

# Low-rank approximation: Applications and algorithms

Ivan Markovsky

School of Electronics and Computer Science  
University of Southampton

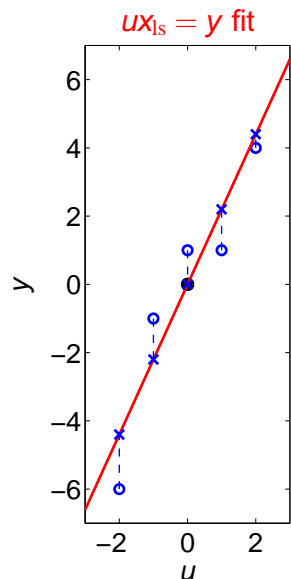
## Outline

Introduction

Applications

Algorithms

## What is a model?



**Classic problem:** Fit the points

$$d_1 = \begin{bmatrix} -2 \\ -6 \end{bmatrix}, d_2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \dots, d_5 = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

by a line passing through the origin.

**Classic solution:** Define  $d_i =: \text{col}(u_i, y_i)$   
and solve the least squares problem

$$\text{col}(u_1, \dots, u_5)x = \text{col}(y_1, \dots, y_5).$$

**The model** is the line

$$\mathcal{B} := \{ d = \text{col}(u, y) \mid ux_{ls} = y \}$$

and not the equation  $ux_{ls} = y$ .

## Model representations

In general,

a linear static model is a subspace  $\mathcal{B}$  of the data space  $\mathbb{R}^q$

**Representations** of a linear static model  $\mathcal{B} \subseteq \mathbb{R}^q$ :

**kernel**  $\mathcal{B} = \ker(R) \quad := \{ d \mid Rd = 0 \}$

**image**  $\mathcal{B} = \text{image}(P) \quad := \{ d = Pv \mid \text{for all } v \}$

**input/output**  $\mathcal{B} = \mathcal{B}_{i/o}(X) \quad := \{ d = \text{col}(u, y) \mid Xu = y \}$

## Links among model representations

$$\begin{array}{c} \text{input/output} \end{array} \quad ux = y \quad \Longleftrightarrow \quad \begin{array}{c} \text{kernel} \\ \underbrace{\begin{bmatrix} x & -1 \end{bmatrix}}_R \end{array} \begin{array}{c} \begin{bmatrix} u \\ y \end{bmatrix} = 0 \end{array} \quad \Longleftrightarrow \quad \begin{array}{c} \text{image} \\ \begin{bmatrix} u \\ y \end{bmatrix} = \underbrace{\begin{bmatrix} 1 \\ x \end{bmatrix}}_P u \end{array}$$

Therefore,  $\mathcal{B} = \{d = \text{col}(u, y) \mid ux = y\} = \ker(R) = \text{image}(P)$

In general:

$$\begin{array}{ccc} \mathcal{B} = \ker(R) & \xleftrightarrow{RP=0} & \mathcal{B} = \text{image}(P) \\ \swarrow \begin{array}{l} X^T = -R_o^{-1} R_i \\ R = [X^T \ -I] \end{array} & & \searrow \begin{array}{l} X^T = P_o P_i^{-1} \\ P^T = [I \ X] \end{array} \\ & \mathcal{B} = \mathcal{B}_{i/o}(X) & \end{array}$$

## Exact modelling

Consider a given data set

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

A model  $\mathcal{B} \subseteq \mathbb{R}^q$  is exact for the data  $\mathcal{D}$  if  $\mathcal{D} \subseteq \mathcal{B}$ .

Exact data modelling problem:

Find a least complex model  $\mathcal{B}$  in a given set of models  $\mathcal{M}$  that fits the data  $\mathcal{D}$  exactly or assert that such a model doesn't exist.

## Model complexity

The dimension of  $\mathcal{B}$ ,  $\dim(\mathcal{B})$  is a measure of  $\mathcal{B}$ 's "complexity".

In the example,

$$\dim(\mathcal{B}) = 1 = \text{rank}(P) = 2 - \text{rank}(R)$$

- $\mathcal{B} = \{0\}$  has  $\dim(\mathcal{B}) = 0$  and is the **least complex model**,
- $\mathcal{B}$  = "line passing through the origin" has  $\dim(\mathcal{B}) = 1$ , and
- $\mathcal{B} = \mathbb{R}^2$  has  $\dim(\mathcal{B}) = 2$  and is the **most complex model**.

## Notes

- The most basic (and simple) data modelling problem.
- The exact case should be considered before (more complicated) approximate and stochastic cases.
- The solution of the exact problem is useful in the solution of approximate and stochastic cases.
- A core question in all sciences. For example,
  - Kepler's laws define a model that fit exactly planetary trajectories,
  - Newton's laws of dynamics define a model that fit exactly the trajectory of any moving body

## Exact linear static model

In the case, when the model class

$\mathcal{M}$  = set of all linear static models

the solution of the exact modelling problem is simple.

- A solution always exists and is unique.
- It is given by  $\mathcal{B} = \text{image}(D)$ , where

$$D := [d_1 \quad \dots \quad d_N] \in \mathbb{R}^{q \times N}$$

- The complexity of  $\mathcal{B}$  is equal to the rank of  $D$

Generically,  $\text{rank}(D) = q$ , so  $\mathcal{B}$  is trivial model (fits everything).

## Low-rank approximation (LRA)

Let  $\hat{D} \in \mathbb{R}^{q \times N}$  be the perturbed data. We want

1.  $\hat{D}$  to be as close as possible to  $D$ , e.g.,  $\min \|D - \hat{D}\|$
2.  $\hat{\mathcal{B}}$  to be an exact model for  $\hat{D}$ , i.e.,  $\hat{D} \in \hat{\mathcal{B}}$
3.  $\hat{\mathcal{B}}$  to have complexity bounded by  $r < q$ , i.e.,  $\dim(\hat{\mathcal{B}}) \leq r$

$$\hat{D} \in \hat{\mathcal{B}} \quad \text{and} \quad \dim(\hat{\mathcal{B}}) \leq r \quad \implies \quad \text{rank}(\hat{D}) \leq r$$

**Approximate modelling problem:** Given  $D \in \mathbb{R}^{q \times N}$ ,  $r$ , and  $\|\cdot\|$ ,

$$\text{minimize over } \hat{D} \quad \|D - \hat{D}\| \quad \text{subject to} \quad \text{rank}(\hat{D}) \leq r$$

## Approximate modelling

In the case, when the model class

$\mathcal{M}$  = set of all linear static models of bounded complexity

a solution to the exact modelling problem may not exist.

**Approximate modelling problem:**

Find a smallest (in specified sense) perturbation of the data  $\mathcal{D}$  that renders exact modelling of the perturbed data solvable.

## Relation to regression problems

The classical approach for data fitting is regression ( $AX \approx B$ ). Regression corresponds to LRA with input/output representation.

$$AX = B \implies \begin{bmatrix} A & B \end{bmatrix} \begin{bmatrix} X \\ -I \end{bmatrix} = 0 \implies \text{rank}(\begin{bmatrix} A & B \end{bmatrix}) \leq \text{col dim}(X)$$

$\implies$  indeed, regression is a way to achieve low-rank approximation.

However,  $AX = B$  does not imply  $\text{rank}(\begin{bmatrix} A & B \end{bmatrix}) \leq \text{col dim}(X)$ .

$\rightsquigarrow$  existence of ill-posed regression problems.

Ill-posedness/conditioning is a consequence of imposing input/output structure on the model that is not corroborated by the data.

## Rank minimization (RM)

Approximate modeling is a trade-off between:

- fitting accuracy ( $\|D - \hat{D}\|$ ) and
- model complexity ( $\text{rank}(\hat{D})$ )

Two possible scalarizations of the **bi-objective optimization** are:

**LRA:** maximize accuracy under a constraint on complexity

**RM:** minimize complexity under a constraint ( $\mathcal{C}$ ) on accuracy

$$\text{minimize over } \hat{D} \quad \text{rank}(\hat{D}) \quad \text{subject to } \hat{D} \in \mathcal{C}$$

RM (as well as LRA) is **NP-hard**, however, there are effective heuristics for RM, e.g., with  $\hat{D} = \text{diag}(\hat{d})$ ,  $\text{rank}(\hat{D}) = \text{card}(\hat{d})$ ,

$$\ell_1 \text{ heuristic: } \min_{\hat{d}} \|\hat{d}\|_1 \quad \text{subject to } \text{diag}(\hat{d}) \in \mathcal{C}$$

## Applications

- **System theory**
  1. Approximate realization
  2. Model reduction
  3. System identification
- **Signal processing**
  4. Linear prediction
  5. FIR modeling
  6. Harmonic retrieval
  7. Array processing
  8. Image deblurring
- **Computer algebra**
  9. Approximate GCD
- **Machine learning**
  10. Data compression
  11. Natural language proc.
  12. Psychometrics
  13. Recommender systems
- **Computer vision**
  14. Structure from motion
- **Chemometrics**
  15. Multivariate calibration

## Generalizations

- **Cost functions**

- weighted norms  $\|\Delta\| = \text{vec}^\top(\Delta) W \text{vec}(\Delta)$
- information criteria  $\|\Delta\| \leftrightarrow \log \det(\Delta)$

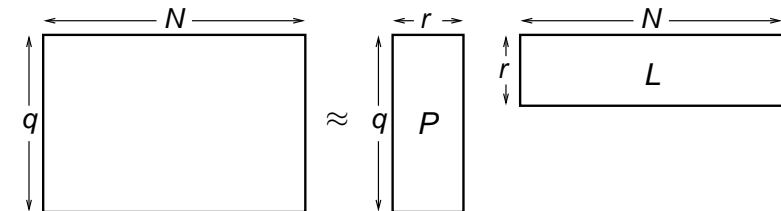
- **Constraints on  $\hat{D}$**

- structured, e.g., Hankel, Sylvester, sparse
- nonnegative
- exact elements

- **Data  $D$**

- tensors
- categorical data
- missing data

## Data compression



$qN$  elements vs  $(q + N)r$  elements

For  $r \ll \max(q, N)$ ,  $qN \gg (q + N)r \implies$  data compression

## Natural language processing

Consider  $N$  documents, involving  $q$  terms and  $r$  concepts.

$d_{ij}$  — frequency of  $i$ th term in  $j$ th document

$\ell_{kj}$  — relevance of  $k$ th concepts to  $j$ th document

$p_{ik}$  — frequency of  $i$ th term in a document of  $k$ th concept only

$$d_j = \begin{bmatrix} d_{1j} \\ \vdots \\ d_{qj} \end{bmatrix} \in \mathbb{R}^q, \quad p_k = \begin{bmatrix} p_{1k} \\ \vdots \\ p_{qk} \end{bmatrix} \in \mathbb{R}^q, \quad \ell_k = \begin{bmatrix} \ell_{1j} \\ \vdots \\ \ell_{rN} \end{bmatrix} \in \mathbb{R}^r$$

Latent semantic analysis model:

$$\underbrace{\begin{bmatrix} d_1 & \cdots & d_N \end{bmatrix}}_D = \underbrace{\begin{bmatrix} p_1 & \cdots & p_r \end{bmatrix}}_P \underbrace{\begin{bmatrix} \ell_1 & \cdots & \ell_N \end{bmatrix}}_L$$

## Approximate latent semantic analysis

The LSA model does not hold exactly because

- the notion of (small number of) concepts is an idealization
- linearity assumption  $D = PL$  is not likely to hold in practice

LRA is used to find a few concepts explaining the data approx.

Document classification:

similarity of documents is evaluated in the concepts space

Synonym discovery:

terms are clustered in the concepts space

Documents search by keywords:

translate first the keywords to a vector in the concepts space and then finding a cluster of documents nearby this vector

## Exact latent semantic analysis

Assuming

- fewer concepts than terms or documents,
- independent concepts, i.e.,  $p_1, \dots, p_r$  linearly independent,
- independent relevance vectors  $\ell_1, \dots, \ell_r$

$\text{rank}(D)$  = the number of concepts related to the documents.

In a rank revealing factorization  $D = PL$ ,

- $P$  indicate relevance of the concepts to the documents
- $L$  indicate the term frequencies related to the concepts

## Psychometrics

The data  $D$  consists of test scores of a group of people

Psychometrics tries to explain the data as a result of a few underlying abilities.

$d_{ij}$  — score in  $i$ th test of  $j$ th person

$\ell_{kj}$  — amount of  $k$ th ability in  $j$ th person

$p_{ik}$  — score in  $i$ th test of a person with  $k$ th ability only

Factor analysis model:

$$\underbrace{\begin{bmatrix} d_1 & \cdots & d_N \end{bmatrix}}_D = \underbrace{\begin{bmatrix} p_1 & \cdots & p_r \end{bmatrix}}_P \underbrace{\begin{bmatrix} \ell_1 & \cdots & \ell_N \end{bmatrix}}_L$$

$\Rightarrow \text{rank}(D) = \#$  of abilities relevant for the tests.

verbal, quantitative, and analytical ability, are believed to be most important in explaining one's academic performance.

## Outline

Introduction

Applications

Algorithms

## Basic low-rank approximation problem

$$\hat{D}^* := \arg \min_{\hat{D}} \|D - \hat{D}\|_F \quad \text{subject to} \quad \text{rank}(\hat{D}) \leq r$$

Theorem (closed form solution)

Let  $D = U\Sigma V^\top$  be the SVD of  $D$  and define

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

An optimal low-rank approximation solution is

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

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

The basic LRA problem is an exception: other approx. criteria and extra constraints lead to NP-hard problems.

Double minimization nature of LRA:

$$\min_{\hat{\mathcal{B}} \in \mathcal{M}} \left( \min_{\hat{D}} \|D - \hat{D}\|_F \quad \text{subject to} \quad \hat{D} \in \hat{\mathcal{B}} \right) \quad (*)$$

If  $\hat{\mathcal{B}}$  is linear, the inner minimization has analytic solution. It gives the distance of the data to the model  $\hat{\mathcal{B}}$ .

Using an image representation of the model:

$$\min_{P \in \mathbb{R}^{q \times r}} \left( \min_{L \in \mathbb{R}^{r \times N}} \|D - PL\|_F \right) = \min_{L \in \mathbb{R}^{r \times N}} \left( \min_{P \in \mathbb{R}^{q \times r}} \|D - PL\|_F \right) \quad (**)$$

For fixed  $P$  the problem is linear in  $L$  and vice versa.

## Variable projection vs. alternating projections

Two ways to approach the double minimization:

- **Variable projections (VARPRO):**  
solve the inner minimization of (\*) analytically  
 $\rightsquigarrow$  nonlinear least squares problem for the model parameters
- **Alternating projections (AP):**  
Alternate between the least squares problems, resulting from (\*\*) with fixed  $P$  and  $L$ , respectively

VARPRO is globally convergent with a super linear conv. rate.

AP is globally convergent with a linear convergence rate.

## Summary

- Linear static models = subspaces. Can be represented as image or kernel of matrix, or graph of map (input/output)
- Exact modeling is not practical but is conceptually useful.
- Approximate modeling is a bi-objective optimization: accuracy vs complexity trade-off.
- LRA — approximate modeling with complexity bound. Regression is special case when input/output is used.
- Most LRA problems have no analytic solution.  
Two basic solution approaches: VARPRO and AP

Thank you