Newton for the matrix sign
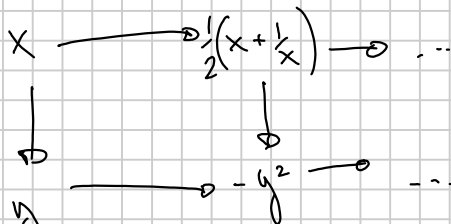
Scalar version:    $X_{k+1} = \frac{1}{2}\left(X_k + X_k^{-1}\right)$     $X_0 = \lambda$

**Theorem:**    $\lim\limits_{k \to \infty} X_k = \text{sign}(\lambda)$

(and the convergence is quadratic).

$y = \dfrac{1+x}{1-x}$                        $(\mathbb{C} \cup \{\infty\}) \longrightarrow (\mathbb{C} \cup \{\infty\})$
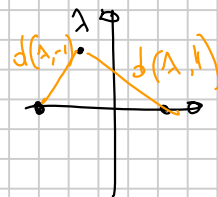
$y_k := \dfrac{1+X_k}{1-X_k}$     $y_{k+1} = \dfrac{1+X_{k+1}}{1-X_{k+1}} = \dfrac{1+\frac{1}{2}\left(X_k + X_k^{-1}\right)}{1-\frac{1}{2}\left(X_k + X_k^{-1}\right)} = \dfrac{X_k^2 + 2X_k + 1}{-X_k^2 + 2X_k - 1}$

$$= -\left(\dfrac{1+X_k}{1-X_k}\right)^2 = -y_k^2$$

$X \longrightarrow \frac{1}{2}\left(x + \frac{1}{x}\right) \longrightarrow \cdots$

$\downarrow$               $\downarrow$

$y \longrightarrow -y^2 \longrightarrow \cdots$

$X_0 = \lambda$     $y_0 = \dfrac{1+\lambda}{1-\lambda}$     $|y_0| > 1$   if   $\text{Re}(\lambda) > 0$

$|y_0| < 1$   if   $\text{Re}(\lambda) < 0$

$|y_0| = \dfrac{|1+\lambda|}{|1-\lambda|} = \dfrac{d(\lambda, -1)}{d(\lambda, 1)}$



$y_1 = -y_0^2$, $y_2 = -y_0^4$, $y_3 = -y_0^8$, $\cdots$

$$|y_k| = \left|y_0^{2^k}\right| \longrightarrow \begin{cases} 0 & \text{if } |y_0| < 1 \\ \infty & \text{if } |y_0| > 1 \end{cases}$$

so   $y_k \longrightarrow \begin{cases} 0 \\ \infty \end{cases}$   quadratically   if   $0 = y = \dfrac{1+x}{1-x} \Leftrightarrow x = -1$

if $\infty = y = \dfrac{1+x}{1-x} \Leftrightarrow x = \infty$

$X_k \longrightarrow \begin{cases} -1 \\ 1 \end{cases}$   quadratically

Behavior of the corresponding matrix iteration

$$X_{k+1} = \frac{1}{2}\left(X_k + X_k^{-1}\right) \qquad X_0 = A$$

$$Y_k = (X_k - S)(X_k + S)^{-1} \qquad S = \text{sign}(A)$$

$$\text{if } A = VJV^{-1}$$

$$Y_0 = (A-S)(A+S)^{-1} = V\left(J - \text{sign}(J)\right)V^{-1}V\left(J + \text{sign}(J)\right)^{-1}V^{-1}$$

$$\underbrace{\phantom{V\left(J + \text{sign}(J)\right)^{-1}V^{-1}}}$$

on diagonal:

$$\lambda_j + \text{sign}(\lambda_j) \neq 0$$

$$= V\begin{bmatrix} \ddots & & \overbrace{\lambda_j^{(1)}}^{\phantom{x}} & \times \\ & (\lambda_j - \text{sign}(\lambda_j))(\lambda_j + \text{sign}(\lambda_j))^{-1} & & \\ & 0 & & \ddots & \times \\ \end{bmatrix} V^{-1}$$

$$\lambda_j = a + bi \qquad \lambda_j + \text{sign}(\lambda_j) =$$
$$(a + \text{sign}(a)) + bi$$
$$\lambda_j - \text{sign}(\lambda_j) = (a - \text{sign}(a)) + bi$$
$$\Rightarrow \text{all eigvls of } Y_0 \text{ are } |\lambda_j| < 1$$

$$Y_{k+1} = (X_{k+1} - S)(X_{k+1} + S)^{-1} = \left(\tfrac{1}{2}(X_k + X_k^{-1}) - S\right)\left(\tfrac{1}{2}(X_k + X_k^{-1}) + S\right)^{-1} =$$

$$= \left(\tfrac{1}{2}(X_k + X_k^{-1}) - S\right)(2X_k)(2X_k)^{-1}\left(\tfrac{1}{2}(X_k + X_k^{-1}) + S\right)^{-1}$$

$$= \left(X_k^2 + I - 2SX_k\right)\left[\left(\tfrac{1}{2}(X_k + X_k^{-1}) + S\right)(2X_k)\right]^{-1}$$

$$= \left(X_k^2 + I - 2SX_k\right)\left(X_k^2 + I + 2SX_k\right)^{-1}$$

$$\underset{\underset{S^2}{\shortparallel}}{\phantom{X_k^2}} \quad \underset{\underset{SX_k = X_k S}{\shortparallel}}{\phantom{2SX_k}} \quad \underset{\underset{S^2}{\shortparallel}}{\phantom{X_k^2}}$$

$$= (X_k - S)^2(X_k + S)^{-2} = \left[(X_k - S)(X_k + S)^{-1}\right]^2 = Y_k^2$$

$$Y_k = Y_0^{2^k} \to 0 \quad \text{because } \rho(Y_0) < 1$$

$$Y_k = (X_k - S)(X_k + S)^{-1} \iff Y_k(X_k + S) = X_k - S$$

$$= (Y_k - I)X_k = (-I - Y_k)S \quad \Rightarrow X_k = (Y_k - I)^{-1}(-I - Y_k)S$$

$$\to S \quad \text{if } Y_k \to 0$$

Algorithm:

$$X_0 = A$$

$$\Rightarrow \text{you } \underline{\text{do}} \text{ need an inverse!}$$
$$\text{inv}(X$$

Repeat $\quad X_{k+1} = \tfrac{1}{2}(X_k + X_k^{-1}).$

Convergence slow when for from $\pm 1$



How to speed up the iteration when this happens? Scaling!

$$sign(X) = sign(\alpha X) \quad \text{for each} \quad \alpha > 0$$

Can we choose $\alpha$ to make sure that all eigenvalues of $A$ have absolute value $\tilde{\sim} 1$

Might be impossible if the eigls of $A$ have different scales, e.g. $eig(A) = 10^5, 2, 10^{-6}(i+1)$

But still, it is better to have $10^{-5}, 2, 10^{-6}(i+1)$

$$\text{than} \quad i+1, 2\cdot 10^6, 10^{11} \quad (\alpha = 10^6)$$

(Q: what happens if you start the iteration $X_{k+1} = \frac{1}{2}(X_k + X_k^{-1})$ from a point on the imaginary axis?) e.g. $X_0 = i$

Idea: we want to balance the orders of magnitudes of eigenvalues so that they are "centered in 1"

($10^5$ is as slow-converging as $10^{-5}$)

Idea 1: make geometric mean of eigenvalues $= 1$ :

Cheaper! $(\lambda_1 \lambda_2 \cdots \lambda_n)^{1/n} = (\det A)^{1/n} = 1 \quad \Longleftrightarrow \quad \det A = 1$

$X^{-1} = U^{-1} L^{-1}$

gives "determinantal scaling".

$\det(..) = prod(diag(U))$

Idea 2: make $[\lambda_{min}, \lambda_{max}]$ "centered around 1"

So you try to make $\alpha \lambda_{min} = (\alpha \lambda_{max})^{-1}$

need some steps of power iteration on $A$, $A^{-1}$

$$\alpha \lambda_{min} = \alpha^{-1} \lambda_{max}^{-1}$$

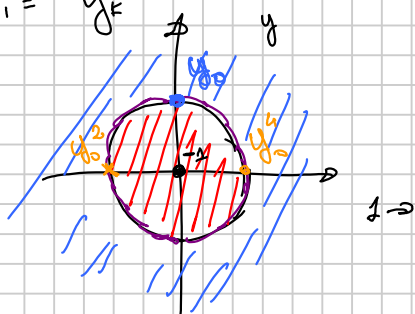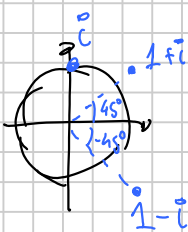$$\alpha^2 = \frac{1}{\lambda_{min} \lambda_{max}}$$
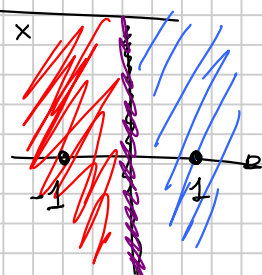
"Spectral scaling"

Idea 3: $\sigma_{min}(\alpha A)\,\sigma_{max}(\alpha A) = 1$    "norm scaling"

need some steps of
power it. on $A^TA$
and $A^{-T}A^{-1}$

$\|\alpha A\|$

---

$X_0 = i$     $X_{k+1} = \frac{1}{2}\left(X_k + X_k^{-1}\right)$     ?

$y_0 = \frac{4x_0}{1-x_0} = \frac{1+i}{1-i} = i$     $y_{k+1} = -y_k^2$

$\cdot \, 1+i$

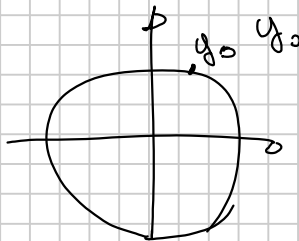$-45°$
$(-45°)$

$1-i$

$y_0 = i$

$y_1 = -y_0^2 = -(-1) = 1$

$y_2 = -1^2 = -1$

$y_3 = -1$ ,     $y_4 = -1, \ldots$

$X_0 = 1.1 i$

---

Stability: complicated. Under small perturbations,

Assume (up change of basis)    $A = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix}$    $\Lambda(A_{11}) \subseteq$ LHP
$\Lambda(A_{22}) \subseteq$ RHP

$\text{sign}(A) = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$     $V_1 = \begin{bmatrix} I \\ 0 \end{bmatrix}$

$\left\| \text{sign}(A+E) - \text{sign}(A) \right\| \lesssim \frac{\|E\|}{\left(\text{sep}(A_{11}, A_{22})\right)^3}$

(Scales as separation$^3$)

However, often one is only interested in $V_1$, $V_2$ stable/antistable
invariant subspaces      $\left( \text{sign}(A) = [V_1 | V_2] \begin{bmatrix} -I & 0 \\ 0 & I \end{bmatrix} [V_1|V_2]^{-1} \right)$

If one uses a computed $\text{sign}(A+E)$ to extract these two

subspaces,
$$\text{Ker}\left(\text{sign}(A+E)+I\right) = \begin{bmatrix} 1 \\ X \end{bmatrix}$$

$$\|X\| \leq \frac{\|E\|}{\text{sep}(A_{11}, A_{22})} \quad \underline{\text{without}} \text{ the cube!}$$

[Bai, Demmel '98, Byers, Mehrmann, He '97]

---

Suppose you are given $M, N$ s.t. $A = M^{-1}N$

Can you compute $\text{sign}(A)$ just from $M, N$ without forming $A$?

(Natural question if you think about generalized eigenvalues / QZ algorithm)

QZ is, essentially, a method to compute eigenvalues of $A = M^{-1}N$
just from $M, N$ without forming $A$)

We can implement the Newton iteration in a similar "inverse-free" fashion

input: $M, N$ s.t. $A = M^{-1}N = X_0$

$$X_1 = \frac{1}{2}\left(A + A^{-1}\right) \overset{?}{=} M_1^{-1}N_1$$

$$\frac{1}{2}\left(M^{-1}N + N^{-1}M^{-1}\cdot\right) = \frac{1}{2}M^{-1}\left(N + M\hat{N}^{-1}M\right)$$

Suppose we find
$\hat{M}, \hat{N}$
$MN^{-1} = \hat{M}^{-1}\hat{N}$

$$= \frac{1}{2}M^{-1}\left(N + \hat{M}^{-1}\hat{N}M\right) = \frac{1}{2}M^{-1}\hat{M}^{-1}\left(\hat{M}N + \hat{N}M\right)$$

$$= \underbrace{\left(\hat{M}M\right)^{-1}}_{M_1^{-1}}\underbrace{\left(\frac{1}{2}\left(\hat{M}N + \hat{N}M\right)\right)}_{N_1}$$

$(M_0, N_0) \longrightarrow (M_1, N_1) \longrightarrow (M_2, N_2)$
$\longrightarrow (M_3, N_3) \longrightarrow \cdots \longrightarrow (M_\infty, N_\infty)$

How to find
$MN^{-1} = \hat{M}^{-1}\hat{N}$ ?

$$\iff \hat{M}M = \hat{N}N \iff \begin{bmatrix} \hat{M} & \hat{N} \end{bmatrix}\begin{bmatrix} M \\ -N \end{bmatrix} = 0$$

$$\begin{bmatrix} \hat{M} & \hat{N} \end{bmatrix}^T = \text{Ker}\left(\begin{bmatrix} M \\ -N \end{bmatrix}^T\right) \quad \text{(computed e.g. with a QR decomposition)}$$

There are cases in which computing this kernel is a lot

more stable than computing $M^{-1}$ or $N^{-1}$:

for instance,

$$\begin{bmatrix} M \\ -N \end{bmatrix} = Q \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \\ \varepsilon & 0 \\ 0 & 1 \end{bmatrix} \begin{matrix} Q^T \\ Q^T \end{matrix} \qquad \text{for a small } \varepsilon > 0$$

$$\begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \\ \varepsilon & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \cdot \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

$$Q_1, Q_2 \in \mathbb{R}^{4\times 2} \qquad R_1 \in \mathbb{R}^{2\times 2}$$

$$Q_2^T Q_1 = 0 \implies Q_2^T = \begin{bmatrix} \hat{M} & \hat{N} \end{bmatrix}$$

("linear algebra on pencils")    matrix pencil $= \left( -M\lambda + N \right)$

---

$$x_{k+1} = \frac{1}{2}\left( x_k + x_k^{-1} \right) = \frac{x_k^2 + 1}{2 x_k} \qquad \leftarrow \text{"even terms" of } (1+x_k)^2$$
$$\qquad \leftarrow \text{"odd terms" of } (1+x_k)^2$$

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

The same works for higher powers:

$$x_{k+1} = \frac{(1+x_k)^3 + (1-x_k)^3}{(1+x_k)^3 - (1-x_k)^3} = \frac{1+3x_k+3x_k^2+x_k^3 + \left(1-3x_k+3x_k^2-x_k^3\right)}{1+3x_k+3x_k^2+x_k^3 - \left(1-3x_k+3x_k^2-x_k^3\right)}$$

$$= \frac{2\left(1+3x_k^2\right)}{2\left(3x_k+x_k^3\right)} \qquad \begin{matrix} \text{has fixed points in } \pm 1 \\ \text{of order 3} \end{matrix}$$

Problem: it is no longer true that $x_k \to \begin{cases} -1 & \text{if } \operatorname{Re} x_0 < 0, \\ 1 & \text{if } \operatorname{Re} x_0 > 0 \end{cases}$