

# An efficient preconditioner for the Riemannian trust-region method on the manifold of fixed-rank matrices

Marco Sutti

National Center for Theoretical Sciences  
Taipei, Taiwan

ICIAM 2023, Waseda University, Tokyo  
August 25, 2023

# Overview

Preprint: [Implicit low-rank Riemannian schemes for the time integration of stiff partial differential equations](#), M. Sutti and B. Vandereycken, submitted, arXiv preprint arXiv:2305.11532.

## Contributions:

- ▶ Preconditioner for the RTR method on the [manifold of fixed-rank matrices](#).
- ▶ Applications within implicit numerical integration schemes to solve stiff, time-dependent PDEs.

## This talk:

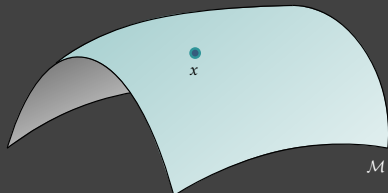
- I. Optimization on matrix manifolds.
- II. The manifold of fixed-rank matrices.
- III. Preconditioner.
- IV. Numerical application.

# Optimization problems on matrix manifolds

- ▶ We can state the **optimization problem** as

$$\min_{x \in \mathcal{M}} f(x),$$

where  $f : \mathcal{M} \rightarrow \mathbb{R}$  is the **objective function** and  $\mathcal{M}$  is some **matrix manifold**.



- ▶ **Matrix manifold**: any manifold that is constructed from  $\mathbb{R}^{n \times p}$  by taking either **embedded submanifolds** or **quotient manifolds**.
  - ▶ **Examples of embedded submanifolds**: orthogonal Stiefel manifold, manifold of symplectic matrices, **manifold of fixed-rank matrices**, ...
  - ▶ **Example of quotient manifold**: the Grassmann manifold.
- ▶ **Motivation**: by exploiting the **underlying geometric structure**, only feasible points are considered!

# Problems considered: variational problems

- ▶ **Variational problem**, called “LYAP” herein,

$$\begin{cases} \min_w \mathcal{F}(w(x, y)) = \int_{\Omega} \frac{1}{2} \|\nabla w(x, y)\|^2 - \gamma(x, y) w(x, y) \, dx \, dy \\ \text{such that } w = 0 \text{ on } \partial\Omega, \end{cases}$$

where  $\nabla = \left( \frac{\partial}{\partial x}, \frac{\partial}{\partial y} \right)$ ,  $\Omega = [0, 1]^2$  and  $\gamma$  is the source term.

- ▶ **Discretization on a uniform grid**: regardless of the specific form of  $\mathcal{F}$ , we obtain the general formulation

$$\min_W F(W) \quad \text{s.t.} \quad W \in \{X \in \mathbb{R}^{n \times n} : \text{rank}(X) = r\},$$

where  $F$  denotes the discretization of the functional  $\mathcal{F}$ .

---

“LYAP” variational problem: [Henson 2003, Gratton/Sartenaer/Toint 2008, Wen/Goldfarb 2009, S./Vandereycken 2021, ...]

# Riemannian manifold and gradient

A manifold  $\mathcal{M}$  endowed with a smoothly-varying inner product (called Riemannian metric  $g$ ) is called Riemannian manifold.

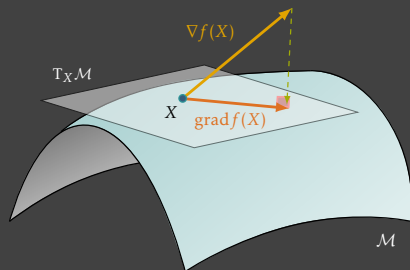
$\leadsto$  A couple  $(\mathcal{M}, g)$ , i.e., a manifold with a Riemannian metric on it.

Let  $f: \mathcal{M} \rightarrow \mathbb{R}$ . E.g., the objective function in an optimization problem.

$\leadsto$  For any embedded submanifold:

- ▶ Riemannian gradient: projection onto  $T_X\mathcal{M}$  of the Euclidean gradient

$$\text{grad } f(X) = P_{T_X\mathcal{M}}(\nabla f(X)).$$



$\leadsto$   $\nabla f(X)$  is the Euclidean gradient of  $f(X)$ .

---

Matrix and vector calculus: [The Matrix Cookbook](http://www.matrixcalculus.org), [www.matrixcalculus.org](http://www.matrixcalculus.org), ...

Automatic differentiation on low-rank manifolds: [Novikov/Rakhuba/Oseledets 2022]

# The manifold of fixed-rank matrices

- ▶ Our optimization problem is defined over

$$\mathcal{M}_r = \{X \in \mathbb{R}^{n \times n} : \text{rank}(X) = r\}.$$

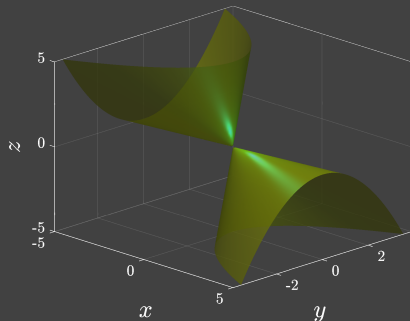
$\leadsto \mathcal{M}_r$  has a smooth structure ...

$2 \times 2$  example:

$$X = \begin{bmatrix} x & -2y \\ y & z \end{bmatrix}.$$

**Parametrization:**

$\text{rank}(X) = 1 \Leftrightarrow xz = -2y^2$  and  
 $x, z \neq 0$ .



- ▶ **Theorem:**  $\mathcal{M}_r$  is a smooth Riemannian submanifold embedded in  $\mathbb{R}^{n \times n}$  of dimension  $r(2n - r)$ .

# Alternative characterization

- ▶ Using the singular value decomposition (SVD), we have the equivalent characterization

$$\mathcal{M}_r = \{U\Sigma V^T : U^T U = I_r, V^T V = I_r, \Sigma = \text{diag}(\sigma_i), \sigma_1 \geq \dots \geq \sigma_r > 0\}.$$

The diagram illustrates the SVD decomposition of a matrix  $X$ . Matrix  $X$  is  $n \times n$ . It is equal to the product of matrix  $U$  ( $n \times r$ ), matrix  $\Sigma$  ( $r \times r$ ), and matrix  $V^T$  ( $r \times n$ ). Matrix  $\Sigma$  is shown as a square with a diagonal line and the Greek letter  $\Sigma$  below it. Matrix  $U$  is shown as a vertical rectangle. Matrix  $V^T$  is shown as a horizontal rectangle.

- ▶ Only  $2nr + r$  coefficients instead of  $n^2$ . If  $r \ll n$ , then big memory savings.
- ▶ Perform the calculations **directly** in the **factorized format**.

# Riemannian Hessian and preconditioning/1

- ▶ In the case of **Riemannian submanifolds**, the full Riemannian Hessian of  $f$  at  $x \in \mathcal{M}$  is given by the projected Euclidean Hessian plus the curvature part

$$\text{Hess } f(x)[\xi] = P_x \nabla^2 f(x) P_x + P_x (\text{“curvature terms”}) P_x.$$

$\leadsto$  Use  $P_x \nabla^2 f(x) P_x$  as a preconditioner in RTR.

- ▶ For **LYAP**, we can get the symmetric  $n^2$ -by- $n^2$  matrix

$$H_X = P_X (A \otimes I + I \otimes A) P_X.$$

- ▶ **Inverse of  $H_X$**   $\leadsto$  good candidate for a preconditioner.

**!** **Not inverted directly**, since this would cost  $\mathcal{O}(n^6)$ .

- ▶ A good preconditioner should reduce the number of iterations of the inner trust-region solver. It has to be **effective** and **cheap** to compute.



## Riemannian Hessian and preconditioning/2

- ▶ Applying the preconditioner in  $X \in \mathcal{M}_r$  means solving for  $\xi \in T_X \mathcal{M}$  the system

$$H_X \text{vec}(\xi) = \text{vec}(\eta),$$

where  $\eta \in T_X \mathcal{M}$  is a known tangent vector.

- ▶ This is equivalent to

$$P_X(A\xi + \xi A) = \eta.$$

- ▶ Using the definition of the orthogonal projector onto  $T_X \mathcal{M}_r$ , we obtain

$$P_U(A\xi + \xi A)P_V + P_U^\perp(A\xi + \xi A)P_V + P_U(A\xi + \xi A)P_V^\perp = \eta,$$

which is equivalent to the system

$$\begin{cases} P_U(A\xi + \xi A)P_V = P_U \eta P_V, \\ P_U^\perp(A\xi + \xi A)P_V = P_U^\perp \eta P_V, \\ P_U(A\xi + \xi A)P_V^\perp = P_U \eta P_V^\perp. \end{cases}$$

⋮

↪ Many (tedious) calculations, but the numerical results are pretty striking!

# “LYAP” variational problem

Table: Effect of preconditioning: dependence on size for LYAP.

Prec.	size	Rank 5						Rank 10					
		10	11	12	13	14	15	10	11	12	13	14	15
No	$n_{\text{outer}}$	51	54	61	59	162	92	300	103	61	63	62	59
	$\sum n_{\text{inner}}$	4561	9431	21066	36556	30069	30096	27867	30025	33818	45760	44467	38392
	$\max n_{\text{inner}}$	1801	3191	7055	9404	1194	1851	2974	3385	8894	24367	24537	25013
Yes	$n_{\text{outer}}$	41	45	50	52	56	60	44	64	62	53	56	56
	$\sum n_{\text{inner}}$	44	45	50	52	56	60	69	104	82	60	69	56
	$\max n_{\text{inner}}$	4	1	1	1	1	1	9	9	8	8	8	1

- ▶ **Stopping criterion:** maximum number of outer iterations  $n_{\text{max outer}} = 300$ . The inner solver is stopped when  $\sum n_{\text{inner}}$  first exceeds 30000.
- ▶ **Impressive reduction** in the number of iterations of the inner solver.
- ▶  $n_{\text{outer}}$  and  $\sum n_{\text{inner}}$  depend (quite mildly) on size, while  $\max n_{\text{inner}}$  is basically constant.

# Allen–Cahn equation/1

- ▶ **Reaction-diffusion equation** that models the process of phase separation in multi-component alloy systems.
  - ▶ Other applications include: mean curvature flows, two-phase incompressible fluids, complex dynamics of dendritic growth, and image segmentation ...
- ▶ In its simplest form, it reads

$$\frac{\partial w}{\partial t} = \varepsilon \Delta w + w - w^3.$$

- ▶ It is a **stiff**, time-dependent PDE.

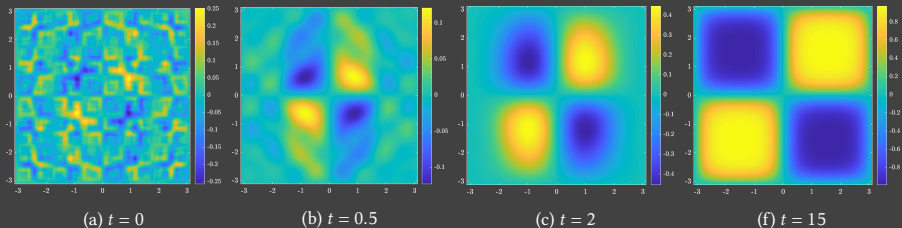


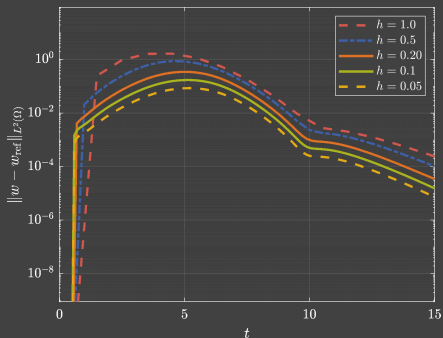
Figure: Time evolution of the solution  $w$  to the Allen–Cahn equation, with ERK4,  $h = 10^{-4}$ .

Allen–Cahn equation: [Allen/Cahn 1972, Allen/Cahn 1973]

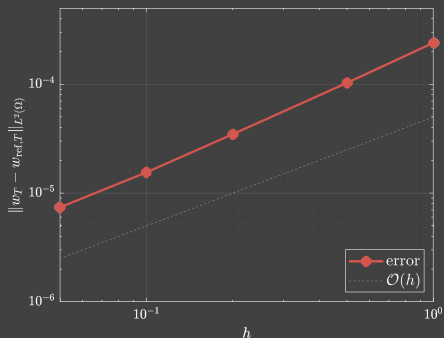
# Allen–Cahn equation/2 - low-rank evolution

- We build the functional

$$\min_w \mathcal{F}(w) := \int_{\Omega} \frac{\varepsilon h}{2} \|\nabla w\|^2 + \frac{(1-h)}{2} w^2 + \frac{h}{4} w^4 - \tilde{w} \cdot w \, dx \, dy.$$



(a)



(b)

**Figure:** Panel (a): error versus time for the preconditioned low-rank evolution of the Allen-Cahn equation. Panel (b): error at  $T = 15$  versus time step  $h$ .

# Conclusions

## Pros and cons:

- ⊕ Efficient preconditioner on the manifold of fixed-rank matrices.
- ⊕ Solid, quite well-understood mathematical theory behind.
- ⊖ If the problem does not admit a low-rank representation, then there is no advantage over using dense matrices.

## Outlook:

- ▶ Go to higher-order numerical integration methods.
- ▶ Other applications in mind, e.g., diffusion problems in mathematical biology or problems with low-rank tensor structure.

Thank you for your attention!

Questions?



# Metric, projection, gradient, retraction

- ▶ The **Riemannian metric** is

$$g_X(\xi, \eta) = \langle \xi, \eta \rangle = \text{Tr}(\xi^\top \eta), \quad \text{with } X \in \mathcal{M}_r \quad \text{and} \quad \xi, \eta \in T_X \mathcal{M}_r,$$

where  $\xi, \eta$  are seen as matrices in the ambient space  $\mathbb{R}^{n \times n}$ .

- ▶ **Orthogonal projection** onto the tangent space at  $X$  is

$$P_{T_X \mathcal{M}_r} : \mathbb{R}^{n \times n} \rightarrow T_X \mathcal{M}_r, \quad Z \rightarrow P_U Z P_V + P_U^\perp Z P_V + P_U Z P_V^\perp.$$

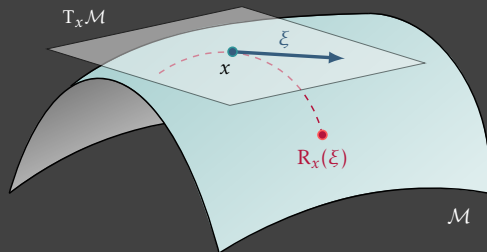
- ▶ **Riemannian gradient**: projection onto  $T_X \mathcal{M}_r$  of the **Euclidean gradient**

$$\text{grad } f(X) = P_{T_X \mathcal{M}_r}(\nabla f(X)).$$

- ▶ **Retraction**  $R_X : T_X \mathcal{M}_r \rightarrow \mathcal{M}_r$ . Typical: **truncated SVD**.

# Retractions

- ▶ Move in the direction of  $\xi$  while remaining constrained to  $\mathcal{M}$ .
- ▶ Smooth mapping  $R_x: T_x\mathcal{M} \rightarrow \mathcal{M}$  with a local condition that preserves gradients at  $x$ .



- ▶ The **Riemannian exponential mapping** is also a retraction, but it is not computationally efficient.
- ▶ **Retractions: first-order approximation of the Riemannian exponential!**



# Riemannian trust-region (RTR) method

---

**Algorithm 1:** Riemannian trust-region (RTR)

---

1 Given  $\bar{\Delta} > 0, \Delta_1 \in (0, \bar{\Delta})$

2 **for**  $i = 1, 2, \dots$  **do**

3     **Define** the second-order model

$$m_i: T_{x_i} \mathcal{M} \rightarrow \mathbb{R}, \xi \mapsto f(x_i) + \langle \text{grad } f(x_i), \xi \rangle + \frac{1}{2} \langle \text{Hess } f(x_i)[\xi], \xi \rangle.$$

4     **Trust-region subproblem:** compute  $\eta_i$  by solving

$$\eta_i = \operatorname{argmin} m_i(\xi) \quad \text{s.t.} \quad \|\xi\| \leq \Delta_i.$$

5     Compute  $\rho_i = (\widehat{f}(0) - \widehat{f}_i(\eta_i)) / (m_i(0) - m_i(\eta_i))$ .

6     **if**  $\rho_i \geq 0.05$  **then**

7         | Accept step and set  $x_{i+1} = R_{x_i}(\eta_i)$ .

8     **else**

9         | Reject step and set  $x_{i+1} = x_i$ .

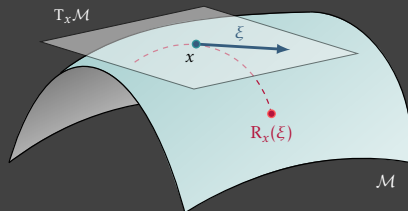
10     **end if**

11     Radius update: set

$$\Delta_{i+1} = \begin{cases} \min(2\Delta_i, \bar{\Delta}) & \text{if } \rho_i \geq 0.75 \text{ and } \|\eta_i\| = \Delta_i, \\ 0.25\|\eta_i\| & \text{if } \rho_i \leq 0.25, \\ \Delta_i & \text{otherwise.} \end{cases}$$

12 **end for**

---



TR method: [Goldfeld/Quandt/Trotter 1966, Sorensen 1982, Fletcher 1980/1987 ...]

RTR method: [Absil/Baker/Gallivan 2007]

## An example of factorized gradient

- ▶ “LYAP” functional:  $\mathcal{F}(w(x, y)) = \int_{\Omega} \frac{1}{2} \|\nabla w(x, y)\|^2 - \gamma(x, y) w(x, y) dx dy$ .
- ▶ The gradient of  $\mathcal{F}$  is the variational derivative  $\frac{\delta \mathcal{F}}{\delta w} = -\Delta w - \gamma$ .
- ▶ The discretized Euclidean gradient in matrix form is given by

$$G = AW + WA - \Gamma.$$

with  $A$  is the second-order periodic finite difference differentiation matrix.

- ▶ The first-order optimality condition  $G = AW + WA - \Gamma = 0$  is a Lyapunov (or Sylvester) equation.

→ Factorized Euclidean gradient:

$$G = \begin{bmatrix} AU & U & U_{\gamma} \end{bmatrix} \text{blkdiag}(\Sigma, \Sigma, \Sigma_{\gamma}) \begin{bmatrix} V & AV & V_{\gamma} \end{bmatrix}^{\top}.$$

$$\begin{bmatrix} AU & U & U_{\gamma} \end{bmatrix} \begin{bmatrix} \square & & \\ & \square & \\ & & \square \end{bmatrix} \begin{bmatrix} V & AV & V_{\gamma} \end{bmatrix}^{\top}$$

# Tangent vectors

- ▶ A **tangent vector**  $\xi$  at  $X = U\Sigma V^\top$  is represented as

$$\xi = UMV^\top + U_p V^\top + UV_p^\top,$$

$$M \in \mathbb{R}^{r \times r}, \quad U_p \in \mathbb{R}^{n \times r}, \quad U_p^\top U = 0, \quad V_p \in \mathbb{R}^{n \times r}, \quad V_p^\top V = 0.$$

- ▶ We can rewrite it as

$$\xi = (UM + U_p)V^\top + UV_p^\top.$$

$\leadsto \xi$  is a rank- $2r$  bounded matrix. Useful in implementation.