

# True Newton vs. Modified Newton for sqrtm

Note Title

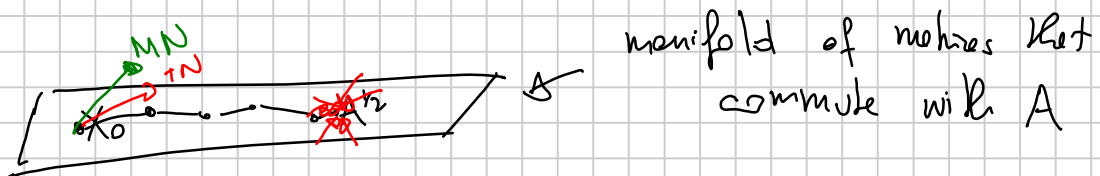
2023-04-20

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

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

$X_0$  commutes with  $A \Rightarrow$  TN, MN coincide in exact arithmetic...  
but not in machine arithmetic

Geometric picture:

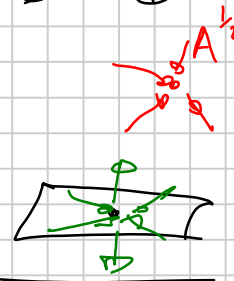


TN, MN coincide on the manifold, but not on all  $\mathbb{C}^{n \times n}$

When one allows for errors, two different iterations arise

TN is quadratically convergent to  $A^{1/2}$  on all  $\mathbb{C}^{n \times n}$

MN is convergent only on the manifold!



Dynamical system (in discrete time)

is defined by a map  $F: \mathbb{C}^m \rightarrow \mathbb{C}^m$

We are interested in studying the iteration  $\begin{cases} X_0 \text{ given} \\ X_{k+1} = F(X_k) \end{cases}$

If  $X_* = F(X_*)$ ,  $X_*$  is said to be a fixed point

We are interested in studying the iteration in a neighborhood of  $x_*$ , using its Jacobian  $J_{F, x_*}$  (which we assume exists)

$$x_0 = x_* + e \quad \text{for a small perturbation } e \in \mathbb{C}^m \quad x_* \in \mathbb{C}^m$$

$$x_1 = F(x_0 + e) = F(x_*) + J_{F, x_*} e + o(\|e\|)$$

$$x_2 = F(x_1) = F(x_* + J_{F, x_*} e + o(\|e\|)) = x_* + J_{F, x_*}^2 e + o(\|e\|)$$

$$\vdots$$

$$x_k = x_* + J_{F, x_*}^k e + o(\|e\|)$$

If  $\rho(J_{F, x_*}) < 1$ ,  $J_{F, x_*}^k \rightarrow 0$ , hence the iteration started from  $x_* + e$  will converge to  $x_*$

$x_*$  is stable

If  $\rho(J_{F, x_*}) > 1$ ,  $\|J_{F, x_*}^k\| \rightarrow \infty$ , and at least for some choices of  $e$  the iteration will not converge to  $x_*$

$x_*$  is unstable

For TN,  $J_{TN, A^{\frac{1}{2}}} = 0$ , and this is what makes the Newton method converge quadratically.

We would like to see if MN has a stable or unstable fixed point in  $A^{\frac{1}{2}}$ , i.e. if  $\rho(J_{MN, A^{\frac{1}{2}}})$  is smaller or larger than 1.

$$\text{MN: } F: X \rightarrow \frac{1}{2}(X + X^{-1}A)$$

The Jacobian of  $F$  (after vectorization) coincides with the Fréchet derivative of  $F$ .

To do this, we need to find the Fréchet derivative of the map  $X \mapsto X^{-1}$

$$(I-M)^{-1} = I+M+M^2+M^3+\dots \text{ if } \rho(M) < 1.$$

$$(X+E)^{-1} = (X(1+X^{-1}E))^{-1} = (1+X^{-1}E)^{-1} X^{-1} = (1 - X^{-1}E + X^{-1}EX^{-1}E - X^{-1}EX^{-1}EX^{-1}E + \dots) X^{-1}$$

$$= X^{-1} + \underbrace{X^{-1}EX^{-1}}_{\downarrow} + O(\|E\|^2)$$

$$\downarrow \\ \mathcal{L}_{\text{inv}, X}[E]$$

$$\mathcal{L}_{F, X}[E]$$

So

$$F(X+E) = \frac{1}{2}(X+E + (X+E)^{-1}A) = \frac{1}{2}(X + X^{-1}A) + \frac{1}{2}(E - X^{-1}EX^{-1}A) + O(\|E\|^2)$$

$$\hat{\mathcal{L}}_{F, X} = \frac{1}{2} \left( I_{n^2} - (X^{-1}A)^T \otimes X^{-1} \right)$$

$$\hat{\mathcal{L}}_{F, A^{1/2}} = \frac{1}{2} \left( I_{n^2} - (A^{1/2})^T \otimes A^{-1/2} \right)$$

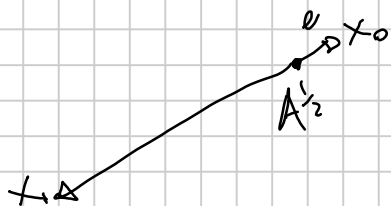
Assuming that  $A^{-1/2}, A^{1/2}$  are put in Schur form, we obtain that if  $\Lambda(A) = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ , then

$$\Lambda(\hat{\mathcal{L}}_{F, A^{1/2}}) = \left\{ \frac{1}{2} \left( 1 - \lambda_i^{1/2} \cdot \lambda_j^{-1/2} \right) : i, j = 1, 2, \dots, n \right\}$$

When  $A$  is ill-conditioned,  $|\lambda_{\max}/\lambda_{\min}|$  is large, and

$$\frac{1}{2} \left( 1 - \left( \frac{\lambda_{\max}}{\lambda_{\min}} \right)^{1/2} \right) \text{ can be large, too.}$$

So MN is unstable: even if started very close to  $X_* = A^{1/2}$ ,



the machine arithmetic error in the iterations can be amplified and make the method diverge.

So MN  $\rightarrow$  can't be used, unstable

TN  $\rightarrow$  can't be used, too expensive

But variants of MN have a very different stability behavior

Variant: Denman-Beavers iteration (DB)

Start from  $X_{k+1} = \frac{1}{2} (X_k + \underbrace{X_k^{-1} A}_{Y_k^{-1}})$

$Y_k^{-1} := X_k^{-1} A$

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

$$\begin{cases} X_{k+1} = \frac{1}{2} (X_k + Y_k^{-1}) & X_0 = A \\ Y_{k+1} = \frac{1}{2} (Y_k + X_k^{-1}) & Y_0 = I \end{cases}$$

DB iteration

EX: show that DB is equivalent to applying the Newton iteration for the sign to compute

$$\text{sign} \left( \begin{bmatrix} 0 & A \\ I & 0 \end{bmatrix} \right) = \begin{bmatrix} 0 & A^{\frac{1}{2}} \\ A^{-\frac{1}{2}} & 0 \end{bmatrix}$$

Sol: Newton for the sign reads

$$z_k = \begin{bmatrix} 0 & X_k \\ Y_k & 0 \end{bmatrix} \quad z_{k+1} = \frac{1}{2} (z_k + z_k^{-1}) = \frac{1}{2} \left( \begin{bmatrix} 0 & X_k \\ Y_k & 0 \end{bmatrix} + \begin{bmatrix} 0 & Y_k^{-1} \\ X_k^{-1} & 0 \end{bmatrix} \right) \\ = \begin{bmatrix} 0 & \frac{1}{2} (X_k + Y_k^{-1}) \\ \frac{1}{2} (Y_k + X_k^{-1}) & 0 \end{bmatrix} \quad \checkmark$$

Let us study the stability of the DB iteration.

We are interested in the fixed point  $(X_*, Y_*) = (A^{\frac{1}{2}}, A^{-\frac{1}{2}})$

$$F \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(X + Y^{-1}) \\ \frac{1}{2}(Y + X^{-1}) \end{pmatrix} \quad \begin{array}{l} X \rightarrow X+E \\ Y \rightarrow Y+F \end{array}$$

$$L_{DB, \begin{pmatrix} * \\ * \end{pmatrix}} \begin{pmatrix} E \\ F \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(E - Y^{-1}FY^{-1}) \\ \frac{1}{2}(F - X^{-1}EX^{-1}) \end{pmatrix} \quad \hat{L}_{DB} \in \mathbb{C}^{2n^2 \times 2n^2}$$

Instead of computing eigenvalues, we note that  $\hat{L}_{DB, \begin{pmatrix} A^{1/2} \\ A^{-1/2} \end{pmatrix}}$  is idempotent:  $L_{DB}^2 = I$

$$\begin{aligned} L_{DB, \begin{pmatrix} A^{1/2} \\ A^{-1/2} \end{pmatrix}} \begin{pmatrix} E \\ F \end{pmatrix} &= \begin{pmatrix} \frac{1}{2} \left( \frac{1}{2}(E - Y^{-1}FY^{-1}) - Y^{-1} \left( \frac{1}{2}(F - X^{-1}EX^{-1}) \right) Y^{-1} \right) \\ \frac{1}{2} \left( \frac{1}{2}(F - X^{-1}EX^{-1}) - X^{-1} \left( \frac{1}{2}(E - Y^{-1}FY^{-1}) \right) X^{-1} \right) \end{pmatrix} \\ &= \begin{pmatrix} \frac{1}{4}E - \frac{1}{4}Y^{-1}FY^{-1} - \frac{1}{4}Y^{-1}FY^{-1} + \frac{1}{4} \overbrace{Y^{-1}X^{-1}}^F \overbrace{EX^{-1}Y^{-1}}^F \\ \text{"} \\ \frac{1}{4}E - \frac{1}{4}Y^{-1}FY^{-1} - \frac{1}{4}Y^{-1}FY^{-1} + \frac{1}{4} \overbrace{Y^{-1}X^{-1}}^F \overbrace{EX^{-1}Y^{-1}}^F \end{pmatrix} \\ &= \begin{pmatrix} \left( \frac{1}{4} + \frac{1}{4} \right) E - \left( \frac{1}{4} + \frac{1}{4} \right) Y^{-1}FY^{-1} \\ \text{"} \end{pmatrix} = L_{DB, \begin{pmatrix} A^{1/2} \\ A^{-1/2} \end{pmatrix}} \begin{pmatrix} E \\ F \end{pmatrix}. \end{aligned}$$

$\lambda = \lambda^2 \rightarrow \lambda \in \{0, 1\}$ , or, more directly, the powers

$\hat{L}_{DB, \begin{pmatrix} A^{1/2} \\ A^{-1/2} \end{pmatrix}}$  are all equal to  $\hat{L}_{DB, \begin{pmatrix} A^{1/2} \\ A^{-1/2} \end{pmatrix}}$  and hence are bounded.

$$X_0 = X_* + e$$

$$X_k = X_* + \left( J_{F, X_*} \right)^k e + o(\|e\|)$$

So DB does not amplify errors due to machine arithmetic  $\Rightarrow$  the iteration is (locally) stable.

DB is indeed an effective iteration, although it has 50%

higher computational cost.

$$\frac{1}{2}(X + X^{-1}A)$$

$$\frac{8}{3}n^3$$

$$O(n^2)$$

vs.

$$\frac{1}{2}(X + Y^{-1})$$

$$\frac{1}{2}(Y + X^{-1})$$

$$2n^3$$

$$O(n^2)$$

$$4n^3 + O(n^2)$$

EX show that Newton for the matrix sign has a stable fixed point.

(Numerical experiments)

Function of large-scale matrices.

Let  $A$  be large and sparse.

Problem:  $f(A)$  is often full, e.g.,  $f(x) = x^{-1}$

Solution: we look for methods that can compute  $f(A) \cdot b$  directly given  $A \in \mathbb{C}^{n \times n}$ ,  $b \in \mathbb{C}^n$ , analogously to linear systems.

1) polynomial or rational approximation:  $f(A)b \approx q^{-1}(A)p(A)b$

$b, Ab, A^2b, \dots \rightarrow \deg(p) \cdot n \cdot \text{nz}(A)$  cost

$q(A)^{-1}v$ , with iterative methods, can still be computed

with an iteration whose cost is  $n_{\text{steps}} \cdot (\text{cost of evaluating } q(A)v)$

$= n_{\text{steps}} \cdot \deg(q) \cdot n \cdot \text{nz}(A)$ .

Problem: need to find good rational approximation, and these require knowledge of  $\Lambda(A)$ .

2) Contour integration:

$$f(A) = \frac{1}{2\pi i} \int_{\Gamma} f(z) (zI - A)^{-1} dz \approx \sum_{k=1}^N w_k f(x_k) (x_k I - A)^{-1}.$$

for certain nodes  $x_1, \dots, x_N$  weights  $w_1, \dots, w_N$

$$f(A)b \approx \sum_{k=1}^N w_k f(x_k) \underbrace{(x_k I - A)^{-1} b}_{\text{linear system sol.}}$$

One can use methods for sparse linear systems, e.g. sparse LU.

Cost:  $N \cdot (\text{cost of solving } (x_k I - A)y = b)$

This is actually also a rational approximation:

$$f(z) \approx \sum_{k=1}^N w_k f(x_k) (x_k - z)^{-1}$$

So it is a special case of method 1, with the function in a particular form.

3) Methods based on the Arnoldi iteration.

Also rational approximations, but include an almost-optimal way to choose them and a framework for their evaluation.

$$f(A)b$$