

$$\text{Sign}(z) = \begin{cases} +1 & z \in \text{right half-plane} \\ -1 & z \in \text{left h. p.} \end{cases}$$



$M$  with no eigenvalues on the imaginary axis

$$\text{sign}(M) \begin{cases} X_0 = M \\ X_{k+1} = \frac{1}{2} (X_k + X_k^{-1}) \end{cases}$$

If  $z$  is large,

$$\frac{1}{2}(z + z^{-1}) \approx \frac{1}{2}z$$

$$\text{If } z \ll 1, \quad \frac{1}{2}(z + z^{-1}) \approx z^{-1}$$

so convergence is slow when far away from  $\pm 1$

Idea to fix it:  $\text{sign}(M) = \text{sign}(\alpha M)$  for each  $\alpha > 0$ .

We would like to take the eigenvalues as close to 1 as possible by choosing  $\alpha$  appropriately

If  $\alpha\lambda_1, \dots, \alpha\lambda_n$  are the eigenvalues of  $\alpha M$ , then

$$\text{we want } \sqrt[n]{|(\alpha\lambda_1)(\alpha\lambda_2) \dots (\alpha\lambda_n)|} = 1$$

$$\Leftrightarrow \alpha \sqrt[n]{|\lambda_1 \lambda_2 \dots \lambda_n|} = 1$$

$$\Leftrightarrow \alpha = \frac{1}{|\det(M)|^{1/n}} \quad M \in \mathbb{C}^{n \times n}$$

This is called determinantal scaling

It comes for free, since the same decomposition (e.g. LU) can be used to compute  $X^{-1}$  and  $\det(X)$

Alternatives: use power methods to approximate  $\lambda_1(X_k)$ ,  $\lambda_n(X_k)$  and balance them.

---

Remarks on stability:

The sign function is sensitive to perturbation:

$$M = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix} \quad \text{sign}(M) = \begin{bmatrix} -I & Z \\ 0 & I \end{bmatrix}$$

$Z$  is obtained via a Sylvester equation, if  $\text{sep}(A, B) \ll 1$ ,  $Z$  is large.

$$\text{If } \tilde{M} = M + E \quad \|E\| = \varepsilon \quad \|\tilde{Z} - Z\| \leq \frac{\varepsilon}{\text{sep}(A, B)^2}$$

However, if we are only interested in computing invariant subspaces,

$$\text{Ker}(\text{sign}(M) + I)$$

Perturbing  $M$  to  $M + E$  changes the invariant subspace only by  $\frac{\varepsilon}{\text{sep}(A, B)}$

In practice, this method produces invariant subspaces as accurate as the reoriented Schur method.

Note that  $X_0$  ill-conditioned  $\Leftrightarrow$  problem is ill-conditioned.

---

Matrix square root

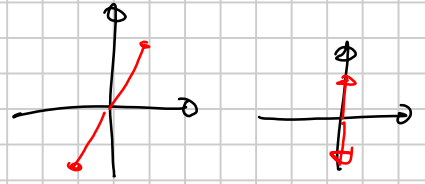
Principal square root:

for  $z \in \mathbb{C}$ , if  $z \notin \{\text{negative reals}\}$ ,

$f(z) = z^{1/2}$  is the square root in the RHP. (right half-plane)

$$f(0) = 0^{1/2} = 0$$

$f(\text{negative real})$  undefined  
(there is no continuous way to do it)



Note that  $(z^2)^{1/2} = z \cdot \text{sign}(z)$

so  $\text{sign}(M) = (M^2)^{1/2} \cdot M^{-1}$   
Also, one can prove that

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

Schur-Parlett for the square root:

Idea:  $M = QUQ^*$

• compute  $S = U^{1/2}$  (triangular factor)

•  $M^{1/2} = QU^{1/2}Q^*$

To compute  $f(S)$ , the standard Schur-Parlett method involves recursions of the form

$$S_{ij} = \frac{(\text{numerator})}{U_{jj} - U_{ii}}$$

and has trouble if  $M$  has close eigenvalues

For the matrix square root, we can avoid the problem:

we compute  $S$  not from  $SU = US$ , but from  $S^2 = U$

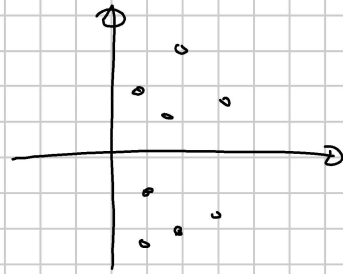
$$U_{ij} = (S^2)_{ij} = S_{ii} S_{ij} + S_{i,i+1} S_{i+1,j} + \dots + S_{i,j-1} S_{j-1,j} + S_{jj} S_{ij}$$

$$\Rightarrow S_{ij} = \frac{U_{ij} - \sum_{i < k < j} S_{ik} S_{kj}}{S_{ii} + S_{jj}}$$

The denominator does not involve  $U_{jj} - U_{ii}$ .

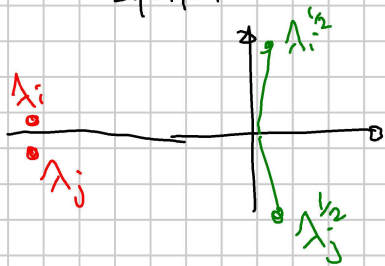
$$S_{ii} + S_{jj} = \lambda_i^{1/2} + \lambda_j^{1/2}$$

$\lambda_{ii}, \lambda_{jj} \in \text{RHP} \Rightarrow \text{Re}(\lambda_{ii} + \lambda_{jj}) > 0$   
and the denominator cannot vanish.



Actually, the case when  $S_{ii} + S_{jj}$  is small is only when

$\|L_{\text{sqrt}, M}\|$  is large: either  $\lambda_i, \lambda_j$  are small, or



$\lambda_i, \lambda_j$  close to the negative real line.

`sqrtm` in Matlab uses blockwise Schur-Parlett in a recursive way:

$$M = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$$

with  $A, B$  square blocks  
with almost the same size,

then  $M^{1/2} = \begin{bmatrix} A^{1/2} & Z \\ 0 & B^{1/2} \end{bmatrix}$   $A^{1/2}, B^{1/2}$  computed recursively  
 $Z$  computed from

$$\begin{bmatrix} A^{1/2} & Z \\ 0 & B^{1/2} \end{bmatrix} \begin{bmatrix} A^{1/2} & Z \\ 0 & B^{1/2} \end{bmatrix} = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$$

$$A^{1/2} Z + Z B^{1/2} = C \quad \leftarrow \text{Sylvester equation.}$$

Theorem: Let  $U$  be upper triangular,  $\tilde{S}$  be the matrix computed with the Schur-Parlett variant in machine arithmetic.

$$\tilde{S}^2 = U + \Delta \quad |\Delta_{ij}| = (|S|^2)_{ij} \cdot O(u)$$

(u = machine precision)

Theorem:  $\tilde{X}$  computed with the Schur-Parlett variant, for a general  $M \in \mathbb{C}^{n \times n}$  (not nec. triangular) satisfies

$$\|\tilde{X}^2 - M\|_F \leq \|X\|_F^2 \cdot O(n^3 u)$$

This is not backward stability, which would be

$$\|\tilde{X}^2 - M\|_F \leq \|M\|_F \cdot O(n^3 u)$$

$$\Leftrightarrow \tilde{X}^2 = M + \Delta \quad \frac{\|\Delta\|_F}{\|M\|_F} = O(n^3 u)$$

Instead of having  $\Delta$  small w.r.t.  $\|M\|_F$ , we only know that it is small w.r.t.  $\|\tilde{X}\|_F^2$

$$\|M\| = \|X^2\| \leq \|X\| \cdot \|X\|.$$

$$\tilde{X}^2 = M + \Delta \quad \frac{\|\tilde{X} - X\|}{\|X\|_F} = O(u) \cdot \text{condition number}$$

$$\frac{\|\Delta\|}{\|M\|} = O(u)$$

$$S_{12} = \frac{U_{12}}{S_{22} + S_{11}}$$

$$\tilde{S}_{12} = \frac{U_{12}}{\tilde{S}_{22} + \tilde{S}_{11}} (1 + \delta_1)(1 + \delta_2)$$

$$\tilde{S}_{11} \tilde{S}_{12} + \tilde{S}_{12} \tilde{S}_{22} = U_{12} (1 + \delta_1)(1 + \delta_2) = U_{12} + \boxed{U_{12}(\delta_1 + \delta_2)}$$

$$\|U \tilde{X} - b\| = O(\text{mach. precision.})$$



Wed May 7 → makeup lecture 9-11

Newton method for the square root:

Multivariable Newton method:  $F(x) = x^2 - M$

$$X_{k+1} = X_k - E \quad E = \text{Jac}_{F, X_k}^{-1} [F(X_k)]$$

↑

Fréchet derivative,  $\mathbb{C}^{n \times n} \rightarrow \mathbb{C}^{n \times n}$

The Fréchet derivative of the map  $F(x) = x^2 - M$  in  $X$  is the operator  $E \mapsto XE + EX$

$$F(X+E) = (X+E)^2 - M = \underbrace{X^2 - M}_{F(X)} + \underbrace{EX + XE}_{L_X[E]} + \underbrace{E^2}_{o(\|E\|)}$$

$$L_{X_k}[E] = F'(X_k) = 2X_k E$$

$E$  solves the Sylvester equation  $X_k E + E X_k = X_k^2 - M$

Algorithm:

1. Choose  $X_0$

2. For  $k=0, 1, 2, 3, \dots$

• solve  $X_k E + E X_k = X_k^2 - M$  to find  $E$

• set  $X_{k+1} = X_k - E$

Problem: expensive, one Schur form at each step

However, we can do something better: suppose  $X_0$  and  $M$  commute. Then,  $E = (2X_0)^{-1}(X_0^2 - M)$  solves the equation:

$E$  is a function of  $X_0, M$  and it commutes with  $X_0, M$ , and we have

$$X_0 (2X_0)^{-1} (X_0^2 - M) + (2X_0)^{-1} (X_0^2 - M) X_0 \\ = (2X_0) (2X_0)^{-1} (X_0^2 - M) = X_0^2 - M \quad \checkmark$$

Now,

$$X_1 = X_0 - E = X_0 - (2X_0)^{-1} (X_0^2 - M)$$

is built with  $X_0, M$ , and one can see that it commutes with  $M$ :  $X_1 M = M X_1$

One can prove by induction:

$$\text{Let } M \in \mathbb{C}^{n \times n}, \quad X_0 \text{ s.t. } X_0 M = M X_0$$

- Then,
1.  $E = (2X_k)^{-1} (X_k^2 - M)$  solves the Sylvester equation appearing in the Newton method
  2.  $X_k$  commutes with  $M$ .

So we can simplify the algorithm:

1. Choose  $X_0$  that commutes with  $M$
2. For  $k=0, 1, 2, \dots$

$$E = (2X_k)^{-1} (X_k^2 - M)$$

$$X_{k+1} = X_k - E$$

Or equivalently:

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

Modified Newton iteration:  $= \frac{1}{2} (X_k + X_k^{-1} M)$

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

with  $X_0 = \alpha I$  or  $X_0 = \alpha M$

However:

