# Errors-in-variables modelling

## Ivan Markovsky

### University of Southampton

---

# Outline

Introduction

Static model identification

LTI model identification

Algorithms

Nonlinear models

---

# Data modelling setup

Given:

- data generating system       the "true" system $\Sigma_{\text{true}}$
- from which data is observed       the data $\{ w^1, \ldots, w^N \}$
- set of candidate models       the model class $\mathcal{M}$

Choose:

- a model in the model class that
- approximates "well" the true system

However, since the true system is unknown,

- the model is chosen to approximate "well" the data

---

# User choices

- data preprocessing  (centering, scaling, filtering, . . . )
- model class
- fitting criterion

often made by heuristic rules or by trail and error (experience)

User choices correspond to

- prior knowledge and/or
- assumptions

about the true system

System identification theory aims to

- justify particular fitting criteria  (statistics)
- derive algorithms  (optimization, numerical methods)

## Classical paradigm

Assumptions:

- "true model in the model class" assumption
  the data generating system belongs to the model class
- input/output partitioning of the variables is a priori given
- the data-model mismatch is a stochastic process

Two variations:

- treat observed inputs as exact — regression
  The uncertainty is attributed to unobserved latent inputs.
- treating all variables as noisy — errors-in-variables model
  The uncertainty is attributed to the measurement noise.

Latency model:    $e$ — process noise (unobserved)



EIV model:    $\widetilde{w} = (\widetilde{u}, \widetilde{y})$ — measurement noise



$\Sigma_{\text{true}}$ — data generating system,  $w_{\text{d}} = (u_{\text{d}}, y_{\text{d}})$ — observed data

## Main results in the classical paradigm
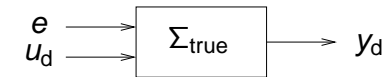
Modelling methods that are

- consistent
- efficient and
- produce confidence bounds

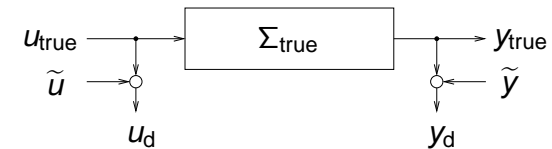under specified assumptions on the true model and the errors.

Are these assumptions reasonable in applications?

The true model in the model class assumption is often not realistic. The model-data mismatch is due to a combination of

1. wrong model class
2. process noise
3. measurement noise

In control and signal processing, 1 often dominates 2 and 3

- the model class consists of low-order LTI systems but the "true" system is high-order nonlinear time-varying
- the effect of the unobserved inputs is not strong
- the measurement devices are accurate

## Model behavior

The behavior $\mathscr{B}$ of a model $\Sigma$ is the set of all trajectories of $\Sigma$

$$\mathscr{B} = \{ (u, y) \mid (u, y) \text{ is trajectory of } \Sigma \}$$

$\mathscr{B}$ completely specifies $\Sigma$ and allow us to postpone the issue of choosing a model representation to a later stage of the analysis

$$w = (u, y) \in \mathscr{B} \quad \Longleftrightarrow \quad y \text{ is an output of } \Sigma \text{ for an input } u$$

The inputs-outputs partition $(u, y)$ of the trajectory $w$ is not as essential as the classical setting implies.

The behavioral approach was introduced by Jan C. Willems

## Deterministic approximation

Latency identification: given data $w_d$ and model class $\mathscr{M}$

$$\text{minimize (over } \widehat{\mathscr{B}} \in \mathscr{M} \text{ and } e) \ \|e\| \text{ subject to } (e, w_d) \in \widehat{\mathscr{B}}$$

EIV identification: given data $w_d$ and model class $\mathscr{M}$

$$\text{minimize (over } \widehat{\mathscr{B}} \in \mathscr{M} \text{ and } \widehat{w}) \ \|w_d - \widehat{w}\| \text{ subject to } \widehat{w} \in \widehat{\mathscr{B}}$$

*". . . the noise model H in (3.1) is from this point of view just an alibi for determining the predictor. . . . This also means that the difference between a "stochastic system" (3.1) and a "deterministic" one (3.35) is not fundamental."*

*L. Ljung,* System identification: Theory for the user
*Second edition, 1999, Page 74*

## Maximum likelihood estimation

The deterministic approximation problems yield maximum likelihood estimates assuming that

1. there is a true model $\mathscr{B}_{\text{true}}$ and

2. the model-data mismatch is a stochastic process

   - Latency model: there is $e_{\text{true}}$, such that $(w_d, e_{\text{true}}) \in \mathscr{B}_{\text{true}}$
     $e_{\text{true}}$ – realization of zero mean, white, Gaussian process

   - EIV model: there is $w_{\text{true}} \in \mathscr{B}_{\text{true}}$, such that $w_d = w_{\text{true}} + \widetilde{w}$
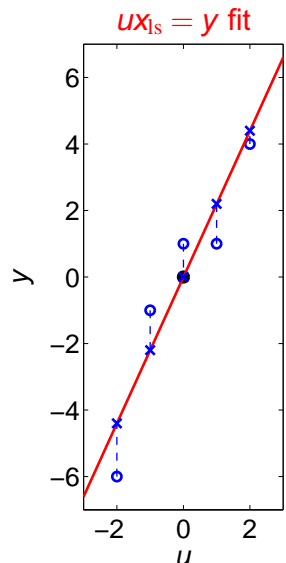     $\widetilde{w}$ – realization of zero mean, white, Gaussian process

## Notes

- choosing representation of the model, makes the identification problems parameter optimization problems

- however, different representations lead to different parameter optimization problems

- latency and EIV are applications of a general principle:
  *impose relevant prior knowledge by the selection of the model class and the data fitting criterion*

- combined with the deterministic point of view, this principle leads to low-rank approximation problems

## Static linear model: an example



$ux_{ls} = y$ fit

**Line fitting problem:** Fit the points

$$w_d^1 = \begin{bmatrix} -2 \\ -6 \end{bmatrix}, \ w_d^2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \ \dots, \ w_d^5 = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

by a line passing through the origin.

**Classic solution:** Define $w_d^i =: \mathrm{col}(u_d^i, y_d^i)$
and solve the least squares problem

$$xu_d^i = y_d^i, \quad \text{for } i = 1, \dots, 5$$

**The model** is the fitting line

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

## Static linear model

A static linear model $\mathscr{B}$ with $q$ variables is a subspace of $\mathbb{R}^q$.

Complexity of $\mathscr{B}$ is defined to be the dimension of $\mathscr{B}$, which is equal to the number of inputs (free variables)

$\mathscr{L}_m$ — set of all static linear models with at most $m$ inputs

Rank constraint on the data matrix

$$w_d^i \in \mathscr{B} \in \mathscr{L}_m, \ i = 1, \dots, N \quad \Longleftrightarrow \quad \mathrm{rank}(\begin{bmatrix} w_d^1 & \cdots & d_d^N \end{bmatrix}) \leq m$$

We will write $D \in \mathscr{B}$

$$D := \begin{bmatrix} w_d^1 & \cdots & d_d^N \end{bmatrix}$$

when each column of $D$ is in $\mathscr{B}$.

Static linear EIV identification:    given $D$ and $m$

   minimize (over $\widehat{\mathscr{B}}$ and $\widehat{D}$) $\|D - \widehat{D}\|$ subject to $\widehat{D} \in \widehat{\mathscr{B}} \in \mathscr{L}_m$

$\Updownarrow$

Low-rank approximation:    given $D$ and $m$

   minimize   (over $\widehat{D}$)   $\|D - \widehat{D}\|_{\mathrm{F}}$   subject to   $\mathrm{rank}(\widehat{D}) \leq m$

- nonconvex, however, an analytic solution exists (SVD)
- also known as the principal component analysis

## Modified low-rank approximation problems

- Weighted low-rank approximation

    minimize   $\sum w_{ij}(D_{ij} - \widehat{D}_{ij})^2$   subject to   $\mathrm{rank}(\widehat{D}) \leq m$

    Allows to treat missing data by setting weights $w_{ij}$ to zero.

- Nonnegative low-rank approximation

    minimize   $\|D - \widehat{D}\|$   subject to   $\mathrm{rank}(\widehat{D}) \leq m$
    
         and   $\widehat{D}_{ij} \geq 0$,   for all $i, j$

- Structured low-rank approximation

    minimize   $\|D - \widehat{D}\|$   subject to   $\mathrm{rank}(\widehat{D}) \leq m$

         and $\widehat{D}$ has the same structure as $D$

    Allows to treat dynamical models $\rightsquigarrow$ using Hankel structure.

## Linear time-invariant (LTI) models

Consider the time series

$$\big(w(1),\ldots,w(T)\big), \qquad w(t) \in \mathbb{R}^q$$

Two special cases:

- $q = 1$ variable, $m = 1$ input — scalar, autonomous model
- $q = 2$ var., $m = 1$ input — single input, single output model

The difference equation

$$r_0 w(t) + r_1 w(t+1) + \cdots + w(t+n) = 0, \quad r_i \in \mathbb{R}^{1 \times q} \qquad \text{(DE)}$$

defines an LTI model with $m = q - 1$ inputs of order at most $n$.

$\mathscr{L}_{m,n}$ — LTI model CLASS $\leq m$ inputs and order $\leq n$

Note that $\mathscr{L}_{m,0} = \mathscr{L}_m$ — the class of static models.

## Linear prediction problem

Future values of $w$ are estimated as linear comb. of past values

$$w(t) = c_1 w(t-1) + c_2 w(t-2) + \cdots + c_n w(t-n) \qquad \text{(LP)}$$

$c_i$ are the linear prediction coefficients

Given an observed signal $w_{\mathrm{d}}$, how do we find the coefficients $c_i$?

There are many methods for doing this:

- Pisarenko, Prony, Kumaresan–Tufts methods
- subspace methods
- frequency domain methods
- maximum likelihood method $\equiv$ Hankel low-rank approx.

## Sum-of-damped-exponentials model

Model the signal $w$ as

$$w(t) = \sum_{i=1}^n a_i e^{d_i t} e^{\mathbf{i}(\omega_i t + \phi_i)} \qquad \text{(SDE)}$$

where $a_i$, $d_i$, $\phi_i$, and $\omega_i$ are parameters of the model

| | | | | |
|---|---|---|---|---|
| $a_i$ | — | amplitudes | $d_i$ | — | dampings |
| $\omega_i$ | — | frequencies | $\phi_i$ | — | initial phases |

For all $\{a_i, d_i, \omega_i, \phi_i\}$ there are $c_i$ and $w(-n+1),\ldots,w(0)$, s.t. the solution of (LP) coincides with (SDE) and vice verse.

the LP problem $\quad\Longleftrightarrow\quad$ modeling by (SDE)

## LTI models and Hankel structured matrices

$$r_0 w(t) + r_1 w(t+1) + \cdots + r_n w(t+n) = 0, \qquad w(t) \in \mathbb{R}^{1 \times q}$$

for $t = 1,\ldots,T-n$, is equivalent to the system of equations

$$\begin{bmatrix} r_0 & \cdots & r_n \end{bmatrix} \underbrace{\begin{bmatrix} w(1) & w(2) & w(3) & \cdots & w(T-n) \\ w(2) & w(3) & w(4) & \cdots & w(T-n+1) \\ w(3) & w(4) & w(5) & \cdots & w(T-n+1) \\ \vdots & \vdots & \vdots & & \vdots \\ w(n+1) & w(n+2) & w(n+3) & \cdots & w(T) \end{bmatrix}}_{\mathscr{H}_{n+1}(w)} = 0$$

$$\Longleftrightarrow \quad \mathrm{rank}\,\big(\mathscr{H}_{n+1}(w)\big) \leq m(n+1) + n, \quad m := q - \mathrm{row\,dim}(r)$$

- the subspace methods are based on the SVD of $\mathscr{H}_{n+1}(w)$ (unstructured low-rank approximation)
- the maximum-likelihood method preserves the structure

## EIV identification

$$w \in \mathscr{B} \in \mathscr{L}_{m,n} \quad \Longleftrightarrow \quad \text{rank}\left(\mathscr{H}_{n+1}(w)\right) \leq m(n+1)+n$$

EIV identification: given data $w_\mathrm{d}$, # of inputs $m$, and order $n$

minimize (over $\widehat{\mathscr{B}} \in \mathscr{M}$ and $\widehat{w}$) $\|w_\mathrm{d} - \widehat{w}\|$ subject to $\widehat{w} \in \widehat{\mathscr{B}}$

$\Updownarrow$

Hankel structured low-rank approximation: given $w_\mathrm{d}$ and $k$

minimize   over $\widehat{w}$   $\|w_\mathrm{d} - \widehat{w}\|$   subject to   $\text{rank}\left(\mathscr{H}_{n+1}(\widehat{w})\right) \leq k$

## Structured low-rank approximation

No closed form solution is known for the general SLRA problem

$$\widehat{w}^* := \arg\min_{\widehat{w}} \|w_\mathrm{d} - \widehat{w}\| \quad \text{subject to} \quad \text{rank}\left(\mathscr{S}(\widehat{w})\right) \leq n$$

NP-hard, consider solution methods based on local optimization

Representing the constraint in a kernel form, the problem is

$$\min_{r, rr^\top = 1} \left( \min_{\widehat{w}} \|w_\mathrm{d} - \widehat{w}\| \quad \text{subject to} \quad r\mathscr{S}(\widehat{w}) = 0 \right)$$

Double minimization with bilinear equality constraint.

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

## Variable projection vs. alternating projections

Two ways to approach the double minimization:

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

$$\min_{r, rr^\top = 1} r\mathscr{S}(w_\mathrm{d})\left(G^\top(r)G(r)\right)^{-1}\mathscr{S}^\top(w_\mathrm{d})r^\top$$

  $\rightsquigarrow$ 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.

## Algorithmic details using the VARPRO approach

The structured low-rank approximation problem is equivalent to

$$\min_{r, rr^\top = 1} r\mathscr{S}(w_\mathrm{d})\left(G^\top(r)G(r)\right)^{-1}\mathscr{S}^\top(w_\mathrm{d})r^\top$$

To evaluate the cost function we need to solve for $z$

$$\left(G^\top(r)G(r)\right)z = \left(r\mathscr{S}(w_\mathrm{d})\right)$$

What special structure does $G^\top G$ have?

Banded Toeplitz for any $\mathscr{S} = \begin{bmatrix} \mathscr{S}_1 & \cdots & \mathscr{S}_q \end{bmatrix}$, where $\mathscr{S}_i$ is Toeplitz, Hankel, Toeplitz+Hankel, unstructured, or fixed.

## Special case: sum-of-damped-exp. modeling

In the sum-of-damped-exp. modeling, the structure is

$$\mathscr{S}(w) = \mathscr{H}_{n+1}(w)$$

What matrix $G$ satisfies

$$r\mathscr{H}_{n+1}(w) = 0 \quad \Longleftrightarrow \quad wG(r) = 0$$

for all $r$ and $w$? What is the structure of $G^\top G$?

## Special case: sum-of-damped-exp. modeling

$$\begin{bmatrix} r_0 & r_1 & \cdots & r_n \end{bmatrix} \underbrace{\begin{bmatrix} w(1) & w(2) & \cdots & w(T-n) \\ w(2) & w(3) & \cdots & w(T-n+1) \\ \vdots & \vdots & & \vdots \\ w(n+1) & w(n+2) & \cdots & w(T) \end{bmatrix}}_{\mathscr{H}_{n+1}(w)}$$

$$= \begin{bmatrix} w_1 & w_2 & \cdots & w_T \end{bmatrix} \underbrace{\begin{bmatrix} r_0 & & & \\ r_1 & r_0 & & \\ \vdots & r_1 & \ddots & \\ r_n & \vdots & \ddots & r_0 \\ & r_n & & r_1 \\ & & \ddots & \vdots \\ & & & r_n \end{bmatrix}}_{G(r)}$$

## Special case: sum-of-damped-exp. modeling

Therefore,

$$G^\top G = \begin{bmatrix} r_0 & r_1 & \cdots & r_n & & & \\ & r_0 & r_1 & \cdots & r_n & & \\ & & \ddots & \ddots & & \ddots & \\ & & & r_0 & r_1 & \cdots & r_n \end{bmatrix} \begin{bmatrix} r_0 & & & \\ r_1 & r_0 & & \\ \vdots & r_1 & \ddots & \\ r_n & \vdots & \ddots & r_0 \\ & r_n & & r_1 \\ & & \ddots & \vdots \\ & & & r_n \end{bmatrix}$$

(All missing elements are zeros.)

## Special case: sum-of-damped-exp. modeling

$$G^\top G = \begin{bmatrix} \sum_{i=0}^{n} r_i r_i & \sum_{i=1}^{n} r_i r_{i-1} & \cdots & r_n r_0 & & & \\ \sum_{i=1}^{n} r_{i-1} r_i & \ddots & & & \ddots & & \\ \vdots & & \ddots & & & \ddots & \\ r_0 r_n & & & \ddots & & & r_n r_0 \\ & \ddots & & & & & \vdots \\ & & \ddots & & & \ddots & \sum_{i=1}^{n} r_i r_{i-1} \\ & & & r_0 r_n & \cdots & \sum_{i=1}^{n} r_{i-1} r_i & \sum_{i=0}^{n} r_i r_i \end{bmatrix}$$

banded Toeplitz, bandwidth $2n+1$

## Data fitting by a second order model

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

Consider first exact data:

$$w \in \mathscr{B}(A,b,c) \iff w^\top A w + b^\top w + c = 0$$
$$\iff \big\langle \underbrace{\text{col}(w \otimes_s w, w, 1)}_{w_{\text{ext}}}, \underbrace{\text{col}\big(\text{vec}_s(A), b, c\big)}_{\theta} \big\rangle = 0$$

$$\{w_1,\dots,w_N\} \in \mathscr{B}(\theta) \iff \theta \in \text{left ker} \underbrace{\big[ w_{\text{ext},1} \quad \cdots \quad w_{\text{ext},N} \big]}_{D_{\text{ext}}}, \quad \theta \neq 0$$

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

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)

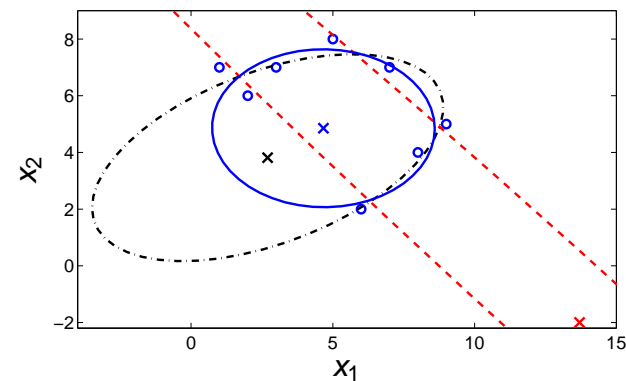○ — data points    × — centers

---

## Conclusions

- User choices should impose prior knowledge for the true systems rather than convenient theoretical assumption.

- Stochastic estimation and deterministic approximation are two sides of the same coin.

- The behavioral setting leads to elegant and useful data modeling philosophy.

- Its algorithmic implementation is low-rank approximation.

- EIV static linear model identification — unstructured LRA.
  EIV dynamic LTI model identification — Hankel LRA.

- Algorithms based on VARPRO and alternating projections.

- Nonlinear model identification via data transformation.

---

## Thank you