

Matrix sign function

Note Title

2023-04-04

Recap: Schur-Parlett method

1) Reduce to diagonal case

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

$$\Lambda(A) \subset \text{LHP}$$

$$\Lambda(B) \subset \text{RHP}$$

2) solve $AZ - ZB = -2C$ (*)

$$3) \text{sign} \left(\begin{bmatrix} A & C \\ 0 & B \end{bmatrix} \right) = \begin{bmatrix} -1 & Z \\ 0 & 1 \end{bmatrix}$$

Conditioning of (*) related to $\text{sep}(A, B)$.

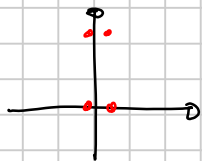
Perturbation theory

From generic results on matrix functions, if $M = VDV^{-1}$

$$\|L_{\text{sign}, M}^{\wedge}\| \leq K(V)^2 \max_{\lambda, \mu \in \Lambda(M)} |f[\lambda, \mu]|$$

$$f[\lambda, \mu] = \begin{cases} f'(\lambda) & \text{if } \lambda = \mu \\ \frac{f(\lambda) - f(\mu)}{\lambda - \mu} & \text{otherwise} \end{cases}$$

For the sign, $f[\lambda, \mu] = \begin{cases} 0 & \lambda = \mu \text{ or in same half-plane} \\ \frac{\pm 2}{\lambda - \mu} & \text{if } \lambda, \mu \text{ on opposite sides of im. axis} \end{cases}$



Theorem: If $M = Q \begin{bmatrix} A & C \\ 0 & B \end{bmatrix} Q^*$, where Q orthogonal,
 $\Lambda(A) \subset \text{LHP}$, $\Lambda(B) \subset \text{RHP}$. Let $\delta = \text{sep}(A, B)$

and consider a perturbation $\tilde{M} = M + E$ of norm $\|E\|_F \leq \varepsilon$ sufficiently small.

Then,

$$1) \quad \|\text{sign}(M)\|_F = O\left(\frac{1}{\delta}\right) \quad \checkmark \quad O\left(\frac{\varepsilon}{\delta^3}\right)$$

$$2) \quad \frac{\|\text{sign}(M+E) - \text{sign}(M)\|_F}{\|\text{sign}(M)\|_F} = O\left(\frac{\varepsilon}{\delta^2}\right) \quad \checkmark \quad O\left(\frac{1}{\delta}\right)$$

3) the Hurwitz stable invariant subspace of $M+E$,
 i.e. its invariant subspace relative to the LHP,
 has a basis of the form $Q \begin{bmatrix} 1 \\ x \end{bmatrix}$, where $\|x\|_F = O\left(\frac{\varepsilon}{\delta}\right)$
 (so it is an $O\left(\frac{\varepsilon}{\delta}\right)$ perturbation of the
 invariant subspace $\text{Im } Q \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ of M)

Proof:

$$1) \quad \text{sign}(M) = Q \begin{bmatrix} -1 & z \\ 0 & 1 \end{bmatrix} Q^* \quad \|\text{sign}(M)\|_F = O(\|z\|_F)$$

$$z \text{ solves } Az - zB = -zC \text{ hence } \|z\|_F \leq \|(I \otimes A - B^T \otimes I)^{-1}\| \cdot \|zC\|_F$$

$$= \frac{1}{\text{sep}(A, B)} \cdot 2\|C\|_F = O\left(\frac{1}{\delta}\right)$$

3) From a result in the beginning of the course:

$$Q^*(M+E)Q = \begin{bmatrix} 1 & 0 \\ x & 1 \end{bmatrix} \begin{bmatrix} \hat{A} & \hat{C} \\ 0 & \hat{B} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -x & 1 \end{bmatrix} \quad \rightarrow \quad \begin{bmatrix} \tilde{A} + \tilde{C}x & \tilde{C} \\ 0 & \tilde{B} - x\tilde{C} \end{bmatrix}$$

for a matrix X with norm $\|X\|_F = O\left(\frac{\varepsilon}{\delta}\right)$

and hence $Q \begin{bmatrix} 1 \\ x \end{bmatrix}$ spans the Hurwitz stable inv. subspace of $M+E$.

(For sufficiently small ε , $\Lambda(\tilde{A} + \tilde{C}x) \subset \text{LHP}$, $\Lambda(\tilde{B} - x\tilde{C}) \subset \text{RHP}$)

$$Q^*(M+E)Q = \begin{bmatrix} \hat{A} & \hat{C} \\ \tilde{D} & \tilde{B} \end{bmatrix} \quad Q^*MQ = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$$

2) Continue using the previous decomposition to compute $\text{sign}(M+E)$:

$$\begin{aligned} \text{sign}(M+E) &= Q \begin{bmatrix} 1 & 0 \\ X & 1 \end{bmatrix} \text{sign} \left(\begin{bmatrix} \hat{A} & \hat{C} \\ 0 & \hat{B} \end{bmatrix} \right) \begin{bmatrix} 1 & 0 \\ -X & 1 \end{bmatrix} Q^* \\ &= Q \begin{bmatrix} 1 & 0 \\ X & 1 \end{bmatrix} \begin{bmatrix} -1 & \hat{Z} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -X & 1 \end{bmatrix} Q^* \end{aligned}$$

where \hat{Z} solves $\hat{A}\hat{Z} - \hat{Z}\hat{B} = -2\hat{C}$

$$\text{i.e. } (\hat{A} + \tilde{C}X)\hat{Z} - \hat{Z}(\tilde{B} - X\tilde{C}) = -2\hat{C}$$

This is a small perturbation of $\text{sign}(M) = Q \begin{bmatrix} -1 & Z \\ 0 & 1 \end{bmatrix} Q^*$

$\|\hat{Z} - Z\|_F = ?$ To bound $\|\hat{Z} - Z\|_F$, we use perturbation theory for linear systems.

$$AZ - ZB = -2C \Leftrightarrow (I \otimes A - B^T \otimes I) \text{vec}(Z) = \text{vec}(-2C)$$

$$(\hat{A} + \tilde{C}X)\hat{Z} - \hat{Z}(\tilde{B} - X\tilde{C}) = -2\hat{C} \Leftrightarrow (I \otimes (\hat{A} - \tilde{C}X) - (\tilde{B} - X\tilde{C})^T \otimes I) \text{vec}(\hat{Z}) = \text{vec}(-2\hat{C})$$

Lemma: if x solves $Tx=c$, and \tilde{x} solves $(T+\delta T)\tilde{x}=(c+\delta c)$
 then $\frac{\|\tilde{x}-x\|}{\|x\|} \leq K(T) \left(\frac{\|\delta c\|}{\|c\|} + \frac{\|\delta T\|}{\|T\|} \right) + O\left(\left(\frac{\|\delta T\|}{\|T\|} \right)^2 \right)$

$$\frac{\|\hat{Z} - Z\|_F}{\|Z\|_F} \leq K(T) \left(\frac{\|2\hat{C} - C\|_F}{\|C\|_F} + \frac{\|\hat{A} - \tilde{C}X - A\|_F + \|\tilde{B} - X\tilde{C} - B\|_F}{\|T\|} \right) + O(\dots)$$

$$O\left(\frac{1}{\delta}\right) \cdot \left(O(\epsilon) + O\left(\frac{\epsilon}{\delta}\right) \right)$$

$$= O\left(\frac{\epsilon}{\delta^2}\right)$$

$$\begin{aligned} \text{sign}(M+E) - \text{sign}(M) &= Q \underbrace{\begin{bmatrix} 1 & 0 \\ x & 1 \end{bmatrix}}_{O\left(\frac{\epsilon}{\delta}\right) \text{ part. of } I} \underbrace{\begin{bmatrix} -1 & \hat{z} \\ 0 & 1 \end{bmatrix}}_{O\left(\frac{\epsilon}{\delta^2}\right) \|z\| \text{ part. of } \begin{bmatrix} -1 & \hat{z} \\ 0 & 1 \end{bmatrix}} \underbrace{\begin{bmatrix} 1 & 0 \\ -x & 1 \end{bmatrix}}_{O\left(\frac{\epsilon}{\delta}\right) \text{ part. of } I} Q^x - Q \begin{bmatrix} 1 & \hat{z} \\ 0 & 1 \end{bmatrix} Q^x \\ &= O\left(\frac{\epsilon}{\delta^2}\right) \|z\|_F = O\left(\frac{\epsilon}{\delta^2}\right) \cdot \|\text{sign}(M)\|_F \quad \square \end{aligned}$$

[Byers, Mehrmann, He 97].

Newton method for the matrix sign

Let us consider the matrix iteration

$$X_0 = M \quad X_{k+1} = \frac{1}{2}(X_k + X_k^{-1}) \quad k=0, 1, 2, \dots$$

We wish to see that $\lim_{k \rightarrow \infty} X_k = \text{sign}(M)$

Let us first study the case $M = V D V^{-1}$ diagonalizable
 \uparrow $\begin{matrix} \text{diag}(\text{eigenvalues}) \\ \text{eigenvectors} \end{matrix}$

$$X_0 = V D V^{-1} \quad X_1 = \frac{1}{2}(X_0 + X_0^{-1}) = V \left(\frac{1}{2}(D + D^{-1}) \right) V^{-1}$$

$$= V \begin{bmatrix} \frac{\lambda_1 + \frac{1}{\lambda_1}}{2} & & & \\ & \frac{\lambda_2 + \frac{1}{\lambda_2}}{2} & & \\ & & \ddots & \\ & & & \frac{\lambda_n + \frac{1}{\lambda_n}}{2} \end{bmatrix} V^{-1} = V \begin{bmatrix} f(\lambda_1) & & & \\ & \ddots & & \\ & & \ddots & \\ & & & f(\lambda_n) \end{bmatrix} V^{-1}$$

set $f(x) = \frac{1}{2}\left(x + \frac{1}{x}\right)$

and analogously $X_k = V \begin{bmatrix} f^k(\lambda_1) \\ \vdots \\ f^k(\lambda_n) \end{bmatrix} V^{-1}$ (by induction).
 f composed k times

$$f(x) = \frac{1}{2} \left(x + \frac{1}{x} \right) = \frac{x^2 + 1}{2x}$$

$f(x)$ is the iterative method obtained by applying Newton's method to search a root of $x^2 - 1 = 0$: indeed,

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

Fixed points of f : $f(a) = a \iff \frac{a^2 + 1}{2a} = a \iff a^2 + 1 = 2a^2 \iff a = \pm 1$

Lemma: the limit of the scalar iteration $x_{k+1} = f(x_k) = \frac{1}{2} \left(x_k + \frac{1}{x_k} \right)$

is

$$\lim_{k \rightarrow \infty} f^k(x_0) = \begin{cases} 1 & \text{if } \operatorname{Re}(x_0) > 0 \\ -1 & \text{if } \operatorname{Re}(x_0) < 0 \\ \text{complicated} & \text{if } \operatorname{Re}(x_0) = 0 \end{cases}$$

Proof Change of variables (Cayley transform)

$$y = \frac{x-1}{x+1}$$

Let us first compute the inverse map:

$$y = \frac{x-1}{x+1} \iff xy + y = x - 1 \iff x(y-1) = -1-y \iff x = \frac{1+y}{1-y}$$

if $y_k = \frac{x_k - 1}{x_k + 1}$ and $x_{k+1} = \frac{1}{2} \left(x_k + \frac{1}{x_k} \right)$, then

$$y_{k+1} = \frac{x_{k+1} - 1}{x_{k+1} + 1} = \frac{\frac{1}{2} \left(x_k + \frac{1}{x_k} \right) - 1}{\frac{1}{2} \left(x_k + \frac{1}{x_k} \right) + 1} = \frac{x_k^2 + 1 - 2x_k}{x_k^2 + 1 + 2x_k} = \left(\frac{x_k - 1}{x_k + 1} \right)^2 = y_k^2$$

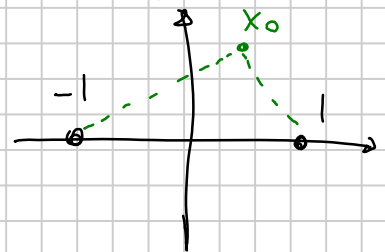
$$\begin{array}{ccccccc}
 x_0 & \xrightarrow{f} & x_1 & \xrightarrow{f} & x_2 & \xrightarrow{f} & x_3 \\
 \downarrow \frac{x-1}{x+1} & & \downarrow \frac{x-1}{x+1} & & \downarrow \frac{x-1}{x+1} & & \downarrow \frac{x-1}{x+1} \\
 y_0 & \xrightarrow{g} & y_1 & \xrightarrow{g} & y_2 & \xrightarrow{g} & y_3
 \end{array}$$

$$g(y) = y^2$$

So the iteration is, in the y variables, repeated squaring.

If $x_0 \in \text{RHP}$, $\text{dist}(x_0, 1) < \text{dist}(x_0, -1)$

$$|y_0| = \left| \frac{x_0 - 1}{x_0 + 1} \right| < 1$$



Then, $y_k = (y_0)^{2^k}$ converges to 0, i.e. $\lim_{k \rightarrow \infty} y_k = 0$

and by the reverse change of variables

$$\lim_{k \rightarrow \infty} x_k = \lim_{k \rightarrow \infty} \frac{1 + y_k}{1 - y_k} = 1$$

If $x_0 \in \text{LHP}$ $\text{dist}(x_0, -1) < \text{dist}(x_0, 1)$

$$|y_0| = \left| \frac{x_0 - 1}{x_0 + 1} \right| > 1$$

Then, $y_k = (y_0)^{2^k}$ converges to ∞ and

$$\lim_{k \rightarrow \infty} x_k = \lim_{k \rightarrow \infty} \frac{1 + y_k}{1 - y_k} = -1$$

$$\hat{U} = \begin{array}{c} \begin{array}{cc} & \begin{array}{c} \text{n_llhp} \\ \downarrow \\ \text{n} \end{array} \\ \begin{array}{c} \downarrow \\ \text{i} \end{array} \\ \begin{array}{|cc|} \hline A & C \\ \hline \text{O} & B \\ \hline \end{array} \end{array}$$