

Methods for $A^{1/2}$

Note Title

2023-04-18

Recap: Schur-Parlett

- 1) $A = QUQ^*$ to reduce to triangular U
- 2) Compute diagonal of $S = f(U)$ $S_{ii} = f(U_{ii})$
- 3) Looking at entry (i, j) of the identity $US = SU$, we obtain formulas $S_{ij} = \frac{\dots}{U_{ii} - U_{jj}}$ which we can use to compute entries S_{ij} one by one

$$4) f(A) = QSQ^*$$

For $f(x) = x^{1/2}$, variant to avoid the denominator:

use $S^2 = U$ instead of $SU = US$

$$U_{ij} = S_{ii} S_{ij} + S_{i,i+1} S_{i+1,j} + S_{i,i+2} S_{i+2,j} + \dots + S_{ij} S_{jj} \quad i < j$$

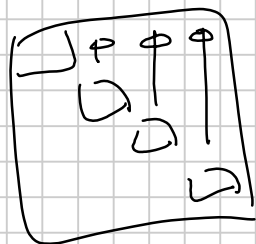
$$S_{ij} = \frac{U_{ij} - S_{i,i+1} S_{i+1,j} - \dots - S_{i,j-1} S_{j-1,j}}{S_{ii} + S_{jj}}$$

g entries below and to the left of S_{ij}

denominator

$$S_{ii} + S_{jj} = U_{ii}^{1/2} + U_{jj}^{1/2} \text{ cannot be } 0 : U_{ii}, U_{jj} \in \mathbb{R}^+ \Rightarrow U_{ii} + U_{jj} \in \mathbb{R}^+$$

Moreover, $\frac{1}{S_{ii} + S_{jj}} = f[U_{ii}, U_{jj}] = \frac{U_{ii}^{1/2} - U_{jj}^{1/2}}{U_{ii} - U_{jj}}$ is large only when $f(A)$ is ill-cond.
(Matlab implementation)



```

function B = sqrtm_schurparlett(A)
[Q, U] = schur(A, 'complex');
n = size(A, 1);
S = zeros(n, n);
for j = 1:n
    S(j,j) = sqrt(U(j,j));
    for i = j-1:-1:1
        num = U(i,j) - S(i,i+1:j-1)*S(i+1:j-1,j);
        S(i,j) = num / (S(i,i)+S(j,j));
    end
end
B = Q*S*Q';
if isreal(A)
    B = real(B);
end

```

$\frac{\|\tilde{B}^2 - A\|}{\|A\|} = O(u) \iff$ the computed \tilde{B} is the exact square root of $A + \delta_A$ $\frac{\|\delta_A\|}{\|A\|} = O(u)$
 i.e., the method is backward stable.

This method is unfortunately not backward stable, but we can prove that

$$\frac{\|\tilde{B}^2 - A\|}{\|\tilde{B}\|^2} = O(u)$$

(Note that $\|A\| = \|\tilde{B}^2\| = \|\tilde{B} \cdot \tilde{B}\| \leq \|\tilde{B}\| \cdot \|\tilde{B}\|$, and there can be major cancellation for non-normal matrices)

Theorem: suppose $\tilde{S} = f(u)$ is computed in floating point with machine precision u , then

$$\tilde{S}^2 = U + \delta_U \quad |\delta_U| \leq |S|^2 O(n)u$$

$$|\delta_{UVij}| \leq \sum_{k=i}^j |S_{ik}| \cdot |S_{kj}|, \text{ elementwise abs. value and inequalities.}$$

From this theorem, with a little more work we obtain a norm bound $\|\tilde{B}^2 - A\| \leq O(u) \|\tilde{B}\|^2$.

Proof: Note first that in the rhs we can freely replace S with \tilde{S} , because $S - \tilde{S} = O(u)$

$$\tilde{S}_{ij} = \left(U_{ij} \ominus \tilde{S}_{i,i+1} \otimes \tilde{S}_{i+1,j} \ominus \dots \ominus \tilde{S}_{i,j-1} \otimes \tilde{S}_{j-1,j} \right) \oslash (\tilde{S}_{ii} \oplus \tilde{S}_{jj})$$

Replace $a \otimes b$ with $(a \otimes b)(1 + \epsilon_k)$ $|\epsilon_k| \leq u$
 + ignore higher-order terms

$$\begin{aligned} & |U_{ij} - \tilde{S}_{i,i+1} \tilde{S}_{i+1,j} - \dots - \tilde{S}_{i,j-1} \tilde{S}_{j-1,j} - \tilde{S}_{ij}(\tilde{S}_{ii} + \tilde{S}_{jj})| \\ &= |\text{error terms}| \leq u \left(|U_{ij}| + |\tilde{S}_{i,i+1} \tilde{S}_{i+1,j}| + \dots + |\tilde{S}_{i,j-1} \tilde{S}_{j-1,j}| + |\tilde{S}_{ij}| (|\tilde{S}_{ii}| + |\tilde{S}_{jj}|) \right) \end{aligned}$$

e.g. with only one term:

$$\begin{aligned} \tilde{S}_{ij} &= \left(U_{ij} \ominus \tilde{S}_{i,i+1} \tilde{S}_{i+1,j} \right) \oslash (\tilde{S}_{ii} \oplus \tilde{S}_{jj}) \\ &= \left(U_{ij} - \tilde{S}_{i,i+1} \tilde{S}_{i+1,j} \right) (1 + \epsilon_1)(1 + \epsilon_3) / \left((\tilde{S}_{ii} + \tilde{S}_{jj})(1 + \epsilon_2) \right) \end{aligned}$$

$$U_{ij}(1 + \epsilon_1)(1 + \epsilon_3) - (\tilde{S}_{i,i+1} \tilde{S}_{i+1,j})(1 + \epsilon_1)(1 + \epsilon_3) - \tilde{S}_{ij}(\tilde{S}_{ii} + \tilde{S}_{jj})(1 + \epsilon_2) = 0$$

$$\begin{aligned} U_{ij} - \tilde{S}_{i,i+1} \tilde{S}_{i+1,j} - \tilde{S}_{ij} \tilde{S}_{ii} - \tilde{S}_{ij} \tilde{S}_{jj} &= -U_{ij}(\epsilon_1 + \epsilon_3) + (\tilde{S}_{i,i+1} \tilde{S}_{i+1,j})(\epsilon_1 + \epsilon_3) \\ &\quad + \tilde{S}_{ij}(\tilde{S}_{ii} + \tilde{S}_{jj})\epsilon_2 + O(u^2) \end{aligned}$$

$$|U - \tilde{S}^2|_{ij} \leq 2u \left(|U_{ij}| + |\tilde{S}_{i,i+1} \tilde{S}_{i+1,j}| + |\tilde{S}_{ij} \tilde{S}_{ii}| + |\tilde{S}_{ij} \tilde{S}_{jj}| \right)$$

$(|U| + |\tilde{S}|^2)_{ij}$

$$|U - \tilde{S}^2|_{ij} \leq (|U|_{ij} + |\tilde{S}|^2_{ij}) O(n)U$$

$$|U|_{ij} = |S_{ii}S_{ij} + S_{i,i+1}S_{i+1,j} + \dots + S_{ij}S_{jj}|$$

$$\leq |S_{ii}| |S_{ij}| + |S_{i,i+1}| |S_{i+1,j}| + \dots + |S_{ij}| |S_{jj}| = |S|^2_{ij}$$

then $|U - \tilde{S}^2|_{ij} \leq |S|^2_{ij} O(n)U$

Matlab's `sqrtm()` uses a similar algorithm:
 a recursive block version of this: split U in blocks

$$U = \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix} \text{ of equal } (\pm 1) \text{ size}$$

then compute $S_{11} = (U_{11})^{\frac{1}{2}}$ $S_{22} = (U_{22})^{\frac{1}{2}}$ recursively and
 S_{12} by solving a Sylvester equation that comes from

$$\begin{bmatrix} S_{11} & S_{12} \\ 0 & S_{22} \end{bmatrix} \begin{bmatrix} S_{11} & S_{12} \\ 0 & S_{22} \end{bmatrix} = \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix}$$

$$S_{11}S_{12} + S_{12}S_{22} = U_{12} \quad \text{a Sylvester equation,}$$

Another method: matrix iterations.

$B = A^{\frac{1}{2}}$ solves $A - B^2 = 0$, i.e., it is a root of ~~$f(x) = A - x^2$~~ .

$$f(x) = x^2 - A$$

Multivariate Newton method: $F: \mathbb{C}^m \rightarrow \mathbb{C}^m$

$$\begin{cases} x_0 \in \mathbb{C}^m \text{ given} \\ x_{k+1} = x_k - (\text{Jac} F_{x_k})^{-1} F(x_k) \end{cases} \quad | \quad x_{k+1} = x_k - h \text{ where } h \text{ solves}$$

$$\begin{matrix} \uparrow & \uparrow \\ m \times m & m \times 1 \end{matrix} \quad |$$

$$(\text{Jac } F_{x_k})h = F(x_k)$$

We can apply this method working on the "vectorized form" of f .

$$m = n^2 \quad F(x) = \text{vec} \left((\text{unvec}(x))^2 - A \right)$$

Jac F_x is the Kronecker form of the Fréchet derivative

$$L_{f,x}[E] = XE + EX$$

So re-seen as a matrix iteration the Newton method becomes

$$\begin{cases} X_0 \in \mathbb{C}^{n \times n} \text{ given} \\ X_{k+1} = X_k - H \text{ where } H \text{ solves } X_k H + H X_k = X_k^2 - A \end{cases}$$

Multivariate Newton method: Starting from given X_0 :

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

1. Solve $X_k H + H X_k = X_k^2 - A$

2. $X_{k+1} = X_k - H$

We expect this to converge quadratically, as Newton methods do.

Problem: to solve the Sylvester equation we need a Schur form of X_k at each step.

However, there is a simpler way to avoid the Sylvester solution

Lemma: If you start from X_0 s.t. $A X_0 = X_0 A$ (e.g. $X_0 = \alpha \cdot I$, for $\alpha > 0$, or $X_0 = \alpha A$), then

1. At each step, X_k and A commute

2. At each step, $H = (2X_k)^{-1} (X_k^2 - A)$

Proof: Combined induction.

For $k=0$, 1. is obvious

Once proved 1. for a certain k , we can prove 2.:

$$\begin{aligned}X_k H + H X_k &= X_k (2X_k)^{-1} (X_k^2 - A) + (2X_k)^{-1} (X_k^2 - A) X_k \\ &= \left(X_k (2X_k)^{-1} + (2X_k)^{-1} X_k \right) (X_k^2 - A) \\ &= \left(\frac{1}{2} I + \frac{1}{2} I \right) (X_k^2 - A) = X_k^2 - A\end{aligned}$$

H, being a simple expression in X_k and A , commutes, too.

$\Rightarrow X_{k+1} = X_k - H$ commutes with A , too.

This gives a cheaper version of the method

$$\begin{pmatrix} \text{Schur} & \sim 30n^3 \\ \text{Inversion} & \sim 2n^3 \end{pmatrix}$$

So we can modify the Newton iteration to

$$\begin{aligned}X_{k+1} &= X_k - (2X_k)^{-1} (X_k^2 - A) = X_k - \frac{1}{2} (X_k^{-1} X_k^2 - X_k^{-1} A) \\ &= \frac{1}{2} (X_k + X_k^{-1} A)\end{aligned}$$

Modified Newton method:

Given X_0 that commutes with A :

for $k=0, 1, 2, \dots$

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

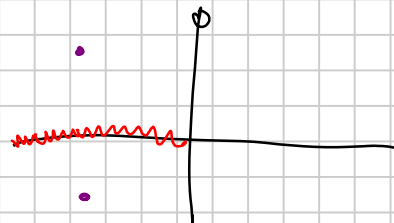
end

MN computes the same iterates as the "true Newton" method, at least in exact arithmetic.

These iterates converge quadratically to a solution of $X^2 - A = 0$.

Theorem: Assume A has no real non-positive eigenvalues, and that the iteration is started from $X_0 = \alpha I$ or $X_0 = \alpha A$ with $\alpha > 0$.

Then, the iterations converge to the principal square root $A^{1/2}$



Proof: set $Z_k := A^{-1/2} X_k$. Then,

$$\begin{aligned} Z_{k+1} &= A^{-1/2} X_{k+1} = A^{-1/2} \frac{1}{2} (X_k + X_k^{-1} A) = \frac{1}{2} (A^{-1/2} X_k + X_k^{-1} A^{1/2}) \\ &= \frac{1}{2} (Z_k + Z_k^{-1}) \end{aligned}$$

This is the Newton iteration for the matrix sign!

Z_k converges quadratically to $\text{sign}(Z_0)$

$$Z_0 = \begin{cases} A^{-1/2} \cdot \alpha I = \alpha A^{-1/2} \\ A^{-1/2} \alpha A = \alpha A^{1/2} \end{cases} \begin{matrix} \swarrow \\ \searrow \end{matrix} \begin{matrix} \text{have all eigenvalues in RHP} \\ \text{since } \Lambda(A^{1/2}) \subset \text{RHP} \end{matrix}$$

Hence, $\text{sign}(Z_0) = I$

So $Z_k \rightarrow I$ and $X_k \rightarrow A^{1/2}$ when $k \rightarrow \infty$.

However, in practice their behavior is very different! $\triangle!$
We will see it in the following lecture

```

rng(0);
M = randn(10); M = M*M';

X = eye(size(M));
Y = eye(size(M));
for k = 1:15
    X = X - lyap(X, X, M-X^2); % True Newton
    Y = 1/2 * (Y + Y\M);
    TNres(k) = norm(X^2-M) / norm(M);
    MNres(k) = norm(Y^2-M) / norm(M);
    TNcom(k) = norm(M*X-X*M) / (norm(M)*norm(X));
    MNcom(k) = norm(M*Y-Y*M) / (norm(M)*norm(Y));
end
[(1:15)' TNres' MNres' TNcom' MNcom']

```

```

>> sqrt_comparison
ans =
 1.0000e+00  1.2716e+01  1.2716e+01  5.0496e-17  5.0496e-17
 2.0000e+00  2.9472e+00  2.9472e+00  2.7535e-16  1.6357e-16
 3.0000e+00  5.5013e-01  5.5013e-01  8.5127e-16  5.6933e-16
 4.0000e+00  4.8810e-02  4.8810e-02  4.6606e-16  3.8247e-15
 5.0000e+00  5.6788e-04  5.6788e-04  9.7668e-17  2.9243e-14
 6.0000e+00  8.0577e-08  8.0577e-08  7.0494e-17  2.2360e-13
 7.0000e+00  1.6465e-15  1.8219e-12  8.2882e-17  1.7098e-12
 8.0000e+00  1.0450e-16  1.3934e-11  9.4525e-17  1.3078e-11
 9.0000e+00  1.1462e-16  1.0658e-10  5.7112e-17  1.0003e-10
 1.0000e+01  8.5601e-17  8.1522e-10  1.1203e-16  7.6519e-10
 1.1000e+01  9.6578e-17  6.2358e-09  8.2755e-17  5.8532e-09
 1.2000e+01  1.0385e-16  4.7700e-08  6.9193e-17  4.4773e-08
 1.3000e+01  9.9863e-17  3.6488e-07  1.0014e-16  3.4249e-07
 1.4000e+01  1.0058e-16  2.7911e-06  7.0488e-17  2.6198e-06
 1.5000e+01  1.0279e-16  2.1350e-05  5.8032e-17  2.0040e-05

```

TN works well, but MN stops converging at residual $\approx 10^{-12}$!