



Computational methods for large-scale matrix equations and application to PDEs

V. Simoncini

Dipartimento di Matematica

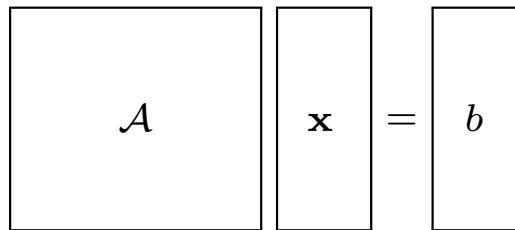
Alma Mater Studiorum - Università di Bologna

`valeria.simoncini@unibo.it`

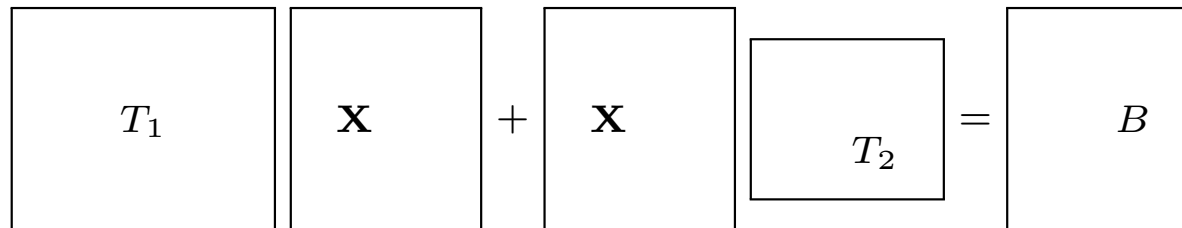
Linear (vector) systems and linear matrix equations

Problem: solve the linear problem

$$\mathcal{A}\mathbf{x} = b \quad \text{or} \quad T_1\mathbf{X} + \mathbf{X}T_2 = B$$



A diagram illustrating the linear system $\mathcal{A}\mathbf{x} = b$. It consists of three vertical rectangular boxes. The first box on the left is wider than it is tall and contains the symbol \mathcal{A} . The second box is narrower and taller, containing the symbol \mathbf{x} . The third box is also narrower and taller, containing the symbol b . An equals sign is placed between the second and third boxes.

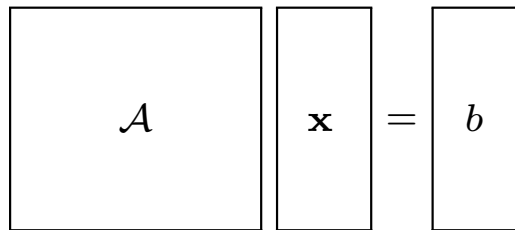


A diagram illustrating the matrix equation $T_1\mathbf{X} + \mathbf{X}T_2 = B$. It consists of five rectangular boxes arranged horizontally. The first box is tall and contains T_1 . The second box is tall and contains \mathbf{X} . A plus sign is between the second and third boxes. The third box is tall and contains \mathbf{X} . The fourth box is wide and contains T_2 . An equals sign is between the fourth and fifth boxes. The fifth box is tall and contains B .

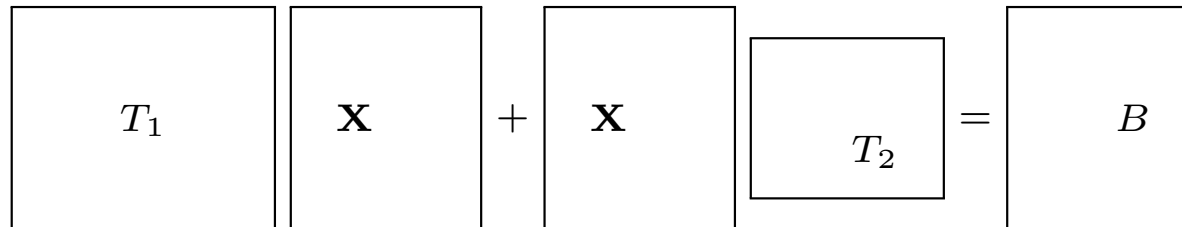
Linear (vector) systems and linear matrix equations

Problem: solve the linear problem

$$\mathcal{A}\mathbf{x} = b \quad \text{or} \quad T_1\mathbf{X} + \mathbf{X}T_2 = B$$



A diagram representing the linear system $\mathcal{A}\mathbf{x} = b$. It consists of three vertical rectangular boxes. The first box on the left is wider than the other two and contains the symbol \mathcal{A} . To its right is a narrower box containing the symbol \mathbf{x} . To the right of the \mathbf{x} box is an equals sign, followed by a third narrow box containing the symbol b .



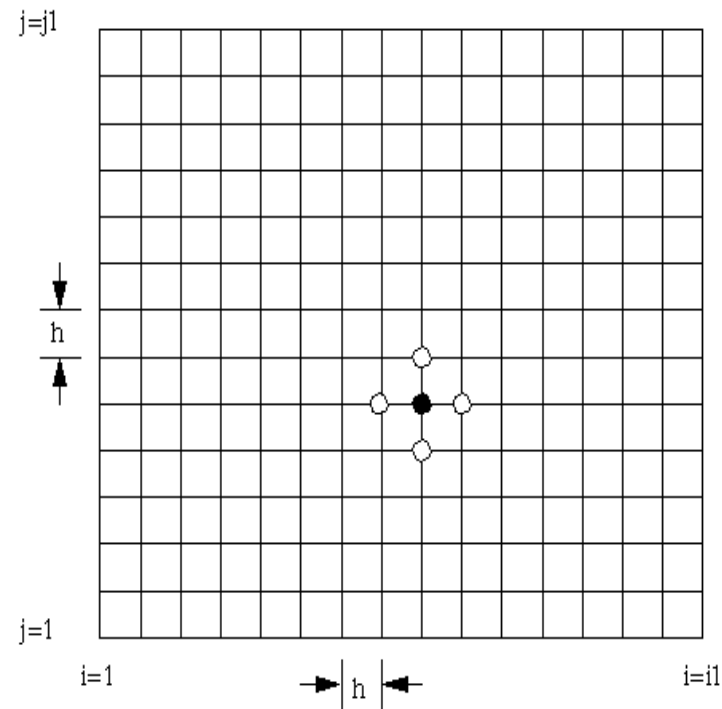
A diagram representing the matrix equation $T_1\mathbf{X} + \mathbf{X}T_2 = B$. It consists of five rectangular boxes arranged horizontally. The first box is tall and contains T_1 . To its right is a shorter box containing \mathbf{X} . To the right of the \mathbf{X} box is a plus sign. To the right of the plus sign is another shorter box containing \mathbf{X} . To the right of this second \mathbf{X} box is a shorter box containing T_2 . To the right of the T_2 box is an equals sign, followed by a tall box containing B .

Remark: In discretizing PDEs with tensor bases, the two problems may be mathematically equivalent !

The Poisson equation

$$-u_{xx} - u_{yy} = f, \quad \text{in } \Omega = (0, 1)^2$$

+ Dirichlet b.c. (zero b.c. for simplicity)



The Poisson equation

$$-u_{xx} - u_{yy} = f, \quad \text{in } \Omega = (0, 1)^2 \quad + \text{ Dirichlet zero b.c.}$$

FD Discretization: $U_{i,j} \approx u(x_i, y_j)$, with (x_i, y_j) interior nodes, so that

$$u_{xx}(x_i, y_j) \approx \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} = \frac{1}{h^2} [1, -2, 1] \begin{bmatrix} U_{i-1,j} \\ U_{i,j} \\ U_{i+1,j} \end{bmatrix}$$

$$u_{yy}(x_i, y_j) \approx \frac{U_{i,j-1} - 2U_{i,j} + U_{i,j+1}}{h^2} = \frac{1}{h^2} [U_{i,j-1}, U_{i,j}, U_{i,j+1}] \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

$$-T_1 \mathbf{U} - \mathbf{U} T_1^\top = F, \quad F_{ij} = f(x_i, y_j), \quad T_1 = \frac{1}{h^2} \text{tridiag}(1, -2, 1)$$

The Poisson equation

$$-u_{xx} - u_{yy} = f, \quad \text{in } \Omega = (0, 1)^2 \quad + \text{Dirichlet zero b.c.}$$

FD Discretization: $U_{i,j} \approx u(x_i, y_j)$, with (x_i, y_j) interior nodes, so that

$$u_{xx}(x_i, y_j) \approx \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} = \frac{1}{h^2} [1, -2, 1] \begin{bmatrix} U_{i-1,j} \\ U_{i,j} \\ U_{i+1,j} \end{bmatrix}$$

$$u_{yy}(x_i, y_j) \approx \frac{U_{i,j-1} - 2U_{i,j} + U_{i,j+1}}{h^2} = \frac{1}{h^2} [U_{i,j-1}, U_{i,j}, U_{i,j+1}] \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

$$-T_1 \mathbf{U} - \mathbf{U} T_1^\top = F, \quad F_{ij} = f(x_i, y_j), \quad T_1 = \frac{1}{h^2} \text{tridiag}(1, -2, 1)$$

Lexicographic ordering:

$$\mathbf{U} \rightarrow \text{vec}(\mathbf{U}) = \mathbf{u} = [\mathbf{U}_{11}, \dots, \mathbf{U}_{n,1}, \mathbf{U}_{1,2}, \dots, \mathbf{U}_{n,2}, \dots]^\top$$

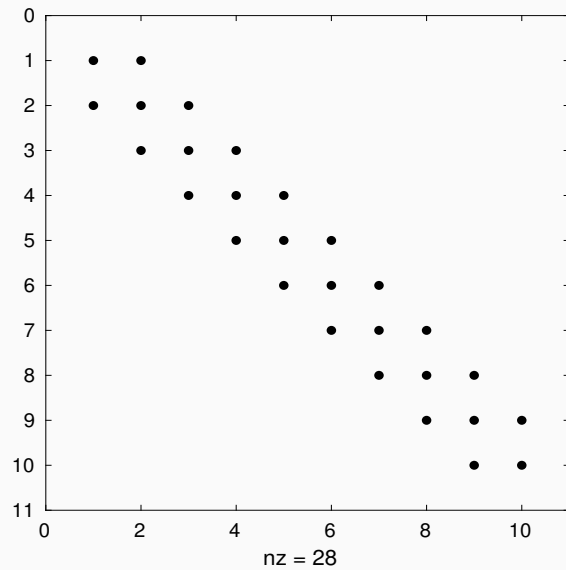
$$\boxed{\mathcal{A} \mathbf{u} = f} \quad \text{with } \mathcal{A} = -I \otimes T_1 - T_1 \otimes I, \quad f = \text{vec}(F)$$

$((M \otimes N)$ Kronecker product, $(M \otimes N) = (M_{i,j} N)$)

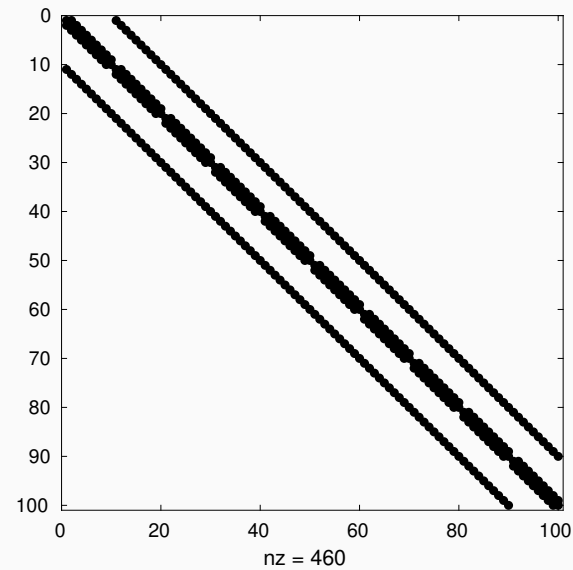
Computational considerations

$$T_1 \mathbf{U} + \mathbf{U} T_2 = F, \quad T_i \in \mathbb{R}^{n_i \times n_i}$$

$$\mathcal{A} \mathbf{u} = f \quad \mathcal{A} = I \otimes T_1 + T_2 \otimes I \in \mathbb{R}^{n_1 n_2 \times n_1 n_2}$$



T_1



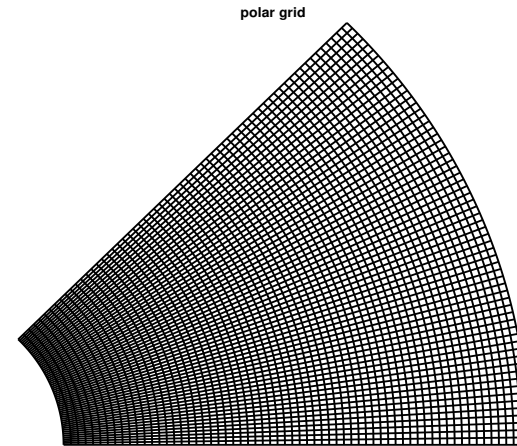
\mathcal{A}

Discretization of more complex domains (with Y. Hao)

$$-u_{xx} - u_{yy} = f, \quad \text{in } \Omega$$

$$(x, y) \in \Omega, \quad x = r \cos \theta, \quad y = r \sin \theta$$

$$(r, \theta) \in [r_0, r_1] \times [0, \frac{\pi}{4}]$$



♣ Transformed equation in polar coordinates:

$$-r^2 \tilde{u}_{rr} - r \tilde{u}_r - \tilde{u}_{\theta\theta} = \tilde{f}, \quad (r, \theta) \in [r_0, r_1] \times [0, \frac{\pi}{4}]$$

Matrix equation after mapping to the rectangle:

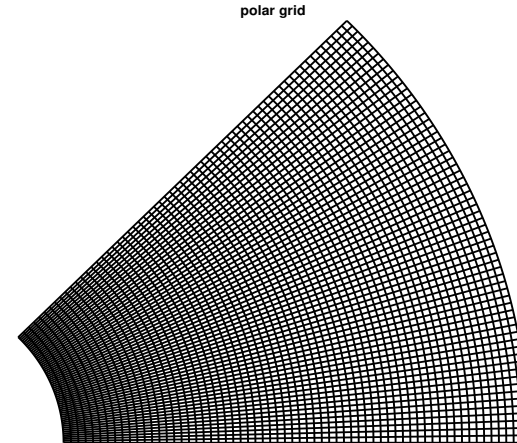
$$\boxed{\Phi^2 T \tilde{U} + \tilde{U} T - \Phi B \tilde{U} = \tilde{F}} \quad \Leftrightarrow \quad \boxed{(\Phi^2 T - \Phi B) \tilde{U} + \tilde{U} T = \tilde{F}}$$

Discretization of more complex domains (with Y. Hao)

$$-u_{xx} - u_{yy} = f, \quad \text{in } \Omega$$

$$(x, y) \in \Omega, \quad x = r \cos \theta, \quad y = r \sin \theta$$

$$(r, \theta) \in [r_0, r_1] \times [0, \frac{\pi}{4}]$$



♣ Transformed equation in log-polar coordinates ($r = e^\rho$):

$$-\hat{u}_{\rho\rho} - \hat{u}_{\theta\theta} = \hat{f}, \quad (\rho, \theta) \in [\rho_0, \rho_1] \times [0, \frac{\pi}{4}]$$

Matrix equation after mapping to the rectangle:

$$\boxed{T\hat{U} + \hat{U}T = \hat{F}}$$

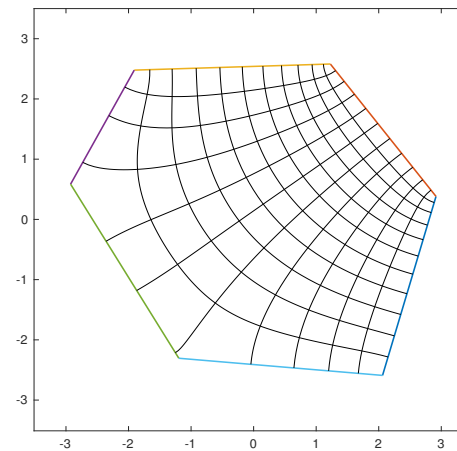
Poisson equation in a polygon with more than 4 edges (with Y. Hao)

♣ Schwarz-Christoffel conformal mappings between polygon Ω and rectangle Π

$$-u_{xx} - u_{yy} = f, \quad (x, y) \in \Omega$$

$$-\tilde{u}_{\xi\xi} - \tilde{u}_{\eta\eta} = \mathcal{J} \tilde{f}, \quad (\xi, \eta) \in \Pi$$

(\mathcal{J} Jacobian det of SC mapping)



With finite difference discretization in Π :

$$\boxed{T_1 U + U T_2 = F}, \quad F = \tilde{F} + b.c., \quad \text{and} \quad \tilde{F}_{i,j} = (\mathcal{J} \tilde{f})(\xi_i, \eta_j), \quad 1 \leq i \leq n_1, \quad 1 \leq j \leq n_2$$

Poisson equation is the ideal setting for SC mappings!

More general settings

- Convection-diffusion eqns in a rectangle
(see, e.g., Palitta & Simoncini, 2016)
- Space-Time discretizations via tensorized high order methods
(see, e.g., joint wrk w/ Henning, Palitta and Urban, 2020)
- Galerkin FE discretization of Stochastic PDEs
(see, e.g., joint wrk w/ Powell and Silvester, 2017)
- Isogeometric Analysis
(see, e.g., Sangalli and Tani, 2016)
- ...

... A classical approach, Bickley & McNamee, 1960, Wachspress, 1963
(Early literature on difference equations)

Numerical solution of the Sylvester equation

$$AU + UB = G$$

Various settings:

- Small A and small B : Bartels-Stewart algorithm
 1. Compute the Schur forms:
 $A^* = URU^*$, $B = VSV^*$ with R, S upper triangular;
 2. Solve $R^*Y + YS = U^*GV$ for Y (element-wise);
 3. Compute $U = UYV^*$

Numerical solution of the Sylvester equation

$$AU + UB = G$$

Various settings:

- Small A and small B : Bartels-Stewart algorithm
 1. Compute the Schur forms:
 $A^* = URU^*$, $B = VSV^*$ with R, S upper triangular;
 2. Solve $R^*Y + YS = U^*GV$ for Y (element-wise);
 3. Compute $U = UYV^*$
- Large A and small B : Column decoupling
 1. Compute the decomposition $B = WSW^{-1}$, $S = \text{diag}(s_1, \dots, s_m)$
 2. Set $\hat{G} = GW$
 3. For $i = 1, \dots, m$ solve $(A + s_i I)(\hat{U})_i = (\hat{G})_i$
 4. Compute $U = \hat{U}W^{-1}$

Numerical solution of the Sylvester equation

$$AU + UB = G$$

Various settings:

- Small A and small B : Bartels-Stewart algorithm
 1. Compute the Schur forms:
 $A^* = URU^*$, $B = VSV^*$ with R, S upper triangular;
 2. Solve $R^*Y + YS = U^*GV$ for Y (element-wise);
 3. Compute $U = UYV^*$.
- Large A and small B : Column decoupling
 1. Compute the decomposition $B = WSW^{-1}$, $S = \text{diag}(s_1, \dots, s_m)$
 2. Set $\hat{G} = GW$
 3. For $i = 1, \dots, m$ solve $(A + s_i I)(\hat{U})_i = (\hat{G})_i$
 4. Compute $U = \hat{U}W^{-1}$
- Large A and large B : Iterative solution (G low rank, or G sparse)

Numerical solution of large scale Sylvester equations

$$AU + UB = G$$

with G low rank

- Projection methods
- ADI (Alternating Direction Iteration)
- Data sparse approaches (structure-dependent)

Projection methods

Seek $U_k \approx U$ of low rank:

$$U_k = \begin{bmatrix} U_k^{(1)} \\ \end{bmatrix} [(U_k^{(2)})^\top]$$

with $U_k^{(1)}, U_k^{(2)}$ tall

Index k “related” to the approximation rank

See, Simoncini, SIREV 2016.

Multiterm linear matrix equation

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

$A_i \in \mathbb{R}^{n \times n}$, $B_i \in \mathbb{R}^{m \times m}$, \mathbf{X} unknown matrix

Possibly large dimensions, structured coefficient matrices

The problem in its full generality is far from tractable, although the transformation to a matrix-vector equation [...] allows us to use the considerable arsenal of numerical weapons currently available for the solution of such problems.

Peter Lancaster, SIAM Rev. 1970

Multiterm linear matrix equation. Classical device

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

Kronecker formulation $(B_1^\top \otimes A_1 + \dots + B_\ell^\top \otimes A_\ell) \mathbf{x} = c \Leftrightarrow \mathcal{A} \mathbf{x} = c$

Iterative methods: matrix-matrix multiplications and rank truncation

(Benner, Breiten, Bouhamidi, Chehab, Damm, Grasedyck, Jbilou, Kressner, Matthies, Nagy, Onwunta, Raydan, Stoll, Tobler, Wedderburn, Zander, ...)

Kronecker product : $M \otimes P = \begin{bmatrix} m_{11}P & \dots & m_{1n}P \\ \vdots & \ddots & \vdots \\ m_{n1}P & \dots & m_{nn}P \end{bmatrix}$ and $\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X)$

Multiterm linear matrix equation. Classical device

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

Kronecker formulation $(B_1^\top \otimes A_1 + \dots + B_\ell^\top \otimes A_\ell) \mathbf{x} = c \Leftrightarrow \mathcal{A} \mathbf{x} = c$

Iterative methods: matrix-matrix multiplications and rank truncation

(Benner, Breiten, Bouhamidi, Chehab, Damm, Grasedyck, Jbilou, Kressner, Matthies, Nagy, Onwunta, Raydan, Stoll, Tobler, Wedderburn, Zander, ...)

Kronecker product : $M \otimes P = \begin{bmatrix} m_{11}P & \dots & m_{1n}P \\ \vdots & \ddots & \vdots \\ m_{n1}P & \dots & m_{nn}P \end{bmatrix}$ and $\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X)$

Alternatives to Kronecker form:

- Fixed point iterations (an “evergreen” ...)
- Projection-type methods \Rightarrow low rank approximation
- Ad-hoc problem-dependent procedures
- etc.

Current very active area of research

Truncated matrix-oriented CG for Kronecker form

Input: $\mathcal{A}(X) = A_1 X B_1 + A_2 X B_2 + \dots + A_\ell X B_\ell$, right-hand side $C \in \mathbb{R}^{n \times n}$ in low-rank format. Truncation operator \mathcal{T} .

Output: Matrix $X \in \mathbb{R}^{n \times n}$ in low-rank format s.t. $\|\mathcal{A}(X) - C\|_F / \|C\|_F \leq \text{tol}$.

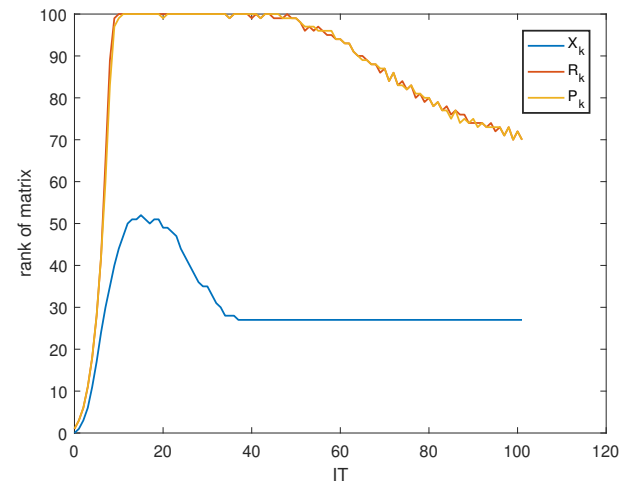
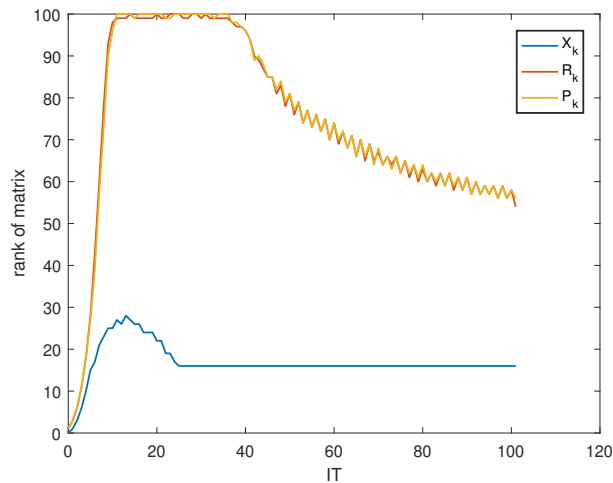
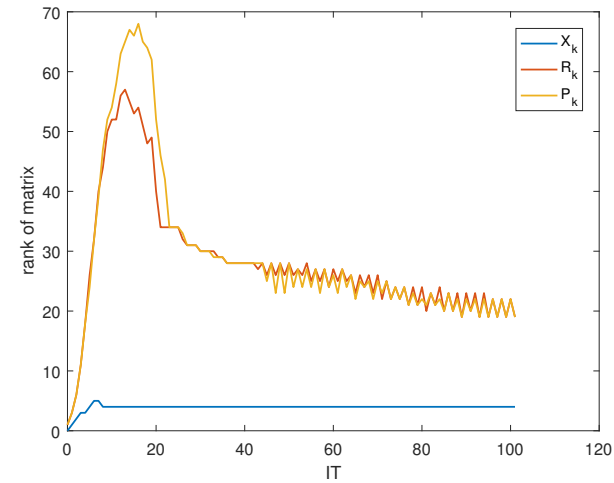
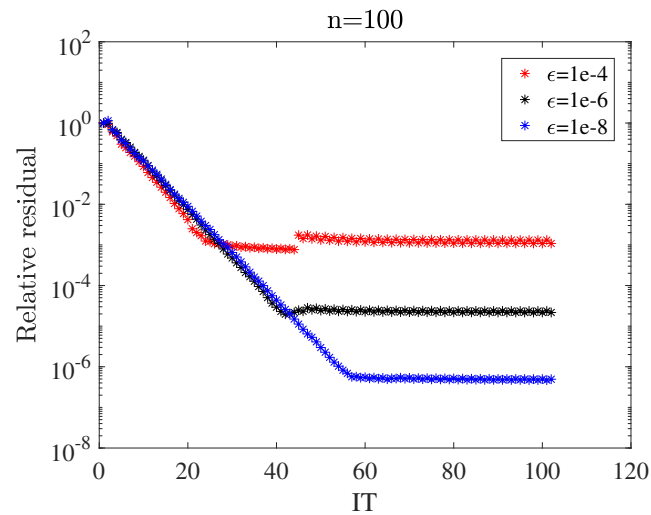
- 1: $X_0 = 0, R_0 = C, P_0 = R_0, Q_0 = \mathcal{A}(P_0)$
- 2: $\xi_0 = \langle P_0, Q_0 \rangle, k = 0$ $\langle X, Y \rangle = \text{tr}(X^\top Y)$
- 3: **while** $\|R_k\|_F > \text{tol}$ **do**
- 4: $\omega_k = \langle R_k, P_k \rangle / \xi_k$
- 5: $X_{k+1} = X_k + \omega_k P_k,$ $X_{k+1} \leftarrow \mathcal{T}(X_{k+1})$
- 6: $R_{k+1} = C - \mathcal{A}(X_{k+1}),$ Optionally: $R_{k+1} \leftarrow \mathcal{T}(R_{k+1})$
- 7: $\beta_k = -\langle R_{k+1}, Q_k \rangle / \xi_k$
- 8: $P_{k+1} = R_{k+1} + \beta_k P_k,$ $P_{k+1} \leftarrow \mathcal{T}(P_{k+1})$
- 9: $Q_{k+1} = \mathcal{A}(P_{k+1}),$ Optionally: $Q_{k+1} \leftarrow \mathcal{T}(Q_{k+1})$
- 10: $\xi_{k+1} = \langle P_{k+1}, Q_{k+1} \rangle$
- 11: $k = k + 1$
- 12: **end while**
- 13: $X = X_k$

♣ Iterates kept in factored form!

Kressner and Tobler, 2011

Threshold based truncated CG. $n = 100$, $\text{tol}=\epsilon \in \{10^{-4}, 10^{-6}, 10^{-8}\}$

$A = \frac{1}{h^2} \text{tridiag}(-1, \underline{2}, -1)$, $M = \text{diag}(a_1)$, $N = \text{diag}(a_2)$, a_1 and a_2 random vectors



Projection-type methods. 1

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

Given approximation spaces $\mathcal{K}_A, \mathcal{K}_B,$

$$\mathbf{X} \approx X_m \quad \text{with} \quad \text{vec}(X_m) \in \mathcal{K}_B \otimes \mathcal{K}_A$$

Projection-type methods. 1

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

Given approximation spaces $\mathcal{K}_A, \mathcal{K}_B$,

$$\mathbf{X} \approx X_m \quad \text{with} \quad \text{vec}(X_m) \in \mathcal{K}_B \otimes \mathcal{K}_A$$

\mathbf{X} is approximated by a low rank matrix !

that is, $X_m := V_m Y_m W_m^\top$, $\mathcal{K}_A = \text{Range}(V_m)$, $\mathcal{K}_B = \text{Range}(W_m)$

Galerkin condition:

$$R := A_1 X_m B_1 + A_2 X_m B_2 + \dots + A_\ell X_m B_\ell - C \quad \perp \quad \mathcal{K}_B \otimes \mathcal{K}_A$$

$$V_m^\top R W_m = 0$$

Projection-type methods. 1

$$A_1 \mathbf{X} B_1 + A_2 \mathbf{X} B_2 + \dots + A_\ell \mathbf{X} B_\ell = C$$

Given approximation spaces $\mathcal{K}_A, \mathcal{K}_B$,

$$\mathbf{X} \approx X_m \quad \text{with} \quad \text{vec}(X_m) \in \mathcal{K}_B \otimes \mathcal{K}_A$$

\mathbf{X} is approximated by a low rank matrix !

that is, $X_m := V_m Y_m W_m^\top$, $\mathcal{K}_A = \text{Range}(V_m)$, $\mathcal{K}_B = \text{Range}(W_m)$

Galerkin condition:

$$R := A_1 X_m B_1 + A_2 X_m B_2 + \dots + A_\ell X_m B_\ell - C \quad \perp \quad \mathcal{K}_B \otimes \mathcal{K}_A$$

$$V_m^\top R W_m = 0$$

Projected matrix equation:

$$V_m^\top (A_1 X_m B_1 + \dots + A_\ell X_m B_\ell - C) W_m = 0$$

$$(V_m^\top A_1 V_m) \mathbf{Y} (W_m^\top B_1 W_m) + \dots + (V_m^\top A_\ell V_m) \mathbf{Y} (W_m^\top B_\ell W_m) - V_m^\top C W_m = 0$$

Projection-type methods. 2

Solve for \mathbf{Y} :

$$(V_m^\top A_1 V_m) \mathbf{Y} (W_m^\top B_1 W_m) + \dots + (V_m^\top A_\ell V_m) \mathbf{Y} (W_m^\top B_\ell W_m) - V_m^\top C W_m = 0$$

Then, implicitly generate $X_m := V_m \mathbf{Y}_m W_m^\top$

Procedure generalizes the case $\ell = 2$, using the classical *Galerkin projection* methodology

Projection-type methods. 2

Solve for \mathbf{Y} :

$$(V_m^\top A_1 V_m) \mathbf{Y} (W_m^\top B_1 W_m) + \dots + (V_m^\top A_\ell V_m) \mathbf{Y} (W_m^\top B_\ell W_m) - V_m^\top C W_m = 0$$

Then, implicitly generate $X_m := V_m \mathbf{Y}_m W_m^\top$

Procedure generalizes the case $\ell = 2$, using the classical *Galerkin projection* methodology

Optimality property:

Palitta and Simoncini, 2020

$$\|X_\star - X_m\|_{\mathcal{A}} = \min_{\substack{Z = V_m Y W_m^\top \\ Y \in \mathbb{R}^{m \times m}}} \|X_\star - Z\|_{\mathcal{A}},$$

where $\|X\|_{\mathcal{A}}^2 = \text{trace} \left(\sum_{j=1}^{\ell} X^\top A_j X B_j \right)$.

Projection-type methods. 2

Solve for \mathbf{Y} :

$$(V_m^\top A_1 V_m) \mathbf{Y} (W_m^\top B_1 W_m) + \dots + (V_m^\top A_\ell V_m) \mathbf{Y} (W_m^\top B_\ell W_m) - V_m^\top C W_m = 0$$

Then, implicitly generate $X_m := V_m \mathbf{Y}_m W_m^\top$

Procedure generalizes the case $\ell = 2$, using the classical *Galerkin projection* methodology

Optimality property:

Palitta and Simoncini, 2020

$$\|X_\star - X_m\|_{\mathcal{A}} = \min_{\substack{Z = V_m Y W_m^\top \\ Y \in \mathbb{R}^{m \times m}}} \|X_\star - Z\|_{\mathcal{A}},$$

where $\|X\|_{\mathcal{A}}^2 = \text{trace} \left(\sum_{j=1}^{\ell} X^\top A_j X B_j \right)$.

Crucial issues for effectiveness:

- Choice of spaces $\mathcal{K}_A, \mathcal{K}_B$ and their construction. Ideally,

$$\text{range}(V_m) \subseteq \text{range}(V_{m+1}), \quad \text{range}(W_m) \subseteq \text{range}(W_{m+1})$$

- Solution of the reduced multiterm equation

A “simple” example

$$A\mathbf{X} + \mathbf{X}A + M\mathbf{X}M = ff^\top, \quad A, M \text{ spd, } f \text{ vector}$$

♣ No available *direct* methods for the generic case, except Kronecker form

A “simple” example

$$A\mathbf{X} + \mathbf{X}A + M\mathbf{X}M = ff^\top, \quad A, M \text{ spd, } f \text{ vector}$$

♣ No available *direct* methods for the generic case, except Kronecker form

Matrix-oriented CG: $X^{(k)} = X_1^{(k)} G^{(k)} (X_1^{(k)})^\top$

$\text{range}(X_1^{(k)}) \subset \mathbb{Q}_k = \text{span}\{f, Af, Mf, A^2f, AMf, MAf, M^2f, \dots\}$, $\dim(\mathbb{Q}_{k+1}) \leq \dim(\mathbb{Q}_k) + 2^k$

A “simple” example

$$A\mathbf{X} + \mathbf{X}A + M\mathbf{X}M = ff^\top, \quad A, M \text{ spd, } f \text{ vector}$$

♣ No available *direct* methods for the generic case, except Kronecker form

Matrix-oriented CG: $X^{(k)} = X_1^{(k)} G^{(k)} (X_1^{(k)})^\top$

$\text{range}(X_1^{(k)}) \subset \mathbb{Q}_k = \text{span}\{f, Af, Mf, A^2f, AMf, MAf, M^2f, \dots\}$, $\dim(\mathbb{Q}_{k+1}) \leq \dim(\mathbb{Q}_k) + 2^k$

Galerkin method: Choose $\mathcal{K}_m = \text{range}(V_m)$ with

$$\begin{aligned} V_0 = f &=: \underline{v_1} & V_1 &= [v_1, Av_1, Mv_1] =: [v_1, \underline{v_2}, v_3] \\ & & V_2 &= [V_1, Av_2, Mv_2] =: [v_1, v_2, \underline{v_3}, v_4, v_5] \\ & & V_3 &= [V_2, Av_3, Mv_3] =: [v_1, v_2, v_3, \underline{v_4}, v_5, v_6, v_7] \\ & & & \text{etc.} \end{aligned}$$

$\Rightarrow \mathcal{K}_m = \text{span}\{f, Af, Mf, A^2f, AMf, MAf, M^2f, \dots\}$, $\dim(\mathcal{K}_m) = 2m + 1$

$\mathbb{Q}_k = \text{range}(V_{2^k-1})$

Hao and Simoncini, work in progress

Computational methods for certain structured problems

A particular case^a:

$$A\mathbf{X} + \mathbf{X}A^\top + M_1\mathbf{X}M_1 + \dots + M_\ell\mathbf{X}M_\ell = F,$$

with $A \in \mathbb{R}^{n \times n}$, M_i s with very low rank s_i , $M_i = U_i V_i^\top$

^aIn fact, terms in the form $M_i\mathbf{X}N_i$ can also be treated

Computational methods for certain structured problems

A particular case^a:

$$A\mathbf{X} + \mathbf{X}A^\top + M_1\mathbf{X}M_1 + \dots + M_\ell\mathbf{X}M_\ell = F,$$

with $A \in \mathbb{R}^{n \times n}$, M_i s with very low rank s_i , $M_i = U_i V_i^\top$

Using the Kronecker form ($\ell = 1$):

$$(A \otimes I + I \otimes A + (U_1 \otimes U_1)(V_1 \otimes V_1)^\top)\mathbf{x} = f$$

that is

$$(\mathcal{A} + \mathcal{U}\mathcal{V}^\top)\mathbf{x} = f$$

with $\mathcal{U} = U_1 \otimes U_1$, $\mathcal{V} = V_1 \otimes V_1$ again of low rank s_1^2

^aIn fact, terms in the form $M_i\mathbf{X}N_i$ can also be treated

Computational methods for certain structured problems

A particular case^a:

$$A\mathbf{X} + \mathbf{X}A^\top + M_1\mathbf{X}M_1 + \dots + M_\ell\mathbf{X}M_\ell = F,$$

with $A \in \mathbb{R}^{n \times n}$, M_i s with very low rank s_i , $M_i = U_i V_i^\top$

Using the Kronecker form ($\ell = 1$):

$$(A \otimes I + I \otimes A + (U_1 \otimes U_1)(V_1 \otimes V_1)^\top)\mathbf{x} = f$$

that is

$$(\mathcal{A} + \mathcal{U}\mathcal{V}^\top)\mathbf{x} = f$$

with $\mathcal{U} = U_1 \otimes U_1$, $\mathcal{V} = V_1 \otimes V_1$ again of low rank s_1^2

Solution method: Sherman-Morrison-Woodbury formula

$$\mathbf{x} = (\mathcal{A} + \mathcal{U}\mathcal{V}^\top)^{-1}f = \mathcal{A}^{-1}f - \mathcal{A}^{-1}\mathcal{U}(I + \mathcal{V}^\top\mathcal{A}^{-1}\mathcal{U})^{-1}\mathcal{V}^\top\mathcal{A}^{-1}f$$

^aIn fact, terms in the form $M_i\mathbf{X}N_i$ can also be treated

Matrix-oriented Sherman-Morrison-Woodbury formula

$$\mathbf{x} = \mathcal{A}^{-1}f - \mathcal{A}^{-1}\mathcal{U}(I + \mathcal{V}^\top \mathcal{A}^{-1}\mathcal{U})^{-1}\mathcal{V}^\top \mathcal{A}^{-1}f$$

1. Solve $\mathcal{A}w = f$
2. Solve $\mathcal{A}p_j = \mathbf{u}_j$ where $\mathcal{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{s^2}]$ to give $\mathcal{P} = [p_1, \dots, p_{s^2}]$;
3. Compute $H = I + \mathcal{V}^\top \mathcal{P} \in \mathbb{R}^{s^2 \times s^2}$
4. Solve $Hg = \mathcal{V}^\top w$
5. Compute $x = w - \mathcal{P}g$.

Matrix-oriented Sherman-Morrison-Woodbury formula

$$\mathbf{x} = \mathcal{A}^{-1}f - \mathcal{A}^{-1}\mathcal{U}(I + \mathcal{V}^\top \mathcal{A}^{-1}\mathcal{U})^{-1}\mathcal{V}^\top \mathcal{A}^{-1}f$$

1. Solve $\mathcal{A}w = f$
2. Solve $\mathcal{A}p_j = u_j$ where $\mathcal{U} = [u_1, \dots, u_{s^2}]$ to give $\mathcal{P} = [p_1, \dots, p_{s^2}]$;
3. Compute $H = I + \mathcal{V}^\top \mathcal{P} \in \mathbb{R}^{s^2 \times s^2}$
4. Solve $Hg = \mathcal{V}^\top w$
5. Compute $x = w - \mathcal{P}g$.

Steps 1. and 2.:

$$w = \mathcal{A}^{-1}f \quad \Leftrightarrow \quad AW + WA^\top = F, \quad f = \text{vec}(F)$$

Analogously for each $p_j = \text{vec}(P_j)$ in step 2

$AW + WA^\top = P_j$ Lyapunov equations, with the same A - cheap “direct” solution

Matrix-oriented Sherman-Morrison-Woodbury formula

$$\mathbf{x} = \mathcal{A}^{-1}f - \mathcal{A}^{-1}\mathcal{U}(I + \mathcal{V}^\top \mathcal{A}^{-1}\mathcal{U})^{-1}\mathcal{V}^\top \mathcal{A}^{-1}f$$

1. Solve $\mathcal{A}w = f$
2. Solve $\mathcal{A}p_j = \mathbf{u}_j$ where $\mathcal{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{s^2}]$ to give $\mathcal{P} = [p_1, \dots, p_{s^2}]$;
3. Compute $H = I + \mathcal{V}^\top \mathcal{P} \in \mathbb{R}^{s^2 \times s^2}$
4. Solve $Hg = \mathcal{V}^\top w$
5. Compute $x = w - \mathcal{P}g$.

Steps 1. and 2.:

$$w = \mathcal{A}^{-1}f \quad \Leftrightarrow \quad AW + WA^\top = F, \quad f = \text{vec}(F)$$

Analogously for each $p_j = \text{vec}(P_j)$ in step 2

$AW + WA^\top = P_j$ Lyapunov equations, with the same A - cheap “direct” solution

Step 3.

$$\mathbf{v}_j^\top \mathcal{A}^{-1}\mathbf{u}_t = v_i^\top P_t v_k, \quad j = (k-1)s + i$$

Analogously for $\mathcal{V}^\top w$ in step 4

A numerical example

Let X_* be a ref. soln (uniformly distr.random), and rhs computed explicitly

We monitor:

$$Err := \frac{\|X - X_*\|_F}{\|X_*\|_F}$$

| n | s_1/s_2 | Matrix form | | Vector Form | |
|-----|-----------|-------------|----------|-------------|----------|
| | | CPU time | Err | CPU time | Err |
| 40 | 3/5 | 0.013 | 3.81e-11 | 0.195 | 2.29e-10 |
| | 6/10 | 0.017 | 9.05e-10 | 0.657 | 4.98e-10 |
| | 12/20 | 0.035 | 5.25e-09 | 2.333 | 1.35e-08 |
| 80 | 3/5 | 0.022 | 2.15e-10 | 5.283 | 1.22e-09 |
| | 6/10 | 0.033 | 8.38e-09 | 15.408 | 1.84e-08 |
| | 12/20 | 0.074 | 2.50e-08 | 56.347 | 3.46e-08 |
| 160 | 3/5 | 0.043 | 1.29e-09 | 129.957 | 6.89e-09 |
| | 6/10 | 0.070 | 1.10e-08 | 281.946 | 2.69e-08 |
| | 12/20 | 0.220 | 2.90e-07 | 1030.242 | 1.20e-06 |

Table 1: Symmetric and dense matrix A and U_1, U_2 ($\ell = 2$) for various ranks s_1, s_2

Hao and Simoncini, 2021. See also Damm, 2008, Massei et al 2018.

Conclusions

- Rich setting for new algorithmic strategies
- Certain approaches appropriate for solving linear *tensor* equations
- Devise more general “direct” solvers, to be used (also) in the projection phase!

Visit: www.dm.unibo.it/~simoncin

Email address: valeria.simoncini@unibo.it

REFERENCES

1. Yue Hao and V. S., *Matrix equation solving of PDEs in polygonal domains using conformal mappings*. Journal of Numerical Mathematics, vol. 29, no. 3, 2021
2. Yue Hao and V. S., *The Sherman-Morrison-Woodbury formula for generalized linear matrix equations and applications*, Numer. Linear Algebra w/Apl. 28 (5), 2021
3. Catherine E. Powell, David Silvester and V. S., *An efficient reduced basis solver for stochastic Galerkin matrix equations*, SIAM J. Scientific Computing, 39 (1), (2017).
4. Davide Palitta and V. S., *Matrix-equation-based strategies for convection-diffusion equations*, BIT Numerical Mathematics, 56-2, (2016).
5. V.S., *Computational methods for linear matrix equations*,SIAM Review, 58-3(2016)