

# Data driven Koopman spectral analysis - a numerical linear algebra perspective

Zlatko Drmač

Department of Mathematics, University of Zagreb, Croatia

Research supported by: DARPA Contract HR0011-16-C-0116 *"On a Data-Driven, Operator-Theoretic Framework for Space-Time Analysis of Process Dynamics"* and the DARPA Contract HR0011-18-9-0033 *"The Physics of Artificial Intelligence"*

joint work with

Igor Mezić and Ryan Mohr, University of California, Santa Barbara

Operator Theoretic Methods in Dynamic Data Analysis and Control  
Institute for Pure and Applied Mathematics 2019

# Overview

- 1 Introduction
- 2 DMD and RRRR-DMD
- 3 Vandermonse+DFT based algorithm
- 4 Snapshot reconstruction
- 5 ...

# Data driven spectral analysis

Suppose we are given a sequence of snapshots  $\mathbf{f}_i \in \mathbb{C}^n$  of an underlying dynamics, that are driven by an inaccessible *black box* linear operator  $\mathbb{A}$ ;

$$\mathbf{f}_{i+1} = \mathbb{A}\mathbf{f}_i, \quad i = 1, \dots, m, \quad m < n, \quad (1)$$

with some initial  $\mathbf{f}_1$  and a time lag  $\delta t$ . No other information is available.

The two basic tasks of the Dynamic Mode Analysis (DMD) are

- 1 Identify approximate eigenpairs  $(\lambda_j, z_j)$  such that

$$\mathbb{A}z_j \approx \lambda_j z_j, \quad \lambda_j = |\lambda_j|e^{i\omega_j\delta t}, \quad j = 1, \dots, k; \quad k \leq m, \quad (2)$$

- 2 Derive a spectral spatio-temporal representation of the snapshots  $\mathbf{f}_i$ :

$$\mathbf{f}_i \approx \sum_{j=1}^{\ell} z_{\varsigma_j} \alpha_j \lambda_{\varsigma_j}^{i-1} \equiv \sum_{j=1}^{\ell} z_{\varsigma_j} \alpha_j |\lambda_{\varsigma_j}|^{i-1} e^{i\omega_{\varsigma_j}(i-1)\delta t}, \quad i = 1, \dots, m. \quad (3)$$

# Data driven spectral analysis – deep connections to Koopman operator theory/applications

The decomposition of the snapshots (3) reveals dynamically relevant spatial structures, the  $z_{\varsigma_j}$ 's, that evolve with amplitudes and frequencies encoded in the corresponding  $\lambda_{\varsigma_j}$ 's. It is desirable to have small number  $\ell$  of the most important modes  $z_{\varsigma_1}, \dots, z_{\varsigma_\ell}$ ,  $\varsigma_j \in \{1, \dots, k\}$ .

Such a sequence of snapshot (vectors of observables) can be obtained e.g. using black-boxed high performance ODE/PDE software, or e.g. by hi speed camera in an analysis of combustion instabilities in flame dynamics.

In a carefully designed framework with reach enough set of properly selected observables, the DMD can be considered as a finite dimensional spectral approximation of the Koopman operator associated with the dynamics under study [Arbabi+Mezić]. This deep theoretical connection gives the DMD a pivotal role in computational study of complex phenomena in fluid dynamics, see e.g. [Rowley], [Williams].

# Data driven spectral analysis - applications and software implementations

Other successful applications of DMD include e.g. aeroacoustics [Lele+Nichols], affective computing (analysis of videos for human emotion recognition [Cat Le Ngo+et al.]), robotics (filtering external perturbation using DMD based prediction [Berger+et al.]), algorithmic trading on financial markets [Mann+Kutz], analysis of infectious disease spread [Proctor+Eckhoff], neuroscience [Brunon+et al.] – just to name a few.

## Computational aspects (for software implementation):

- matrix multiplication and other simple matrix/vector operations
- SVD decomposition, Moore-Penrose pseudo-inverse, least squares
- eigenvalue of matrices of moderate dimensions

All necessary software implementations available in state of the art packages such as Matlab, Python (NumPy, SciPy) – LAPACK based.

**So is there anything left to do for a numerical analyst?**

# Tool: Krylov subspaces

- For  $i = 1, 2, \dots, m$ , define the Krylov matrices

$$\mathbf{X}_i = (\mathbf{f}_1 \quad \mathbf{f}_2 \quad \dots \quad \mathbf{f}_{i-1} \quad \mathbf{f}_i), \quad \mathbf{Y}_i = (\mathbf{f}_2 \quad \mathbf{f}_3 \quad \dots \quad \mathbf{f}_i \quad \mathbf{f}_{i+1}) \equiv \mathbb{A}\mathbf{X}_i,$$

and the corresponding Krylov subspaces  $\mathcal{X}_i = \text{range}(\mathbf{X}_i) \subset \mathbb{C}^n$ .

- Assume that at the index  $m$ ,  $\mathbf{X}_m$  is of full column rank. This implies

$$\mathcal{X}_1 \subsetneq \mathcal{X}_2 \subsetneq \dots \subsetneq \mathcal{X}_i \subsetneq \mathcal{X}_{i+1} \subsetneq \dots \subsetneq \mathcal{X}_m \subsetneq \dots \subsetneq \mathcal{X}_\ell = \mathcal{X}_{\ell+1}, \quad ,$$

i.e.  $\dim(\mathcal{X}_i) = i$  for  $i = 1, \dots, m$ , and there must be the smallest saturation index  $\ell$  at which  $\mathcal{X}_\ell = \mathcal{X}_{\ell+1}$ .

- $\mathbb{A}\mathcal{X}_\ell \subseteq \mathcal{X}_\ell$ , It is well known that then  $\mathcal{X}_\ell$  is the smallest  $\mathbb{A}$ -invariant subspace that contains  $\mathbf{f}_1$ .
- The action of  $\mathbb{A}$  on  $\mathcal{X}_m$  is known,  $\mathbb{A}(\mathbf{X}_m v) = \mathbf{Y}_m v$  for any  $v \in \mathbb{C}^m$ . Hence, useful spectral information can be obtained using the computable restriction  $\mathbf{P}_{\mathcal{X}_m} \mathbb{A}|_{\mathcal{X}_m}$ , that is, the Ritz values and vectors extracted using the Rayleigh quotient of  $\mathbb{A}$  with respect to  $\mathcal{X}_m$ .

# Tool: Krylov decomposition and Rayleigh-Ritz extraction

- To that end, let the vector  $c = (c_i)_{i=1}^m$  be computed from the least squares approximation

$$\|\mathbf{f}_{m+1} - \mathbf{X}_m c\|_2 \longrightarrow \min_c \quad (1)$$

and let  $r_{m+1} = \mathbf{f}_{m+1} - \mathbf{X}_m c$  be the corresponding residual. Recall that, by virtue of the theorem of projection,  $\mathbf{X}_m c = \mathbf{P}_{\mathcal{X}_m} \mathbf{f}_{m+1}$  and that  $r_{m+1}$  is orthogonal to the range of  $\mathbf{X}_m$ ,  $\mathbf{X}_m^* r_{m+1} = 0$ .

- Let  $E_{m+1} = r_{m+1} e_m^T$ ,  $e_m = (0, \dots, 0, 1)^T$ . The Krylov decomposition reads:

$$\mathbb{A} \mathbf{X}_m = \mathbf{X}_m C_m + E_{m+1}, \quad C_m = \begin{pmatrix} 0 & 0 & \dots & 0 & c_1 \\ 1 & 0 & \dots & 0 & c_2 \\ 0 & 1 & \dots & 0 & c_3 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & c_m \end{pmatrix},$$

# Rayleigh–Ritz extraction – basic properties

- ①  $C_m = (\mathbf{X}_m^* \mathbf{X}_m)^{-1} (\mathbf{X}_m^* \mathbb{A} \mathbf{X}_m) \equiv \mathbf{X}_m^\dagger \mathbb{A} \mathbf{X}_m = (\mathbf{X}_m^* \mathbf{X}_m)^{-1} (\mathbf{X}_m^* \mathbf{Y}_m)$  is the Rayleigh quotient, i.e. the matrix representation of  $\mathbf{P}_{\mathcal{X}_m} \mathbb{A}|_{\mathcal{X}_m}$
- ② If  $r_{m+1} = 0$  (and thus  $E_{m+1} = 0$  and  $m = \ell$ ) then  $\mathbb{A} \mathbf{X}_m = \mathbf{X}_m C_m$  and each eigenpair  $C_m w = \lambda w$  of  $C_m$  yields an eigenpair of  $\mathbb{A}$ ,  $\mathbb{A}(\mathbf{X}_m w) = \lambda(\mathbf{X}_m w)$ .
- ③ If  $r_{m+1} \neq 0$ , then  $(\lambda, z \equiv \mathbf{X}_m w)$  is an approximate eigenpair,  $\mathbb{A}(\mathbf{X}_m w) = \lambda(\mathbf{X}_m w) + r_{m+1}(e_m^T w)$ , i.e.  $\mathbb{A}z = \lambda z + r_{m+1}(e_m^T w)$ . The Ritz pair  $(\lambda, z)$  is acceptable if the residual

$$\frac{\|\mathbb{A}z - \lambda z\|_2}{\|z\|_2} = \frac{\|r_{m+1}\|_2}{\|z\|_2} |e_m^T w|$$

is sufficiently small. It holds that  $z^* r_{m+1} = 0$ , and

$$\lambda = \frac{z^* \mathbb{A} z}{z^* z} = \operatorname{argmin}_{\zeta \in \mathbb{C}} \|\mathbb{A} z - \zeta z\|_2$$

( $\lambda z$  is the orthogonal projection of  $\mathbb{A} z$  onto the span of  $z$ ).



# Beautiful structure and bad news

The spectral decomposition of  $C_m$  has beautiful structure. Assume for simplicity that the eigenvalues  $\lambda_i$ ,  $i = 1, \dots, m$ , are algebraically simple. It is easily checked that the spectral decomposition of  $C_m$  reads

$$C_m = \mathbb{V}_m^{-1} \Lambda_m \mathbb{V}_m, \quad \text{where } \Lambda_m = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{pmatrix}, \quad \mathbb{V}_m = \begin{pmatrix} 1 & \lambda_1 & \dots & \lambda_1^{m-1} \\ 1 & \lambda_2 & \dots & \lambda_2^{m-1} \\ \vdots & \vdots & \dots & \vdots \\ 1 & \lambda_m & \dots & \lambda_m^{m-1} \end{pmatrix}.$$

The Ritz vectors are the columns of  $Z_m = \mathbf{X}_m \mathbb{V}_m^{-1}$ .

**Bad news: The Vandermonde matrix  $\mathbb{V}_m$  is ill-conditioned!**

The condition number  $\kappa_2(\mathbb{V}_m) \equiv \|\mathbb{V}_m\|_2 \|\mathbb{V}_m^{-1}\|_2$  of any  $100 \times 100$  real Vandermonde matrix is **larger than  $3 \cdot 10^{28}$** ,  
 $(\kappa_2(\mathbb{V}_m) \geq 2^{m-2}/\sqrt{m}, m = 100, [\text{Gautschi}])$ .

Better: Schmid's DMD – compute the Rayleigh quotient using a POD (truncated SVD) basis of  $\mathbf{X}_m$ .

# SVD and best low rank approximation

## Theorem

*Eckart-Young-Mirsky Let the SVD of  $M \in \mathbb{C}^{n \times m}$  be*

$$M = U\Sigma V^*, \quad \Sigma = \text{diag}(\sigma_i)_{i=1}^{\min(m,n)}, \quad \sigma_1 \geq \dots \geq \sigma_{\min(m,n)} \geq 0.$$

*For  $k \in \{1, \dots, \min(m, n)\}$ , define  $U_k = U(:, 1 : k)$ ,  $\Sigma_k = \Sigma(1 : k, 1 : k)$ ,  $V_k = V(:, 1 : k)$ , and  $M_k = U_k \Sigma_k V_k^*$ . Then*

$$\min_{\text{rank}(N) \leq k} \|M - N\|_2 = \|M - M_k\|_2 = \sigma_{k+1} \quad (2)$$

$$\min_{\text{rank}(N) \leq k} \|M - N\|_F = \|M - M_k\|_F = \sqrt{\sum_{i=k+1}^{\min(n,m)} \sigma_i^2}. \quad (3)$$

Hence, if  $\sigma_m \ll \sigma_1$ , the condition number  $\kappa_2(\mathbf{X}_m) = \|\mathbf{X}_m\|_2 \|\mathbf{X}_m^\dagger\|_2 = \frac{\sigma_1}{\sigma_m}$  is large,  $\mathbf{X}_m$  can be made singular with a perturbation  $\delta \mathbf{X}_m$  such that  $\|\delta \mathbf{X}_m\|_2 / \|\mathbf{X}_m\|_2 = \sigma_m / \sigma_1 = 1 / \kappa_2(\mathbf{X}_m) \ll 1$ .

# Schmid's DMD

To avoid the ill-conditioning, Schmid used the thin truncated SVD  $\mathbf{X}_m = U\Sigma V^* \approx U_k \Sigma_k V_k^*$ , where  $U_k = U(:, 1:k)$  is  $n \times k$  orthonormal ( $U_k^* U_k = I_k$ ),  $V_k = V(:, 1:k)$  is  $m \times k$ , also orthonormal ( $V_k^* V_k = I_k$ ), and  $\Sigma_k = \text{diag}(\sigma_i)_{i=1}^k$  contains the largest  $k$  singular values of  $\mathbf{X}_m$ . In brief,  $U_k$  is the POD basis for the snapshots  $\mathbf{f}_1, \dots, \mathbf{f}_m$ . Since

$$\mathbf{Y}_m = \mathbb{A}\mathbf{X}_m \approx \mathbb{A}U_k \Sigma_k V_k^*, \quad \text{and} \quad \mathbb{A}U_k = \mathbf{Y}_m V_k \Sigma_k^{-1}, \quad (4)$$

the Rayleigh quotient  $S_k = U_k^* \mathbb{A}U_k$  with respect to the range of  $U_k$  can be computed as

$$S_k = U_k^* \mathbf{Y}_m V_k \Sigma_k^{-1}, \quad (5)$$

which is suitable for data driven setting because it does not use  $\mathbb{A}$  explicitly. Clearly, (4, 5) only require that  $\mathbf{Y}_m = \mathbb{A}\mathbf{X}_m$ ; it is not necessary that  $\mathbf{Y}_m$  is shifted  $\mathbf{X}_m$ . Each eigenpair  $(\lambda, w)$  of  $S_k$  generates the corresponding Ritz pair  $(\lambda, U_k w)$  for  $\mathbb{A}$ .

# Schmid's DMD

Algorithm  $[Z_k, \Lambda_k] = \text{DMD}(\mathbf{X}_m, \mathbf{Y}_m)$

**Input:** •  $\mathbf{X}_m = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ ,  $\mathbf{Y}_m = (\mathbf{y}_1, \dots, \mathbf{y}_m) \in \mathbb{C}^{n \times m}$  that define a sequence of snapshots pairs  $(\mathbf{x}_i, \mathbf{y}_i \equiv \mathbb{A}\mathbf{x}_i)$ . (Tacit assumption is that  $n$  is large and that  $m \ll n$ .)

- 1:  $[U, \Sigma, V] = \text{svd}(\mathbf{X}_m)$  ; { *The thin SVD:  $\mathbf{X}_m = U\Sigma V^*$ ,  $U \in \mathbb{C}^{n \times m}$ ,  $\Sigma = \text{diag}(\sigma_i)_{i=1}^m$ ,  $V \in \mathbb{C}^{m \times m}$* }
- 2: Determine numerical rank  $k$ .
- 3: Set  $U_k = U(:, 1 : k)$ ,  $V_k = V(:, 1 : k)$ ,  $\Sigma_k = \Sigma(1 : k, 1 : k)$
- 4:  $S_k = ((U_k^* \mathbf{Y}_m) V_k) \Sigma_k^{-1}$ ; { *Schmid's formula for the Rayleigh quotient  $U_k^* \mathbb{A} U_k$* }
- 5:  $[W_k, \Lambda_k] = \text{eig}(S_k)$  {  $\Lambda_k = \text{diag}(\lambda_i)_{i=1}^k$ ;  $S_k W_k(:, i) = \lambda_i W_k(:, i)$ ;  $\|W_k(:, i)\|_2 = 1$ }
- 6:  $Z_k = U_k W_k$  { *Ritz vectors*}

**Output:**  $Z_k, \Lambda_k$

# Data driven (computable) residual

Not all computed Ritz pairs will provide good approximations of eigenpairs of the underlying  $\mathbb{A}$ , and it is desirable that each pair is accompanied with an error estimate that will determine whether it can be accepted and used in the next steps of a concrete application. The residual is computationally feasible and usually reliable measure of fitness of a Ritz pair. With a simple modification, the DMD Algorithm can be enhanced with residual computation, without using  $\mathbb{A}$  explicitly.

## Proposition

For the Ritz pairs  $(\lambda_i, Z_k(:, i) \equiv U_k W_k(:, i))$ ,  $i = 1, \dots, k$ , computed in the DMD Algorithm, the residual norms can be computed as follows:

$$r_k(i) = \|\mathbb{A}Z_k(:, i) - \lambda_i Z_k(:, i)\|_2 = \|(\mathbf{Y}_m V_k \Sigma_k^{-1})W_k(:, i) - \lambda_i Z_k(:, i)\|_2. \quad (6)$$

Further, if  $\mathbb{A} = S \text{diag}(\alpha_i)_{i=1}^n S^{-1}$ , then  $\min_{\alpha_j} |\lambda_i - \alpha_j| \leq \kappa_2(S) r_k(i)$  (by the Bauer–Fike Theorem).

# Data driven (computable) residual

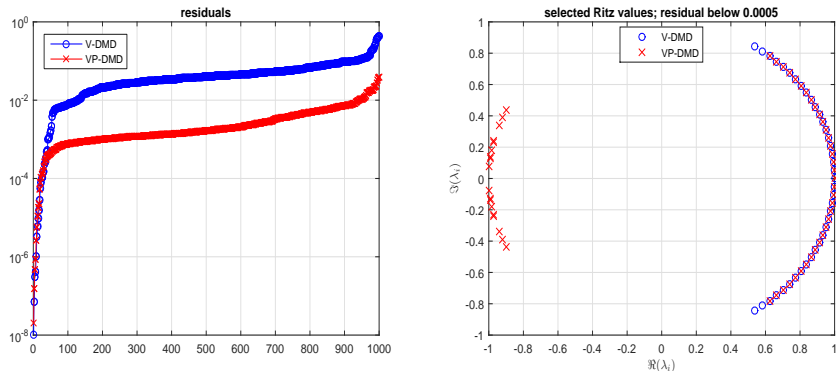
## Example

test The well studied and understood model of laminar flow around a cylinder is based on the two-dimensional incompressible Navier-Stokes equations

$$\frac{\partial \mathbf{v}}{\partial t} = -(\mathbf{v} \cdot \nabla) \mathbf{v} + \nu \Delta \mathbf{v} - \frac{1}{\rho} \nabla p, \quad \nabla \cdot \mathbf{v} = 0, \quad (7)$$

where  $\mathbf{v} = (v_x, v_y)$  is velocity field,  $p$  is pressure,  $\rho$  is fluid density and  $\nu$  is kinematic viscosity. The flow is characterized by the Reynolds number  $\mathfrak{Re} = v^* D / \nu$  where, for flow around circular cylinder, the characteristic quantities are the inlet velocity  $v^*$  and the cylinder diameter  $D$ . For a detailed analytical treatment of the problem see [Bagheri], [Glaz+et al.]; for a more in depth description of the Koopman analysis of fluid flow we refer to [Mezić], [Rowley].

# Data driven (computable) residual



**Figure:** Left: Comparison of the residuals of the Ritz pairs computed by the DMD\_RRR Algorithm with velocities as observables (V-DMD, circles  $\circ$ ) and with both velocities and pressures (VP-DMD, crosses,  $\times$ ). Right: Selected Ritz values with velocities as observables ( $\circ$ ) and with both velocities and pressures ( $\times$ ).

# Refined Ritz vectors

The Ritz vectors are not optimal eigenvectors approximations from a given subspace  $\mathcal{U}_k = \text{range}(U_k)$ . Hence, for a computed Ritz value  $\lambda$ , instead of the associated Ritz vector, we can choose a vector  $z \in \mathcal{U}_k$  that minimizes the residual. From the variational characterization of the singular values, it follows that

$$\begin{aligned} \min_{z \in \mathcal{U}_k \setminus \{0\}} \frac{\|\mathbb{A}z - \lambda z\|_2}{\|z\|_2} &= \min_{w \neq 0} \frac{\|\mathbb{A}U_k w - \lambda U_k w\|_2}{\|U_k w\|_2} \\ &= \min_{\|w\|_2=1} \|(\mathbb{A}U_k - \lambda U_k)w\|_2 = \sigma_{\min}(\mathbb{A}U_k - \lambda U_k), \end{aligned}$$

where  $\sigma_{\min}(\cdot)$  denotes the smallest singular value of a matrix, and the minimum is attained at the right singular vector  $w_\lambda$  corresponding to  $\sigma_\lambda \equiv \sigma_{\min}(\mathbb{A}U_k - \lambda U_k)$ . As a result, the refined Ritz vector corresponding to  $\lambda$  is  $U_k w_\lambda$  and the optimal residual is  $\sigma_\lambda$ . Detailed analysis by Jia.



# Data driven refinement of Ritz vectors

The minimization of the residual can be replaced with computing the smallest singular value with the corresponding right singular vector of  $B_k - \lambda U_k$ , where  $B_k \equiv \mathbb{A}U_k = \mathbf{Y}_m V_k \Sigma_k^{-1}$ . Compute the QRF

$$(U_k \quad B_k) = QR, \quad R = \begin{matrix} & k & k \\ k' & \begin{pmatrix} R_{[11]} & R_{[12]} \\ 0 & R_{[22]} \end{pmatrix} \end{matrix}, \quad k' = \min(n - k, k)$$

and write the pencil  $B_k - \lambda U_k$  as

$$B_k - \lambda U_k = Q \left( \begin{pmatrix} R_{[12]} \\ R_{[22]} \end{pmatrix} - \lambda \begin{pmatrix} R_{[11]} \\ 0 \end{pmatrix} \right) \equiv QR_\lambda, \quad R_\lambda = \begin{pmatrix} R_{[12]} - \lambda R_{[11]} \\ R_{[22]} \end{pmatrix}.$$

## Theorem

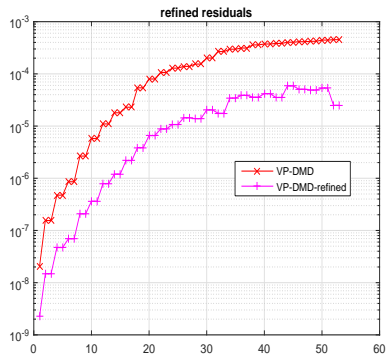
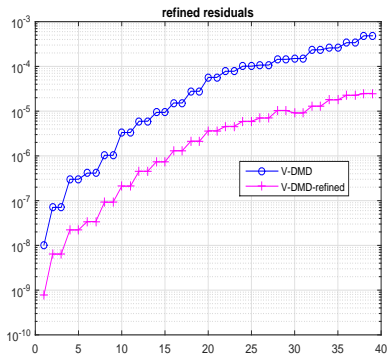
*Let for the Ritz value  $\lambda = \lambda_i$ ,  $w_{\lambda_i}$  denote the right singular vector of the smallest singular value  $\sigma_{\lambda_i}$  of the matrix  $R_{\lambda_i}$ . Then  $z = z_{\lambda_i} \equiv U_k w_{\lambda_i}$  minimizes the residual, whose minimal value equals  $\sigma_{\lambda_i} = \|R_{\lambda_i} w_{\lambda_i}\|_2$ .*

## Enhanced DMD Algorithm

$$[Z_k, \Lambda_k, r_k, \rho_k] = \text{DMD\_RRR}(\mathbf{X}_m, \mathbf{Y}_m; \epsilon) \text{ \{Refined Rayleigh-Ritz DMD\}}$$

- 1:  $D_x = \text{diag}(\|\mathbf{X}_m(:, i)\|_2)_{i=1}^m$ ;  $\mathbf{X}_m^{(1)} = \mathbf{X}_m D_x^\dagger$ ;  $\mathbf{Y}_m^{(1)} = \mathbf{Y}_m D_x^\dagger$
- 2:  $[U, \Sigma, V] = \text{svd}(\mathbf{X}_m^{(1)})$ ; numerical rank:  $k = \max\{i : \sigma_i \geq \sigma_1 \epsilon\}$ .
- 3: Set  $U_k = U(:, 1:k)$ ,  $V_k = V(:, 1:k)$ ,  $\Sigma_k = \Sigma(1:k, 1:k)$
- 4:  $B_k = \mathbf{Y}_m^{(1)}(V_k \Sigma_k^{-1})$ ; *\{Schmid's formula for  $\mathbb{A}U_k\}$*
- 5:  $[Q, R] = \text{qr}((U_k, B_k))$ ; *\{The thin QR factorization\}*
- 6:  $S_k = \text{diag}(\overline{R_{ii}})_{i=1}^k R(1:k, k+1:2k)$  *\{ $S_k = U_k^* \mathbb{A} U_k\}$*
- 7:  $\Lambda_k = \text{eig}(S_k)$  *\{\Lambda\_k = \text{diag}(\lambda\_i)\_{i=1}^k; Ritz values, i.e. eigenvalues of  $S_k\}$*
- 8: **for**  $i = 1, \dots, k$  **do**
- 9:    $[\sigma_{\lambda_i}, w_{\lambda_i}] = \text{svd}_{\min}(\begin{pmatrix} R(1:k, k+1:2k) - \lambda_i R(1:k, 1:k) \\ R(k+1:2k, k+1:2k) \end{pmatrix})$ ;
- 10:    $W_k(:, i) = w_{\lambda_i}$ ;  $r_k(i) = \sigma_{\lambda_i}$  *\{Optimal residual,  $\sigma_{\lambda_i} = \|R_{\lambda_i} w_{\lambda_i}\|_2\}$*
- 11:    $\rho_k(i) = w_{\lambda_i}^* S_k w_{\lambda_i}$  *\{Rayleigh quotient,  $\rho_k(i) = (U_k w_{\lambda_i})^* \mathbb{A} (U_k w_{\lambda_i})\}$*
- 12: **end for**
- 13:  $Z_k = U_k W_k$  *\{Refined Ritz vectors\}*

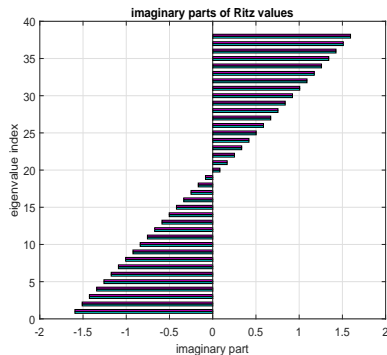
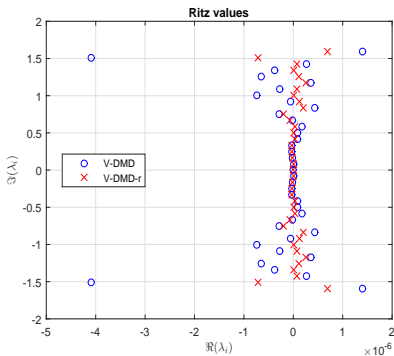
# Residuals of refined selected pairs



**Figure:** Comparison of the refined residuals of the Ritz pairs computed by the DMD\_RRR Algorithm with velocities as observables (top 39 pairs in V-DMD, circles  $\circ$ ) and with both velocities and pressures (top 53 pairs in VP-DMD, crosses,  $\times$ ). The noticeable staircase structure on the graphs corresponds to complex conjugate Ritz pairs.

# Koopman eigenvalues ... cyclic group structure

$$\mathfrak{K}_j = \frac{\log \lambda_j}{2\pi\Delta t} \equiv \frac{\log |\lambda_j| + i \arg \lambda_j}{2\pi\Delta t} \approx \frac{1}{2\pi\Delta t} i \arg \lambda_j. \quad (8)$$



**Figure:** The cyclic group structure of the eigenvalues on the unit circle is nicely visualized in the  $\mathfrak{K}_j$ 's; cf. (8). These values nicely correspond to the analytically derived formulas for the Koopman eigenvalues of the flow [Bagheri].

# A synthetic example

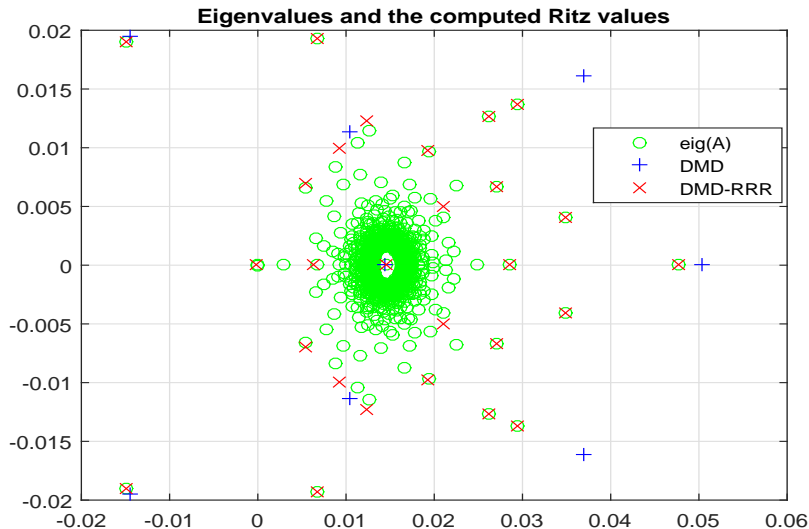
## Goal: DMD black-box software

The main goal of the modifications of the DMD algorithm is to provide a reliable black-box, data driven software device that can estimate part of the spectral information of the underlying linear operator, and that also can provide an error bound.

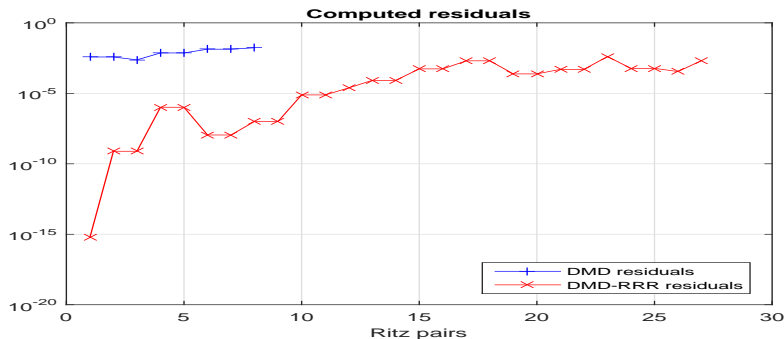
## Example (A case study)

The test matrix is generated as  $A = e^{-B^{-1}}$  where  $B$  is pseudo-random with entries uniformly distributed in  $[0, 1]$ , and then  $\mathbb{A} = A/\|A\|_2$ . Although this example is purely synthetic, it may represent a situation with the spectrum entirely in the unit disc, such as e.g. in the case of an off-attractor analysis of a dynamical system, after removing the peripheral eigenvalues, see e.g. Mohr & Mezić 2014.

# Accuracy of the computed Ritz values

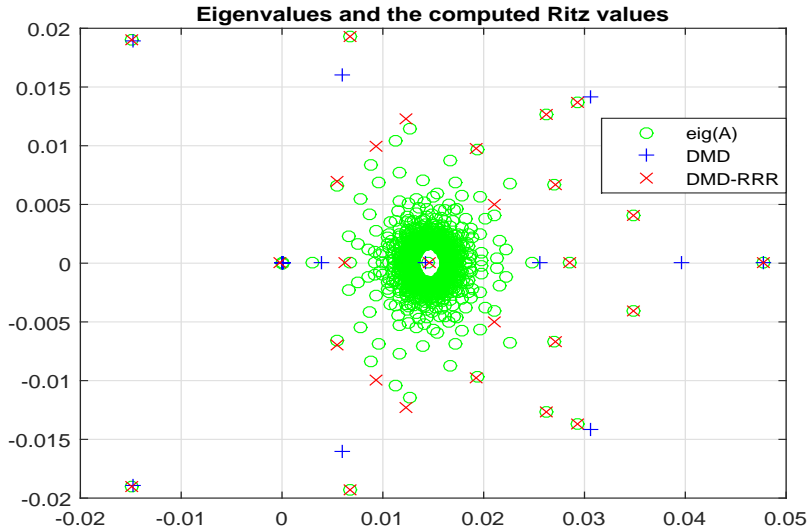


# Comparing residuals

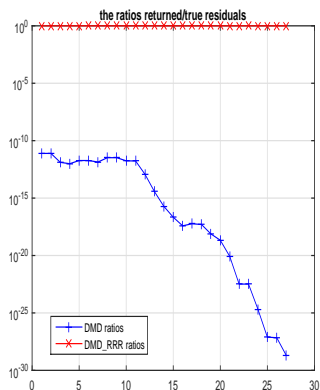
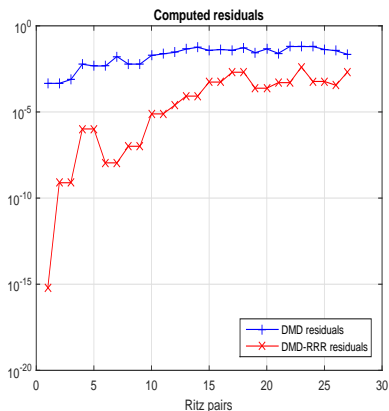


**Figure:** Comparison of the residuals of the Ritz pairs computed by the DMD Algorithm (pluses +) and the DMD\_RRR Algorithm (crosses, ×), with the same threshold in the truncation criterion for determining the numerical rank.

# Ritz values wit $k = 27$ (hard coded)

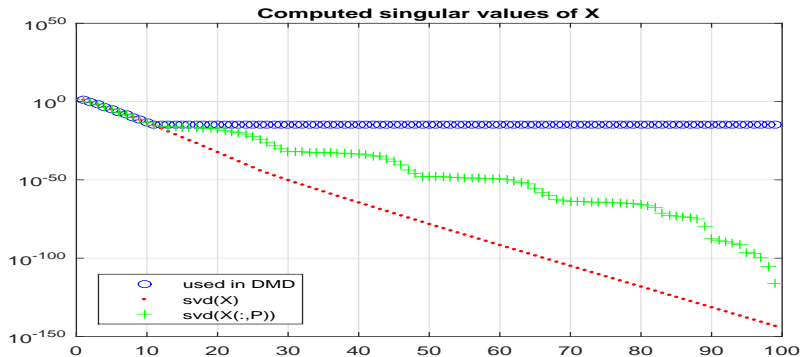




Residuals wit  $k = 27$  (hard coded)

$$\eta_i = \frac{\|(\mathbf{Y}_m \mathbf{V}_k \Sigma_k^{-1}) \mathbf{W}_k(:, i) - \lambda_i (\mathbf{U}_k \mathbf{W}_k(:, i))\|_2}{\|\mathbf{A}(\mathbf{U}_k \mathbf{W}_k(:, i)) - \lambda_i (\mathbf{U}_k \mathbf{W}_k(:, i))\|_2} \equiv 1, \quad i = 1, \dots, k.$$

# Singular values of $\mathbf{X}_m$ computed three times



**Figure:** The blue circles ( $\circ$ ) are the values used in the DMD Algorithm and are computed as  $[U, \Sigma, V] = \text{svd}(\mathbf{X}_m, 'econ')$ . The red dots ( $\cdot$ ) are computed as  $\Sigma = \text{svd}(\mathbf{X}_m)$ , and the pluses ( $+$ ) are the results of  $\Sigma = \text{svd}(\mathbf{X}_m(:, P))$ , where  $P$  is randomly generated permutation.

# Floating point SVD

If  $\mathbf{X}_m \approx \tilde{U}\tilde{\Sigma}\tilde{V}^*$  is the computed SVD of  $\mathbf{X}_m$ , then there exist unitary matrices  $\hat{U}$ ,  $\hat{V}$ , and a perturbation  $\delta\mathbf{X}_m$  (*backward error*) such that  $\|\hat{U} - \tilde{U}\|_2 \leq \epsilon_1$ ,  $\|\hat{V} - \tilde{V}\|_2 \leq \epsilon_2$ , and

$$\mathbf{X}_m + \delta\mathbf{X}_m = \hat{U}\tilde{\Sigma}\hat{V}^*, \quad \|\delta\mathbf{X}_m\|_2 \leq \epsilon\|\mathbf{X}_m\|_2. \quad (9)$$

## Theorem (Weyl and Wieland-Hoffman)

Let the singular values of  $\mathbf{X}_m$  and  $\mathbf{X}_m + \delta\mathbf{X}_m$  be  $\sigma_1 \geq \dots \geq \sigma_{\min(m,n)}$  and  $\tilde{\sigma}_1 \geq \dots \geq \tilde{\sigma}_{\min(m,n)}$ , respectively. Then

$$\max_i |\tilde{\sigma}_i - \sigma_i| \leq \|\delta\mathbf{X}_m\|_2; \quad \sqrt{\sum_{i=1}^{\min(m,n)} |\tilde{\sigma}_i - \sigma_i|^2} \leq \|\delta\mathbf{X}_m\|_F.$$

Hence, if we combine this Theorem with the backward stability (9), we have that for each computed singular value  $\tilde{\sigma}_i = \sigma_i + \delta\sigma_i$

$$|\delta\sigma_i| \leq \|\delta\mathbf{X}_m\|_2 \leq \epsilon\|\mathbf{X}_m\|_2; \quad |\delta\sigma_i|/\sigma_i \leq \epsilon\|\mathbf{X}_m\|_2/\sigma_i. \quad (10)$$

$$\|\delta \mathbf{f}_i\|_2 \leq \|\delta \mathbf{X}_m\|_2 \leq \epsilon \|\mathbf{X}_m\|_2 \leq \epsilon \sqrt{m} \max_i \|\mathbf{f}_i\|_2$$

Bad news for small  $\sigma_i$ 's:  $\max_i |\delta \sigma_i| / \sigma_i \leq \epsilon \kappa_2(\mathbf{X}_m)$ .

Suppose we have backward error  $\delta \mathbf{X}_m$  such that

$$\|\delta \mathbf{X}_m(:, i)\|_2 \leq \epsilon \|\mathbf{X}_m(:, i)\|_2, \quad i = 1, \dots, m. \quad (11)$$

In terms of the snapshots, this reads  $\|\delta \mathbf{f}_i\|_2 \leq \epsilon \|\mathbf{f}_i\|_2$ , for all snapshots.

### Theorem (Eisenstat and Ipsen)

Let  $\sigma_1 \geq \dots \geq \sigma_n$  and  $\tilde{\sigma}_1 \geq \dots \geq \tilde{\sigma}_n$ .  $A + \delta A = \Xi_1 A \Xi_2$  and let  $\xi = \max\{\|\Xi_1 \Xi_1^T - I\|_2, \|\Xi_2^T \Xi_2 - I\|_2\}$ . Then

$$|\tilde{\sigma}_i - \sigma_i| \leq \xi \sigma_i, \quad i = 1, \dots, n.$$

Hence, if  $\mathbf{X}_m$  is of full column rank,  $\mathbf{X}_m + \delta \mathbf{X}_m = (\mathbb{I}_n + \delta \mathbf{X}_m \mathbf{X}_m^\dagger) \mathbf{X}_m$  and n application of this theorem yields

$$\max_i \frac{|\sigma_i - \tilde{\sigma}_i|}{\sigma_i} \leq 2 \|\delta \mathbf{X}_m \mathbf{X}_m^\dagger\|_2 + \|\delta \mathbf{X}_m \mathbf{X}_m^\dagger\|_2^2.$$

$\|\delta \mathbf{X}_m \mathbf{X}_m^\dagger\|_2$  invariant under column scalings!

# Dicussion on the SVD

Matlab uses different algorithms in the `svd()` function, depending on whether the singular vectors are requested on output.

- The faster but less accurate method is used in the call  $[U, S, V] = \text{svd}(\mathbf{X}_m, 'econ')$ . It is very likely that the full SVD, including the singular vectors, is computed using the divide and conquer algorithm, `xGESDD()` in LAPACK.
- For computing only the singular values  $S = \text{svd}(X)$  calls the QR SVD, `xGESVD()` in LAPACK.

Note that the same fast `xGESDD()` subroutine is (to our best knowledge) under the hood of the Python function `numpy.linalg.svd`.

Numerical robustness of both `xGESVD()`, `xGESDD()` depends on  $\kappa_2(\mathbf{X}_m)$ , and if one does not take advantage of the fact that scaling is allowed, the problems illustrated in this example are likely to happen.

Better: Jacobi SVD (`xGEJSV()`, `xGESVJ()` in LAPACK, Drmač 2009.) and preconditioned QR (`xGESVDQ()`, LAPACK, Drmač 2018.).

# Extension to weighted DMD using energy $(\cdot, \cdot)_M$

$$[Z_k, \Lambda_k] = \text{Weighted\_DMD}(\mathbf{X}_m, \mathbf{Y}_m; M; \epsilon)$$

## Input:

- $\mathbf{X}_m = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ ,  $\mathbf{Y}_m = (\mathbf{y}_1, \dots, \mathbf{y}_m) \in \mathbb{C}^{n \times m}$  that define a sequence of snapshots pairs  $(\mathbf{x}_i, \mathbf{y}_i \equiv \mathbb{A}\mathbf{x}_i)$ .
- Hermitian  $n \times n$  positive definite  $M$  that defines the inner product.  $M = LL^*$ ,  $L = \text{chol}(M)$ .
- Tolerance level  $\epsilon$  for numerical rank

- 1:  $[\tilde{U}_k, L, [\hat{U}_k]] = \text{Weighted\_POD}(\mathbf{X}_m; M; \epsilon)$
- 2:  $\hat{S}_k = U_k^* L^* \mathbf{Y}_m V_k \Sigma_k^{-1} \{ \text{Weighted Rayleigh quotient} \}$
- 3:  $[W_k, \Lambda_k] = \text{eig}(\hat{S}_k) \{ \Lambda_k = \text{diag}(\lambda_i)_{i=1}^k; \hat{S}_k W_k(:, i) = \lambda_i W_k(:, i); \|W_k(:, i)\|_2 = 1 \}$
- 4:  $Z_k = \hat{U}_k W_k \{ \text{Ritz vectors} \}$

**Output:**  $Z_k, \Lambda_k$

Changing physical units changes the condition number and the rank!? For more details see [Drmač+Mezić+Mohr].

## ... rewind ... Beautiful structure and bad news

The spectral decomposition of  $C_m$  has beautiful structure. Assume for simplicity that the eigenvalues  $\lambda_i$ ,  $i = 1, \dots, m$ , are algebraically simple. It is easily checked that the spectral decomposition of  $C_m$  reads

$$C_m = \mathbb{V}_m^{-1} \Lambda_m \mathbb{V}_m, \text{ where } \Lambda_m = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{pmatrix}, \quad \mathbb{V}_m = \begin{pmatrix} 1 & \lambda_1 & \dots & \lambda_1^{m-1} \\ 1 & \lambda_2 & \dots & \lambda_2^{m-1} \\ \vdots & \vdots & \dots & \vdots \\ 1 & \lambda_m & \dots & \lambda_m^{m-1} \end{pmatrix}.$$

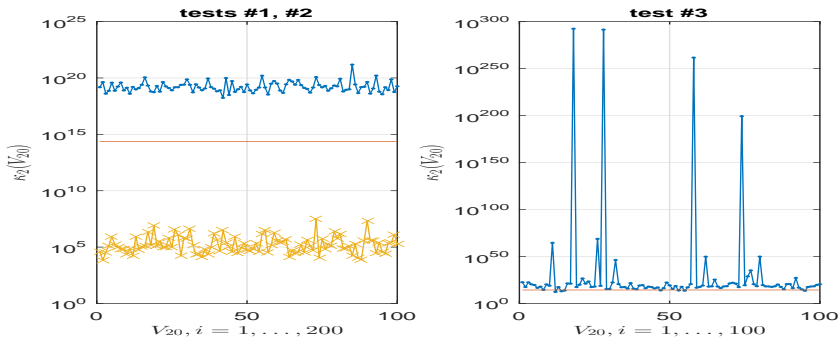
The Ritz vectors are the columns of  $Z_m = \mathbf{X}_m \mathbb{V}_m^{-1}$ .

**Bad news: The Vandermonde matrix  $\mathbb{V}_m$  is ill-conditioned!**

The condition number  $\kappa_2(\mathbb{V}_m) \equiv \|\mathbb{V}_m\|_2 \|\mathbb{V}_m^{-1}\|_2$  of any  $100 \times 100$  real Vandermonde matrix is larger than  $3 \cdot 10^{28}$ ,  
 $(\kappa_2(\mathbb{V}_m) \geq 2^{m-2}/\sqrt{m}, m = 100, [\text{Gautschi}])$ .

Better: Schmid's DMD – compute the Rayleigh quotient using a POD (truncated SVD) basis of  $\mathbf{X}_m$ .

# Ill-conditioning of Vandermonde matrices: examples



**Figure:** The spectral condition number over three sets of the total of 300 Vandermonde matrices of dimension  $m = 20$ ,  $V_{20}(\lambda_i)$ ;  $(\lambda_i) = \text{eig}(A)$ . *Left panel:* First, 100 matrices are generated in Matlab as  $A = \text{rand}(m, m)$ ,  $A = A/\max(\text{abs}(\text{eig}(A)))$ . Then, 100 matrices are generated as  $A = \text{randn}(m, m)$ ,  $A = A/\max(\text{abs}(\text{eig}(A)))$ . *Right panel:* 100 samples of  $V(\lambda_i)$  are generated using the eigenvalues of  $A = \exp(-\text{inv}(\text{rand}(m, m)))$ ,  $A = A/\max(\text{abs}(\text{eig}(A)))$ . The horizontal line marks  $1/(m\epsilon) \approx 2.25 \times 10^{14}$ .



# Vandermonde $\times$ DFT = Cauchy; DFT = Vandermonde

Let  $\mathbb{F}$  denote the DFT matrix,  $\mathbb{F}_{ij} = \omega^{(i-1)(j-1)} / \sqrt{m}$ , where  $\omega = e^{2\pi i/m}$ ,  $i = \sqrt{-1}$ . Now, recall that DFT transforms Vandermonde into Cauchy matrices as follows: If  $\lambda_i^m \neq 1$ , then

$$(\mathbb{V}_m \mathbb{F})_{ij} = \left[ \frac{\lambda_i^m - 1}{\sqrt{m}} \right] \left[ \frac{1}{\lambda_i - \omega^{1-j}} \right] [\omega^{1-j}] \equiv (\mathcal{D}_1)_{ii} \mathcal{C}_{ij} (\mathcal{D}_2)_{jj}, \quad 1 \leq j \leq m. \quad (1)$$

If  $\lambda_i = \omega^{1-j}$  for some index  $j$ , write  $\lambda_i^m - 1 = \prod_{k=1, k \neq j}^m (\lambda_i - \omega^{1-k})$  and replace (1) with the equivalent formula for the  $i$ -th row

$$(\mathbb{V}_m \mathbb{F})_{ij} = \underbrace{\frac{1}{\sqrt{m}}}_{(\mathcal{D}_1)_{ii}} \prod_{\substack{k=1 \\ k \neq j}}^m (\lambda_i - \omega^{1-k}) \underbrace{\omega^{1-j}}_{(\mathcal{D}_2)_{jj}}, \quad (\mathbb{V}_m \mathbb{F})_{ik} = 0 \text{ for } k \neq j. \quad (2)$$

This is the starting point for accurate computation of the SVD of  $\mathbb{V}_m$  [Demmel]. We use it for  $Z_m = \mathbf{X}_m \mathbb{V}_m^{-1} \equiv (\mathbf{X}_m \mathbb{F})(\mathbb{V}_m \mathbb{F})^{-1}$ .

# How this transforms $\mathbb{A}\mathbf{X}_m = \mathbf{X}_m C_m + r_{m+1} e_m^T$ ?

It is interesting to see how this transformation  $\mathbb{V}_m \mapsto \mathbb{V}_m \mathbb{F}$  fits the framework of the Krylov decomposition  $\mathbb{A}\mathbf{X}_m = \mathbf{X}_m C_m + r_{m+1} e_m^T$ , where  $C_m = \mathbb{V}_m^{-1} \Lambda_m \mathbb{V}_m$ . Post-multiply this with  $\mathbb{F}$  to obtain

$$\mathbb{A}(\mathbf{X}_m \mathbb{F}) = (\mathbf{X}_m \mathbb{F}) \mathbb{F}^* (\mathbb{V}_m^{-1} \Lambda_m \mathbb{V}_m) \mathbb{F} + r_{m+1} e_m^T \mathbb{F} \quad (3)$$

and then, using  $\mathbb{V}_m \mathbb{F} = \mathcal{D}_1 \mathcal{C} \mathcal{D}_2$ ,

$$\mathbb{F}^* C_m \mathbb{F} = \mathbb{F}^* (\mathbb{V}_m^{-1} \Lambda_m \mathbb{V}_m) \mathbb{F} = \mathcal{D}_2^* \mathcal{C}^{-1} \mathcal{D}_1^{-1} \Lambda_m \mathcal{D}_1 \mathcal{C} \mathcal{D}_2 = \mathcal{D}_2^* \mathcal{C}^{-1} \Lambda_m \mathcal{C} \mathcal{D}_2 \text{ and}$$

$$\mathbb{A}(\mathbf{X}_m \mathbb{F}) = (\mathbf{X}_m \mathbb{F}) ((\mathcal{C} \mathcal{D}_2)^{-1} \Lambda_m (\mathcal{C} \mathcal{D}_2)) + r_{m+1} e_m^T \mathbb{F} \iff (4)$$

$$\mathbb{A}(\mathbf{X}_m \mathbb{F} \mathcal{D}_2^*) = (\mathbf{X}_m \mathbb{F} \mathcal{D}_2^*) (\mathcal{C}^{-1} \Lambda_m \mathcal{C}) + r_{m+1} e_m^T \mathbb{F} \mathcal{D}_2^*. \quad (5)$$

If we think of each row  $\mathbf{X}_m(i, :)$  as a time trajectory of the corresponding observable, then  $\mathbf{X}_m(i, :)\mathbb{F}$  represents its image in the frequency domain, and (4, 5) is the corresponding Krylov decomposition. Possible insightful connections (?) to the Laskar algorithm [Laskar, Arbabi+Mezić] and (data centered) DMD and temporal DFT see [Chen+Tu+Rowley].

# But $\kappa_2(\text{Vandermonde} \times \text{DFT}) = \kappa_2(\text{Vandermonde})$ ?!

Note that the matrices  $\mathcal{D}_1$ ,  $\mathcal{C}$ ,  $\mathcal{D}_2$  are given implicitly by the parameters  $\lambda_i$  (eigenvalues, available on input) and the  $m$ -th roots of unity  $\zeta_j = \omega^{1-j}$ ,  $j = 1, \dots, m$  (easily precomputed to any desired precision and tabulated), so that the DFT  $\mathbb{V}_m \mathbb{F}$  is not done by actually running an FFT. It suffices to make a note that the  $\lambda_i$ 's and the roots of unity are the parameters (original data) that define  $\mathbb{V}_m \mathbb{F}$  as in (1), (2).

## Besides nice matrix-theoretical connection, what is the gain?

Applying the DFT to  $\mathbb{V}_m$  in order to avoid the ill-conditioning of  $\mathbb{V}_m$  may seem a futile effort – since  $\mathbb{F}$  is unitary,  $\kappa_2(\mathcal{D}_1 \mathcal{C} \mathcal{D}_2) = \kappa_2(\mathbb{V}_m \mathbb{F}) = \kappa_2(\mathbb{V}_m)$ . Further, Cauchy matrices are also notoriously ill-conditioned, so, in essence, we have traded one badly conditioned structure to another one.

How bad it can be? What is the meaning of bad, ill-conditioned anyway? Is Cauchy matrix really badly conditioned?

# Caveat ill-conditioning: Hilbert matrix example

Ill-conditioning is not always obvious in the sizes of its entries – the entries of the innocuous-looking  $100 \times 100$  Hilbert matrix  $H_{100}$  range from  $1/199 \approx 5.025 \cdot 10^{-3}$  to 1, and  $\kappa_2(H_{100}) > 10^{150}$ .

condition number(condition number)=condition number [Higham]

```
>> cond(hilb(100))
ans = 4.6226e+19
```

If the computed/estimated condition number is above  $1/\text{eps}$  (in Matlab,  $1/\text{eps}=4.503599627370496\text{e}+15$ ), it might be a severe underestimate. This may lead to an underestimate of extra precision needed to handle the ill-conditioning.

High accuracy numerical linear algebra :)

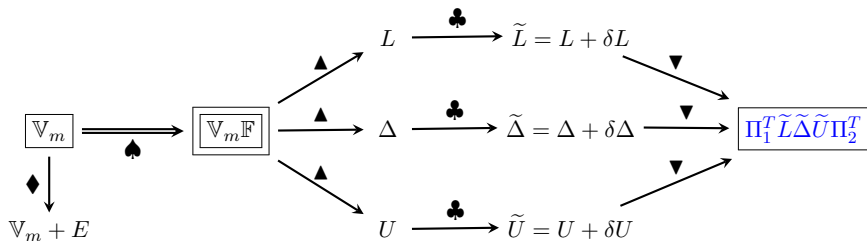
Accurate LU, QR, SVD of any Cauchy or Vandermonde matrix is feasible without higher precision arithmetic. Good algorithms are available!

# SVD( $D_1 \times \text{Cauchy} \times D_2$ )

Given Cauchy matrix  $\mathcal{C}$  and any two diagonal matrices  $D_1, D_2$ , the SVD of  $G = D_1 \mathcal{C} D_2$  can be computed to nearly full precision as follows:

- 1 Compute the LDU,  $P_1 G P_2 = L D U$  using explicit determinant based formulas to update the Schur complement [Demmel]. This is entry wise forward stable computation of  $L, D, U$ . Moreover,  $\kappa(L), \kappa(U)$  are moderate. (Small  $\|\delta L\|/\|L\|, \|\delta U\|/\|U\|, |\delta D_{ii}|/|D_{ii}|$  is also OK) ( $G = \text{hilb}(100), \kappa_2(G) > 10^{150}, \kappa_2(L) = \kappa_2(U) \approx 72.24, \kappa_2(D) \approx 2.32 \cdot 10^{149}$ )
- 2 Compute the SVD of the product  $L D U$  using a Jacobi type SVD algorithm (product SVD [Drmač]). The forward error is determined by  $\max(\kappa(L), \kappa(U))$ , independent of  $\kappa_2(D)$ . The backward errors  $\|\Delta L\|/\|L\|, \|\Delta U\|/\|U\|, \Delta D_{ii}/D_{ii}$  are small.

The key is in forward stable reparametrization, so that the new representation is well-conditioned (for particular algorithm).



**Figure:** Legend:  $\spadesuit$  = the DFT of  $\mathbb{V}_m$  using the explicit formulas (1);  $\blacktriangle$  = the forward stable pivoted LDU [Demmel];  $\clubsuit$  = forward errors in the computed factors  $\tilde{L}$ ,  $\tilde{\Delta}$ ,  $\tilde{U}$ ;  $\blacktriangledown$  = implicit representation of  $\mathbb{V}_m \mathbb{F}$  as the product  $\Pi_1^T \tilde{L} \tilde{\Delta} \tilde{U} \Pi_2^T$ ;  $\blacklozenge$  = direct computation with  $\mathbb{V}_m$ , using standard algorithms, produces backward error  $E$  that is small in matrix norm, and the condition number is  $\kappa_2(\mathbb{V}_m)$ .

$$|\tilde{L}_{ij} - L_{ij}| \leq \epsilon |L_{ij}|, \quad |\tilde{\Delta}_{ii} - \Delta_{ii}| \leq \epsilon |\Delta_{ii}|, \quad |\tilde{U}_{ij} - U_{ij}| \leq \epsilon |U_{ij}|, \quad (6)$$

$$\widehat{W}_m = (((((\mathbf{X}_m \mathbb{F}) \Pi_2) U^{-1}) \Delta^{-1}) L^{-1}) \Pi_1. \quad (7)$$

See [Demmel], [Demmel+Koev], [Dopico+Molera] for error analysis.

# Matlab code for $Z_m = X_m V_m^{-1}$

```
function Z = X_inv_Vandermonde( lambda, X )
% X_inv_Vandermonde computes Z = X*inv(V(lambda)), where X has m
% columns and V(lambda)=fliplr(vander(lambda)) is the m x m
% Vandermonde matrix defined by the m x 1 vector lambda;
% V(lambda)_{ij} = lambda(i)^(j-1), i,j=1,...,m.
%.....
% Coded by Zlatko Drmac, drmac@math.hr.
%.....
%
m = length(lambda) ;
% .. pivoted LDU of V(lambda)*DFT ; p1, p2 permutations
[ L, D, U, p1, p2 ] = Vand_DFT_LDU( lambda, m, 'LDU' ) ;
Z = ifft(X,[],2) ;
Z = ( ( Y(:,p2) / U ) * diag(sqrt(m)./D) ) / L ;
pli(p1) = 1:m ; Z = Z(:,pli) ; % pli is the inverse of p1
end
```

- Not as simple as  $Z_m = X_m/V_m$
- More accurate than  $Z_m = X_m/V_m$ ; independent of the distribution of the  $\lambda_i$ 's

# Numerical stress test drive: Q2D Kolmogorov flow

## Example

We use the simulation data of a 2D model obtained by depth averaging the Navier–Stokes equations for a shear flow in a thin layer of electrolyte suspended on a thin lubricating layer of a dielectric fluid; see [Tithof+et al], [Suri+et al] for more detailed description of the experimental setup.<sup>a</sup> The (scalar) vorticity field data consists of  $n_t$  snapshots of dimensions  $n_x \times n_y$ ; in this particular example  $n_t = 1201 \equiv m + 1$ ,  $n_x = n_y = 128$ . The  $n_x \times n_y \times n_t$  tensor is matricized into  $n_x \cdot n_y \times n_t$  matrix of snapshots  $(\mathbf{f}_1, \dots, \mathbf{f}_{n_t})$ , and  $\mathbf{X}_m$  is of dimensions  $16384 \times 1200$ .

<sup>a</sup>We thank Michael Schatz, Balachandra Suri, Roman Grigoriev and Logan Kageorge from the Georgia Institute of Technology for providing us with the data.

This is a good stress test because  $\kappa_2(\mathbb{V}_m) > 10^{76}$

Want to show that we can use  $C_m, \mathbb{V}_m$  in working precision (IEEE 64 bit) despite the fact that  $\kappa_2(\mathbb{V}_m) > 10^{76} \gg 1/\text{roundoff}_{64} \approx 4.5 \cdot 10^{15}$



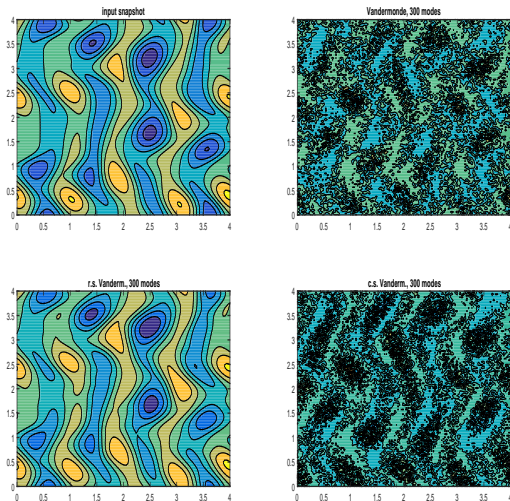
# Test the reconstruction potential

## Reconstructing snapshots using selected modes

$$\mathbf{f}_i \approx \sum_{j=1}^{\ell} z_{\zeta_j} \alpha_j \lambda_{\zeta_j}^{i-1} \equiv \sum_{j=1}^{\ell} z_{\zeta_j} \alpha_j |\lambda_{\zeta_j}|^{i-1} e^{i\omega_{\zeta_j} (i-1)\delta t}, \quad i = 1, \dots, m.$$

where, for simplicity, the modes are selected by taking given number of modes with absolutely largest amplitudes  $|\alpha_j|$  (*dominant modes*).

- ❶ inversion of the Vandermonde matrix by the backslash operator in Matlab ;  $\kappa_2(\mathbb{V}_m) > 10^{76} \gg 1/\text{roundoff}_{64} \approx 4.5 \cdot 10^{15}$
- ❷ inversion of the row scaled Vandermonde matrix by the backslash operator in Matlab:  $\mathbb{V}_m = D_r \mathbb{V}_m^{(r)}$ ,  $\widehat{\mathbb{W}}_m = (\mathbf{X}_m (\mathbb{V}_m^{(r)})^{-1}) D_r^{-1}$ , where  $D_r = \text{diag}(\|\mathbb{V}_m(i, :)\|)_{i=1}^m$  ;  $\kappa_2(\mathbb{V}_m^{(r)}) \approx 3.1 \cdot 10^7 \approx 0.45 \cdot 1/\text{roundoff}_{32}$
- ❸ inversion of the column scaled Vandermonde matrix by the backslash operator in Matlab:  $\mathbb{V}_m = \mathbb{V}_m^{(c)} D_c$ ,  $\widehat{\mathbb{W}}_m = (\mathbf{X}_m D_c^{-1}) (\mathbb{V}_m^{(c)})^{-1}$ , where  $D_c = \text{diag}(\|\mathbb{V}_m(:, i)\|)_{i=1}^m$ .  $\kappa_2(\mathbb{V}_m^{(c)}) \approx 3.0 \cdot 10^{21}$

$f_{321}$  with 300 modes; similar results for other  $f_i$ 's

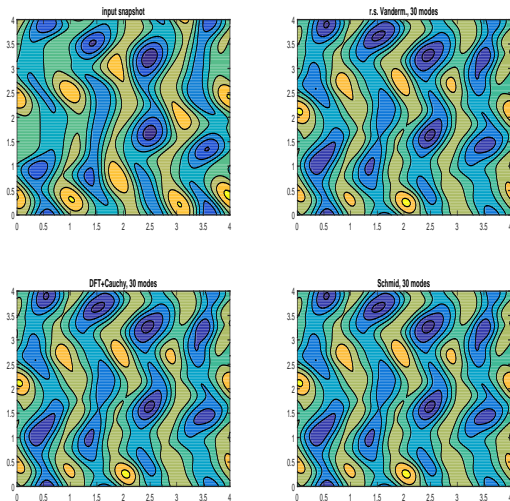
**Figure:** Reconstruction of  $f_{321}$  using 300 dominant modes. The linear systems are solved in Matlab using the backslash operator. Note how using the row scaled  $\mathbb{V}_m^{(r)}$  improves the reconstruction (first plot in the second row), while backslashing the original  $\mathbb{V}_m$  and the column scaled matrix  $\mathbb{V}_m^{(c)}$  yields poor results (plots in the second column on the Figure).

# Björck-Pereyera, DMD, DFT+Cauchy

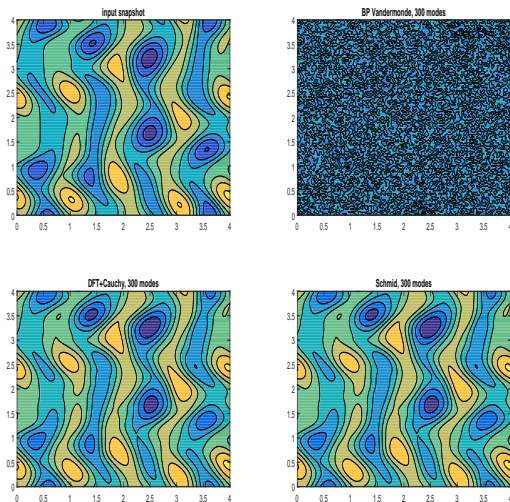
We now test the following three methods:

- ① Companion matrix formulation with the Björck-Pereyera method for Vandermonde systems. Although forward stable in the special case of real and ordered  $\lambda_i$ 's, this method may be very sensitive in the case of general complex  $\lambda_i$ 's and relatively large dimension  $m$ .
- ② Companion matrix formulation with the DFT and inversion of the Cauchy matrix. Since  $\mathbb{F}$  and  $\mathcal{D}_2$  in (1) are unitary, the algorithm solves linear system with the matrix  $\mathcal{D}_1\mathcal{C} = \mathbb{V}_m\mathbb{F}\mathcal{D}_2^*$  of condition number bigger than  $10^{76}$ . No additional scaling is used; we want to illustrate the claim that such high condition number cannot spoil the result.
- ③ Schmid's DMD method. Here we expect good reconstruction results, provided it is feasible for given data and the parameters. The SVD is not truncated because  $\kappa_2(\mathbf{X}_m) \approx 5.5 \cdot 10^{10}$  ( $\sigma_{\max}(\mathbf{X}_m) \approx 4.2 \cdot 10^3$ ,  $\sigma_{\min}(\mathbf{X}_m) \approx 7.5 \cdot 10^{-8}$ ).

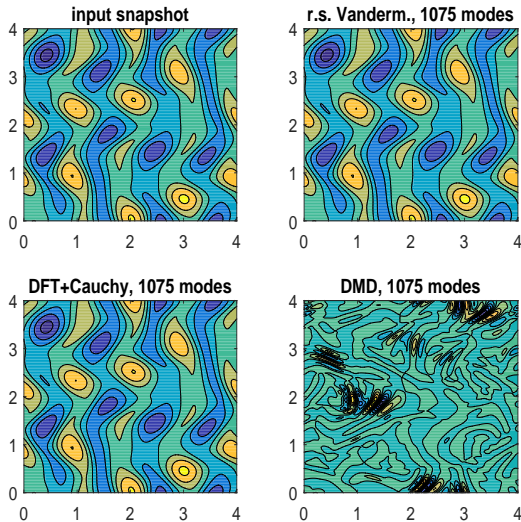
# Reconstructing $f_{321}$ with 30 modes



**Figure:** Reconstruction of  $f_{321}$  using 30 dominant modes. The row scaled Vandermonde and the DFT+Cauchy inversion, as well as the Schmid's DMD reconstruction (second row) succeeded in reconstructing  $f_{321}$  using 30 modes with dominant amplitudes.



**Figure:** Reconstruction of  $f_{321}$  using 300 dominant modes. Björck-Pereyera method (second plot in the first row) failed to produce any useful data. The DFT+Cauchy inversion and the Schmid's DMD reconstruction (second row) succeeded in reconstructing  $f_{321}$  pretty much using 300 modes with dominant amplitudes.



**Figure:** Reconstruction of  $f_{1111}$  using 1075 dominant modes. The DFT+Cauchy inversion did well, and the Schmid's DMD reconstruction (second row) unfortunately failed. (With 1074 modes DMD performed well, but starting with 1075 all reconstruction failed, including the one with all 1200 modes.)

# Snapshot reconstruction – some theory

Wanted are the coefficients  $\alpha_1, \dots, \alpha_\ell$  that minimize

$$\sum_{i=1}^m \left\| \mathbf{f}_i - \sum_{j=1}^{\ell} z_j \alpha_j \lambda_j^{i-1} \right\|_2^2 \longrightarrow \min. \quad (1)$$

If we set  $Z_\ell = (z_1, \dots, z_\ell)$ ,  $\vec{\alpha} = (\alpha_1, \dots, \alpha_\ell)^T$ ,  $\Delta_\alpha = \text{diag}(\vec{\alpha})$ ,  $\Lambda_j = (\lambda_1^{j-1}, \dots, \lambda_\ell^{j-1})^T$ , and  $\Delta_{\Lambda_j} = \text{diag}(\Lambda_j)$  then the objective (1) reads

$$\Omega^2(\alpha) \equiv \left\| \mathbf{X}_m - Z_\ell \Delta_\alpha \begin{pmatrix} \Lambda_1 & \Lambda_2 & \dots & \Lambda_m \end{pmatrix} \right\|_F^2 \longrightarrow \min. \quad (2)$$

Compute the tall QR factorization  $Z_\ell = QR$  and define projected snapshots  $\mathbf{g}_i = Q^* \mathbf{f}_i$ , then the LS problem can be compactly written as

$$\|\vec{\mathbf{g}} - S\vec{\alpha}\|_2 \longrightarrow \min, \quad \vec{\mathbf{g}} = \begin{pmatrix} \mathbf{g}_1 \\ \vdots \\ \mathbf{g}_m \end{pmatrix}, \quad S = (I_m \otimes R) \begin{pmatrix} \Delta_{\Lambda_1} \\ \vdots \\ \Delta_{\Lambda_m} \end{pmatrix} \equiv \begin{pmatrix} R\Delta_{\Lambda_1} \\ \vdots \\ R\Delta_{\Lambda_m} \end{pmatrix}. \quad (3)$$

# Solution(s) by generalized inverse(s)

The optimal solution is obtained as  $\vec{\alpha} = S^\dagger \vec{g}$  using the explicit normal equations formula, DMDSP [Jovanović+Schmid+Nichols]. Here  $S^\dagger$  is the Moore-Penrose generalized inverse.

On the other hand, deploying the reflexive g-inverse of  $S$ ,

$$S^- = \begin{pmatrix} \Delta_{\Lambda_1} \\ \vdots \\ \Delta_{\Lambda_m} \end{pmatrix}^\dagger (I \otimes R^{-1}),$$

we obtain interesting explicit formulas for  $\vec{\alpha}_* = S^- \vec{g}$  that reveal a relation to the Generalized Laplace Analysis, GLA, [Mezić+Mohr].

$$S^- \neq S^\dagger$$

Recall that, by definition,  $S^-$  satisfies

$SS^-S = S$ ,  $S^-SS^- = S^-$  and  $S^-S = (S^-S)^*$ . In fact, since  $SS^- \neq (SS^-)^*$ , we have in general that  $S^- \neq S^\dagger$ , so  $\vec{\alpha}_* \neq S^\dagger \vec{g}$ .



$\vec{\alpha}_\star = S^{-1}\vec{g}$  solves a weighted LS problem

$$\begin{aligned}\vec{\alpha}_\star &= \begin{pmatrix} \Delta_{\Lambda_1} \\ \vdots \\ \Delta_{\Lambda_m} \end{pmatrix}^\dagger (I \otimes R^{-1}) \begin{pmatrix} \mathbf{g}_1 \\ \vdots \\ \mathbf{g}_m \end{pmatrix} = \left( \sum_{k=1}^m \Delta_{\Lambda_k}^* \Delta_{\Lambda_k} \right)^{-1} \sum_{i=1}^m \Delta_{\Lambda_i}^* (R^{-1} \mathbf{g}_i) \\ &= \sum_{i=1}^m \begin{pmatrix} \frac{\bar{\lambda}_1^{i-1}}{\sum_{k=1}^m |\lambda_1|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_1 \\ \frac{\bar{\lambda}_2^{i-1}}{\sum_{k=1}^m |\lambda_2|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_2 \\ \vdots \\ \frac{\bar{\lambda}_\ell^{i-1}}{\sum_{k=1}^m |\lambda_\ell|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_\ell \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m \frac{\bar{\lambda}_1^{i-1}}{\sum_{k=1}^m |\lambda_1|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_1 \\ \sum_{i=1}^m \frac{\bar{\lambda}_2^{i-1}}{\sum_{k=1}^m |\lambda_2|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_2 \\ \vdots \\ \sum_{i=1}^m \frac{\bar{\lambda}_\ell^{i-1}}{\sum_{k=1}^m |\lambda_\ell|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_\ell \end{pmatrix}\end{aligned}$$

## Theorem

Let  $M = I \otimes (RR^*)^{-1}$ ,  $(x, y)_M = y^* M x$ , and let  $\|x\|_M = \sqrt{x^* M x}$ . Then  $\vec{\alpha}_\star$  is the minimum  $\|\cdot\|_2$ -norm solution of the weighted least squares problem

$$\|\vec{g} - S\vec{\alpha}\|_M \longrightarrow \min. \quad (4)$$

# Comparison with GLA reconstruction [Mezić+Mohr]

$$\vec{\alpha}_\star = \begin{pmatrix} \sum_{i=1}^m \frac{\bar{\lambda}_1^{i-1}}{\sum_{k=1}^m |\lambda_1|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_1 \\ \sum_{i=1}^m \frac{\bar{\lambda}_2^{i-1}}{\sum_{k=1}^m |\lambda_2|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_2 \\ \vdots \\ \sum_{i=1}^m \frac{\bar{\lambda}_\ell^{i-1}}{\sum_{k=1}^m |\lambda_\ell|^{2(k-1)}} (R^{-1} \mathbf{g}_i)_\ell \end{pmatrix}$$

$$\vec{\alpha}_{(GLA)} = \frac{1}{m} \sum_{i=1}^m \Lambda_\ell^{-i+1} R^{-1} \mathbf{g}_i = \begin{pmatrix} \frac{1}{m} \sum_{i=1}^m \lambda_1^{-i+1} (R^{-1} \mathbf{g}_i)_1 \\ \frac{1}{m} \sum_{i=1}^m \lambda_2^{-i+1} (R^{-1} \mathbf{g}_i)_2 \\ \vdots \\ \frac{1}{m} \sum_{i=1}^m \lambda_\ell^{-i+1} (R^{-1} \mathbf{g}_i)_\ell \end{pmatrix}.$$

## Proposition

When the spectrum of  $\mathbb{A}$  lies on the unit circle,  $\vec{\alpha}_\star = \vec{\alpha}_{(GLA)}$ .

For more see recent paper [Drmač+Mezić+Mohr].

# Concluding remarks

- We have presented modifications of the DMD algorithm, together with theoretical analysis and justification, discussion of the potential weaknesses of the original method, and examples that illustrate the advantages of the new proposed method. From the point of view of numerical linear algebra, the deployed techniques are not new; however, the novelty is in adapting them to the data driven setting and turning the DMD into a more powerful tool.
- Using high accuracy numerical linear algebra techniques we were able to curb the ill-conditioning of the companion matrix's associated Vandermonde matrix allowing us to invert it and find the DMD modes.
- In addition to the inherent elegance in terms of the companion matrix formulation of DMD, we provide a numerical linear algebra framework that has a close connection to Koopman operator theory, as well as to the results coming from Generalized Laplace Analysis theory.

# For more on this and complete references list see



Z. Drmač, I. Mezić, and R. Mohr.

Data driven modal decompositions: analysis and enhancements.

*SIAM Journal on Scientific Computing*, 40(4):A2253–A2285, 2018.



Z. Drmač, I. Mezić, and R. Mohr.

Data driven Koopman spectral analysis in Vandermonde–Cauchy form via the DFT: numerical method and theoretical insights.

*ArXiv e-prints*, August 2018.



Z. Drmač, I. Mezić, and R. Mohr.

On least squares problem with certain Vandermonde–Khatri–Rao structure with applications to DMD.

*ArXiv e-prints*, 2018.

# Few more comments on the numerical aspects of $S^\dagger \vec{g}$

For given  $(\lambda_j, z_j)$ 's and nonnegative weights  $\mathfrak{w}_i$ , find the  $\alpha_j$ 's to achieve

$$\sum_{i=1}^m \mathfrak{w}_i^2 \|\mathbf{f}_i - \sum_{j=1}^{\ell} z_j \alpha_j \lambda_j^{i-1}\|_2^2 \longrightarrow \min. \quad (1)$$

Set  $\mathbf{W} = \text{diag}(\mathfrak{w}_i)_{i=1}^m$ . The weights  $\mathfrak{w}_i > 0$  are used to emphasize snapshots whose reconstruction is more important. Let  $\mathbf{\Lambda} = \text{diag}(\lambda_j)_{j=1}^{\ell}$ ,

$$\Delta_{\alpha} = \begin{pmatrix} \alpha_1 & 0 & \cdot & 0 \\ 0 & \alpha_2 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 \\ 0 & \cdot & 0 & \alpha_{\ell} \end{pmatrix}, \quad \Lambda_i = \begin{pmatrix} \lambda_1^{i-1} \\ \lambda_2^{i-1} \\ \cdot \\ \lambda_{\ell}^{i-1} \end{pmatrix}, \quad \Delta_{\Lambda_i} = \begin{pmatrix} \lambda_1^{i-1} & 0 & \cdot & 0 \\ 0 & \lambda_2^{i-1} & \cdot & \cdot \\ \cdot & \cdot & \cdot & 0 \\ 0 & \cdot & 0 & \lambda_{\ell}^{i-1} \end{pmatrix} \equiv \mathbf{\Lambda}^{i-1},$$

and write the objective (1) as the function of  $\alpha = (\alpha_1, \dots, \alpha_{\ell})^T$ ,

$$\Omega^2(\alpha) \equiv \|\mathbf{X}_m - Z_{\ell} \Delta_{\alpha} (\Lambda_1 \quad \Lambda_2 \quad \dots \quad \Lambda_m)\|_F^2 \longrightarrow \min, \quad (2)$$

$$(\Lambda_1 \quad \Lambda_2 \quad \dots \quad \Lambda_m) = \begin{pmatrix} 1 & \lambda_1 & \dots & \lambda_1^{m-1} \\ 1 & \lambda_2 & \dots & \lambda_2^{m-1} \\ \vdots & \vdots & \dots & \vdots \\ 1 & \lambda_{\ell} & \dots & \lambda_{\ell}^{m-1} \end{pmatrix} \equiv \mathbb{V}_{\ell, m} \in \mathbb{C}^{\ell \times m}. \quad (3)$$

# Explicit normal equations solution

$$\|(\mathbf{W} \otimes \mathbb{I}_\ell) [\vec{\mathbf{g}} - S\boldsymbol{\alpha}]\|_2 \longrightarrow \min, \quad \text{where } \vec{\mathbf{g}} = \begin{pmatrix} \mathbf{g}_1 \\ \vdots \\ \mathbf{g}_m \end{pmatrix}, \quad S = (\mathbb{I}_m \otimes R) \begin{pmatrix} \Delta_{\Lambda_1} \\ \vdots \\ \Delta_{\Lambda_m} \end{pmatrix}$$

## Theorem

*With the notation as above, the unique solution  $\boldsymbol{\alpha}$  of the LS problem (1) is*

$$\boldsymbol{\alpha} = [(R^*R) \circ (\overline{\mathbb{V}_{\ell,m} \mathbf{W}^2 \mathbb{V}_{\ell,m}^*})]^{-1} [(\overline{\mathbb{V}_{\ell,m} \mathbf{W}} \circ (R^*G\mathbf{W}))\mathbf{e}], \quad (4)$$

*where  $G = (\mathbf{g}_1 \quad \dots \quad \mathbf{g}_m)$ ,  $\mathbf{e} = (1 \quad \dots \quad 1)^T$ . In terms of  $\mathbf{X}_m, Z_\ell$ ,*

$$\boldsymbol{\alpha} = [(Z_\ell^* Z_\ell) \circ (\overline{\mathbb{V}_{\ell,m} \mathbf{W}^2 \mathbb{V}_{\ell,m}^*})]^{-1} [(\overline{\mathbb{V}_{\ell,m} \mathbf{W}} \circ (Z_\ell^* \mathbf{X}_m \mathbf{W}))\mathbf{e}]. \quad (5)$$

This includes the DMDSP of [Jovanović+et al] and solution for scattering coefficients in multistatic antenna array processing [Lev-Ari] as unweighted cases. **Are normal equations safe to use? Let us experiment with a small dimension example.**

# Squaring the conditon number – loosing definiteness

Let  $\mathbf{W} = \mathbf{I}$ . Let  $\ell = 3$ ,  $m = 4$ ,  $\xi = \sqrt{\epsilon}$ ,  $\lambda_1 = \xi$ ,  $\lambda_2 = 2\xi$ ,  $\lambda_3 = 0.2$ , so that the Vandermonde section  $\mathbb{V}_{\ell,m}$  equals

$$\mathbb{V}_{\ell,m} = \begin{pmatrix} 1 & 1.490116119384766e-08 & 2.220446049250313e-16 & 3.308722450212111e-24 \\ 1 & 2.980232238769531e-08 & 8.881784197001252e-16 & 2.646977960169689e-23 \\ 1 & 2.000000000000000e-01 & 4.000000000000001e-02 & 8.000000000000002e-03 \end{pmatrix},$$

$$R = \begin{pmatrix} 1 & 1 & 1 \\ 0 & \xi/2 & \xi \\ 0 & 0 & \xi \end{pmatrix} = \begin{pmatrix} 1 & 1.000000000000000e+00 & 1.000000000000000e+00 \\ 0 & 7.450580596923828e-09 & 1.490116119384766e-08 \\ 0 & 0 & 1.490116119384766e-08 \end{pmatrix}.$$

Here  $\kappa_2(\mathbb{V}_{\ell,m}) \approx 10^9$ ,  $\kappa_2(R) \approx 10^9 \ll 1/\text{roundoff}_{64} \approx 4.5 \cdot 10^{15}$ .

<code>&gt;&gt; chol(Vlm*Vlm')</code>	<code>&gt;&gt; chol((R'*R).*(Vlm*Vlm'))</code>
Error using chol	Error using chol
Matrix must be ....	Matrix must be positive definite.

<code>&gt;&gt; chol(R'*R)</code>	Normal equations matrix is
Error using chol	not definite!
Matrix must be positive definite.	

# Indefinite $\circ$ Indefinite = Positive Definite ?!

Use the same  $\mathbb{V}_{\ell,m}$  but change the definition of  $R$  to

$$R = \begin{pmatrix} 1 & 1 & 1 \\ 0 & \xi & \xi \\ 0 & 0 & \xi/2 \end{pmatrix} = \begin{pmatrix} 1 & 1.0000000000000000e+00 & 1.0000000000000000e+00 \\ 0 & 1.490116119384766e-08 & 1.490116119384766e-08 \\ 0 & 0 & 7.450580596923828e-09 \end{pmatrix}.$$

If we repeat the experiment with the Cholesky factorizations, we obtain

```
>> chol(Vlm*Vlm')
```

Error using chol

Matrix must be positive definite.

```
>> chol(R'*R)
```

Error using chol

Matrix must be positive definite.

```
>> TC = chol((R'*R).*(Vlm*Vlm'))
```

TC =

1	1.0000000000000000e+00	1.0000000002980232e+00
0	1.490116119384765e-08	1.999999880790710e-01
0	0	4.079214149695062e-02



# How accurately we can solve with

$$C = (R' * R) .* (V_{lm} * V_{lm}')$$

Based on [Demmel], we know that floating point Cholesky factorization  $C = LL^*$  ( $L$  lower triangular with positive diagonal) of  $C$  is feasible if the matrix  $C_s = (c_{ij} / \sqrt{c_{ii}c_{jj}})_{i,j=1}^{\ell}$  is well conditioned. Further, if we solve the linear system  $Cx = b \neq 0$  using the Cholesky factor in the forward and backward substitutions, then the computed solution  $\tilde{x}$  satisfies

$$\frac{\|D_C(\tilde{x} - C^{-1}b)\|_2}{\|D_C\tilde{x}\|_2} \leq g(\ell)\epsilon\kappa_2(C_s), \quad (6)$$

where  $g(\ell)$  is modest function of the dimension,  $D_C = \text{diag}(\sqrt{c_{ii}})_{i=1}^{\ell}$ . Note that this implies component-wise error bound for each  $\tilde{x}_i \neq 0$ :

$$\frac{|\tilde{x}_i - (C^{-1}b)_i|}{|\tilde{x}_i|} \leq \underbrace{\left[ \frac{\|D_C\tilde{x}\|_2}{\sqrt{c_{ii}}|\tilde{x}_i|} \right]}_{\geq 1} g(\ell)\epsilon\kappa_2(C_s). \quad (7)$$

In the Q2D Kolmogorov flow example,  $\kappa_2(C) > 10^{87} \gg \kappa_2(C_s) \approx 8.5+01$ .

## Theorem

Let  $A$  and  $B$  be Hermitian positive semidefinite matrices with positive diagonal entries, and let  $C = A \circ B$ . If  $A_s = (a_{ij}/\sqrt{a_{ii}a_{jj}})$ ,  $B_s = (b_{ij}/\sqrt{b_{ii}b_{jj}})$ ,  $C_s = (c_{ij}/\sqrt{c_{ii}c_{jj}})$ , then

$$\max(\lambda_{\min}(A_s), \lambda_{\min}(B_s)) \leq \lambda_i(C_s) \leq \min(\lambda_{\max}(A_s), \lambda_{\max}(B_s)). \quad (8)$$

In particular,  $\|C_s^{-1}\|_2 \leq \min(\|A_s^{-1}\|_2, \|B_s^{-1}\|_2)$  and  $\kappa_2(C_s) \leq \min(\kappa_2(A_s), \kappa_2(B_s))$ . If  $A$  or  $B$  is diagonal, all inequalities in this theorem become equalities.

## Corollary

Let  $C \equiv (R^*R) \circ (\overline{\mathbb{V}_{\ell,m} \mathbf{W}^2 \mathbb{V}_{\ell,m}^*})$ ,  $C_s = (c_{ij}/\sqrt{c_{ii}c_{jj}})$ . Further, let  $R = R_c \Delta_r$  and  $\mathbb{V}_{\ell m} \mathbf{W} = \Delta_v (\mathbb{V}_{\ell m} \mathbf{W})_r$  with diagonal scaling matrices  $\Delta_r$  and  $\Delta_v$  such that  $R_c$  has unit columns and  $(\mathbb{V}_{\ell m} \mathbf{W})_r$  has unit rows (in Euclidean norm). Then

$$\kappa_2(C_s) \leq \min(\kappa_2(R_c)^2, \kappa_2((\mathbb{V}_{\ell,m} \mathbf{W})_r)^2).$$

$$\|\vec{g} - S\alpha\|_2 \longrightarrow \min; S = Q_S R_S, \alpha = R_S^{-1}(Q_S^* \vec{g})$$

$$\alpha = R_S^{-1}(R_S^{-*}(S^* \vec{g})), r = \vec{g} - S\alpha$$

$$\delta\alpha = R_S^{-1}(R_S^{-*}(S^* r)), \alpha_* = \alpha + \delta\alpha \quad (9)$$

### Algorithm Corrected semi-normal solution

**Input:**  $R, \Lambda, G, S$

**Output:** Corrected solution  $\alpha_*$

- 1: Compute the triangular factor  $R_S$  in the QR factorization of  $S$ .
- 2:  $g_S = [(\overline{\nabla_{\ell,m}} \circ (R^* G))\mathbf{e}]$  {Note,  $g_S = S^* \vec{g}$ . Use xTRMM from BLAS 3.}
- 3:  $\alpha = R_S^{-1}(R_S^{-*} g_S)$  {Use xTRSM or xTRTRS or xTRSV from LAPACK.}
- 4:  $r_{\square} = G - R(\alpha \quad \Lambda\alpha \quad \Lambda^2\alpha \quad \dots \quad \Lambda^{m-1}\alpha) \equiv G - R\text{diag}(\alpha)\nabla_{\ell,m}$
- 5:  $r_S = [(\overline{\nabla_{\ell,m}} \circ (R^* r_{\square}))\mathbf{e}]$  {Note,  $r_S = S^* r$ . Use xTRMM from BLAS 3.}
- 6:  $\delta\alpha = R_S^{-1}(R_S^{-*} r_S)$  {Use xTRSM or xTRTRS or xTRSV from LAPACK.}
- 7:  $\alpha_* = \alpha + \delta\alpha$

Considerably improves over normal equations, but needs QRF of  $S$ .

# Algorithm: Recursive QR factorization of $S$ for $m = 2^p$

**Input:** Upper triangular  $R \in \mathbb{C}^{\ell \times \ell}$ ; diagonal  $\Lambda \in \mathbb{C}^{\ell \times \ell}$ ; number of snapshots  $m = 2^p$

**Output:** Upper triangular QR factor  $R_S = T_p$  of  $S \in \mathbb{C}^{2^p \ell \times \ell}$

$\boxed{T_4}$	$\leftarrow \boxed{T_3}$	$\leftarrow \boxed{T_2}$	$\leftarrow \boxed{T_1}$	$\leftarrow R\Lambda^0$
0	0	0	0	$\leftarrow R\Lambda^1$
0	0	0	$\leftarrow T_1\Lambda^2$	$R\Lambda^2$
0	0	0	0	$R\Lambda^3$
0	0	$\leftarrow T_2\Lambda^4$	$T_1\Lambda^4$	$R\Lambda^4$
0	0	0	0	$R\Lambda^5$
0	0	0	$T_1\Lambda^6$	$R\Lambda^6$
0	0	0	0	$R\Lambda^7$
0	$\leftarrow T_3\Lambda^8$	$T_2\Lambda^8$	$T_1\Lambda^8$	$R\Lambda^8$
0	0	0	0	$R\Lambda^9$
0	0	0	$T_1\Lambda^{10}$	$R\Lambda^{10}$
0	0	0	0	$R\Lambda^{11}$
0	0	$T_2\Lambda^{12}$	$T_1\Lambda^{12}$	$R\Lambda^{12}$
0	0	0	0	$R\Lambda^{13}$
0	0	0	$T_1\Lambda^{14}$	$R\Lambda^{14}$
0	0	0	0	$R\Lambda^{15}$

```

1 :  $T_0 = R$ 
2 : for  $i = 1 : p$  do
3 :    $\begin{pmatrix} \boxed{T_i} \\ 0 \end{pmatrix} = \text{qr}\left(\begin{pmatrix} \boxed{T_{i-1}} \\ T_{i-1}\Lambda^{2^{i-1}} \end{pmatrix}\right)$ 
4 : end for

```

# Matlab code for $S = \mathbb{V}_{\ell,m}^T \odot R$ (Khatri-Rao product $\odot$ )

```
function T = QR_Khatri_Rao_VTR_2p( R, Lambda, p )
% QR_Khatri_Rao_VTR_2p computes the upper triangular factor
% in the QR factorization of the Khatri-Rao product
% S=Khatri_Rao(Vlm.',R), where R is an <ell x ell> upper
% triangular matrix, and Vlm is an <ell x m> Vandermonde
% matrix V, whose columns are V(:,i) = Lambda.^(i-1),
% i = 1,...,m, and m=2^p.
% Input:
% R      upper triangular matrix
% Lambda vector, defines Vlm = Vandermonde matrix
% p      integer >=0      defines m = 2^p
% Output:
% T      triangular QR factor of Khatri_Rao(Vlm.',R)
T = R ; D = Lambda ;
%
for i = 1 : p
[~, T] = qr( [ T ; T*diag(D)], 0 ) ;
D = D.^2 ;
end
end
```

**Input:** Upper triangular  $R \in \mathbb{C}^{\ell \times \ell}$ ; diagonal  $\Lambda \in \mathbb{C}^{\ell \times \ell}$ ;  $m$

**Output:** Upper triangular QR factor  $R_S = \mathbb{T}_{j-1}$  of  $S$  in (3)

```

1: Compute the binary representation of  $m$ :  $m \equiv \mathbf{b} =$ 
   ( $\mathbf{b}_{\lfloor \log_2 m \rfloor}, \dots, \mathbf{b}_1, \mathbf{b}_0$ )2,  $m \equiv \sum_{j=1}^{j^*} 2^{i_j}$ 
2: Let  $\lfloor \log_2 m \rfloor = i_{j^*} > i_{j^*-1} > \dots > i_2 > i_1 \geq 0$ 
3:  $T_0 = R$ 
4: if  $i_1 = 0$  then
5:    $\mathbb{T}_1 = T_0$ ;  $j = 2$ ;  $\wp = 1$ 
6: else
7:    $\mathbb{T}_0 = []$ ;  $j = 1$ ;  $\wp = 0$ 
8: end if
9: for  $k = 1 : i_{j^*}$  do
10:   $\begin{pmatrix} \boxed{T_k} \\ \mathbf{0} \end{pmatrix} = \text{qr} \left( \begin{pmatrix} \boxed{T_{k-1}} \\ T_{k-1} \Lambda^{2^{k-1}} \end{pmatrix} \right)$  {Local factor.}
11:  if  $k = i_j$  then
12:    if  $\mathbb{T}_{j-1} \neq []$  then
13:       $\begin{pmatrix} \boxed{\mathbb{T}_j} \\ \mathbf{0} \end{pmatrix} = \text{qr} \left( \begin{pmatrix} \boxed{\mathbb{T}_{j-1}} \\ T_k \Lambda^{\wp} \end{pmatrix} \right)$  {Global factor.}
14:    else
15:       $\boxed{\mathbb{T}_j} = \boxed{T_k}$ 
16:    end if
17:     $j := j + 1$ ;  $\wp := \wp + 2^k$ 
18:  end if
19: end for
```

# Comments and concluding remarks

- Provably small backward error

$$\|\delta S(:, j)\|_2 \leq \eta \|S(:, j)\|_2, \quad j = 1, \dots, \ell; \quad \eta \leq f(\ell, m)\epsilon,$$

- The relevant condition number is of the column scaled  $S$ :

## Corollary

$$\begin{aligned} \kappa_2(S_c) = \sqrt{\kappa_2(C_s)} &\leq \min(\kappa_2(R_c), \kappa_2((\mathbb{V}_{\ell, m})_r)) \\ &\leq \sqrt{\ell} \min\left(\min_{D=\text{diag}} \kappa_2(RD), \min_{D=\text{diag}} \kappa_2(D\mathbb{V}_{\ell, m})\right). \end{aligned}$$

- If the data is real, can work in real arithmetic even if the eigenvalues are complex (conjugate pairs)