

Slides available:

<http://www.ibis.t.u-tokyo.ac.jp/ryotat/tensor12kyoto.pdf>

# Statistical Performance of Convex Tensor Decomposition

Ryota Tomioka

2012/01/26 @ Kyoto University

Perspectives in Informatics 4B

Collaborators: Taiji Suzuki, Kohei Hayashi, Hisashi Kashima

# Netflix challenge (2006-2009)

---

- \$1,000,000 prize
- Goal: Improve the performance of a **video recommendation system**

(predict who likes which movies)

- Example:



Likes “Star Wars” and “E.T.”,  
Doesn’t like “Minority Report”.

Does he like “Blade Runner”?

# Matrix completion view

---

	Star Wars	E.T.	Minority Report	Blade Runner	Monsters Inc.
User A	+1	+1	-1	?	?
User B	+1	?	?	+1	?
User C	?	+1	-1	?	+1
User D	+1	?	?	?	+1
.	.	.	.	.	.
.	.	.	.	.	.

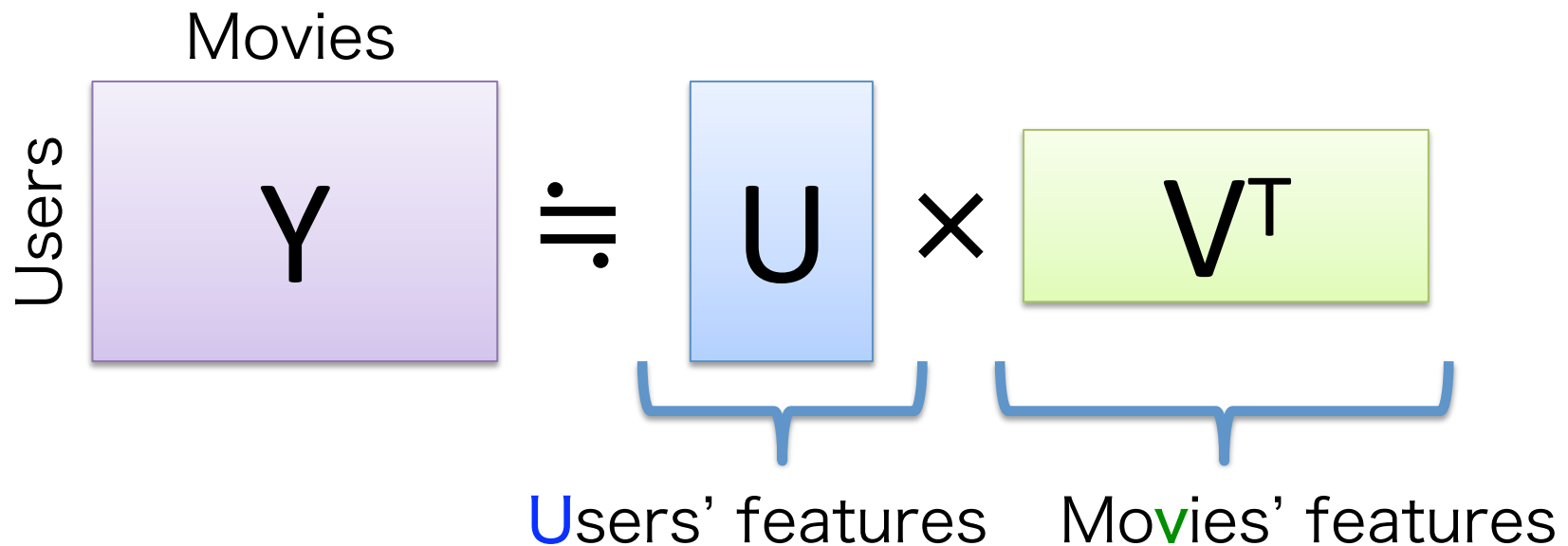
Goal: fill the missing entries!

# Matrix completion

---

- Impossible without an assumption. (Missing entries can be arbitrary) --- **problem is ill-posed**
- Most common assumption:

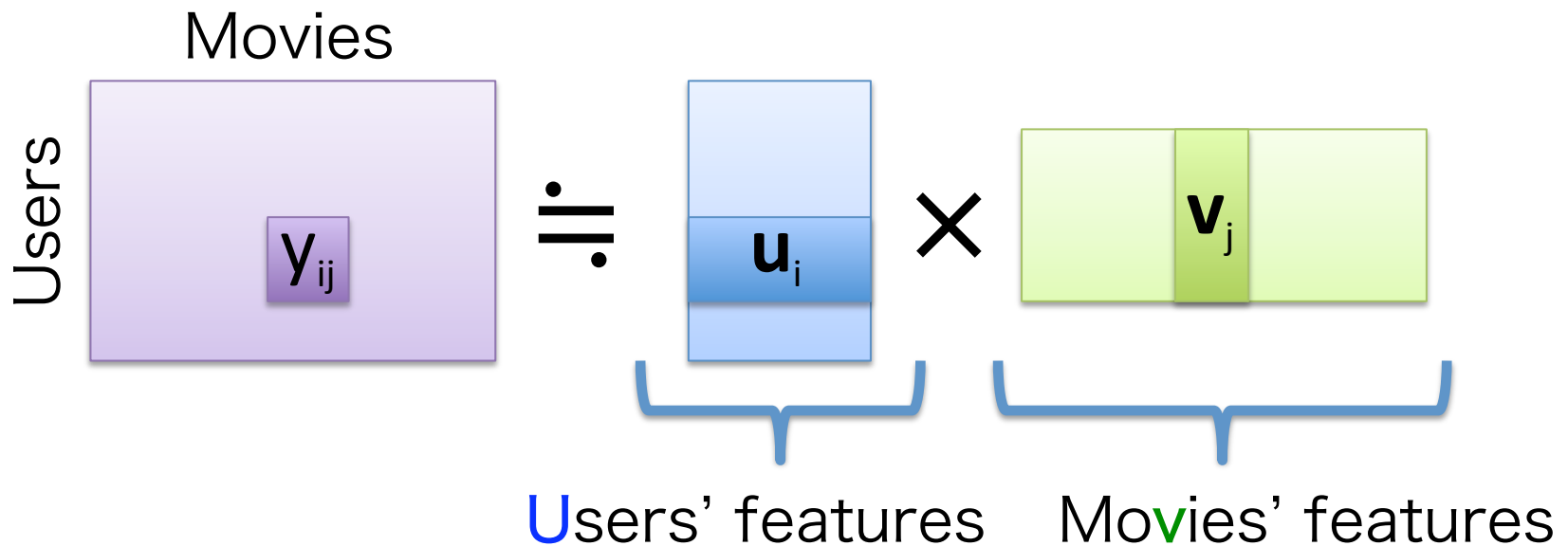
Low-rank decomposition



# Matrix completion

- Most common assumption:

Low-rank decomposition (rank  $r$ )

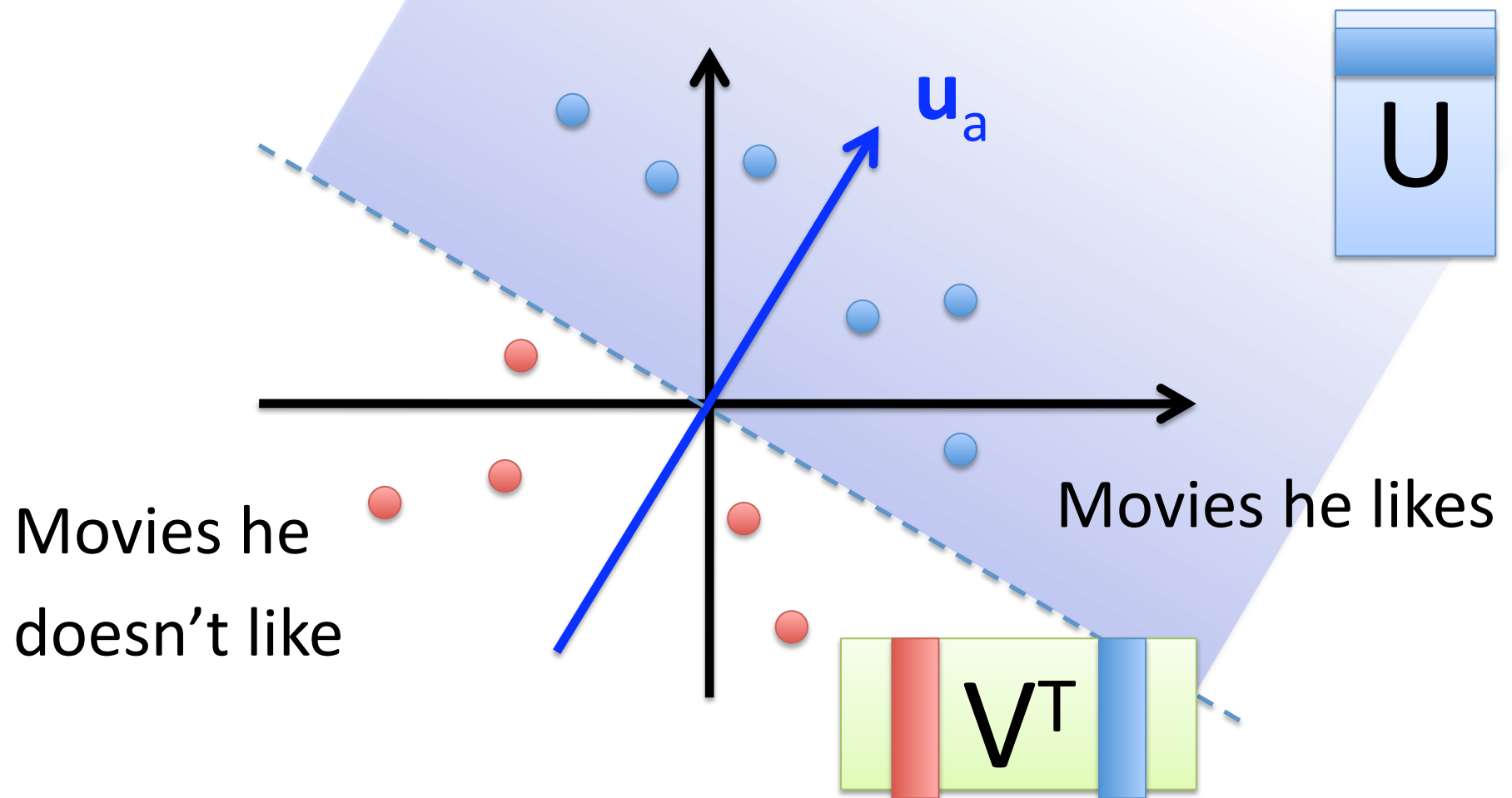


$$y_{ij} = u_i^T v_j \quad \left( \begin{array}{l} \text{dot product} \\ \text{in } r\text{-dim space} \end{array} \right)$$

# Geometric Intuition

r-dimensional space

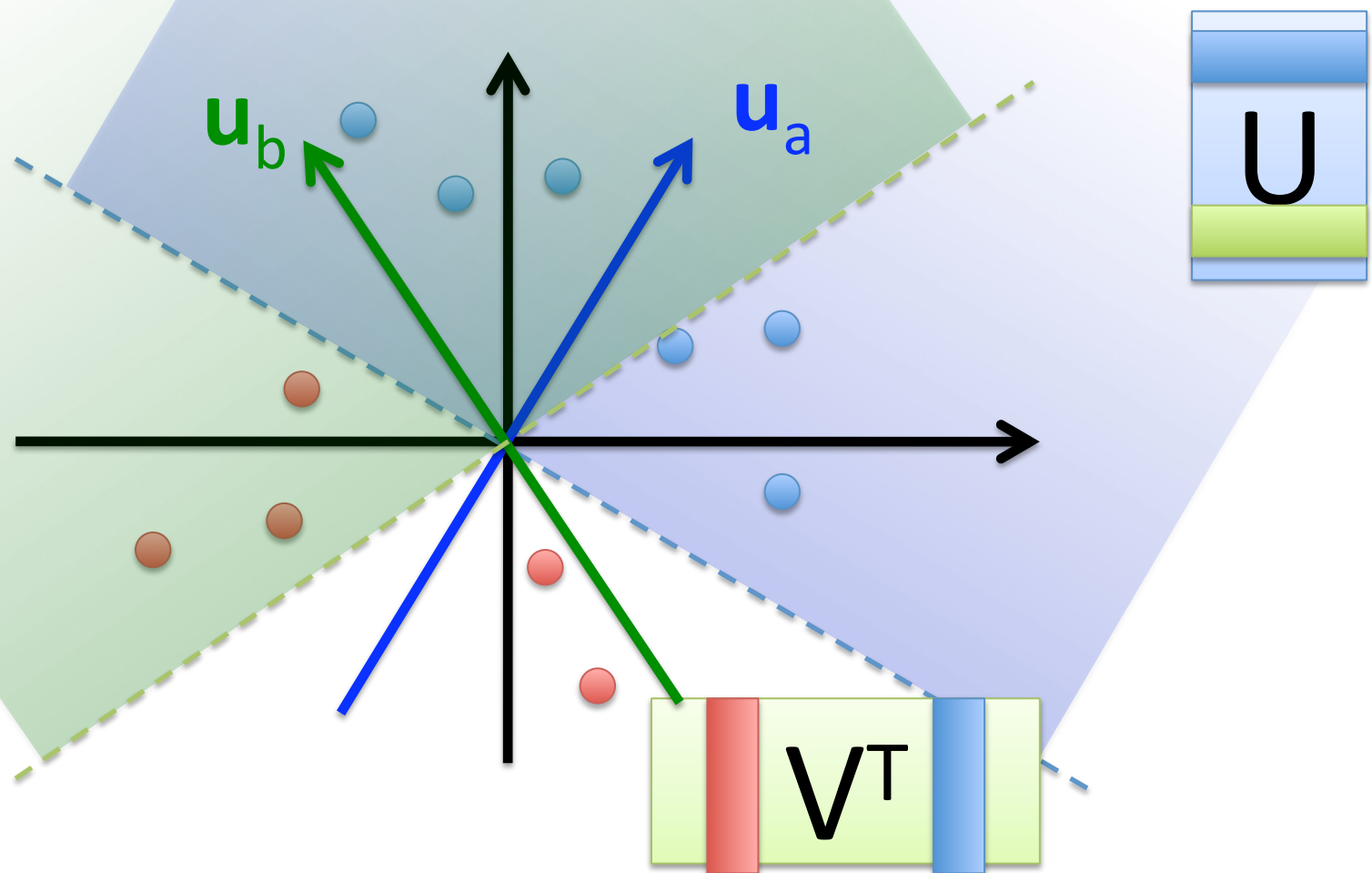
(r: the rank of the decomposition)



# Geometric Intuition

r-dimensional space

(r: the rank of the decomposition)



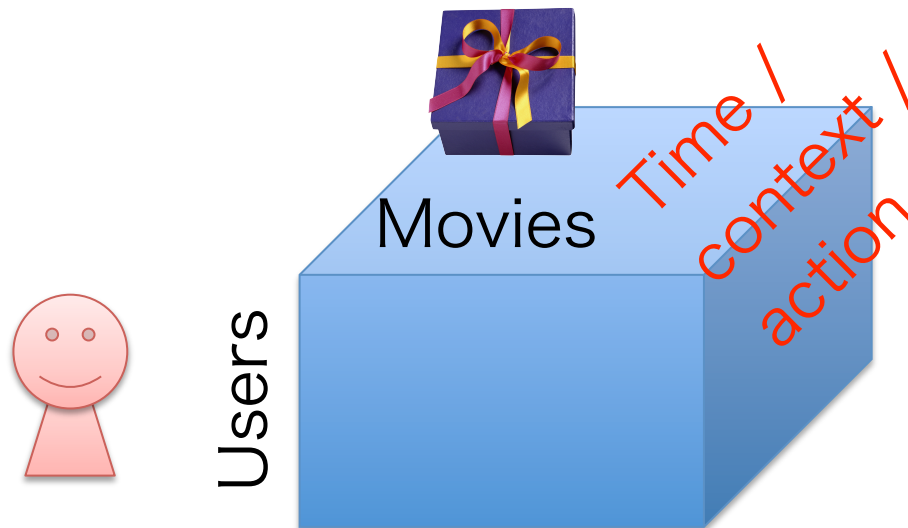
# Tensor data completion

---

- Tensor = Multi-dimensional array
- Beyond 2D

Movie preference

+ time / context / action



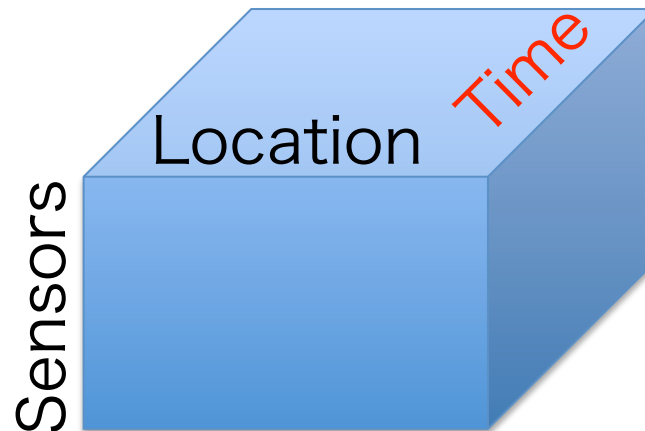
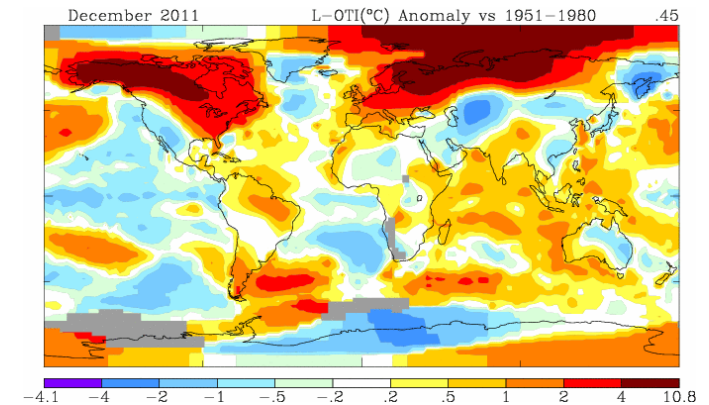


# Tensor data completion

- Tensor = Multi-dimensional array
- Beyond 2D

Climate monitoring

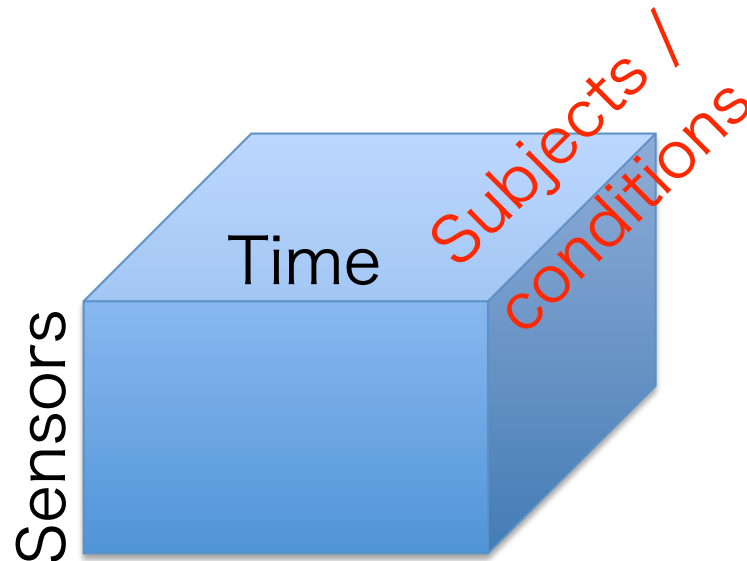
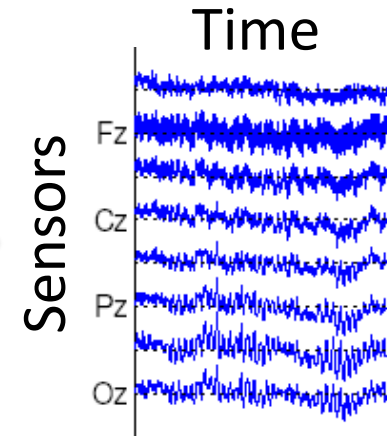
- temperature
- humidity
- rainfall



# Tensor data completion

- Tensor = Multi-dimensional array
- Beyond 2D

Neuroscience  
(brain imaging)



# Rest of this talk

---

- Computing low-rank matrix decomposition
- Generalizing from matrix to tensor
- Analyzing the performance
  - Statistical learning theory

# Computing low-rank matrix decomposition

# Computing low-rank decomposition

- If all entries are observed (no missing entries)
  - Given  $Y$ , compute singular value decomposition (SVD)

The diagram illustrates the SVD decomposition of a matrix  $Y$ . Matrix  $Y$  is represented by a purple box with dimensions  $m$  (height) and  $n$  (width). It is equated to the product of three matrices:  $U$  (blue box,  $m \times r$ ),  $\Sigma$  (red box,  $r \times r$ ), and  $V^T$  (green box,  $r \times n$ ). The dimensions are explicitly labeled around each matrix box.

where  $U, V$ : Orthogonal ( $U^T U = I, V^T V = I$ )

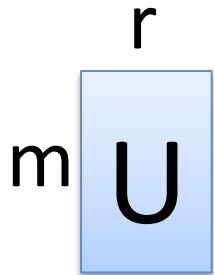
$$\Sigma = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{pmatrix} \quad \sigma_j: j\text{th largest singular value}$$

# Tolerating missings

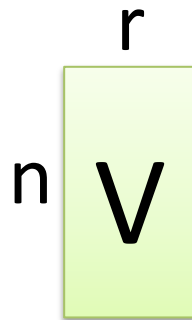
Optimization problem

$$\underset{U, V}{\text{minimize}} \sum_{(ij) \in \Omega} (y_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2$$

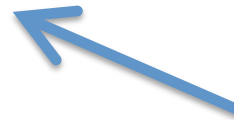
Users'  
features



Movies'  
features



Set of observed  
index pairs



# Tolerating missings

Optimization problem

Non-convex!

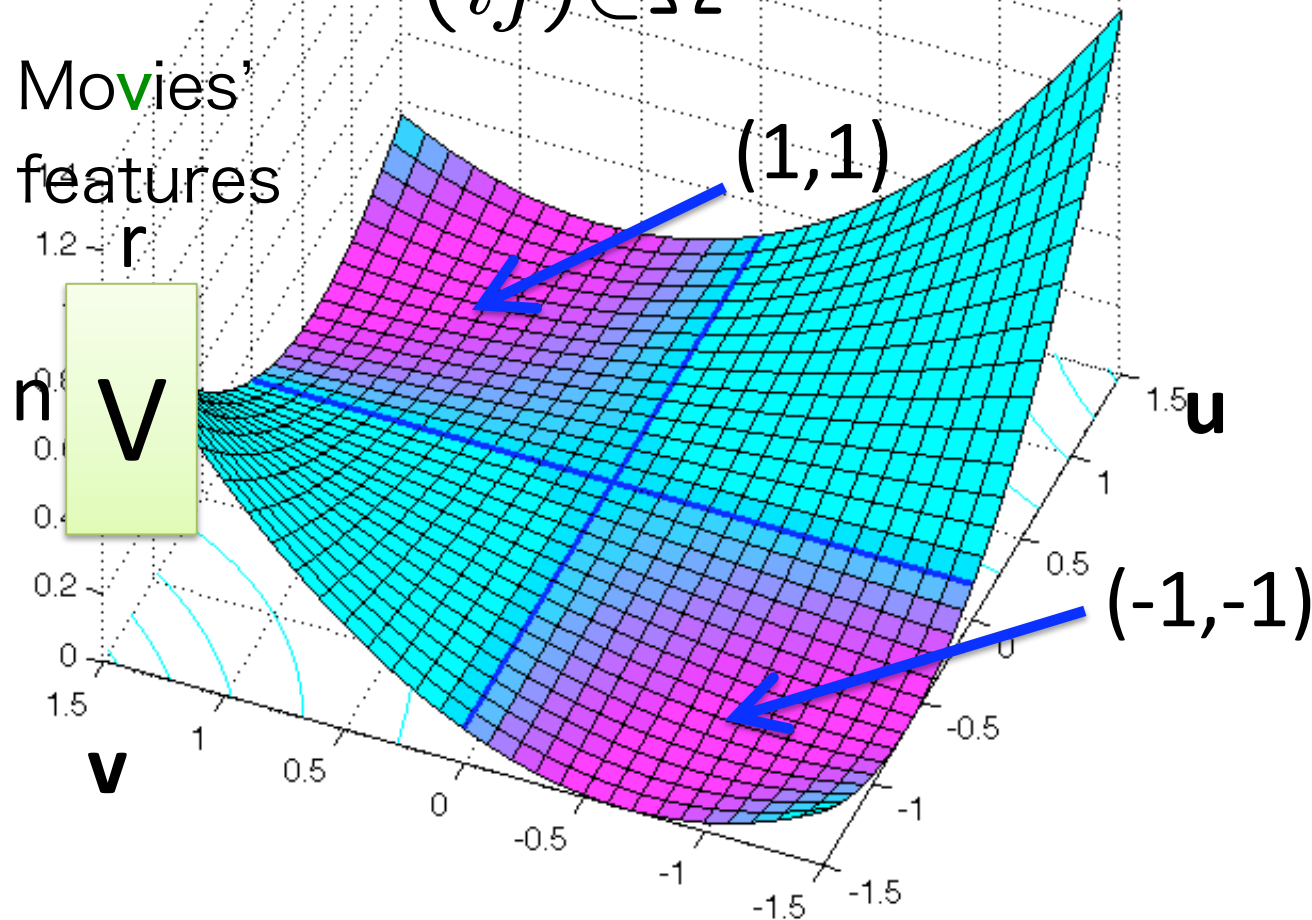
$$\underset{U, V}{\text{minimize}} \sum_{(i,j) \in \Omega} (y_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2$$

Users' features

$r$   
 $m$   $U$

Movies' features

$n$   $V$



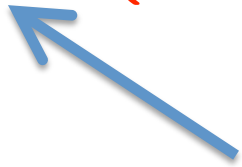
# Tolerating missings

---

Optimization problem

Still non-convex!

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \sum_{(ij) \in \Omega} (y_{ij} - w_{ij})^2, \\ \text{subject to} & \text{rank}(\mathbf{W}) \leq r \end{array}$$

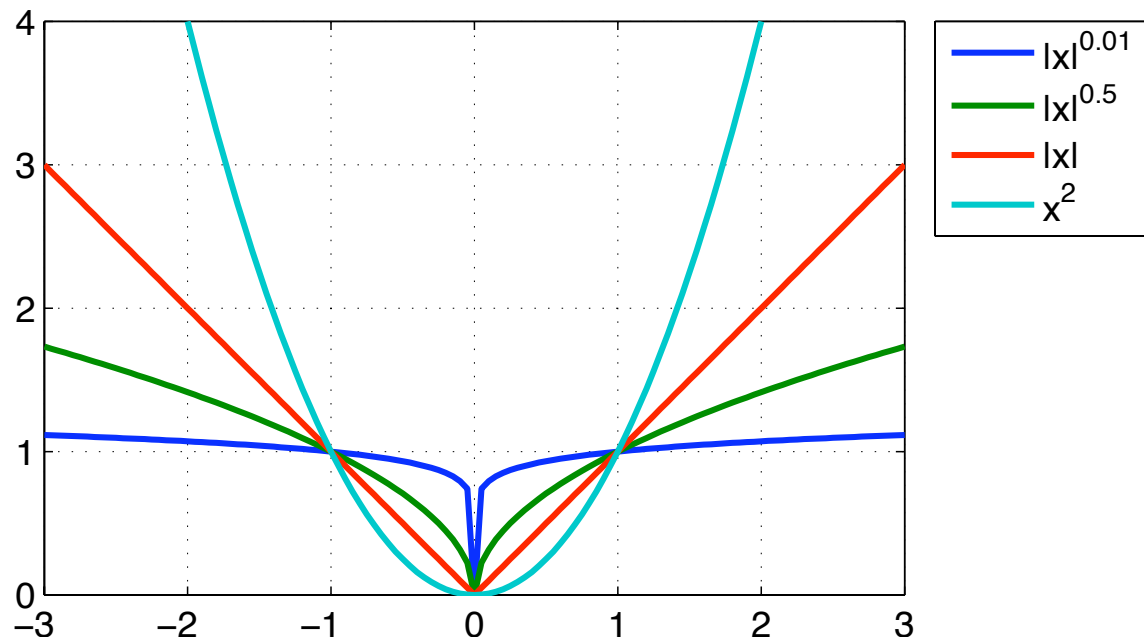
 Rank constraint  
is NP hard



# Convex relaxation of rank

Schatten  $p$ -norm  
(to the  $p$ th power)  $\|\mathbf{W}\|_{S_p}^p := \sum_{j=1}^r \sigma_j^p(\mathbf{W})$   
 $\sigma_j(\mathbf{W})$  :  $j$ th largest singular value

$$\|\mathbf{W}\|_{S_p}^p \xrightarrow{p \rightarrow 0} \text{rank}(\mathbf{W})$$



$p=1$  is the tightest  
convex relaxation  
(also known as  
trace norm /  
nuclear norm)

# Tolerating missings

Optimization problem

Convex relaxation

$$\begin{aligned} & \underset{\mathbf{W}}{\text{minimize}} && \sum_{(ij) \in \Omega} (y_{ij} - w_{ij})^2, \\ & \text{subject to} && \|\mathbf{W}\|_{S_1} \leq \tau \end{aligned}$$

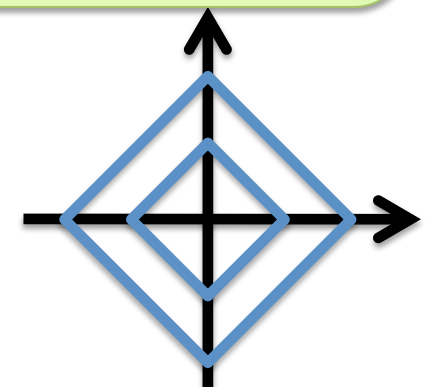
Schatten 1-norm

(nuclear norm,  
trace norm)

$$\|\mathbf{W}\|_{S_1} = \sum_{j=1}^r \sigma_j(\mathbf{W})$$

$\sigma_j(\mathbf{W})$  :  $j$ th largest singular value

Cf. Lasso ( $L_1$  norm) for variable selection  
= linear sum of abs. coefficients



# Take home messages

---

- Rank constrained minimization is hard to solve (non-convex and NP hard)
- Can be relaxed into a tractable convex problem using Schatten 1-norm.

# How about tensors?

- How to define tensor rank?
- How related to matrix rank?

# Ranf of a tensor

Definition. Let  $\mathcal{X} \in \mathbb{R}^{n_1 \times \cdots \times n_K}$  ( $K$ th order tensor)

The smallest number  $R$  such that the given tensor  $X$  is written as

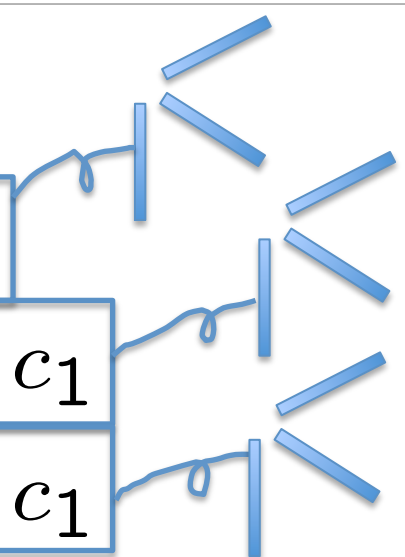
$$\mathcal{X} = \sum_{r=1}^R \mathcal{A}_r \quad \text{where} \quad \mathcal{A}_r = \begin{array}{c} \diagup \\ | \diagdown \end{array} \quad \text{is rank one.}$$

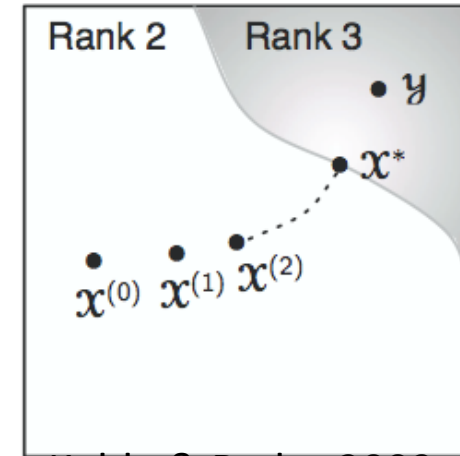
(can be written as an outer product of  $K$  vectors)

- Called CP (CANDECOMPO/PARAFAC) decomposition
- Bad news: NP hard to compute the rank  $R$  even for a fully observed  $X$ .

# Bad news 2: Tensor rank is not closed

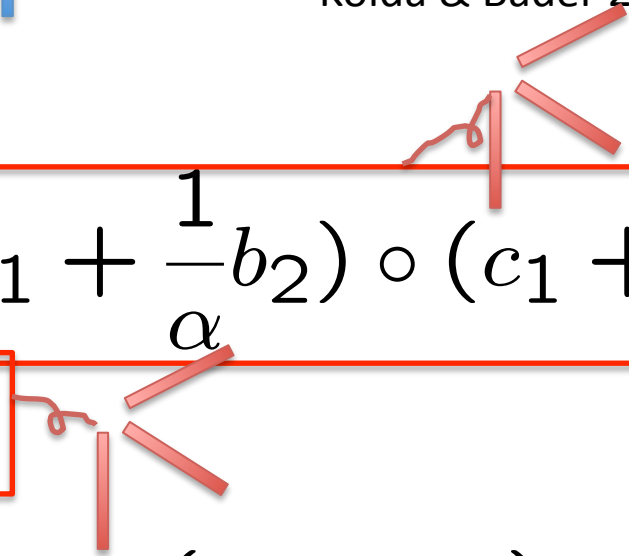
X is rank 3

$$\mathcal{X} = \boxed{a_1 \circ b_1 \circ c_2} + \boxed{a_1 \circ b_2 \circ c_1} + \boxed{a_2 \circ b_1 \circ c_1}$$




Kolda & Bader 2009

Y is rank 2

$$\mathcal{Y} = \boxed{\alpha \left( a_1 + \frac{1}{\alpha} a_2 \right) \circ \left( b_1 + \frac{1}{\alpha} b_2 \right) \circ \left( c_1 + \frac{1}{\alpha} c_2 \right)} - \boxed{\alpha a_1 \circ b_1 \circ c_1}$$


$$\|\mathcal{X} - \mathcal{Y}\|_F \rightarrow 0 \quad (\alpha \rightarrow \infty)$$

# Tucker decomposition [Tucker 66]

Diagram illustrating the Tucker decomposition of a 3D tensor  $X$  (size  $n_1 \times n_2 \times n_3$ ) into a Core tensor  $C$  (size  $r_1 \times r_2 \times r_3$ ) and three orthogonal factor matrices  $U^{(1)}$ ,  $U^{(2)}$ , and  $U^{(3)}$  (sizes  $n_1 \times r_1$ ,  $n_2 \times r_2$ , and  $n_3 \times r_3$  respectively).

The decomposition is represented by the equation:

$$X = C \times_1 U^{(1)} \times_2 U^{(2)} \times_3 U^{(3)}$$

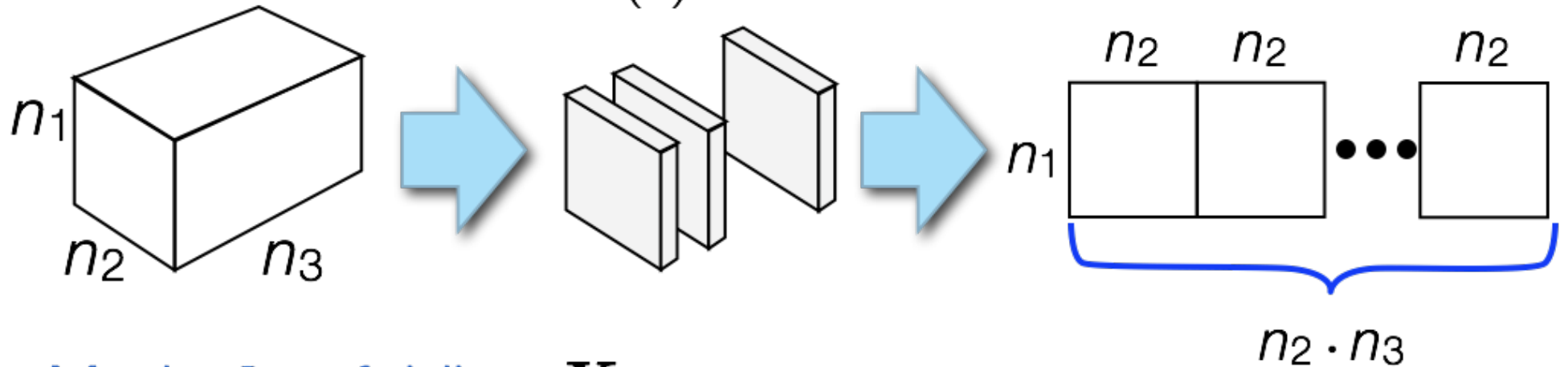
Below the diagram, the tensor equation is given:

$$\left( X_{ijk} = \sum_{a=1}^{r_1} \sum_{b=1}^{r_2} \sum_{c=1}^{r_3} C_{abc} U_{ia}^{(1)} U_{jb}^{(2)} U_{kc}^{(3)} \right)$$

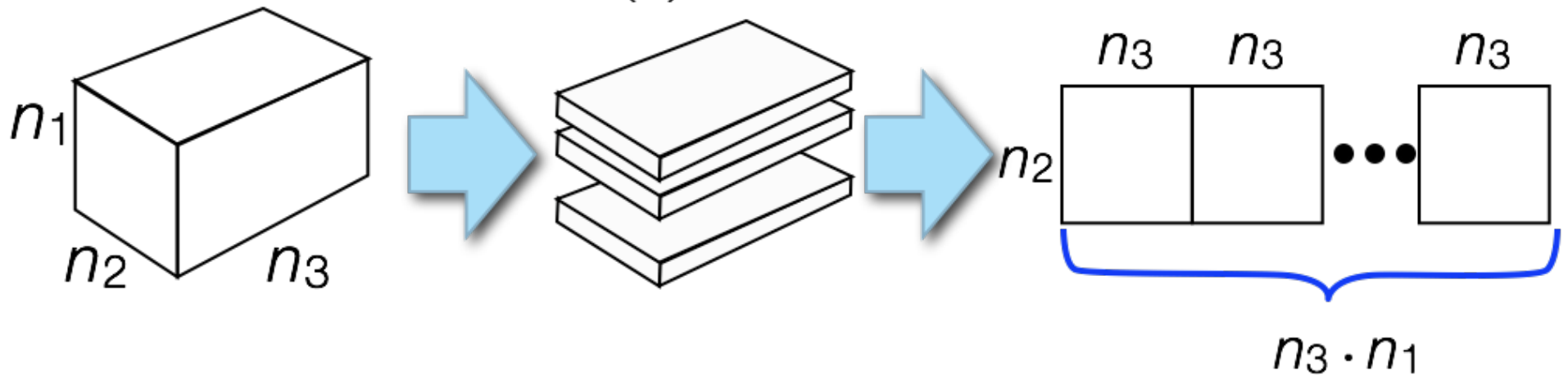
- Also known as higher-order SVD [De Lathauwer+00]
- Rank  $(r_1, r_2, r_3)$  can be computed in polynomial time using unfolding operations.

# Mode-k unfoldings (matricization)

Mode-1 unfolding  $\mathbf{X}_{(1)}$



Mode-2 unfolding  $\mathbf{X}_{(2)}$

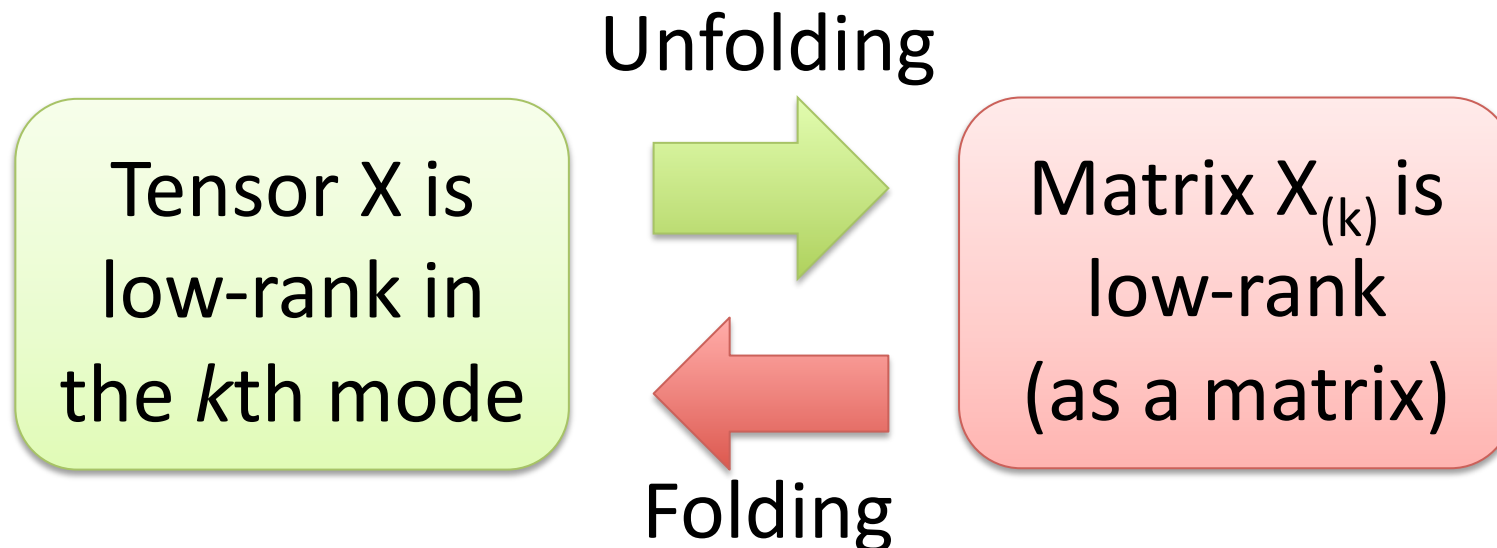




# Computing Tucker rank

---

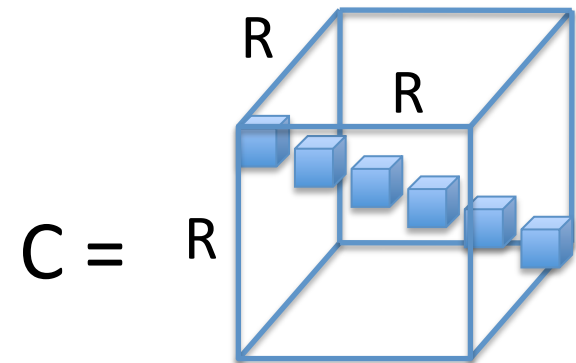
- For each  $k=1,\dots,K$ 
    - Compute the **mode-k unfolding**  $X_{(k)}$
    - Compute the (matrix) rank of  $X_{(k)}$
- $$r_k = \text{rank}(\mathbf{X}_{(k)})$$



# Computing Tucker rank

- For each  $k=1,\dots,K$ 
    - Compute the **mode-k unfolding**  $X_{(k)}$
    - Compute the (matrix) rank of  $X_{(k)}$
- $$r_k = \text{rank}(\mathbf{X}_{(k)})$$
- Difference between Tensor rank and Tucker rank
    - Tensor rank is a single number  $R$  (may not be easy to compute)
    - Tucker rank is defined for each mode (easy to compute)

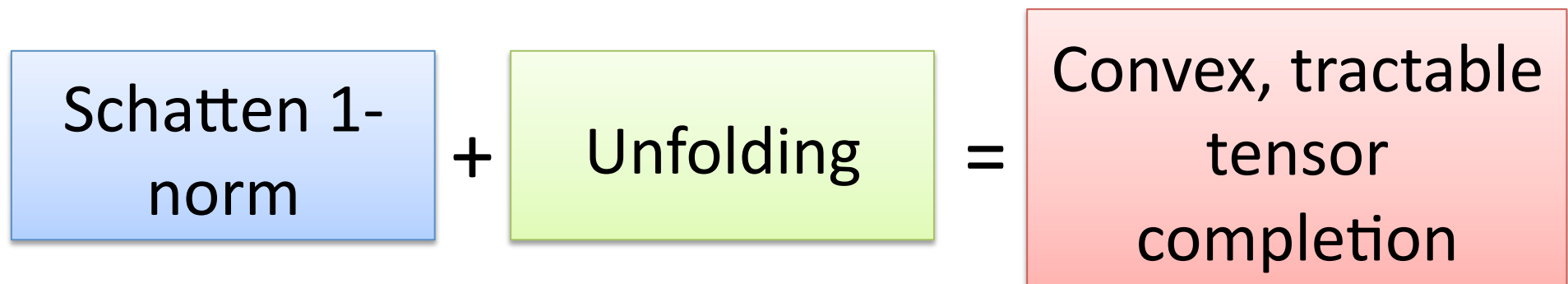
- CP decomp is a special case of Tucker decomp with  
 $R=r_1=r_2=\dots=r_K$  and diagonal core



# Basic idea

---

- We know how to do **matrix completion with Schatten 1-norm (tractable convex optimization)**
- We know how to compute Tucker rank (=the rank of the mode-k unfolding)



# Overlapped Schatten 1-norm for Tensors

$$\left| \left| \left| \mathcal{W} \right| \right|_{S_1} := \frac{1}{K} \sum_{k=1}^K \left\| \mathbf{W}_{(k)} \right\|_{S_1}$$

Schatten 1-norm of  
the [mode-k unfolding](#)

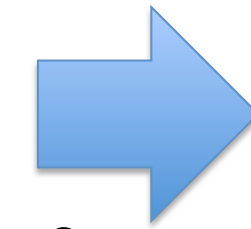
Measures the overall low-rank-ness  
(not just a single mode)

# Convex Tensor Estimation

---

## Matrix

Estimation of *low-rank* matrix  
(hard)



Convex  
relaxation

Schatten 1-norm  
minimization  
(tractable)

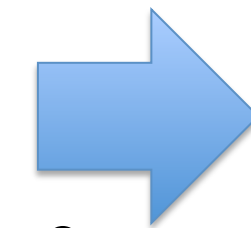
[Fazel, Hindi, Boyd  
01]



Generalize

## Tensor

Estimation of *low-rank* **tensor**  
(hard)



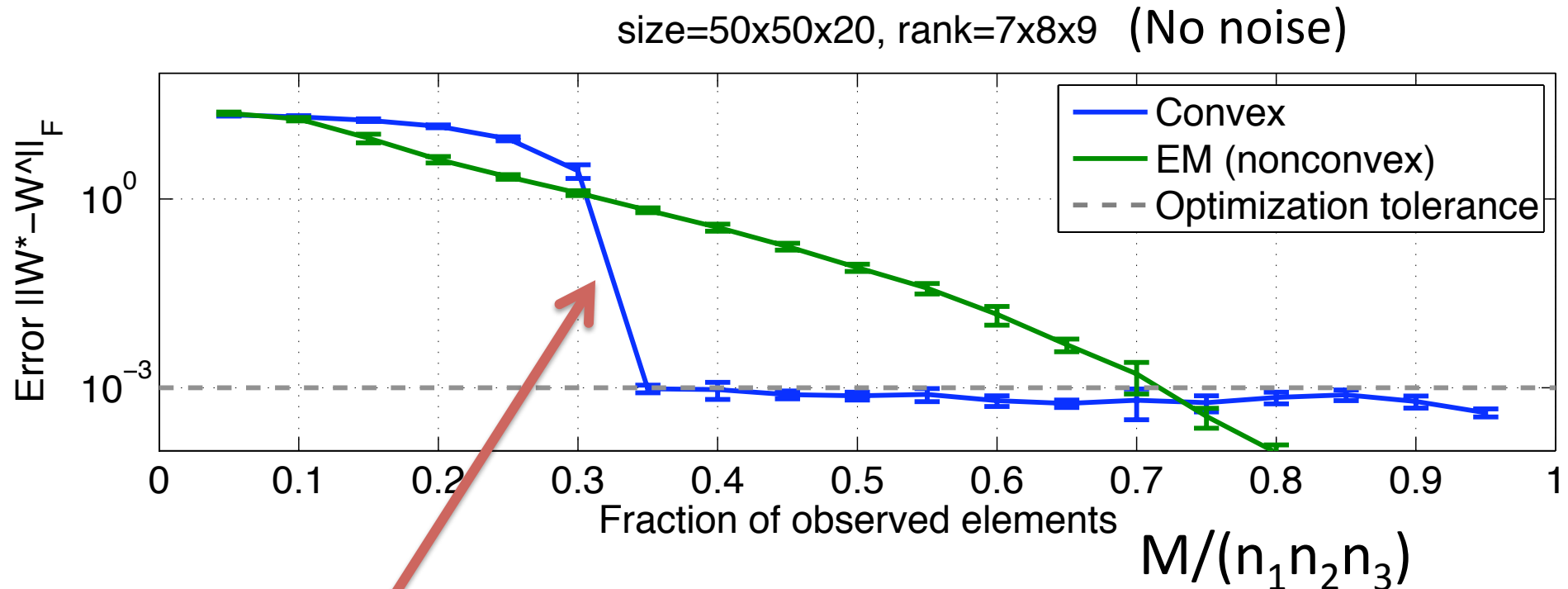
Convex  
relaxation

**Overlapped  
Schatten 1-norm  
minimization**

[Liu+09, Signoretto+10,  
Tomioka+10, Gandy+11]

# Empirical performance

Tensor completion result [Tomika et al. 2010]



Phase transition!!

Can we predict this theoretically?

# Analyzing the performance of convex tensor decomposition

# Problem setting

---

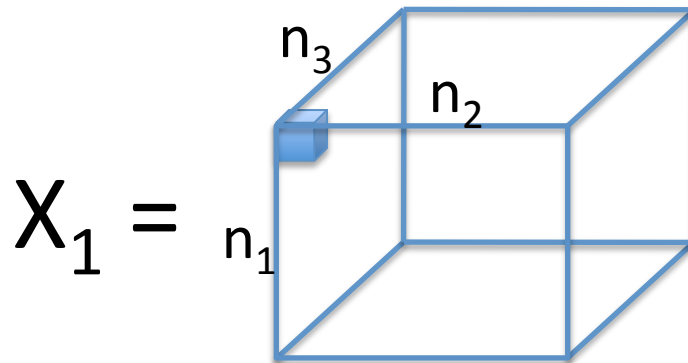
## Observation model

$$y_i = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + \epsilon_i \quad (i = 1, \dots, M)$$

$\mathbf{w}^*$  true tensor rank- $(r_1, \dots, r_K)$

$\epsilon_i$  Gaussian noise

## Example (tensor completion)





# Problem setting

---

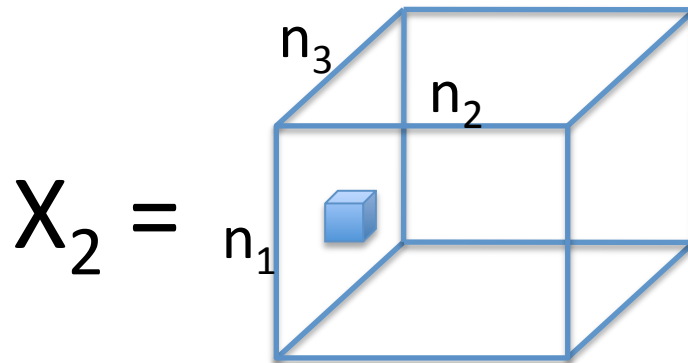
## Observation model

$$y_i = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + \epsilon_i \quad (i = 1, \dots, M)$$

$\mathbf{w}^*$  true tensor rank- $(r_1, \dots, r_K)$

$\epsilon_i$  Gaussian noise

## Example (tensor completion)



# Problem setting

---

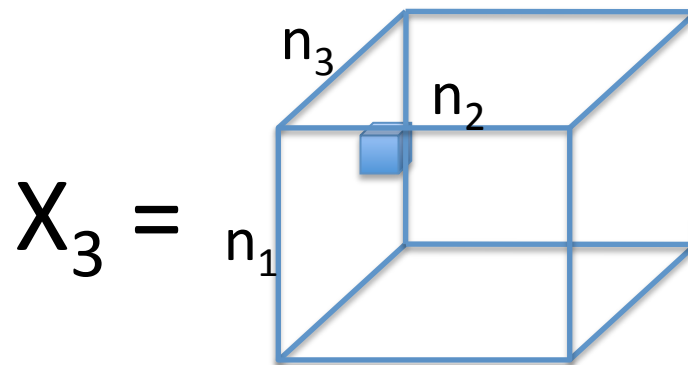
## Observation model

$$y_i = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + \epsilon_i \quad (i = 1, \dots, M)$$

$\mathbf{w}^*$  true tensor rank- $(r_1, \dots, r_K)$

$\epsilon_i$  Gaussian noise

## Example (tensor completion)



# Problem setting

---

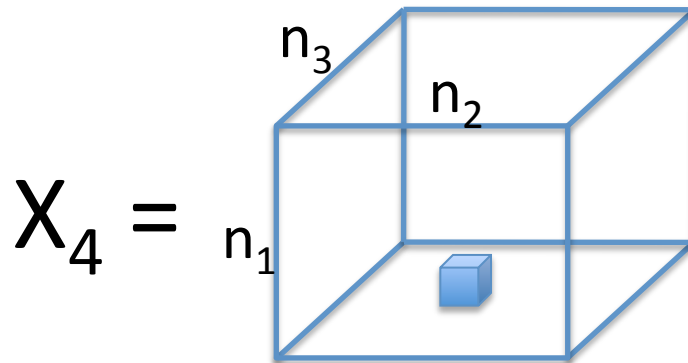
## Observation model

$$y_i = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + \epsilon_i \quad (i = 1, \dots, M)$$

$\mathbf{w}^*$  true tensor rank- $(r_1, \dots, r_K)$

$\epsilon_i$  Gaussian noise

## Example (tensor completion)



and so on...

# Problem setting

## Observation model

$$y_i = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + \epsilon_i \quad (i = 1, \dots, M)$$

$\mathbf{w}^*$  true tensor rank- $(r_1, \dots, r_K)$

$\epsilon_i$  Gaussian noise  $N(0, \sigma^2)$

## Optimization

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\operatorname{argmin}} \left( \underbrace{\frac{1}{2M} \|\mathbf{y} - \mathfrak{X}(\mathbf{w})\|_2^2}_{\text{Empirical error}} + \lambda_M \underbrace{\|\mathbf{w}\|_{S_1}}_{\text{Regularization}} \right)$$

$$(N = \prod_{k=1}^K n_k)$$

Observation  
model

Reg. Const.

$$\mathfrak{X} : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

$$\mathfrak{X}(\mathbf{w}) = (\langle \mathbf{x}_1, \mathbf{w} \rangle, \dots, \langle \mathbf{x}_M, \mathbf{w} \rangle)^\top$$

# Analysis objective

---

- We would like to show something like

Estimated tensor

True low-rank tensor

Mean squared error

$$\frac{\|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2}{N} \leq O_p \left( \frac{c(\mathbf{n}, \mathbf{r})}{M} \right)$$

The size  $\mathbf{n} = (n_1, \dots, n_K)$

The rank  $\mathbf{r} = (r_1, \dots, r_K)$

Number of samples  $M$

# Theorem: random Gauss design

---

Assume elements of  $X_i$  are drawn iid from standard normal distribution. Moreover

$$\frac{\text{\#samples } (M)}{\text{\#variables } (N)} \geq c_1 \underbrace{\|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}_{\text{Normalized rank}} \approx \frac{r}{n}$$

$$\|\mathbf{n}^{-1}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{1/n_k} \right)^2, \quad \|\mathbf{r}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2$$

# Theorem: random Gauss design

Assume elements of  $X_i$  are drawn iid from standard normal distribution. Moreover

$$\frac{\text{\#samples } (M)}{\text{\#variables } (N)} \geq c_1 \underbrace{\|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}_{\text{Normalized rank}} \approx \frac{r}{n}$$

Convergence!

$$\frac{\|\hat{\mathbf{w}} - \mathbf{w}^*\|_F^2}{N} \leq O_p \left( \frac{\sigma^2 \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}{M} \right)$$

$$\|\mathbf{n}^{-1}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{1/n_k} \right)^2, \quad \|\mathbf{r}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2$$

# Proof idea

Since  $\hat{\mathcal{W}}$  minimizes the objective,

Estimated  
tensor

True low-  
rank tensor

$$\text{Obj}(\hat{\mathcal{W}}) \leq \text{Obj}(\mathcal{W}^*)$$

It is not so hard to see:

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \left\langle \mathfrak{X}^*(\epsilon)/M, \hat{\mathcal{W}} - \mathcal{W}^* \right\rangle + \lambda_M \left| \left| \hat{\mathcal{W}} - \mathcal{W}^* \right| \right|_{S_1}$$

$\mathfrak{X}^*(\epsilon) = \sum_{i=1}^M \epsilon_i \mathcal{X}_i$

What we want to derive:

$$\frac{\left| \left| \hat{\mathcal{W}} - \mathcal{W}^* \right| \right|_F^2}{N} \leq O_p \left( \frac{c(n, r)}{M} \right)$$



# Proof outline (1/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \langle \mathfrak{X}^*(\epsilon)/M, \hat{\mathcal{W}} - \mathcal{W}^* \rangle + \lambda_M \|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1}$$

Inequality 1: upper-bound the dot product

$$\langle \mathfrak{X}^*(\epsilon)/M, \hat{\mathcal{W}} - \mathcal{W}^* \rangle \leq O_p \left( \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1} \right)$$

(optimization duality / random matrix theory)

# Proof outline (1/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \left( \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}} + \lambda_M \right) \|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1}$$

Inequality 1: upper-bound the dot product

$$\langle \mathfrak{X}^*(\epsilon)/M, \hat{\mathcal{W}} - \mathcal{W}^* \rangle \leq O_p \left( \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1} \right)$$

Trade-off between  $\sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}}$  and  $\lambda_M$

➡ Optimal reg. const  $\lambda_M \simeq O_p \left( \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}} \right)$

# Proof outline (2/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2}}{M}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1}$$

Inequality 2: relate the Schatten 1-norm  
with the Frobenius norm

$$\|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1} \leq \sqrt{\|\mathbf{r}\|_{1/2}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F$$

(relation between L1- and L2-norm)

# Proof outline (2/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}{M}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F$$

Inequality 2: relate the Schatten 1-norm  
with the Frobenius norm

$$\|\hat{\mathcal{W}} - \mathcal{W}^*\|_{S_1} \leq \sqrt{\|\mathbf{r}\|_{1/2}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F$$

(relation between L1- and L2-norm)

# Proof outline (3/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\frac{1}{2M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2 \leq \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}{M}} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F$$

Inequality 3: lower-bound the left hand-side

$$\kappa \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq \frac{1}{M} \|\mathfrak{X}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2$$

$$\text{If } \frac{\text{\#samples } (M)}{\text{\#variables } (N)} \geq c_1 \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}$$

(Gordon-Slepian Theorem in Gaussian process theory)

# Proof outline (3/3)

---

Estimated  
tensor

True low-  
rank tensor

$$\kappa \left\| \hat{\mathcal{W}} - \mathcal{W}^* \right\|_F^2 \leq \sqrt{\frac{\sigma^2 N \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}{M}} \left\| \hat{\mathcal{W}} - \mathcal{W}^* \right\|_F$$

Inequality 3: lower-bound the left hand-side

$$\kappa \left\| \hat{\mathcal{W}} - \mathcal{W}^* \right\|_F^2 \leq \frac{1}{M} \|\mathfrak{x}(\hat{\mathcal{W}} - \mathcal{W}^*)\|_2^2$$

$$\text{If } \frac{\text{\#samples } (M)}{\text{\#variables } (N)} \geq c_1 \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}$$

(Gordon-Slepian Theorem in Gaussian process theory)

# Back to the theorem statement

Assume elements of  $X_i$  are drawn iid from standard normal distribution. Moreover

$$\frac{\text{\#samples } (M)}{\text{\#variables } (N)} \geq c_1 \underbrace{\|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}_{\text{Normalized rank}} \approx \frac{r}{n}$$

Convergence!

$$\frac{\|\hat{\mathbf{w}} - \mathbf{w}^*\|_F^2}{N} \leq O_p \left( \frac{\sigma^2 \|\mathbf{n}^{-1}\|_{1/2} \|\mathbf{r}\|_{1/2}}{M} \right)$$

$$\|\mathbf{n}^{-1}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{1/n_k} \right)^2, \quad \|\mathbf{r}\|_{1/2} := \left( \frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2$$

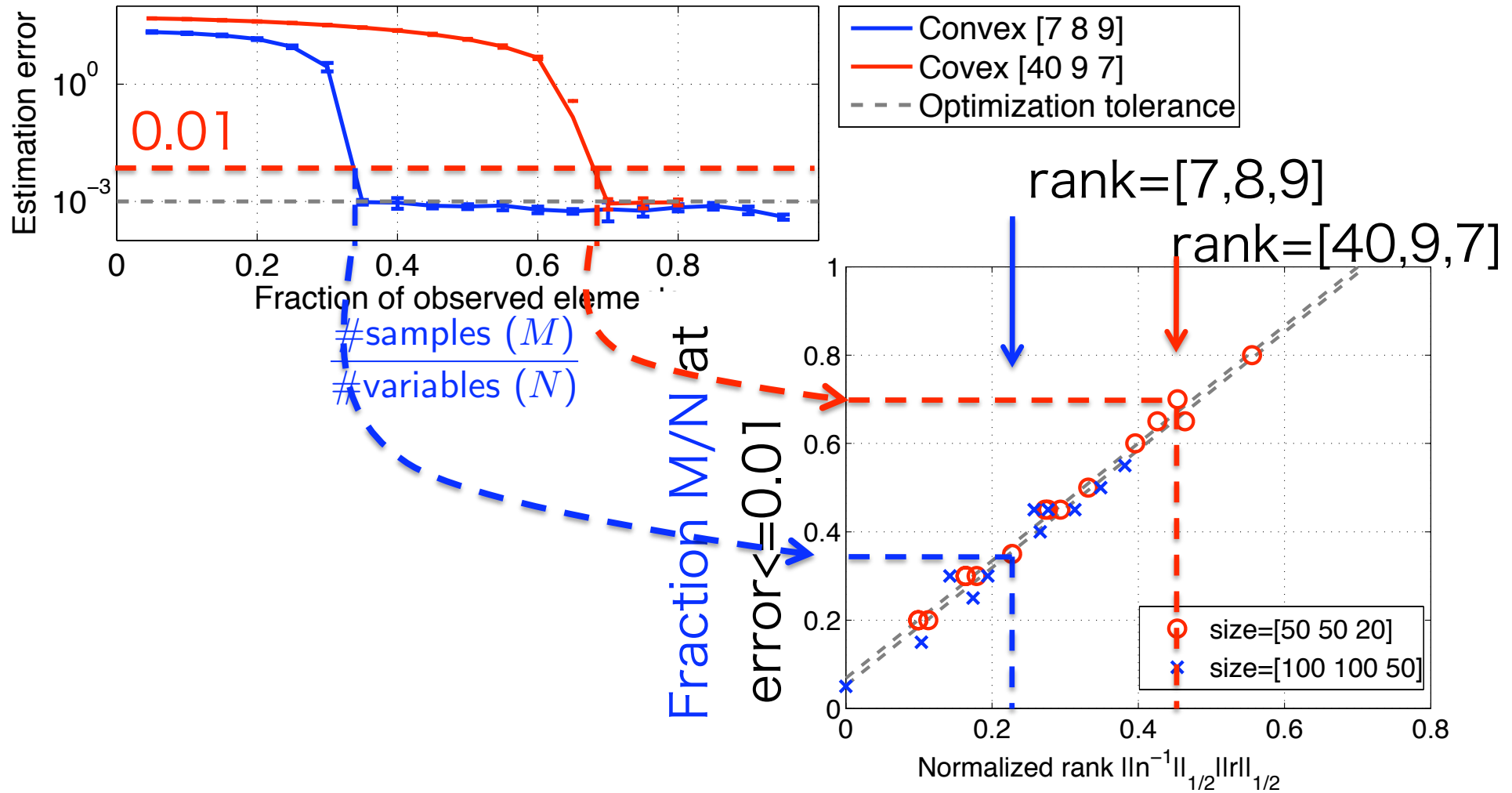
Notice:

- Sample-size condition independent of noise  $\sigma^2$ .
- Bound RHS proportional to  $\sigma^2$ .

Threshold behavior in the limit  $\sigma^2 \rightarrow 0$

# Tensor completion results

size = 50x50x20 true rank 7x8x9 or 40x9x7

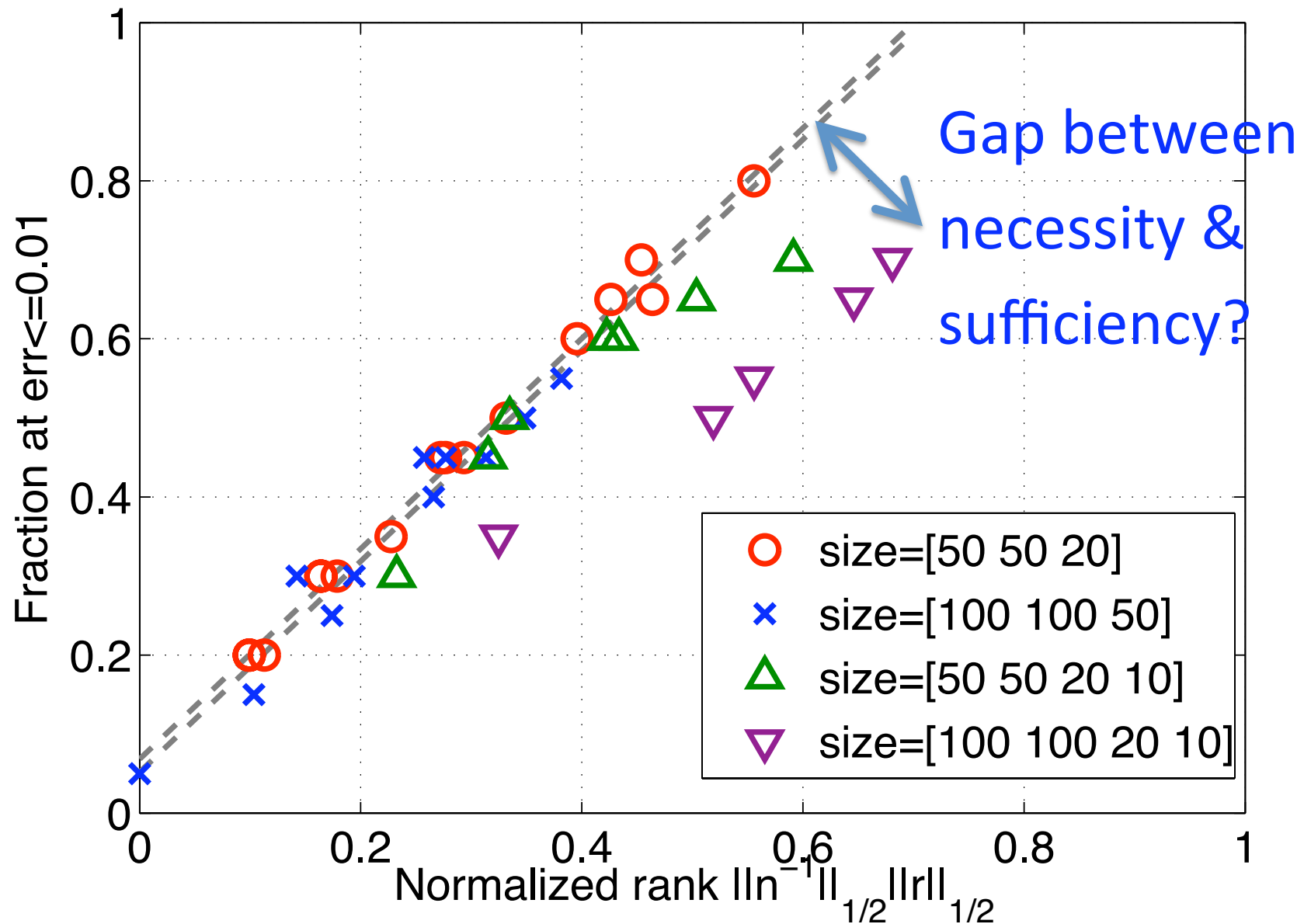


No observation noise

Normalized rank



# Including 4<sup>th</sup> order tensors

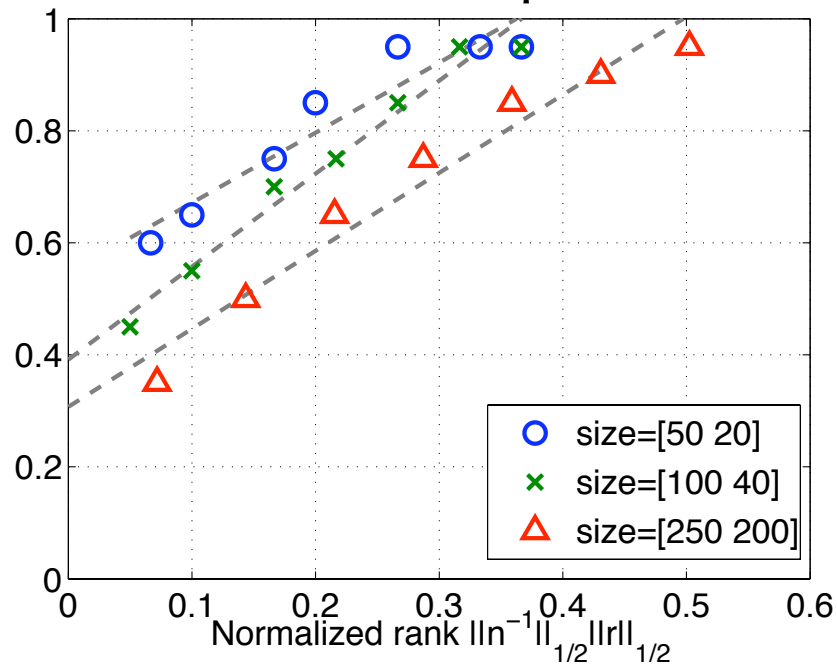


# Matrix / tensor completion

Fraction  $M/N$  at error  $\leq 0.01$

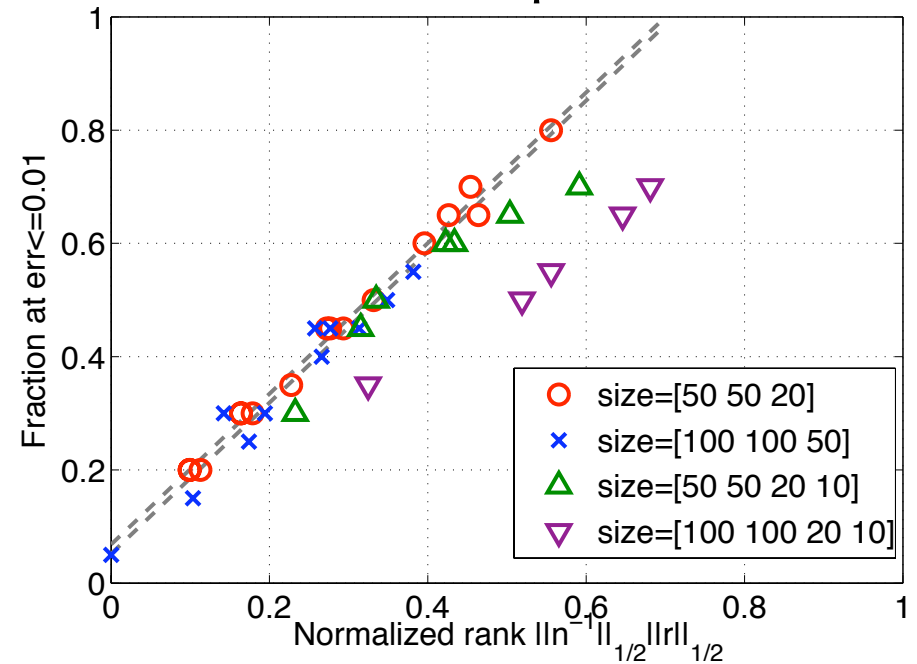
$K = 2$

Matrix completion



$K \geq 3$

Tensor completion



Tensor completion *easier* than matrix completion!?

# Conclusion

---

- Many real world problems can be cast into the form of tensor data analysis.
- Convex optimization is a useful tool also for the analysis of higher order tensors.
- Proposed a convex tensor decomposition algorithm with performance guarantee
- Normalized rank predicts empirical scaling behavior well

# Issues

---

- Why matrix completion more difficult than tensor completion?
- How big the gap between necessity and sufficiency?
- Random Gaussian design  $\neq$  tensor completion
  - $\Rightarrow$  Incoherence (Candes & Recht 09)
  - $\Rightarrow$  Spikiness (Negahban et al 10)
- When only some modes are low-rank
  - Schatten 1-norm is too strong  $\Rightarrow$  Mixture approach
  - E.g. Mode 1, 4 is low rank but the rest is not (combinatorial problem)

T h a n k y o u !