

Approximate system identification

Ivan Markovsky

University of Southampton

Outline

- From exact to approximate identification
- Misfit vs latency
- Misfit minimization
- Misfit computation
 - kernel
 - image
 - input/state/output

From exact to approximate identification

Exact system identification and the MPUM

Exact identification problem:

Given a vector time series

$$w_d = (w_d(1), \dots, w_d(T)) \in (\mathbb{R}^w)^T$$

find the smallest $m \in \mathbb{N}$ and $\ell \in \mathbb{N}$ and LTI system $\mathcal{B}_{\text{mpum}} \in \mathcal{L}_{m,\ell}^w$, s.t.

$$w_d \in \mathcal{B}_{\text{mpum}}.$$

The model $\mathcal{B}_{\text{mpum}} = \mathcal{B}_{\text{mpum}}(w_d)$ is unique and is called the MPUM of w_d in the model class $\mathcal{L}_{m,\ell}^w$.

There are effective algorithms for computing representations of $\mathcal{B}_{\text{mpum}}(w_d)$ from given $w_d = (u_d, y_d)$ and an upper bound ℓ_{\max} on ℓ .

Identifiability

Provided that $w_d \in \overline{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w$, find conditions under which

$$\mathcal{B}_{\text{mpum}}(w_d) = \overline{\mathcal{B}}.$$

Main theoretical result in the exact identification setting:

$\overline{\mathcal{B}}$ is identifiable from $w_d = (u_d, y_d)$ in $\mathcal{L}_{m, \ell_{\max}}^w$ if

1. w_d is exact, i.e., $w_d \in \overline{\mathcal{B}}$,
2. the model class is correct, i.e., $\overline{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w$,
3. u_d is persistently exciting of order $\ell_{\max} + n$, and
4. $\overline{\mathcal{B}}$ is controllable

Notes

- the conditions are only sufficient
- conditions 1, 2, and 4 are not verifiable from the given data and should therefore be postulated
- a known input/output partitioning of the data is assumed
- from a practical point of view, condition 1 is a strong assumption that limits the applicability of the exact SYSID methods

approaches for relaxing condition 1 are described next

MPUM in the case of “noisy” data

Q: What is the MPUM of a noisy trajectory

$$w_d = \bar{w} + \tilde{w} \quad \text{where} \quad \bar{w} \in \overline{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w \quad \text{and} \quad \tilde{w} \text{ is random and zero mean?} \quad (\text{EIV})$$

A: With probability 1,

$$\mathcal{B}_{\text{mpum}}(w_d) = (\mathbb{R}^w)^{\mathbb{Z}_+} \quad (\text{all variables are inputs})$$

This is a trivial model because it fits every trajectory.

Alternatively, $\mathcal{B}_{\text{mpum}}(w_d)$ **does not exist** in a model class $\mathcal{M} = \mathcal{L}_{m, \ell_{\max}}^w$ of bounded ($m < w$, $\ell_{\max} \ll T$) complexity.

In what follows we assume a given bounded complexity model class $\mathcal{M} = \mathcal{L}_{m, \ell_{\max}}^w$, so we refer to lack of existence rather than trivial MPUM.

In practice, w_d is often generated by a

nonlinear, **infinite** dimensional, **time-varying** system $\overline{\mathcal{B}}$

possibly with process and measurement noises, i.e., $\overline{\mathcal{B}} \notin \mathcal{L}_{m, \ell_{\max}}^w$

\implies even without noise, identifying exact model is often not possible

It can be argued that in practice the approximation aspect ($\hat{\mathcal{B}} \approx \overline{\mathcal{B}}$) is often more important than the stochastic estimation ($\hat{\mathcal{B}} \rightarrow \overline{\mathcal{B}}$ as $T \rightarrow \infty$)

An approximate $\hat{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w$ is what is anyway needed:

Many prediction and control methods are based on LTI models

\implies even if it was possible to identify $\overline{\mathcal{B}}$, it would be necessary to approximate it by $\hat{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w$

Misfit vs latency

Modifications of the MPUM concept

Unless the model class $\mathcal{M} = \mathcal{L}_{m, \ell_{\max}}^w$ is enlarged, *i.e.*, (m, ℓ_{\max}) is increased, until the MPUM exists in \mathcal{M}

we have to accept falsified models in \mathcal{M}
 \rightsquigarrow approximate SYSID

The main question in approximate SYSID is:

Q: Which (falsified) model in \mathcal{M} to choose?

A: In some sense the “least falsified” (“approximately unfalsified”) one.

Two major notions of “least falsified” are small **misfit** and small **latency**.

They quantify the discrepancy between the model and the data.

Misfit approach for approximate SYSID

Consider given data w_d and model class $\mathcal{M} = \mathcal{L}_{m, \ell_{\max}}^w$.

If the MPUM does not exist in \mathcal{M} , i.e.,

$$\mathcal{B}_{\text{mpum}}(w_d) \notin \mathcal{M}$$

we aim to find an approximate model $\hat{\mathcal{B}}$ for w_d in \mathcal{M} .

The misfit approach modifies w_d as little as possible, so that the modified data, say \hat{w} , has MPUM in \mathcal{M} , i.e.,

$$\mathcal{B}_{\text{mpum}}(\hat{w}) \in \mathcal{M}$$

The approximate model for w_d in \mathcal{M} is defined as $\hat{\mathcal{B}}_{\text{misfit}} := \mathcal{B}_{\text{mpum}}(\hat{w})$.

The modification of the data is measured by the **misfit** $\|w_d - \hat{w}\|$

Latency

The latency approach augments w_d by, as small as possible, variable e , so that the augmented data $w_{\text{ext}} := \text{col}(e, w_d)$ has MPUM in the augmented model class, *i.e.*,

$$\mathcal{B}_{\text{mpum}}(\text{col}(e, w_d)) \in \mathcal{M}_{\text{ext}} := \mathcal{L}_{m+e, \ell_{\max}}^{w+e}$$

Let Π_w be the projection of $\text{col}(e, w)$ on w . The approximate model for w_d in \mathcal{M} is defined as

$$\hat{\mathcal{B}}_{\text{latency}} := \Pi_w \mathcal{B}_{\text{mpum}}(\text{col}(e, w_d))$$

The size of e is measured by the **latency** $\|e\|$

Notes

- Both the misfit and latency approaches reduce the approximate SYSID problem to (different) exact SYSID problems:
 - $\hat{\mathcal{B}}_{\text{misfit}}$ is exact for the modified data \hat{w}
 - $\hat{\mathcal{B}}_{\text{latency}}$ is obtained from an exact model for $\text{col}(e, w_d)$
- If $\mathcal{B}_{\text{mpum}}(w_d) \in \mathcal{M}$ (the data is exact),

$$\hat{\mathcal{B}}_{\text{misfit}} = \hat{\mathcal{B}}_{\text{latency}} = \mathcal{B}_{\text{mpum}}(w_d) \quad (\hat{w} = w_d \text{ and } e = 0)$$

So, misfit and latency are indeed extensions of the MPUM.

- The misfit approach modifies w_d but does not change \mathcal{M} .
- The latency approach modifies \mathcal{M} but does not change w_d .

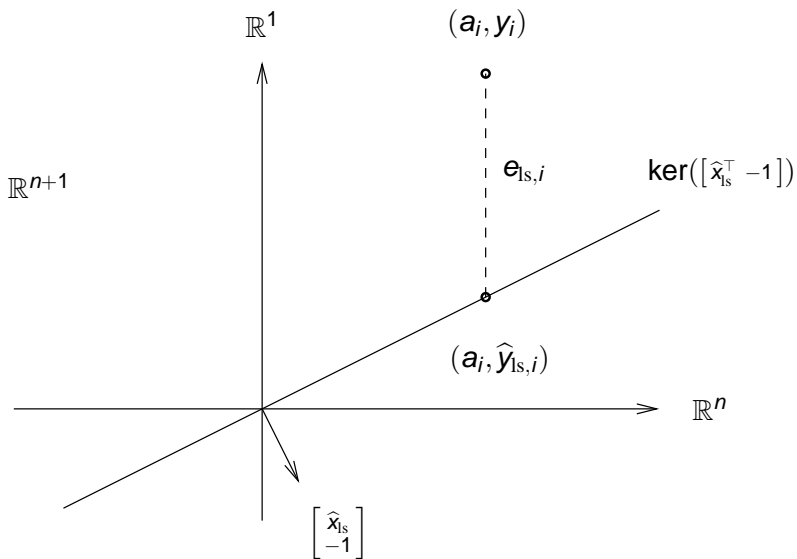
Static case: latency \leftrightarrow LS misfit \leftrightarrow TLS

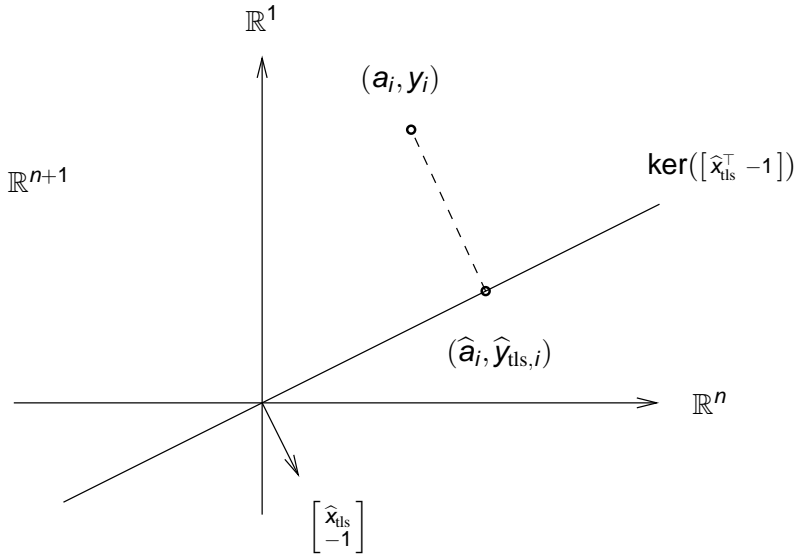
LS: minimize $_{e,x} \|e\|_2$ subject to $Ax = y + e$ e latent variable

$$\text{latency}((A, y), x) := \left(\min_e \|e\|_2^2 \text{ s.t. } Ax = y + e \right) = \|Ax - y\|_2^2$$

TLS: minimize $_{\Delta A, \Delta y, x} \left\| \begin{bmatrix} \Delta A & \Delta y \end{bmatrix} \right\|_F$ $\Delta A, \Delta y$ data corrections
subject to $(A + \Delta A)x = y + \Delta y$

$$\begin{aligned} \text{misfit}((A, b), x) &:= \min_{\Delta A, \Delta b} \left\| \begin{bmatrix} \Delta A & \Delta b \end{bmatrix} \right\|_F \text{ s.t. } (A + \Delta A)x = b + \Delta b \\ &= \frac{\|Ax - b\|_2}{\sqrt{1 + \|x\|_2^2}} \end{aligned}$$



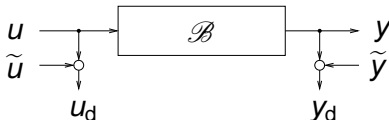


Statistical interpretation of misfit and latency

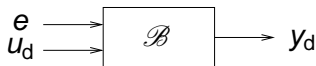
misfit \leftrightarrow errors-in-variables (EIV) model

latency \leftrightarrow ARMAX model

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



ARMAX model: e — process noise



Assumptions: \tilde{w} , e — zero mean, stationary, white, ergodic, Gaussian

- $\Pi_w \mathcal{B}_{\text{ext}}$ — deterministic part of the model
- $\Pi_e \mathcal{B}_{\text{ext}}$ — stochastic part of the model

Notes:

- The Kalman filter and LQG control are based on the latency model
- \implies stochastic part is used by the KF and LQG controller

Misfit minimization

Identification problems

Misfit minimization (GTLS): given $w_d \in (\mathbb{R}^w)^T$ and $\ell_{\max} \in \mathbb{N}$, find

$$\hat{\mathcal{B}}_{\text{gtls}}^* := \arg \min_{\hat{\mathcal{B}}, \hat{w}} \|w_d - \hat{w}\| \quad \text{subject to} \quad \hat{w} \in \hat{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}$$

Latency minimization (PEM): given $w_d \in (\mathbb{R}^w)^T$ and $\ell_{\max} \in \mathbb{N}$, find

$$\hat{\mathcal{B}}_{\text{pem}}^* := \arg \min_{\hat{\mathcal{B}}_{\text{ext}}, e} \|e\| \quad \text{subject to} \quad (e, \hat{w}) \in \hat{\mathcal{B}}_{\text{ext}} \in \mathcal{L}_{m+e, \ell_{\max}}$$

Notes:

- nonconvex optimization problems
- solution methods based on local optimization
- initial approximation obtained from subspace methods

Misfit minimization

Define the misfit between w_d and \mathcal{B} as follows

$$\text{misfit}(w_d, \mathcal{B}) := \underset{\hat{w}}{\text{minimize}} \quad \|w_d - \hat{w}\|_{\ell_2} \quad \text{subject to} \quad \hat{w} \in \mathcal{B}$$

the minimizer \hat{w}^* is projection of w_d on \mathcal{B} (best ℓ_2 **approx. of w_d in \mathcal{B}**)

alternatively, \hat{w}^* is the **smoothed estimate of w_d , given \mathcal{B}**

our goal is to find the model $\hat{\mathcal{B}}$ that minimizes $\text{misfit}(w_d, \mathcal{B})$, i.e.,

$$\hat{\mathcal{B}} := \underset{\hat{\mathcal{B}}}{\arg \min} \quad \text{misfit}(w_d, \mathcal{B}) \quad \text{subject to} \quad \mathcal{B} \in \mathcal{M}$$

a double minimization problem: inner minimization is projection on a subspace (easy), outer minimization is a nonconvex problem (difficult)

Maximum likelihood estimator in the EIV setup

Assuming that the data is generated according to the model

$$w_d = \bar{w} + \tilde{w}, \quad \text{where } \bar{w} \in \bar{\mathcal{B}} \in \mathcal{M} \quad \text{and} \quad \tilde{w} \sim \mathcal{N}(0, s^2 I)$$

$\hat{\mathcal{B}}$ is the **maximum likelihood estimator** of the true model $\bar{\mathcal{B}}$

$\hat{\mathcal{B}}$ is a **consistent** estimator of the true model $\bar{\mathcal{B}}$, i.e., $\hat{\mathcal{B}} \rightarrow \bar{\mathcal{B}}$ as $T \rightarrow \infty$

The log-likelihood function is (“const” does not depend on \hat{w} and $\hat{\mathcal{B}}$)

$$\ell(\hat{\mathcal{B}}, \hat{w}) = \begin{cases} \text{const} - \frac{1}{2s^2} \|w_d - \hat{w}\|_{\ell_2}^2, & \text{if } \hat{w} \in \hat{\mathcal{B}} \in \mathcal{M} \\ -\infty, & \text{otherwise,} \end{cases}$$

likelihood evaluation



misfit computation

Misfit computation

Computation of the misfit

Given w_d and $\mathcal{B} \in \mathcal{L}_{m, \ell_{\max}}^w$, find $\text{misfit}(w_d, \mathcal{B}) := \min_{\hat{w} \in \mathcal{B}} \|w_d - \hat{w}\|_{\ell_2}$

\mathcal{B} is subspace \implies the constraint is linear
 \implies ordinary LS problem

Using general purpose LS solvers, the comput. complexity is $O(T^3)$.

Time-invariance of \mathcal{B} , however, implies Toeplitz structure of the LS prob.

Structure exploiting misfit computation methods have **complexity $O(T)$** .

They are based on:

1. structured matrix computations
(e.g., using displacement rank theory and the gen. Schur alg.)
2. Riccati recursions (Kalman smoother)

Misfit computation using image representation

$$\text{minimize}_{\hat{w}} \quad \|w_d - \hat{w}\|_{\ell_2} \quad \text{subject to} \quad \hat{w} \in \mathcal{B} := \text{image}(M(\sigma)) \quad (\text{M})$$

Recall from Lecture 5 that

$$w \in \mathcal{B} \iff w = \underbrace{\begin{bmatrix} M_0 & M_1 & \cdots & M_\ell \\ & M_0 & M_1 & \cdots & M_\ell \\ & & \ddots & \ddots & \\ & & & M_0 & M_1 & \cdots & M_\ell \end{bmatrix}}_{\mathcal{I}_M} \begin{bmatrix} v(1) \\ v(2) \\ \vdots \\ v(T+\ell) \end{bmatrix}$$

So that (M) is an ordinary least squares problem

$$\text{minimize}_v \quad \|w_d - \mathcal{I}_M v\| \quad (\text{M}')$$

and

$$\text{misfit}(w_d, \text{image}(M(\sigma))) = w_d^\top \mathcal{I}_M (\mathcal{I}_M^\top \mathcal{I}_M)^{-1} \mathcal{I}_M^\top w_d$$

Misfit computation using kernel representation

$$\text{minimize}_{\hat{w}} \quad \|w_d - \hat{w}\|_{\ell_2} \quad \text{subject to} \quad \hat{w} \in \mathcal{B} := \ker(R(\sigma)) \quad (\text{R})$$

Recall from Lecture 5 that

$$w_d \in \mathcal{B} \iff \underbrace{\begin{bmatrix} R_0 & R_1 & \cdots & R_\ell \\ & R_0 & R_1 & \cdots & R_\ell \\ & & \ddots & \ddots & \\ & & & R_0 & R_1 & \cdots & R_\ell \end{bmatrix}}_{\mathcal{T}_T(R)} \begin{bmatrix} w_d(1) \\ w_d(2) \\ \vdots \\ w_d(T) \end{bmatrix} = 0$$

So that (R) is an equality constrained least squares problem

$$\text{minimize}_{\hat{w}} \quad \|w_d - \hat{w}\|_{\ell_2} \quad \text{subject to} \quad \mathcal{T}_R \hat{w} = 0 \quad (\text{R}')$$

In order to solve (R') explicitly, we need a basis for $\text{null}(\mathcal{T}_R)$. Let

N be such that $\mathcal{T}_R N = 0$ and N is full column rank

Then, $\mathcal{T}_R \hat{w} = 0 \iff \exists z \in \mathbb{R}^{\text{col dim}(N)}$ s.t. $\hat{w} = Nz$, and

$$\text{misfit}(w_d, \ker(R(\sigma))) = w_d^\top N (N^\top N)^{-1} N^\top w_d.$$

Note that the columns of \mathcal{T}_M form a particular basis for the null space of \mathcal{T}_R . This can be seen algebraically from

$$\mathcal{T}_R \mathcal{T}_M = 0 \quad \text{and} \quad \mathcal{T}_M \text{ is full column rank}$$

or (better) from a system theoretic point of view:

$$\text{colspan}(\mathcal{T}_R) = \text{null}(\mathcal{T}_R) = \mathcal{B}_T.$$

Efficient reduction of (R) to (M)

Computing N , reduces (R) to (M), however, N need not be Toeplitz.

Moreover, computing N by general purpose methods is expensive, while the path $R \mapsto M \mapsto \mathcal{T}_M$ is cheap.

Let $R(z) =: \begin{bmatrix} Q(z) & P(z) \end{bmatrix}$ with $P(z) \in \mathbb{R}^{p \times p}[z]$ nonsingular (this is equivalent to assuming existence of I/O partition $w = \text{col}(u, y)$)

Compute the right matrix fraction of $G(z) := P^{-1}(z)Q(z)$

$$G(z) = Q_l(z)P_l^{-1}(z)$$

Then

$$M(z) = \begin{bmatrix} P_l(z) \\ Q_l(z) \end{bmatrix}$$

which also reduces (R) to (M).

Misfit computation using I/S/O representation

$$\text{minimize}_{\hat{w}} \quad \|w_d - \hat{w}\|_{\ell_2} \quad \text{subject to} \quad \hat{w} \in \mathcal{B} := \mathcal{B}_{\text{i/s/o}}(A, B, C, D) \quad (\text{SS})$$

Recall from Lecture 5 that

$$w = (u, y) \in \mathcal{B}_{\text{i/s/o}}(A, B, C, D) \quad \Longleftrightarrow \quad \text{there exists } x_{\text{ini}} \in \mathbb{R}^n, \text{ such that}$$

$$y = \underbrace{\begin{bmatrix} C \\ CA \\ CA^2 \\ \vdots \\ CA^{T-1} \end{bmatrix}}_{\mathcal{O}} x_{\text{ini}} + \underbrace{\begin{bmatrix} H(0) & & & \\ H(1) & H(0) & & \\ H(2) & H(1) & H(0) & \\ \vdots & \ddots & \ddots & \ddots \\ H(T-1) & \dots & H(2) & H(1) & H(0) \end{bmatrix}}_{\mathcal{I}_H} u$$

where $H(0) = D$ and $H(t) = CA^{t-1}B$, for $t = 1, 2, \dots$

Then (SS) is equivalent to the ordinary LS problem:

$$\text{minimize}_{x_{\text{ini}}, \hat{u}} \quad \left\| \begin{bmatrix} u_d \\ y_d \end{bmatrix} - \begin{bmatrix} I & 0 \\ \mathcal{O} & \mathcal{T}_H \end{bmatrix} \begin{bmatrix} x_{\text{ini}} \\ \hat{u} \end{bmatrix} \right\| \quad (\text{SS}')$$

Efficient solution via Riccati recursion \rightsquigarrow Kalman smoother

References

1. B. Roorda and C. Heij,
Global total least squares modeling of multivariate time series,
IEEE-AC, 40(1):50–63, 1995
2. I. Markovsky et al.,
Application of structured total least squares for system
identification and model reduction,
IEEE-AC, 50(10):1490–1500, 2005
3. P. Lemmerling and B. De Moor,
Misfit versus latency,
Automatica, 37:2057–2067, 2001.
4. I. Markovsky, J. C. Willems, S. Van Huffel, and B. De Moor.
Exact and Approximate Modeling of Linear Systems
SIAM, 2006

Software

A Matlab toolbox:

`ftp.esat.kuleuven.be/pub/SISTA/markovsky/
abstracts/04-221a.html`

Exercises 3 applies a simple GTLS algorithm on benchmark problem and compares the results with the ones of the latency minimization approach, implemented in the System Identification Toolbox.