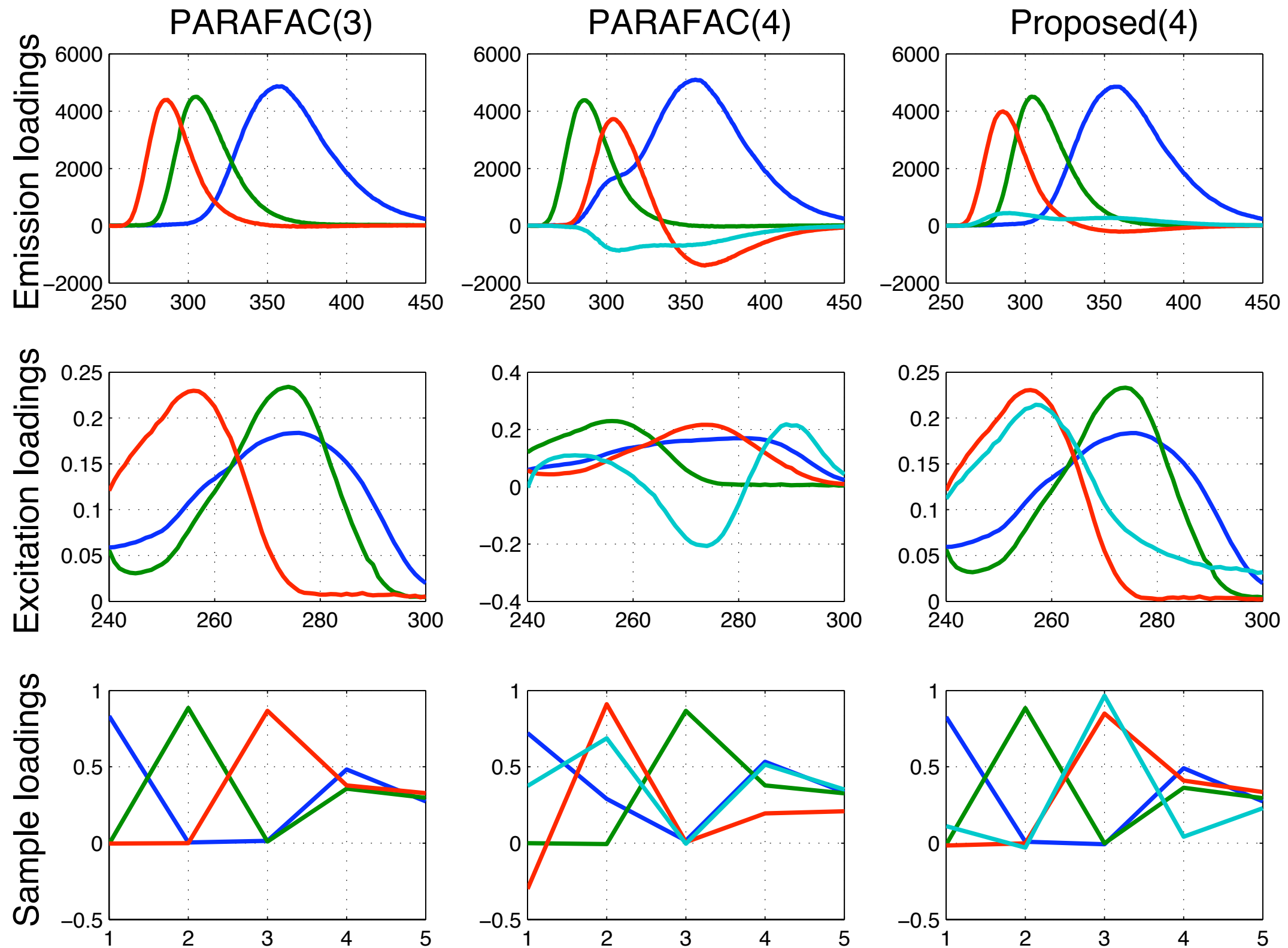
Improving Interpretability

- Apply PARAFAC on the core (4x4x5) obtained by the proposed “constraint” approach.
- Separate imputation problem and interpretation problem.

$$\begin{aligned}\mathcal{X} &= \mathcal{C} \times_1 U_1 \times_2 U_2 \times_3 U_3 \\ &= (\mathbf{A}^{(1)} \odot \mathbf{A}^{(2)} \odot \mathbf{A}^{(3)}) \times_1 U_1 \times_2 U_2 \times_3 U_3 \\ &\stackrel{\text{PARAFAC}}{=} (U_1 \mathbf{A}^{(1)}) \odot (U_2 \mathbf{A}^{(2)}) \odot (U_3 \mathbf{A}^{(3)})\end{aligned}$$

PARAFAC

Obtained factors



Summary

- Low-rank tensor completion can be computed in a **convex optimization problem** using the trace norm regularization.
 - No need to specify the rank beforehand.
- Convex formulation is more **accurate** and **faster** than conventional EM-based Tucker decomposition.
- Curious “phase transition” found → compressive-sensing-type analysis is an on-going work.
- Combination of proposed+PARAFAC is useful.
- Code:
 - <http://www.ibis.t.u-tokyo.ac.jp/RyotaTomiooka/Softwares/Tensor>

Acknowledgment

- This work was supported in part by MEXT KAKENHI 22700138, 80545583, JST PRESTO, and NTT Communication Science Laboratories.

ADMM convergence

- Step 1: ADMM is equivalent to Douglas-Rachford Splitting in the dual

$$\alpha^{t+1} = \text{prox}_{g^*} \left(\text{prox}_{f^*}(-\mathbf{A}^\top \cdot) (\alpha^t - z^t) + z^t \right)$$

$$z^{t+1} = \text{prox}_g \left(\text{prox}_{f^*}(-\mathbf{A}^\top \cdot) (\alpha^t - z^t) + z^t \right)$$