

Newton method for the matrix square root:

$$X_0 = \alpha I \text{ or } X_0 = \alpha M$$

$$F(x) = x^2 - M$$

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

Modified Newton:

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

Using the fact that $X_0 M = M X_0$, one can prove that $X_k M = M X_k$ for each k , and that the two iterations produce the same sequence X_k in exact arithmetic.

Lemma: suppose M has no negative eigenvalues, then $X_k \rightarrow M^{\frac{1}{2}}$, the principal square root of M

Multiply the modified Newton method equation by $M^{\frac{1}{2}}$ to get

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

$$= \frac{1}{2} (M^{\frac{1}{2}} X_k + X_k^{-1} M^{\frac{1}{2}}) \text{ set } Y_k := M^{-\frac{1}{2}} X_k \text{ to get}$$

$$Y_{k+1} = \frac{1}{2} (Y_k + Y_k^{-1}) \quad \leftarrow \text{matrix sign iteration!}$$

$$\text{So } \lim_{k \rightarrow \infty} Y_k = \text{sign}(Y_0) = \text{sign}(M^{-\frac{1}{2}} X_0)$$

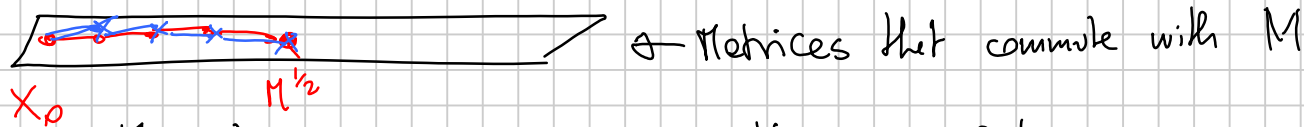
If $X_0 = \alpha M$, with $\alpha > 0$, then $M^{-\frac{1}{2}} X_0 = \alpha M^{\frac{1}{2}}$

Since $M^{\frac{1}{2}}$ has all eigenvalues in the RHP, $\text{sign}(M^{-\frac{1}{2}} X_0) = I$

$$\text{So } I = \lim_{k \rightarrow \infty} Y_k = \lim_{k \rightarrow \infty} M^{-\frac{1}{2}} X_k$$

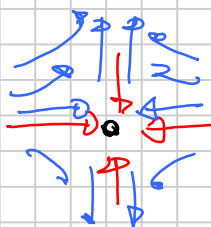
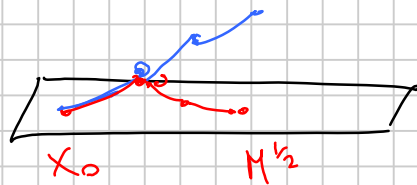
$$\text{And } \lim_{k \rightarrow \infty} X_k = M^{\frac{1}{2}}$$

If $X_0 = \alpha I$, $Y_0 = \alpha M^{-\frac{1}{2}}$ and $\text{sign}(Y_0) = I$ anyway.



Matrices that commute with M

The iterations coincide on this manifold



But they have very different behaviour outside of the manifold: TN is attractive, MN is repulsive, so the iterations diverge even when started very close to $M^{1/2}$.

TN is quadratically convergent thanks to the general theory of Newton methods.

Discrete-time dynamical system associated to the map

$$F: \mathbb{C}^N \rightarrow \mathbb{C}^N$$

$$x_{k+1} = F(x_k) \quad k=0, 1, 2, 3, \dots$$

Suppose we start close to a fixed point $x_* = F(x_*)$

$x_0 = x_* + e$, with $e \in \mathbb{C}^N$ small.

$$x_1 = F(x_0) = F(x_* + e) = F(x_*) + \underbrace{J_{F, x_*}}_{N \times N \text{ matrix}} \cdot e + O(\|e\|^2)$$

$N \times N$ matrix

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

$$x_k = x_* + (J_{F, x_*})^k e + O(\|e\|^2)$$

If $\rho(J_{F, X_*}) < 1$, the iterates get closer to X_*
 $\Rightarrow X_*$ is an attractive fixed point

If $\rho(J_{F, X_*}) > 1$, then for most starting e the powers
 $(J_{F, X_*})^k e$ get larger and larger and
the iterates do not approach X_*
 $\Rightarrow X_*$ is a repulsive fixed point.

(cfr. analysis of 1-variable fixed-point methods).



For modified Newton, we have

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

The Fréchet derivative of $F(x) = \frac{1}{2}(x + x^{-1}M)$

is obtained by

$$\begin{aligned} F(x+E) &= \frac{1}{2}(x+E + (x+E)^{-1}M) \\ &= \frac{1}{2}(x+E + (x^{-1} - x^{-1}E x^{-1} + x^{-1}E x^{-1}E x^{-1} + \dots)M) \\ &= \frac{1}{2}(x + x^{-1}M) + \frac{1}{2}(E - x^{-1}E x^{-1}M) + o(\|E\|) \end{aligned}$$

$$L_{F, X} [E] = \frac{1}{2}(E - X^{-1}E X^{-1}M)$$

$$L_{F, M^{\frac{1}{2}}} [E] = \frac{1}{2}(E - M^{-\frac{1}{2}} E M^{\frac{1}{2}})$$

Using Kronecker products, we can see it as an $n^2 \times n^2$ matrix:

$$K = \frac{1}{2} \left(I_{n^2} - (M^{\frac{1}{2}})^T \otimes M^{-\frac{1}{2}} \right)$$

$$\text{vec}(A \times B) = (B^T \otimes A) \text{vec } X$$

I can change basis so that this matrix becomes upper triangular:

take Schur forms $(M^{\frac{1}{2}})^T = Q_1 U_1 Q_1^*$, $M^{-\frac{1}{2}} = Q_2 U_2 Q_2^*$

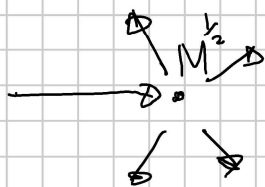
$$K = \frac{1}{2} (Q_1 \otimes Q_2) (I - U_1 \otimes U_2) (Q_1 \otimes Q_2)^*$$

On the diagonal, I can read off

$$\Lambda(K) = \left\{ \frac{1}{2} (1 - \lambda_i^{\frac{1}{2}} \cdot \lambda_j^{-\frac{1}{2}}) : i, j = 1, \dots, n \right\} \text{ where } \lambda_1, \dots, \lambda_n = \Lambda(M)$$

If M has two eigenvalues such that $\frac{\lambda_i}{\lambda_j}$ is large,

then $\rho(K) > 1 \Rightarrow M^{\frac{1}{2}}$ is a repulsive fixed point.



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

Denman-Beavers iteration: set $X_k^{-1} M = Y_k^{-1}$ $Y_k = M^{-1} X_k$
to get

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

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

$$= \frac{1}{2} (Y_k + X_k^{-1})$$

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

$$F: \mathbb{C}^{2n^2} \rightarrow \mathbb{C}^{2n^2}$$

The Jacobian of this iteration is idempotent: $J^2 = J$.

Functions of sparse matrices

$A \in \mathbb{C}^{n \times n}$ large and sparse

In general, $f(A)$ is going to be dense, but we can hope to compute $f(A)b$ for a vector b

e.g. to solve $\begin{cases} \dot{x} = Ax \\ x(0) = b \end{cases}$ we need to compute $x(t) = \exp(tA)b$

Problem: given f , A large and sparse, b , compute $f(A)b$.

1. Schur-based methods won't work

2. methods based on $f(A) \approx p(A)$ or $f(A) \approx r(A)$ (polynomial or rational approximation) work, because we can evaluate $p(A)b$, $q(A)^{-1}p(A)b$ even for large and sparse A .

$$A^3 b = A(A(Ab))$$

3. Cauchy formula:

$$f(A)b = \frac{1}{2\pi i} \int f(z)(zI - A)^{-1}b \, dz$$

$$\approx \sum_{k=1}^N w_k \underbrace{(z_k I - A)^{-1}b}_{\substack{\downarrow \\ \text{sparse linear system}}} \quad \text{for certain nodes } z_k \text{ and weights } w_k$$

This also reduces to a rational approximation:

$$r(z) = \sum_{k=1}^N w_k \frac{1}{z_k - z}$$

Problem in both cases: find $p(z)$ or $r(z)$ that approximates $f(z)$

in $\Lambda(A)$ (which itself is difficult to compute).

Now idea: use Arnoldi.

Def: Given $A \in \mathbb{C}^{m \times m}$, $b \in \mathbb{C}^m$ $m > n$

$$K_n(A, b) = \text{span}(b, Ab, A^2b, \dots, A^{n-1}b) \\ = \{p(A)b : p \in \text{polynomials of degree } d < n\}$$

In many problems, we can get good approximations of the solution by projecting it onto a Krylov subspace.

Let $V_n \in \mathbb{C}^{m \times n}$ be a matrix whose columns are an orthonormal basis of $K_n(A, b)$, then to solve $Ax = b$,

we project it to solve instead $(V_n^* A V_n) y = V_n^* b$

and then recover $x = V_n y$

Also, the eigenvalues of A are usually well approximated by those of $A_n = V_n^* A V_n \in \mathbb{C}^{n \times n}$ (Ritz values)

Usually $\Lambda(A_n)$ is a good approximation of the outer eigenvalues of A (those with largest modulus).

Arnoldi produces an orthonormal basis V_n for $K_n(A, b)$

$$V_n^* b = \begin{bmatrix} \beta \\ 0 \\ \vdots \\ 0 \end{bmatrix} = e_1 \beta \quad \beta = \|b\|.$$

Lemma: For all polynomials of degree $d < n$,

$$p(A)b = V_n p(\underbrace{V_n^* A V_n}_{A_n}) \underbrace{V_n^* b}_{e_1 \beta}.$$

$m \times n$ matrix

$n \times n$ smaller matrix

Proof: it is enough to prove that $A^j b = V_n A_n^j V_n^* b$
for all $j = 0, 1, 2, \dots, n-1$.

$\boxed{j=0}$: $b = V_n V_n^* b$? $V_n V_n^*$ is the orthogonal projection
onto $K_n(A, b)$

Since $b \in K_n(A, b)$, the projection does nothing,

so $b = V_n V_n^* b$

$\boxed{j=1}$ $Ab = V_n V_n^* A V_n V_n^* b$?

$V_n V_n^* b = b$, as above

$V_n V_n^* Ab = Ab$ because $Ab \in K_n(A, b)$

Similarly,

$V_n \underbrace{V_n^* A V_n V_n^* A \dots V_n^* A}_{j \text{ times } V_n^* A V_n} V_n b = A^j b$

as long as $j < n$, because we can cancel out these projections
one by one: $V_n V_n^* A^j b = A^j b$ if $A^j b \in K_n(A, b)$. \square

$f(A)b = V_n \underbrace{f(A_n)}_{n \times n} e, \beta$

$A_n = V_n^* A V_n \in \mathbb{C}^{n \times n}$

For a generic matrix function, we can hope that

$f(A)b \approx V_n f(A_n) e, \beta =: c$

is a good approximation.

Note that this is also a polynomial approximation:

$f(A_n) = p_n(A_n)$, where p_n is the interpolating

polynomial for f on $\Lambda(A_n)$ of degree $d < n$.

This is not the interpolating polynomial $p(z)$ of degree m on $\Lambda(A)$ for which $f(A) = p(A)$, but one with smaller degree.

$$c = V_n f(A_n) e_i \beta = V_n p_n(A_n) e_i \beta = p_n(A) b$$

$$f(A) b = p(A) b$$

$$p_n(\mu) = f(\mu) \text{ for Ritz values } \mu \in \Lambda(A_n)$$

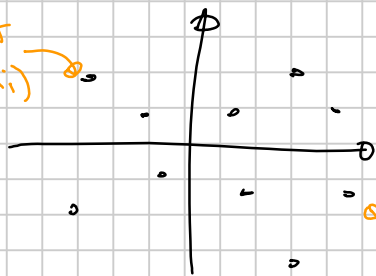
The Ritz values approximate well outer eigenvalues, so this gives us a good approximation on those

$$f(A) b = V \begin{bmatrix} f(\lambda_1) \\ \vdots \\ f(\lambda_m) \end{bmatrix} V^* b$$

$$p_n(A) b = V \begin{bmatrix} p_n(\lambda_1) \\ \vdots \\ p_n(\lambda_m) \end{bmatrix} V^* b$$

This will be a good approximation if the eigenvalues for which $f(\lambda_i)$ is larger are outer eigenvalues

smallest $\exp(\lambda_i)$



$$\text{E.g. } f(\lambda_i) = \exp(\lambda_i)$$

is larger when λ_i is rightmost in the complex plane

(Just an intuition for now, proofs in the next lecture.)