

Tensor decompositions: old, new, and beyond

Ryota Tomioka

@MLSS, Kyoto 2015

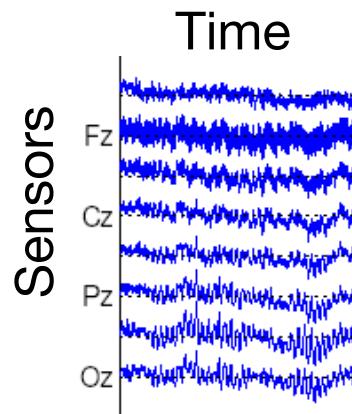
Toyota Technological Institute at Chicago



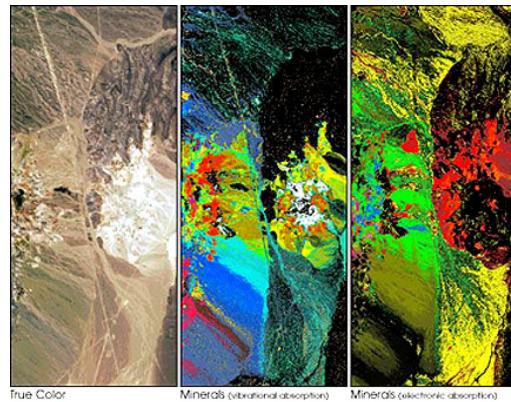
Tensors appear everywhere

Matrices
Tensors

Multivariate time-series



Spatio-temporal data

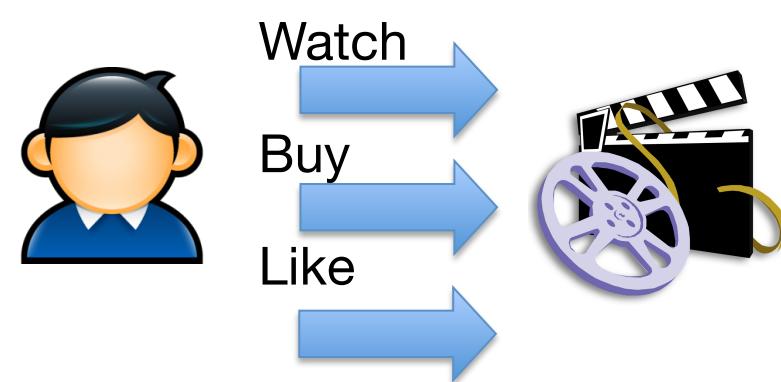


Collaborative filtering

Movies

	Star Wars	Titanic	Blade Runner
User 1	5	2	4
User 2	1	4	2
User 3	5	?	?

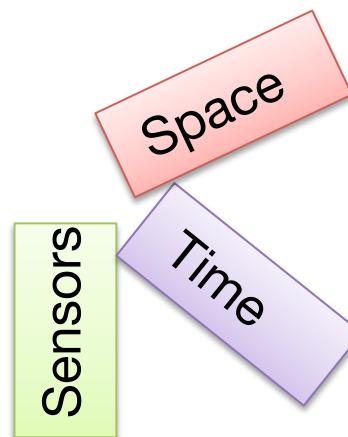
Multiple relations



Tensor decomposition

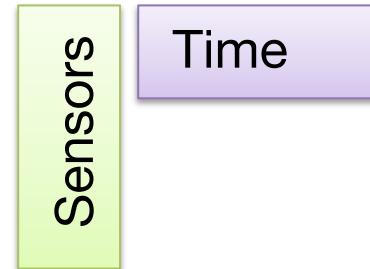
Tensors

Spatio-temporal data

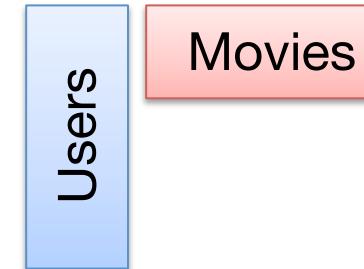


Matrices

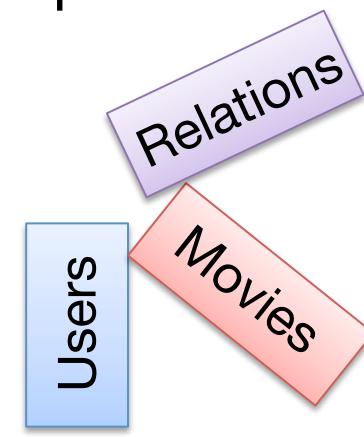
Multivariate time-series



Collaborative filtering

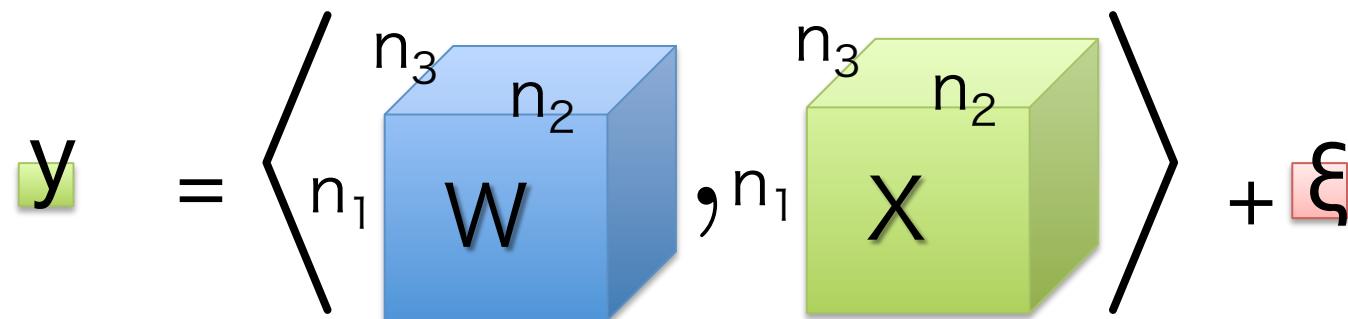


Multiple relations



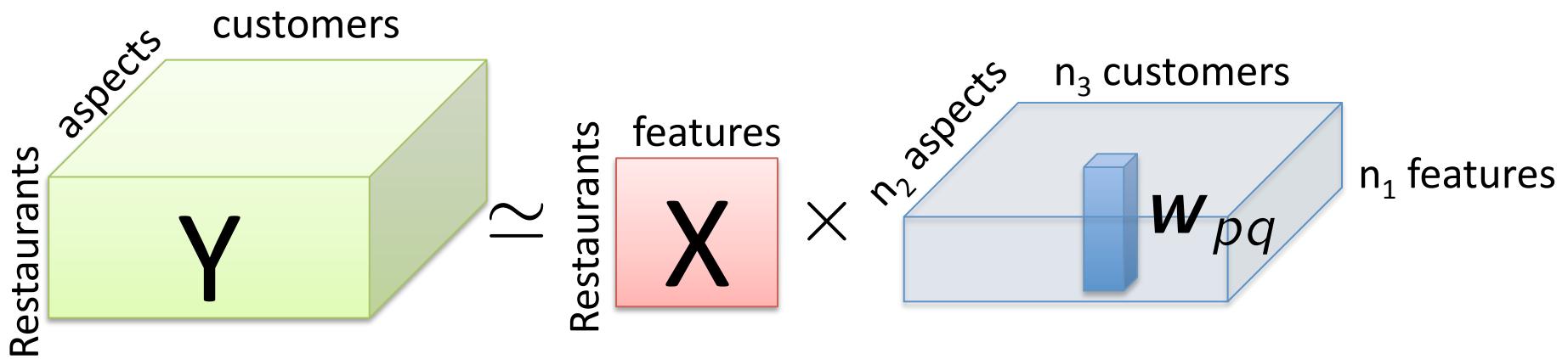
Learning with tensors

- Input: tensor; output: scalar

$$y = \langle \text{W} \rangle_{n_1, n_2, n_3} + \xi$$


– input X may not be low-rank but the weight W may be.

- Multi-task learning with tensors

$$\text{Y} \sim \text{X} \times \text{W}_{pq}$$


Learning objectives

Part 1: Matrices

- What is rank? Different views.
- Nuclear norm and factor-based regularization.
- Learning bounds

Part 2: Tensors

- What is rank? Which view generalizes, and how?
- Nuclear norm? Factor-based regularization?

Part 3: Beyond

- Convex relaxation for Tucker model
- Interaction between computation and statistics

Lecture format

- I will try to be as interactive as possible.
 - Ask questions when you don't understand.
 - Please interrupt.
- I will use iPython Notebooks to give further details.
 - URL:
<https://github.com/ryotat/mlss15/tree/master/python>
 - Python environment?
 - Download Enthought Canopy if you don't

Let's start from matrices

- Rank
- Decomposition
- Regularization
- Theory
 - and see how these generalize (or does not generalize) to tensors

$$\mathbf{A} \in \mathbb{R}^{m \times n}$$

(with $m \leq n$, w.l.o.g.)

Rank

- What are the ranks of these matrices?

$$\begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \cdot (1 \quad 2 \quad 3)$$



Rank 1

$$\begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & 1 \\ 3 & 0 & 3 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 3 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$$



Rank 2

Rank

- The rank of a matrix A

$$= \min R \text{ such that } A = \sum_{r=1}^R u_r v_r^\top$$

$= \max C$ such that A_C is linearly independent
(A_C is the submatrix of A indexed by C)

$=$ number of non-zero singular values

Singular values/vectors

- are solutions of

$$A\mathbf{v} = \sigma\mathbf{u}$$

$$A^\top \mathbf{u} = \sigma\mathbf{v}$$

σ : singular value; (\mathbf{u}, \mathbf{v}) : singular vectors

Fact: for $(\sigma_j, \mathbf{u}_j, \mathbf{v}_j)$ $j=1,\dots,m$, as long as $\sigma_j \neq \sigma_k$

$$\mathbf{u}_j^\top \mathbf{u}_k = 0, \quad \text{and} \quad \mathbf{v}_j^\top \mathbf{v}_k$$

Singular values and eigenvalues

- Eigenvalue equation

$$\begin{bmatrix} 0 & A \\ A^\top & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \sigma \begin{bmatrix} u \\ v \end{bmatrix}$$

Symmetric matrix

$$\begin{bmatrix} 0 & A \\ A^\top & 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} U & U \\ V & -V \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix} \begin{bmatrix} U^\top & V^\top \\ U^\top & -V^\top \end{bmatrix}$$

Negative eigenvalues correspond to flipping the sign of V .

Singular values and eigenvalues

- Eigenvalue equation

$$\begin{bmatrix} 0 & A \\ A^\top & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \sigma \begin{bmatrix} u \\ v \end{bmatrix}$$

Symmetric matrix

$$\begin{bmatrix} 0 & A \\ A^\top & 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} U & U \\ V & -V \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix} \begin{bmatrix} U^\top & V^\top \\ U^\top & -V^\top \end{bmatrix}$$

That is, $A = U\Sigma V^\top$

Algorithmic view

Fixed-point algorithm

1. Initialize u^0 and v^0 randomly
2. For $t=0,1,2,\dots$

$$u^{t+1} = \frac{Av^t}{\|Av^t\|},$$

$$v^{t+1} = \frac{A^\top u^t}{\|A^\top u^t\|}$$

Does this converge? – Yes, essentially power iteration for $\begin{bmatrix} 0 & A \\ A^\top & 0 \end{bmatrix}$

Singular value decomposition

$$A = U\Sigma V^\top$$

$$U^\top U = I_m, \quad V^\top V = I_m, \quad \Sigma = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{pmatrix}$$

Usual convention: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$.

This gives a decomposition

$$A = \sum_{j=1}^m \sigma_j u_j v_j^\top$$

Is this minimal? Yes – keep only $\sigma_j > 0$.

Experiment

- Open USPS.ipynb

Summary: Rank

1. $\min R$ such that $A = \sum_{r=1}^R \mathbf{u}_r \mathbf{v}_r^\top$
2. $\max C$ such that A_C is linearly independent
(A_C is the submatrix of A indexed by C)
3. number of non-zero singular values

Summary: Rank

- Which definition generalizes to tensors?

$$1. \min R \text{ such that } A = \sum_{r=1}^R \mathbf{u}_r \mathbf{v}_r^\top$$

Yes - Rank

2. $\max C$ such that A_C is linearly independent
(A_C is the submatrix of A indexed by C)

Yes – Multilinear rank

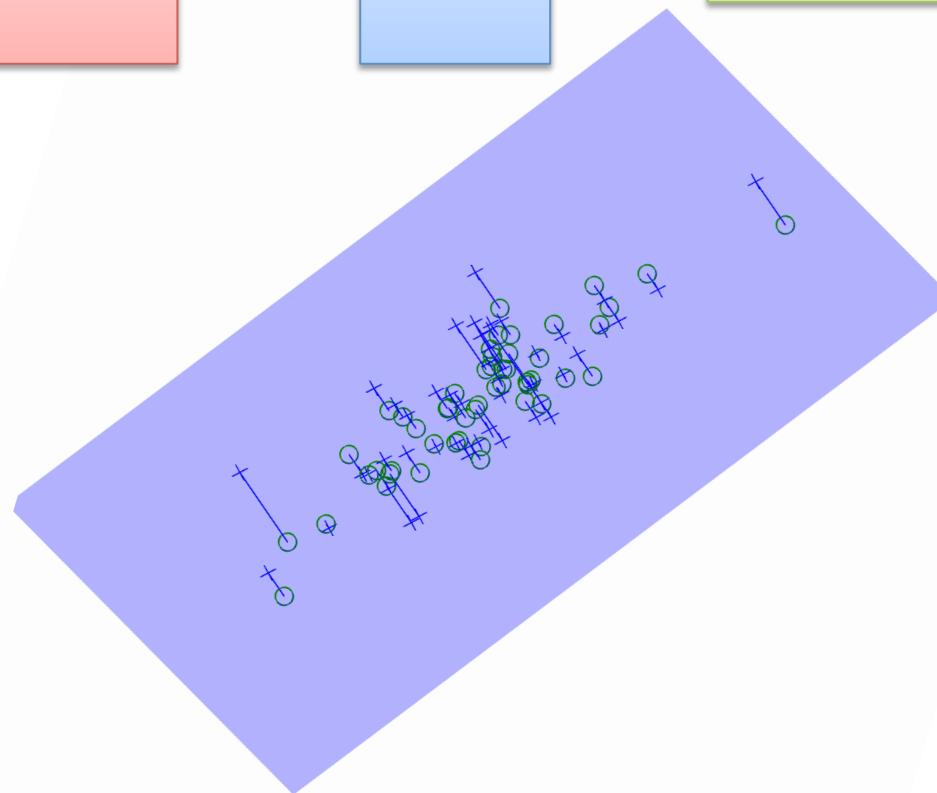
3. number of non-zero singular values

???

Best rank- r approximation

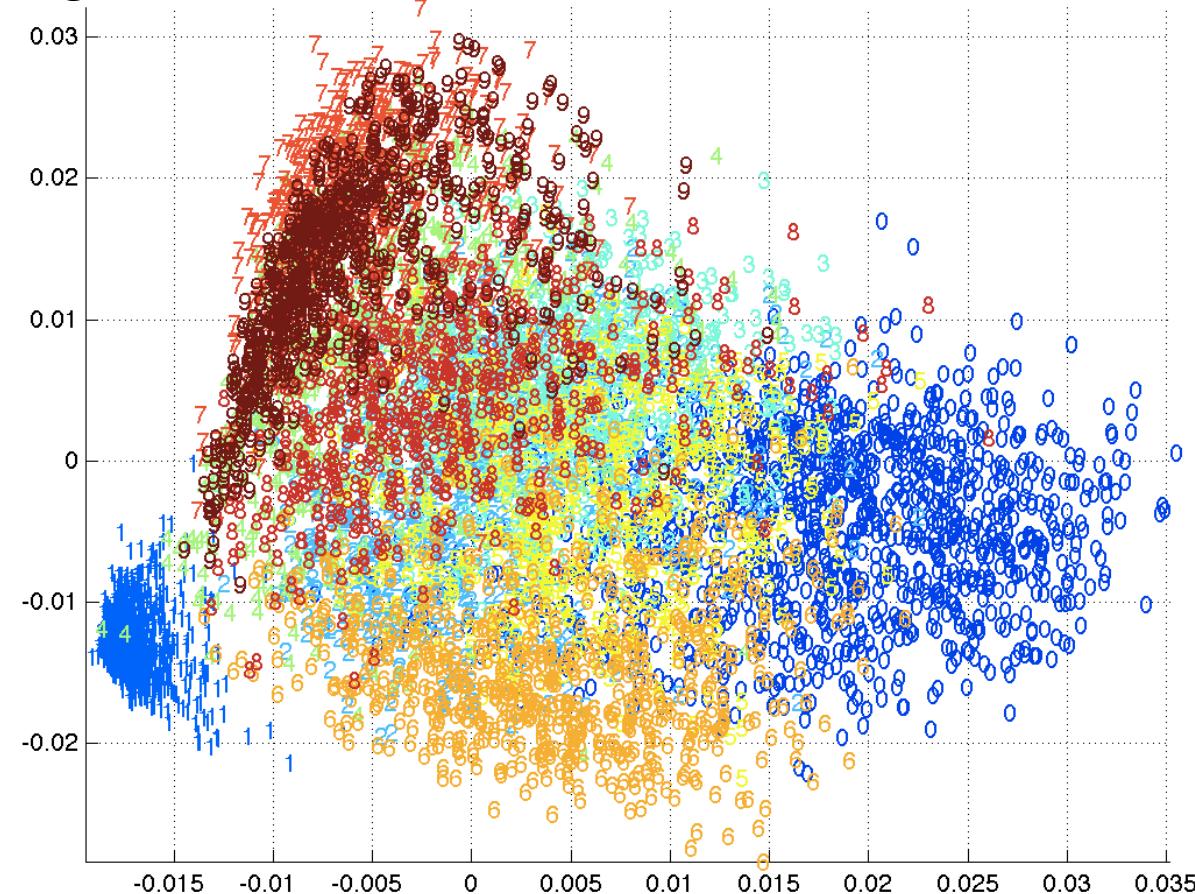
#dimensions $r = \#$ embedding dimensions

$$A \doteq U \times V^T$$



Best rank-2 approximation

USPS digits dataset



Data from: <http://statweb.stanford.edu/~tibs/ElemStatLearn/data.html>

Best rank-r approximation

- Given an $m \times n$ matrix A , we want to find

$$\hat{A}_r = \underset{\substack{\mathbf{X} \in \mathbb{R}^{m \times n}: \\ \text{rank}(\mathbf{X}) \leq r}}{\operatorname{argmin}} \|A - \mathbf{X}\|_F$$

Best rank-1 approximation

- The simplest case

$$\hat{\mathbf{A}}_1 = \underset{\substack{\mathbf{X} \in \mathbb{R}^{m \times n}: \\ \text{rank}(\mathbf{X}) \leq 1}}{\operatorname{argmin}} \| \mathbf{A} - \mathbf{X} \|_F$$

Fact:

$$\hat{\mathbf{A}}_1 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top$$

σ_1 : the largest singular value

Proof

- Let $X = uv^\top$

$$\|A - X\|_F^2 = \|A\|_F^2 - 2u^\top Av + \|u\|_2^2 \cdot \|v\|_2^2$$

Setting the derivative to zero, we have

$$\begin{aligned} Av &= \|v\|_2^2 \cdot u, \\ A^\top u &= \|u\|_2^2 \cdot v \end{aligned} \quad \left. \begin{array}{l} \text{Equations} \\ \text{defining singular} \\ \text{values/vectors} \end{array} \right\}$$

Solution: $u = \sqrt{\sigma_1}u_1, \quad v = \sqrt{\sigma_1}v_1$

Viewing it as a maximization

- Finding the best rank-one approximation

$$\hat{A}_1 = \underset{\substack{\mathbf{X} \in \mathbb{R}^{m \times n}: \\ \text{rank}(\mathbf{X}) \leq 1}}{\operatorname{argmin}} \|A - \mathbf{X}\|_F$$

is equivalent to finding (u, v) that maximizes

$$\sigma_1 = \max_{\substack{\mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^n, \\ \|\mathbf{u}\|=1, \|\mathbf{v}\|=1}} \mathbf{u}^\top A \mathbf{v}$$

Note: both problems are non-convex.

Still admit poly-time algorithm (singular value decomp.)

Spectral norm $\|A\|$

- The largest singular value

$$\sigma_1 = \max_{\substack{\mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^n, \\ \|\mathbf{u}\|=1, \|\mathbf{v}\|=1}} \mathbf{u}^\top A \mathbf{v}$$

is a norm. We denote it by $\|A\| = \sigma_1$

- A norm needs to satisfy

1) Positive homogeneity $\|\alpha A\| = |\alpha| \cdot \|A\| \quad (\forall \alpha \in \mathbb{R})$

2) Subadditivity $\|A + B\| \leq \|A\| + \|B\|$

3) Separation $\|A\| = 0 \quad \Rightarrow \quad A = 0$

Best rank r approximation

- For $r > 1$, $\hat{\mathbf{A}}_r = \mathbf{U}_r \Sigma_r \mathbf{V}_r^\top$

(take the first r columns of \mathbf{U} and \mathbf{V} and corresponding singular values)

PSYCHOMETRIKA—VOL. 1, NO. 3
SEPTEMBER, 1936

THE APPROXIMATION OF ONE MATRIX BY ANOTHER OF LOWER RANK

CARL ECKART AND GALE YOUNG
University of Chicago, Chicago, Illinois



Best rank r approximation

- For $r > 1$, $\hat{\mathbf{A}}_r = \mathbf{U}_r \boldsymbol{\Sigma}_r \mathbf{V}_r^\top$

(take the first r columns of \mathbf{U} and \mathbf{V} and corresponding singular values)

PSYCHOMETRIKA—VOL. 1, NO. 3
SEPTEMBER, 1936

Implication: Subtracting the best rank-one approximation reduces the rank by one.
Does this generalize to tensors? – No!

Low-rank regression

- Input $x \in \mathbb{R}^d$, output $y \in \mathbb{R}^C$

$$y = \mathbf{W}^\top x + \xi$$

where $\mathbf{W} \in \mathbb{R}^{d \times C}$, $\text{rank}(\mathbf{W}) \leq r$

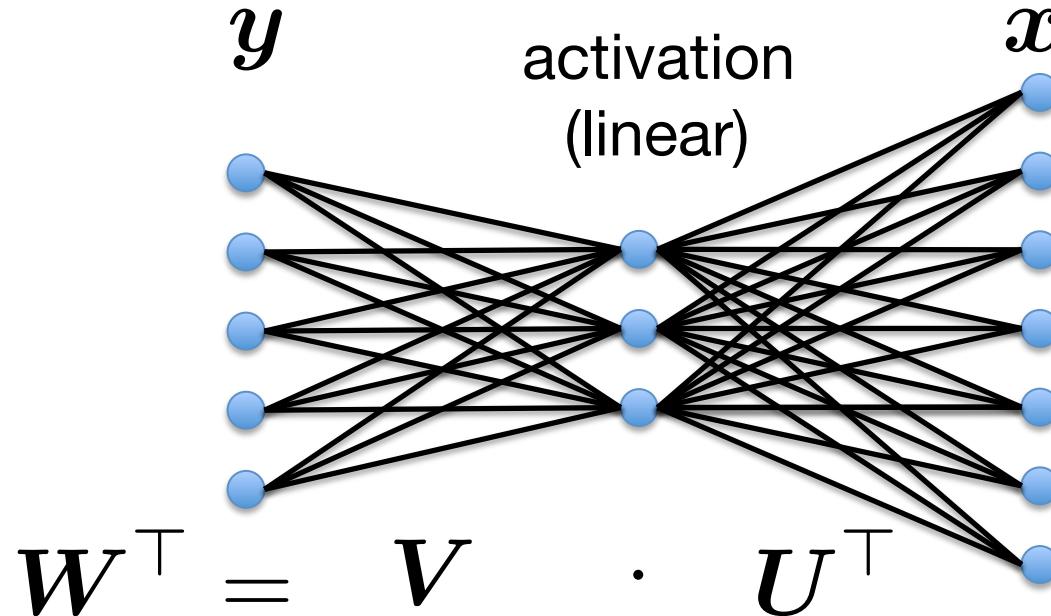
- Given training inputs (x_i, y_i) , $i=1, \dots, m$

$$\hat{\mathbf{W}} = \underset{\substack{\mathbf{W} \in \mathbb{R}^{d \times C}: \\ \text{rank}(\mathbf{W}) \leq r}}{\operatorname{argmin}} \|Y - \mathbf{W}^\top X\|_F^2$$

$$Y = [y_1, \dots, y_m], \quad X = [x_1, \dots, x_m]$$

Connection to linear neural network

- Linear neural network (Baldi & Hornik, 1989)



Matrix factorization = Linear neural network

Analytic solution

$$\mathbf{W}^\top = \mathbf{P}_V \boldsymbol{\Sigma}_{XY}^{-\top} \boldsymbol{\Sigma}_{XX}^{-1}$$

$$\mathbf{P}_V = \mathbf{U}_r \mathbf{U}_r^\top,$$

\mathbf{U}_r : the top r s.v. of $\boldsymbol{\Sigma}_{XY}^{-\top} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY}$

$$\boldsymbol{\Sigma}_{XX} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top, \quad \boldsymbol{\Sigma}_{XY} = \sum_{i=1}^m \mathbf{x}_i \mathbf{y}_i^\top$$

Neural Networks and Principal Component Analysis: Learning from Examples Without Local Minima

PIERRE BALDI AND KURT HORNIK*

University of California, San Diego

(Received 18 May 1988; revised and accepted 16 August 1988)

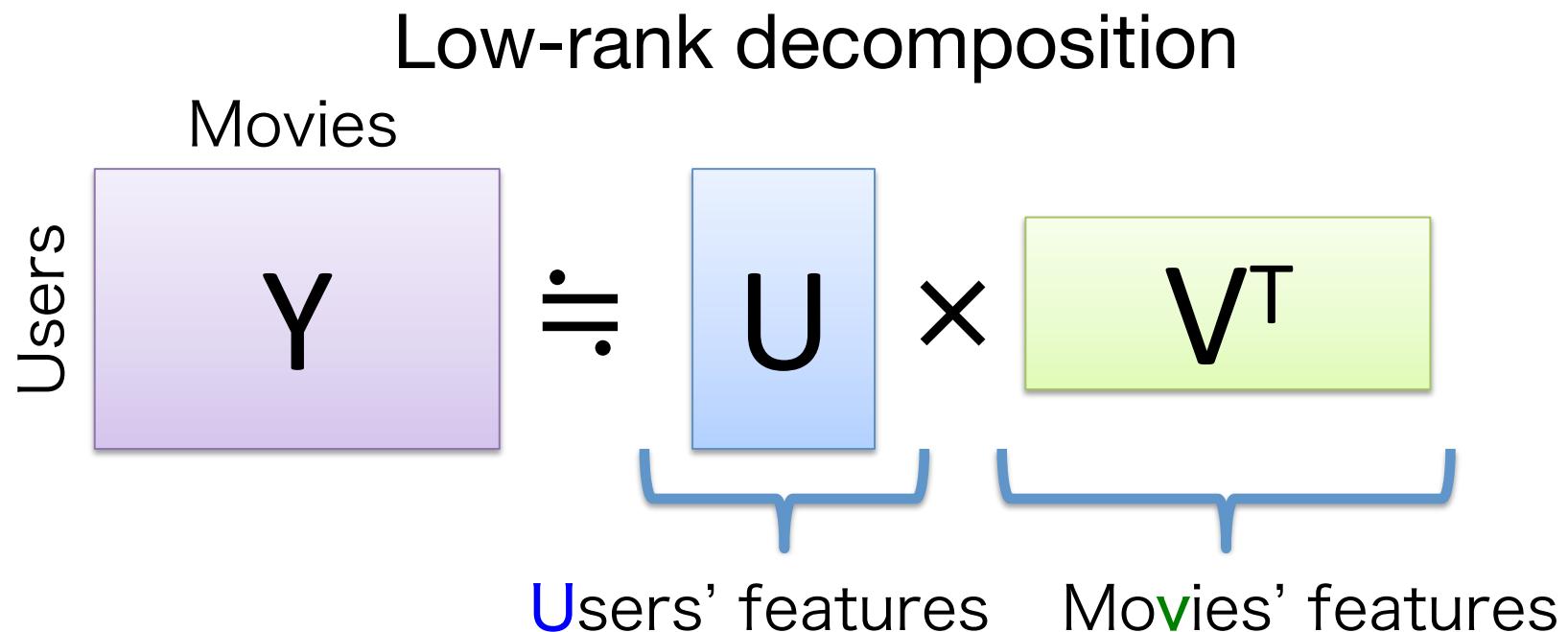
Matrix completion

- “Netflix problem” – Fill in the missing entries

$$Y = \begin{array}{c} \begin{matrix} & \text{Star wars} & \text{Titanic} & \text{Blade runner} & \text{Lost in} \\ \text{User 1} & 5 & 2 & 4 & 5 \\ \text{User 2} & 1 & 4 & 2 & 4 \\ \text{User 3} & 5 & 1 & ? & ? \\ \text{User 4} & ? & ? & ? & 4 \end{matrix} \\ \hline \end{array}$$

Matrix completion

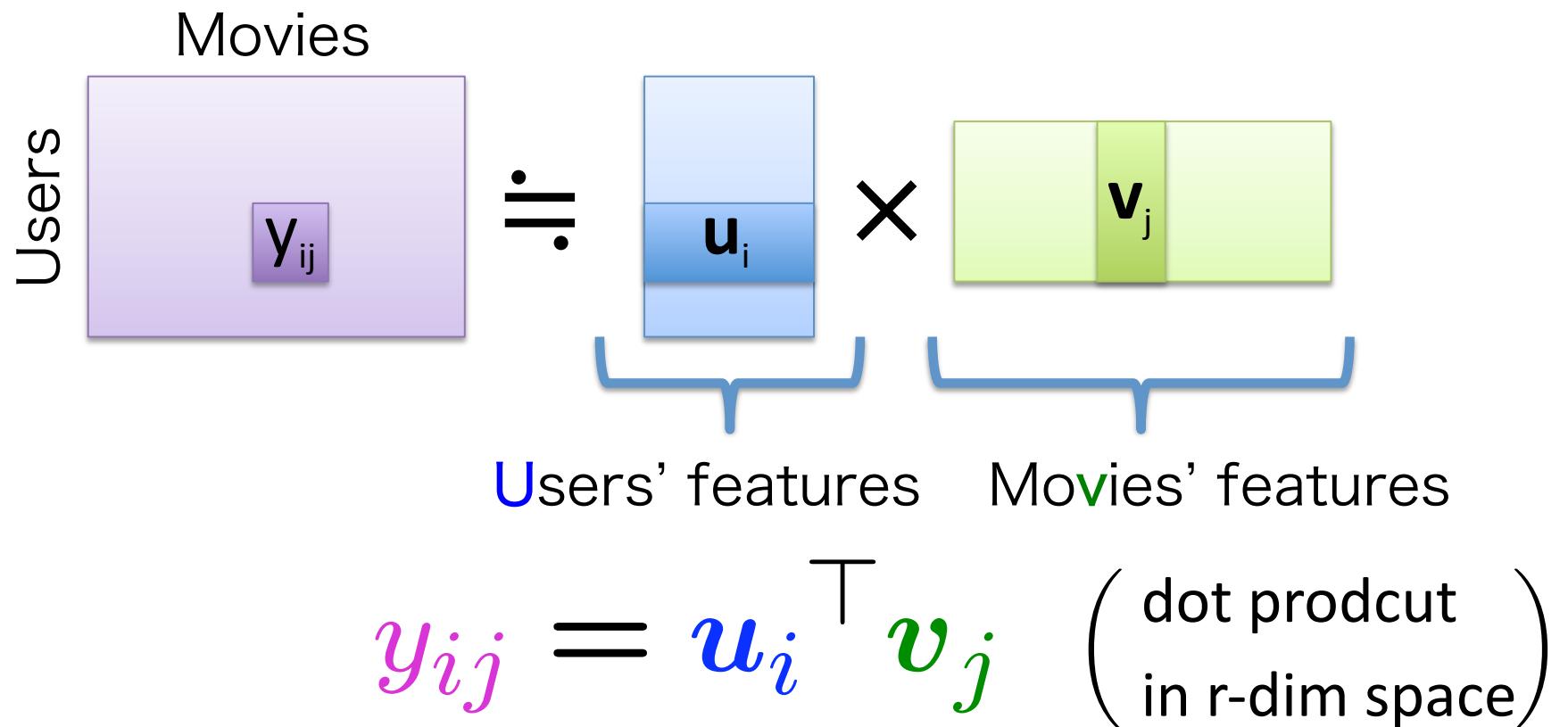
- Impossible without an assumption. (Missing entries can be arbitrary) --- problem is ill-posed
- Most common assumption:



Matrix completion

- Most common assumption:

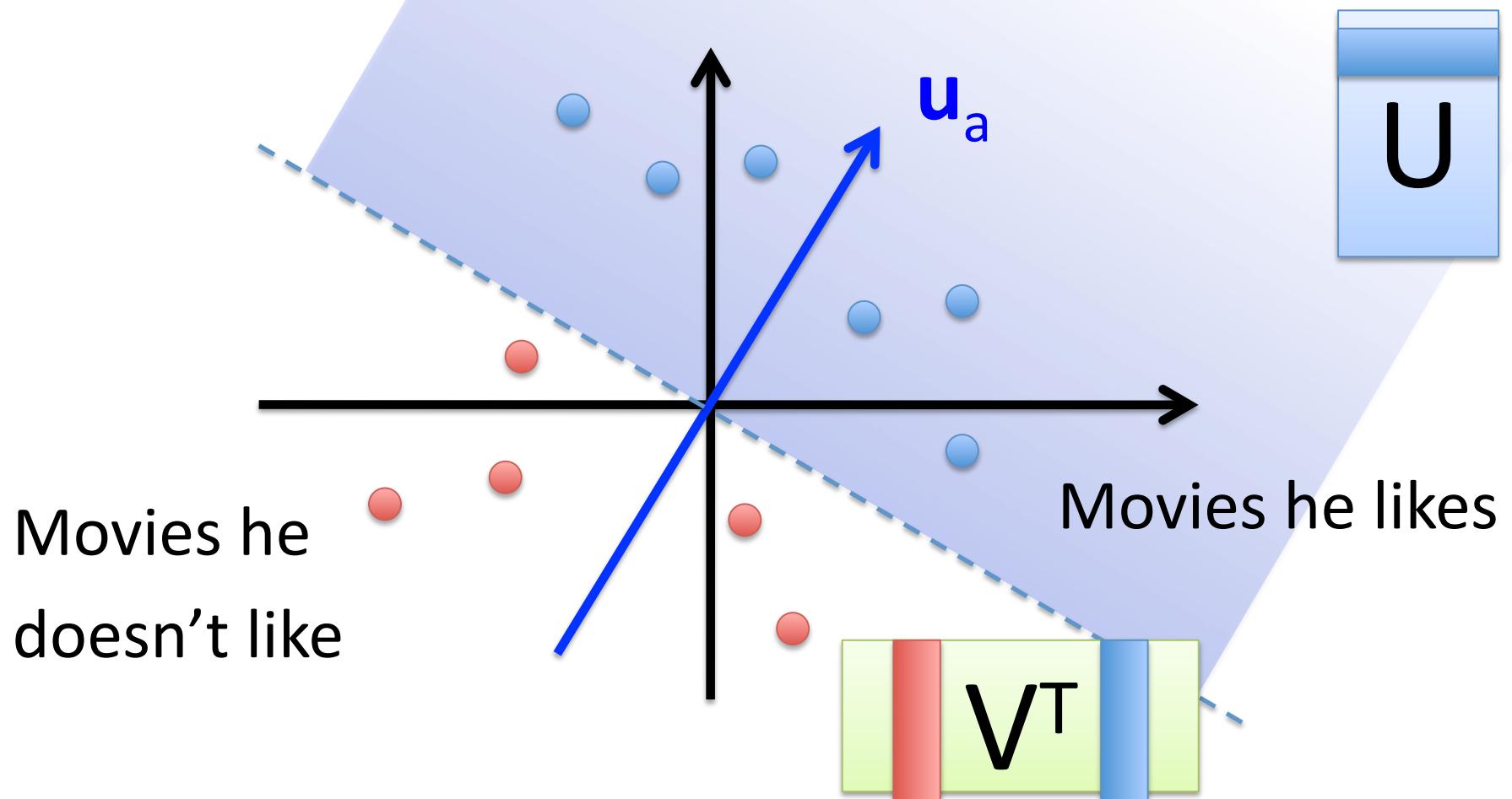
Low-rank decomposition (rank r)



Geometric Intuition

r-dimensional space

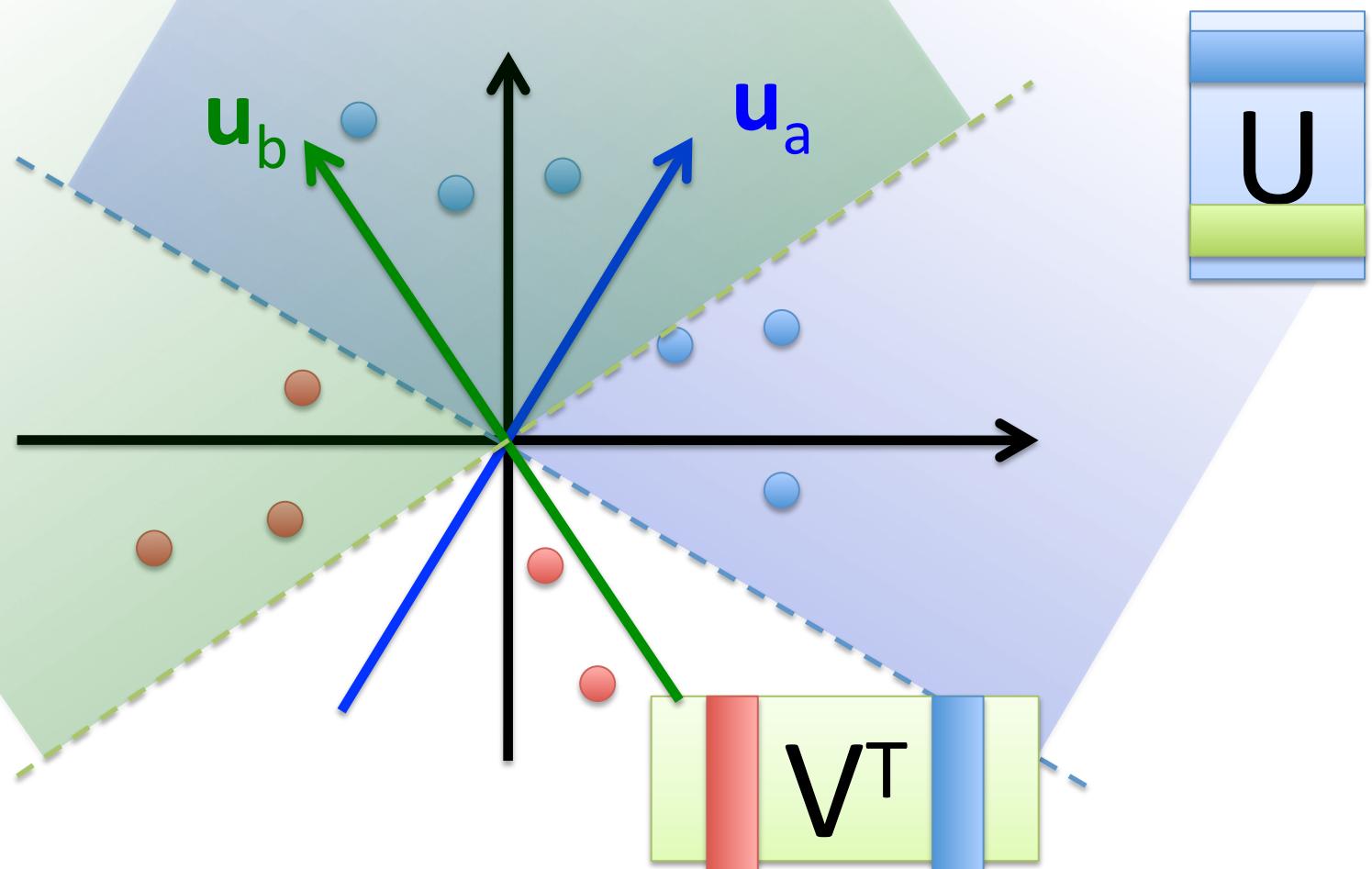
(r: the rank of the decomposition)



Geometric Intuition

r-dimensional space

(r: the rank of the decomposition)



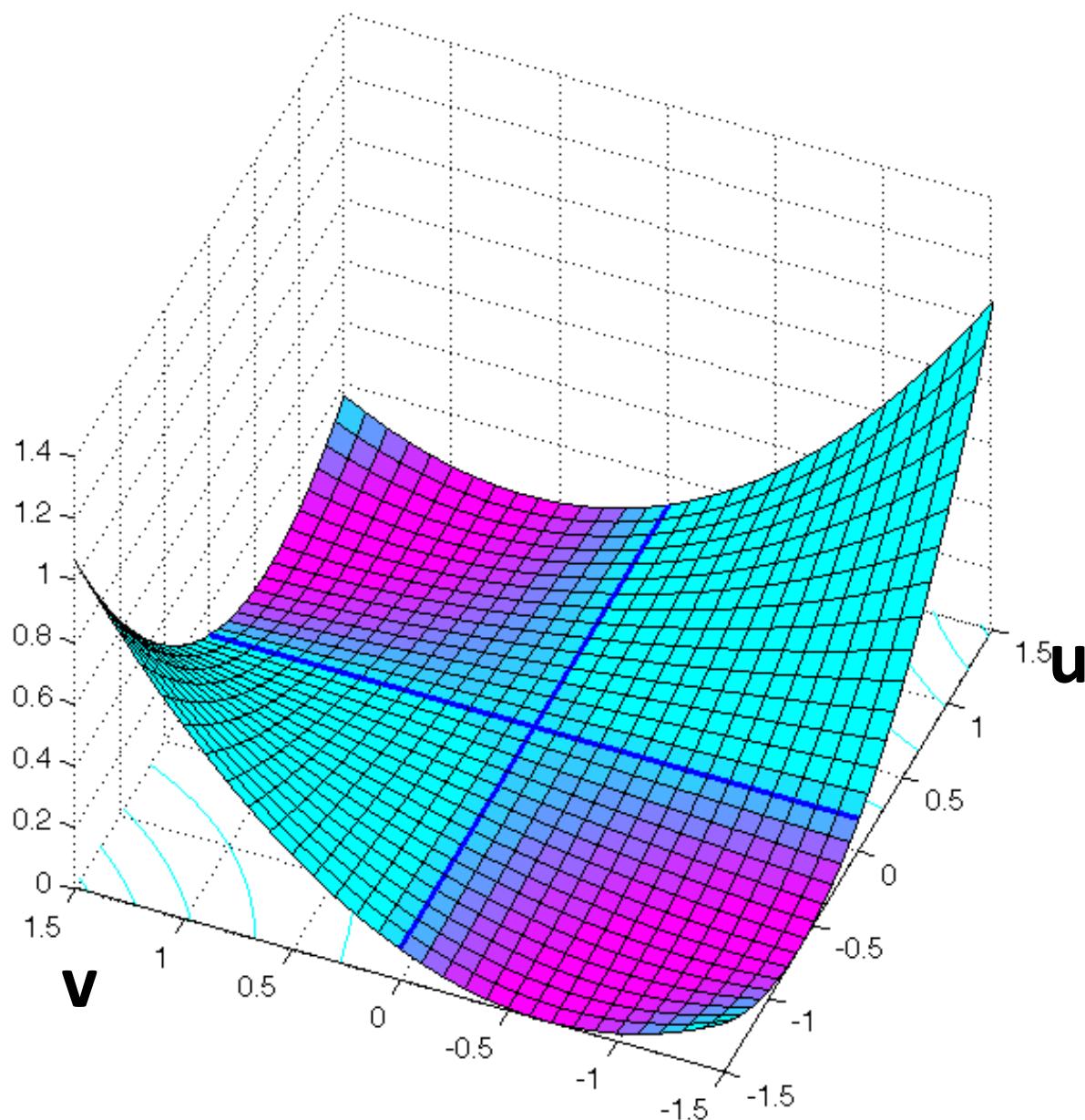
Netflix problem –best rank- r approximation from partial observation

$$\underset{\mathbf{U} \in \mathbb{R}^{m \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Y_{i,j} - \mathbf{u}_i^\top \mathbf{v}_j)^2$$

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}_1^\top \\ \mathbf{u}_2^\top \\ \vdots \\ \mathbf{u}_m^\top \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix}$$

Note: \mathbf{u}_i and \mathbf{v}_j are rows of \mathbf{U} and \mathbf{V} , respectively.

Non-convexity



Regularization

- Add L2 regularization terms

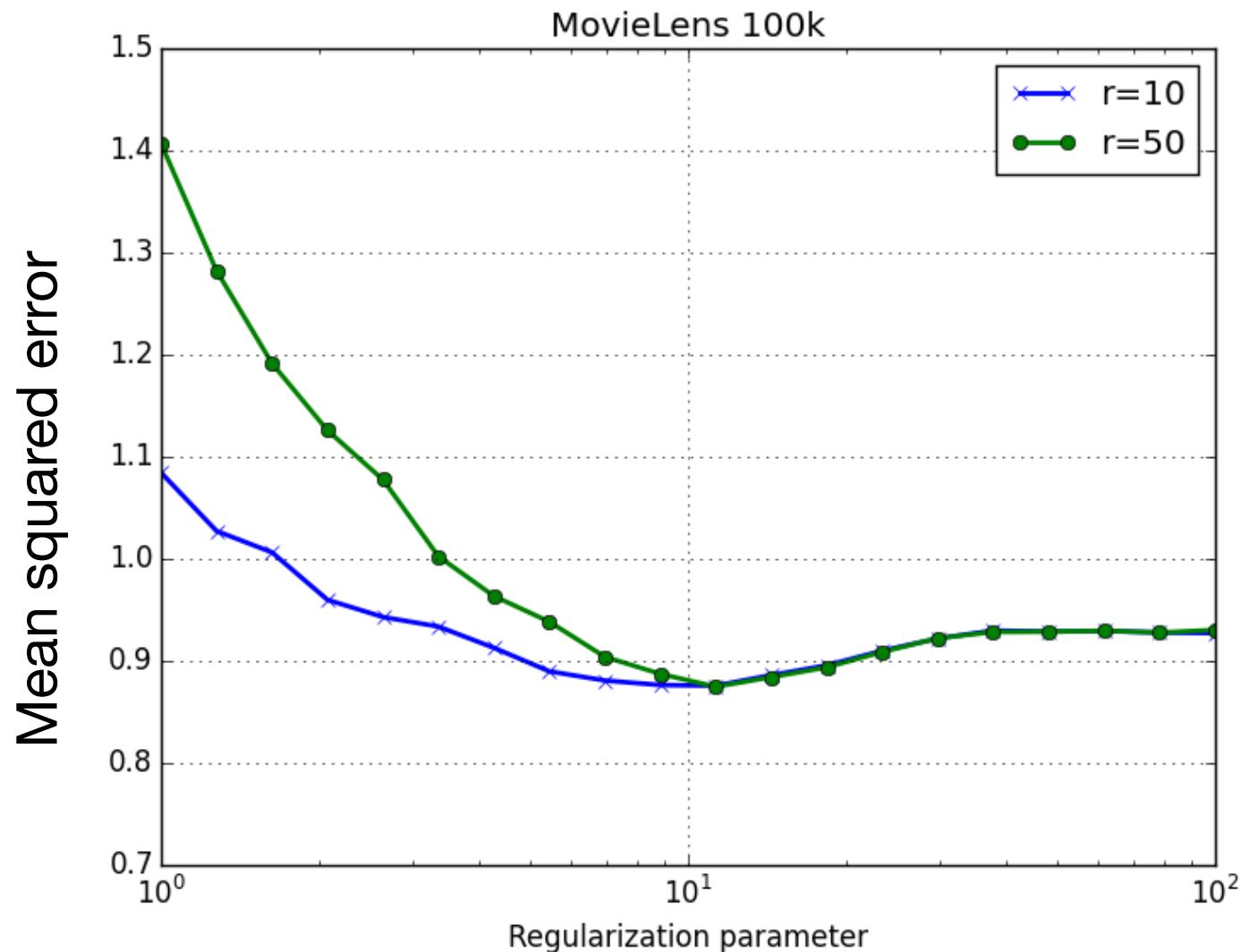
$$\underset{\substack{\mathbf{U} \in \mathbb{R}^{m \times r}, \\ \mathbf{V} \in \mathbb{R}^{n \times r}}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Y_{i,j} - \mathbf{u}_i^\top \mathbf{v}_j)^2 + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$$

How do we minimize this? –Gradient descent!

Let's try it with $r=10$ and $r=50$

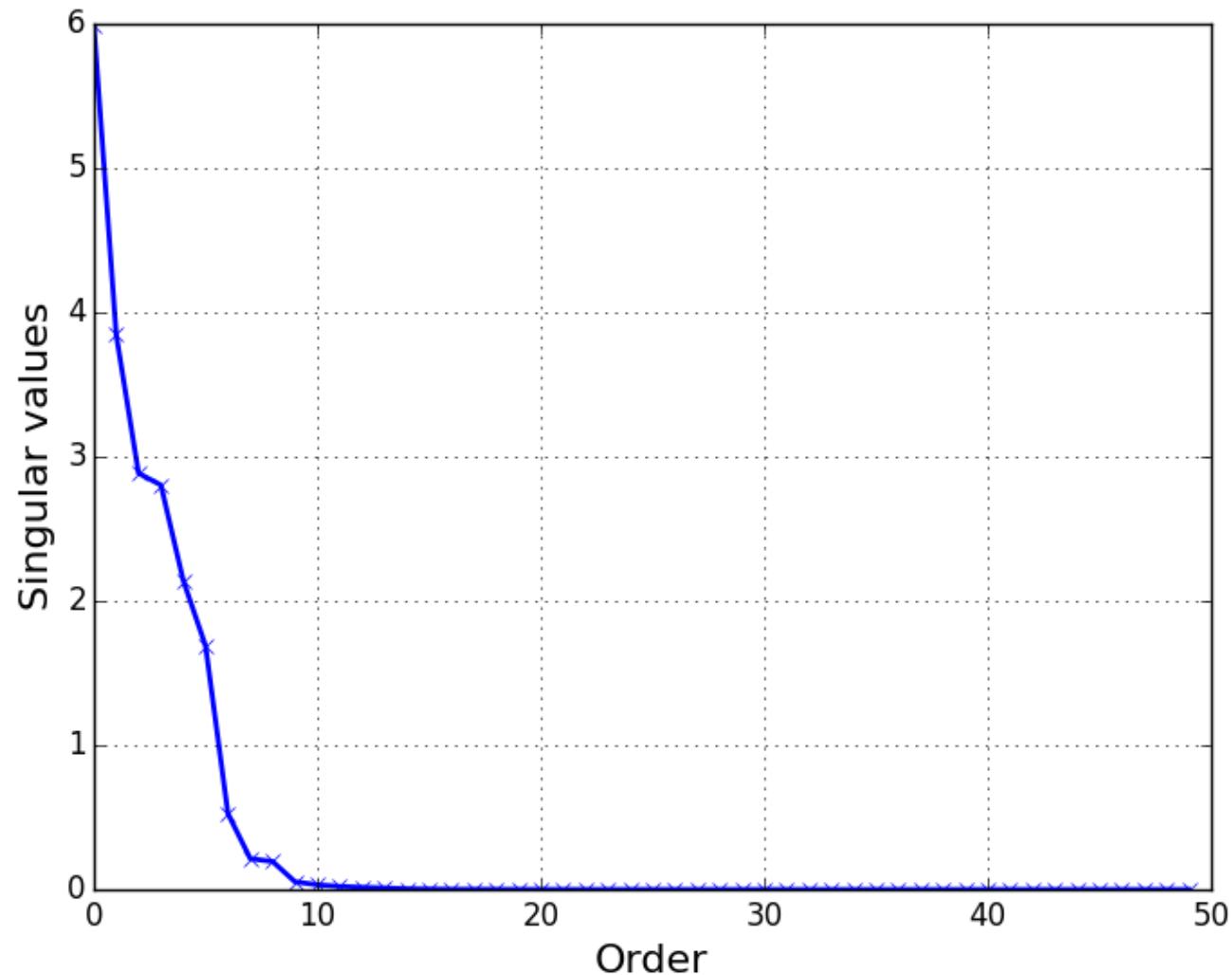
Results for $r=10$ and $r=50$

are surprisingly similar



Singular values

$\lambda = 23.4$



What is going on?

- Define $\mathbf{W} = \mathbf{U}\mathbf{V}^\top = (\mathbf{u}_i^\top \mathbf{v}_j)_{i,j}$
- Loss term: only depends on the product \mathbf{W} .
- Regularizer: let's define

$$g(\mathbf{W}) = \min_{\substack{\mathbf{U} \in \mathbb{R}^{m \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}: \\ \mathbf{W} = \mathbf{U}\mathbf{V}^\top}} \frac{1}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$$

Note: this is a function of \mathbf{W} (\mathbf{U} and \mathbf{V} are minimized-out)!

Effective regularizer

- Now our objective is

$$\underset{\mathbf{W} \in \mathbb{R}^{m \times n}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Y_{i,j} - W_{i,j})^2 + \lambda g(\mathbf{W})$$

with
$$g(\mathbf{W}) = \min_{\substack{\mathbf{U} \in \mathbb{R}^{m \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}: \\ \mathbf{W} = \mathbf{U}\mathbf{V}^\top}} \frac{1}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$$

Is this a convex minimization problem in \mathbf{W} ?

Convexity

Claim: $g(\mathbf{w})$ is a norm (and therefore convex), if r is large enough.

Proof:

1. Positive homogeneity:

(for $\mathbf{W}' = \alpha \mathbf{W}$, take $\mathbf{U}' = \text{sign}(\alpha) \sqrt{|\alpha|} \mathbf{U}$ and $\mathbf{V}' = \sqrt{|\alpha|} \mathbf{V}$.)

2. Triangular inequality:

(for $\mathbf{W} = \mathbf{X} + \mathbf{Y}$, take $\mathbf{U} = [\mathbf{U}_X, \mathbf{U}_Y]$ and $\mathbf{V} = [\mathbf{V}_X, \mathbf{V}_Y]$.)

3. Separation:

$g(\mathbf{W}) = 0$ implies $\mathbf{U}=\mathbf{V}=0$, which implies $\mathbf{W}=0$.

Nuclear norm

- The norm $\|\mathbf{W}\|_* = \min_{\mathbf{W}=\mathbf{U}\mathbf{V}^\top} \frac{1}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$

is called the **nuclear norm**.

We use the norm notation instead of $g(\mathbf{W})$ because we know that it is a norm.

- Moreover,

$$\|\mathbf{W}\|_* = \max_{\substack{\mathbf{X} \in \mathbb{R}^{m \times n}, \\ \|\mathbf{X}\| \leq 1}} \langle \mathbf{X}, \mathbf{W} \rangle$$

Spectral norm

$$\langle \mathbf{X}, \mathbf{W} \rangle = \sum_{j=1}^m \sigma_j$$

Linear sum of s.v.

Proof

- Easy to see that $\|W\|_* \leq \sum_{j=1}^m \sigma_j \leq \max_{\|\mathbf{X}\| \leq 1} \langle \mathbf{X}, W \rangle$
(given $W = \tilde{\mathbf{U}}\Sigma\tilde{\mathbf{V}}^\top$, take $\mathbf{U} = \tilde{\mathbf{U}}\sqrt{\Sigma}$, $\mathbf{V} = \tilde{\mathbf{V}}\sqrt{\Sigma}$, and $\mathbf{X} = \tilde{\mathbf{U}}\tilde{\mathbf{V}}^\top$.)
- For any $W = \sum_{j=1}^m \mathbf{u}_j \mathbf{v}_j^\top$
$$\begin{aligned}\langle \mathbf{X}, W \rangle &= \sum_{j=1}^m \mathbf{u}_j^\top \mathbf{X} \mathbf{v}_j \\ &\leq \sum_{j=1}^m \|\mathbf{u}_j\|_2 \cdot \|\mathbf{v}_j\|_2 \quad (\text{because } \|\mathbf{X}\| \leq 1) \\ &\leq \frac{1}{2} \sum_{j=1}^m (\|\mathbf{u}_j\|_2^2 + \|\mathbf{v}_j\|_2^2) \quad (\text{AM-GM})\end{aligned}$$

Nuclear norm and spectral norm

- Nuclear norm and spectral norm are **dual**

$$\|\mathbf{W}\|_* = \sum_{j=1}^m \sigma_j, \quad \|\mathbf{X}\| = \max_j \sigma_j$$

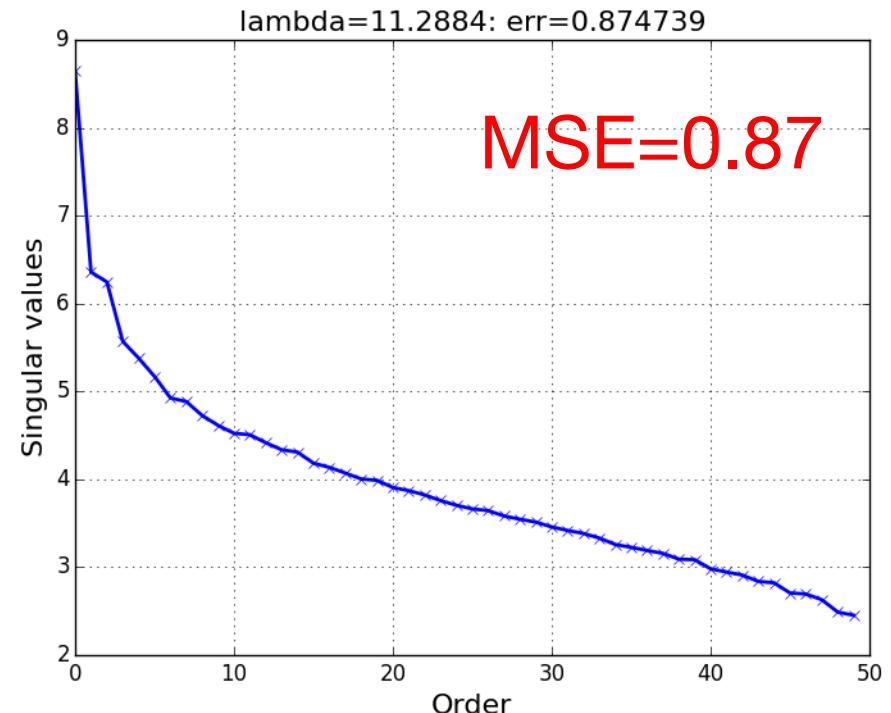
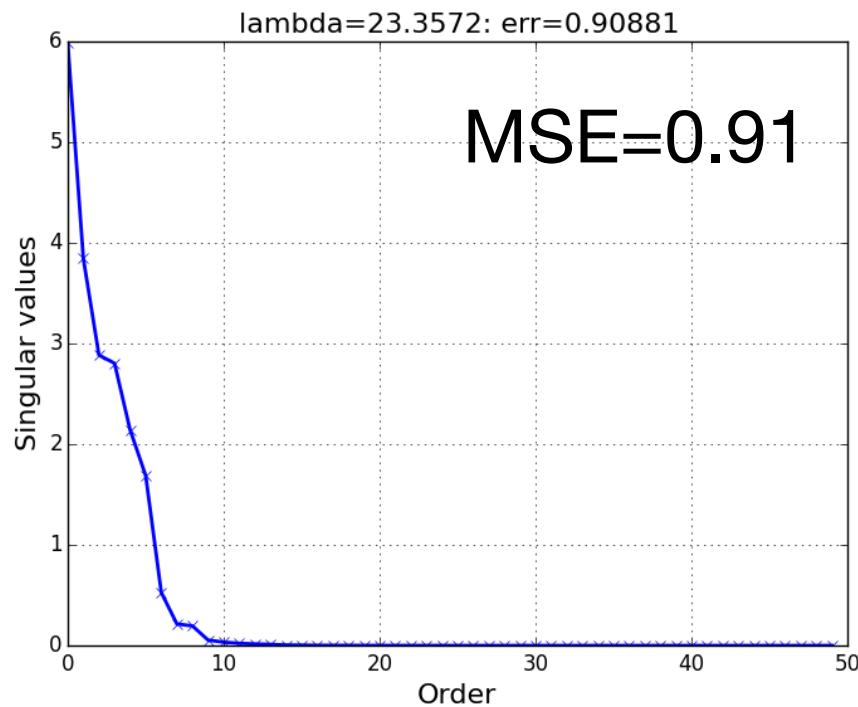
to each other. Just like L_1 and L_∞

$$\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|, \quad \|\mathbf{x}\|_\infty = \max_j |x_j|$$

When do we use the L1 norm? What can we expect from using the nuclear norm?

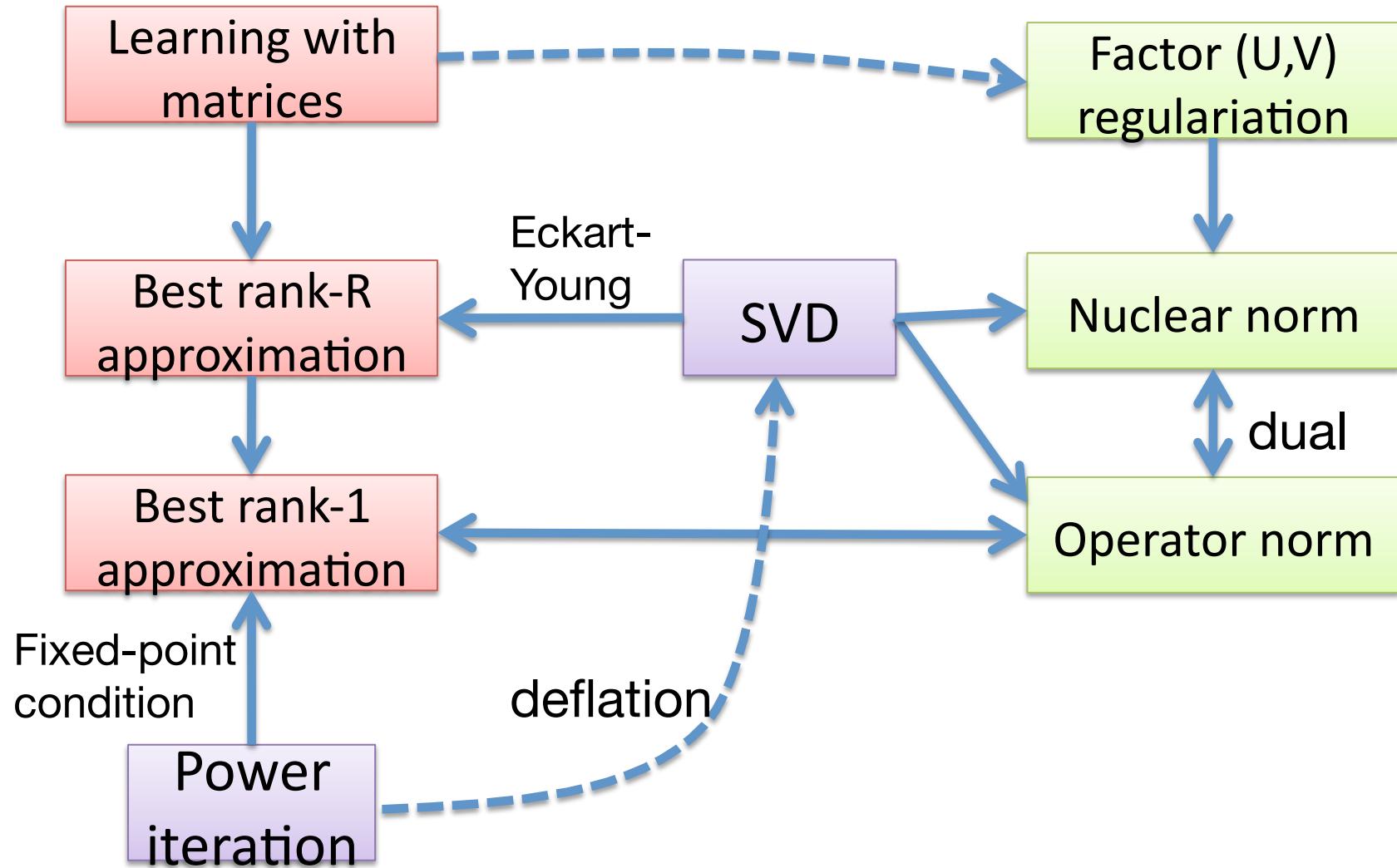
Nuclear norm and sparsity

- Nuclear norm promotes W to be **low rank**



but maybe that's not the end of the story...

Summary for matrix decomposition



Little bit of theory

VC dimension

Rademacher complexity

Complexity of low-rank matrices

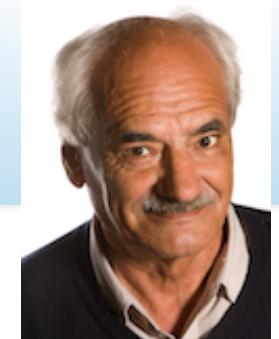
- Let $f_{\mathbf{U}, \mathbf{V}, b}(\mathbf{X}) = \langle \mathbf{X}, \mathbf{U}\mathbf{V}^\top \rangle + b$ and define

$$\mathcal{F}_r = \{f_{\mathbf{U}, \mathbf{V}, b} : \mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}, b \in \mathbb{R}\}$$

- Examples
 - Low-rank regression: $\mathbf{X} = \mathbf{x} \cdot \mathbf{e}_j^\top$
 - Matrix completion: $\mathbf{X} = \mathbf{e}_i \cdot \mathbf{e}_j^\top$
 - Matrix classification: general \mathbf{X} .

What is the complexity of \mathcal{F}_r ? How many samples do we need?

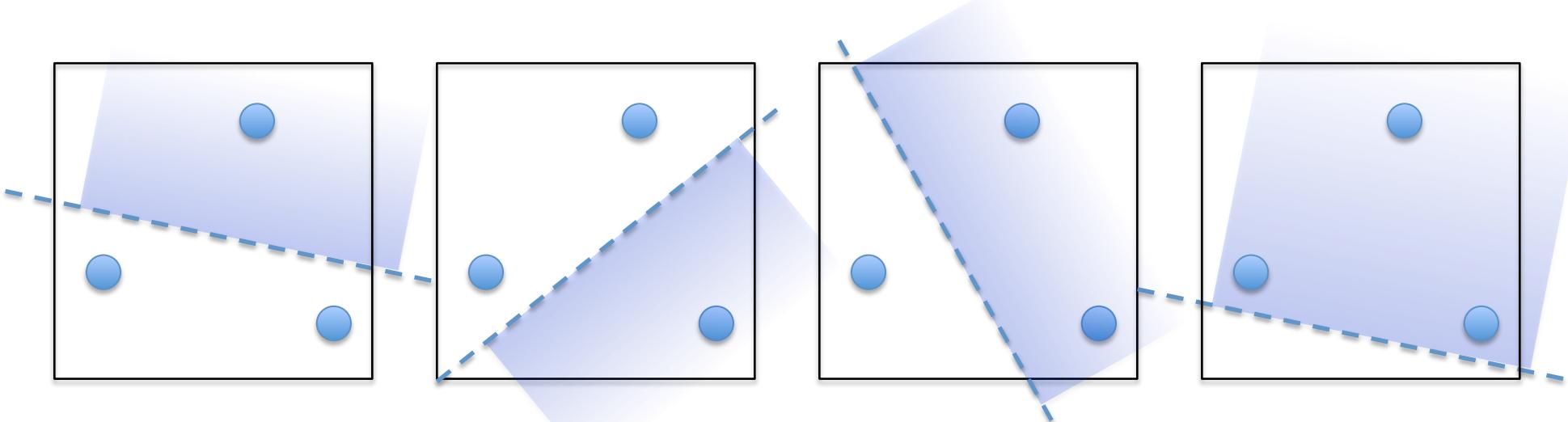
Background



- VC dimension of \mathcal{F} :

$$d(\mathcal{F}) = \max m$$

$$\text{s.t. } |\{(\operatorname{sgn}(f(\mathbf{X}_1)), \dots, \operatorname{sgn}(f(\mathbf{X}_m)) : f \in \mathcal{F}\}| = 2^m$$



Note: standard asymptotic statistics does not work here because the model is not regular (degenerate Fisher information).

Bound on the VC dimension

$$d(\mathcal{F}_r) \leq (r(n_1 + n_2) + 1) \log(r(n_1 + n_2) + 1)$$

For sufficiently large n_1 and n_2 . No more than #parameters!

- This implies that for classification

$$R(f) - R_{\text{emp}}(f) \leq \tilde{O} \left(\sqrt{\frac{(r(n_1 + n_2) + 1)}{m}} \right)$$

for any $f \in \mathcal{F}_r$ with high probability (ignoring the log term)

where $R(f) = \Pr(y \neq f(\mathbf{X}))$, $R_{\text{emp}}(f) = \frac{1}{m} \sum_{i=1}^m I(y_i \neq f(\mathbf{X}_i))$

- Similar bound for real-valued output using **pseudodimension**.

Proof of the VC dimension

- Warren's theorem

THEOREM 3. *Let p_1, \dots, p_m be real polynomials in n variables, each of degree at most $d \geq 1$. If $m \geq n$, the number of sign sequences $\text{sgn } \mathbf{p}(\mathbf{x}) = (\text{sgn } p_1(\mathbf{x}), \dots, \text{sgn } p_m(\mathbf{x}))$ that consist of terms $+1, -1$ does not exceed $(4edm/n)^n$.*

Warren (1968) “Lower bounds for approximation by nonlinear manifolds”

Note that

$$f_{\mathbf{U}, \mathbf{V}, b}(\mathbf{X}_i) = \langle \mathbf{X}_i, \mathbf{U}\mathbf{V}^\top \rangle + b$$

is a polynomial of degree at most 2 in $r(n_1+n_2)+1$ variables.

Note: (U, V, b) are variables and X_i are coefficients.

Low-nuclear norm matrices

- Let $f_W(\mathbf{X}) = \langle \mathbf{X}, \mathbf{W} \rangle$ and define

$$\mathcal{F}_W = \{f_W : \mathbf{W} \in \mathbb{R}^{n_1 \times n_2}, \|\mathbf{W}\|_* \leq W\}$$

(no bias term for simplicity)

- Examples

– Low-rank regression: $\mathbf{X} = \mathbf{x} \cdot \mathbf{e}_j^\top$

– Matrix completion: $\mathbf{X} = \mathbf{e}_i \cdot \mathbf{e}_j^\top$

– Matrix classification: general \mathbf{X} .

Low-nuclear norm matrices

- Let $f_{\mathbf{W}}(\mathbf{X}) = \langle \mathbf{X}, \mathbf{W} \rangle$ and define

$$\mathcal{F}_{\mathbf{W}} = \{f_{\mathbf{W}} : \mathbf{W} \in \mathbb{R}^{n_1 \times n_2}, \|\mathbf{W}\|_* \leq W\}$$

(no bias term for simplicity)

- Rademacher complexity

$$\hat{\mathcal{R}}_n(\mathcal{F}) = \mathbb{E}_{\boldsymbol{\sigma}} \sup_{f \in \mathcal{F}} \left| \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{X}_i) \right|$$

(σ_i takes +1 or -1 with probability 1/2)

- measures how well F can fit +1/-1 noise.
- is **scale sensitive**.

Rademacher complexity bound

[see Kakade, Shalev-Shwartz, Tewari 12]

If $\|X_i\| \leq X$ (for $i = 1, \dots, m$)

Spectral norm

$$\hat{\mathcal{R}}_m(\mathcal{F}_W) \leq \sqrt{\frac{2}{m} \log((n_1 + n_2)e)} \cdot X \cdot W.$$

Consequence: for any Lipschitz continuous loss function ℓ ,

$$\mathbb{E} \ell(y, f(\mathbf{X})) - \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(\mathbf{X}_i)) \leq O \left(\sqrt{\frac{X^2 W^2 \log(n_1 + n_2)}{m}} \right)$$

How large are X and W ? If \mathbf{X}_i ($i=1, \dots, m$) are drawn from standard normal distribution, $X^2 = O(n_1 + n_2)$, and $W^2 = O(r)$, if \mathbf{W} has a constant Frobenius norm.

Proof of RC bound

- From the duality btwn. the spectral & nuclear norms

$$\frac{1}{m} \sum_{i=1}^m \sigma_i \langle \mathbf{X}_i, \mathbf{W} \rangle \leq \left\| \frac{1}{m} \sum_{i=1}^m \sigma_i \mathbf{X}_i \right\| \cdot \|\mathbf{W}\|_*$$

 $\leq W$

- Thus

$$\hat{\mathcal{R}}_m(\mathcal{F}_W) \leq \mathbb{E}_{\boldsymbol{\sigma}} \left\| \frac{1}{m} \sum_{i=1}^m \sigma_i \mathbf{X}_i \right\| \cdot W$$

and using random matrix theory,

$$\mathbb{E}_{\boldsymbol{\sigma}} \left\| \frac{1}{m} \sum_{i=1}^m \sigma_i \mathbf{X}_i \right\| \leq \sqrt{\frac{2}{m} \log((n_1 + n_2)e)} \cdot X$$

if $\|\mathbf{X}_i\| \leq X$ (for $i = 1, \dots, m$).

(Using Corollary 4.2 from Tropp (2011) “User-Friendly Tail Bounds for Sums of Random Matrices”)

Tightening by expectation

[Foygel & Srebro, 2011]

- Previous analysis does not assume any distribution for inputs X_1, \dots, X_m .
- Considering the matrix completion setting and taking expectation over the random draw of the observed positions,

$$\mathcal{R}_m(\mathcal{F}_W) := \mathbb{E} \hat{\mathcal{R}}_m(\mathcal{F}_W) = O \left(\sqrt{\frac{(n_1 + n_2) \log(n_1 + n_2)}{n_1 n_2 m}} \cdot W \right)$$

Expectation over
the draw of the
observed positions

W can be as large as $(rn_1 n_2)^{1/2}$
and sample complexity
 $O(r(n_1 + n_2)/\varepsilon^2)$

Approach from high-dimensional statistics

[Negahban & Wainwright 2011]

- Let the entries of X_1, \dots, X_m drawn i.i.d. from standard Gaussian distribution and **assume**

$$y_i = \langle \mathbf{X}_i, \mathbf{W}^* \rangle + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2), \quad (i = 1, \dots, m)$$

truth (rank r)

then the estimator

$$\hat{\mathbf{W}} = \underset{\mathbf{W} \in \mathbb{R}^{n_1 \times n_2}}{\operatorname{argmin}} \left(\frac{1}{2m} \sum_{i=1}^m (y_i - \langle \mathbf{X}_i, \mathbf{W} \rangle)^2 + \lambda_m \|\mathbf{W}\|_* \right)$$

with $\lambda_m = c \cdot \sigma \sqrt{(n_1 + n_2)/m}$ satisfies

$$\|\hat{\mathbf{W}} - \mathbf{W}^*\|_F^2 \leq O \left(\frac{\sigma^2 r(n_1 + n_2)}{m} \right)$$

with high probability; $r = \operatorname{rank}(\mathbf{W}^*)$.

Summary so far

- Rank and nuclear norm: both measure complexity of a matrix-based hypothesis class.
- L₂ regularization on the factors (U and V) induces nuclear norm regularization on the matrix UV^T .
- Nuclear norm promotes the matrix to be low-rank. However, it can be independently justified even when it doesn't.

Tensors

Rank and multilinear rank

Non-closedness and other weirdness

Tucker decomposition

Graphical representation

Recap: Rank for matrices

1. $\min R$ such that $A = \sum_{r=1}^R \mathbf{u}_r \mathbf{v}_r^\top$
2. $\max C$ such that A_C is linearly independent
(A_C is the submatrix of A indexed by C)
3. number of non-zero singular values

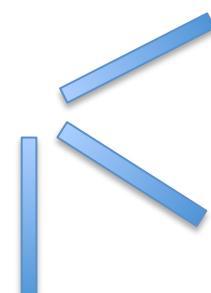
Rank for tensors

- A K th-order tensor

$$\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_K}$$

(we will take $K=3$ below to simplify notation)

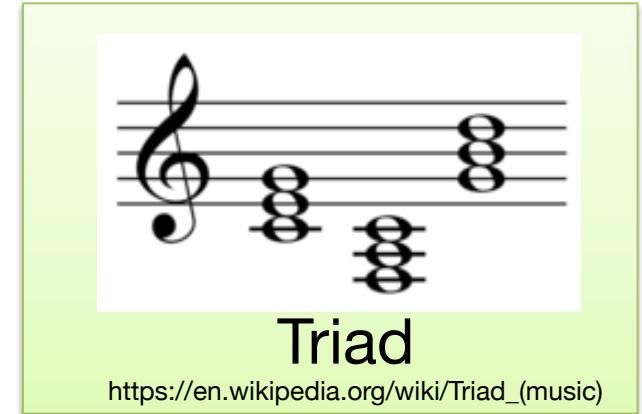
- Definition: Rank is the min R such that

$$\mathcal{A} = \sum_{r=1}^R \underbrace{\mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r}_{\text{outer product of 3 vectors}}$$


Bibliographical note

- The decomposition

$$\mathcal{A} = \sum_{r=1}^R \underbrace{\mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r}_{\text{Triad}}$$



is called **CP (Canonical Polyadic)** decomposition.

[Hitchcock (1927) “The expression of a tensor or a polyadic as a sum of products”]

- Also known as CANDECOMP [Carroll & Chang, 1970] and PARAFAC [Harshman, 1970].

Rank quiz

- What is the rank of this tensor?

$$\mathcal{A} = \begin{bmatrix} 1 & 1 | 1 & 1 \\ 1 & 1 | 1 & 1 \end{bmatrix}$$


1st slice 2nd slice

$$= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



$$\mathcal{A} = \begin{bmatrix} 0 & 1 | 1 & 0 \\ 1 & 0 | 0 & 100 \end{bmatrix}$$

$$= \begin{bmatrix} 0.0995 \\ 0.995 \end{bmatrix} \circ \begin{bmatrix} 0.5025 \\ 5.025 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 10 \end{bmatrix} + \begin{bmatrix} 0.0995 \\ -0.995 \end{bmatrix} \circ \begin{bmatrix} -0.5025 \\ 5.025 \end{bmatrix} \circ \begin{bmatrix} 1 \\ -10 \end{bmatrix}$$



$$\mathcal{A} = \begin{bmatrix} 0 & 1 | 1 & 0 \\ 1 & 0 | 0 & 0 \end{bmatrix}$$



Tensor rank is not closed

- Consider the limit $\alpha \rightarrow \infty$ of a rank 2 tensor

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

$$\lim_{\alpha \rightarrow \infty} \mathcal{Y} = \underbrace{a_1 \circ b_1 \circ c_2}_{\text{rank 1}} + \underbrace{a_1 \circ b_2 \circ c_1}_{\text{rank 1}} + \underbrace{a_2 \circ b_1 \circ c_1}_{\text{rank 1}}$$

- The limit is rank 3 if $[a_1, a_2]$, $[b_1, b_2]$, $[c_1, c_2]$ are linearly independent.

(Example from Kolda & Bader 2009)

Rank of $2 \times 2 \times 2$ tensors

Table 1

Orbits of $2 \times 2 \times 2$ tensors under the action of invertible multilinear transformation (**S**, **T**, **U**) over the real field. The letters *D* and *G* stand for “degenerate” (zero volume set in the eight-dimensional space of $2 \times 2 \times 2$ tensors) and “typical” (positive volume set), respectively.

Canonical form	Tensor rank	Multilinear rank	Sign Δ
$D_0 :$ $\begin{bmatrix} 0 & 0 & & 0 & 0 \\ 0 & 0 & & 0 & 0 \end{bmatrix}$	0	(0,0,0)	0
$D_1 :$ $\begin{bmatrix} 1 & 0 & & 0 & 0 \\ 0 & 0 & & 0 & 0 \end{bmatrix}$	1	(1,1,1)	0
$D_2 :$ $\begin{bmatrix} 1 & 0 & & 0 & 0 \\ 0 & 1 & & 0 & 0 \end{bmatrix}$	2	(2,2,1)	0
$D'_2 :$ $\begin{bmatrix} 1 & 0 & & 0 & 1 \\ 0 & 0 & & 0 & 0 \end{bmatrix}$	2	(1,2,2)	0
$D''_2 :$ $\begin{bmatrix} 1 & 0 & & 0 & 0 \\ 0 & 0 & & 1 & 0 \end{bmatrix}$	2	(2,1,2)	0
$G_2 :$ $\begin{bmatrix} 1 & 0 & & 0 & 0 \\ 0 & 0 & & 0 & 1 \end{bmatrix}$	2	(2,2,2)	+
$D_3 :$ $\begin{bmatrix} 0 & 1 & & 1 & 0 \\ 1 & 0 & & 0 & 0 \end{bmatrix}$	3	(2,2,2)	0
$G_3 :$ $\begin{bmatrix} -1 & 0 & & 0 & 1 \\ 0 & 1 & & 1 & 0 \end{bmatrix}$	3	(2,2,2)	-

Stegeman & Comon (2010) “Subtracting a best rank-1 approximation may increase tensor rank”

Bounding the rank

- A trivial bound

$$\text{rank}(\mathcal{A}) \leq n_1 n_2 n_3$$

(consider each entry as a rank-one tensor)

- A slightly better bound

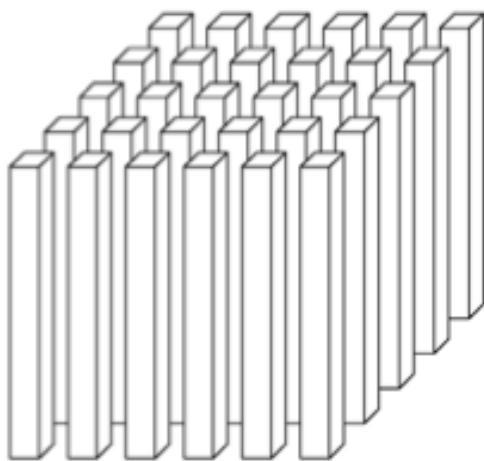
$$\text{rank}(\mathcal{A}) \leq \min(n_1 n_2, n_2 n_3, n_1 n_3)$$

(consider each fiber as a rank-one tensor)

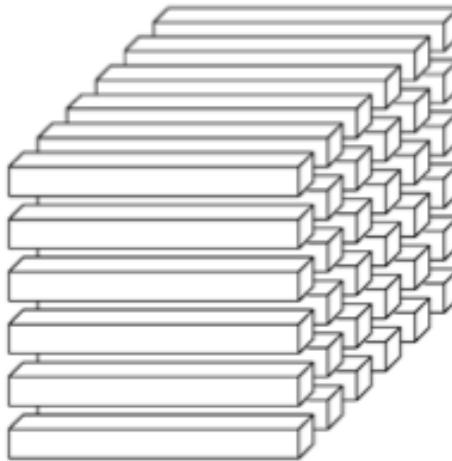
Rank = number of linearly independent columns?

- What is a column vector of a tensor? –**fiber**

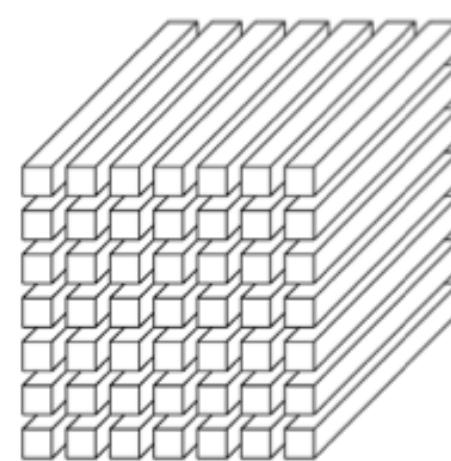
A **mode- k fiber** is a column vector obtained by fixing all but the k th index.



(a) Mode-1 (column) fibers: $\mathbf{x}_{:jk}$



(b) Mode-2 (row) fibers: $\mathbf{x}_{i:k}$



(c) Mode-3 (tube) fibers: $\mathbf{x}_{ij:}$

Fig. 2.1 Fibers of a 3rd-order tensor.

(Kolda & Bader 2009)

Unfolding and mode- k rank

- mode- k unfolding

$$A_{(k)} = \begin{array}{c} \text{mode-}k \text{ fibers} \\ \hline \text{---} \\ \text{---} \end{array} \quad \dots \quad \dots \quad \dots \quad \in \mathbb{R}^{n_k \times \prod_{k' \neq k} n_{k'}}$$

- mode- k rank $\text{rank}_k(\mathcal{A}) = \text{rank}(A_{(k)})$
- Multilinear rank $(\text{rank}_1(\mathcal{A}), \dots, \text{rank}_K(\mathcal{A}))$

Questions

- How do these two notions of rank relate to each other? When are they equal?
- What kind of decomposition does low-multilinear rank imply?

Equivalence when rank=1

Detecting a rank-one tensor is easy!

- Suppose A is rank one

$$\mathcal{A} = \mathbf{u} \circ \mathbf{v} \circ \mathbf{w} = (u_i v_j w_k)_{i,j,k}$$

(again assuming $K=3$ for simplicity)



Rank 1

Mode-1 unfolding

$$\begin{aligned}\mathbf{A}_{(1)} &= [uv_1w_1 \quad uv_2w_1 \quad \cdots \quad uv_{n_2}w_{n_3}] \\ &= \mathbf{u}(\mathbf{v} \otimes \mathbf{w})^\top \quad (\otimes: \text{Kronecker product})\end{aligned}$$

Similarly $\mathbf{A}_{(2)}$ and $\mathbf{A}_{(3)}$ are also rank 1.

Thus multilinear rank=(1,1,1). The other way is also true.

General relation

$$\text{rank}_k(\mathcal{A}) \leq \text{rank}(\mathcal{A}) \leq \prod_{k=1}^K \text{rank}_k(\mathcal{A})$$

- Left hand side:

$$\mathbf{A}_{(1)} = \sum_{r=1}^R \mathbf{u}_r (\mathbf{v}_r \otimes \mathbf{w}_r)^\top$$



- Right hand side: given by the so-called **Tucker decomposition**

Preliminaries

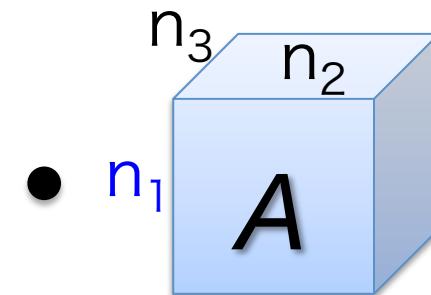
- Define the mode- k product

$$\mathcal{A} \times_k \mathbf{U} := \left(\sum_{j=1}^{n_k} A_{i_1, \dots, i_{k-1}, j, i_{k+1}, \dots, i_K} U_{i_k j} \right)_{i_1, \dots, i_K}$$

for $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_K}$ and $\mathbf{U} \in \mathbb{R}^{n'_1 \times n_k}$

For example,

$$\mathcal{A} \times_1 \mathbf{U} = {}^{n'_1} \begin{matrix} n_1 \\ \mathbf{U} \end{matrix}$$



In other words,

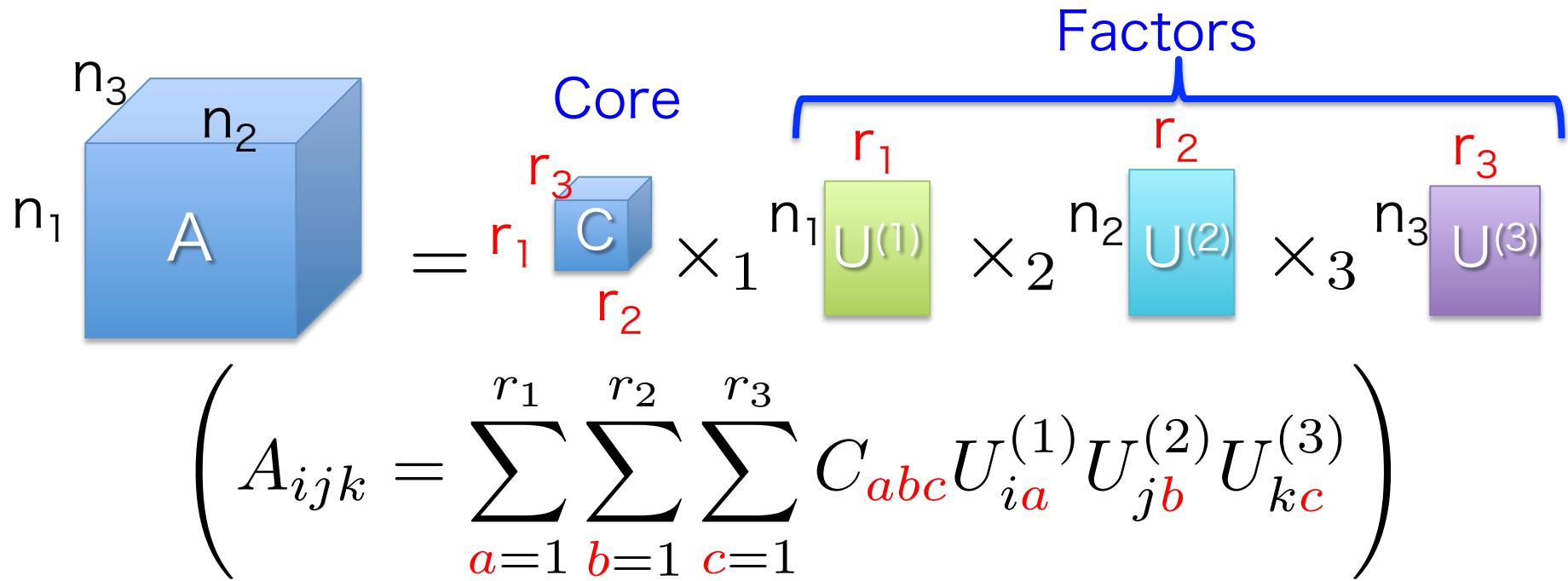
$$(\mathcal{A} \times_k \mathbf{U})_{(k)} = \mathbf{U} \mathcal{A}_{(k)}$$

U applies to each mode-k fiber independently

Tucker decomposition

[Tucker (1966) "Some mathematical notes on three-mode factor analysis"]

- Tensor A has multilinear rank (r_1, \dots, r_K) implies



- Also known as higher-order SVD [De Lathauwer+00]

Intuition

- Recall that $(\mathcal{A} \times_k \mathbf{U})_{(k)} = \mathbf{U} \mathbf{A}_{(k)}$
- Take $\mathbf{U}^{(k)}$, the (full) left singular vectors of $\mathbf{A}_{(k)}$

$$(\mathcal{A} \times_k \mathbf{U}^{(k)\top})_{(k)} = r_k \begin{matrix} n_k \\ \mathbf{U}^{(k)\top} \end{matrix} \bullet \begin{matrix} \prod_{k' \neq k} n_{k'} \\ \mathbf{A} \end{matrix}$$

- Repeat

$$\mathcal{A} \times_1 \mathbf{U}^{(1)\top} \cdots \times_K \mathbf{U}^{(K)\top} = \begin{matrix} r_3 \\ r_1 \\ C \\ r_2 \end{matrix}$$

- Reconstruct

$$\mathcal{A} = \mathcal{C} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_K \mathbf{U}^{(K)}$$

Properties

- Tensor A has multilinear rank (r_1, \dots, r_K) implies

$$\mathcal{A} = \sum_{a=1}^{r_1} \sum_{b=1}^{r_2} \sum_{c=1}^{r_3} C_{abc} u_a^{(1)} \circ u_b^{(2)} \circ u_c^{(3)}$$

Thus

$$\text{rank}(\mathcal{A}) \leq \prod_{k=1}^K \text{rank}_k(\mathcal{A})$$

- In fact,

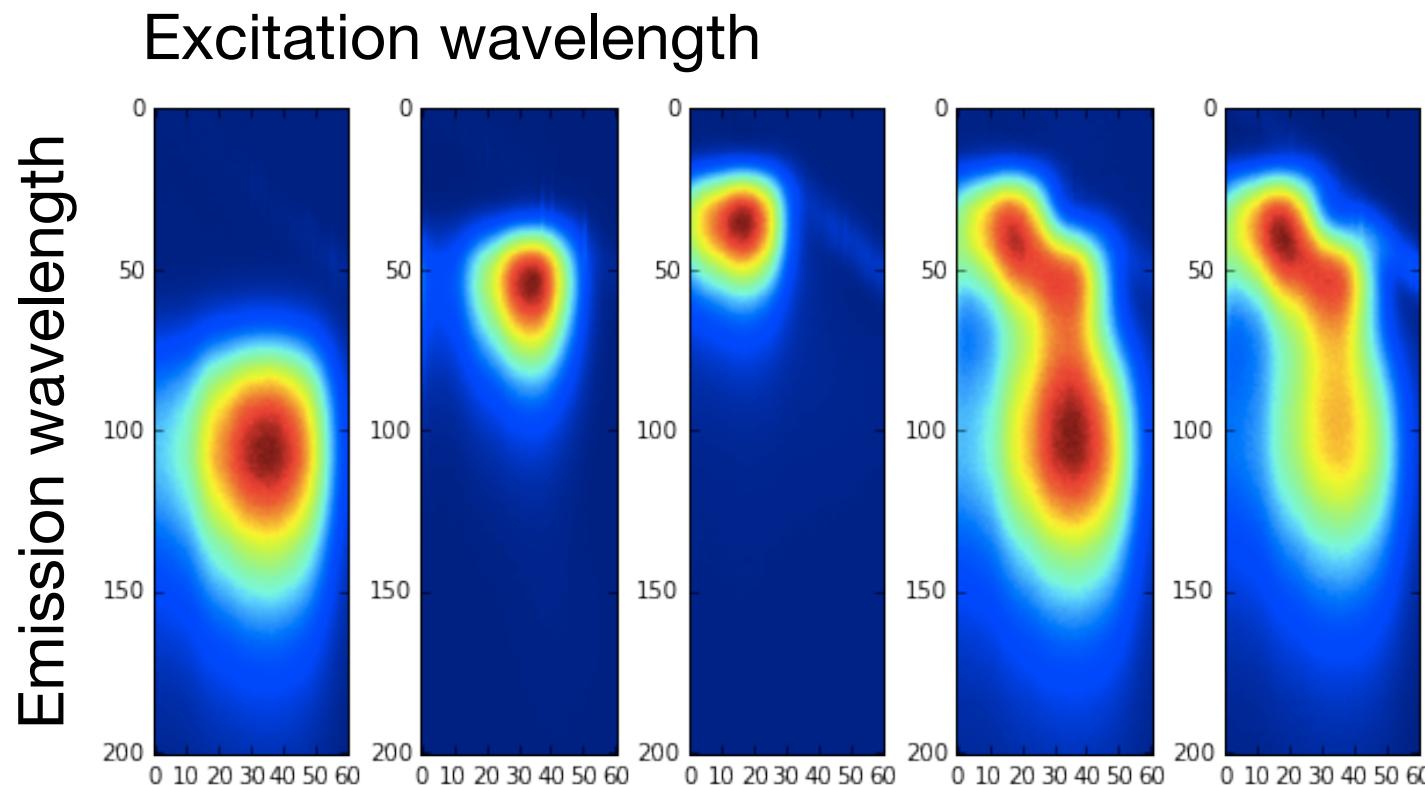
$$\text{rank}(\mathcal{A}) = \text{rank}(\mathcal{C})$$

if the columns of $U^{(1)}, \dots, U^{(K)}$ are independent.

Amino acids fluorescence data

[Bro (1997) “PARAFAC: Tutorial and applications”]

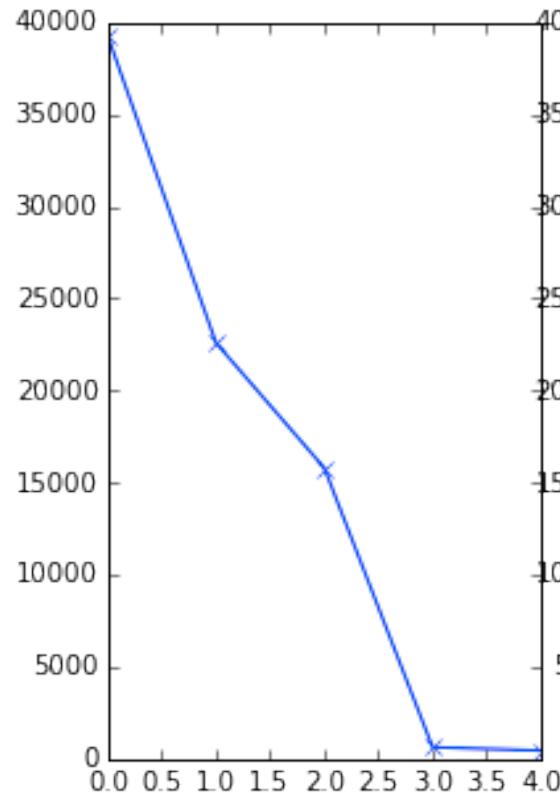
- Five samples containing different amounts of three types of amino acids (Tyr, Trp, and Phe) measured by fluorescence spectrometer



Mode-wise SVD

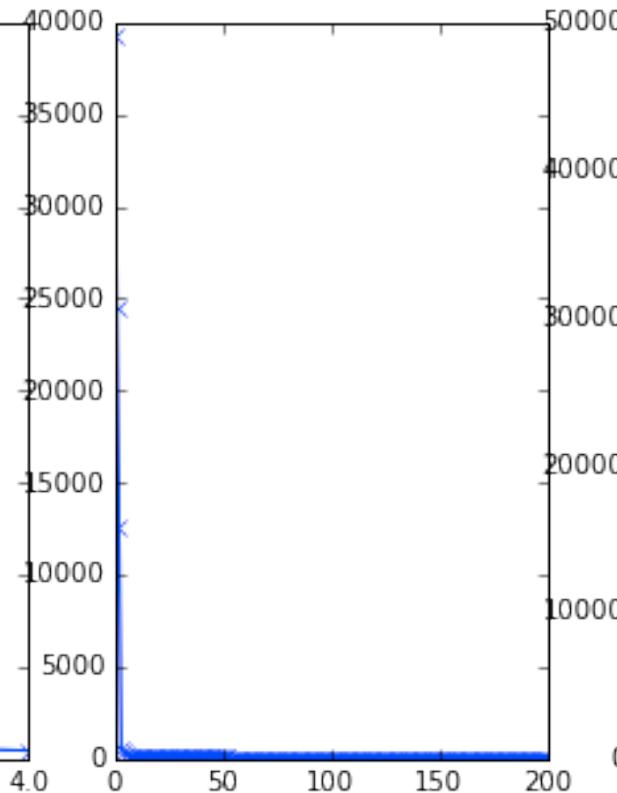
Sample dim.

SVD($A_{(1)}$)



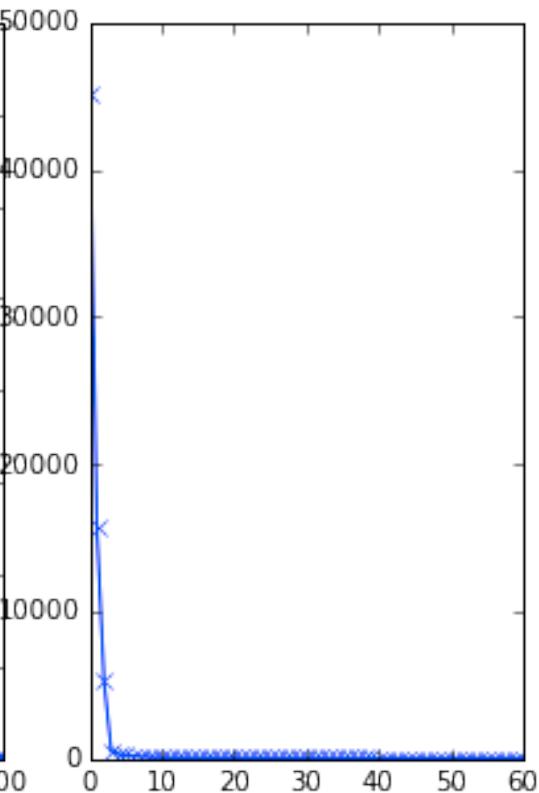
Emission dim.

SVD($A_{(2)}$)



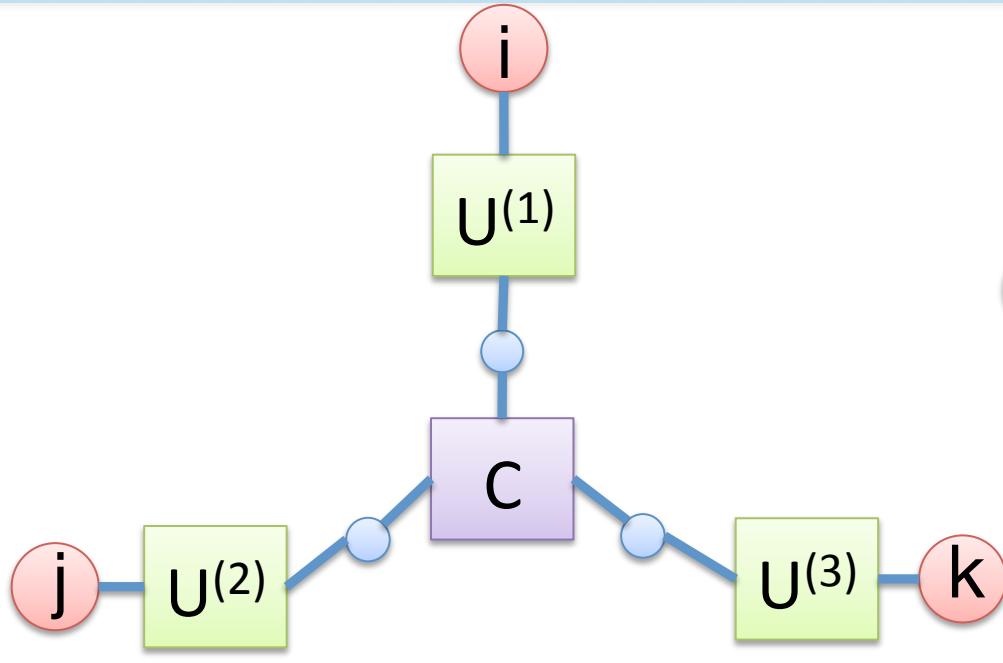
Excitation dim.

SVD($A_{(3)}$)

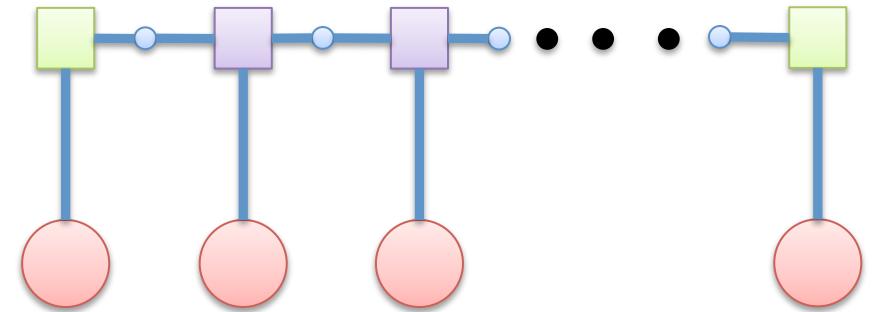


Multilinear rank $\doteq (3, 3, 3) \Rightarrow 3 \leq \text{rank}(A) \leq 9$

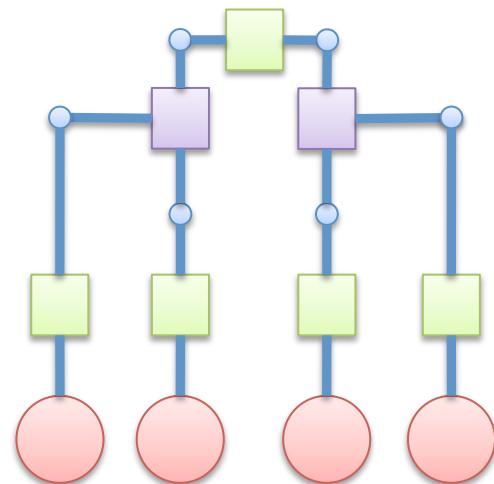
Graphical representation



Tucker decomposition



Tensor train decomposition
[Oseledets, 2011]



Hierarchical Tucker decomposition
[Grasedyck, 2010]

Short summary

- Tensor rank is a much more delicate concept than matrix rank
 - difficulty of computing (**NP-hard** [Håstad, 90])
 - non-closedness
 - may increase by subtracting the best rank-1 approximation [Stegeman & Comon 2010]
- Rank can be bounded by the multilinear rank, which can be obtained easily by mode-wise SVD.

Tensor decomposition

Best rank-k approximation

Singular values and eigenvalues

Power method

Random projection

Best rank- R approximation

- Minimization problem

$$\underset{\mathbf{U}, \mathbf{V}, \mathbf{W}}{\text{minimize}} \left\| \mathcal{A} - \sum_{r=1}^R \mathbf{u}_r \circ \mathbf{v}_r \circ \mathbf{w}_r \right\|_F$$

(Frobenius norm defined as $\|\mathcal{A}\|_F = \sqrt{\langle \mathcal{A}, \mathcal{A} \rangle}$)

- Hardness
 - Known to be NP-hard even for R=1 (3-SAT)
- Questions: for R=1
 - Does it lead to equations defining SVD?
 - Does it relate to the operator norm?

Connection to spectral norm

- Finding

$$(\hat{\mathbf{u}}, \hat{\mathbf{v}}, \hat{\mathbf{w}}) = \operatorname{argmin}_{\mathbf{u}, \mathbf{v}, \mathbf{w}} \|\mathcal{A} - \mathbf{u} \circ \mathbf{v} \circ \mathbf{w}\|_F$$

is equivalent to computing

$$\|\mathcal{A}\| = \max_{\begin{array}{l} \|\tilde{\mathbf{u}}\| \leq 1, \\ \|\tilde{\mathbf{v}}\| \leq 1, \\ \|\tilde{\mathbf{w}}\| \leq 1 \end{array}} \langle \mathcal{A}, \tilde{\mathbf{u}} \circ \tilde{\mathbf{v}} \circ \tilde{\mathbf{w}} \rangle$$

In fact, $\hat{\mathbf{u}} = \sigma \tilde{\mathbf{u}}$, $\hat{\mathbf{v}} = \tilde{\mathbf{v}}$, $\hat{\mathbf{w}} = \tilde{\mathbf{w}}$ with $\sigma = \|\mathcal{A}\|$ is a best rank-one approximation.

Connection to singular values

- Stationary condition implies

$$\mathcal{A} \times_2 \tilde{\mathbf{v}}^\top \times_3 \tilde{\mathbf{w}}^\top = \sigma \tilde{\mathbf{u}}$$

$$\mathcal{A} \times_1 \tilde{\mathbf{u}}^\top \times_3 \tilde{\mathbf{w}}^\top = \sigma \tilde{\mathbf{v}}$$

$$\mathcal{A} \times_1 \tilde{\mathbf{u}}^\top \times_2 \tilde{\mathbf{v}}^\top = \sigma \tilde{\mathbf{w}}$$

generalizes equations defining s.v. for matrices.

- We can define any triplet of unit vectors satisfying the above equations **singular vectors** [Lim, 2005]
- Difference:
 - LHS is quadratic, whereas RHS is linear (motivation for defining L_p singular vectors with $p=K$.)
 - Generally, not orthogonal.

Symmetrization

- Define a super-symmetric $n' \times n' \times n'$ tensor

$$(n' = n_1 + n_2 + n_3)$$

$$\bar{\mathcal{A}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \mathcal{A}^{(2,3,1)} \\ 0 & \mathcal{A}^{(3,2,1)} & 0 \end{bmatrix} \left| \begin{array}{ccc} 0 & 0 & \mathcal{A}^{(1,3,2)} \\ 0 & 0 & 0 \\ \mathcal{A}^{(3,1,2)} & 0 & 0 \end{array} \right| \begin{bmatrix} 0 & \mathcal{A}^{(1,2,3)} & 0 \\ \mathcal{A}^{(2,1,3)} & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

- Eigenvectors of $\bar{\mathcal{A}}$, solutions of

$$\bar{\mathcal{A}} \times_2 \bar{\mathbf{u}}^\top \times_3 \bar{\mathbf{u}}^\top = \lambda \bar{\mathbf{u}}$$

include vectors of the form $\bar{\mathbf{u}} = \begin{bmatrix} \tilde{\mathbf{u}} \\ \tilde{\mathbf{v}} \\ \tilde{\mathbf{w}} \end{bmatrix}$

E.V.P. for symmetric tensors

- Power iteration (detailed below) will quickly converge to one of the solutions –but which solution?
- The solutions are **not orthogonal** to each other.
Cannot use the usual deflation approach (it may actually increase the rank...)
- The **orthogonally decomposable case**

$$\mathcal{A} = \sum_{r=1}^R \lambda_r \mathbf{u}_r \circ \mathbf{u}_r \circ \mathbf{u}_r, \quad (\mathbf{u}_r \perp \mathbf{u}_{r'}, \forall r \neq r')$$

has attracted considerable attention [Anandkumar, Ge, Hsu, Kakade, Telgarsky, 2012]

Power method for symmetric tensors

Fixed-point algorithm

1. Initialize u^0 randomly
2. For $t=0,1,2,\dots$

$$\boldsymbol{u}^{t+1} = \frac{\mathcal{A} \times_2 (\boldsymbol{u}^t)^\top \times_3 (\boldsymbol{u}^t)^\top}{\|\mathcal{A} \times_2 (\boldsymbol{u}^t)^\top \times_3 (\boldsymbol{u}^t)^\top\|_2}$$

Why does it work?

- For symmetric orthogonally decomposable tensors,

$$\mathcal{A} = \sum_{r=1}^R \lambda_r \mathbf{u}_r \circ \mathbf{u}_r \circ \mathbf{u}_r, \quad (\mathbf{u}_r \perp \mathbf{u}_{r'}, \forall r \neq r')$$

$$\mathbf{u}^{t+1} \propto \mathcal{A} \times_2 (\mathbf{u}^t)^\top \times_3 (\mathbf{u}^t)^\top$$

$$= \sum_{r=1}^R \lambda_r (\langle \mathbf{u}_r, \mathbf{u}^t \rangle)^2 \mathbf{u}_r$$

$$\langle \mathbf{u}_r, \mathbf{u}^{t+1} \rangle \propto \lambda_r (\langle \mathbf{u}_r, \mathbf{u}^t \rangle)^2 \quad (\text{Used orthogonality})$$

Notes on power iteration for orthogonally decomposable tensors

$$\langle \mathbf{u}_r, \mathbf{u}^{t+1} \rangle \propto \lambda_r (\langle \mathbf{u}_r, \mathbf{u}^t \rangle)^2$$

- The inner product $\langle \mathbf{u}_r, \mathbf{u}^t \rangle$ will converge to zero **exponentially fast** except for the one with the largest absolute value.
- The rate is $K-1$: power iteration converges **faster for higher order tensors**.
- Odd order K : can take $\lambda_r \geq 0$. RHS ≥ 0 . Even order K : RHS can be negative (oscillation).
- **Deflation works** (because the factors are orthogonal).
- Robustness to noise (see Anandkumar et al. 12)

Power method for asymmetric tensors

1. Initialize u^0, v^0, w^0 randomly
2. For $t=0,1,2,\dots$

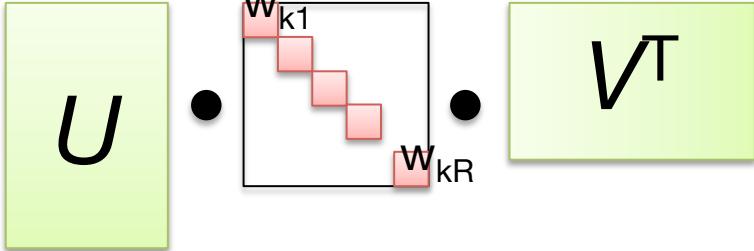
$$\mathbf{u}^{t+1} = \frac{\mathcal{A} \times_2 (\mathbf{v}^t)^\top \times_3 (\mathbf{w}^t)^\top}{\|\mathcal{A} \times_2 (\mathbf{v}^t)^\top \times_3 (\mathbf{w}^t)^\top\|_2},$$

$$\mathbf{v}^{t+1} = \frac{\mathcal{A} \times_1 (\mathbf{u}^t)^\top \times_3 (\mathbf{w}^t)^\top}{\|\mathcal{A} \times_1 (\mathbf{u}^t)^\top \times_3 (\mathbf{w}^t)^\top\|_2},$$

$$\mathbf{w}^{t+1} = \frac{\mathcal{A} \times_1 (\mathbf{u}^t)^\top \times_2 (\mathbf{v}^t)^\top}{\|\mathcal{A} \times_1 (\mathbf{u}^t)^\top \times_2 (\mathbf{v}^t)^\top\|_2}$$

Tensor and simultaneous diagonalization

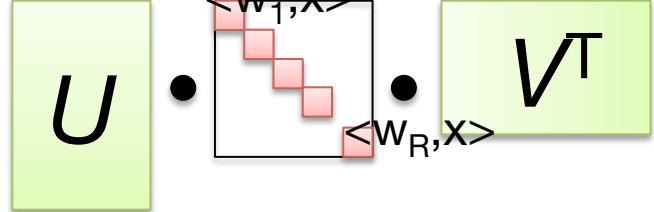
- If $R < \min(n_1, n_2)$, slices are low-rank

$$A_{::,k} = \sum_{r=1}^R w_{kr} u_r v_r^\top = \begin{matrix} U \\ \bullet \\ \begin{matrix} w_{k1} \\ \vdots \\ w_{kR} \end{matrix} \\ \bullet \\ V^\top \end{matrix}$$
The diagram illustrates the decomposition of a tensor slice $A_{::,k}$ into three components. On the left, the equation shows the sum of R terms, each consisting of a weight w_{kr} multiplied by a column vector u_r and its transpose v_r^\top . To the right of the equation is a diagram where the first term U is represented by a green box. The second term is a square matrix with red blocks on the diagonal, labeled with weights $w_{k1}, w_{k2}, \dots, w_{kR}$. The third term V^\top is represented by a green box.

- Factors are common to all slices.
- If we can simultaneously diagonalize all slices, we are done.

Random projection

- Consider the random projection

$$A \times_3 x^\top = \sum_{r=1}^R \langle w_r, x \rangle u_r v_r^\top = U \cdot \begin{matrix} & \begin{matrix} \langle w_1, x \rangle \\ \vdots \\ \langle w_R, x \rangle \end{matrix} \end{matrix} \bullet V^\top$$


x is a random vector. The projection is again a low-rank matrix if $R < \min(n_1, n_2)$.

- Random slice is a special case of random projection
- How many random projections/slices do we need?
 - Two! (if the factors are linearly independent)

[Harshman (1972) Determination and Proof of Minimum Uniqueness Conditions for PARAFAC1]

Random projection algorithm

[Harshman 1972; Goyal, Vempala, Xiao, 2014; Moitra et al.]

1. Let A_1 and A_2 be two independent random projections of rank-R tensor A . (see Random Projection.ipynb)
2. Compute truncated rank-R SVD $A_1 = U_1 S_1 V_1^\top$
3. Compute eigen-decomposition

$$S_1^{-1} U_1^\top A_2 V_1 = P \Lambda P^{-1}$$

4. Recover $U = U_1 S_1 P$

$$V = V_1 P^{-\top}$$

5. Note that

$$A_1 = U V^\top, \quad A_2 = U \Lambda V^\top$$

Simultaneous
diagonalization!

Learning with tensors

Regularization

Nuclear norm

Nuclear norm as an atomic norm

Hardness

Learning with tensors

- Examples
 - Input: pair of vectors (x , z); output: vector y .

$$y = n_1 \begin{matrix} n_3 \\ n_2 \\ W \end{matrix} \times_2 x \times_3 z + \xi$$

The diagram shows the equation $y = n_1 \begin{matrix} n_3 \\ n_2 \\ W \end{matrix} \times_2 x \times_3 z + \xi$. To the left of the equation is a green vertical bar labeled y . To the right is a blue cube labeled W with dimensions $n_1 \times n_2 \times n_3$. Below the cube is a green vertical bar labeled x . To the right of x is a green vertical bar labeled z . To the right of z is a red vertical bar labeled ξ with the word "noise" written above it.

- Input: tensor; output: scalar

$$y = \left\langle n_1 \begin{matrix} n_3 \\ n_2 \\ W \end{matrix}, n_1 \begin{matrix} n_3 \\ n_2 \\ X \end{matrix} \right\rangle + \xi$$

The diagram shows the equation $y = \left\langle n_1 \begin{matrix} n_3 \\ n_2 \\ W \end{matrix}, n_1 \begin{matrix} n_3 \\ n_2 \\ X \end{matrix} \right\rangle + \xi$. To the left of the equation is a green vertical bar labeled y . To the right of the equation is a red vertical bar labeled ξ . Between the two bars are two blue cubes, one labeled W and one labeled X , each with dimensions $n_1 \times n_2 \times n_3$. Brackets on either side of the cubes indicate they are vectors in a space, and the angle brackets indicate the calculation of their dot product.

Regularization

- In learning setting, we are interested in minimizing

$$\underset{\mathcal{W} \in \mathbb{R}^{n_1 \times \cdots \times n_K}}{\text{minimize}} \quad \frac{1}{m} \sum_{i=1}^m \ell(y_i, \langle \mathcal{X}_i, \mathcal{W} \rangle) + R(\mathcal{W})$$

- Examples
 - Low-rank regression: $\mathcal{X} = e_{\text{out}} \circ \mathbf{x} \circ \mathbf{z}$
 - Tensor completion: $\mathcal{X} = e_i \circ e_j \circ e_k$
 - Tensor classification: (general tensor X)

Regularization

- In learning setting, we are interested in minimizing

$$\underset{\mathcal{W} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{m} \sum_{i=1}^m \ell(y_i, \langle \mathcal{X}_i, \mathcal{W} \rangle) + R(\mathcal{W})$$

- Questions:
 - What regularization term $R(\mathcal{W})$ should we use?
 - How is it related to rank?
 - How does regularizing the factors $(\mathcal{U}^{(1)}, \dots, \mathcal{U}^{(K)})$ relate to a regularization for \mathcal{W} ?

Regularizing the factors

- Recall that the nuclear norm

$$\|W\|_* = \min_{\substack{U, V: \\ W=UV^\top}} \frac{1}{2} (\|U\|_F^2 + \|V\|_F^2)$$

minimum over all factorizations $W=UV^\top$.

- A naïve idea

$$R_{2,2}(W) = \min_{\substack{U^{(1)}, U^{(2)}, U^{(3)}: \\ W=\sum_{r=1}^R u_r^{(1)} \circ u_r^{(2)} \circ u_r^{(3)}}} \frac{1}{3} (\|U^{(1)}\|_F^2 + \|U^{(2)}\|_F^2 + \|U^{(3)}\|_F^2)$$

Is this convex? Is this a norm?

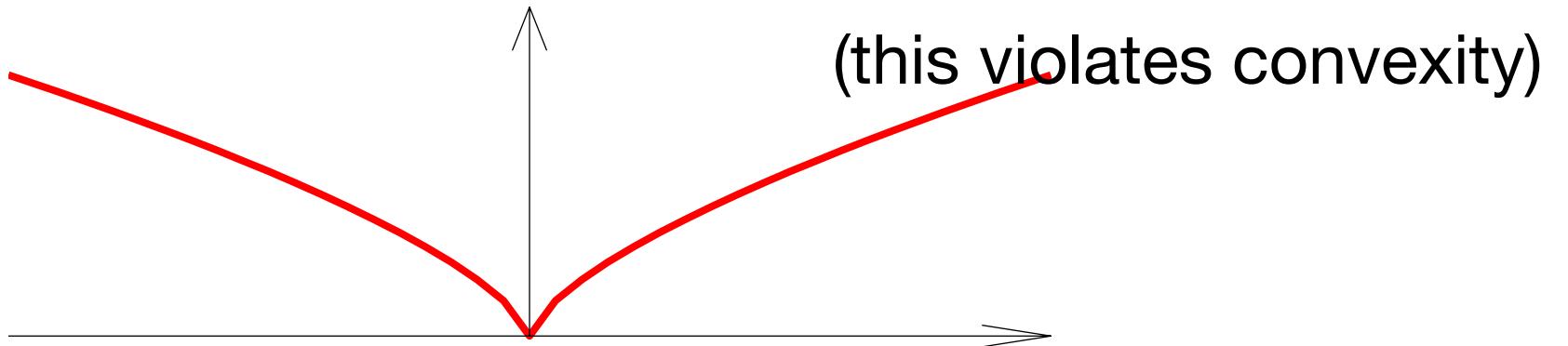
Non-convexity

- Given any decomposition

$$R_{2,2}(\mathcal{W}) \leq \left(\|U^{(1)}\|_F \cdot \|U^{(2)}\|_F \cdot \|U^{(3)}\|_F \right)^{\frac{2}{3}} \quad (\text{AM-GM inequality})$$

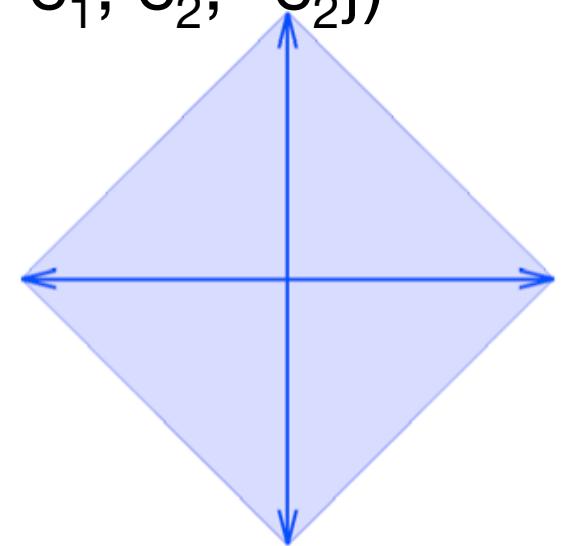
- Thus taking $(\hat{U}^{(1)}, \hat{U}^{(2)}, \hat{U}^{(3)}) = \underset{U^{(1)}, U^{(2)}, U^{(3)}}{\operatorname{argmin}} (\dots)$

$$R_{2,2}(\alpha \mathcal{W}) \leq \left(|\alpha| \|\hat{U}^{(1)}\|_F \cdot \|\hat{U}^{(2)}\|_F \cdot \|\hat{U}^{(3)}\|_F \right)^{\frac{2}{3}} = |\alpha|^{\frac{2}{3}} R_{2,2}(\mathcal{W})$$



Atomic norm

- Atomic norm [Chandrasekaran+2012] is a norm whose unit ball is the convex hull of an **atomic set**.
- Example:
 - L_1 norm unit ball in 2D = $\text{conv}(\{e_1, -e_1, e_2, -e_2\})$
 - Group lasso
 - Nuclear norm
 - etc.



(Tensor) nuclear norm

- Atomic norm with respect to the atomic set

$$\mathcal{A}_{\text{rank}1} = \left\{ \underbrace{\mathbf{u}^{(1)} \circ \cdots \circ \mathbf{u}^{(K)}}_{\text{rank-one tensor}} : \|\mathbf{u}^{(1)}\|_2 = \cdots = \|\mathbf{u}^{(K)}\|_2 = 1 \right\}$$

- Can be written as a minimum

$$\begin{aligned} \|\mathcal{W}\|_{\text{nuc}} &= \min_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(K)}} \sum_{r=1}^R \|\mathbf{u}_r^{(1)}\|_2 \cdot \|\mathbf{u}_r^{(2)}\|_2 \cdots \|\mathbf{u}_r^{(K)}\|_2 \\ \text{s.t. } \mathcal{W} &= \sum_{r=1}^R \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \cdots \circ \mathbf{u}_r^{(K)} \end{aligned}$$

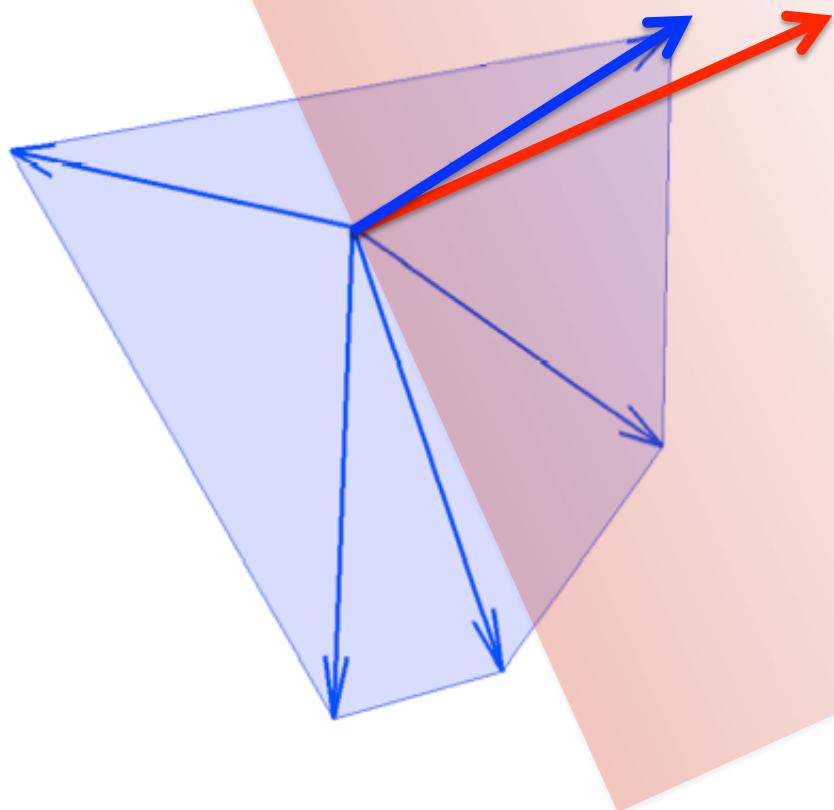
- In other words,

$$\mathcal{W} \in \text{conv}(\mathcal{A}_{\text{rank}1}) \Leftrightarrow \exists \text{decomposition } \sum_{r=1}^R \|\mathbf{u}_r^{(1)}\|_2 \cdots \|\mathbf{u}_r^{(K)}\|_2 \leq 1$$

Dual of nuclear norm is the spectral norm

- Dual of nuclear norm

$$\|\mathcal{A}\|_{\text{nuc}^*} = \max_{\mathcal{W}} \langle \mathcal{A}, \mathcal{W} \rangle, \quad \text{s.t.} \quad \mathcal{W} \in \text{conv}(\mathcal{A}_{\text{rank } 1})$$



$$\|\mathcal{A}\|_{\text{nuc}^*} = \|\mathcal{A}\|$$

Summary

Fact:

$$\max_{\substack{\mathcal{X}: \\ \|\mathcal{X}\| \leq 1}} \langle \mathcal{X}, \mathcal{W} \rangle = \min_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}} \sum_{r=1}^R \|\mathbf{u}_r^{(1)}\|_2 \cdot \|\mathbf{u}_r^{(2)}\|_2 \cdot \|\mathbf{u}_r^{(3)}\|_2$$

s.t. $\mathcal{W} = \sum_{r=1}^R \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \mathbf{u}_r^{(3)}$

Dual of (tensor)
spectral norm = (tensor)
 = nuclear norm = Atomic norm
 wrt $\mathcal{A}_{\text{rank}1}$

How do we compute this norm?

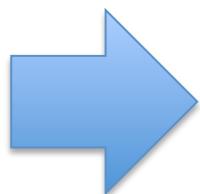
- Unfortunately it is NP hard to compute both the nuclear norm and the spectral norm for tensors.
- Wait, they are norms (convex functions) aren't they easy to deal with?
- Equivalently written as convex problem with $|A_{\text{rank}1}|$ many constraints. Just finding the most violated constraint is NP hard.

Generative model?

- By AM-GM inequality

$$\|\mathcal{W}\|_{\text{nuc}} = \min_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(K)}} \sum_{r=1}^R \frac{1}{K} \left(\|\mathbf{u}_r^{(1)}\|_2^K + \dots + \|\mathbf{u}_r^{(K)}\|_2^K \right)$$

$$\text{s.t. } \mathcal{W} = \sum_{r=1}^R \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \dots \circ \mathbf{u}_r^{(K)}$$

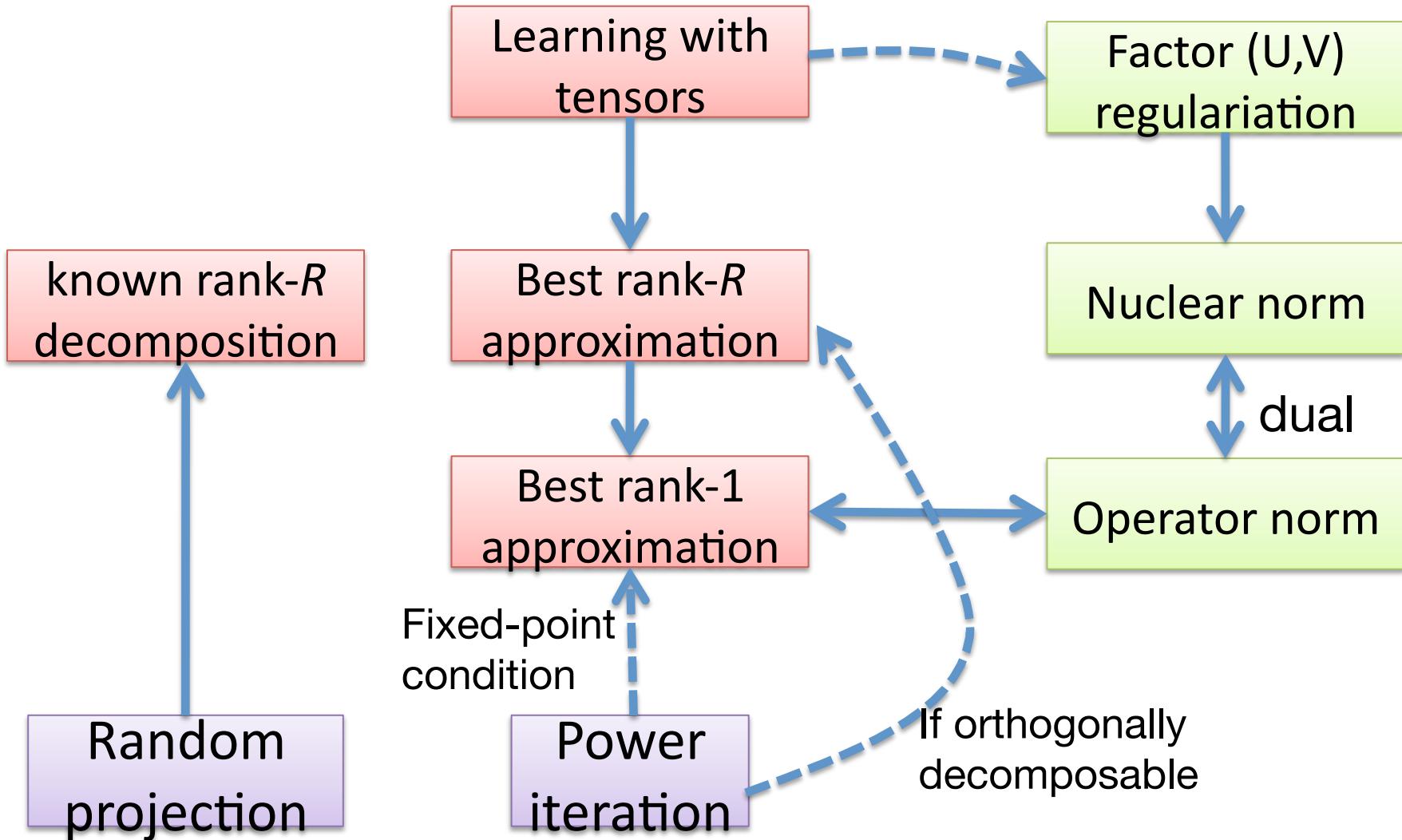

$$P(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(K)}) \propto \prod_{k=1}^K \exp \left(-\frac{1}{K} \sum_{r=1}^R \|\mathbf{u}_r^{(k)}\|_2^K \right)$$

- If we consider instead

$$\mathcal{A}_{\text{rank1}, \mathbf{K}} = \left\{ \mathbf{u}^{(1)} \circ \dots \circ \mathbf{u}^{(K)} : \|\mathbf{u}^{(1)}\|_{\mathbf{K}} = \dots = \|\mathbf{u}^{(K)}\|_{\mathbf{K}} = 1 \right\}$$

$$\|\mathcal{W}\|_{\text{nuc}} = \min_{\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(K)}} \sum_{r=1}^R \frac{1}{K} \left(\|\mathbf{u}_r^{(1)}\|_{\mathbf{K}}^K + \dots + \|\mathbf{u}_r^{(K)}\|_{\mathbf{K}}^K \right)$$

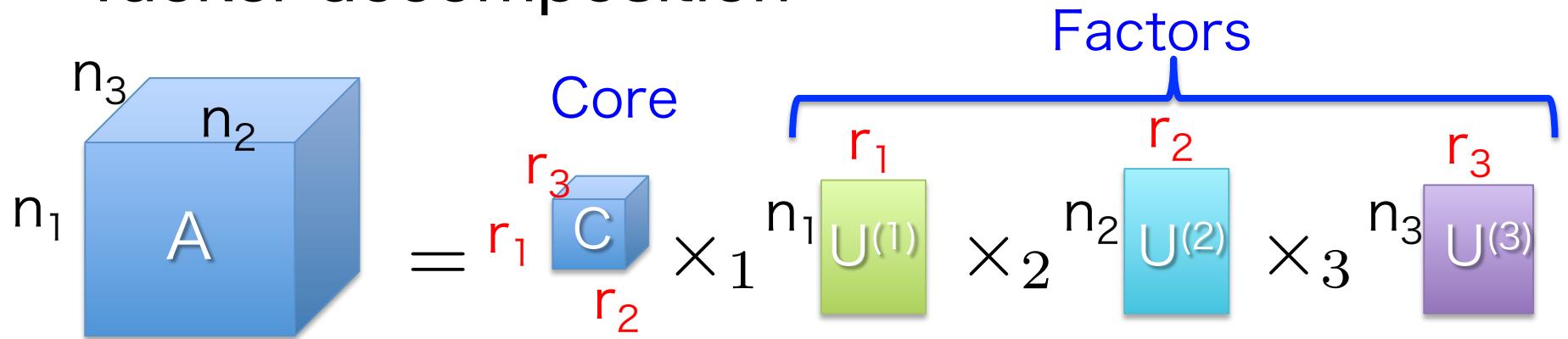
Summary for tensor decomposition



Convex relaxation of Tucker

Overlapped regularization

- If a tensor W is low-rank in the sense of Tucker decomposition



- All unfoldings $W_{(1)}$, $W_{(2)}$, $W_{(3)}$ are simultaneously low-rank.
- How can we encourage W to be simultaneously low-rank?

Approach 1: As a matrix

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

$$\underset{\mathcal{W} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2m} \|\mathbf{y} - \mathfrak{X}(\mathcal{W})\|_2^2 + \lambda_m \|\mathbf{W}_{(k)}\|_*$$

Observation operator $\mathfrak{X} : \mathbb{R}^{n_1 \times \dots \times n_K} \rightarrow \mathbb{R}^m$
 $\mathfrak{X}(\mathcal{W}) = (\langle \mathcal{X}_1, \mathcal{W} \rangle, \dots, \langle \mathcal{X}_m, \mathcal{W} \rangle)^\top$

Pro: Basically a matrix problem

⇒ Easy to carry out, theory available

Con: Have to be lucky to pick the right mode

Approach 2: Overlapped trace norm

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

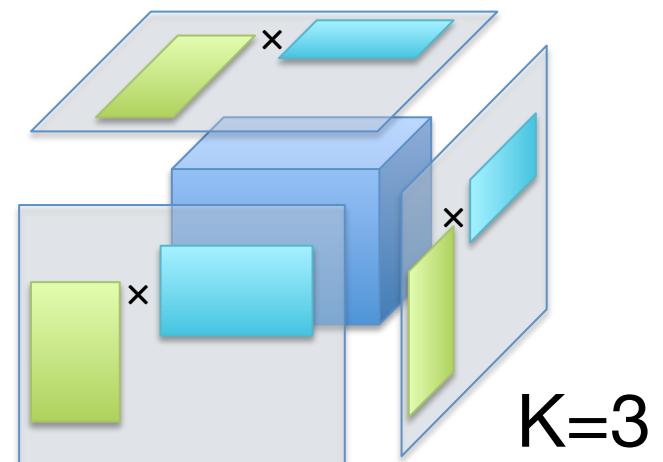
- Regularize by the sum of mode-wise nuclear norms

$$\underset{\mathcal{W} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2m} \|\mathbf{y} - \mathfrak{X}(\mathcal{W})\|_2^2 + \lambda_m \sum_{k=1}^K \|\mathbf{W}_{(k)}\|_*$$

- Intuition:

- the same tensor is regularized to be **simultaneously low-rank** w.r.t. all modes

Nuclear norm of unfoldings



Overlapped nuclear norm

[T+10; Signoretto+10; Gandy+11; Liu+09]

- Regularize the sum of nuclear norms

$$\| \mathcal{W} \|_{S_1/1} := \sum_{k=1}^K \underbrace{\| W_{(k)} \|_*}_{\text{Nuclear norm of unfoldings}}$$

- Is this a norm?

Yes –convex because sum of convex functions and positive homogeneous.

- Is this easy to compute?

Yes – K -unfoldings and SVDs.

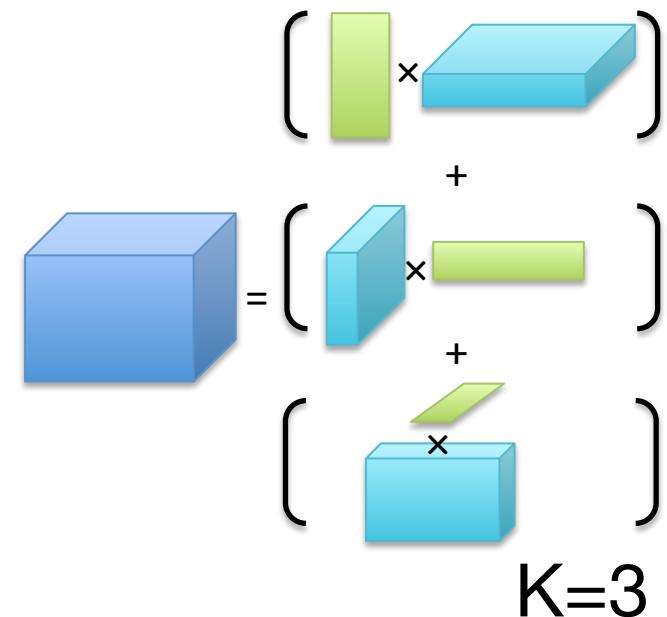
Approach 3: Mixture of low-rank tensors

[T, Hayashi, Kashima10]

- Each mixture component $W^{(k)}$ is regularized to be low-rank in the corresponding mode.

$$\underset{W \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2m} \left\| \mathbf{y} - \mathfrak{X} \left(\sum_{k=1}^K W^{(k)} \right) \right\|_2^2 + \lambda_m \sum_{k=1}^K \|W^{(k)}\|_*$$

- Each $W^{(k)}$ takes care of each mode
- Sparse solution (group lasso like effect)
- Note: the sum may not be low rank



Latent nuclear norm

- Effective regularization term

$$\|\mathcal{W}\|_{\overline{S_1/1}} := \min_{\mathcal{W}^{(1)}, \dots, \mathcal{W}^{(K)}} \sum_{k=1}^K \left\| \mathbf{W}_{(k)}^{(k)} \right\|_* \quad \text{s.t.} \quad \mathcal{W} = \sum_{k=1}^K \mathcal{W}^{(k)}$$

- Is this a norm?

Yes –it is the dual of the overlapped S_∞/∞ -norm:

$$\|\mathcal{X}\|_{\underline{S_\infty/\infty}} := \max_k \left\| \mathbf{X}_{(k)} \right\|$$

spectral norm

- Is it easy to compute?

Yes –convex minimization problem involving matrix nuclear norm

Property of the norms

- Nuclear norm

$$\|\mathbf{W}\|_* \leq \sqrt{r} \|\mathbf{W}\|_F \quad \text{if } \text{rank}(\mathbf{W}) \leq r,$$

- Overlapped trace norm

$$\|\mathcal{W}\|_{\underline{S_1/1}} \leq \sum_{k=1}^K \sqrt{r_k} \|\mathcal{W}\|_F$$

if \mathcal{W} has multilinear rank (r_1, \dots, r_K)

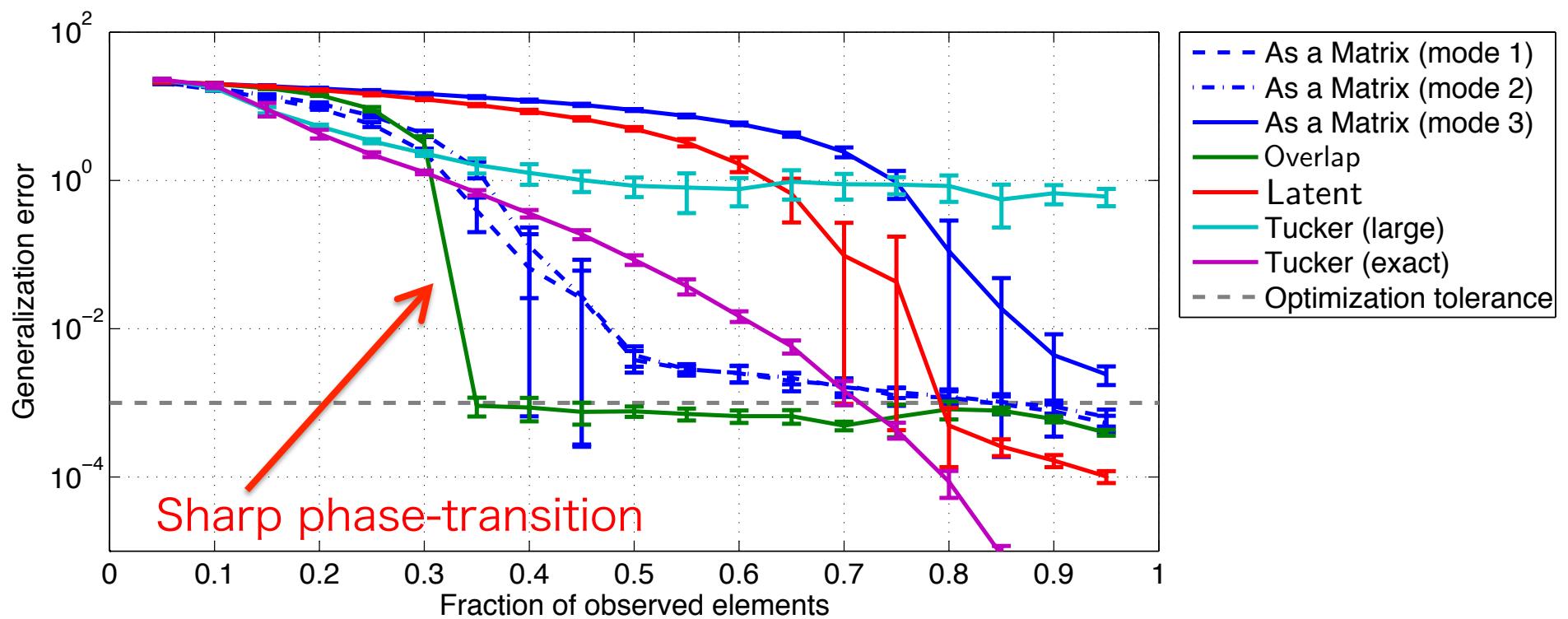
- Latent trace norm

$$\|\mathcal{W}\|_{\overline{S_1/1}} \leq \min_k \sqrt{r_k} \|\mathcal{W}\|_F$$

if \mathcal{W} has multilinear rank (r_1, \dots, r_K)

Tensor completion result

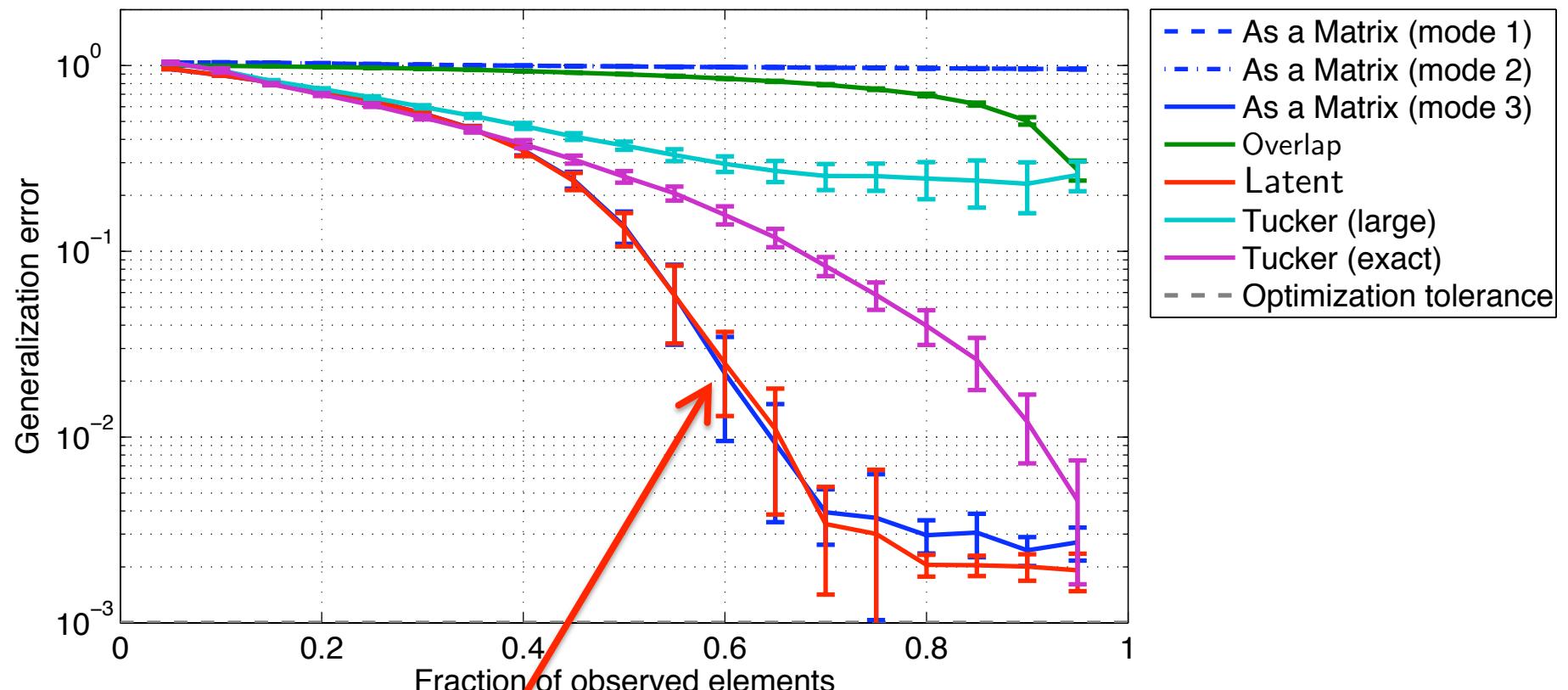
- True tensor: $50 \times 50 \times 20$, rank $7 \times 8 \times 9$. No noise ($\lambda=0$).
- Random train/test split.



Tucker implemented in n-way toolbox. [Andersson & Bro, 00]

Almost full-rank case

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



Latent approach can automatically figure out which mode is low-rank

Overlapped norm: key properties

- How it relates to the multilinear rank

$$\|\mathcal{W}\|_{\underline{S_1/1}} \leq \underbrace{\sum_{k=1}^K \sqrt{r_k}}_{\text{average of the ranks}} \|\mathcal{W}\|_F$$

- How the dual norm behaves

$$\mathbb{E} \|\mathcal{X}\|_{\underline{(S_1/1)^*}} \leq O \left(\sigma \frac{1}{K^2} \sum_{k=1}^K \sqrt{N/n_k} \right), \quad (N = \prod_k n_k)$$

Random Gaussian tensor $\mathcal{N}(0, \sigma^2)$

Overlapped norm: denoising guarantee

[T,Suzuki,Hayashi,Kashima,11]

- Let \mathcal{Y} be a noisy observation of \mathcal{W}^*

$$\mathcal{Y} = \mathcal{W}^* + \mathcal{E} \quad (\mathcal{E} \sim \mathcal{N}(0, \sigma^2))$$

- Then the estimator

$$\hat{\mathcal{W}} = \operatorname{argmin}_{\mathcal{W}} \left(\frac{1}{2} \|\mathcal{Y} - \mathcal{W}\|_F^2 + \lambda \|\mathcal{W}\|_{S_1/1} \right)$$

with $\lambda = c_0 \sigma (\sum_{k=1}^K \sqrt{N/n_k})/K$ gives

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_1 \sigma^2 \left(\frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2 \left(\frac{1}{K} \sum_{k=1}^K 1/\sqrt{n_k} \right)^2$$

with prob. at least $1 - \exp(-\text{poly}(n))$

Latent norm: key properties

- How it relates to the multilinear rank

$$\|\mathcal{W}\|_{\overline{S_1/1}} \leq (\underbrace{\min_k \sqrt{r_k}}_{\text{minimum rank}}) \|\mathcal{W}\|_F$$

- How the dual norm behaves

$$\mathbb{E} \|\mathcal{X}\|_{(\overline{S_1/1})^*} \leq \tilde{O} \left(\sigma \max_k \sqrt{N/n_k} \right), \quad (N = \prod_k n_k)$$

 Random Gaussian tensor $\mathcal{N}(0, \sigma^2)$

Denoising guarantees

- Overlapped nuclear norm: sensitive to the average rank

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_1 \sigma^2 \left(\frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2 \left(\frac{1}{K} \sum_{k=1}^K 1/\sqrt{n_k} \right)^2$$

- Latent nuclear norm: sensitive to the minimum rank

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_2 \sigma^2 \min_k r_k / \min_k n_k$$

- Matters only if $r_k \neq r_k'$ or $n_k \neq n_k'$ (typical many tensor problems)

Could we obtain a bound of $\min_k (r_k/n_k)$?

Scaled latent nuclear norm

[Wimalawarne, Sugiyama, T, 2014]

- Normalize by $\sqrt{n_k}$

$$\|\mathcal{W}\|_{\text{scaled}} = \inf_{\mathcal{W}^{(1)}, \dots, \mathcal{W}^{(K)}} \sum_{k=1}^K \frac{1}{\sqrt{n_k}} \|\mathbf{W}_{(k)}^{(k)}\|_{S_1} \quad \text{s.t.} \quad \sum_{k=1}^K \mathcal{W}^{(k)} = \mathcal{W}$$

- relation to the multilinear rank

$$\|\mathcal{W}\|_{\text{scaled}} \leq \left(\min_k \sqrt{\frac{r_k}{n_k}} \right) \|\mathcal{W}\|_F$$

- behavior of the dual norm

$$\mathbb{E} \|\mathcal{X}\|_{\text{scaled}^*} \leq \tilde{O} \left(\sigma \sqrt{N} \right), \quad (N = \prod_k n_k)$$

Denoising guarantees

- Overlapped nuclear norm: sensitive to the average rank

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_1 \sigma^2 \left(\frac{1}{K} \sum_{k=1}^K \sqrt{r_k} \right)^2 \left(\frac{1}{K} \sum_{k=1}^K 1/\sqrt{n_k} \right)^2$$

- Latent nuclear norm: sensitive to the minimum rank

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_2 \sigma^2 \min_k r_k / \min_k n_k$$

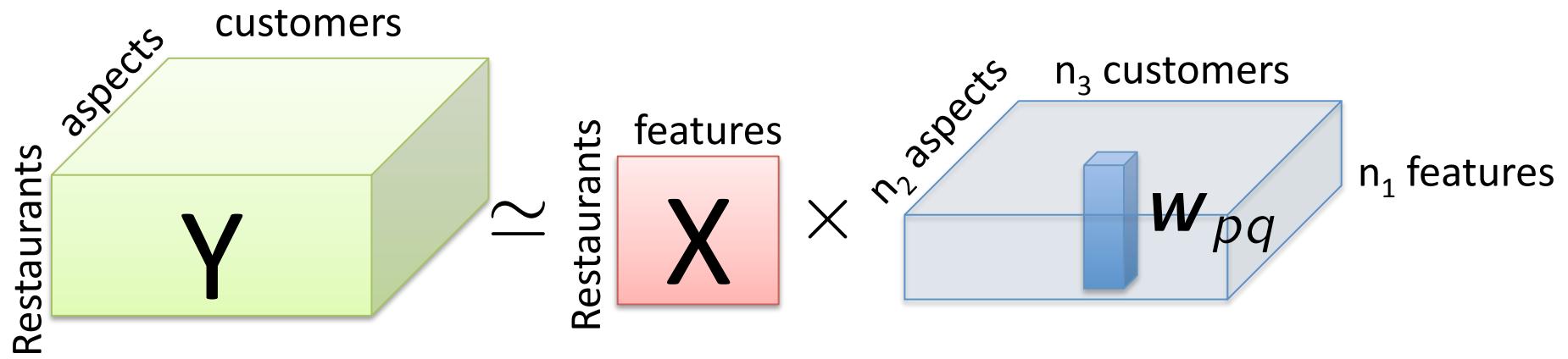
- Scaled latent norm: sensitive to the minimum rank-to-dimension ratio

$$\frac{1}{N} \|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_3 \sigma^2 \min_k \frac{r_k}{n_k}$$

Multi-task learning with tensors

[Romera-Paredes+13; Signoretto+13]

- Example: multi-aspect restaurant recommendation



- Issue: heterogeneous dimensions and ranks
 - Number of aspects is usually much smaller than the number of customers.
 - Rank (amount of sharing of information) may be different for aspect, customers, and features.

Theory for multi-linear MTL

- $n_2 \times n_3$ tasks. Observed tasks $S \subseteq [n_2] \times [n_3]$
- Training samples
 - many tasks may have no samples (zero-shot learning) $(\mathbf{x}_{ipq}, y_{ipq})_{i=1}^{m_{pq}}, (p, q) \in S$
- Training

$$\hat{\mathcal{W}} = \underset{\mathcal{W} \in \mathbb{R}^{n_1 \times n_2 \times n_3}}{\operatorname{argmin}} \hat{L}(\mathcal{W}), \text{ s.t. } \|\mathcal{W}\|_* \leq C$$

where $\hat{L}(\mathcal{W}) = \frac{1}{|S|} \sum_{(p,q) \in S} \frac{1}{m_{pq}} \sum_{i=1}^{m_{pq}} \ell(\langle \mathbf{x}_{ipq}, \mathbf{w}_{pq} \rangle, y_{ipq})$

Analysis

[Wimalawarne, Sugiyama, T, 14]

- Lemma: Let \mathcal{W}^* be any tensor with rank (r_1, r_2, r_3) and element-wise bounded by B . For $C = B\phi(\mathbf{r})\sqrt{N}$

$$L(\hat{\mathcal{W}}) - L(\mathcal{W}^*) \leq O \left(\Lambda B \underbrace{\phi(\mathbf{r})}_{\text{How the norm relates to the rank}} \cdot \frac{\sqrt{N}}{|S|} \cdot \underbrace{\mathbb{E} \left\| \sum_{(p,q) \in S} \frac{1}{m_{pq}} \sum_{i=1}^{m_{pq}} \sigma_{ipq} \mathcal{X}_{ipq} \right\|_{\star^*}}_{\text{How the dual norm behaves}} \right)$$

How the norm relates to the rank

$$\|\mathcal{W}\|_{\star} \leq \phi(\mathbf{r}) \|\mathcal{W}\|_F$$

How the dual norm behaves

σ_{ipq} : Rademacher RV

Λ : Lipschitz constant of the loss ℓ

κ : Measure of feature correlation $\mathbb{E} [\mathbf{x}_{ipq} \mathbf{x}_{ipq}^\top] = C_{pq} \preceq \kappa/n_1 \mathbf{I}_{n_1}$

$L(\mathcal{W})$: Expected loss $L(\mathcal{W}) = \frac{1}{n_2 n_3} \sum_{(p,q) \in [n_2] \times [n_3]} \mathbb{E} \ell (\langle \mathbf{x}_{pq}, \mathbf{w}_{pq} \rangle, y_{pq})$

Theorem (overlap norm)

[Wimalawarne, Sugiyama, T, 14]

- Let W^* be any tensor with rank (r_1, r_2, r_3) and element-wise bounded by B . For $C = B \sqrt{\|\mathbf{r}\|_{1/2} n_1 n_2 n_3}$

$$L(\hat{\mathcal{W}}) - L(\mathcal{W}^*) \leq O \left(\Lambda B \sqrt{\frac{\kappa}{m|S|} \|\mathbf{r}\|_{1/2} \min_k (D_k \log D_k)} \right)$$

Λ : Lipschitz constant of the loss ℓ

κ : Measure of feature correlation $\mathbb{E} [\mathbf{x}_{ipq} \mathbf{x}_{ipq}^\top] = C_{pq} \preceq \kappa/n_1 \mathbf{I}_{n_1}$

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

$$D_1 := n_1 + n_2 n_3, \quad D_2 := n_2 + n_1 n_3, \quad D_3 := n_3 + n_1 n_2$$

Theorem (latent norm)

[Wimalawarne, Sugiyama, T, 14]

- Let W^* be any tensor with rank (r_1, r_2, r_3) and element-wise bounded by B . For $C = B \sqrt{\min_k r_k n_1 n_2 n_3}$

$$L(\hat{W}) - L(W^*) \leq O \left(\Lambda B \sqrt{\frac{\kappa}{m|S|} \min_k r_k \max_k (D_k \log D_k)} \right)$$

Λ : Lipschitz constant of the loss ℓ

κ : Measure of feature correlation $\mathbb{E} [x_{ipq} x_{ipq}^\top] = C_{pq} \preceq \kappa/n_1 I_{n_1}$

$D_1 := n_1 + n_2 n_3, \quad D_2 := n_2 + n_1 n_3, \quad D_3 := n_3 + n_1 n_2$

Theorem (scaled latent norm)

[Wimalawarne, Sugiyama, T, 14]

- Let W^* be any tensor with rank (r_1, r_2, r_3) and element-wise bounded by B . For $C = B \sqrt{\min_k (r_k/n_k) n_1 n_2 n_3}$

$$L(\hat{W}) - L(W^*) \leq O \left(\Lambda B \sqrt{\frac{\kappa}{m|S|} \min_k \left(\frac{r_k}{n_k} \right) n_1 n_2 n_3 \log(\max_k D_k)} \right)$$

Λ : Lipschitz constant of the loss ℓ

κ : Measure of feature correlation $\mathbb{E} [x_{ipq} x_{ipq}^\top] = C_{pq} \preceq \kappa/n_1 I_{n_1}$

$D_1 := n_1 + n_2 n_3, \quad D_2 := n_2 + n_1 n_3, \quad D_3 := n_3 + n_1 n_2$

Sample complexities

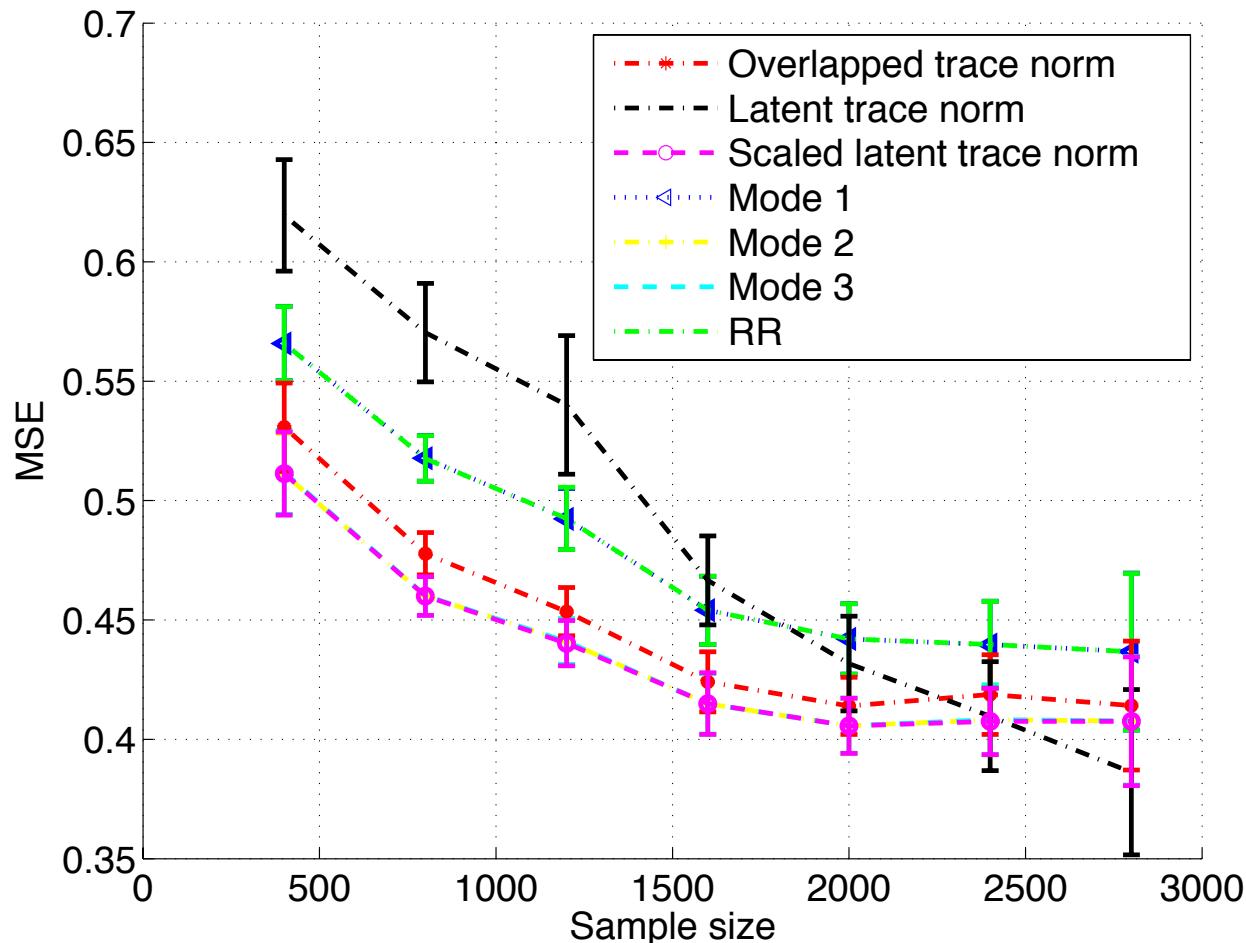
(ignoring log terms and per $1/\varepsilon^2$)

	Overlap	Latent	Scaled latent
Multi-task learning (homogeneous)	$\kappa \ \mathbf{r}\ _{1/2} n^2$	$\kappa (\min_k r_k) n^2$	$\kappa (\min_k r_k) n^2$
Multi-task learning (heterogeneous)	$\kappa \ \mathbf{r}\ _{1/2} n_2 n_3$	$\kappa n_1 n_2 n_3$	$\kappa \min(r_1 n_2 n_3, n_1 n_2 r_3)$
Tensor completion	$\ \mathbf{r}\ _{1/2} \min_k D_k$	$\min_k r_k \max_k D_k$	$\min_k \frac{r_k}{n_k} N$

Restaurant recommendation data

[Vargas-Gómez+2011; Romera-Paredes+2013]

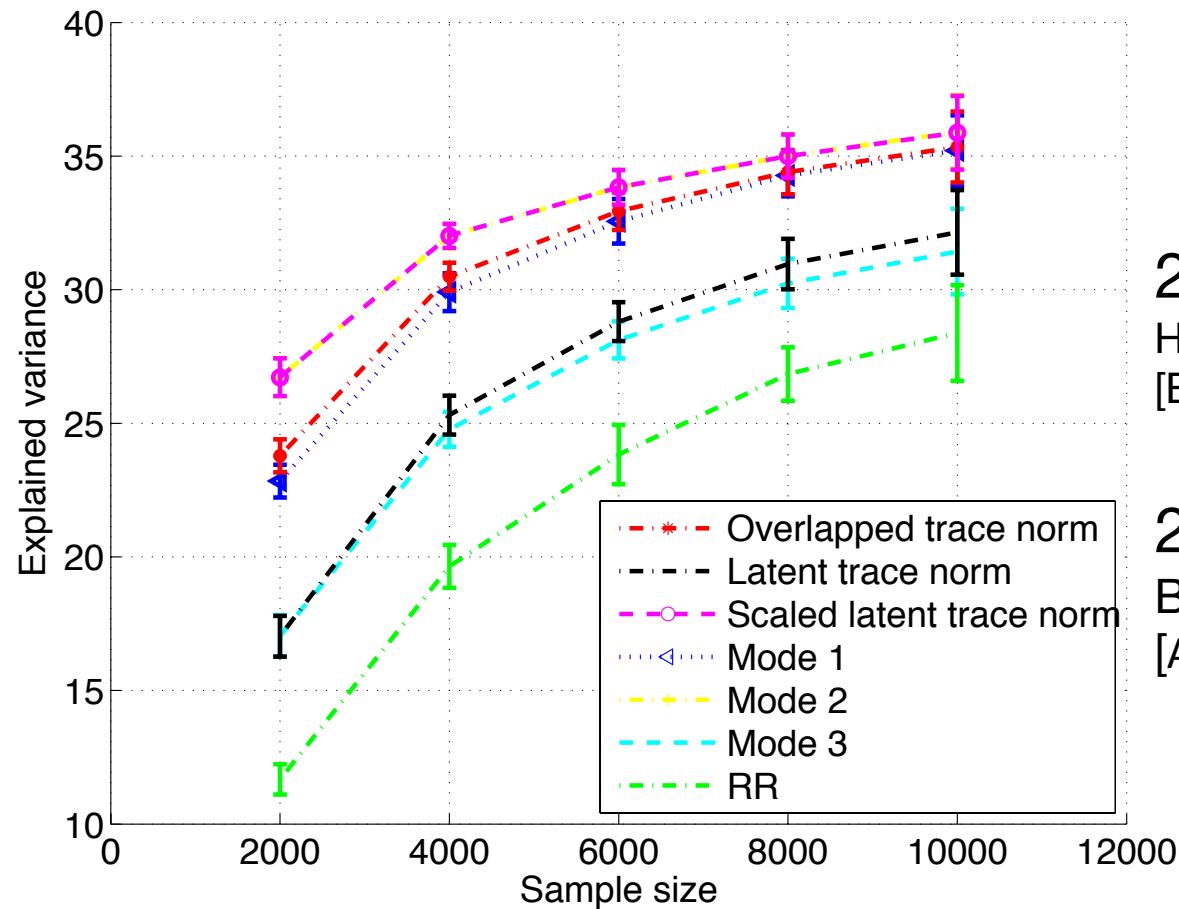
- 45 features x 3 aspects x 138 customers



ILEA School data

[ILEA; Bakker & Heskes 2003]

- 24 features x 139 schools x 3 years



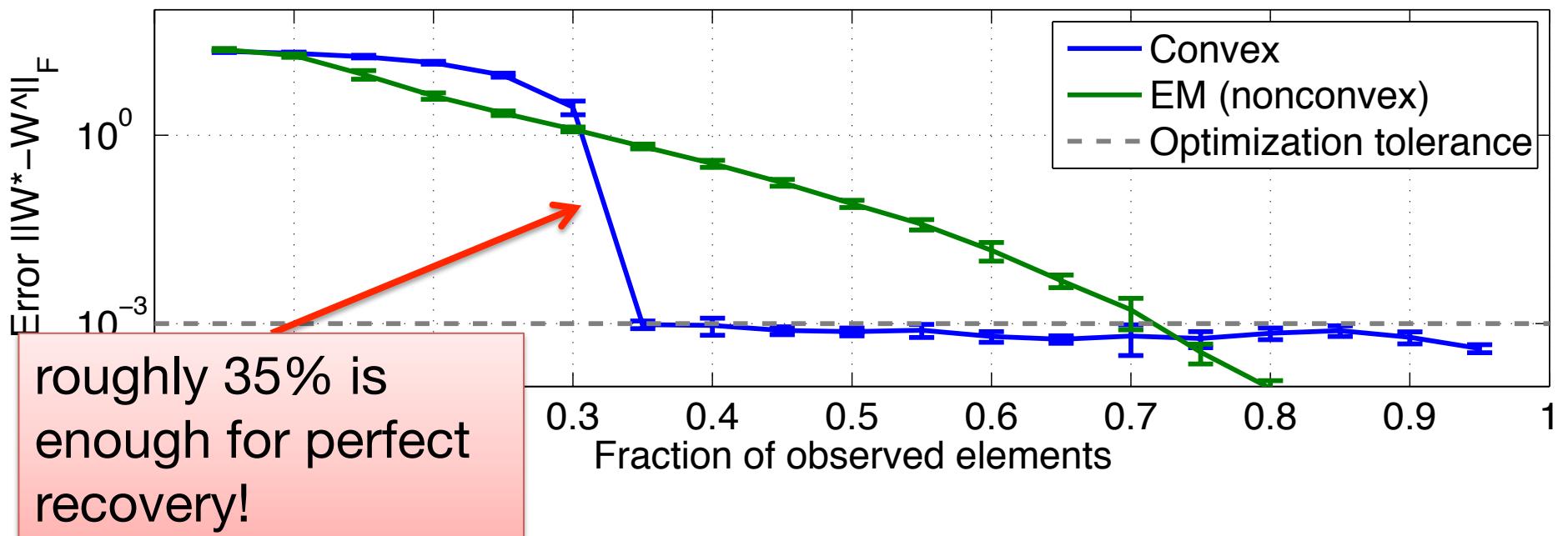
29.5%
Hierarchical Bayes
[Bakker & Heskes 03]

26.7%
Bilinear MTL
[Argyriou 08]

Tensor completion

$$\begin{aligned} \min_{\mathcal{W}} \quad & \|\mathcal{W}\|_{S_1/1} \\ \text{s.t.} \quad & W_{i,j,k} = Y_{i,j,k} \end{aligned}$$

size=50x50x20, rank=7x8x9



EM algorithm implemented in n-way toolbox. [Andersson & Bro, 00]

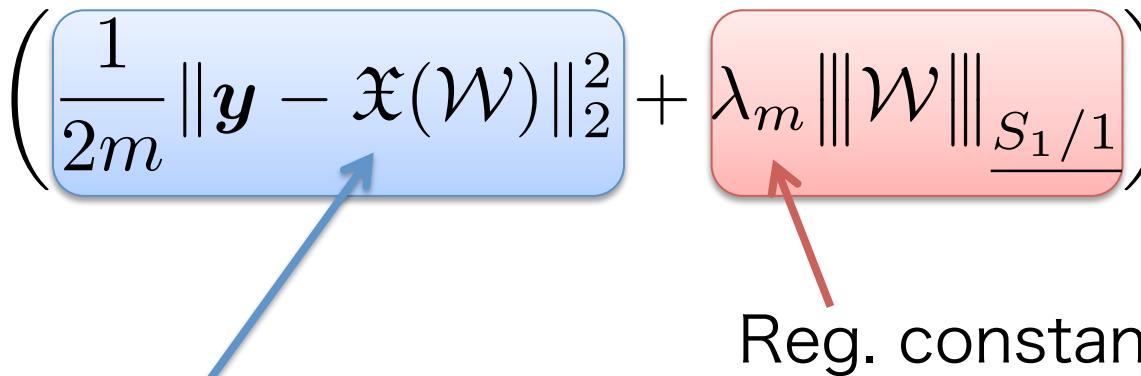
Problem setting

Observation \mathcal{W}^* : true K-way tensor with rank (r_1, \dots, r_K)

$$y_i = \langle \mathcal{X}_i, \mathcal{W}^* \rangle + \epsilon_i \quad (i = 1, \dots, m)$$

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

Optimization	Likelihood	Regularization
$\hat{\mathcal{W}} = \underset{\mathcal{W} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\operatorname{argmin}} \left(\frac{1}{2m} \ \mathbf{y} - \mathfrak{X}(\mathcal{W})\ _2^2 + \lambda_m \ \mathcal{W}\ _{S_1/1} \right)$		



Reg. constant

Observation operator $\mathfrak{X} : \mathbb{R}^{n_1 \times \dots \times n_K} \rightarrow \mathbb{R}^m$

$$\mathfrak{X}(\mathcal{W}) = (\langle \mathcal{X}_1, \mathcal{W} \rangle, \dots, \langle \mathcal{X}_m, \mathcal{W} \rangle)^\top$$

Theorem (random Gaussian design)

There are constants c_1, c_2 such that

$$\|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_2 \sigma^2 \frac{rn^{K-1}}{M}$$

holds with high probability if

$$\frac{\#\text{samples}(m)}{\#\text{variables}(\prod_k n_k)} \geq c_1 \frac{r}{n} \quad \text{and} \quad \lambda_m = 8\sigma \sqrt{n^{K-1}/m}$$

Condition for the sample size m :

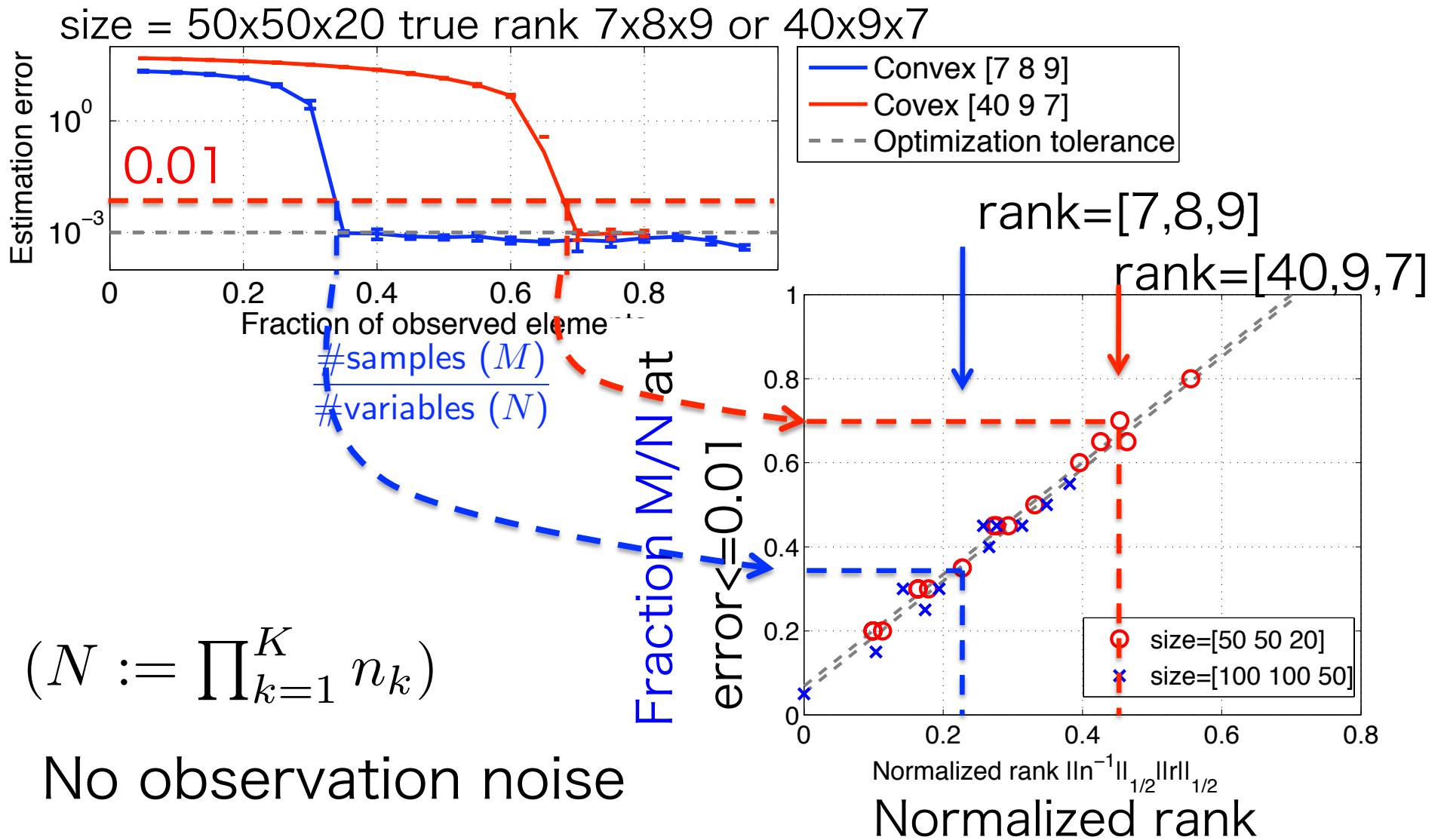
- Depends on the **normalized rank r/n**
- Independent of the noise σ

Condition for the reg. parameter λ_m :

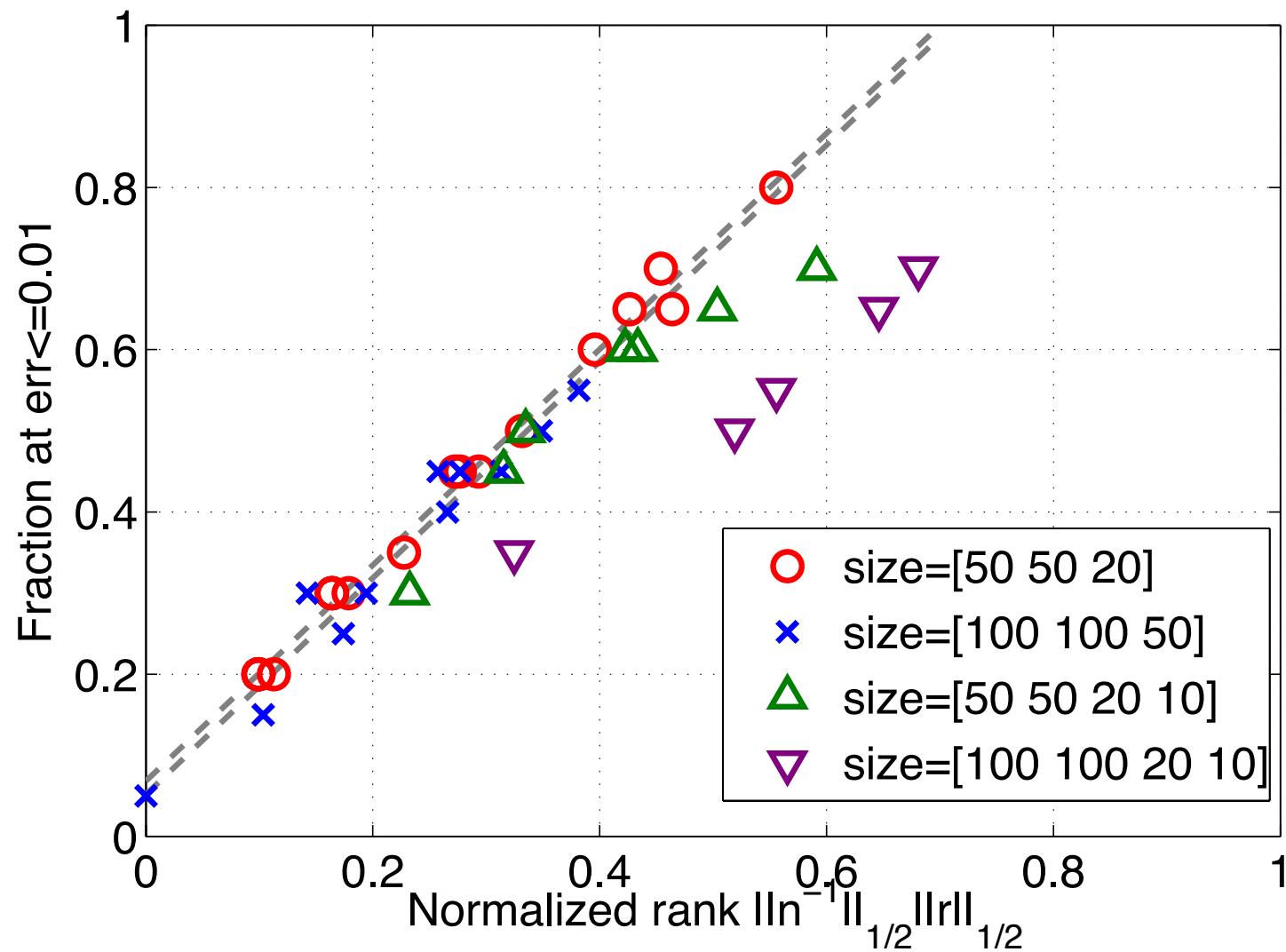
- Independent of the rank r
- Depends on the **noise σ**

For simplicity $n_1, \dots, n_K = n, r_1, \dots, r_K = r$

Tensor completion



Theory vs. Experiments (4th order)



Limitation: too many samples required!

- Sample complexity = number of samples required to obtain error ϵ

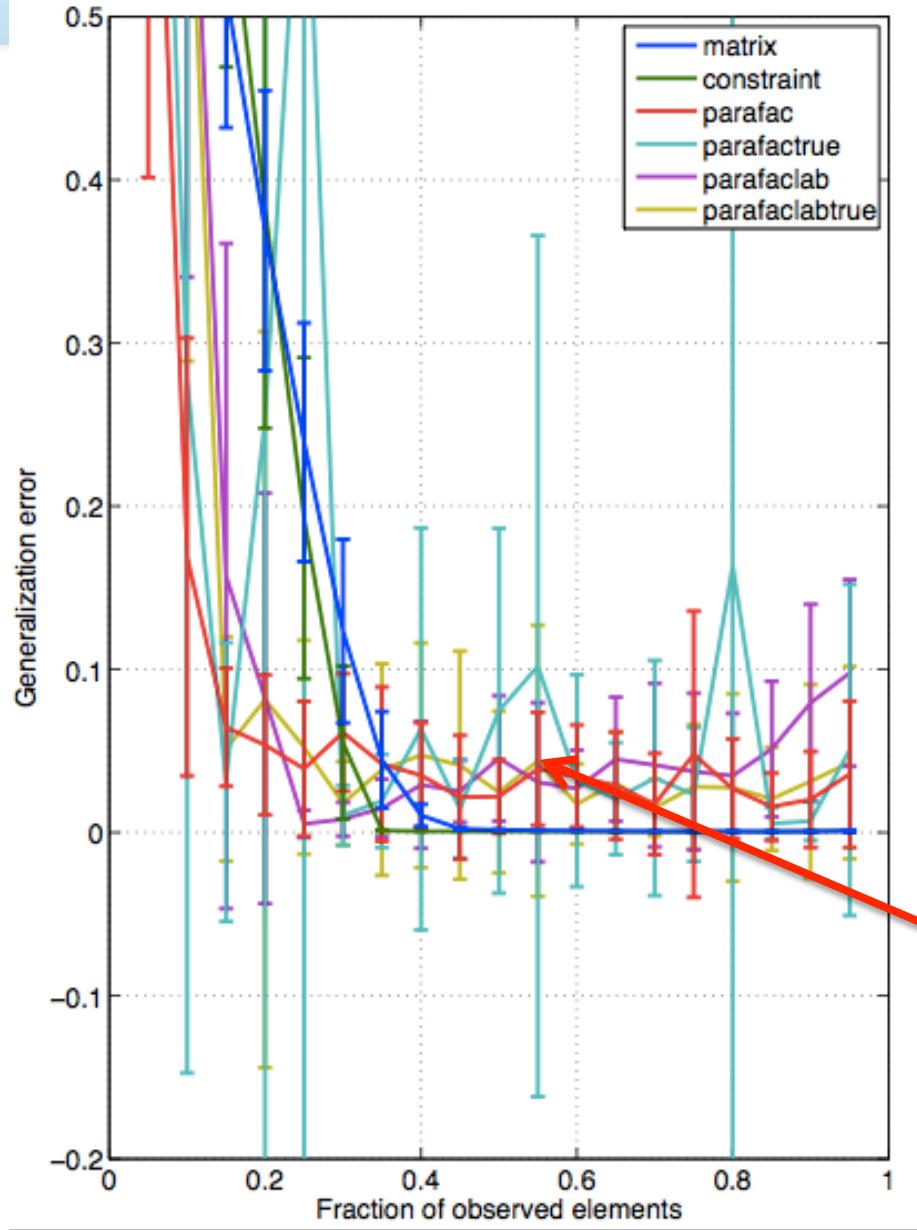
$$\|\hat{\mathcal{W}} - \mathcal{W}^*\|_F^2 \leq c_2 \sigma^2 \frac{rn^{K-1}}{M} \Leftrightarrow M = O(rn^{K-1}/\epsilon)$$

- For $K=2$: $O(rn) =$ number of parameters rank r matrix has.
- For $K>2$: too many. In fact, for the $50 \times 50 \times 20$ rank (7,8,9) case, about 1500 ($\doteq 50^7 + 50^8 + 20^9 + 7^8 \cdot 9$) samples ($3\% \ll 35\%$) should be right.

Tensorlab

Sorber, Van Barel and De Lathauwer (2014).

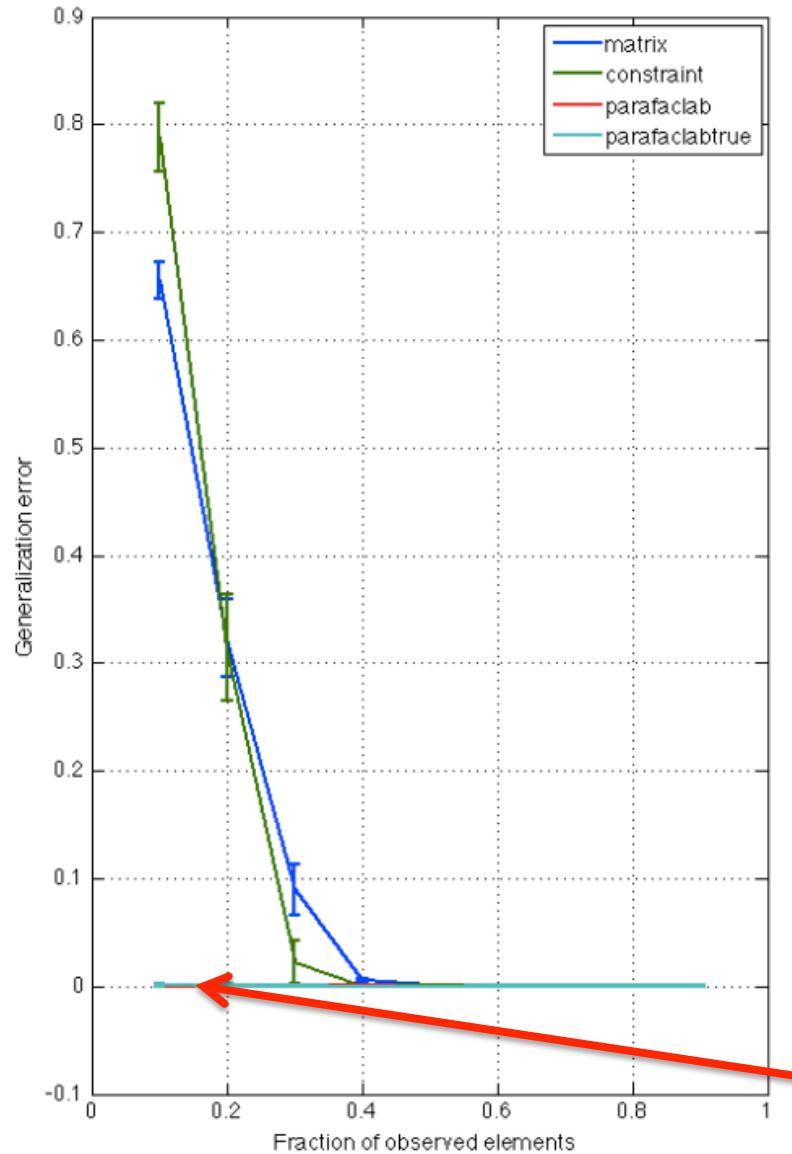
- RT “Hi, I was trying cpd with missing entries and found that the error does not go to zero as the number of observed entries increases...”



Tensorlab

Tensorlab

Sorber, Van Barel and De Lathauwer (2014).



- Laurent Sorber's reply
“Thank you for your detailed example ... I have been experimenting further with your code, and below is a method that gets good performance with Tensorlab ...”

Tensorlab

Lesson learned

- When you are doing something non-convex, details matter a lot!
 - Random initialization turns out to be better than the default intelligent initialization.
 - Make sense to try several short restarts (4×5).
 - Default max iteration too small.
- Although global guarantee may be absent, nonlinear optimization could work very well in practice.

Computation-statistics trade-off

Simplest case [Montanari & Richard 14]

- Noise corrupted symmetric rank-one tensor

$$\begin{matrix} n \\ n \end{matrix} \text{---} Y = \beta \begin{matrix} u^* \\ u^* \\ u^* \end{matrix} + \begin{matrix} n \\ n \end{matrix} \text{---} Z$$

Gaussian $N(0,1)$

β : signal-to-noise ratio

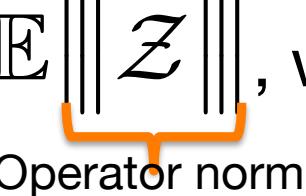
- How large does β need to be to recover u^* from Y ?

Optimal: maximum likelihood

[Richard & Montanari 14]

- Consider the maximum likelihood estimator

$$\hat{\mathbf{u}}_{\text{ML}} = \underset{\mathbf{u} \in \mathbb{R}^n}{\operatorname{argmax}} \langle \mathcal{Y}, \mathbf{u} \circ \mathbf{u} \circ \mathbf{u} \rangle$$

- If $\beta \geq \mathbb{E} \|\mathcal{Z}\|$, with high probability


$$1 - |\langle \hat{\mathbf{u}}_{\text{ML}}, \mathbf{u}^* \rangle| = O\left(\frac{\mathbb{E} \|\mathcal{Z}\|}{\beta}\right)$$

- Note that $\mathbb{E} \|\mathcal{Z}\| = O\left(\sqrt{nK \log(K)}\right)$ [T&Suzuki 14]
- This is optimal (no estimator can do better)
- But computationally interactable (best rank-1 approx.)!

Power method

[Anandkumar+12; Richard & Montanari 14]

- Is a computationally tractable algorithm
- With $\beta \geq c \cdot \sqrt{n^K K \log(K)}$ randomly initialized power iteration converges to a solution with

$$1 - |\langle \hat{u}, u^* \rangle| = O\left(\frac{\mathbb{E} \|\mathcal{Z}\|}{\beta}\right)$$

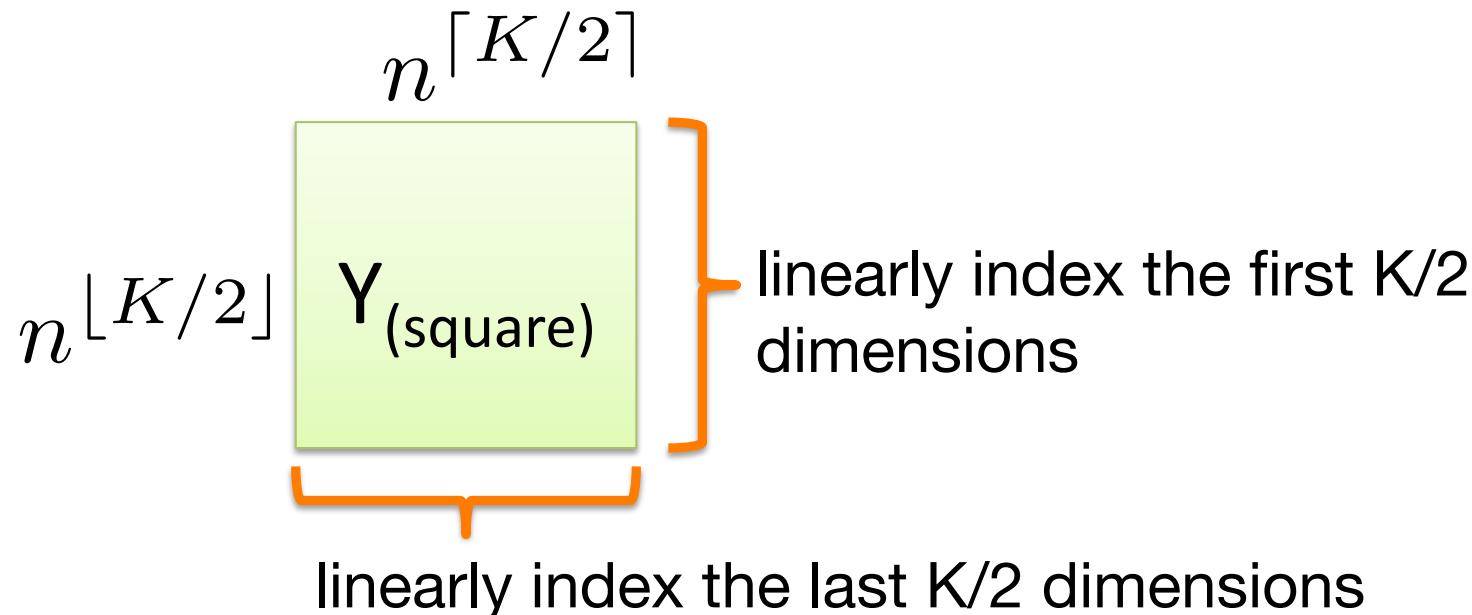
with high probability.

- Note the sharp increase in the required β (trading statistical performance for tractability)

Recursive unfolding

[Richard & Montanari 14]

1. Matricize Y so that it looks as square as possible



2. Compute the leading singular vector v . Note that v is $n^{K/2}$ dimensional.
3. Reshape v into an $n \times n^{K/2-1}$ matrix and return the top-left singular vector.

Performance of (recursive) unfolding

- Richard & Montanari: recursive unfolding attains

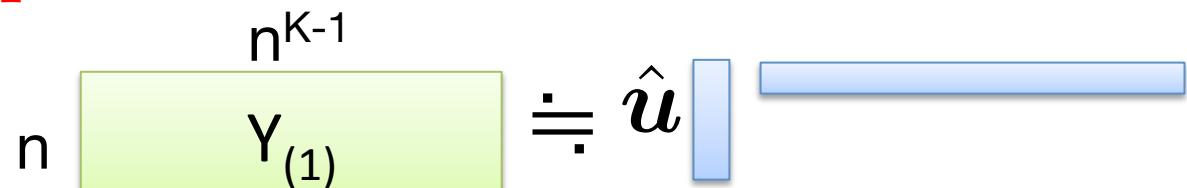
$$1 - |\langle \hat{u}, u^* \rangle| = O\left(\frac{Kn^{\lceil K/2 \rceil}}{\beta^2}\right)$$

under $\beta \geq c \cdot \sqrt{n^{\lceil K/2 \rceil} K}$ with high probability.

– loose for odd-order tensors –conjectured $n^{K/4}$ suffice

- A tighter analysis [Zheng & T, 2015]: taking the top-left singular vector of the ordinary (rectangular) unfolding is sufficient to obtain a threshold

$$\beta \geq c' \cdot n^{K/4}$$



Why simpler method gets better bound

- R&M analyzes the perturbation of

$$Y_{\text{(square)}} = \beta \underbrace{\begin{pmatrix} u^* & \otimes & \cdots & \otimes & u^* \end{pmatrix}^\top}_{\lceil K/2 \rceil} + Z$$
$$\mathbb{E}\|Z\| = O(\sqrt{n^{\lceil K/2 \rceil}})$$

- Z&T analyzes the perturbation of

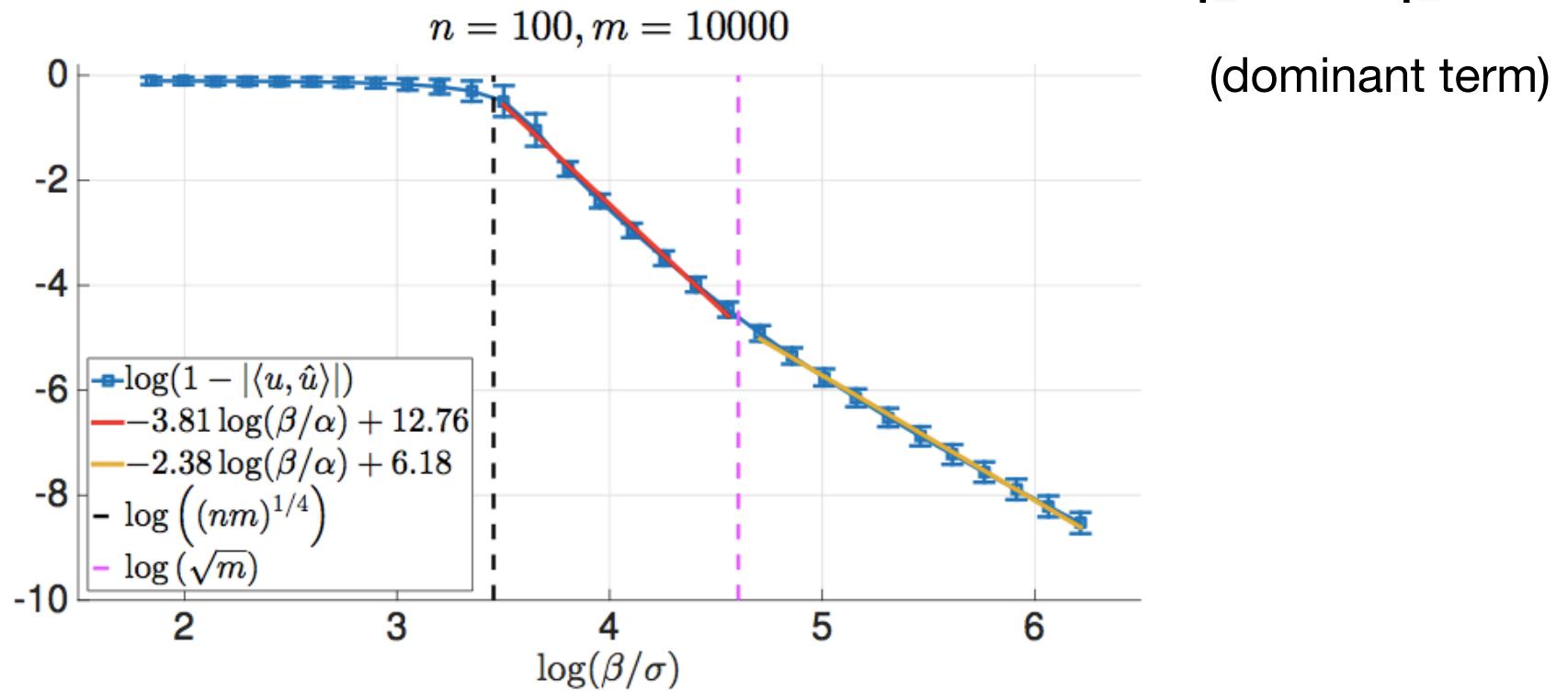
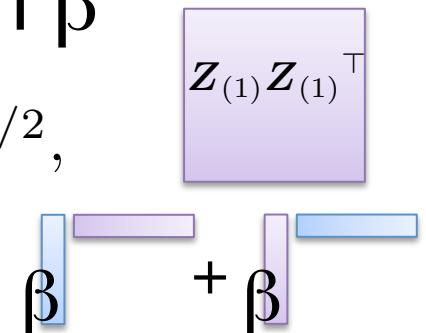
$$n \begin{pmatrix} Y_{(1)} & Y_{(1)}^\top \end{pmatrix} = \beta^2 \begin{pmatrix} u^* & (u^*)^\top \end{pmatrix} + \beta \text{ (cross terms)} + \beta \begin{pmatrix} Z_{(1)} & Z_{(1)}^\top \end{pmatrix}$$

(better controlled perturbation!)

Two-phase behavior

- Because the cross term depends on β

$$1 - |\langle \hat{u}, u^* \rangle| = \begin{cases} O\left(\frac{n^K}{\beta^4}\right), & \text{if } c \cdot n^{K/4} \leq \beta \leq n^{(K-1)/2}, \\ O\left(\frac{n}{\beta^2}\right), & \text{if } n^{(K-1)/2} \leq \beta \end{cases}$$



Summary of rank-one tensor recovery

- Ideal ML estimator: $\beta \geq c \cdot \sqrt{nK \log(K)}$
- Randomly initialized PI: $\beta \geq c' \cdot \sqrt{n^K K \log(K)}$
 - Note: in practice, it seems to be much more robust!
- (Recursive) unfolding: $\beta \geq c'' \cdot n^{K/4}$
- Questions:
 - Is there a fundamental limit how small β can be and still computationally tractable?
 - How about more general problems?

Barak & Moitra (2015) “Tensor Prediction, Rademacher Complexity and Random 3-XOR”

- Studied the atomic norm w.r.t

$$\mathcal{A}_{\text{rank}1,\infty} = \left\{ \mathbf{u} \circ \mathbf{u} \circ \mathbf{u} : \|\mathbf{u}\|_\infty \leq \frac{C}{\sqrt{n}} \right\}$$

atoms are “balanced”

for symmetric $n \times n \times n$ tensors.

- For tensor completion, the atomic norm has almost optimal $\mathcal{O}(n/\varepsilon^2)$ samples complexity but intractable.
- Sixth order sum-of-squares relaxation yields $\mathcal{O}(n^{K/2}/\varepsilon^2)$ sample complexity and polytime – $\mathcal{O}(n^6)$ size SDP.
- Suggested no efficient algorithm can beat $\mathcal{O}(n^{K/2}/\varepsilon^2)$.

Jain & Oh (2014) “Provable Tensor Factorization with Missing Data”

- Showed that alternating minimization algorithm with careful initialization achieves $O(r^5 n^{3/2} (\log(n))^4 \log(r \|W^*\|_F / \varepsilon))$ sample complexity for $K=3$.
 - Slightly high scaling with the rank (probably not optimized)
 - The same $O(n^{K/2})$ scaling
 - Note that if we consider $\sigma^2 = 1/\beta^2 = 1/m$ (inverse sample size), the same $m \geq c'' n^{K/2}$ scaling for (recursive) unfolding.

Conclusion

- Tensors lie in the intersection of a lot of interesting areas –theory, data analysis, machine learning, and numerical linear algebra.
- Many tensor operations are known to be hard but often gradient descent/power iteration works in practice
 - Interaction between computation and statistics: can we trade statistical performance for tractability?
 - Need to understand hardness in contexts.
- Learning with tensors is an area already explored empirically (relational data, neural networks) but not much theoretically

T h a n k y o u !