

Lecture 1: Classical vs behavioral paradigms for data modeling

Ivan Markovsky

ELEC doctoral school 2013

Vrije Universiteit Brussel

Outline

Course mechanics

Motivating example

Low-rank approximation

Overview of applications

Overview of algorithms

Summary

Exercises

Course mechanics

- ▶ 4 sessions on Monday and Tuesday 13:00–17:00

- ▶ session = 1h lecture + 40' exercises + 20' break

"I hear, I forget; I see, I remember; I do, I understand."

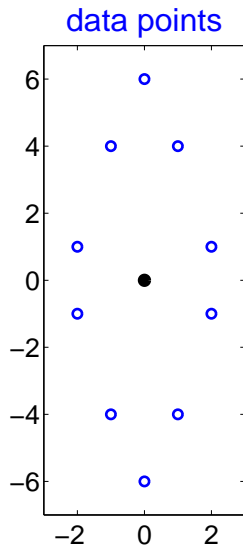
- ▶ topics

1. the behavioral paradigm for data modeling
2. exact system identification
3. approximate system identification

- ▶ book + lecture slides

homepages.vub.ac.be/~imarkovs/doctoral-school.html

A line fitting example



- ▶ classic problem: fit points

$$d_1 = \begin{bmatrix} 0 \\ 6 \end{bmatrix}, d_2 = \begin{bmatrix} 1 \\ 4 \end{bmatrix}, \dots, d_{10} = \begin{bmatrix} -1 \\ 4 \end{bmatrix}$$

by a line passing through the origin

- ▶ terminology and notation

data space

$$\mathcal{U} = \mathbb{R}^2$$

data

$$\mathcal{D} = \{d_1, \dots, d_{10}\}$$

model

$$\mathcal{B} \subset \mathcal{U}$$

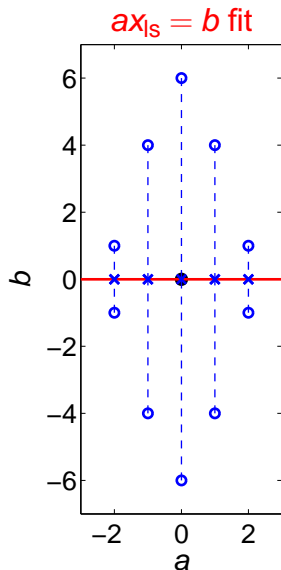
model class

\mathcal{M} : lines through 0

fitting criterion

$$\text{dist}(\mathcal{D}, \mathcal{B})$$

A line fitting example (cont.)



- ▶ classic solution: define

$$d_i =: \begin{bmatrix} a_i \\ b_i \end{bmatrix}, \quad A := \text{col}(a_1, \dots, a_{10}) \\ B := \text{col}(b_1, \dots, b_{10})$$

and solve a least squares problem

$$Ax \approx B$$

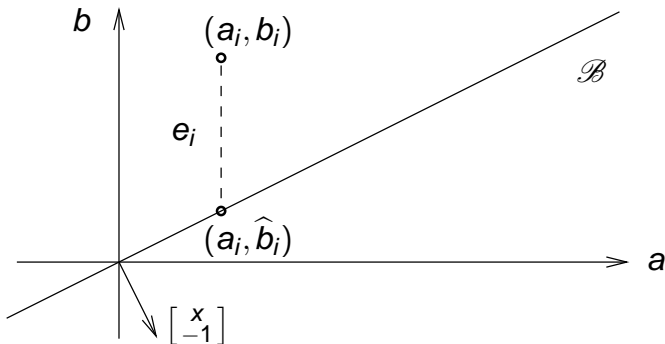
- ▶ the model is given by

$$\mathcal{B}_{ls} := \{ (a, b) \in \mathbb{R}^2 \mid ax_{ls} = b \}$$

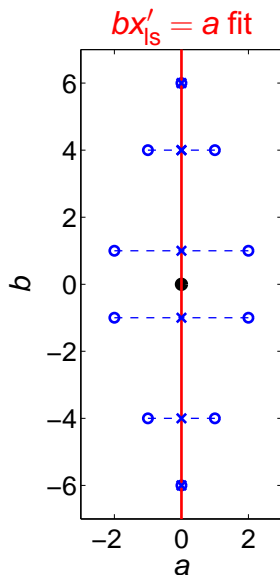
representation of a line through 0

$\mathcal{B}_{\text{ls}}(x_{\text{ls}})$ minimizes the **vertical distances** from \mathcal{D} to $\mathcal{B} \in \mathcal{M}$

$$\begin{aligned}\text{dist}_{\text{ls}}(\mathcal{D}, \mathcal{B}) &= \min_{\hat{B}} \|B - \hat{B}\|_2 \quad \text{s.t.} \quad (a_i, \hat{b}_i) \in \mathcal{B} \text{ for all } i \\ &= \|B - Ax_{\text{ls}}\|_2\end{aligned}$$



A line fitting example (cont.)



- ▶ we can also solve

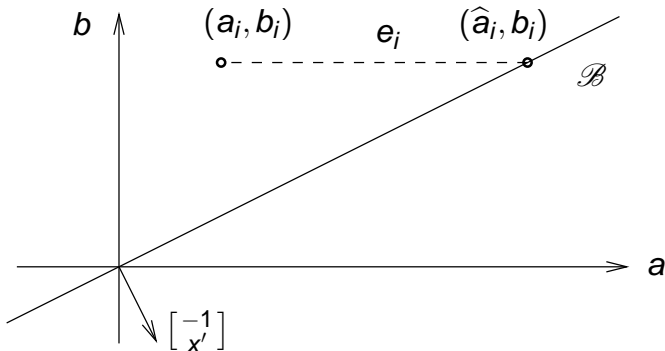
$$Bx' \approx A$$

- ▶ and obtain **another fitting line**

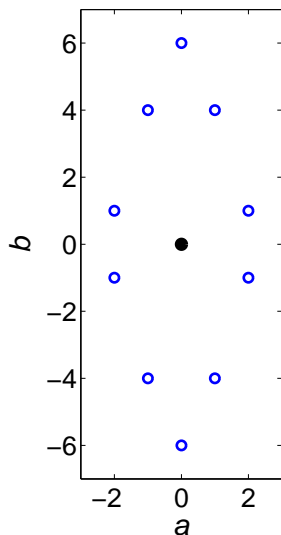
$$\mathcal{B}'_{\text{ls}} := \{ (a, b) \in \mathbb{R}^2 \mid bx'_{\text{ls}} = a \}$$

$\mathcal{B}'_{\text{ls}}(\mathbf{x}'_{\text{ls}})$ minimizes the **horizontal distances** from \mathcal{D} to $\mathcal{B} \in \mathcal{M}$

$$\begin{aligned}\text{dist}'_{\text{ls}}(\mathcal{D}, \mathcal{B}) &= \min_{\hat{A}} \|A - \hat{A}\|_2 \quad \text{s.t.} \quad (\hat{a}_i, b_i) \in \mathcal{B} \text{ for all } i \\ &= \|A - B\mathbf{x}'_{\text{ls}}\|_2\end{aligned}$$



A line fitting example (cont.)



- ▶ total least squares problem:

$$\min_{x, \hat{A}, \hat{B}} \sqrt{\|A - \hat{A}\|_2^2 + \|B - \hat{B}\|_2^2}$$

subject to $\hat{A}x = \hat{B}$

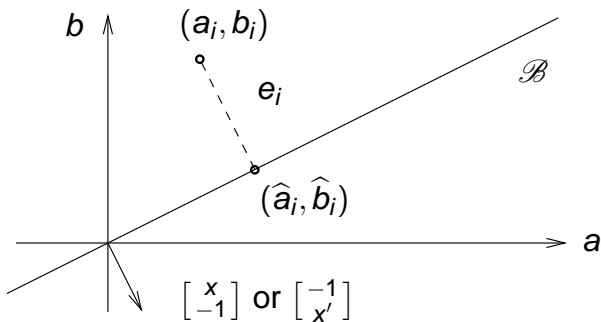
- ▶ x_{tls} does not exist ($x_{\text{tls}} = \infty$)
- ▶ however, $x'_{\text{tls}} = 0$ exists

$$\min_{x', \hat{A}, \hat{B}} \sqrt{\|A - \hat{A}\|_2^2 + \|B - \hat{B}\|_2^2}$$

subject to $\hat{B}x' = \hat{A}$

$\mathcal{B}_{\text{tls}}(\mathbf{x}_{\text{tls}})$ minimizes **orthogonal distances** from \mathcal{D} to $\mathcal{B} \in \mathcal{M}$

$$\begin{aligned} \text{dist}_{\text{tls}}(\mathcal{D}, \mathcal{B}) &= \min_{\hat{\mathcal{D}}} \sqrt{\sum_{i=1}^N \|\mathbf{d}_i - \hat{\mathbf{d}}_i\|_2^2} \quad \text{s.t.} \quad \hat{\mathcal{D}} \subset \mathcal{B} \\ &= \frac{\|\mathbf{A}\mathbf{x}_{\text{tls}} - \mathbf{B}\|_2^2}{1 + \|\mathbf{x}_{\text{tls}}\|_2^2} \end{aligned}$$



Conclusions

- ▶ least squares is representation dependent
- ▶ total least squares is representation invariant
- ▶ total least squares may have no solution

What are the issues?

- ▶ a representation is nonunique
- ▶ its choice should be independent of a fitting criterion
- ▶ a representation should cover all models in \mathcal{M}

Inputs and outputs

- ▶ the model considered is linear static with var. (a, b)
- ▶ two different representations:

$$\{(a, b) \in \mathbb{R}^2 \mid ax = b\} \quad (*)$$

$$\{(a, b) \in \mathbb{R}^2 \mid bx' = a\} \quad (**)$$

- ▶ $(*)$ and $(**)$ define \mathcal{B} by **functions**
 - ▶ $a \mapsto b$ in $(*)$
 - ▶ $b \mapsto a$ in $(**)$
- ▶ input/output representations
 - ▶ in $(*)$, a is input, b is output (a causes b)
 - ▶ in $(**)$, b is input, a is output (b causes a)

Input/output representation

- ▶ separately $(*)$, $(**)$ don't represent all models in \mathcal{M}
- ▶ **I/O representation:** any $\mathcal{B} \in \mathcal{M}$ is representable as

$$\mathcal{B} = \mathcal{B}(\mathbf{x}, \Pi) = \{ \Pi \begin{bmatrix} a \\ b \end{bmatrix} \mid \mathbf{a}\mathbf{x} = \mathbf{b} \}$$

for some $\mathbf{x} \in \mathbb{R}$ and a permutation matrix Π

- ▶ link to system of linear equations

$$\mathcal{D} \subset \mathcal{B}(\mathbf{x}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}) \iff \mathbf{A}\mathbf{x} = \mathbf{B}$$

$$\mathcal{D} \subset \mathcal{B}(\mathbf{x}, \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}) \iff \mathbf{B}\mathbf{x}' = \mathbf{A}$$

where

$$\mathbf{A} := \text{col}(\mathbf{a}_1, \dots, \mathbf{a}_{10}) \quad \text{and} \quad \mathbf{B} := \text{col}(\mathbf{b}_1, \dots, \mathbf{b}_{10})$$

Kernel representation

- ▶ any $\mathcal{B} \in \mathcal{M}$ can be represented as

$$\mathcal{B} = \ker(R) := \{ d \in \mathbb{R}^2 \mid Rd = R_1 a + R_2 b = 0 \}$$

for some nonzero vector $R \in \mathbb{R}^{1 \times 2}$

- ▶ $Rd = 0$ defines a **relation** between a and b
- ▶ $\mathcal{D} \subset \ker(R)$ implies that

$$R \underbrace{\begin{bmatrix} d_1 & \cdots & d_N \end{bmatrix}}_D = 0$$

Image representation

- ▶ any $\mathcal{B} \in \mathcal{M}$ can be represented as

$$\mathcal{B} = \text{image}(P) := \{ d = P\ell \mid \ell \in \mathbb{R} \}$$

for some vector $P \in \mathbb{R}^{2 \times 1}$

- ▶ $\mathcal{D} \subset \text{image}(P)$ implies that

$$[d_1 \quad \cdots \quad d_N] = PL$$

for some $L \in \mathbb{R}^{1 \times N}$

Important observation

- ▶ common feature of the representations considered

$$\left. \begin{array}{lll} \exists \mathbf{x} \in \mathbb{R} & A\mathbf{x} = B & \implies \\ \exists \mathbf{x}' \in \mathbb{R} & B\mathbf{x}' = A & \implies \\ \exists \mathbf{x} \in \mathbb{R}, \Pi \text{ permut.} & \begin{bmatrix} \mathbf{x} & -1 \end{bmatrix} \Pi D = 0 & \iff \\ \exists R \in \mathbb{R}^{1 \times 2}, R \neq 0 & RD = 0 & \iff \\ \exists P \in \mathbb{R}^{2 \times 1}, L \in \mathbb{R}^{1 \times N} & D = PL & \iff \end{array} \right\} \text{rank}(D) = 1$$

- ▶ representation free characterization of the exact data

$$\begin{array}{c} \mathcal{D} \subset \mathcal{B} \text{ and} \\ \mathcal{B} \text{ is a line through } 0 \\ \Updownarrow \\ \text{rank}(D) = 1 \end{array}$$

Low-rank approximation

- ▶ representation free formulation
- ▶ exact modeling problem:

\exists exact model for $\mathcal{D} \iff D$ is rank deficient

- ▶ approximate modeling problem:

minimize over $\hat{\mathcal{D}}$ $\text{dist}(\mathcal{D}, \hat{\mathcal{D}})$
subject to \exists exact model for $\hat{\mathcal{D}}$



minimize over \hat{D} $\text{dist}(D, \hat{D})$
subject to \hat{D} is rank deficient

Generalizations

1. multivariable data fitting $\mathcal{U} = \mathbb{R}^q$

- ▶ linear static model \leftrightarrow subspace
- ▶ model complexity \leftrightarrow subspace dimension
- ▶ $\text{rank}(D) \leftrightarrow$ upper bound on the model complexity

2. nonlinear static modeling

- ▶ $\mathcal{D} \mapsto D$ — nonlinear function
- ▶ nonlinearly structured low-rank approximation

3. linear time-invariant dynamical models

- ▶ $\mathcal{D} \mapsto$ Hankel matrix D
- ▶ Hankel structured low-rank approximation

4. nonlinear dynamic (2. + 3.)

Related frameworks

- ▶ behavioral approach in systems and control theory
 - ▶ representation free: model = set (the behavior)
 - ▶ no a priori separation of inputs and outputs
- ▶ errors-in-variables modeling
 - ▶ all variables are perturbed by noise
 - ▶ maximum likelihood estimation \leftrightarrow LRA
- ▶ principal component analysis
 - ▶ another statistical setting for LRA
- ▶ factor analysis
 - ▶ factors \leftrightarrow latent variables in image repr.

Overview of applications

- ▶ **systems and control**
 - ▶ approximate realization / model reduction
 - ▶ system identification
- ▶ **signal processing**
 - ▶ approximate deconvolution
 - ▶ image deblurring
- ▶ **machine learning**
 - ▶ dimensionality reduction
 - ▶ recommender systems
- ▶ **computer algebra**
 - ▶ approximate common divisor

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

unstructured

Hankel

$q \times 1$ Hankel

$q \times N$ Hankel

mosaic Hankel

[Hankel unstructured]

block-Hankel Hankel-block

linear static

scalar LTI

q -variate LTI

N equal length traj.

N general traj.

finite impulse response

2D linear shift-invariant

Problems with analytic solutions

- ▶ unstructured, unweighted $(\|\cdot\|_F := \|\text{vec}(\cdot)\|_2)$

$$\begin{array}{ll} \text{minimize} & \text{over } \hat{D} \quad \|D - \hat{D}\|_F \\ \text{subject to} & \text{rank}(\hat{D}) \leq r \end{array} \quad (\text{LRA})$$

- ▶ unstructured, with some fixed rows
- ▶ unstructured, with left/right weighting matrices

$$\begin{array}{ll} \text{minimize} & \text{over } \hat{D} \quad \|W_l(D - \hat{D})W_r\|_F \\ \text{subject to} & \text{rank}(\hat{D}) \leq r \end{array}$$

- ▶ circulant structure, unweighted

Theorem (Eckart–Young–Mirsky)

Let $D = U\Sigma V^\top$ be the (thin) SVD of $D \in \mathbb{R}^{q \times N}$ and define

$$U = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{matrix} r & q-r \\ q & \end{matrix}, \quad \Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{matrix} r & q-r \\ q-r & r \end{matrix}, \quad V = \begin{bmatrix} V_1 & V_2 \end{bmatrix} \begin{matrix} r & q-r \\ N & \end{matrix}$$

An optimal low-rank approximation (a solution of (LRA)) is

$$\hat{D}^* = U_1 \Sigma_1 V_1^\top, \quad \hat{\mathcal{B}}^* = \ker(U_2^\top) = \text{image}(U_1).$$

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

- ▶ “truncated SVD”
- ▶ $\hat{\mathcal{B}}^*$ depends only on the left singular vectors
- ▶ in general
 - ▶ structures other than circular
 - ▶ norms other than 2-norm
 - ▶ weights other than “left/right” multiplication of $D - \hat{D}$lead to hard non-convex optimization problems
- ▶ there are many (heuristic) solution methods

Overview of algorithms

- ▶ **global solution methods**
 - ▶ SDP relaxations of rational function min. problem
 - ▶ systems of polynomial equations (computer algebra)
 - ▶ branch-and-bound, simulating annealing, ...
- ▶ **local optimization methods**
 - ▶ variable projections (Lecture 3)
 - ▶ alternating projections
 - ▶ variations (parameterization + optimization method)
- ▶ **convex relaxations / multistage methods**
 - ▶ subspace methods (Lecture 2)
 - ▶ nuclear norm heuristic

Summary

- ▶ linear static model = subspace
- ▶ model representations
 - ▶ input/output (a function, system $AX = B$)
 - ▶ kernel (implicit function, relation)
 - ▶ image (introduces latent variables)
- ▶ representation invariant problem formulation \leadsto LRA
- ▶ different representations \leadsto different solution methods

*“... most linear resistors let us treat current as a function of voltage or voltage as a function of current, since R is neither zero nor infinite. But in the two limiting cases - the short circuit and the open circuit - that's not true. To fit these cases neatly in a unified framework, we shouldn't think of the relation between current and voltage as defining a function. **It's just a relation!**”*

John Baez

Instructions

- ▶ individual work (1–5 min.)
- ▶ discussion within groups (\sim 4 people)
- ▶ Konstantin, Mariya, and Ivan are around to help
- ▶ give us a sign when you are ready
- ▶ class discussion
- ▶ there are problems for homework

Exercise 1: Constant approximation

- ▶ the approximate modeling problem:

$$\begin{array}{ll} \text{minimize} & \text{over } \hat{D} \quad \|D - \hat{D}\|_F \\ \text{subject to} & \{\hat{d}_1, \dots, \hat{d}_N\} \subset \text{line through } 0 \end{array} \quad (*)$$

is equivalent to low-rank approximation

- ▶ the solution is given by the EYM theorem (SVD of D)
- ▶ solve $(*)$ with additional constraint

$$\hat{d}_i = \hat{d}, \quad \text{for some } \hat{d} \in \mathbb{R}^q \quad (**)$$

Exercise 2: Line fitting

- ▶ we considered the problem of fitting

$$\mathcal{D} = \{d_1, \dots, d_N\} \subset \mathbb{R}^2$$

to a line passing through the origin

- ▶ the solution is rank-1 approximation of the matrix

$$D = [d_1 \quad \dots \quad d_N]$$

- ▶ consider now the problem of fitting \mathcal{D} to any line in \mathbb{R}^2
- ▶ is it also equivalent to low-rank approximation?

Exercise 3: Conic section fitting

- ▶ conic section model

$$\mathcal{B}(S, u, v) = \{ d \in \mathbb{R}^2 \mid d^\top S d + u^\top d + v = 0 \}$$

where $S = S^\top$, u , v are model parameters

- ▶ express the exact fitting condition $\mathcal{D} \subset \mathcal{B}(S, u, v)$ as a rank constraint on a matrix $D(\mathcal{D})$
- ▶ find a conic section fitting the points

$$d_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \quad d_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad d_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad d_4 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Exercise 4: Subspace clustering

- ▶ this problem is a special case of the **Generalized PCA**
- ▶ union of two lines model

$$\mathcal{B}(R^1, R^2) = \{d \in \mathbb{R}^2 \mid (R^1 d)(R^2 d) = 0\}$$

where $R^1, R^2 \in \mathbb{R}^{1 \times 2}$, $R^1, R^2 \neq 0$ are model parameters

- ▶ express the exact fitting condition $\mathcal{D} \subset \mathcal{B}(R^1, R^2)$ as a rank constraint on a matrix $D(\mathcal{D})$
- ▶ find a union of two lines model fitting the points

$$d_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \quad d_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad d_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad d_4 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Exercise 5: LTI autonomous system fitting

- ▶ linear time-invariant autonomous system

$$\begin{aligned} \mathcal{B}|_T(R) = \{ (w(1), \dots, w(T)) \mid \\ R_0 w(t) + R_1 w(t+1) + \dots + R_\ell w(t+\ell) = 0, \\ \text{for } t = 1, \dots, T - \ell \} \end{aligned}$$

- ▶ express the exact fitting condition $w \in \mathcal{B}(R)$ as a rank constraint on a matrix constructed from w
- ▶ find the smallest ℓ , for which $\exists R \in \mathbb{R}^{1 \times (\ell+1)}$, such that

$$w_d := (1, 2, 4, 7, 13, 24, 44, 81) \in \mathcal{B}_8(R)$$

Exercise 6: Polynomial common divisor

- ▶ the polynomials

$$p(z) = p_0 + p_1 z + \cdots + p_{\ell_p} z^{\ell_p}$$

$$q(z) = q_0 + q_1 z + \cdots + q_{\ell_q} z^{\ell_q}$$

have a common divisor

$$c(z) = c_0 + c_1 z + \cdots + c_{\ell_c} z^{\ell_c}$$

iff $p = ca$ and $q = cb$ for some polynomials a and b

- ▶ express the common divisor condition as a rank constraint on a matrix constructed from p , q