

Structured low-rank approximation with missing data

Ivan Markovsky and Konstantin Usevich

Department ELEC
Vrije Universiteit Brussel
{imarkovs, kusevich}@vub.ac.be

Abstract

We consider low-rank approximation of affinely structured matrices with missing elements. The method proposed is based on reformulation of the problem as inner and outer optimization. The inner minimization is a singular linear least-norm problem and admits an analytic solution. The outer problem is a nonlinear least squares problem and is solved by local optimization methods: minimization subject to quadratic equality constraints and unconstrained minimization with regularized cost function. The method is generalized to weighted low-rank approximation with missing values and is illustrated on approximate low-rank matrix completion, system identification, and data-driven simulation problems. An extended version of the paper is a literate program, implementing the method and reproducing the presented results.

Keywords: low-rank approximation, missing data, variable projections, system identification, approximate matrix completion.

AMS subject classifications: 93A30, 93B30, 93B40, 37M10, 41A29, 62J12, 49N30, 15A83.

1 Introduction

The paper describes a solution method for matrix structured low-rank approximation, *i.e.*, approximation of a given matrix by another matrix whose elements satisfy certain predefined relations (matrix structure) and whose rank is less than or equal to a predefined value. The combination of matrix and low-rank structure makes structured low-rank approximation a tool for data modeling. Low-rank property of a matrix is equivalent to existence of an exact low-complexity linear model for the data. Moreover, the rank of the matrix is related to the complexity of the model. The structure, imposed on the approximation, is related to properties of the model. For example, Hankel structure corresponds to time-invariance of a linear dynamical model for the data. A tutorial exposition of the subject with an emphasis on applications in systems theory, control, and signal processing, is given in [12, 14].

A novel feature of the low-rank approximation problem, considered in this paper, is that elements of the data matrix can be missing (not specified). Missing data may occur in practical applications due to malfunctioning of measurement device, communication channel, or storing device. In such cases, the most common strategy is to collect a complete data record by repeating the data collection experiment. In other applications, however, the missing data problem is intrinsic and can not be avoided by repeated experiments. Two such applications reviewed in Subsection 1.2 and further illustrated by numerical examples in Section 5 are recommender systems and data-driven simulation.

Although structured low-rank approximation and approximate low-rank matrix completion (missing data estimation in low-rank matrices) are independently active research topics, the combined problem of missing data estimation in affine structured low-rank matrices has not been considered before. Both the structured low-rank approximation and approximate matrix completion problems are nonconvex optimization problems that in general admit no analytic solution. Therefore, in both domains local optimization and convex relaxation heuristics are used as solution techniques. In this paper, we use the local optimization approach.

Structured low-rank approximation has been studied in the literature from different viewpoints: numerical algorithm for computing locally optimal or suboptimal solutions, statistical properties of the resulting estimators, and applications. From a numerical point of view, the main challenge is to achieve fast and robust computational methods that can deal effectively with large data sets. From a statistical point of view, the main challenge is to establish conditions for consistency and efficiency of the methods. Our objective in this paper is different: we aim to unify as

many data modeling applications as possible and derive a single algorithm that solves them. Of course, this goal can be achieved by brute force optimization. The challenge is to discover and use effectively the structure of the problem. In the general problem considered, this structure is a separation of variables with analytic solution over one set of variables. This approach is related to the variable projections method [8] used in [18].

The general solution method proposed in the paper has computational complexity $O(n^3)$ per iteration, where n is the number of columns of the data matrix. This makes it unsuitable for large scale applications. In special cases, such as unstructured matrix with missing data and Hankel structured matrix, there are efficient $O(n)$ methods. Modification of the methods in the paper for efficient computation in the case of mosaic-Hankel [9] matrices, is possible and will be presented elsewhere.

1.1 Problem formulation

We denote missing data values by the symbol NaN (“not a number”). The considered low-rank approximation problem is: Given a data vector $p \in (\mathbb{R} \cup \{\text{NaN}\})^{n_p}$,

$$\begin{aligned} & \text{minimize} \quad \text{over } \hat{p} \in \mathbb{R}^{n_p} \quad \sum_{\{i \mid p_i \neq \text{NaN}\}} (p_i - \hat{p}_i)^2 \\ & \text{subject to} \quad \text{rank}(\mathcal{S}(\hat{p})) \leq r, \end{aligned} \tag{SLRA}$$

where

$$\mathcal{S} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{m \times n}, \quad \text{defined by} \quad \mathcal{S}(\hat{p}) = S_0 + \sum_{i=1}^{n_p} S_i \hat{p}_i, \tag{\mathcal{S}}$$

is the matrix structure—an affine function from the structure parameter space \mathbb{R}^{n_p} to the set of matrices $\mathbb{R}^{m \times n}$. With \mathcal{G} denoting the vector of indices of the given values $\{i \mid p_i \neq \text{NaN}\}$ and $p_{\mathcal{G}}$ denoting the subvector of p with indices in \mathcal{G} , the approximation criterion can be written as

$$\sum_{i \in \mathcal{G}} (p_i - \hat{p}_i)^2 = \|p_{\mathcal{G}} - \hat{p}_{\mathcal{G}}\|_2^2.$$

Without loss of generality, we assume throughout the paper that $r < m \leq n$. Using the kernel representation of the rank constraint

$$\text{rank}(\mathcal{S}(\hat{p})) \leq r \iff \text{there is } R \in \mathbb{R}^{(m-r) \times m}, \text{ such that } R\mathcal{S}(\hat{p}) = 0 \text{ and } R \text{ has full row rank}, \tag{KER}$$

the following equivalent problem to (SLRA) is obtained

$$\begin{aligned} & \text{minimize} \quad \text{over } \hat{p} \in \mathbb{R}^{n_p} \text{ and } R \in \mathbb{R}^{(m-r) \times m} \quad \|p_{\mathcal{G}} - \hat{p}_{\mathcal{G}}\|_2^2 \\ & \text{subject to} \quad R\mathcal{S}(\hat{p}) = 0 \quad \text{and} \quad R \text{ has full row rank.} \end{aligned} \tag{SLRA}_R$$

Problem (SLRA)_R is a double minimization over the parameters R and \hat{p}

$$\text{minimize} \quad \text{over } R \in \mathbb{R}^{(m-r) \times m} \quad M(R) \quad \text{subject to} \quad R \text{ has full row rank}, \tag{SLRA}'_R$$

where

$$M(R) := \min_{\hat{p}} \|p_{\mathcal{G}} - \hat{p}_{\mathcal{G}}\|_2^2 \quad \text{subject to} \quad R\mathcal{S}(\hat{p}) = 0. \tag{INNER}$$

The evaluation of the cost function M , *i.e.*, solving (INNER) for a given value of R , is referred to as the *inner minimization problem*. This problem is solved analytically in Section 2. The remaining problem of minimizing M over R is referred to as the *outer minimization problem*. It is a nonlinear least-squares problem, which, in general, admits no analytic solution. General purpose local optimization methods are used in Section 3 for its numerical solution. In Section 4, the approach is generalized to weighted 2-norm approximation criteria with missing values. Numerical examples of solving approximation problems with missing data by the proposed methods are shown in Section 5.

1.2 Applications

1.2.1 Linear static modeling with missing data: An approximate low-rank matrix completion problem

Consider a set of vectors $\mathcal{D} = \{d_1, \dots, d_N\}$ in \mathbb{R}^q . A linear model \mathcal{B} for the data \mathcal{D} is a subspace of the data space \mathbb{R}^q , and the dimension of \mathcal{B} is a measure of the model's complexity. Linear static modeling is the problem of finding a low-complexity model (low-dimensional subspaces) that fits the data as close as possible. Existence of an exact linear model \mathcal{B} for the data \mathcal{D} , i.e., $\mathcal{D} \subset \mathcal{B}$, with complexity at most r is equivalent to the data matrix $D = [d_1 \ \dots \ d_N]$ having rank at most r . Therefore, measuring the fit between the data point d_i and the model \mathcal{B} by the orthogonal distance

$$\text{dist}(d_i, \mathcal{B}) = \min_{\hat{d}_i \in \mathcal{B}} \|d_i - \hat{d}_i\|_2,$$

the approximate fitting problem becomes a rank r matrix approximation problem (SLRA) with unstructured data matrix $\mathcal{S}(p) = D$.

Suppose now that some elements d_{ij} , $(i, j) \in \mathcal{J}_{\text{missing}}$ of the data matrix are missing. Equivalently, only the elements d_{ij} , $(i, j) \in \mathcal{J}_{\text{given}}$ of D are specified. The exact linear static modeling problem becomes a low-rank matrix completion problem [6]:

$$\text{find } \hat{D} \text{ such that } \text{rank}(\hat{D}) \leq r \text{ and } \hat{D}_{\mathcal{J}_{\text{given}}} = D_{\mathcal{J}_{\text{given}}}.$$

Here $D_{\mathcal{J}}$ denotes the vector of elements of D with indices in \mathcal{J} . In the context of approximate data fitting by linear static model and missing data values, the relevant problem is approximate low-rank matrix completion [5]:

$$\text{minimize over } \hat{D} \quad \|\hat{D}_{\mathcal{J}_{\text{given}}} - D_{\mathcal{J}_{\text{given}}}\|_2^2 \quad \text{subject to} \quad \text{rank}(\hat{D}) \leq r. \quad (\text{AMC})$$

The approximate low-rank matrix completion problem (AMC) is used for building recommender systems. In recommender system applications, there is a set of users and a set of products. Some users rate some products. The goal is to predict the user ratings on products that they have not rated. The underlying assumption that makes the solution of this problem possible is that the full user-ratings matrix is low-rank. The low-rank property is observed empirically and can be explained intuitively as existence of a small number of groups of users with the same “taste” (i.e., users that like or dislike the same products). In practice, the low-rank assumption is satisfied only approximately, which makes the approximation aspect of the problem essential.

The main issue in building real-life recommender systems is the high dimensionality and sparsity of the data matrix. Additional important issues in building practical recommender systems is the fact that the given and missing ratings are discrete and that apart from the users' ratings, there is demographic information about the users. Taking into account this prior information may improve significantly the accuracy of the missing values estimation. These issues, however, are outside the scope of the present paper.

1.2.2 System identification with missing data: An approximate block-Hankel structured low-rank matrix completion problem

A discrete-time linear time-invariant dynamical model is a set of time series

$$\mathcal{B}(R) := \{w \mid R_0 w(t) + R_1 w(t+1) + \dots + R_\ell w(t+\ell) = 0, \text{ for all } t\} \quad (\mathcal{B})$$

that satisfy a constant coefficients difference equation. The matrices $R_0, R_1, \dots, R_\ell \in \mathbb{R}^{p \times q}$ are parameters specifying the model. Note that the linear static model is a special case of a linear time-invariant dynamical model when the lag ℓ of the difference equation representing the system is equal to zero.

A finite time series

$$w = (w(1), \dots, w(T)), \quad \text{where } w(t) \in \mathbb{R}^q,$$

is an exact trajectory of the system defined in (\mathcal{B}) if the following matrix equation is satisfied

$$\underbrace{\begin{bmatrix} R_0 & R_1 & \cdots & R_\ell \end{bmatrix}}_R \begin{bmatrix} w(1) & w(2) & w(3) & \cdots & w(T-\ell) \\ w(2) & w(3) & \ddots & & w(T-\ell+1) \\ w(3) & \ddots & & & \vdots \\ \vdots & & & & \\ w(\ell+1) & w(\ell+2) & \cdots & & w(T) \end{bmatrix} = 0.$$

$\mathcal{H}_{\ell+1}(w)$

Without loss of generality, we assume that the parameter matrix R has full row rank p , which implies that

$$\text{rank}(\mathcal{H}_{\ell+1}(w)) \leq q(\ell+1) - p.$$

We showed above that the data w is an exact trajectory of a system $\mathcal{B}(R)$, if the block-Hankel matrix $\mathcal{H}_{\ell+1}(w)$ is rank deficient. Therefore, as in the static case, approximate modeling by a linear time-invariant system is a low-rank approximation problem

$$\text{minimize over } \hat{w} \quad \|w - \hat{w}\|_2^2 \quad \text{subject to} \quad \text{rank}(\mathcal{H}_{\ell+1}(\hat{w})) \leq q(\ell+1) - p. \quad (\text{SYSID})$$

Note, however, that the linear time-invariant model class imposes a block-Hankel structure constraint on the approximation matrix.

Identification from a trajectory with missing elements is therefore a block-Hankel structured low-rank matrix completion problem. A special system identification problems for the class of auto-regressive exogenous systems with missing data is considered in [10] and a method based on frequency domain techniques is proposed in [24]. These papers do not link the system identification problem to the block-Hankel low-rank approximation problem (SYSID), so that their approaches are different from ours. The method developed in this paper, when specialized to block-Hankel structure can be used for general multivariable system identification in the time domain.

1.2.3 Data-driven simulation and control

The trajectory w of a dynamical system can be partitioned into inputs u , *i.e.*, free variables, and outputs y , *i.e.*, variables that are determined by the inputs, initial conditions, and the model. Let $w = (u, y)$ be such a partition. The output $y = (y(1), \dots, y(T))$ of \mathcal{B} is uniquely determined by the input $u = (u(1), \dots, u(T))$ and the initial conditions

$$w_p = (w(-\ell+1), w(-\ell+1), \dots, w(0)).$$

This gives a “signal processor” interpretation of a dynamical system.

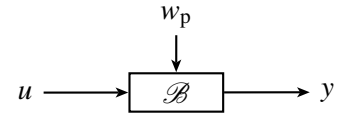
The simulation problem aims to find the output y_f of a system \mathcal{B} , corresponding to a given input u_f and initial conditions w_p , *i.e.*,

$$\text{find } y_f \text{ such that } w_p \wedge (u_f, y_f) \in \mathcal{B}.$$

($w_p \wedge w_f$ denotes the concatenation of w_p and w_f .) This is a classic problem in system theory and numerical linear algebra, for which many solutions exist, *e.g.*, for systems with no inputs, the problem is related to the computation of the matrix exponential [22]. The classical simulation methods require a representation (state space, transfer function, convolution kernel, *etc.*) of the model. Such a representation is often obtained from data by system identification. The question occurs of solving the simulation problem directly from the data without identifying a model representation as a byproduct and using it in a model based solution. We call the direct problem data-driven simulation [15].

The data-driven simulation problem is a mosaic-Hankel structured low-rank approximation problem with fixed (exact) and missing data. To see this, let w' denotes the data that implicitly specifies the model and $w'' = w_p'' \wedge (u_f'', y_f'')$ be the to-be-simulated trajectory. We express the condition that w' and w'' are trajectories of a linear time-invariant system with lag ℓ in matrix language as

$$\text{rank} \left(\begin{bmatrix} \mathcal{H}_{\ell+1}(w') & \mathcal{H}_{\ell+1}(w'') \end{bmatrix} \right) \leq q(\ell+1) - p.$$



This is a model-free description of the simulation problem—the existence of a model is implicit in the rank constraint. The fixed data are the initial conditions w_p'' and the input u_f'' and the missing data are the to-be computed response y_p'' . When the data w' is not an exact trajectory of the model, the matrix $\mathcal{H}_{\ell+1}(w')$ is generically full rank, so that an approximation is needed. The resulting data-driven simulation problem is a mosaic-Hankel structured low-rank approximation with missing data:

$$\begin{aligned} & \text{minimize} \quad \text{over } \hat{w}' \text{ and } \hat{w}'' \quad \|w' - \hat{w}'\|_2^2 \\ & \text{subject to} \quad \text{rank} \left(\begin{bmatrix} \mathcal{H}_{\ell+1}(\hat{w}') & \mathcal{H}_{\ell+1}(\hat{w}'') \end{bmatrix} \right) \leq q(\ell-1) - p, \\ & \quad \hat{w}_p'' = w_p'', \quad \text{and} \quad \hat{u}_f'' = u_f''. \end{aligned} \tag{DDSIM}$$

Notation

In the rest of the paper, we use the following notation.

- $A_{\mathcal{I}, \mathcal{J}}$ is the submatrix of A with rows in \mathcal{I} and columns in \mathcal{J} . The row/column index can be replaced by the symbol “:”, in which case all rows/columns are selected.
- $\mathcal{M} / \mathcal{G}$ is the vector of indices of p that are missing / given, and n_m / n_g is the number of missing / given elements.
- A^+ is the pseudo inverse of A and A^\perp is a matrix which rows form a basis for the left null space of A .
- $\text{vec}(\cdot)$ is the column-wise vectorization operator.

2 Analytical solution of the inner minimization problem

In this section, we consider the inner minimization problem (INNER).

Problem 1. Given affine structure \mathcal{S} , structure parameter vector $p \in \{\mathbb{R} \cup \{\text{NaN}\}\}^{n_p}$, and a kernel parameter $R \in \mathbb{R}^{(m-r) \times m}$, evaluate the cost function $M(R)$, defined in (INNER), and find a point \hat{p} that attains the minimum.

For a given structure \mathcal{S} and $R \in \mathbb{R}^{(m-r) \times m}$, we define the matrix

$$G := [\text{vec}(RS_1) \quad \cdots \quad \text{vec}(RS_{n_p})] \in \mathbb{R}^{(m-r)n \times n_p}. \tag{G}$$

Theorem 2. Under the following assumptions:

1. $G_{:, \mathcal{M}}$ is full column rank,
2. $1 \leq (m-r)n - n_m \leq n_g$, and
3. $\bar{G} := G_{:, \mathcal{M}}^\perp G_{:, \mathcal{G}}$ is full row rank,

Problem 1 has a unique global minimum

$$\hat{p}_{\mathcal{G}} = p_{\mathcal{G}} - \bar{G}^\top (\bar{G} \bar{G}^\top)^{-1} s \quad \text{and} \quad \hat{p}_{\mathcal{M}} = -G_{:, \mathcal{M}}^+ (\text{vec}(RS_0) + G_{:, \mathcal{G}} \hat{p}_{\mathcal{G}}), \tag{\hat{p}}$$

with objective function value

$$M(R) = s^\top (\bar{G} \bar{G}^\top)^{-1} s, \quad \text{where} \quad s := \bar{G} p_{\mathcal{G}} + G_{:, \mathcal{M}}^\perp \text{vec}(RS_0). \tag{M}$$

Proof. Defining

$$\Delta p_{\mathcal{G}} := p_{\mathcal{G}} - \hat{p}_{\mathcal{G}}$$

and using the identity

$$R\mathcal{S}(\hat{p}) = 0 \quad \Longleftrightarrow \quad G\hat{p} = -\text{vec}(RS_0),$$

we have

$$R\mathcal{S}(\hat{p}) = 0 \quad \Longleftrightarrow \quad \begin{bmatrix} G_{:,g} & G_{:,m} \end{bmatrix} \begin{bmatrix} p_g - \Delta p_g \\ \hat{p}_m \end{bmatrix} = -\text{vec}(RS_0).$$

Therefore, (INNER) is equivalent to

$$M(R) := \min_{\Delta p_g \in \mathbb{R}^{n_g}, \hat{p}_m \in \mathbb{R}^{n_m}} \|\Delta p_g\|_2^2 \quad \text{subject to} \quad \begin{bmatrix} G_{:,g} & G_{:,m} \end{bmatrix} \begin{bmatrix} \Delta p_g \\ -\hat{p}_m \end{bmatrix} = G_{:,g} p_g + \text{vec}(RS_0),$$

which is a generalized linear least norm problem. The solution follows from Lemma 3. \square

Generalized least norm problem

Lemma 3. *Consider the generalized linear least norm problem*

$$f = \min_{x,y} \|x\|_2^2 \quad \text{subject to} \quad Ax + By = c, \quad (\text{GLN})$$

with $A \in \mathbb{R}^{m \times n_x}$, $B \in \mathbb{R}^{m \times n_y}$, and $c \in \mathbb{R}^m$. Under the following assumptions:

1. B is full column rank,
2. $1 \leq m - n_y \leq n_x$, and
3. $\bar{A} := B^\perp A$ is full row rank,

problem (GLN) has a unique solution

$$\begin{aligned} f &= c^\top (B^\perp)^\top (\bar{A} \bar{A}^\top)^{-1} B^\perp c, \\ x &= \bar{A}^\top (\bar{A} \bar{A}^\top)^{-1} B^\perp c \quad \text{and} \quad y = B^+(c - Ax). \end{aligned} \quad (\text{SOL})$$

Proof. Under assumption 1, B has a nontrivial left kernel of dimension $m - n_y$. Therefore for the nonsingular matrix

$$T = \begin{bmatrix} B^+ \\ B^\perp \end{bmatrix} \in \mathbb{R}^{m \times m}$$

$$TB = \begin{bmatrix} B^+ \\ B^\perp \end{bmatrix} B = \begin{bmatrix} B^+ B \\ B^\perp B \end{bmatrix} = \begin{bmatrix} I_{n_y} \\ 0 \end{bmatrix}.$$

Pre-multiplying both sides of the constraint of (GLN) by T , we have the following equivalent constraint

$$\begin{bmatrix} B^+ Ax \\ B^\perp Ax \end{bmatrix} + \begin{bmatrix} y \\ 0 \end{bmatrix} = \begin{bmatrix} B^+ c \\ B^\perp c \end{bmatrix}.$$

The first equation

$$y = B^+(c - Ax)$$

uniquely determines y , given x . The second equation

$$B^\perp Ax = B^\perp c \quad (*)$$

defines a linear constraint for x only. By assumption 2, it is an underdetermined system of linear equations. Therefore, (GLN) is equivalent to the following standard least norm problem

$$f = \min_x \|x\|_2^2 \quad \text{subject to} \quad B^\perp Ax = B^\perp c. \quad (\text{GLN}')$$

By assumption 3 the solution is unique and is given by (SOL). \square

Note 4 (About assumptions 1–3). Assumption 1 is a necessary condition for uniqueness of the solution. Relaxing assumption 1 implies that any vector in the affine space

$$\mathcal{Y} = B^+(c - Ax) + \text{null}(B)$$

is a solution to (GLN). Assumption 2 ensures that the problem is a least norm problem and has a nontrivial solution. In the case $m = n_y$, the problem has a trivial solution $f = 0$. In the case $m - n_y > n_x$, the problem generically has no solution because the constraint (*) is an overdetermined system of equations. Assumption 3 is also required for uniqueness of the solution. It can also be relaxed, making y nonunique.

Note 5 (Link to weighted least norm problems with singular weight matrix). Consider the weighted least norm problem

$$\min_z z^\top W z \quad \text{subject to} \quad Dz = c,$$

with singular positive semidefinite weight matrix W . Using a change of variables $\bar{z} = T^{-1}z$, where T is a nonsingular matrix, we obtain the equivalent problem

$$\min_{\bar{z}} \bar{z}^\top T^\top W T \bar{z} \quad \text{subject to} \quad DT \bar{z} = c.$$

There exists a nonsingular matrix T , such that

$$T^\top W T = \begin{bmatrix} I_{n_x} & \\ & 0 \end{bmatrix}.$$

Partitioning \bar{z} and $\bar{D} := DT^{-1}$ conformably as

$$\bar{z} = \begin{bmatrix} x \\ y \end{bmatrix} \quad \text{and} \quad \bar{D} = \begin{bmatrix} A & B \end{bmatrix}$$

we obtain problem (GLN).

3 Outer minimization problem

The outer minimization problem (SLRA_R) is a nonlinear least-squares problem, which we solve by general purpose local optimization methods. In order to apply standard optimization methods, however, we need first to replace the rank constraint with equivalent equality or inequality constraints.

The full row rank constraint on R is equivalent to and can be enforced in the parameter optimization method by the equality constraint

$$RR^\top = I_{m-r}. \quad (\text{f.r.r. } R)$$

Then, the outer minimization problem becomes a constrained nonlinear least squares problems

$$\text{minimize over } R \in \mathbb{R}^{(m-r) \times m} \quad M(R) \quad \text{subject to} \quad RR^\top - I_{m-r} = 0. \quad (\text{SLRA}'_R)$$

(SLRA'_R) is an optimization problem on a Stiefel manifold [1] and can be solved by specialized methods, *e.g.*, the GenRTR package [2], or by general purpose penalty methods for constrained optimization [23]. Next, we consider a penalty method, *i.e.*, reformulation of (SLRA'_R) as a regularized unconstrained nonlinear least squares problem by adding the regularization term $\gamma \|RR^\top - I_{m-r}\|_F^2$, where $\|\cdot\|_F$ is the Frobenius norm, to the cost function

$$\text{minimize over } R \in \mathbb{R}^{(m-r) \times m} \quad M(R) + \gamma \|RR^\top - I_{m-r}\|_F^2. \quad (\text{SLRA}''_R)$$

The parameter γ should be chosen “large enough” in order to enforce the constraint (f.r.r. R). A corollary of the following theorem shows that $\gamma = \|p_{\mathcal{G}}\|_2^2$ is sufficiently large for linear structures.

Theorem 6. Let $M : \mathbb{R}^{(m-r) \times m} \rightarrow \mathbb{R}_+$ be a homogeneous function, *i.e.*, $M(R) = M(TR)$, for any R and a nonsingular $m \times m$ matrix T . Assume that γ satisfies

$$\gamma > \min_{\{R \mid \text{rank}(R)=m-r\}} M(R). \quad (\gamma)$$

Then, the optimal solutions of problem (SLRA''_R) coincide with the optimal solutions of (SLRA'_R).

Proof. We call a set $\mathcal{R} \subset \mathbb{R}^{d \times m}$ a “homogeneous set” if for all $R \in \mathcal{R}$ and for all nonsingular matrices $T \in \mathbb{R}^{d \times d}$, $TR \in \mathcal{R}$. Let R be a solution to (SLRA_R'') with the constraint $R \in \mathcal{R}$, where \mathcal{R} is a homogeneous set. We will show that

$$\|RR^\top - I_{m-r}\|_F^2 = m - r - \text{rank}(R). \quad (*)$$

There exists an orthogonal matrix U diagonalizing RR^\top . We have,

$$\begin{aligned} \|RR^\top - I_{m-r}\|_F^2 &= \|URR^\top U^\top - I_{m-r}\|_F^2 \\ &= \|\text{diag}(a_1, \dots, a_{\text{rank}(R)}, 0, \dots, 0) - I_{m-r}\|_F^2, \quad \text{where } a_i > 0 \\ &= \sum_{i=1}^{\text{rank}(R)} (a_i - 1)^2 + m - r - \text{rank}(R). \end{aligned}$$

Suppose that $a_i \neq 1$ for some i . The matrix

$$R' = \text{diag}(1, \dots, 1, 1/\sqrt{a_i}, 1, \dots, 1)R$$

has the same row span and rank as R , so that by the homogeneity property of M , $M(R) = M(R')$. However, we have

$$\|RR^\top - I_{m-r}\|_F^2 > \|R'R'^\top - I_{m-r}\|_F^2,$$

so that $R' \in \mathcal{R}$ achieves smaller value of the cost function of (SLRA_R'') than R . This is a contradiction. Therefore, $a_i = 1$ for all i . This concludes the proof of $(*)$.

So far we showed that minimization of the cost function in (SLRA_R'') on homogeneous sets is equivalent to minimization of

$$M(R) + \gamma(m - r - \text{rank}(R)). \quad (M'')$$

The set of full row rank matrices

$$\mathcal{R}_f := \{R \in \mathbb{R}^{(m-r) \times m} \mid \text{rank}(R) = m - r\}$$

and the set of rank-deficient matrices

$$\mathcal{R}_d := \{R \in \mathbb{R}^{(m-r) \times m} \mid \text{rank}(R) < m - r\}$$

are homogeneous. Denote the solutions of (SLRA_R'') on these sets as

$$\begin{aligned} M_f^* &:= \inf_{R \in \mathcal{R}_f} M(R) + \gamma \|RR^\top - I_{m-r}\|_F^2 \stackrel{(*)}{=} \inf_{R \in \mathcal{R}_f} M(R) < \gamma, \\ M_d^* &:= \inf_{R \in \mathcal{R}_d} M(R) + \gamma \|RR^\top - I_{m-r}\|_F^2 \stackrel{(*)}{=} \inf_{R \in \mathcal{R}_d} \underbrace{M(R)}_{\geq 0} + \underbrace{\gamma(m - r - \text{rank}(R))}_{\geq \gamma}. \end{aligned}$$

Then, $M_f^* < \gamma \leq M_d^*$ and

$$M^* := \inf_{R \in \mathbb{R}^{(m-r) \times m}} M(R) + \gamma \|RR^\top - I_{m-r}\|_F^2 = M_f^*.$$

In addition, the minimum of (SLRA_R') coincides with M_f^* by the homogeneity of M . Therefore, the solutions of (SLRA_R'') and (SLRA_R') coincide if one of them exists. \square

Note 7 (Choice of γ). $\gamma = \max_{R \in \mathcal{R}_f} M(R)$ always satisfies condition (γ) . In particular, for a linear structure \mathcal{S} , it is sufficient to take $\gamma = \|p_{\mathcal{G}}\|_2^2$, because $\mathcal{S}(0)$ has zero rank, and $R\mathcal{S}(0) = 0$ holds for any R .

Note 8 (Initial approximation). Solving the outer minimization problem by local minimization requires an initial approximation for the parameter R , i.e., a suboptimal solution of the structured low-rank approximation problem. Such a solution can be computed from a heuristic that ignores the data matrix structure \mathcal{S} and fills in the missing values with initial estimates. Rigorous analysis of the missing values imputation question is done in [11]. Theorem 1.1 of [11] gives theoretical justification for the zero imputation in the case of unstructured \mathcal{S} . The resulting unstructured low-rank approximation problem can then be solved analytically in terms of the singular value decomposition.

Note 9 (Efficient computation and software implementation). Efficient evaluation of the cost function and its derivatives in the special case of mosaic-Hankel matrix structure is presented in a companion paper [26]. The method, presented in this paper (general affine structure) and the efficient methods of [26] are implemented in Matlab (using Optimization Toolbox) and in C++ (using by the Levenberg-Marquardt algorithm [20] from the GNU Scientific Library [7]), respectively. Description of the software and overview of its applications is given in [16].

4 Weighted approximation

Problem (SLRA) is generalized in this section to the weighted structured low-rank approximation problem

$$\begin{aligned} & \text{minimize} && \text{over } \hat{p} \in \mathbb{R}^{n_p} && (p_{\mathcal{G}} - \hat{p}_{\mathcal{G}})^\top W_g (p_{\mathcal{G}} - \hat{p}_{\mathcal{G}}) \\ & \text{subject to} && \text{rank}(\mathcal{S}(\hat{p})) \leq r, \end{aligned} \quad (\text{WSLRA})$$

where W_g is a positive definite matrix.

In case of a diagonal weight matrix

$$W_g = \text{diag}(w_g) = \text{diag}(w_1, \dots, w_{n_g}), \quad (w_g)$$

the weights can be specified by a positive vector w_g .

The change of variables

$$p'_{\mathcal{G}} = W_g^{1/2} p_{\mathcal{G}} \quad \text{and} \quad \hat{p}'_{\mathcal{G}} = W_g^{1/2} \hat{p}_{\mathcal{G}} \quad (p \mapsto p')$$

reduces Problem (WSLRA) to an equivalent unweighted problem (SLRA). We have

$$\mathcal{S}(\hat{p}) = S_0 + \text{vec}^{-1}(\mathbf{S}\hat{p}), \quad \text{where } \mathbf{S} := [\text{vec}(S_1) \ \cdots \ \text{vec}(S_{n_p})] \in \mathbb{R}^{mn \times n_p}. \quad (\mathbf{S})$$

The structure \mathcal{S}' of the equivalent problem is defined by the matrices S_0 and

$$\mathbf{S}' = [\text{vec}(S'_1) \ \cdots \ \text{vec}(S'_{n_p})], \quad \text{where } \mathbf{S}'_{:, \mathcal{G}} = \mathbf{S}_{:, \mathcal{G}} W_g^{-1/2} \text{ and } \mathbf{S}'_{:, \mathcal{M}} = \mathbf{S}_{:, \mathcal{M}}. \quad (\mathcal{S} \mapsto \mathcal{S}')$$

We showed that problem (WSLRA) is solved by:

1. preprocessing the data p and the structure \mathcal{S} , as in $(p \mapsto p')$ and $(\mathcal{S} \mapsto \mathcal{S}')$,
2. solving the equivalent unweighted problem with structure parameter vector p' , structure specification \mathcal{S}' , and rank specification r , and
3. postprocessing the solution \hat{p}' , obtained in step 2, in order to obtain the solution $\hat{p}_{\mathcal{G}} = W_g^{-1/2} \hat{p}'_{\mathcal{G}}$ of the original problem.

Using the transformation $(p \mapsto p')$, $(\mathcal{S} \mapsto \mathcal{S}')$ and the solution (M) of (SLRA), we obtain the following explicit expression for the cost function of (WSLRA)

$$M(R) = (\bar{G} p_{\mathcal{G}} - G_{:, \mathcal{M}}^\perp \text{vec}(RS_0))^\top W_g^{-1} \bar{G}^\top (\bar{G} W_g^{-1} \bar{G}^\top)^{-1} \bar{G} W_g^{-1} (\bar{G} p_{\mathcal{G}} - G_{:, \mathcal{M}}^\perp \text{vec}(RS_0)), \quad (\mathbf{M}_W)$$

where $\bar{G} = G_{:, \mathcal{M}}^\perp G_{:, \mathcal{G}}$ and G is defined in (G).

Note 10 (Specification of fixed parameter values by infinite weights). In the case of a diagonal weight matrix (w_g) , an infinite weight $w_j = \infty$ specifies a fixed parameter value $\hat{p}_j = p_j$. A problem with infinite weights is equivalent to a regular structured low-rank approximation problem with fixed parameters assigned to the constant term S_0 of the structure specification. Let \mathcal{J}_f be the set of indices of the fixed structure parameters and $\overline{\mathcal{J}}_f$ its complement

$$\mathcal{J}_f = \{j \in \{1, \dots, n_p\} \mid \hat{p}_j = p_j\} \quad \text{and} \quad \overline{\mathcal{J}}_f = \{j \in \{1, \dots, n_p\} \mid j \notin \mathcal{J}_f\}.$$

The equivalent problem has structure, defined by

$$\mathcal{S}'(\hat{p}') = S_0 + \underbrace{\sum_{i \in \mathcal{J}_f} S_i p_i}_{S'_0} + \sum_{i \in \overline{\mathcal{J}}_f} S_i \hat{p}_i, \quad \text{where } \hat{p}' := \hat{p}|_{\overline{\mathcal{J}}_f}.$$

The estimated vector \hat{p} is recovered from the parameter vector \hat{p}' of the equivalent problem by

$$\hat{p}|_{\overline{\mathcal{J}}_f} = \hat{p}' \quad \text{and} \quad \hat{p}|_{\mathcal{J}_f} = p|_{\mathcal{J}_f}.$$

Note 11 (Solving (SLRA) as weighted unstructured problem). Consider an instance of problem (SLRA), referred to as problem P1, with structure $\mathcal{S} = \mathcal{S}_1$ and an instance of problem (WSLRA), referred to as problem P2, with unstructured correction ($\mathcal{S}_2 = \text{vec}^{-1}$, $n_{p_2} = mn$) and weight matrix

$$W_2^{-1} = \mathbf{S}_1 \mathbf{S}_1^\top. \quad (\mathcal{S}_1 \mapsto W_2)$$

It can be verified by inspection that the cost functions (M) and (M_W) of problems P1 and P2, respectively, coincide. The weight matrix $W_2 \in \mathbb{R}^{mn \times mn}$, defined in $(\mathcal{S}_1 \mapsto W_2)$, however is singular ($\text{rank}(W_2)$ is equal to the number of structure parameters of problem P1, which is less than mn). In the derivation of the cost function (M_W) it is assumed that W_g is positive definite, so that minimization of (M_W) is not equivalent to problem P2.

5 Numerical examples

In this section, we present numerical examples with the three problems covered in the introduction:

- unstructured noisy matrix completion,
- system identification with missing data, and
- data-driven simulation.

The correctness of the results and the effectiveness of the methods in the paper is validated by comparison with alternative methods, specifically developed for these applications. All simulations are done in Matlab and are reproducible in the sense of [4]. An extended version [17] of this paper is a literate program (in noweb format [25]), implementing the methods in the paper and generating the presented numerical results. The necessary m-files can be downloaded from

<http://homepages.vub.ac.be/~imarkovs/publications.html>

5.1 Approximate matrix completion

In the approximate matrix completion problem (AMC) of Section 1.2.1, the methods in the paper are compared with the following alternative methods:

- `wlra` — the alternating projections method of [13],
- `optspace` — a method based on spectral techniques and manifold optimization [11],
- `lmafit` — the successive over-relaxation algorithm of [27], and
- `rttrmc` — the Riemannian trust-region method of [3].

These methods use an image representation

$$\hat{D} = \mathcal{S}(\hat{p}) = PL, \quad \text{where } P \in \mathbb{R}^{m \times r} \text{ and } L \in \mathbb{R}^{r \times n},$$

of the unstructured $m \times n$ rank r matrix \hat{D} . The number of optimization variables in the image representation is $r(m+n)$ (rm in the case of `rttrmc`), so that the methods are suitable for problems with small rank r . The above cited methods are developed for recommender system applications, where the data matrix is sparse in the given elements and this sparsity is effectively used in the computations.

In contrast, the kernel representation (KER) has $m(m-r)$ variables, so that it is suitable for problems with small co-rank $m-r$. The implementation [17] of the methods in the paper is applicable only for small size problems (say, $m < 10$, $n < 100$, and $m-r < 3$). For larger problems, the efficient C implementation [16] (denoted by `slra-c` below), of the variable projection approach can be used by setting small values for the weights corresponding to the missing values. In the reported results, `slra-m` corresponds to (SLRA'_R) and `slra-r` corresponds to (SLRA''_R) with $\gamma = \|p_{\mathcal{G}}\|_2^2$ (see Note 7). `slra-c` is applied by setting the weights of the missing values to 10^{-6} .

The data matrices, used in the simulation examples, are generated as random rank r matrices $\mathcal{S}(\bar{p})$ plus noise, where the noise matrix is zero mean Gaussian with independent identically distributed elements. A fraction of randomly selected elements of the data matrix are missing. The simulation parameters for the experiments are:

- matrix dimensions and rank, defined by variables m , n , and r , respectively;
- noise level, defined by a variable nl , and
- fraction of given elements, defined by a variable eps .

The relative approximation errors

$$e_g = \frac{\|p_g - \hat{p}_g\|_2}{\|p_g\|_2} \quad \text{and} \quad e_m = \frac{\|\bar{p}_{\mathcal{M}} - \hat{\bar{p}}_{\mathcal{M}}\|_2}{\|\bar{p}_{\mathcal{M}}\|_2}$$

(rows eg and em) and execution time (row t) are shown below for the methods compared in three simulation problems. In a problem with exact 10×100 matrix of rank 8 with 8 missing values (exact matrix completion problem):

11a $\langle \text{AMC, example 1 11a} \rangle \equiv$ (?0—1)

`ex = 'amc_ex1'; m = 10; n = 100; r = 8; nl = 0; eps = 0.995; test_amc`
the methods in the paper, `wlra`, and `rtrmc` recover exactly the missing values:

	'slra-m'	'slra-r'	'slra-c'	'wlra'	'optspace'	'lmafit'	'rtrmc'
'eg'	[2e-16]	[2e-16]	[1.2410e-04]	[5e-11]	[0.0213]	[2.5156e-04]	[2e-12]
'em'	[1e-14]	[1e-14]	[0.0013]	[3e-09]	[0.9204]	[0.1923]	[2e-10]
't'	[1.3632]	[3.5105]	[0.2710]	[0.1614]	[0.2992]	[0.0440]	[0.2823]

For the same setup with 10% noise (approximate matrix completion):

11b $\langle \text{AMC, example 2 11b} \rangle \equiv$ (?0—1)

`ex = 'amc_ex2'; m = 10; n = 100; r = 8; nl = 0.1; eps = .995; test_amc`
the methods in the paper, `wlra`, and `rtrmc` obtain the same approximation error (around 8% for the missing values).

	'slra-m'	'slra-r'	'slra-c'	'wlra'	'optspace'	'lmafit'	'rtrmc'
'eg'	[0.0117]	[0.0117]	[0.0154]	[0.0117]	[0.0231]	[0.0118]	[0.0117]
'em'	[0.0796]	[0.0797]	[0.1439]	[0.0803]	[0.9376]	[0.1233]	[0.0796]
't'	[10.0445]	[36.9907]	[0.2400]	[0.0968]	[0.3364]	[0.0465]	[0.2638]

For a problem with 10×1000 matrix of rank 9 with 101 missing values and 5% noise

11c $\langle \text{AMC, example 3 11c} \rangle \equiv$ (?0—1)

`ex = 'amc_ex3'; m = 10; n = 1000; r = 9; nl = 0.05; eps = .99; test_amc`
`slra-m` and `slra-r` are not applicable (they require too much time and memory). However, the C implementation `slra-c` is competitive with the alternative methods in terms of execution time and obtains the same approximation error as `rtrmc`. In this example, `lmafit` achieves the smallest error and is also the fastest of all compared methods.

	'slra-m'	'slra-r'	'slra-c'	'wlra'	'optspace'	'lmafit'	'rtrmc'
'eg'	[NaN]	[NaN]	[0.0041]	[0.0062]	[0.0209]	[0.0041]	[0.0041]
'em'	[NaN]	[NaN]	[0.2633]	[86.1488]	[0.9021]	[0.2116]	[0.2634]
't'	[NaN]	[NaN]	[0.0732]	[1.5501]	[1.5773]	[0.1047]	[0.2978]

5.2 System identification with missing data

Consider the system identification problem (SYSID), described in Section 1.2.2. The data w is a noisy $T = 100$ samples long random trajectory of a single-input single-output linear time-invariant system $\mathcal{B}(\bar{R})$ with lag $\ell = 2$. Samples $w(t)$, for $t \in \mathcal{T}_m$, are missing. The noise is zero mean white Gaussian process with covariance matrix $\sigma^2 I_q$, i.e., both the inputs $u = w_1$ and the outputs $y = w_2$ are perturbed and the input and the output noise variances are equal.

The true model parameters

$$\bar{R}_0 = [-1 \quad 0.81], \quad \bar{R}_1 = [1 \quad -1.456], \quad \bar{R}_2 = [\bar{Q}_2 \quad \bar{P}_2] = [-1 \quad 1] \quad (*)$$

are normalized with $\bar{P}_2 = 1$ and the same normalization is used for the identified model parameter \hat{R} . This ensures that the parameters are unique and the systems $\mathcal{B}(\bar{R})$ and $\mathcal{B}(\hat{R})$ can be compared by the relative parameter error

$$e_R = \frac{\|\bar{R} - \hat{R}\|_2}{\|\bar{R}\|_2}.$$

The identification problem is solved by the methods `slra-m` and `slra-r`, developed in the paper; the efficient software `slra-c` of [16] (with zero weights substituted by small constants); and the method `sysid` of [24].

The simulation parameters in the experiments are the

- number of samples T ,
- set of missing values \mathcal{T}_m , specified by a variable T_m , and
- noise variance interval, specified by a vector NL .

The reported results show the estimation error e_R for the compared methods and for the different noise levels specified in NL . In the case of uniformly distributed missing data samples:

12a $\langle \text{SYSID example 1 12a} \rangle \equiv$ (? 0—1)
`ex = 'sysid_ex1'; T = 100; N = 8; NL = linspace(0, 0.1, N); Tm = 30:3:70; test_sysid`
 the developed methods perform slightly better than the alternative method for small noise level and slightly worse for high noise level:

'nl'	[0]	[0.0143]	[0.0286]	[0.0429]	[0.0571]	[0.0714]	[0.0857]	[0.1000]
'slra-m'	[1.5624e-15]	[0.0085]	[0.0134]	[0.0390]	[0.0325]	[0.0693]	[0.0347]	[0.0778]
'slra-r'	[1.5624e-15]	[0.0085]	[0.0134]	[0.0390]	[0.0325]	[0.0693]	[0.0347]	[0.0778]
'slra-c'	[2.9202e-06]	[0.0085]	[0.0134]	[0.0390]	[0.0325]	[0.0693]	[0.0347]	[0.0778]
'sysid'	[5.6315e-15]	[0.0096]	[0.0148]	[0.0434]	[0.0356]	[0.0697]	[0.0322]	[0.0758]

Similar results are obtained in the case of consecutive missing samples:

12b $\langle \text{SYSID example 2 12b} \rangle \equiv$ (? 0—1)
`ex = 'sysid_ex2'; T = 100; N = 8; NL = linspace(0, 0.1, N); Tm = 40:60; test_sysid`

'nl'	[0]	[0.0143]	[0.0286]	[0.0429]	[0.0571]	[0.0714]	[0.0857]	[0.1000]
'slra-m'	[1.6123e-15]	[0.0063]	[0.0142]	[0.0366]	[0.0529]	[0.0738]	[0.0410]	[0.0850]
'slra-r'	[1.6123e-15]	[0.0063]	[0.0142]	[0.0366]	[0.0529]	[0.0738]	[0.0411]	[0.0851]
'slra-c'	[6.5227e-07]	[0.0063]	[0.0142]	[0.0366]	[0.0529]	[0.0738]	[0.0411]	[0.0851]
'sysid'	[3.3311e-16]	[0.0080]	[0.0191]	[0.0431]	[0.0578]	[0.0766]	[0.0417]	[0.0829]

This latter problem can be solved also by treating the data as two independent trajectories without missing data.

5.3 Data-driven simulation

Consider the data-driven simulation problem (DDSIM), described in Section 1.2.3. The to-be-simulated system is $\bar{\mathcal{B}} = \mathcal{B}(\bar{R})$, with parameter \bar{R} given in (*) and with an input/output partition $w = (u, y)$ of the variables. The given trajectory $w' = (u', y') \in (\mathbb{R}^2)^{T'}$ is a noise corrupted random trajectory of $\bar{\mathcal{B}}$ (the same simulation setup as in Section 5.2) and the to-be-simulated trajectory $w'' = (u'', y'') \in (\mathbb{R}^2)^{T''}$ is the impulse response \bar{h} of $\bar{\mathcal{B}}$, i.e., the response of $\bar{\mathcal{B}}$ to pulse input under zero initial conditions:

$$u'' = (\underbrace{0, \dots, 0}_{\ell}, \underbrace{1, 0, \dots, 0}_{\text{pulse input}}) \quad \text{and} \quad y'' = (\underbrace{0, \dots, 0}_{\ell}, \underbrace{\hat{h}(0), \hat{h}(1), \dots, \hat{h}(T_2 - \ell - 1)}_{\text{impulse response}}).$$

The methods in the paper are compared with an alternative subspace-type method `ddsim` [19, 15] in terms of the relative approximation error

$$e_h = \frac{\|\bar{h} - \hat{h}\|_2}{\|\bar{h}\|_2}.$$

Subspace methods are multi stage methods, i.e., they split the nonconvex optimization problem in several steps that are individually convex, but do not guarantee global or local optimality with respect to the overall problem. In general, subspace-type methods are more efficient but less accurate than local optimization-based methods. The results for $T' = 30$, $T'' = 52$, and noise level in the range 0–40%:

12c $\langle \text{DDSIM example 12c} \rangle \equiv$ (? 0—1)
`T1 = 30; T2 = 52; N = 6; NL = linspace(0, 0.1, N); test_ddsim`
 confirm this rule of thumb:

'nl'	[0]	[0.0200]	[0.0400]	[0.0600]	[0.0800]	[0.1000]
'slra-m'	[1.8920e-15]	[0.0451]	[0.1149]	[0.1652]	[0.2678]	[0.2301]
'slra-r'	[1.8920e-15]	[0.0451]	[0.1150]	[0.1652]	[0.2679]	[0.2298]
'slra-c'	[1.0123e-05]	[0.0451]	[0.1149]	[0.1652]	[0.2679]	[0.2298]
'ddsim'	[3.2215e-15]	[0.0572]	[0.1727]	[0.2772]	[0.3012]	[0.5311]

The true impulse response and the approximations for the 40% noise run are plotted in Figure 1.

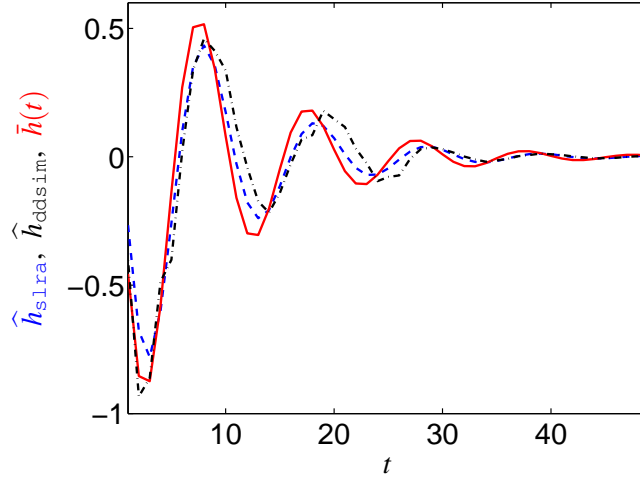


Figure 1: Data-driven simulation of impulse response: true impulse response — solid line, approximation by the methods in the paper — dashed line, approximation by the methods of [15] — dashed-dotted line.

6 Conclusions

A variable-projection-like method for structured low-rank approximation with missing data was developed. The approach was furthermore generalized to weighted structured low-rank approximation with missing values. After elimination of the approximation \hat{p} , the remaining nonlinear least-squares problem subject to quadratic equality constraints was solved as an equivalent regularized unconstrained optimization problem.

The problem and solution methods developed have applications in matrix completion (unstructured problems), system identification with missing data, and data-driven simulation and control (mosaic-Hankel structured problems). The performance of the methods in the paper was illustrated on small-size simulation examples and was compared with the performance of problem specific methods. Efficient computation for large scale problems appearing in applications such as recommender systems and system identification is a topic of future research.

Acknowledgements

The implementation of the method `sysid` for system identification with missing data, used in the numerical examples of Section 5.2, was kindly provided to us by Rik Pintelon. The research leading to these results has received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007-2013) / ERC Grant agreement number 258581 “Structured low-rank approximation: Theory, algorithms, and applications”.

References

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, Princeton, NJ, 2008.
- [2] C. Baker, P.-A. Absil, and K. Gallivan. GenRTR Riemannian optimization package. <http://www.math.fsu.edu/~cbaker/GenRTR>.
- [3] N. Boumal and P.-A. Absil. RTRMC: A Riemannian trust-region method for low-rank matrix completion. In J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F.C.N. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24 (NIPS)*, pages 406–414. 2011.
- [4] J. Buckheit and D. Donoho. *Wavelets and statistics*, chapter "Wavelab and reproducible research". Springer-Verlag, Berlin, New York, 1995.
- [5] E. Candes and Y. Plan. Matrix completion with noise. *arXiv:0903.3131*, March 2009.

- [6] E. Candés and B. Recht. Exact matrix completion via convex optimization. *Found. of Comput. Math.*, 9:717–772, 2009.
- [7] M. Galassi et al. *GNU Scientific Library Reference Manual*. 3rd edition.
- [8] G. Golub and V. Pereyra. Separable nonlinear least squares: the variable projection method and its applications. *Institute of Physics, Inverse Problems*, 19:1–26, 2003.
- [9] G. Heinig. Generalized inverses of Hankel and Toeplitz mosaic matrices. *Linear Algebra Appl.*, 216(0):43–59, February 1995.
- [10] A. Isaksson. Identification of ARX-models subject to missing data. *IEEE Trans. Automat. Control*, 38(5):813–819, May 1993.
- [11] R. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. *J. Mach. Learn. Res.*, 11:2057–2078, August 2010.
- [12] I. Markovsky. Structured low-rank approximation and its applications. *Automatica*, 44(4):891–909, 2008.
- [13] I. Markovsky. *Algorithms and iterate programs for weighted low-rank approximation with missing data*, volume 3 of *Springer Proc. Mathematics*, pages 255–273. Springer, 2011.
- [14] I. Markovsky. *Low Rank Approximation: Algorithms, Implementation, Applications*. Springer, 2012.
- [15] I. Markovsky and P. Rapisarda. Data-driven simulation and control. *Int. J. Control*, 81(12):1946–1959, 2008.
- [16] I. Markovsky and K. Usevich. Software for weighted structured low-rank approximation. Technical Report 339974, Univ. of Southampton, <http://eprints.soton.ac.uk/339974>, 2012.
- [17] I. Markovsky and K. Usevich. Structured low-rank approximation with missing values. Technical report, Vrije Univ. Brussel, homepages.vub.ac.be/~imarkovs/publications.html, 2012.
- [18] I. Markovsky, S. Van Huffel, and R. Pintelon. Block-Toeplitz/Hankel structured total least squares. *SIAM J. Matrix Anal. Appl.*, 26(4):1083–1099, 2005.
- [19] I. Markovsky, J. C. Willems, P. Rapisarda, and B. De Moor. Algorithms for deterministic balanced subspace identification. *Automatica*, 41(5):755–766, 2005.
- [20] D. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *SIAM J. Appl. Math.*, 11:431–441, 1963.
- [21] C. Moler and Ch. Van Loan. Nineteen dubious ways to compute the exponential of a matrix. *SIAM Review*, 20(4):801–836, 1978.
- [22] J. Nocedal and S. Wright. *Numerical optimization*. Springer-Verlag, 1999.
- [23] R. Pintelon and J. Schoukens. Frequency domain system identification with missing data. *IEEE Trans. Automat. Control*, 45(2):364–369, February 2000.
- [24] N. Ramsey. Iterate programming simplified. *IEEE Software*, 11:97–105, 1994.
- [25] K. Usevich and I. Markovsky. Variable projection for affinely structured low-rank approximation in weighted 2-norm, 2012. Available from <http://arxiv.org/abs/1211.3938>.
- [26] Z. Wen, W. Yin, and Y. Zhang. Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. Technical report, Rice University, 2010. CAAM Technical Report TR10–07.