

Exact and approximate linear system identification

Ivan Markovsky

University of Southampton

- Exact identification: the most powerful unfalsified model
- Identifiability conditions and algorithms
- From exact to approximate identification: misfit vs latency
- Misfit computation and minimization

Exact identification: $w_d \mapsto \mathcal{B}$, $w_d \in \mathcal{B}$

An exact identification problem

Problem P1 (Exact identification)

Given two vector time series

$$u_d = (u_d(1), \dots, u_d(T)) \in (\mathbb{R}^m)^T \quad \text{“inputs”}$$

$$y_d = (y_d(1), \dots, y_d(T)) \in (\mathbb{R}^p)^T \quad \text{“outputs”}$$

find $n \in \mathbb{N}$ and LTI system \mathcal{B} of order n , with m inputs and p outputs, s.t.

$$w_d := (u_d, y_d) \in \mathcal{B},$$

i.e., w_d is a trajectory of \mathcal{B} .

How can we check that “ $w_d \in \mathcal{B}$ ”?

Checking that $w_d \in \mathcal{B}$

Let \mathcal{B} be defined by a minimal input/state/output representation

$$\mathcal{B} = \mathcal{B}_{i/s/o}(A, B, C, D) := \{ (u, y) \mid \sigma x = Ax + Bu, y = Cx + Du \}$$

$(u_d, y_d) \in \mathcal{B}_{i/s/o}(A, B, C, D) \iff$ **there exists $x_{ini} \in \mathbb{R}^n$** , such that

$$y_d = \underbrace{\begin{bmatrix} C \\ CA \\ CA^2 \\ \vdots \\ CA^{T-1} \end{bmatrix}}_{\theta_T(A, C)} x_{ini} + \begin{bmatrix} D & & & \\ CB & D & & \\ CAB & CB & D & \\ \vdots & \ddots & \ddots & \ddots \\ CA^{T-1}B & \dots & CAB & CB & D \end{bmatrix} u_d$$

(y_d is the response of \mathcal{B} under input u_d and initial condition x_{ini})

Revised exact identification problem

Problem P1' (Exact identification)

Given two vector time series

$$\begin{aligned} u_d &= (u_d(1), \dots, u_d(T)) \in (\mathbb{R}^m)^T && \text{"inputs"} \\ y_d &= (y_d(1), \dots, y_d(T)) \in (\mathbb{R}^p)^T && \text{"outputs"} \end{aligned}$$

find the **smallest $n \in \mathbb{N}$** and LTI system \mathcal{B} of order n , with m inputs and p outputs, such that

$$w_d = (u_d, y_d) \in \mathcal{B}.$$

Comments

- P1 is an **exact fitting problem**, a most basic SYSID problem
- easily generalizable to a **set of N time series**
 $u_{d,1}, \dots, u_{d,N} \in (\mathbb{R}^m)^T$ and $y_{d,1}, \dots, y_{d,N} \in (\mathbb{R}^p)^T$
- the **realization problem**

impulse response $\mapsto (A, B, C, D)$

is a special case of P1 for a set of m time series

- while m is given, **finding n is part of the problem**
any observable system of order $n \geq pT$ is a (trivial) solution
- we are actually interested is a **solution of a minimal order**

Set of LTI systems with a bounded complexity

Notation: $\mathcal{L}_{m,\ell}^{w,n}$ is the set of all LTI systems with

- w (external) variables
- **at most m** inputs
- minimal state dimension **at most n** and
- lag (= observability index) **at most ℓ**

For $t \geq n$, **the set $\mathcal{B}|_t$ of all t samples long traj. of \mathcal{B}** has dimension

$$\dim(\mathcal{B}|_t) \leq tm + n \leq tm + p\ell$$

(where $p(\ell - 1) \leq n \leq p\ell$)

$\implies (m, n)$ and (m, ℓ) specify the **complexity** of the model class $\mathcal{L}_{m,\ell}^{w,n}$

Another exact identification problem

Problem P2 (Exact identification)

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 $\mathcal{B} \in \mathcal{L}_{m,\ell}^w$, s.t. $w_d \in \mathcal{B}$.

Comments:

- no separation between inputs and outputs
- the complexity is defined by (m, ℓ)

Identifiability

Most powerful unfalsified model

The most powerful unfalsified model in the model class $\mathcal{L}_{m,\ell}^w$ of a time series $w_d \in (\mathbb{R}^w)^T$ is the system $\mathcal{B}_{\text{mpum}}$ that is

1. in the model class, i.e., $\mathcal{B}_{\text{mpum}} \in \mathcal{L}_{m,\ell}^w$,
2. unfalsified, i.e., $w_d \in \mathcal{B}_{\text{mpum}}|_T$, and
3. most powerful among all LTI unfalsified systems, i.e.,

$$\mathcal{B}' \in \mathcal{L}_{m,\ell}^w \text{ and } w_d \in \mathcal{B}'|_T \implies \mathcal{B}_{\text{mpum}}|_T \subseteq \mathcal{B}'|_T.$$

MPUM may not exist, but if it does, then it is unique

Identifiability question

P2 is the problem of computing the MPUM of w_d in \mathcal{L}^w

The following related question is of interest:

Suppose that

$$w_d \in \overline{\mathcal{B}} \in \mathcal{L}^w$$

and upper bounds n_{\max} , ℓ_{\max} of the order n and lag ℓ of $\overline{\mathcal{B}}$ are given.

Under what conditions $\mathcal{B}_{\text{mpum}}(w_d)$ is equal to the system \mathcal{B} ?

the answer is given by the following lemma

Fundamental Lemma

Let $\overline{\mathcal{B}} \in \mathcal{L}_m^{w,n}$ be controllable and let $w_d := (u_d, y_d) \in \overline{\mathcal{B}}|_T$.

Then, if u_d is persistently exciting of order $L+n$,

$$\text{image} \left(\begin{bmatrix} w_d(1) & w_d(2) & w_d(3) & \cdots & w_d(T-L+1) \\ w_d(2) & w_d(3) & w_d(4) & \cdots & w_d(T-L+2) \\ w_d(3) & w_d(4) & w_d(5) & \cdots & w_d(T-L+3) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_d(L) & w_d(L+1) & w_d(L+2) & \cdots & w_d(T) \end{bmatrix} \right) = \overline{\mathcal{B}}|_L$$

\Rightarrow under the conditions of the FL, any L samples long response y of \mathcal{B} can be obtained as $y = \mathcal{H}_L(y_d)g$, for certain $g \rightsquigarrow$ **algorithms**

\Rightarrow with $L = \ell_{\max} + 1$, the FL gives **conditions for identifiability**

Persistency of excitation

$u_d = (u_d(1), \dots, u_d(T))$ is **persistently exciting of order L** if

$$\mathcal{H}_L(u_d) := \begin{bmatrix} u_d(1) & u_d(2) & u_d(3) & \cdots & u_d(T-L+1) \\ u_d(2) & u_d(3) & u_d(4) & \cdots & u_d(T-L+2) \\ u_d(3) & u_d(4) & u_d(5) & \cdots & u_d(T-L+3) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ u_d(L) & u_d(L+1) & u_d(L+2) & \cdots & u_d(T) \end{bmatrix} \quad \text{is full row rank}$$

System theoretic interpretation:

u_d is persistently exciting of order L

\iff

there is no LTI system with # of inputs $< m$ and lag $< L$ for which u_d is a trajectory

Overview of algorithms

1. $w_d \mapsto R(\xi)$
2. $w_d \mapsto$ **impulse response H**
3. $w_d \mapsto (A, B, C, D)$ (possibly **balanced**)
 - 3.1 $w_d \mapsto R(\xi) \mapsto (A, B, C, D)$ or $w_d \mapsto H \mapsto (A, B, C, D)$
 - 3.2 $w_d \mapsto \mathcal{O}_{\ell_{\max}+1}(A, C) \mapsto (A, B, C, D)$
 - 3.3 $w_d \mapsto (x_d(1), \dots, x_d(\ell_{\max} + m + 1)) \mapsto (A, B, C, D)$

MATLAB toolbox:

<ftp.esat.kuleuven.be/pub/SISTA/markovsky/abstracts/05-122.html>

References

1. J. C. Willems.
From time series to linear system—Part II. Exact modelling.
Automatica, 22(6):675–694, 1986.
2. J. C. Willems, P. Rapisarda, I. Markovsky, and B. De Moor.
A note on persistency of excitation.
Systems & Control Letters, 54(4):325–329, 2005.
3. I. Markovsky, J. C. Willems, P. Rapisarda, and B. De Moor.
Algorithms for deterministic balanced subspace identification.
Automatica, 41(5):755–766, 2005.
4. I. Markovsky, J. C. Willems, S. Van Huffel, and B. De Moor.
Exact and Approximate Modeling of Linear Systems
SIAM, 2006

From exact to approximate identification

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. $\overline{\mathcal{B}}$ is controllable, and
4. u_d is persistently exciting of order $\ell_{\max} + n$

Notes

- the conditions are only sufficient
- conditions 1, 2, and 3 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, conditions 1 and 2 are strong and limit 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

$$w_d = \overline{w} + \tilde{w} \quad \text{where} \quad \overline{w} \in \overline{\mathcal{B}} \in \mathcal{L}_{m, \ell_{\max}}^w \quad \text{(EIV)}$$

and \tilde{w} is random and zero mean?

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.*, $\overline{\mathcal{B}} \notin \mathcal{L}_{m, \ell_{\max}}^w$

\Rightarrow 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

\Rightarrow 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 .

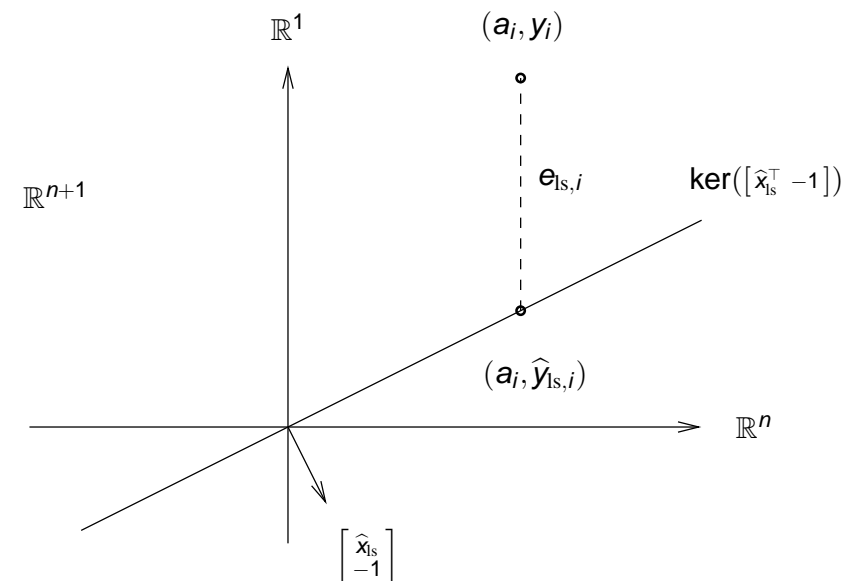
Static case: **latency** \leftrightarrow LS **misfit** \leftrightarrow TLS

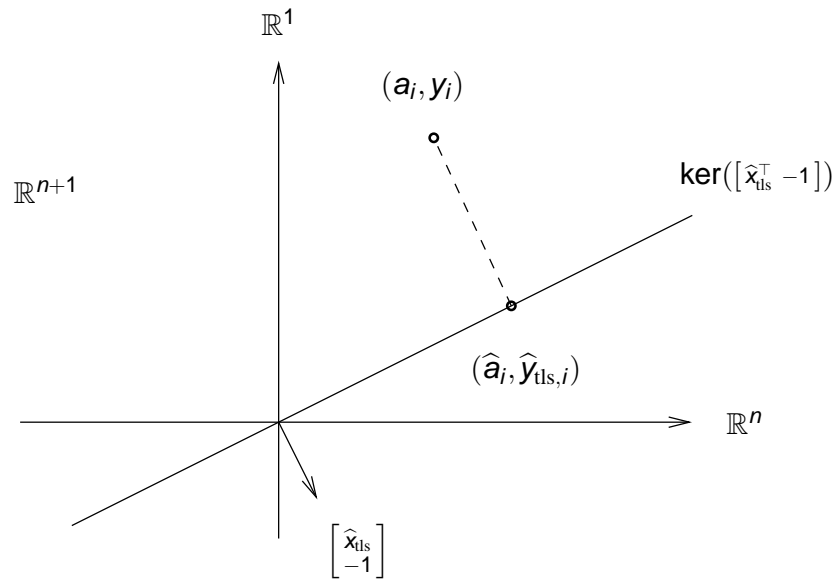
LS: minimize $_{e,x} \|e\|_2$ subject to $Ax = y + e$ **e latent variable**

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

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

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



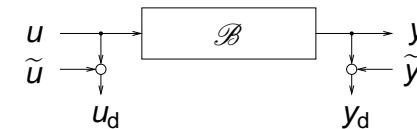


Statistical interpretation of misfit and latency

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

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

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

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.,

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

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

Computation of the misfit

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

$$\begin{aligned} \mathcal{B} \text{ is subspace} &\implies \text{the constraint is linear} \\ &\implies \text{ordinary LS problem} \end{aligned}$$

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. In addition, $\ell_{\max} \ll T$, implies banded structure with bandwidth ℓ_{\max} .

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

They are based on:

1. structured matrix computations (e.g., generalized Schur alg.)
2. Riccati recursions (Kalman smoother)

Maximum likelihood estimator in the EIV setup

Assuming that the data is generated according to the model

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

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

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

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

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

$$\text{likelihood evaluation} \iff \text{misfit computation}$$

Misfit computation using I/S/O representation

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

Recall from Lecture 5 that

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

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

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:

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

Efficient solution via Riccati recursion \rightsquigarrow Kalman smoother

Theorem The solution of (SS) with $D = 0$ and given x_{ini} is

$$\hat{u}(t) = -(B^\top P(t+1)B + I)^{-1} (B^\top P(t+1)A\hat{x}(t) + B^\top s(t+1) - u_d(t))$$

where \hat{x} is given by the **forward recursion**

$$\begin{aligned} \hat{x}(t+1) &= A\hat{x}(t) + B\hat{u}(t), & \hat{x}(0) &= x_{\text{ini}} \\ \hat{y}(t) &= C\hat{x}(t) \end{aligned}$$

and P, s are given by the **backward recursion**

$$\begin{aligned} P(t) &= -A^\top P(t+1)B(B^\top P(t+1)B + I)^{-1} \times \\ &\quad B^\top P(t+1)A + A^\top P(t+1)A + C^\top C, & P(T) &= 0 \end{aligned}$$

$$\begin{aligned} s(t) &= -A^\top P(t+1)B(B^\top P(t+1)B + I)^{-1} \times \\ &\quad (B^\top s(t+1) - u_d(t)) + A^\top s(t+1) - C^\top y_d(t), & s(T) &= 0 \end{aligned}$$

The proof uses dynamic programming

I. Markovsky and B. De Moor, Linear dynamic filtering with noisy input and output, Automatica, 41(1):167–171, 2005

Complete solution in the continuous-time case

I. Markovsky and J. C. Willems and B. De Moor, Continuous-time errors-in-variables filtering CDC 2002, pages 2576–2581

A Matlab toolbox for misfit identification:

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

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