



MAX PLANCK INSTITUTE
FOR DYNAMICS OF COMPLEX
TECHNICAL SYSTEMS
MAGDEBURG



COMPUTATIONAL METHODS IN
SYSTEMS AND CONTROL THEORY

A new low-rank solver for algebraic Riccati equations based on the matrix sign function and principal pivot transforms

Peter Benner
joint work with Federico Poloni (Università di Pisa)

ILAS 2023
25th Conference of the International Linear Algebra Society
MSI04 “Matrix equations”
Madrid, June 12–16, 2023

1. Introduction
2. The Sign Function Method for Algebraic Riccati Equations
3. Another Motivating Application
Closed-loop Balanced Truncation
4. Principal Pivot Transforms
Backward Stability of PPT-based Inversion of SQSD Matrices
Structure-preserving Inversion of SQSD Matrices
5. Factored Form of Sign Function Iteration
6. Conclusions

Algebraic Riccati equation (ARE)

For $A, G = G^T, W = W^T \in \mathbb{R}^{n \times n}$ given and $X \in \mathbb{R}^{n \times n}$ unknown:

$$0 = \mathcal{R}(X) := A^T X + X A - X G X + W.$$

Algebraic Riccati equation (ARE)

For $A, G = G^T, W = W^T \in \mathbb{R}^{n \times n}$ given and $X \in \mathbb{R}^{n \times n}$ unknown:

$$0 = \mathcal{R}(X) := A^T X + X A - X G X + W.$$

Many applications:

- model reduction of (unstable) linear time-invariant (LTI) systems,
- linear-quadratic optimal control problems for LTI systems,
- H_∞ -control, ...

Algebraic Riccati equation (ARE)

For $A, G = G^T, W = W^T \in \mathbb{R}^{n \times n}$ given and $X \in \mathbb{R}^{n \times n}$ unknown:

$$0 = \mathcal{R}(X) := A^T X + X A - X G X + W.$$

Many applications:

- model reduction of (unstable) linear time-invariant (LTI) systems,
- linear-quadratic optimal control problems for LTI systems,
- H_∞ -control, ...

Typical situation in model reduction and control:

- G, W low-rank with $G, W \in \{BB^T, C^T C\}$, where $B \in \mathbb{R}^{n \times m}$, $m \ll n$, and $C \in \mathbb{R}^{p \times n}$, $p \ll n$.
- **Want:** solution with $X = X^T \geq 0$ (and $\Lambda(A - GX) \subset \mathbb{C}^-$), notation: X_{\geq} .

Algebraic Riccati equation (ARE)

For $A, G = G^T, W = W^T \in \mathbb{R}^{n \times n}$ given and $X \in \mathbb{R}^{n \times n}$ unknown:

$$0 = \mathcal{R}(X) := A^T X + X A - X G X + W.$$

Many applications:

- model reduction of (unstable) linear time-invariant (LTI) systems,
- linear-quadratic optimal control problems for LTI systems,
- H_∞ -control, ...

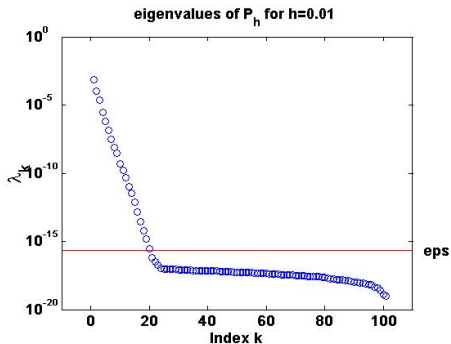
Typical situation in model reduction and control:

- G, W low-rank with $G, W \in \{BB^T, C^T C\}$, where $B \in \mathbb{R}^{n \times m}$, $m \ll n$, and $C \in \mathbb{R}^{p \times n}$, $p \ll n$.
- **Want:** solution with $X = X^T \geq 0$ (and $\Lambda(A - GX) \subset \mathbb{C}^-$), notation: X_{\geq} .
- $n = 10^3 - 10^6$
 $\implies X$ has $10^6 - 10^{12}$ unknowns
 \implies as X is dense in general, we face a storage problem!

Consider spectrum of ARE solution.

Example:

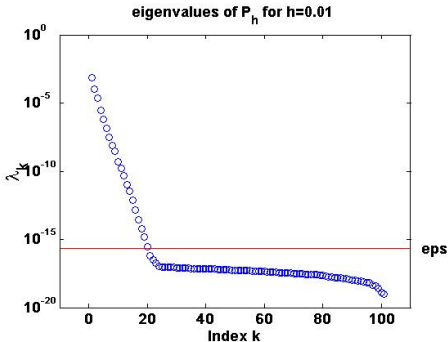
- Linear 1D heat equation with point control,
- $\Omega = [0, 1]$,
- FEM discretization using linear B-splines,
- $h = 1/100 \implies n = 101$.



Consider spectrum of ARE solution.

Example:

- Linear 1D heat equation with point control,
- $\Omega = [0, 1]$,
- FEM discretization using linear B-splines,
- $h = 1/100 \implies n = 101$.



Idea: $X = X^T \geq 0 \implies$

$$X = ZZ^T = \sum_{k=1}^n \lambda_k z_k z_k^T \approx \sum_{k=1}^r \lambda_k z_k z_k^T = \sum_{k=1}^r \left(\sqrt{\lambda_k} z_k \right) \left(\sqrt{\lambda_k} z_k \right)^T =: Z^{(r)} (Z^{(r)})^T.$$

\implies **Goal:** compute $Z^{(r)} \in \mathbb{R}^{n \times r}$ directly w/o ever forming X !



∃ many numerical methods to solve AREs.

Here: revisit the **matrix sign function method**.

∃ many numerical methods to solve AREs.

Here: revisit the **matrix sign function method**.

Definition

For $Z \in \mathbb{R}^{n \times n}$ with $\Lambda(Z) \cap i\mathbb{R} = \emptyset$ and Jordan canonical form

$$Z = S \begin{bmatrix} J^+ & 0 \\ 0 & J^- \end{bmatrix} S^{-1}$$

the **matrix sign function** is

$$\text{sign}(Z) := S \begin{bmatrix} I_k & 0 \\ 0 & -I_{n-k} \end{bmatrix} S^{-1}.$$

∃ many numerical methods to solve AREs.

Here: revisit the **matrix sign function method**.

Definition

For $Z \in \mathbb{R}^{n \times n}$ with $\Lambda(Z) \cap i\mathbb{R} = \emptyset$ and Jordan canonical form

$$Z = S \begin{bmatrix} J^+ & 0 \\ 0 & J^- \end{bmatrix} S^{-1}$$

the **matrix sign function** is

$$\text{sign}(Z) := S \begin{bmatrix} I_k & 0 \\ 0 & -I_{n-k} \end{bmatrix} S^{-1}.$$

Lemma

Let $T \in \mathbb{R}^{n \times n}$ be nonsingular and Z as before, then

$$\text{sign}(T Z T^{-1}) = T \text{sign}(Z) T^{-1}.$$

∃ many numerical methods to solve AREs.

Here: revisit the **matrix sign function method**.

Computation of $\text{sign}(Z)$

$\text{sign}(Z)$ is root of $I_n \implies$ use Newton's method to compute it:

$$Z_0 \leftarrow Z, \quad Z_{j+1} \leftarrow \frac{1}{2} \left(c_j Z_j + \frac{1}{c_j} Z_j^{-1} \right), \quad j = 1, 2, \dots$$

$$\implies \text{sign}(Z) = \lim_{j \rightarrow \infty} Z_j.$$

$c_j > 0$ is scaling parameter for convergence acceleration and rounding error minimization, e.g.

$$c_j = \sqrt{\frac{\|Z_j^{-1}\|_F}{\|Z_j\|_F}},$$

based on “equilibrating” the norms of the two summands [HIGHAM 1986].



Solving AREs with the Matrix Sign Function Method

$$\text{ARE: } 0 = A^T X + XA - XGX + W$$

Key observations:

Key observations:

- Let $H = \begin{bmatrix} A & G \\ W & -A^T \end{bmatrix}$ be the **Hamiltonian matrix** associated to the ARE and X_{\geq} the desired symmetric positive semidefinite solution. Then

$$H \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} = \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} (A - GX_{\geq}),$$

i.e., X_{\geq} defines an n -dimensional invariant subspace of H corresponding to its left-half-plane eigenvalues.

Key observations:

- Let $H = \begin{bmatrix} A & G \\ W & -A^T \end{bmatrix}$ be the **Hamiltonian matrix** associated to the ARE and X_{\geq} the desired symmetric positive semidefinite solution. Then

$$H \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} = \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} (A - G X_{\geq}),$$

i.e., X_{\geq} defines an n -dimensional invariant subspace of H corresponding to its left-half-plane eigenvalues.

- $\frac{1}{2} (I_{2n} + \text{sign}(H))$ is a projector onto the H -invariant subspace corresponding to its left-half-plane eigenvalues

$$\implies (I_{2n} + \text{sign}(H)) \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} = 0.$$

Key observations:

- Let $H = \begin{bmatrix} A & G \\ W & -A^T \end{bmatrix}$ be the **Hamiltonian matrix** associated to the ARE and X_{\geq} the desired symmetric positive semidefinite solution. Then

$$H \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} = \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} (A - GX_{\geq}),$$

i.e., X_{\geq} defines an n -dimensional invariant subspace of H corresponding to its left-half-plane eigenvalues.

- $\frac{1}{2} (I_{2n} + \text{sign}(H))$ is a projector onto the H -invariant subspace corresponding to its left-half-plane eigenvalues

$$\implies (I_{2n} + \text{sign}(H)) \begin{bmatrix} I_n \\ -X_{\geq} \end{bmatrix} = 0.$$

- Hence, X_{\geq} is determined by overdetermined, but consistent linear system of equations once $\text{sign}(H)$ is known.



So far:

**So far:**

- Solution of AREs via $\text{sign}(H)$, fully dense computations!

So far:

- Solution of AREs via $\text{sign}(H)$, fully dense computations!
- Newton iteration for $\text{sign}(H)$ preserves structure, as inversion and addition preserve Hamiltonian structure.

So far:

- Solution of AREs via $\text{sign}(H)$, fully dense computations!
- Newton iteration for $\text{sign}(H)$ preserves structure, as inversion and addition preserve Hamiltonian structure.
- But:
 - ① off-diagonal blocks are not treated in low-rank format,
 - ② X_{\geq} cannot be determined in factored form directly from this.

So far:

- Solution of AREs via $\text{sign}(H)$, fully dense computations!
- Newton iteration for $\text{sign}(H)$ preserves structure, as inversion and addition preserve Hamiltonian structure.
- But:
 - ① off-diagonal blocks are not treated in low-rank format,
 - ② X_{\geq} cannot be determined in factored form directly from this.

Goals

- ① Keep the off-diagonal blocks in H in low-rank form — this would save a significant amount of memory, i.e., working with A, B, C directly would reduce memory requirements by a factor of $\sim 3 - 4$.

So far:

- Solution of AREs via $\text{sign}(H)$, fully dense computations!
- Newton iteration for $\text{sign}(H)$ preserves structure, as inversion and addition preserve Hamiltonian structure.
- But:
 - ① off-diagonal blocks are not treated in low-rank format,
 - ② X_{\geq} cannot be determined in factored form directly from this.

Goals

- ① Keep the off-diagonal blocks in H in low-rank form — this would save a significant amount of memory, i.e., working with A, B, C directly would reduce memory requirements by a factor of $\sim 3 - 4$.
- ② Obtain X_{\geq} in low-rank factored form directly.

Theorem (Kenney/Laub/Jonckheere 1989, B. 2019/22)

Let (A, B) be *stabilizable*, (A, C) be *detectable*, and define the *Hamiltonian matrix*

$$\begin{bmatrix} A & -BB^T \\ -C^T C & -A^T \end{bmatrix}.$$

Theorem (Kenney/Laub/Jonckheere 1989, B. 2019/22)

Let (A, B) be *stabilizable*, (A, C) be *detectable*, and define the *Hamiltonian matrix*

$$\begin{bmatrix} A & -BB^T \\ -C^T C & -A^T \end{bmatrix}.$$

Then the unique stabilizing solution X_{\geq} to the LQR Riccati equation exists and is symmetric positive semidefinite.

Theorem (Kenney/Laub/Jonckheere 1989, B. 2019/22)

Let (A, B) be *stabilizable*, (A, C) be *detectable*, and define the *Hamiltonian matrix*

$$\begin{bmatrix} A & -BB^T \\ -C^T C & -A^T \end{bmatrix}.$$

Then the unique stabilizing solution X_{\geq} to the LQR Riccati equation exists and is symmetric positive semidefinite.

Hence, $A - BB^T X_s$ is stable, the closed-loop Lyapunov equations

$$(A - BB^T X_{\geq})P + P(A - BB^T X_{\geq})^T + BB^T = 0,$$

$$(A - BB^T X_{\geq})^T Q + Q(A - BB^T X_{\geq}) + C^T C = 0,$$

have unique solutions $P = P^T \geq 0$, $Q = Q^T \geq 0$, resp., and it holds

$$\text{sign}(H) = \begin{bmatrix} -I + 2PX_{\geq} & -2P \\ 2X_{\geq}PX_{\geq} - 2X_{\geq} & I - 2X_{\geq}P \end{bmatrix}.$$

Theorem (Kenney/Laub/Jonckheere 1989, B. 2019/22)

Let (A, B) be *stabilizable*, (A, C) be *detectable*, and define the *Hamiltonian matrix*

$$\begin{bmatrix} A & -BB^T \\ -C^T C & -A^T \end{bmatrix}.$$

Then the unique stabilizing solution X_{\geq} to the LQR Riccati equation exists and is symmetric positive semidefinite.

Hence, $A - BB^T X_s$ is stable, the closed-loop Lyapunov equations

$$\begin{aligned} (A - BB^T X_{\geq})P + P(A - BB^T X_{\geq})^T + BB^T &= 0, \\ (A - BB^T X_{\geq})^T Q + Q(A - BB^T X_{\geq}) + C^T C &= 0, \end{aligned}$$

have unique solutions $P = P^T \geq 0$, $Q = Q^T \geq 0$, resp., and it holds

$$\text{sign}(H) = \begin{bmatrix} -I + 2PX_{\geq} & -2P \\ 2X_{\geq}PX_{\geq} - 2X_{\geq} & I - 2X_{\geq}P \end{bmatrix}.$$

Hence, P (and by duality, Q) can be obtained from $\text{sign}(H)$ directly, without solving the AREs at all, and in factored form if sign iterates preserve the off-diagonal low-rank structure!



Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Important observation: M, \tilde{M} are **symmetric quasi-semidefinite** (SQSD).

Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Important observation: M, \tilde{M} are **symmetric quasi-semidefinite** (SQSD).

Here: inversion of SQSD matrices using **principal pivot transforms**

Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Important observation: M, \tilde{M} are **symmetric quasi-semidefinite** (SQSD).

Here: inversion of SQSD matrices using **principal pivot transforms**

- 1 is numerically more robust than standard inversion of symmetric matrices,

Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Important observation: M, \tilde{M} are **symmetric quasi-semidefinite** (SQSD).

Here: inversion of SQSD matrices using **principal pivot transforms**

- 1 is numerically more robust than standard inversion of symmetric matrices,
- 2 allows to work directly with A, B, C without ever forming $2n \times 2n$ -matrices,

Re-write the sign function iteration for $H = \begin{bmatrix} A & BB^T \\ C^T C & -A^T \end{bmatrix}$ in "symmetrized" form:

$$\frac{1}{2} (H + H^{-1}) = \frac{1}{2} (H + (J^T J H)^{-1}) = \frac{1}{2} (H J^T + (J H)^{-1}) J =: \frac{1}{2} (\tilde{M} + M^{-1}) J,$$

where

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \quad M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}, \quad \tilde{M} = \begin{bmatrix} BB^T & -A \\ -A^T & -C^T C \end{bmatrix}.$$

Important observation: M, \tilde{M} are **symmetric quasi-semidefinite** (SQSD).

Here: inversion of SQSD matrices using **principal pivot transforms**

- 1 is numerically more robust than standard inversion of symmetric matrices,
- 2 allows to work directly with A, B, C without ever forming $2n \times 2n$ -matrices,
- 3 yields a sign function iteration for AREs using A, B, C without ever forming $2n \times 2n$ -matrices!

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = [M^{(0)} \quad -I_n]$.

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = [M^{(0)} \quad -I_n]$.

Select a $u \times u$ pivot block $M_{11}^{(0)}$ ($u \in \{1, 2\}$) and permute $W^{(0)}$ accordingly, call the result again $W^{(0)}$.

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = [M^{(0)} \quad -I_n]$.

Select a $u \times u$ pivot block $M_{11}^{(0)}$ ($u \in \{1, 2\}$) and permute $W^{(0)}$ accordingly, call the result again $W^{(0)}$.

Define $K_0 = \begin{bmatrix} M_{11}^{(0)} & 0 \\ M_{21}^{(0)} & I_{n-u} \end{bmatrix}$ and compute

$$W^{(1)} := K_0^{-1}W^{(0)} = \begin{bmatrix} I_u & M_{12}^{(1)} & M_{11}^{(1)} & 0 \\ 0 & M_{22}^{(1)} & M_{21}^{(1)} & -I_{n-u} \end{bmatrix}.$$

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = \begin{bmatrix} M^{(0)} & -I_n \end{bmatrix}$.

Select a $u \times u$ pivot block $M_{11}^{(0)}$ ($u \in \{1, 2\}$) and permute $W^{(0)}$ accordingly, call the result again $W^{(0)}$.

Define $K_0 = \begin{bmatrix} M_{11}^{(0)} & 0 \\ M_{21}^{(0)} & I_{n-u} \end{bmatrix}$ and compute

$$W^{(1)} := K_0^{-1}W^{(0)} = \begin{bmatrix} I_u & M_{12}^{(1)} & M_{11}^{(1)} & 0 \\ 0 & M_{22}^{(1)} & M_{21}^{(1)} & -I_{n-u} \end{bmatrix}.$$

Then

$$M^{(1)} := \begin{bmatrix} M_{11}^{(1)} & M_{12}^{(1)} \\ M_{21}^{(1)} & M_{22}^{(1)} \end{bmatrix} := \begin{bmatrix} -M_{11}^{-1} & M_{11}^{-1}M_{12} \\ M_{21}M_{11}^{-1} & M_{22} - M_{21}M_{11}^{-1}M_{12} \end{bmatrix}.$$

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = [M^{(0)} \quad -I_n]$.

Select a $u \times u$ pivot block $M_{11}^{(0)}$ ($u \in \{1, 2\}$) and permute $W^{(0)}$ accordingly, call the result again $W^{(0)}$.

Define $K_0 = \begin{bmatrix} M_{11}^{(0)} & 0 \\ M_{21}^{(0)} & I_{n-u} \end{bmatrix}$ and compute

$$W^{(1)} := K_0^{-1}W^{(0)} = \begin{bmatrix} I_u & M_{12}^{(1)} & M_{11}^{(1)} & 0 \\ 0 & M_{22}^{(1)} & M_{21}^{(1)} & -I_{n-u} \end{bmatrix}.$$

Then

$$M^{(1)} := \begin{bmatrix} M_{11}^{(1)} & M_{12}^{(1)} \\ M_{21}^{(1)} & M_{22}^{(1)} \end{bmatrix} := \begin{bmatrix} -M_{11}^{-1} & M_{11}^{-1}M_{12} \\ M_{21}M_{11}^{-1} & M_{22} - M_{21}M_{11}^{-1}M_{12} \end{bmatrix}.$$

The mapping $M^{(0)} \rightarrow M^{(1)}$ is called **principal pivot transform (PPT)**.

Definition (Principal Pivot Transform)

Let M be symmetric and invertible, and set $M^{(0)} := M$, $W^{(0)} = [M^{(0)} \quad -I_n]$.

Select a $u \times u$ pivot block $M_{11}^{(0)}$ ($u \in \{1, 2\}$) and permute $W^{(0)}$ accordingly, call the result again $W^{(0)}$.

Define $K_0 = \begin{bmatrix} M_{11}^{(0)} & 0 \\ M_{21}^{(0)} & I_{n-u} \end{bmatrix}$ and compute

$$W^{(1)} := K_0^{-1}W^{(0)} = \begin{bmatrix} I_u & M_{12}^{(1)} & M_{11}^{(1)} & 0 \\ 0 & M_{22}^{(1)} & M_{21}^{(1)} & -I_{n-u} \end{bmatrix}.$$

Then

$$M^{(1)} := \begin{bmatrix} M_{11}^{(1)} & M_{12}^{(1)} \\ M_{21}^{(1)} & M_{22}^{(1)} \end{bmatrix} := \begin{bmatrix} -M_{11}^{-1} & M_{11}^{-1}M_{12} \\ M_{21}M_{11}^{-1} & M_{22} - M_{21}M_{11}^{-1}M_{12} \end{bmatrix}.$$

The mapping $M^{(0)} \rightarrow M^{(1)}$ is called **principal pivot transform (PPT)**.

Repeating this m times with pivots of size u_k , so that $u_0 + \dots + u_{m-1} = n$, yields $M^{(m)} = -M^{-1}$, i.e., a **Gauß-Jordan-type inversion** procedure for symmetric matrices.



Inversion of Symmetric Matrices

Most software packages compute inverses of symmetric matrices M using LDL^T factorization with **Bunch-Kaufman** (diagonal, partial) or **Bunch-Parlett** (complete) pivoting, e.g., xSYTRI from LAPACK and the MATLAB function `inv` based on this. SQSD structure is usually ignored, but turns out to be beneficial!

Most software packages compute inverses of symmetric matrices M using LDL^T factorization with **Bunch-Kaufman** (diagonal, partial) or **Bunch-Parlett** (complete) pivoting, e.g., xSYTRI from LAPACK and the MATLAB function `inv` based on this. SQSD structure is usually ignored, but turns out to be beneficial!

Theorem (Bunch-Parlett)

Let $LDL^T = \Pi M \Pi^T$ be the LDL^T factorization with Bunch-Parlett pivoting of a symmetric matrix M , with pivoting threshold $\tau = \frac{1+\sqrt{17}}{8} \approx 0.64$. Then,

$$\|D\|_{\max} \leq (2.57)^{n-1} \|M\|_{\max}, \quad \text{and} \quad \|L\|_{\max} \leq 2.78.$$

Here: a scalar pivot is chosen if $\max_{k=1, \dots, n} |M[k, k]| \geq \tau \max_{i \neq j} |M[i, j]|$, if such k exists; otherwise maximum 2×2 pivot is chosen.

Most software packages compute inverses of symmetric matrices M using LDL^T factorization with **Bunch-Kaufman** (diagonal, partial) or **Bunch-Parlett** (complete) pivoting, e.g., xSYTRI from LAPACK and the MATLAB function `inv` based on this. SQSD structure is usually ignored, but turns out to be beneficial!

Theorem (Bunch-Parlett)

Let $LDL^T = \Pi M \Pi^T$ be the LDL^T factorization with Bunch-Parlett pivoting of a symmetric matrix M , with pivoting threshold $\tau = \frac{1+\sqrt{17}}{8} \approx 0.64$. Then,

$$\|D\|_{\max} \leq (2.57)^{n-1} \|M\|_{\max}, \quad \text{and} \quad \|L\|_{\max} \leq 2.78.$$

Here: a scalar pivot is chosen if $\max_{k=1, \dots, n} |M[k, k]| \geq \tau \max_{i \neq j} |M[i, j]|$, if such k exists; otherwise maximum 2×2 pivot is chosen.

Worst-case element growth can be slightly improved for SQSD matrices:

Theorem (B./Poloni 2019)

Let $LDL^T = \Pi M \Pi^T$ be the LDL^T factorization with Bunch-Parlett pivoting of a SQSD matrix M , with pivoting threshold $\tau = 1$. Then,

- 1 $\|D\|_{\max} \leq 2^{n-1} \|M\|_{\max}$, and $\|L\|_{\max} \leq 2$.
- 2 $\| |D| |D^{-1}| \|_{\max} \leq 2$, and $\| |D| |D^{-1}| \| \leq 3$.

Theorem (Backward stability of symmetric GJE)

Let \hat{X} be the approximation of $X = -M^{-1}$ computed by the PPT-based symmetric Gauss–Jordan elimination algorithm. Then, each column $\hat{x}_j = \hat{X}e_j$ satisfies

$$-e_j = (M + \Delta_j)\hat{x}_j, \quad |\Delta_j| \leq |M| \left|L^{-T}\right| \left|L^T\right| \varepsilon_n,$$

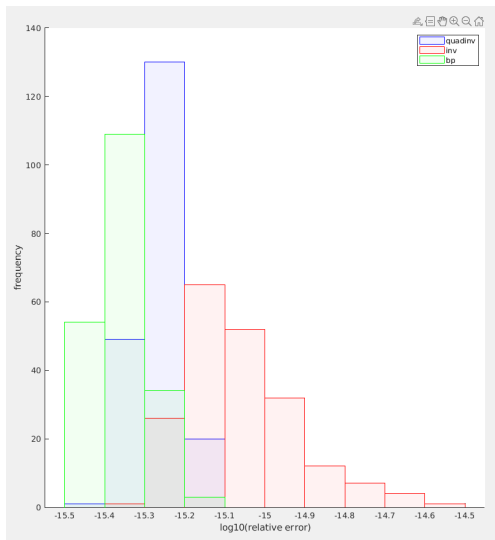
where $\varepsilon_n := \frac{cn\mathbf{u}}{1-cn\mathbf{u}}$ with a constant c independent of n .

Inversion of 200 random 200×200
SQD matrices with

- MATLAB `inv`, based on
DSYTRI from LAPACK
- Bunch-Parlett with complete
pivoting,
- structured inversion using
PPTs.

Inversion of 200 random 200×200 SQD matrices with

- **MATLAB `inv`**, based on **DSYTRI** from **LAPACK**
- **Bunch-Parlett** with **complete pivoting**,
- **structured inversion** using **PPTs**.





Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.



Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.

- For inversion of $M = \begin{bmatrix} C^T C & -A^T \\ -A & -B B^T \end{bmatrix}$, use symmetric GJE based on PPT.



Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.

- For inversion of $M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}$, use symmetric GJE based on PPT.
- The inverse of a SQSD matrix is again SQSD, i.e.,

$$M^{-1} = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}^{-1} = \begin{bmatrix} \hat{C}^T \hat{C} & -\hat{A}^T \\ -\hat{A} & -\hat{B} \hat{B}^T \end{bmatrix}.$$

Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.

- For inversion of $M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}$, use symmetric GJE based on PPT.
- The inverse of a SQSD matrix is again SQSD, i.e.,

$$M^{-1} = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}^{-1} = \begin{bmatrix} \hat{C}^T \hat{C} & -\hat{A}^T \\ -\hat{A} & -\hat{B} \hat{B}^T \end{bmatrix}.$$

- Inversion can be implemented using A, B, C only using again PPT-variant applied to **generator matrix**

$$\mathfrak{G} = \begin{bmatrix} B & A \\ * & C \end{bmatrix},$$

i.e., compute $\mathfrak{X} = \begin{bmatrix} \hat{B} & \hat{A} \\ * & \hat{C} \end{bmatrix}$ representing M^{-1} using A, B, C only without ever forming M
 [POLONI/STRABIĆ 2016]!



Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.

- For inversion of $M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}$, use symmetric GJE based on PPT.
- The inverse of a SQSD matrix is again SQSD, i.e.,

$$M^{-1} = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}^{-1} = \begin{bmatrix} \hat{C}^T \hat{C} & -\hat{A}^T \\ -\hat{A} & -\hat{B} \hat{B}^T \end{bmatrix}.$$

- Inversion can be implemented using A, B, C only using again PPT-variant applied to **generator matrix**

$$\mathfrak{G} = \begin{bmatrix} B & A \\ * & C \end{bmatrix},$$

i.e., compute $\mathfrak{X} = \begin{bmatrix} \hat{B} & \hat{A} \\ * & \hat{C} \end{bmatrix}$ representing M^{-1} using A, B, C only without ever forming M
[POLONI/STRABIĆ 2016]!

- Update $M \rightarrow M_+$ can then be performed also on generators (potentially using rank truncation for offline blocks).

Recall: $M_+ =: \frac{1}{2} (\tilde{M} + M^{-1}) J$ with M, \tilde{M} SQSD.

- For inversion of $M = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}$, use symmetric GJE based on PPT.
- The inverse of a SQSD matrix is again SQSD, i.e.,

$$M^{-1} = \begin{bmatrix} C^T C & -A^T \\ -A & -BB^T \end{bmatrix}^{-1} = \begin{bmatrix} \hat{C}^T \hat{C} & -\hat{A}^T \\ -\hat{A} & -\hat{B} \hat{B}^T \end{bmatrix}.$$

- Inversion can be implemented using A, B, C only using again PPT-variant applied to **generator matrix**

$$\mathfrak{G} = \begin{bmatrix} B & A \\ * & C \end{bmatrix},$$

i.e., compute $\mathfrak{X} = \begin{bmatrix} \hat{B} & \hat{A} \\ * & \hat{C} \end{bmatrix}$ representing M^{-1} using A, B, C only without ever forming M
 [POLONI/STRABIĆ 2016]!

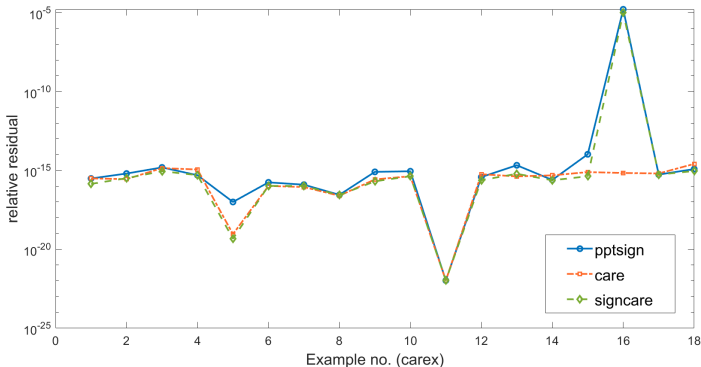
- Update $M \rightarrow M_+$ can then be performed also on generators (potentially using rank truncation for offline blocks).

↪ new version of sign function iteration working directly on generators!

Test the new sign function iteration for AREs based on PPTs ([pptsign](#))

- vs. MATLAB function `care` (based on Schur vector method) and classical sign function method (`signcare`) from MORLab [B./WERNER 2006–2023],
- using 18 examples from carex benchmark collection [B./LAUB/MEHRMANN 1995].

Measure accuracy by $\frac{\|\mathcal{R}(\tilde{X})\|_F}{\|C^T C\|_F + 2\|A\|_F \|\tilde{X}\|_F + \|B B^T\|_F \|\tilde{X}^2\|_F}$.



- Symmetric quasi-semidefinite matrices can be inverted using PPT-based Gauß-Jordan type elimination in a structure-preserving and numerically robust way.
- Sign function iteration for AREs can be reformulated in terms of SQSD matrix inversions and summations, allowing to work with generator matrices (A, B, C) only, without ever forming $2n \times 2n$ matrices.
- Leads to much lower storage requirements and potentially to faster algorithms (fewer flops).
- Application: closed-loop balanced truncation without ever solving AREs.
- **Future work:** sophisticated implementation to really test performance.



P. Benner.

An alternative algorithm for unstable balanced truncation.

MATHMOD 2022 DISCUSSION CONTRIBUTIONS, ARGESIM Report 17, pp. 71–72, 2022.

DOI: [10.11128/arep.17.a17178](https://doi.org/10.11128/arep.17.a17178)



P. Benner, P. Ezzatti, E.S. Quintana-Ortí, A. Remón.

A factored variant of the Newton iteration for the solution of algebraic Riccati equations via the matrix sign function.

NUMERICAL ALGORITHMS, 66:363–377, 2014.



P. Benner, M. Köhler, and J. Saak.

Matrix equations, sparse solvers: M-M.E.S.S.-2.0.1 – philosophy, features and application for (parametric) model order reduction.

In MODEL REDUCTION OF COMPLEX DYNAMICAL SYSTEMS, P. Benner, T. Breiten, H. Faßbender, M. Hinze, T. Stykel, and R. Zimmermann, eds., ISNM 171, Birkhäuser, Cham, pp. 369–392, 2021.



P. Benner and S. W. R. Werner.

MORLAB – A model order reduction framework in MATLAB and Octave.

In MATHEMATICAL SOFTWARE – ICMS 2020, A. M. Bigatti, J. Carette, J. H. Davenport, M. Joswig, and T. de Wolff, eds., LNCS 12097, Springer International Publishing, Cham, pp. 432–441, 2020.



C. Kenney, A. J. Laub, and E. A. Jonckheere.

Positive and negative solutions of dual Riccati equations by matrix sign function iteration.

SYSTEMS CONTROL LETTERS, 13:109–116, 1989.



F. Poloni and N. Strabić.

Principal pivot transforms of quasidefinite matrices and semidefinite Lagrangian subspaces.

ELECTRONIC JOURNAL OF LINEAR ALGEBRA, 31:200–231, 2016.

METT X

10th Workshop on Matrix Equations and Tensor Techniques

September 13–15, 2023

RWTH Aachen University (main building)

<https://www.igpm.rwth-aachen.de/workshop/mett2023>



Special Issue of ETNA (Electronic Transactions on Numerical Analysis),
open for participants only!

Fully (diamond) Open Access without OA charges!