

## The matrix square root

Next (and last, for us) matrix function:  $A^{1/2}$ , principal square root.

$A^{1/2}$  is well defined unless  $A$  has:

- ▶ Real eigenvalues  $\lambda_i < 0$ , or
- ▶ Non-trivial Jordan blocks at  $\lambda_i = 0$  (because  $g(x) = x^{1/2}$  is not differentiable).

## Condition number / sensitivity

The Fréchet derivative of  $f(X) = X^2$  is

$$L_{f,X}(E) = XE + EX, \quad \hat{L} = I \otimes X + X^T \otimes I.$$

The Fréchet derivative of  $g(Y) = Y^{1/2}$  is its inverse,

$$\hat{L}_{g,Y} = (I \otimes Y^{1/2} + (Y^{1/2})^T \otimes I)^{-1}$$

with eigenvalues  $\frac{1}{\lambda_i^{1/2} + \lambda_j^{1/2}}$ ,  $i, j = 1, \dots, n$ .

In particular,  $g$  is ill-conditioned for matrices that either:

- ▶ have a small eigenvalue (taking  $i = j$ ), or
- ▶ have two complex conjugate eigenvalues close to the negative real axis (because then  $\lambda_i^{1/2} \approx ai$ ,  $\lambda_j^{1/2} \approx -ai$ ).

## Modified Schur method

Recall: Schur method:

1. Reduce to a triangular  $T$  using a Schur form;
2. Compute diagonal of  $S = f(T)$ ;
3. Compute off-diagonal entries from  $ST = TS$   
Involves a denominator  $t_{ii} - t_{jj}$ : if it is 0, we must work on blocks.

In the case of  $A^{1/2}$ , we can use  $S^2 = T$  to get the off-diagonal entries instead:

$$s_{ii}s_{ij} + s_{i,i+1}s_{i+1,j} + \cdots + s_{ij}s_{jj} = t_{ij}.$$

This involves a denominator  $s_{ii} + s_{jj}$ : always invertible because  $s_{ii} + s_{jj} \in RHP$ .

This is (more or less) what Matlab uses, by the way (it does it in a divide-and-conquer way).

## Stability of Schur method for `sqrtrm`

Rounding error analysis for

$$s_{ij} = \frac{t_{ij} - s_{i,i+1}s_{i+1,j} - \cdots - s_{i,j-1}s_{j-1,j}}{s_{ii} + s_{jj}}.$$

(replacing each  $s_{ij}$  with the computed  $\tilde{s}_{ij}$ , and considering errors in all operations) leads to

$$\tilde{S}^2 = T + \delta T, \quad |\delta T| \leq |S|^2 \mathcal{O}(nu).$$

Combining it with a (backward-stable) Schur factorization and switching to norm, we get that  $X = A^{1/2}$  is computed with **backward** error

$$\|\hat{X}^2 - A\|_F \leq \mathcal{O}(n^3 u) \|X\|_F^2,$$

This is weaker than backward stability: there could be cancellation in the product  $X^2$ , so  $\|X\|_F^2$  is not really the same thing as  $\|A\|_F$ .

## Newton method

Newton method on  $X^2 - A$ :

$$X_{k+1} = X_k - E, \quad \text{where } E \text{ solves } EX_k + X_k E = X_k^2 - A.$$

Much more expensive than the Schur method: we solve a Sylvster equation at each step (and this requires a Schur form).

**Trick:** If  $X_0$  commutes with  $A$  (for instance, taking  $X_0 = \alpha I$ ), then  $E = (2X_0)^{-1}(X_0^2 - A)$  solves that equation; then  $E, X_1$  commute with  $A$ , too, and so on.

Resulting iteration:

(Modified) Newton iteration (MN)

$$X_{k+1} = \frac{1}{2}(X_k + X_k^{-1}A), \quad X_0 = \alpha I.$$

At each step,  $X_k A = A X_k$ .

## Square root and sign

### Theorem

Assume  $A$  has no eigenvalues in  $\mathbb{R}^-$ . Then, the MN iteration converges to the principal square root  $A^{1/2}$  for each starting point  $X_0 = \alpha I$  or  $X_0 = \alpha A$ , with  $\alpha > 0$ .

**Proof** Pre-multiply by  $A^{-1/2}$ , and use commutativity:

$$A^{-1/2}X_{k+1} = \frac{1}{2} \left( A^{-1/2}X_k + (A^{-1/2}X_k)^{-1} \right).$$

This is the sign iteration! Hence  $A^{-1/2}X_k \rightarrow \text{sign}(A^{-1/2}X_0) = I$ .

**Remark:** if  $A$  has a negative eigenvalue  $\lambda < 0$ , there is another obstruction: neither version of Newton can converge, because  $A^{1/2}$  is non-real. To restore convergence, we need to add a small imaginary part to  $X_0$ .

## Theory and practice

**Problem** All of this holds in **exact arithmetic**, but the method won't work in practice in machine arithmetic! Try `rng(4); A = randn(5);`. These two iterations behave quite differently:

### True Newton

$$X_{k+1} = X_k - E, \quad \text{where } E \text{ solves } EX_k + X_k E = X_k^2 - A.$$

### Modified Newton

$$X_{k+1} = \frac{1}{2}(X_k + X_k^{-1}A).$$

TN converges, but MN diverges after an initial “pseudo-convergence”.

Numerically, the two sequences behave quite differently, and the commutativity property  $X_k A = A X_k$  is lost in MN in a few iterations.

## Local stability

The **geometric picture** The two iterations coincide on the manifold of matrices that commute with  $A$ ,  $\{X \in \mathbb{C}^{n \times n} : AX = XA\}$ , but not on the rest of  $\mathbb{C}^n$ .

Numerical perturbations take us outside of the manifold, and then they do not coincide anymore.

While TN is quadratically convergent, MN does not even have an **stable fixed point** in  $A^{1/2}$ : even when started very close to  $A^{1/2}$ , the iteration diverges.

We can prove this formally.



## Local stability

Local stability of a fixed point of  $X_{k+1} = h(X_k)$  depends on the eigenvalues of its Jacobian.

The Jacobian / Fréchet derivative of  $h(X) = \frac{1}{2}(X + X^{-1}A)$  is

$$L_{h,X}(E) = \frac{1}{2}(E + X^{-1}EX^{-1}A)$$

(use  $(X + E)^{-1} - X^{-1} = (X + E)^{-1}EX^{-1} = X^{-1}EX^{-1} + o(\|E\|)$ ).

Hence  $L_{h,A^{1/2}} = \frac{1}{2}(E + A^{-1/2}EA^{1/2})$ , or

$$\hat{L}_{h,A^{1/2}} = \frac{1}{2} \left( I + (A^{1/2})^T \otimes A^{-1/2} \right).$$

It has eigenvalues  $\frac{1}{2} + \frac{1}{2}\lambda_i^{1/2}\lambda_j^{-1/2}$ , where  $\lambda_i$  are the eigenvalues of  $A$ .

It is easy to construct examples in which  $L_{h,A^{1/2}}$  has eigenvalues with modulus  $> 1$ , hence  $A^{1/2}$  is an **unstable fixed point** of  $h(X)$ .

## Denman–Beavers iteration

However, the stability properties are significantly different for slight variations of the modified Newton's method.

Setting  $Y_k = A^{-1}X_k$ , we can get

Denman–Beavers iteration [Denman–Beavers, '76]

$$\begin{aligned}X_{k+1} &= \frac{1}{2}(X_k + Y_k^{-1}), \\Y_{k+1} &= \frac{1}{2}(Y_k + X_k^{-1}),\end{aligned}$$

(It corresponds to using the relation with the matrix sign and using Newton for the matrix sign to compute  $\text{sign}\left(\begin{bmatrix} 0 & A \\ I & 0 \end{bmatrix}\right)$ .)

## Local stability of the DB iteration

### Theorem

The DB iteration satisfies  $\lim(X_k, Y_k) = (A^{1/2}, A^{-1/2})$ , and it is locally stable.

We have

$$L_{DB,(X,Y)}\left(\begin{bmatrix} E \\ F \end{bmatrix}\right) = \frac{1}{2} \begin{bmatrix} E - Y^{-1}FY^{-1} \\ F - X^{-1}EX^{-1} \end{bmatrix}$$

All  $(X, Y) = (M, M^{-1})$  are fixed points, and in these the Jacobian is **idempotent**, i.e.,  $(K_{DB,(B,B^{-1})})^2 = K_{DB,(B,B^{-1})}$ .

Hence its eigenvalues are 0, 1, and all the Jordan blocks are simple  $\implies$  bounded powers  $\implies$  local stability.

Other variants are available [Higham book, Ch. 6].