

## Low-rank approximation: a tool for data modeling

Ivan Markovsky

University of Southampton

1 / 42

## Exact line fitting

the points  $w_i = (x_i, y_i)$ ,  $i = 1, \dots, N$  lie on a line (\*)

$\Leftrightarrow$

there is  $(a, b, c) \neq 0$ , such that  $ax_i + by_i + c = 0$ , for  $i = 1, \dots, N$

$\Leftrightarrow$

there is  $(a, b, c) \neq 0$ , such that 
$$\begin{bmatrix} a & b & c \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_N \\ y_1 & \cdots & y_N \\ 1 & \cdots & 1 \end{bmatrix} = 0$$

$\Leftrightarrow$

$$\text{rank} \left( \begin{bmatrix} x_1 & \cdots & x_N \\ y_1 & \cdots & y_N \\ 1 & \cdots & 1 \end{bmatrix} \right) \leq 2 \quad (**)$$

3 / 42

## Outline

Examples

A setting for data modeling

Solution methods

2 / 42

- restatement of problem (\*) as an equivalent problem (\*\*)
- however, (\*\*) is a standard problem in linear algebra
- the solution generalizes to
  1. **multivariable data** (points in  $\mathbb{R}^q$ ) fitted by an affine set
  2. **time-series fitting** by linear time-invariant dynamical models
  3. data fitting by **nonlinear models**

4 / 42

## Exact conic section fitting

the points  $w_i = (x_i, y_i)$ ,  $i = 1, \dots, N$  lie on a conic section



there are  $A = A^\top$ ,  $b$ ,  $c$ , at least one of them nonzero, such that

$$w_i^\top A w_i + b^\top w_i + c = 0, \text{ for } i = 1, \dots, N$$

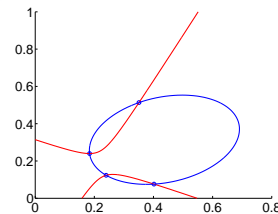


there is  $(a_{11}, a_{12}, a_{22}, b_1, b_2, c) \neq 0$ , such that

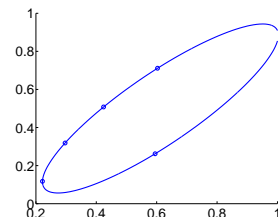
$$\begin{bmatrix} a_{11} & 2a_{12} & b_1 & a_{22} & b_2 & c \end{bmatrix} \begin{bmatrix} x_1^2 & \cdots & x_N^2 \\ x_1 y_1 & \cdots & x_N y_N \\ x_1 & \cdots & x_N \\ y_1^2 & \cdots & y_N^2 \\ y_1 & \cdots & y_N \\ 1 & \cdots & 1 \end{bmatrix} = 0$$

5/42

- $N < 5 \rightsquigarrow$  nonunique fit



- $N = 5$  (different points)  $\rightsquigarrow$  unique fit



- $N > 5 \rightsquigarrow$  generically no conic section fits the data exactly

7/42

the points  $w_i = (x_i, y_i)$ ,  $i = 1, \dots, N$  lie on a conic section



$$\text{rank} \begin{bmatrix} x_1^2 & \cdots & x_N^2 \\ x_1 y_1 & \cdots & x_N y_N \\ x_1 & \cdots & x_N \\ y_1^2 & \cdots & y_N^2 \\ y_1 & \cdots & y_N \\ 1 & \cdots & 1 \end{bmatrix} \leq 5$$

6/42

## Exact fitting by linear homogeneous recurrence relations with constant coefficients

the sequence  $w = (w_1, \dots, w_T)$  is generated by linear recurrence relations with lag  $\leq \ell$



there is  $a = (a_0, a_1, \dots, a_\ell) \neq 0$ , such that

$$a_0 w_i + a_1 w_{i+1} + \cdots + a_\ell w_{i+\ell} = 0, \text{ for } i = 1, \dots, T - \ell$$



there is  $a = (a_0, a_1, \dots, a_\ell) \neq 0$ , such that

$$a^\top \begin{bmatrix} w_1 & w_2 & \cdots & w_{T-\ell} \\ w_2 & w_3 & \cdots & w_{T-\ell+1} \\ \vdots & \vdots & & \vdots \\ w_{\ell+1} & w_{\ell+2} & \cdots & w_T \end{bmatrix} = a^\top \mathcal{H}_\ell(w) = 0$$

8/42

the sequence  $w = (w_1, \dots, w_T)$  is a linear recursion with lag  $\leq \ell$

$$\begin{array}{c} \Downarrow \\ \text{rank} \left( \begin{bmatrix} w_1 & w_2 & \cdots & w_{T-\ell} \\ w_2 & w_3 & \cdots & w_{T-\ell+1} \\ \vdots & \vdots & & \vdots \\ w_{\ell+1} & w_{\ell+2} & \cdots & w_T \end{bmatrix} \right) \leq \ell \end{array}$$

- $T \leq 2\ell \rightsquigarrow$  there is exact fit (independent of  $w$ )
- $T > 2\ell \rightsquigarrow$  generically there is no exact fit

9/42

## Data, model class, and exact fitting test

	line fitting	conic section fitting	linear recurrence with lag $\leq \ell$	GCD
data	points (in $\mathbb{R}^2$ )	points (in $\mathbb{R}^2$ )	sequence	pair of polynomials
model class	lines (in $\mathbb{R}^2$ )	conic sections	autonomous LTI systems	polynomials with nontrivial GCD
exact fitting test	rank condition	rank condition	rank condition	rank condition

11/42

## Existence of greatest common divisor

$$p(z) := p_0 + p_1 z + \cdots + p_m z^m \quad \text{and} \quad q(z) := q_0 + q_1 z + \cdots + q_n z^n$$

have a GCD of degree  $\geq \ell$

$$\begin{array}{c} \Downarrow \\ \dots \\ \Downarrow \\ \text{rank} \left( \begin{bmatrix} p_0 & & & q_0 & & \\ & \ddots & & \vdots & \ddots & \\ & & p_m & q_n & & q_0 \\ & & & \vdots & \ddots & \\ & & & p_m & & q_n \end{bmatrix} \right) \leq m + n - \ell \end{array}$$

10/42

exact fitting test  $\iff$  rank condition

12/42

## Outline

### Examples

### A setting for data modeling

### Solution methods

13/42

## Exact vs approximate models

- $\mathcal{B}$  is an exact model for  $\mathcal{D}$  if  $\mathcal{D} \subset \mathcal{B}$   
otherwise  $\mathcal{B}$  is an approximate model for  $\mathcal{D}$
- $\mathcal{B} = \mathcal{U}$  is a (trivial) exact model for any  $\mathcal{D} \subset \mathcal{U}$   
 $\rightsquigarrow$  we want nontrivial model  
 $\rightsquigarrow$  **notion of model complexity**
- any model is approximate model for any data set  
 $\rightsquigarrow$  we need to quantify the approximation accuracy  
 $\rightsquigarrow$  **notion of model accuracy (w.r.t. the data)**

15/42

## Abstract setting for data modeling

- **data space**  $\mathcal{U}$   
examples:  $\mathbb{R}^q$ ,  $(\mathbb{R}^q)^T$ ,  $\mathbb{R}[\mathbf{z}] \times \mathbb{R}[\mathbf{z}]$ ,  $\{\text{true}, \text{false}\}$
- **data**  $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_N\} \subset \mathcal{U}$   
 $\mathcal{D}_i \in \mathcal{U}$  — observation, realization, or outcome
- **model**  $\mathcal{B} \subset \mathcal{U}$   
an exclusion rule, declares what outcomes are possible
- **model class**  $\mathcal{M} \subset 2^{\mathcal{U}}$

14/42

## Summary

- **data set**  $\mathcal{D} \subset \mathcal{U}$   $\xrightarrow{\text{data modeling problem}}$  **model**  $\mathcal{B} \in \mathcal{M}$ 
  - set of all possible observations  $\mathcal{U}$
  - model class  $\mathcal{M}$
- basic criteria in any data modeling problem are:
  - “simple” model and
  - “good” fit of the data by the model

contradicting objectives
- core issue in data modeling complexity–accuracy trade-off

16/42

## Notes

- in the classical setting, models are viewed as **equations** and a model class is a **parameterized equation**
- in our setting, models are **subsets** of the data space  $\mathcal{U}$  and equations are used as representations of models
- allows us to define equivalence of **model representations**
- establish links among data modeling methods
- model complexity and misfit (lack of fit) b/w data and model have appealing geometrical definitions

17/42

## Linear model complexity

- a linear model  $\mathcal{B}$  is a subspace of  $\mathcal{U}$  ( $\mathcal{U}$  is a vector space)
- the complexity of  $\mathcal{B}$  is defined as its dimension
- in the linear case
 
$$\mathcal{D} \subset \mathcal{B} \implies \text{span}(\mathcal{D}) \subset \mathcal{B}$$
 and the rank of the data matrix is  $\leq \dim(\mathcal{B})$
- $\text{span}(\mathcal{D})$  — the smallest linear model, consistent with  $\mathcal{D}$

19/42

## Model complexity

- the “smaller” a model is the more powerful/useful it is
- the “bigger” a model is the more complex it is
- we prefer simple models over complex ones
- **exact modeling problem:**  
**find the least complex model that fits the data exactly**

18/42

## Model accuracy

- let  $\mathcal{U}$  be a normed vector space with norm  $\|\cdot\|$
- the distance between the data  $\mathcal{D}$  and a model  $\mathcal{B}$

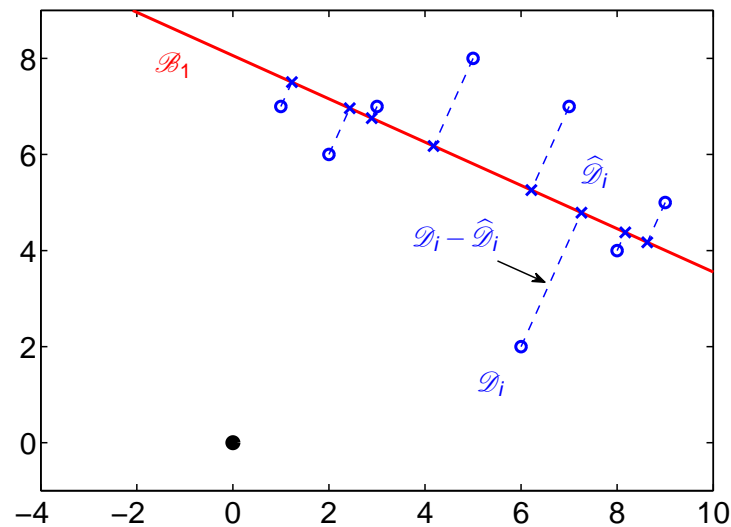
$$\text{dist}(\mathcal{D}, \mathcal{B}) := \min_{\hat{\mathcal{D}} \subset \mathcal{B}} \|\mathcal{D} - \hat{\mathcal{D}}\| \quad (1)$$

measures the lack of fit (misfit) between  $\mathcal{D}$  and  $\mathcal{B}$

- (1) is the projection of the data on the model

20/42

## Example: $\mathcal{U} = \mathbb{R}^2$ , $\mathcal{B}$ linear, Euclidean norm



21/42

## Complexity–accuracy trade-off

- a linear model  $\mathcal{B}$  is a subspace of  $\mathcal{U}$
- a complexity measure of  $\mathcal{B}$  is its dimension —  $\dim(\mathcal{B})$
- misfit — distance from  $\mathcal{D}$  to  $\mathcal{B}$

$$M(\mathcal{D}, \mathcal{B}) := \text{dist}(\mathcal{D}, \mathcal{B}) := \min_{\hat{\mathcal{D}} \in \mathcal{B}} \|\mathcal{D} - \hat{\mathcal{D}}\|_{\mathcal{U}}$$

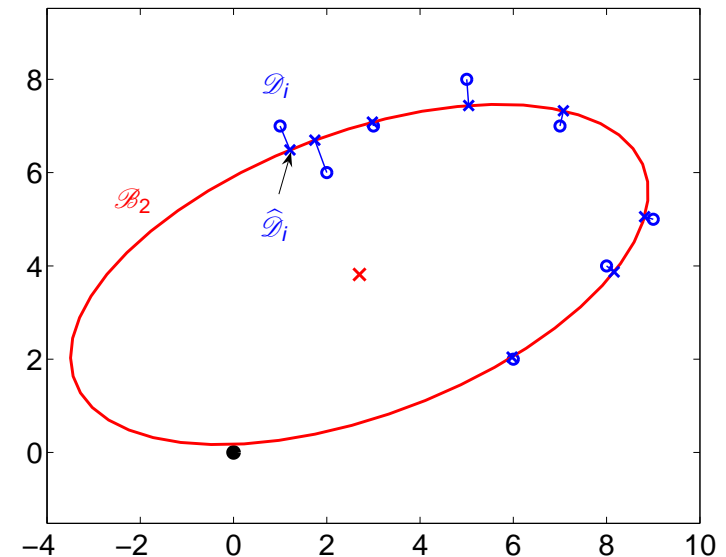
- **data modeling problem:** given  $\mathcal{D} \subset \mathcal{U}$  and  $\|\cdot\|_{\mathcal{U}}$

$$\text{minimize over all linear models } \mathcal{B} \quad \begin{bmatrix} \dim(\mathcal{B}) \\ M(\mathcal{D}, \mathcal{B}) \end{bmatrix} \quad (\text{DM})$$

- a bi-objective optimization problem

23/42

## Example: $\mathcal{U} = \mathbb{R}^2$ , $\mathcal{B}$ quadratic, Euclidean norm



22/42

## The data matrix $\mathcal{S}(p)$

- the data set  $\mathcal{D}$  can be parameterized by a real vector  $p \in \mathbb{R}^{n_p}$  via a map  $\mathcal{S} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{m \times n}$
- $\mathcal{S}$  depends on the application ( $\mathcal{S}$  is affine in case of linear models)
- in static linear modeling problems,  $\mathcal{S}(p)$  is unstructured
- in dynamic LTI modeling problems,  $\mathcal{S}(p)$  is block-Hankel
- fact

$$\dim(\mathcal{B}) \geq \text{rank}(\mathcal{S}(p)) \quad (*)$$

24/42

## The approximation criterion

- $\|\mathcal{D} - \hat{\mathcal{D}}\|_{\mathcal{U}} = \|\mathbf{p} - \hat{\mathbf{p}}\| = \|\tilde{\mathbf{p}}\|$
- weighted 1-, 2-, and  $\infty$ -(semi)norms:
 
$$\|\tilde{\mathbf{p}}\|_{w,1} := \|\mathbf{w} \odot \tilde{\mathbf{p}}\|_1 := \sum_{i=1}^{n_p} |w_i \tilde{p}_i|$$

$$\|\tilde{\mathbf{p}}\|_{w,2} := \|\mathbf{w} \odot \tilde{\mathbf{p}}\|_2 := \sqrt{\sum_{i=1}^{n_p} (w_i \tilde{p}_i)^2}$$

$$\|\tilde{\mathbf{p}}\|_{w,\infty} := \|\mathbf{w} \odot \tilde{\mathbf{p}}\|_{\infty} := \max_{i=1,\dots,n_p} |w_i \tilde{p}_i|$$
- $\mathbf{w}$  — nonnegative vector, specifying the weights
- $\odot$  — element-wise product
- in the stochastic setting of errors-in-variables modeling,  $\|\cdot\|$  corresponds to the distribution of the measurement noise

25 / 42

- (LRA) — low-rank approximation problem
- (RM) — rank minimization problem
- method for solving (RM) can solve (LRA) (using bisection) and vice versa
- varying  $r, e \in [0, \infty)$  the solutions of (LRA) and (RM) sweep the trade-off curve (Pareto optimal solutions of (DM))
- $r$  is discrete and “small”  
 $e$  is continuous and generally unknown
- in applications, an upper bound for  $r$  is often specified

27 / 42

## Low-rank approximation and rank minimization

- (DM) becomes a matrix approximation problem:

$$\text{minimize over } \hat{\mathbf{p}} \quad \begin{bmatrix} \text{rank}(\mathcal{S}(\hat{\mathbf{p}})) \\ \|\mathbf{p} - \hat{\mathbf{p}}\| \end{bmatrix} \quad (\text{DM}')$$

- two possible scalarizations:

1. misfit minimization with a bound  $r$  on the model complexity

$$\text{minimize over } \hat{\mathbf{p}} \quad \|\mathbf{p} - \hat{\mathbf{p}}\| \quad \text{subject to} \quad \text{rank}(\mathcal{S}(\hat{\mathbf{p}})) \leq r \quad (\text{LRA})$$

2. model complexity minimization with a bound  $e$  on the misfit

$$\text{minimize over } \hat{\mathbf{p}} \quad \text{rank}(\mathcal{S}(\hat{\mathbf{p}})) \quad \text{subject to} \quad \|\mathbf{p} - \hat{\mathbf{p}}\| \leq e \quad (\text{RM})$$

26 / 42

## Example: approximate line fitting in $\mathbb{R}^2$

$$\text{minimize over } \mathcal{B} \in \{\text{lines}\} \quad \text{dist}(\mathcal{D}, \mathcal{B})$$

$$\Updownarrow$$

$$\text{minimize over } \hat{x}_i, \hat{y}_i, i = 1, \dots, N \quad \sum_{i=1}^N \left\| \begin{bmatrix} x_i \\ y_i \end{bmatrix} - \begin{bmatrix} \hat{x}_i \\ \hat{y}_i \end{bmatrix} \right\|_2^2$$

$$\text{subject to} \quad \text{rank} \left( \begin{bmatrix} \hat{x}_1 & \cdots & \hat{x}_N \\ \hat{y}_1 & \cdots & \hat{y}_N \\ 1 & \cdots & 1 \end{bmatrix} \right) \leq 2$$

can be solved globally using the singular value decomposition of the data matrix

28 / 42

## Example: approximate conic section fitting in $\mathbb{R}^2$

$$\begin{aligned}
 & \text{minimize} \quad \text{over } \mathcal{B} \in \{\text{conic sections}\} \quad \text{dist}(\mathcal{D}, \mathcal{B}) \\
 & \quad \quad \quad \Updownarrow \\
 & \text{minimize} \quad \text{over } \hat{x}_i, \hat{y}_i, i = 1, \dots, N \quad \sum_{i=1}^N \left\| \begin{bmatrix} x_i \\ y_i \end{bmatrix} - \begin{bmatrix} \hat{x}_i \\ \hat{y}_i \end{bmatrix} \right\|_2^2 \\
 & \text{subject to} \quad \text{rank} \left( \begin{bmatrix} \hat{x}_1^2 & \dots & \hat{x}_N^2 \\ \hat{x}_1 \hat{y}_1 & \dots & \hat{x}_N \hat{y}_N \\ \hat{y}_1^2 & \dots & \hat{y}_N^2 \\ \hat{y}_1 & \dots & \hat{y}_N \\ 1 & \dots & 1 \end{bmatrix} \right) \leq 5
 \end{aligned}$$

29/42

## Algorithms

- with a few exceptions (LRA) and (RM) are non-convex optimization problems
- all general methods are heuristics
- main classes of methods for solving (LRA) and (RM) are:
  - global optimization
  - local optimizations
  - convex relaxations
    - subspace methods and
    - methods based on nuclear norm heuristics

31/42

## Outline

Examples

A setting for data modeling

Solution methods

30/42

## Unstructured low-rank approximation

$$\hat{D}^* := \arg \min_{\hat{D}} \|D - \hat{D}\|_F \quad \text{subject to} \quad \text{rank}(\hat{D}) \leq r$$

Theorem (closed form solution)

Let  $D = U\Sigma V^\top$  be the SVD of  $D$  and define

$$U =: \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{matrix} r & n-r \\ m & \end{matrix}, \quad \Sigma =: \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{matrix} r & n-r \\ n-r & \end{matrix} \quad \text{and} \quad V =: \begin{bmatrix} V_1 & V_2 \end{bmatrix} \begin{matrix} r & n-r \\ m & \end{matrix}$$

An optimal low-rank approximation solution is

$$\hat{D}^* = U_1 \Sigma_1 V_1^\top, \quad (\hat{\mathcal{B}}^* = \ker(U_2^\top) = \text{colspan}(U_1)).$$

It is unique if and only if  $\sigma_r \neq \sigma_{r+1}$ .

32/42



## Structured low-rank approximation

No closed form solution is known for the general SLRA problem

$$\hat{p}^* := \arg \min_{\hat{p}} \|p - \hat{p}\| \quad \text{subject to} \quad \text{rank}(\mathcal{S}(\hat{p})) \leq r.$$

NP-hard, consider solution methods based on local optimization

Representing the constraint in a kernel form, the problem is

$$\min_{R, RR^T = I_{m-r}} \left( \min_{\hat{p}} \|p - \hat{p}\| \quad \text{subject to} \quad R\mathcal{S}(\hat{p}) = 0 \right)$$

Note: Double minimization with bilinear equality constraint.

There is a matrix  $G(R)$ , such that  $R\mathcal{S}(\hat{p}) = 0 \iff G(R)\hat{p} = 0$ .

33/42

## Nuclear norm heuristics

- leads to a semidefinite optimization problem
- existing algorithms with provable convergence properties and readily available high quality software packages
- additional advantage is flexibility: affine inequality constraints in the data modeling problem still leads to semidefinite optimization problems
- disadvantage: the number of optimization variables depends quadratically on the number of data points
- in my experience, the nuclear norm heuristics is less effective than alternative heuristics

35/42

## Variable projection vs. alternating projections

Two ways to approach the double minimization:

- Variable projections (VARPRO):  
solve the inner minimization analytically

$$\min_{R, RR^T = I_{m-r}} \text{vec}^T(R\mathcal{S}(\hat{p})) \left( G(R)G^T(R) \right)^{-1} \text{vec}(R\mathcal{S}(\hat{p}))$$

$\rightsquigarrow$  a nonlinear least squares problem for  $R$  only.

- Alternating projections (AP):  
alternate between solving two least squares problems

VARPRO is globally convergent with a super linear conv. rate.

AP is globally convergent with a linear convergence rate.

34/42

## Nuclear norm heuristics for SLRA

- nuclear norm:  $\|M\|_* = \text{sum of the singular values of } M$
- regularized nuclear norm minimization

$$\begin{aligned} &\text{minimize over } \hat{p} \quad \|\mathcal{S}(\hat{p})\|_* + \gamma \|p - \hat{p}\| \\ &\text{subject to} \quad G\hat{p} \leq h \end{aligned}$$

- using the fact

$$\|M\|_* < \mu \iff \frac{1}{2}(\text{trace}(U) + \text{trace}(V)) < \mu \quad \text{and} \quad \begin{bmatrix} U & M^T \\ M & V \end{bmatrix} \succeq 0$$

we obtain an equivalent SDP problem

$$\begin{aligned} &\text{minimize over } \hat{p}, U, V, v \quad \frac{1}{2}(\text{trace}(U) + \text{trace}(V)) + \gamma v \\ &\text{subject to} \quad \begin{bmatrix} U & \mathcal{S}(\hat{p})^T \\ \mathcal{S}(\hat{p}) & V \end{bmatrix} \succeq 0, \quad \|p - \hat{p}\| < v, \quad G\hat{p} \leq h \end{aligned}$$

36/42

## Nuclear norm heuristics for SLRA

- convex relaxation of (LRA)

$$\text{minimize over } \hat{p} \quad \|p - \hat{p}\| \quad \text{subject to} \quad \|\mathcal{S}(\hat{p})\|_* \leq \mu \quad (\text{RLRA})$$

- motivation: approx. with appropriately chosen bound on the nuclear norm tends to give solutions  $\mathcal{S}(\hat{p})$  of low rank

- (RLRA) can also be written in the equivalent form

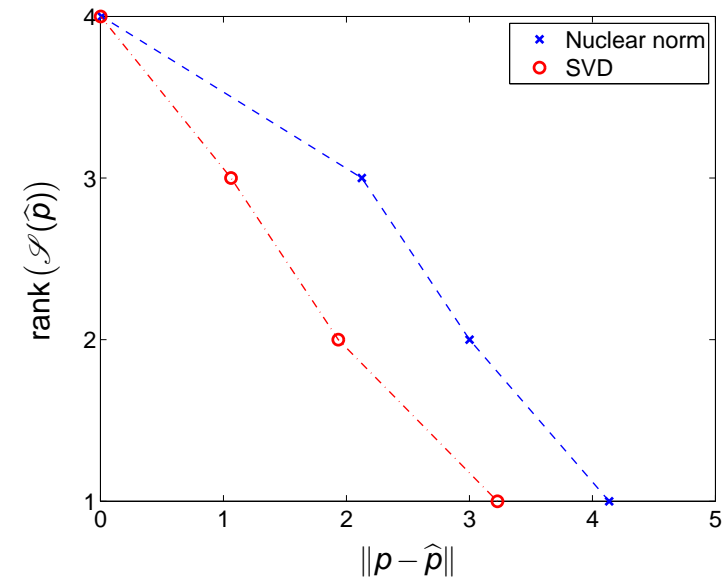
$$\text{minimize over } \hat{p} \quad \|\mathcal{S}(\hat{p})\|_* + \gamma \|p - \hat{p}\| \quad (\text{RLRA}')$$

$\gamma$  — regularization parameter related to  $\mu$  in (RLRA)

- this is a regularized nuclear norm minimization problem

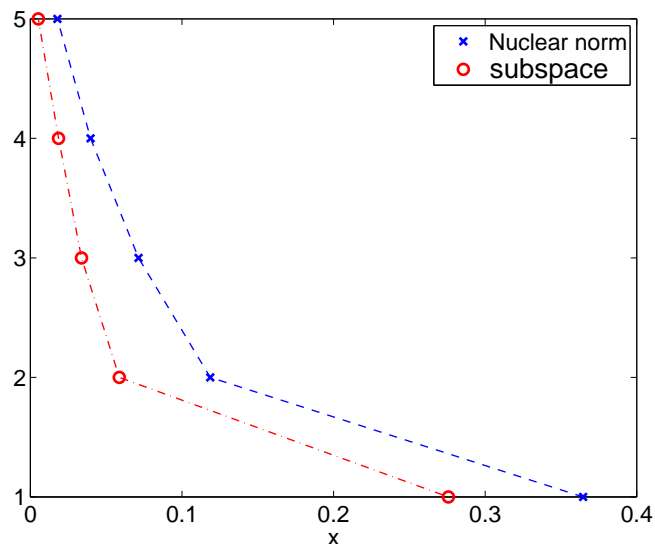
37 / 42

## Unstructured problem's trade-off curve



38 / 42

## Hankel structured problem's trade-off curves



39 / 42

## Conclusions

- common pattern in data modeling

data is exact for a model of bounded complexity



matrix constructed from the data is rank deficient

- exact modeling  $\approx$  rank computation
- approximate modeling is a biobjective opt. problem  
accuracy vs complexity trade-off
- computationally approx. modeling leads to SLRA and RM

40 / 42

- regularized nuclear norm min. is a general and flexible tool
- can be used as a relaxation for low-rank approximation problems with the following desirable features:
  - arbitrary affine structure
  - any weighted 2-norm or even a weighted semi-norm
  - affine inequality constraints
  - regularization
- issues:
  - effectiveness in comparison with other heuristics
  - currently applicable to small sample sizes problems only

## Questions?