

Structured low-rank approximation as optimization on a Grassmann manifold

Konstantin Usevich and Ivan Markovsky
Department ELEC, Vrije Universiteit Brussel

Leuven, 10 July 2013

International Workshop on Advances in Regularization, Optimization, Kernel Methods
and Support Vector Machines: theory and applications

Structured low-rank approximation problem

$$\begin{bmatrix} a & b & c & d \\ b & c & d & e \\ c & d & e & f \end{bmatrix} \approx \begin{bmatrix} \hat{a} & \hat{b} & \hat{c} & \hat{d} \\ \hat{b} & \hat{c} & \hat{d} & \hat{e} \\ \hat{c} & \hat{d} & \hat{e} & \hat{f} \end{bmatrix} = \begin{bmatrix} * \\ * \\ * \end{bmatrix} \begin{bmatrix} * & * & * & * \end{bmatrix} r$$

$$(a, b, c, d, e, f) \approx (\hat{a}, \hat{b}, \hat{c}, \hat{d}, \hat{e}, \hat{f}) \in \mathfrak{M}_r$$

Data
model with complexity r

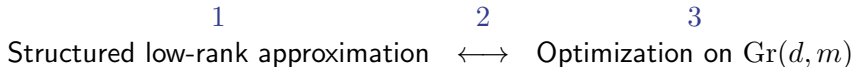
Structured low-rank approximation problem

$$\begin{bmatrix} a & b & c & d \\ b & c & d & e \\ c & d & e & f \end{bmatrix} \approx \begin{bmatrix} \hat{a} & \hat{b} & \hat{c} & \hat{d} \\ \hat{b} & \hat{c} & \hat{d} & \hat{e} \\ \hat{c} & \hat{d} & \hat{e} & \hat{f} \end{bmatrix} = \begin{bmatrix} * \\ * \\ * \end{bmatrix} \left[\begin{array}{cccc} * & * & * & * \end{array} \right] r$$

$$(a, b, c, d, e, f) \approx (\hat{a}, \hat{b}, \hat{c}, \hat{d}, \hat{e}, \hat{f}) \in \mathfrak{M}_r$$

Data
model with complexity r

Structure of [this talk](#):



Hankel matrices

$$\mathcal{H}_m(p) := \begin{bmatrix} p(1) & p(2) & \cdots & p(T-m+1) \\ p(2) & \ddots & & p(T-m+2) \\ \vdots & \ddots & \ddots & \vdots \\ p(m) & p(m+1) & \cdots & p(T) \end{bmatrix}, \quad \begin{array}{l} p = [p(1) \ p(2) \ \cdots p(T)] \\ \text{univariate time series} \end{array}$$

Theorem. (Heinig, 1984)

(evident if $r = m - 1$)

$$\text{rank } \mathcal{H}_m(p) \leq r < m \iff \theta_0 p(t) + \theta_1 p(t+1) + \cdots + \theta_r p(t+r) = 0,$$

linear recurrence $t = 1 : T - r$

Hankel matrices

$$\mathcal{H}_m(p) := \begin{bmatrix} p(1) & p(2) & \cdots & p(T-m+1) \\ p(2) & \ddots & & p(T-m+2) \\ \vdots & \ddots & \ddots & \vdots \\ p(m) & p(m+1) & \cdots & p(T) \end{bmatrix}, \quad \begin{array}{l} p = [p(1) \ p(2) \ \cdots p(T)] \\ \text{univariate time series} \end{array}$$

Theorem. (Heinig, 1984)

(evident if $r = m - 1$)

$$\text{rank } \mathcal{H}_m(p) \leq r < m \iff \theta_0 p(t) + \theta_1 p(t+1) + \cdots + \theta_r p(t+r) = 0, \\ t = 1 : T - r$$


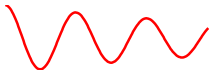

$$\iff p(t) = \sum_{k=1}^d \underbrace{P_k(t)}_{\text{polynomial}} \cdot \lambda_k^t \quad \text{exponential}$$

where

- $\lambda_1, \dots, \lambda_d$ — distinct roots of $\theta(z) = \sum_{j=0}^r \theta_j z^j$
- $\deg P_k(t) = (\text{multiplicity of } \lambda_k) - 1$ ($P_k = \text{const}$ if λ_k is simple)

Low-rank Hankel matrices: examples

$$p(t) = \sum_{k=1}^d P_k(t) \lambda_k^t \iff \theta_0 p(t) + \dots + \theta_r p(t+r) = 0$$

Real p :	picture	formula	r	$\sum_{j=0}^r \theta_j z^j$
exponential		$c\rho^t$	1	$(z - \rho)$
damped sine		$c\rho^t \cos(\omega t + \phi)$	2	$(z - \rho e^{i\omega}) \cdot (z - \rho e^{-i\omega})$
polynomial		$\sum_{k=0}^3 c_k t^k$	4	$(z - 1)^4$

Low-rank approximation

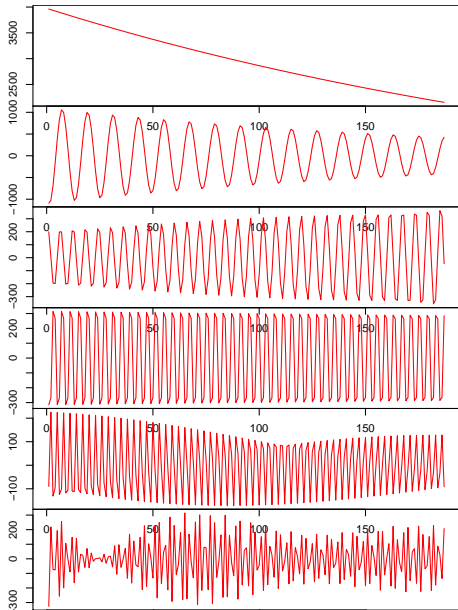
$$\mathcal{H}_m(p) \approx \text{l.r. } \mathcal{H}_m(\hat{p})$$

\leftrightarrow

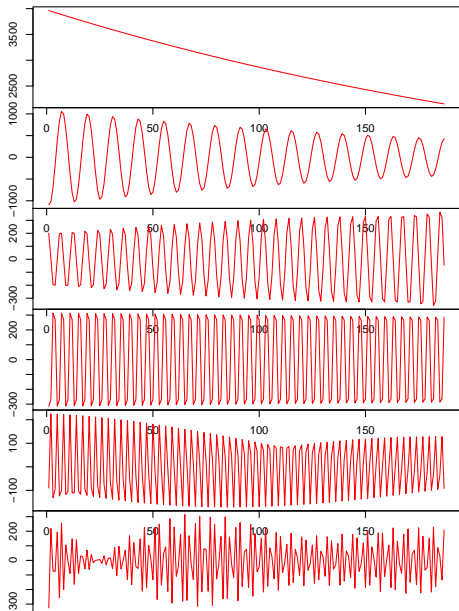
Sparse approximation with ∞ dictionary

$$p \approx \sum_{k=1}^d c_k f_k, \quad f_k(t) \in \left\{ t^\alpha \lambda^t \right\}_{\substack{\alpha \in \mathbb{N}, \\ \lambda \in \mathbb{C}}}$$

Structured low-rank approximation

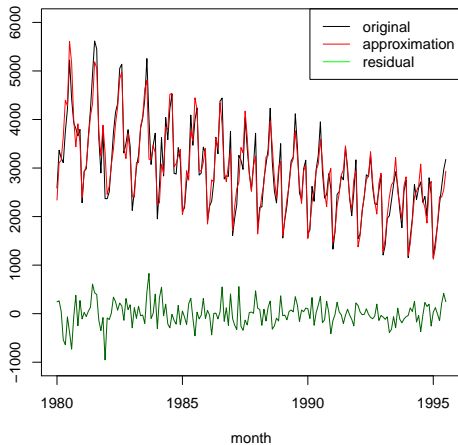


Structured low-rank approximation



5 of 22

Monthly Australian fortified wine sales.



Block-Hankel matrices

$$\mathcal{H}_{\ell+1}(w) := \begin{bmatrix} w(1) & w(2) & \cdots & w(T-\ell) \\ w(2) & \ddots & & w(T-\ell+1) \\ \vdots & \ddots & \ddots & \vdots \\ w(\ell+1) & w(\ell+2) & \cdots & w(T) \end{bmatrix}, \quad \begin{array}{l} w = [w(1) \ w(2) \ \cdots \ w(T)] \in \mathbb{R}^{q \times T} \\ q\text{-variate time series} \end{array}$$

$$\text{rank } \mathcal{H}_{\ell+1}(w) \leq \underbrace{(\ell+1)q}_{\text{number of rows}} - p \iff [R_0 \ \cdots \ R_\ell] \mathcal{H}_{\ell+1}(w) = 0, \quad R_k \in \mathbb{R}^{p \times q}$$

$$\iff R_0 w(t) + \cdots + R_\ell w(t+\ell) = 0, \quad t=1:T-\ell$$

Block-Hankel matrices

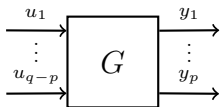
$$\mathcal{H}_{\ell+1}(w) := \begin{bmatrix} w(1) & w(2) & \cdots & w(T-\ell) \\ w(2) & \ddots & & w(T-\ell+1) \\ \vdots & \ddots & \ddots & \vdots \\ w(\ell+1) & w(\ell+2) & \cdots & w(T) \end{bmatrix}, \quad w = [w(1) \ w(2) \ \cdots \ w(T)] \in \mathbb{R}^{q \times T}$$

q-variate time series

$$\text{rank } \mathcal{H}_{\ell+1}(w) \leq \underbrace{(\ell+1)q}_{\text{number of rows}} - p \iff [R_0 \ \cdots \ R_\ell] \mathcal{H}_{\ell+1}(w) = 0, \quad R_k \in \mathbb{R}^{p \times q}$$

$$\iff R_0 w(t) + \cdots + R_\ell w(t + \ell) = 0, \quad t=1:T-\ell$$

\iff (Willems, 1986)



$w = \Pi \text{col}(u, y)$ — trajectory of a dynamical system with lag $\leq \ell$ and $\leq q-p$ inputs

Low-rank approximation

$$\mathcal{H}_m(w) \approx \text{l.r. } \mathcal{H}_m(\hat{w})$$

\iff

System identification from w

G — ?

Sylvester matrices

Two polynomials: $a(z) = \sum_{k=0}^m a_k z^k$ and $b(z) = \sum_{k=0}^m b_k z^k$

Sylvester matrix: $S(a, b) := \begin{bmatrix} a_0 & & & b_0 & & & \\ \vdots & \ddots & & \vdots & \ddots & & \\ \vdots & & a_0 & \vdots & & b_0 & \\ a_m & & \vdots & b_m & & \vdots & \\ & \ddots & \vdots & & \ddots & \vdots & \\ & & a_m & & & b_m & \end{bmatrix} \in \mathbb{R}^{2m \times 2m}$

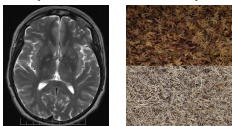
Theorem. $\deg \gcd(a, b) \geq d \iff \text{rank } S(a, b) \leq 2m - d,$

<p>Low-rank approximation $S(a, b) \approx \text{l.r. } S(\hat{a}, \hat{b})$</p>	\leftrightarrow	<p>Approximate common divisor problem $(a(z), b(z)) \approx (p(z)h(z), q(z)h(z))$</p>
--	-------------------	---

Other structures

- Multipolynomial Sylvester matrices
 - approximate GCD of **multiple polynomials**
- Multivariate Sylvester-like matrices (Macaulay matrices)
 - **multivariate polynomials**: approximate Gröbner bases, GCD
- Multilevel (nested) Hankel matrices
 - processing of ***n*-way arrays**

image processing
(MRI, textures)



Afiljler at en.wikipedia /
CC-BY-SA-3.0

symmetric tensor decomposition
(ICA, nonlinear SYSID)

$$\text{Cube} \approx \underset{a}{\overset{a}{\swarrow}} \underset{a}{\searrow} + \underset{b}{\overset{b}{\swarrow}} \underset{b}{\searrow} + \dots + \underset{z}{\overset{z}{\swarrow}} \underset{z}{\searrow}$$

Structured low-rank approximation: formulation

Linear structure: linear map $\mathcal{S} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{m \times n}$

Structured low-rank approximation: Given \mathcal{S} , $\|\cdot\|$, $p \in \mathbb{R}^{n_p}$, $r < m$

$$\underset{\hat{p} \in \mathbb{R}^{n_p}}{\text{minimize}} \quad \|p - \hat{p}\| \quad \text{subject to} \quad \text{rank } \mathcal{S}(\hat{p}) \leq r,$$

Structured low-rank approximation: formulation

Linear structure: linear map $\mathcal{S} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{m \times n}$

Structured low-rank approximation: Given \mathcal{S} , $\|\cdot\|$, $p \in \mathbb{R}^{n_p}$, $r < m$

$$\underset{\hat{p} \in \mathbb{R}^{n_p}}{\text{minimize}} \|p - \hat{p}\| \quad \text{subject to} \quad \text{rank } \mathcal{S}(\hat{p}) \leq r,$$

Weighted Euclidean **semi-norm**:

$$\|p\|_w^2 = \sum_{k=1}^{n_p} w_k p_k^2, \quad w_k \in [0; +\infty]$$

- $w_k = +\infty \iff$ constraint $p_k = \hat{p}_k$ — **fixed values**
- $w_k = 0 \iff$ p_k and \hat{p}_k do not matter — **missing values**

Weighted semi-norm: examples

1. fixed values

$$\begin{bmatrix} & a_0 & & b_0 \\ & \ddots & \vdots & \ddots & \vdots \\ a_0 & & \vdots & b_0 & \vdots \\ \vdots & & a_m & \vdots & b_m \\ \vdots & \ddots & & \vdots & \ddots \\ a_m & & b_m & & \end{bmatrix}$$

flipped Sylvester matrix



Hankel matrices with fixed zeros

2. missing values

$$\circ \begin{bmatrix} p_{1,1} & p_{1,2} & ? & p_{1,4} \\ p_{2,1} & p_{2,2} & p_{2,3} & p_{2,4} \\ p_{3,1} & p_{3,2} & p_{3,3} & ? \\ p_{4,1} & ? & p_{4,3} & p_{4,4} \end{bmatrix} \quad \text{— approximate matrix completion}$$

- System identification with missing data, in particular:

$$\begin{bmatrix} \mathcal{H}_\ell(u_1) & \mathcal{H}_\ell(?) \\ \mathcal{H}_\ell(y_1) & \mathcal{H}_\ell(y_{\text{ref}}) \end{bmatrix} \quad \text{— data-driven control}$$

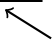
Reparameterization of the problem

$$\underset{\hat{p} \in \mathbb{R}^{n_p}}{\text{minimize}} \quad \|p - \hat{p}\|_w^2 \quad \text{subject to} \quad \text{rank } \mathcal{S}(\hat{p}) \leq r \quad (\text{SLRA})$$

rank constraint

kernel form

$$\text{rank } \mathcal{S}(\hat{p}) \leq r \quad \iff \quad d \begin{matrix} m \\ \boxed{R} \end{matrix} \cdot \begin{matrix} n \\ \boxed{\mathcal{S}(p)} \end{matrix} = 0, \quad \begin{matrix} \text{corank (at least)} \\ d := m - r \end{matrix}$$



 a full row rank matrix

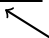
Reparameterization of the problem

$$\underset{\hat{p} \in \mathbb{R}^{n_p}}{\text{minimize}} \quad \|p - \hat{p}\|_w^2 \quad \text{subject to} \quad \text{rank } \mathcal{S}(\hat{p}) \leq r \quad (\text{SLRA})$$

rank constraint

kernel form

$$\text{rank } \mathcal{S}(\hat{p}) \leq r \quad \iff \quad d \begin{matrix} m \\ \boxed{R} \end{matrix} \cdot \begin{matrix} n \\ \boxed{\mathcal{S}(p)} \end{matrix} = 0, \quad \begin{matrix} \text{corank (at least)} \\ d := m - r \end{matrix}$$



 a full row rank matrix

$$(\text{SLRA}) \iff \underset{R \in \mathbb{R}^{d \times m}, \text{rank } R = d}{\text{minimize}} \quad f(R),$$

outer minimization

$$f(R) := \left(\underset{\hat{p} \in \mathbb{R}^{n_p}}{\text{minimize}} \quad \|p - \hat{p}\|_w^2 \quad \text{subject to} \quad R \mathcal{S}(\hat{p}) = 0 \right)$$

inner minimization

Inner minimization problem

$$f(R) := \left(\min_{\hat{p} \in \mathbb{R}^{n_p}} \|p - \hat{p}\|_w^2 \quad \text{subject to} \quad R\mathcal{S}(\hat{p}) = 0 \right), \quad w_k \in [0; +\infty]$$

↑ see (Markovsky, Usevich, 2013)

$$\min_{x, u} \|x\|_2^2 \quad \text{subject to} \quad A(R)x + \underbrace{B(R)u}_{\substack{\text{missing} \\ \text{data}}} = s(R) \quad \text{generalized least-norm problem}$$

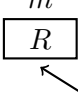
$$f(R) = s^\top B_\perp^\top (B_\perp A (B_\perp A)^\top)^{-1} B_\perp s, \quad \text{where } B_\perp : \begin{array}{l} \text{rowspan}(B_\perp) = \\ \text{(colspan}(B))_\perp \end{array}$$

(Usevich, Markovsky, 2013): for $\mathcal{S}(p) \in \mathbb{R}^{m \times n}$ mosaic Hankel, $w_k > 0$, f , ∇f and Hessian can be evaluated in $O(m^2 n)$ flops.

$$\text{Also, } f(R) = \frac{P(R)}{Q(R)}$$

Outer minimization problem

$$f(R) := \left(\min_{\hat{p} \in \mathbb{R}^{n_p}} \|p - \hat{p}\|_w^2 \text{ subject to } d \begin{array}{|c|} \hline R \\ \hline \end{array} \begin{array}{|c|} \hline \mathcal{S}(p) \\ \hline \end{array} = 0, \right),$$



 full row rank

Note. f depends only on the row space of R :

$$\text{rowspan}(R_1) = \text{rowspan}(R_2) \Rightarrow f(R_1) = f(R_2)$$

$$\Rightarrow \minimize_{R \in \mathbb{R}^{d \times m}, \text{rank } R = d} f(R) \quad \longleftrightarrow \quad \minimize_{\mathcal{L} \in \text{Gr}(d, m)} f(\mathcal{L})$$

Grassmann manifold: $\text{Gr}(d, m) := \{d\text{-dim. subspaces of } \mathbb{R}^m\}$

Constrained minimization

$$\underset{\mathcal{L} \in \text{Gr}(d, m)}{\text{minimize}} f(\mathcal{L}) \text{ --- ? ,}$$

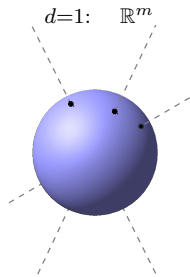
$$\text{Gr}(d, m) := \{d\text{-dim. subspaces of } \mathbb{R}^m\}$$

Constrained minimization:

$$\underset{R \in \mathbb{R}^{d \times m}}{\text{minimize}} f(R) \quad \text{subject to} \quad \underbrace{RR^T = I}_{\text{orthonormal basis}}$$

For example, use **penalty**:

$$\underset{R \in \mathbb{R}^{d \times m}}{\text{minimize}} f(R) + \gamma \|RR^T - I\|_F^2 \quad \text{--- exact}$$



Retraction-based methods

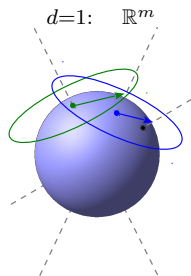
$$\underset{\mathcal{L} \in \text{Gr}(d, m)}{\text{minimize}} f(\mathcal{L}) \text{--- ?}$$

(Absil, Mahoney, Sepulchre, 2008), and others ...


1. From $x_k \in \text{Gr}(d, m)$
choose direction ξ_k in the tangent space
2. Set $x_{k+1} = R_{x_k}(\xi_k)$ (**retraction**)
3. Go to step 1 (is stopping criteria not satisfied).

Optimization methods

on manifolds: gradient descent, trust-region, ...


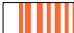


Parametrizations with permutation matrices


For any full-rank $R \in \mathbb{R}^{d \times m}$ there exist d lin. indep. columns: 



For any $\mathcal{L} \in \text{Gr}(d, m)$ there exist $X \in \mathbb{R}^{d \times (m-d)}$ and permutation Π such that $\mathcal{L} = \text{rowspan}([X \ I_d]\Pi)$



Take the permutation matrix:  $\Pi =$ 

Parametrizations with permutation matrices

For any full-rank $R \in \mathbb{R}^{d \times m}$ there exist d lin. indep. columns: 

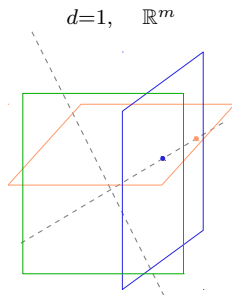


For any $\mathcal{L} \in \text{Gr}(d, m)$ there exist $X \in \mathbb{R}^{d \times (m-d)}$ and permutation Π such that $\mathcal{L} = \text{rowspan}([X \ I_d]\Pi)$

Take the permutation matrix:  $\Pi =$ 

$$\underset{\mathcal{L} \in \text{Gr}(d, m)}{\text{minimize}} f(\mathcal{L})$$

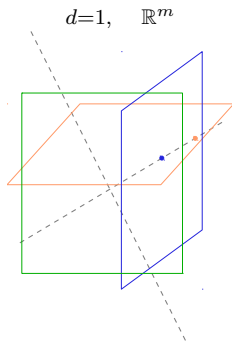
$$\iff \underset{\Pi \text{ — perm.}}{\text{minimize}} \quad \min_{X \in \mathbb{R}^{d \times (m-d)}} f([X \ I_d]\Pi)$$



Optimization with permutations

$$\underbrace{\text{minimize}}_{\substack{\Pi \text{ --- perm.} \\ \binom{m}{d} \text{ possibilities}}} \quad \underbrace{\min}_{X \in \mathbb{R}^{d \times (m-d)}} f([X \ I_d] \Pi)$$

unbounded



Optimization with permutations

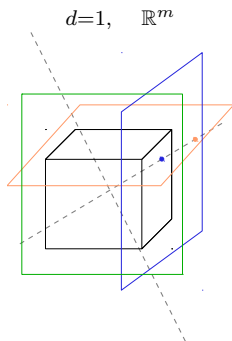
$$\underbrace{\text{minimize}}_{\substack{\Pi \text{ --- perm.} \\ \binom{m}{d} \text{ possibilities}}} \underbrace{\min}_{\substack{X \in [-1;1]^{d \times (m-d)} \\ \text{unbounded}}} f([X \ I_d]\Pi)$$

Theorem. (Knuth, 1985)

For any subspace $\mathcal{L} \in \text{Gr}(d, m)$ there exists a representation $[X \ I_d]\Pi$ with $|X_{k,l}| \leq 1$.

Easy to prove for $d = 1$: if $R = [r_1 \ \cdots \ r_m]$,

$$\text{then } \frac{R}{\max_k r_k} = [* \ \overset{k}{\pm 1} \ *]$$

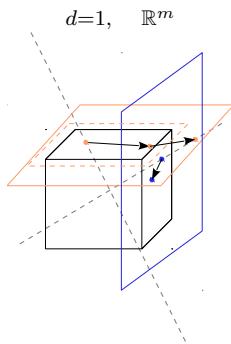


Optimization with switching permutations

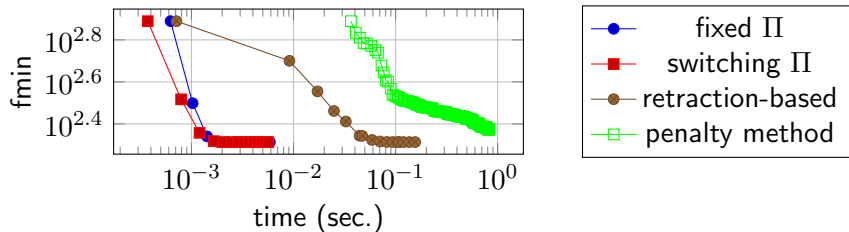
$$\underbrace{\text{minimize}_{\Pi \text{ — perm.}}}_{\binom{m}{d} \text{ possibilities}} \quad \underbrace{\min_{X \in [-1;1]^{d \times (m-d)}}}_{\text{unbounded}} f([X \ I_d]\Pi)$$

Switching permutations:

1. Perform local optimization of $f([X \ I_d]\Pi)$ until convergence, and unless $|X_{k,l}| \leq \Delta$ (where $\Delta > 1$)
2. If $|X_{k,l}| > \Delta$, switch the permutation, and go to step 1; otherwise stop.

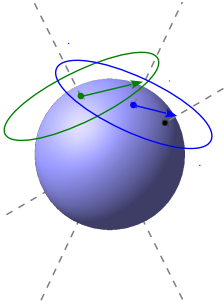
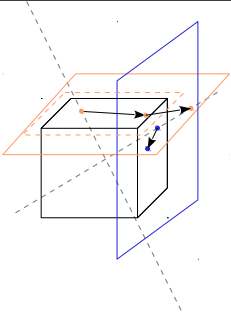


Comparison of the methods



Structured low-rank approximation for system identification: an example from DAISY database, 6×801 block-Hankel matrix with 2×1 blocks (example #3 from the abstract).

Conclusions

	Retraction-based	permutation-based
		
+++	more adapted to the local geometry	simple, any optimization method can be used
---	complicated, every method needs to be adapted/tuned	no bounds on the number of switches

- SLRA as optimization on a Grassmann manifold
K. Usevich and I. Markovsky. (2013).
Optimization on a Grassmann manifold with application
to system identification.
Preprint. <http://homepages.vub.ac.be/~kusevich/preprints.html>
- SLRA with missing values
I. Markovsky and K. Usevich (2013).
Structured low-rank approximation with missing data.
SIAM J. Matrix Anal. Appl. 34(2), 814-830.
- Fast cost function evaluation
K. Usevich and I. Markovsky. (2013).
Variable projection for affinely structured low-rank
approximation in weighted 2-norms.
J. Comput. Appl. Math. (doi:10.1016/j.cam.2013.04.034)
- Software (Matlab/R)
<http://github.com/slra/slra/>

Thank you!