

Data-driven modeling: A low-rank approximation problem

Ivan Markovsky

Vrije Universiteit Brussel

Outline

Setup: data-driven modeling

Problems: system identification, machine learning, ...

Behavioral paradigm \leftrightarrow low-rank approximation

Algorithms: optimization, multistage, convex relaxations

Applications: missing data, data-driven simulation

Connections: TLS, EIV, PCA, rank minimization, ...

General setup



- \mathcal{D} — data, e.g., a vector time series $(\mathbb{R}^q)^{\mathbb{N}}$
- \mathcal{B} — model (behavior): a (sub)set of the data space \mathcal{U}
- \mathcal{M} — model class: a set of models

work plan:

1. define a modeling problem
2. find an algorithm that solves the problem
3. implement the algorithm in software
4. use the software in applications

The problem

prior knowledge, assumptions, and/or prejudices

about what the true or desirable model is

- **model class** — imposes hard constraints
e.g., bound on the model complexity
- **optimization criteria** — impose soft constraints
e.g., small misfit between the model and the data
- real-life problems are vaguely formulated
- often it is not clear what is the “best” problem formulation

“A well defined problem is a half solved problem.”

System identification problems

$$\mathcal{U} = \underbrace{(\mathbb{R}^q \times \dots \times \mathbb{R}^q)}_{T_1} \times \dots \times \underbrace{(\mathbb{R}^q \times \dots \times \mathbb{R}^q)}_{T_N} \quad \text{--- } N, q\text{-variable time series}$$

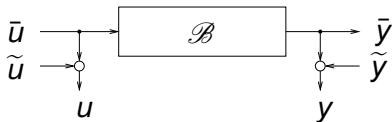
\mathcal{M} is, e.g., bounded complexity (# inputs and lags), LTI systems

- latency (ARMAX):



$$\text{minimize } \|e\| \quad \text{subject to } ((e, u), y) \in \widehat{\mathcal{B}}_{\text{ext}} \in \mathcal{M}$$

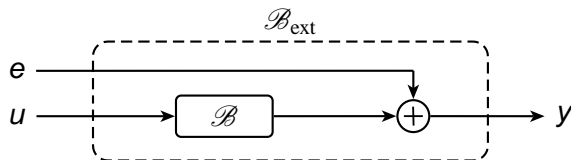
- misfit (EIV):



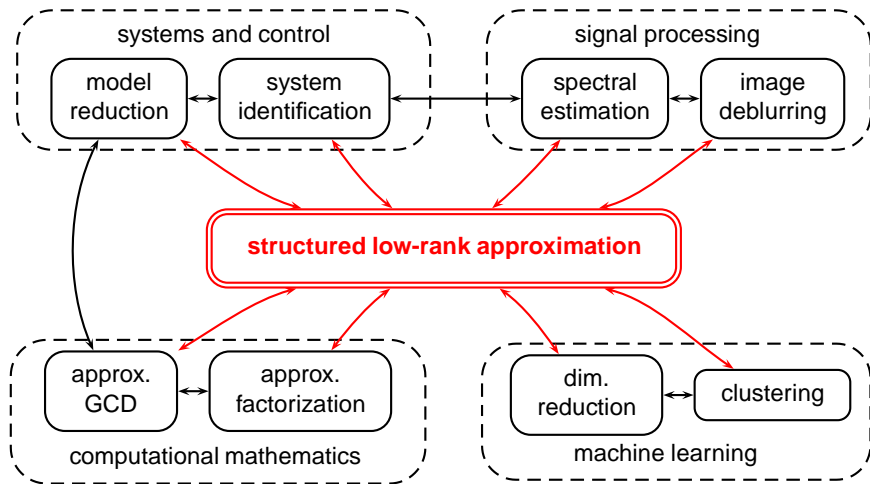
$$\text{minimize } \|(\Delta u, \Delta y)\| \quad \text{subject to } \underbrace{(u + \Delta u)}_{\hat{u}}, \underbrace{(y + \Delta y)}_{\hat{y}} \in \widehat{\mathcal{B}} \in \mathcal{M}$$

Special cases

- \mathcal{M} with lag = 0 \rightsquigarrow static modeling
- \mathcal{M} with # inputs = 0 \rightsquigarrow sum-of-damped-exp. modeling
- FIR systems \rightsquigarrow approximate deconvolution
- EIV with $\Delta u = 0$ or special ARMAX \rightsquigarrow output error



A unifying setting for data modeling



Desirable features of a paradigm

- simple:** can be introduced in 1 slide
- flexible:** applies to a rich class of problems
- practical:** leads to solution methods and algorithms
- optimal:** in theory, finds the "best" solution
- effective:** in practice, can "solve" real-life problems
- automatic:** hyper param. correspond to prior knowledge
- compact:** software implementation requires short code

Structured low-rank approximation

- structure specification $\mathcal{S} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{m \times n}$
- vector of structure parameters $p \in \mathbb{R}^{n_p}$
- weighted 2-norm $\|p\|_w^2 := p^\top W p$
- rank specification r

$$\begin{array}{ll} \text{minimize} & \text{over } \hat{p} \in \mathbb{R}^{n_p} \quad \|p - \hat{p}\|_w^2 \\ \text{subject to} & \text{rank}(\mathcal{S}(\hat{p})) \leq r \end{array} \quad (\text{SLRA})$$

Structure \mathcal{S} \leftrightarrow Model class \mathcal{M}

unstructured

 \leftrightarrow

linear static

Hankel

 \leftrightarrow

scalar LTI

 $q \times 1$ Hankel \leftrightarrow q -variate LTI $q \times N$ Hankel \leftrightarrow N equal length traj.

mosaic Hankel

 \leftrightarrow N general traj.

[Hankel unstructured]

 \leftrightarrow

finite impulse response

block-Hankel Hankel-block

 \leftrightarrow

2D linear shift-invariant

(SLRA) \leftrightarrow approximate data modeling

- $p \leftrightarrow \text{vec}(\mathcal{D})$
- $r \leftrightarrow$ model complexity
- $W \leftrightarrow$ prior knowledge about the data accuracy

(SLRA) is a maximum likelihood estimator in the EIV setting

Singular weight matrix \leftrightarrow fixed and missing values

- consider the special case of element-wise weights

$$\|p - \hat{p}\|_w = \sqrt{\sum_{i=1}^{n_p} w_i (p_i - \hat{p}_i)^2}$$

specified by a vector $w \in \mathbb{R}^{n_p}$

- $w_i = \infty$ imposes equality constraint $\hat{p}_i = p_i$ on (SLRA)

$$w_i = \infty \quad \implies \quad \hat{p}_i = p_i$$

- $w_i = 0$ makes the problem (SLRA) independent of p_i

$$w_i = 0 \quad \implies \quad p_i \text{ is ignored}$$

alternatively, problem (SLRA) is solved with p_i missing

Solution methods

- global solution methods

- SDP relaxations of rational function minimization problem
- systems of polynomial equations (computer algebra)
 - resultant-based methods
 - Stetter-Moller methods
 - subdivision methods
 - homotopy continuation

- local optimization methods

- variable projections
- alternating projections
- variations

parameterization
+
optimization method
=
method

- heuristics

- multistage methods
- nuclear norm heuristic

VARPRO-like solution method

- using the kernel parameterization

$$\text{rank}(\mathcal{J}(\hat{p})) \leq r \iff R\mathcal{J}(\hat{p}) = 0, \quad \text{rank}(R) = m - r$$

- (SLRA) becomes

$$\begin{aligned} &\text{minimize} && \text{over } \hat{p} \text{ and } R && \|p - \hat{p}\|_w^2 \\ &\text{subject to} && R\mathcal{J}(\hat{p}) = 0, \text{rank}(R) = m - r && \end{aligned} \quad (\text{SLRA}_R)$$

- (SLRA_R) is separable in \hat{p} and R , *i.e.*, it is equivalent to

$$\begin{aligned} &\text{minimize} && \text{over } R && f(R) \\ &\text{subject to} && \text{rank}(R) = m - r && \end{aligned} \quad (\text{OUTER})$$

where

$$f(R) := \min_{\hat{p}} \|p - \hat{p}\|_w^2 \quad \text{subject to} \quad R\mathcal{J}(\hat{p}) = 0 \quad (\text{INNER})$$

- \hat{p} is eliminated (projected out) of (SLRA_R)

- evaluation of $f(R)$, *i.e.*, solving (INNER), is least norm prob.
- in SYSID, evaluation of $f(R)$ is a **data smoothing operation**
- in a stochastic setting, it is the **likelihood evaluation**
- efficient computation using **Riccati recursion**
(**Kalman smoothing**)
- in other applications, $f(R)$ can also be evaluate efficiently, by exploiting the matrix structure
- software implementation for mosaic Hankel-like matrices, with fixed and missing data, and linearly structured kernel

<http://github.com/slra/slra>

Structured kernel

- (OUTER) is a nonlinear least-squares problem
- it can be solved with additional constraints
- e.g., linear structure of the kernel

$$R = \mathcal{R}(\theta) := \text{vec}^{-1}(\theta\Psi)$$

- applications requiring structured kernel:

- harmonic retrieval \rightsquigarrow

R palindromic

- SYSID with fixed poles \rightsquigarrow

$$R = R_{\text{fixed}} \star R_{\text{free}}$$

- SYSID with fixed observ. indices \rightsquigarrow

$$R = \begin{bmatrix} \times & \dots & \times & 1 & 0 & 0 \\ \vdots & \ddots & & \ddots & \ddots & 0 \\ \times & \dots & \times & \dots & \times & 1 \end{bmatrix}$$

- common dynamics estimation \rightsquigarrow

\mathcal{R} nonlinear

Autonomous system identification with missing data

- $\mathcal{M} = \mathcal{L}_{0,\ell}$ — LTI systems with 0 inputs and lag $\leq \ell$
- data $y \in \underbrace{\mathbb{R}_{\text{ext}}^{\mathbb{P}} \times \cdots \times \mathbb{R}_{\text{ext}}^{\mathbb{P}}}_T$, where $\mathbb{R}_{\text{ext}} = \mathbb{R} \cup \text{NaN}$
- **problem:** given y and ℓ ,

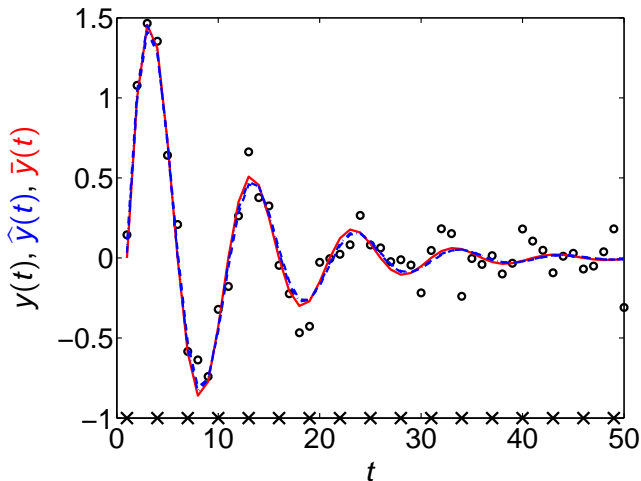
$$\begin{aligned} & \text{minimize} && \text{over } \hat{y} \in (\mathbb{R}^{\mathbb{P}})^T \text{ and } \hat{\mathcal{B}} && \|y - \hat{y}\|_w^2 \\ & \text{subject to} && \hat{y} \in \hat{\mathcal{B}} \in \mathcal{L}_{0,\ell} \end{aligned}$$

- w assigns zeros to the missing data ($y_i(t) = \text{NaN}$)
- $\exists \hat{\mathcal{B}}$, such that $\hat{y} \in \hat{\mathcal{B}} \in \mathcal{L}_{0,\ell} \iff \text{rank}(\mathcal{H}_{\ell+1}(\hat{y})) \leq \ell_{\mathbb{P}}$
- the problem is Hankel structured low-rank approximation

Simulation example

- $p = 1$, $\ell = 2$, $T = 50$, $y = \bar{y} + \text{white noise}$, where
$$\bar{y}(t) = 1.456\bar{y}(t-1) - 0.81\bar{y}(t-2), \quad \bar{y}(0) = 0, \quad \bar{y}(1) = 1$$
- missing values distributed periodically with period 3
- solved with the algorithm based on the VARPRO approach

System identification with periodically missing data



true — solid line

optimal approximation — dashed blue

circles — data points

crosses — location of missing data

Classical simulation problem

given

- LTI system \mathcal{B} (specified by some representation)
- initial condition w_{ini} (specified by trajectory of \mathcal{B})
- input u

find the output y of \mathcal{B} , corresponding to w_{ini} and u

- there are many ways to solve the problem
- the algorithms depend on the model representation (state-space, transfer function, impulse response, ...)

Data-driven simulation

given

- trajectory w' of LTI system \mathcal{B} and the lag ℓ of \mathcal{B}
- initial condition $w_p'' = (w''(1), \dots, w''(\ell))$
- input $u_f'' = (u''(\ell+1), \dots, u''(T_2))$

find the output y_f'' of \mathcal{B} , corresponding to w_p'' and u''

$$y_f'' = (y''(\ell+1), \dots, y''(T_2))$$



find y_f'' and $\mathcal{B} \in \mathcal{L}_{m,\ell}$

such that $w' \in \hat{\mathcal{B}}$ and $\underbrace{w_p'' \wedge (u_f'', y_f'')}_{w''} \in \mathcal{B}$

- there is $\widehat{\mathcal{B}} \in \mathcal{L}_{m,l}$, such that $w' \in \widehat{\mathcal{B}}$ and $w'' \in \widehat{\mathcal{B}}$



$$\text{rank} \left(\begin{bmatrix} \mathcal{H}_{l+1}(w') & \mathcal{H}_{l+1}(w'') \end{bmatrix} \right) \leq 2l + 1$$

mosaic Hankel matrix completion

- with noisy w' , the problem is

$$\begin{aligned} & \text{minimize} && \text{over } \widehat{w}', \widehat{w}'', \widehat{\mathcal{B}} \in \mathcal{L}_{m,l} && \|w' - \widehat{w}'\|_2^2 \\ & \text{subject to} && \widehat{w}', \widehat{w}'' \in \widehat{\mathcal{B}}, && \widehat{w}''_p = w''_p, \quad \widehat{u}''_f = u''_f \end{aligned}$$

mosaic Hankel low-rank approximation
with exact and missing data

Simulation example

- second order SISO system, defined by difference equation

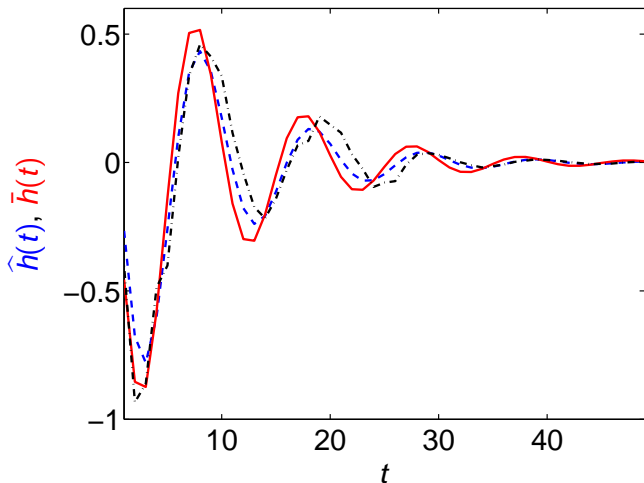
$$\bar{y}(t) = 1.456\bar{y}(t-1) - 0.81\bar{y}(t-2) + \bar{u}(t) - \bar{u}(t-1)$$

- w' is noisy trajectory generated from random input
- y_f'' is the impulse response \bar{h} , i.e.,

$$u'' = (\underbrace{0, \dots, 0}_\ell, \underbrace{1, 0, \dots, 0}_{\text{pulse input}})$$

$$y'' = (\underbrace{0, \dots, 0}_\ell, \underbrace{\hat{h}(0), \hat{h}(1), \dots, \hat{h}(T_2 - \ell - 1)}_{\text{impulse response}})$$

Data-driven simulation of impulse response



true — solid line

optimal approximation — dashed blue

Related frameworks

- **behavioral approach:** representation free modeling
- **total least squares:** (SLRA) with I/O representation

$$R\mathcal{S}(\hat{p}) = [X^T \quad -I] \begin{bmatrix} \hat{A}^T \\ \hat{B}^T \end{bmatrix} = 0 \quad \iff \quad \hat{A}X = \hat{B} \quad (\text{TLS})$$

- **errors-in-variables:** statistical setup for (TLS)
- **principal component analysis:** another statistical setup
- **rank minimization:** “dual” to (SLRA)
(soft constraint on complexity, hard constraint on accuracy)

Current/future work

- comparison of different optimization methods for

$$\text{minimize over } R \quad M(R) \quad \text{subject to} \quad RR^{\top} = I$$

- fast misfit computation \leftrightarrow Kalman smoothing

$$M(R) = \text{vec}^{\top}(w)\Gamma^{-1}(R)\text{vec}(w)$$

efficient computation in the case of missing data

- singularity of Γ (poles/zeros on the unit circle)
- solution of ARMAX identification problems (latency min.)
- static nonlinear modeling is nonlinear SLRA (kernel PCA)
- nD system identification (block-Hankel Hankel-block SLRA)