# Low-Rank Approximation and Its Applications

Ivan Markovsky

University of Southampton

---

# The simplest data modelling example



$ux_{\text{ls}} = y$ fit

Line fitting problem: Fit the points

$$d_1 = \begin{bmatrix} -2 \\ -6 \end{bmatrix}, \; d_2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \; \ldots, \; d_5 = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

by a line passing through the origin.

Classic solution: Define $d_i =: \text{col}(u_i, y_i)$ and solve the least squares problem

$$x \, \text{col}(u_1, \ldots, u_5) = \text{col}(y_1, \ldots, y_5).$$

The model is the fitting line

$$\mathscr{B} := \{ \, d = \text{col}(u, y) \mid x_{\text{ls}} u = y \, \}$$

---

# Data modelling  $\;\not\Longleftrightarrow\;$  Regression

Obviously,

$\mathscr{B}$ is a line passing through the origin    $\not\Longleftrightarrow$    There is $x \in \mathbb{R}$, such that $\mathscr{B} = \{ \, d = \text{col}(u, y) \mid xu = y \, \}$
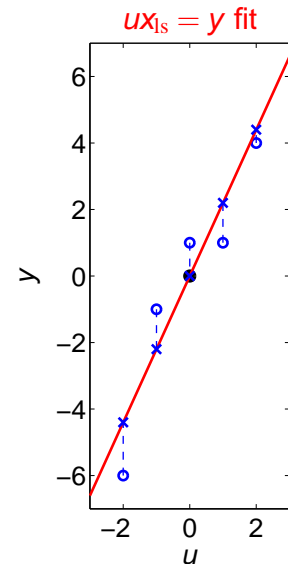
which implies that

Fit the points $d_i = \text{col}(u_i, y_i)$ by a line passing through the origin    $\not\Longleftrightarrow$    Regression $xu \approx y$

Note: Ill-conditioning, a main problem in regression, is a consequence of inadequacy of doing data modelling by regression.

---

# Data modelling  $\;\Longleftrightarrow\;$  low-rank approximation

$\mathscr{B}$ is a line passing through the origin    $\Longleftrightarrow$    $\mathscr{B}$ is a subspace of dimension 1

so that

Fit $d_1, \ldots, d_N$ by a line passing through the origin    $\Longleftrightarrow$    Rank-1 approximation of $D := \begin{bmatrix} d_1 & \cdots & d_N \end{bmatrix}$

An alternative to regression, also known as:

- principal component analysis
- errors-in-variables modeling
- total least squares

# Outline

Introduction

**Applications**

Algorithms

Related problems

# System realisation

The sequence

$$h := \big(h(0), h(1), \dots\big), \qquad h(t) \in \mathbb{R}^{\mathrm{p} \times \mathrm{m}}$$

is realisable by a finite dimensional, linear time-invariant (LTI) system, if and only if

$$\mathscr{H}(h) := \begin{bmatrix} h(1) & h(2) & h(3) & \cdots \\ h(2) & h(3) & \cdot^{\cdot^{\cdot}} & \\ h(3) & \cdot^{\cdot^{\cdot}} & & \\ \vdots & & & \end{bmatrix}$$

has finite rank. Moreover,

$$\mathrm{rank}\big(\mathscr{H}(h)\big) = \text{state dim. of a minimal realisation of } h$$
$$= \text{complexity of an exact LTI model for } h.$$

# Approximate realisation = Model reduction

However, rank deficiency is a nongeneric property (in $\mathbb{Z}_+ \to \mathbb{R}^{\mathrm{p} \times \mathrm{m}}$).

Rank is computed numerically most reliably by the SVD.

From a system theoretic point of view

the SVD does model reduction (Kung's algorithm).

The truncated SVD gives (2-norm) optimal unstructured approx.

Instead, we are aiming at a

structured rank-$\mathrm{n}$ approximation of $\mathscr{H}(h)$:

Find $\widehat{h}$, such that $\|h - \widehat{h}\|$ is minimized and $\mathrm{rank}\big(\mathscr{H}(\widehat{h})\big) = \mathrm{n}$.

# Approximate realisation (model reduction)
# ⇕
# Hankel structured low-rank approximation

The approximate realisation (model reduction) problem is

> Given $h := \big(h(0), h(1), \dots\big)$ and $\mathrm{n} \in \mathbb{N}$, find
>
> $\min_{\widehat{h}} \|h - \widehat{h}\|$  subject to  $\mathrm{rank}\big(\mathscr{H}(\widehat{h})\big) \le \mathrm{n}$

a Hankel structured low-rank approximation (SLRA) problem.

Unfortunately, this problem is NP-complete.

## Deconvolution

Consider the finite sequences

$$h := \big(h(0), h(1), \ldots, h(\mathrm{n})\big), \quad \text{where} \quad h \in \mathbb{R}^{\mathrm{p} \times \mathrm{m}}$$

$$u := \big(u(-\mathrm{n}), \ldots, u(0), u(1) \ldots, u(T)\big) \quad \text{and} \quad y := \big(y(1), \ldots, y(T)\big).$$

Define $\mathrm{row}(y) := \begin{bmatrix} y(1) & \cdots & y(T) \end{bmatrix}$ and the Toeplitz matrix

$$\mathscr{T}_{\mathrm{n}+1}(u) := \begin{bmatrix} u(1) & u(2) & u(3) & \ldots & u(T) \\ u(0) & u(1) & u(2) & \ldots & u(T-1) \\ \vdots & \vdots & \vdots & & \vdots \\ u(-\mathrm{n}) & u(1-\mathrm{n}) & u(2-\mathrm{n}) & \cdots & u(T-\mathrm{n}) \end{bmatrix}$$

With this notation,

$$\begin{array}{ccc} y = h \star u & & \mathrm{row}(y) = \mathrm{row}(h)\mathscr{T}_{\mathrm{n}+1}(u) \\ \text{(convolution)} & \Longleftrightarrow & \text{(linear algebra)} \end{array}$$

## Exact and approximate deconvolution

Exact deconv. problem: Given $u$ and $y$, find $h$, such that $y = h \star u$.

Solution exists if and only if the system of equations

$$\mathrm{row}(y) = \mathrm{row}(h)\mathscr{T}_{\mathrm{n}+1}(u)$$

is solvable for $h$. However with $T > (\mathrm{n}+1)\mathrm{m}$, generically solution does not exist $\rightsquigarrow$ approximate deconvolution problem:

Given $u$, $y$, and $\mathrm{n} \in \mathbb{N}$, find

$$\min_{\widehat{u}, \, \widehat{y}, \, \widehat{h}} \| \mathrm{col}(u, y) - \mathrm{col}(\widehat{u}, \widehat{y}) \| \quad \text{subject to}$$

$$\mathrm{row}(\widehat{y}) = \mathrm{row}(\widehat{h})\mathscr{T}_{\mathrm{n}+1}(\widehat{u})$$

## Deconvolution = FIR system identification

We can interpret

$$y = h \star u$$

as the response of an FIR system with impulse response $h$ to
- initial conditions $\big(u(-\mathrm{n}), \ldots, u(0)\big)$, and
- input $\big(u(1) \ldots, u(T)\big)$.

Then the deconvolution problem has the meaning of an
FIR system identification problem:

Given initial condition, input, and output, find an FIR model.

- exact deconvolution $\implies$ exact FIR fitting model
- approx. deconvolution $\implies$ approx. FIR fitting model

The parameter $\mathrm{n}$ bounds the FIR model complexity.

## Approximate deconvolution $\rightsquigarrow$ SLRA

Assuming that $\mathscr{T}_{\mathrm{n}+1}(\widehat{u})$ is full rank (persistency of excitation),

$$\mathrm{row}(\widehat{y}) = \mathrm{row}(\widehat{h})\mathscr{T}_{\mathrm{n}+1}(\widehat{u}) \quad \Longleftrightarrow \quad \mathrm{rank}\left( \begin{bmatrix} \mathscr{T}_{\mathrm{n}+1}(\widehat{u}) \\ \mathrm{row}(\widehat{y}) \end{bmatrix} \right) = (\mathrm{n}+1)\mathrm{m}$$

Then the approximate deconvolution problem can be written as

Given $u$, $y$, and $\mathrm{n} \in \mathbb{N}$, find

$$\min_{\widehat{u}, \, \widehat{y}} \| \mathrm{col}(u, y) - \mathrm{col}(\widehat{u}, \widehat{y}) \| \quad \text{subject to}$$

$$\mathrm{rank}\left( \begin{bmatrix} \mathscr{T}_{\mathrm{n}+1}(\widehat{u}) \\ \mathrm{row}(\widehat{y}) \end{bmatrix} \right) \leq (\mathrm{n}+1)\mathrm{m}$$

a SLRA problem with structure composed of two blocks:
Toeplitz block and an unstructured block.

## Greatest common divisor (GCD)

Consider the polynomials

$$a(z) := a_0 + a_1 z + \cdots + a_{n_a} z^{n_a}, \quad b(z) := b_0 + b_1 z + \cdots + b_{n_b} z^{n_b}$$

and define the Sylvester matrix

$$S(a,b) := \begin{bmatrix} a_0 & & & b_0 & & \\ \vdots & \ddots & & \vdots & \ddots & \\ a_{n_a} & & a_0 & b_{n_b} & & b_0 \\ & \ddots & \vdots & & \ddots & \vdots \\ & & a_{n_a} & & & b_{n_b} \end{bmatrix} \in \mathbb{R}^{(n_a+n_b)\times(n_a+n_b)}$$

The GCD of $a(z)$ and $b(z)$, has degree $n$, if and only if

$$\text{rank}\left(S(a,b)\right) = n_a + n_b - n.$$

## Approximate GCD $\Longleftrightarrow$ Sylvester SLRA

Given $a(z)$, $b(z)$, and $n \in \mathbb{N}$, find

$$\min_{\widehat{a},\,\widehat{b}} \ \|\text{col}(a,b) - \text{col}(\widehat{a},\widehat{b})\| \quad \text{subject to}$$

$$\text{rank}\left(S(a,b)\right) \leq n_a + n_b - n$$

## Data matrix being low-rank

an exact property $\quad\Longleftrightarrow\quad$ a matrix constructed
holds on the data $\quad\quad\quad\quad$ from data is low-rank

- $h$ is realisable by an $\quad\Longleftrightarrow\quad$ $\text{rank}\left(\mathscr{H}(h)\right) \leq \text{n}$
  LTI system of order $\text{n}$

- $(u,y)$ is fitted by an $\quad\Longleftrightarrow\quad$ $\text{rank}\left(\begin{bmatrix}\mathscr{T}_{\text{n}+1}(u)\\ \text{row}(y)\end{bmatrix}\right) \leq (\text{n}+1)\text{m}$
  $\text{n}$ taps FIR system

- $a(z), b(z)$ have $\quad\Longleftrightarrow\quad$ $\text{rank}\left(S(a,b)\right) \leq n_a + n_b - n$
  GCD of deg. $\geq n$

## Rank of the data matrix

complexity of an exact $\quad\leftrightarrow\quad$ rank of the
model fitting the data $\quad\quad\quad\quad$ data matrix

- order of the realization $\quad=\quad$ $\text{rank}\left(\mathscr{H}(h)\right)$

- number of taps $\quad=\quad$ $\text{rank}\left(\begin{bmatrix}\mathscr{T}_{\text{n}+1}(u)\\ \text{row}(y)\end{bmatrix}\right)/m - 1$
  of an FIR system

- degree of the GCD $\quad=\quad$ rank deficiency of $S(a,b)$

# Main issue: Low-rank approximation

With a bounding on the model complexity,

> generically in the data space, exact property does not hold

$\implies$ an approximation is needed.

## Approximation paradigm:

> modify the data as little as possible, so that the exact property holds for the modified data.

This paradigm leads to structured low-rank approximation.

# Structured low-rank approximation

## Given

- a vector $p \in \mathbb{R}^{n_p}$,
- a mapping $\mathscr{S} : \mathbb{R}^{n_p} \to \mathbb{R}^{m \times n}$ (structure specification)
- a vector norm $\| \cdot \|$, and
- an integer $r$, $0 < r < \min(m, n)$,

## find

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

## Interpretation:

$\widehat{D}^* := \mathscr{S}(\widehat{p}^*)$ is optimal rank-$r$ (or less) approx. of $D := \mathscr{S}(p)$, within the class of matrices with the same structure as $D$.

# Outline

Introduction

Applications

Algorithms

Related problems

# Unstructured low-rank approximation

$$\widehat{D}^* := \arg\min_{\widehat{D}} \| D - \widehat{D} \|_{\mathrm{F}} \quad \text{subject to} \quad \text{rank}(\widehat{D}) \leq r$$

## Theorem (closed form solution)

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

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

An optimal low-rank approximation solution is

$$\widehat{D}^* = U_1 \Sigma_1 V_1^\top, \qquad (\widehat{\mathscr{B}}^* = \ker(U_2^\top) = \text{col span}(U_1)).$$

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

## Structured low-rank approximation

No closed form solution is known for the general SLRA problem

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

NP-hard, consider solution methods based on local optimization

Representing the constraint in a kernel form, the problem is

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

Note: Double minimization with bilinear equality constraint.

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

## Variable projection vs. alternating projections

Two ways to approach the double minimization:

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

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

  $\rightsquigarrow$ 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.

## Variations on low-rank approximation

- Cost functions
  - weighted norms $\quad (\text{vec}^\top(D)W\text{vec}(D))$
  - information criteria $\quad (\log\det(D))$

- Constraints and structures
  - nonnegative
  - sparse

- Data structures
  - nonlinear models
  - tensors

- Optimization algorithms
  - convex relaxations

## Weighted low-rank approximation

In the measurement error model,

$$d_i = \overline{d}_i + \widetilde{d}_i, \quad \overline{d}_i \in \overline{\mathscr{B}}, \quad \widetilde{d}_i \sim \text{Normal}(0, \sigma^2 V_i)$$

the basic low-rank approximation is maximum likelihood estimator assuming $V_i = I$.

Motivation: incorporate prior knowledge $V$ about $\text{cov}(\text{vec}(\widetilde{D}))$

$$\min_{\widehat{D}} \text{vec}^\top(D - \widehat{D}) V^{-1} \text{vec}(D - \widehat{D}) \quad \text{subject to} \quad \text{rank}(\widehat{D}) \leq r$$

Known in chemometrics as maximum likelihood PCA.

NP-hard problem, alternating projections is effective heuristic

## Nonnegative low-rank approximation

Constrained LRA arise in Markov chains and image mining

$$\min_{\widehat{D}} \|D - \widehat{D}\| \quad \text{subject to} \quad \text{rank}(\widehat{D}) \leq r \text{ and } \widehat{D}_{ij} \geq 0 \text{ for all } i,j.$$

Using an image representation, an equivalent problem is

$$\min_{P \in \mathbb{R}^{m \times r}, \, L \in \mathbb{R}^{r \times n}} \|D - PL\| \quad \text{subject to} \quad P_{ik}, L_{kj} \geq 0 \text{ for all } i,k,j.$$
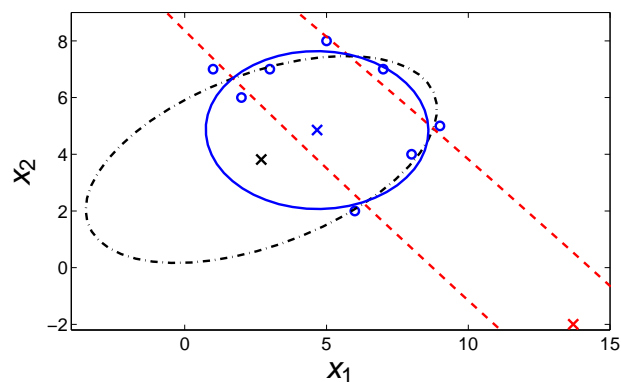
Alternating projections algorithm:

- Choose an initial approximation $P^{(0)}$ and set $k := 0$.
- Solve: $L^{(k)} = \arg\min_L \|D - P^{(k)}L\|$ subject to $L \geq 0$.
- Solve: $P^{(k+1)} = \arg\min_P \|D - PL^{(k)}\|$ subject to $P \geq 0$.
- Repeat until convergence.

---

## Data fitting by a second order model

$$\mathscr{B}(A,b,c) := \{ d \in \mathbb{R}^{\text{d}} \mid d^\top A d + b^\top d + c = 0 \}, \quad \text{with } A = A^\top$$

Consider first exact data:

$$d \in \mathscr{B}(A,b,c) \iff d^\top A d + b^\top d + c = 0$$
$$\iff \big\langle \underbrace{\text{col}(d \otimes_{\text{s}} d, d, 1)}_{d_{\text{ext}}}, \underbrace{\text{col}\big(\text{vec}_{\text{s}}(A), b, c\big)}_{\theta} \big\rangle = 0$$

$$\{ d_1, \ldots, d_N \} \in \mathscr{B}(\theta) \iff \theta \in \text{left ker} \underbrace{\big[ d_{\text{ext},1} \;\; \cdots \;\; d_{\text{ext},N} \big]}_{D_{\text{ext}}}, \quad \theta \neq 0$$

$$\iff \text{rank}(D_{\text{ext}}) \leq \text{d} - 1$$

Therefore, for measured data ⤳ LRA of $D_{\text{ext}}$.

Notes:

- Special case $\mathscr{B}$ an ellipsoid (for $A > 0$ and $4c < b^\top A^{-1} b$).
- Related to kernel PCA

---

## Example: ellipsoid fitting

benchmark example of (Gander *et al.* 94), called "special data"



dashed — LRA     solid — modified LRA

dashed-dotted — orthogonal regression (geometric fitting)

○ — data points     × — centers

---

## Rank minimization

Approximate modeling is a trade-off between:

- fitting accuracy and
- model complexity

Two possible scalarizations of the bi-objective optimization are:

LRA: minimize misfit under a constraint on complexity

RM: minimize complexity under a constraint ($\mathscr{C}$) on misfit

$$\min_X \text{rank}(X) \quad \text{subject to} \quad X \in \mathscr{C}$$

RM is also NP-hard, however, there are effective heuristics, *e.g.*,

with $X = \text{diag}(x)$, $\text{rank}(X) = \text{card}(x)$,

$$\ell_1 \text{ heuristic:} \quad \min_x \|x\|_1 \quad \text{subject to} \quad \text{diag}(x) \in \mathscr{C}$$

## Summary

- SLRA is a generic problem for data modeling.

  search for more applications (pole placement, $\mu$-analysis, . . . )

- In general, SLRA is an NP-complete problem.

  search for special cases that have "nice" solutions
  *e.g.,* circulant SLRA can be computed by DFT.

- The SLRA framework leads to conceptual unification.

## Summary

- Efficient local solution methods

- Different rank representations (kernel, image, $AX = B$)
  lead to equivalent parameter optimization problems.

  Computationally, however, these problems are different.

  For example, the kernel representation leads to
  optimization on a Grassman manifold.

  Currently, it is unexplored which parameterization is
  computational most beneficial.

## Summary

- Effective heuristics, based on convex relaxations

- Practical advantage: one algorithm (and a piece of
  software) can solve a variety of problems

- Extensions of SLRA for tensors and nonlinear models

  A framework with a potential for much to be done.

# Thank you