

Let A be a $n \times n$ matrix and \mathcal{L} some space of matrices of the same order, but of lower complexity. This means, for instance in case A is non singular, that solving linear systems $X\mathbf{z} = \mathbf{b}$, $X \in \mathcal{L}$, is much easier than solving systems $A\mathbf{z} = \mathbf{b}$.

We associate to A an element \mathcal{A} in \mathcal{L} such that for any $\varepsilon > 0$ there exist n_ε , r_ε and a splitting of A of type $A = \mathcal{A} + R + E$, where the matrices R and E satisfy the conditions $\text{rank}(R) < r_\varepsilon$, $\|E\| < \varepsilon$ for all $n > n_\varepsilon$ (PROPERTY).

If at least one such element \mathcal{A} exists, then we have the problem of determining the best \mathcal{A} , or an as good as possible \mathcal{A} .

In other words we have an approximation problem (of A) with matrices of the form $\mathcal{A} + R$, $\mathcal{A} \in \mathcal{L}$, $\text{rank}(R)$ small (or of the form $D + R$, D diagonal, $\text{rank}(R)$ small, if $\mathcal{L} = \{UDU^{-1} : D_{ij} = 0, i \neq j\}$ with U non singular). See [OT],[ZOT] (one needs to know SVD to read [OT]) where also a *black dot* algorithm is proposed to solve it. (Possible research: black dot to the coefficient matrix of google system? extend black dot from circulants to more general \mathcal{L} ? f.i. to spaces in pages 3–12 of [CDDFFZ] ? look for alternatives to black dot algorithm, ok also for \mathcal{L} ! see my e-mail).

Once a good \mathcal{A} has been determined, a system $A\mathbf{x} = \mathbf{b}$ can be efficiently solved by applying iterative methods to a *preconditioned* version of the same system, obtained using \mathcal{A} as preconditioner. (This by an obvious generalization of the Theorem(clusterTC) stated below). In fact, as a consequence of PROPERTY the eigenvalues of $\mathcal{A}^{-1}A$ cluster around 1 (check it!). (Note that one should also have that the condition number of the new system is bounded uniformly on n). (Note also that if \mathcal{A} is singular, as we will see it may happen, then is introduced a matrix $\hat{\mathcal{A}}$ with the same eigenvalues of \mathcal{A} except the null ones which are replaced by 1 and \mathcal{A} is set equal to $\hat{\mathcal{A}}$ [OT]).

Let us see an EXAMPLE: $A = \text{real symmetric Toeplitz}$, $\mathcal{L} = \text{Circulants}$

Set $A = [t_{|i-j|}]_{i,j=1}^n$, where t_k , $k = 0, \dots$, are real parameters, and $\mathcal{L} = \{C(\mathbf{z}) : \mathbf{z} \in \mathbb{C}^n\}$, where $C(\mathbf{z})$ is the circulant matrix with first row \mathbf{z}^T .

Observe that if A is non singular, then a system $A\mathbf{x} = \mathbf{b}$ can be solved via direct methods with $O(n^2)$ a.o. (arithmetic operations) [L,Tr], which reduce to $O(n(\log n)^2)$ if A is p.d. [AG] (a further reduction for particular A is obtained in [Dick]).

Solving a circulant system $C(\mathbf{z})\mathbf{x} = \mathbf{b}$ is much cheaper. Via the known representation of $C(\mathbf{z})$

$$C(\mathbf{z}) = \sum_{k=1}^n z_k P^{k-1} = \sqrt{n} F d(F\mathbf{z}) F^*, \quad P = \begin{bmatrix} 0 & 1 & & \\ & & \ddots & \\ & & & 1 \\ 1 & & & \end{bmatrix},$$

($d(\mathbf{v}) = \text{diag}(z_i)$) in terms of the Fourier matrix $F = [\omega^{(i-1)(j-1)}]_{i,j=1}^n / \sqrt{n}$, $\omega = e^{-i2\pi/n}$, two (or three) FFT are enough.

Circulants make iterative CG/GMRES methods an alternative to direct methods, in solving a Toeplitz system. In fact, each step of CG, for example, requires a matrix-vector multiplication $A \cdot \mathbf{v}$, which can be computed via

order- $2n$ FFTs through the formula

$$\begin{bmatrix} A & S \\ S & A \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} A\mathbf{v} \\ S\mathbf{v} \end{bmatrix}, \quad S = \begin{bmatrix} 0 & t_{n-1} & \cdots & t_1 \\ t_{n-1} & & & \\ \vdots & & & \\ t_1 & & & \end{bmatrix}$$

(compute it via order- $(2n+l)$ FFTs !). Moreover, the number of steps required by such methods to converge often reduces drastically if suitable circulants are used as preconditioners. In these cases iterative outperform direct methods.

T.CHAN. Set $\mathcal{A} = C(\mathbf{z}_A^{(TC)})$ where, if $\mathbf{z}^{TC} = \mathbf{z}_A^{(TC)}$, then

$$z_{i+1}^{TC} = \frac{1}{n}((n-i)t_i + it_{n-i}), \quad i = 0, \dots, n-1$$

($t_n = 0$). Observe that $\|A - \mathcal{A}\|_F = \min_{\mathbf{z}} \|A - C(\mathbf{z})\|_F$ [TC]. Moreover, \mathcal{A} is real symmetric, like A .

The proof of the following result shows that under suitable assumptions on the real sequence $\{t_k\}$, which make the matrix A p.d., such \mathcal{A} (besides being p.d. as A) has exactly the required PROPERTY, i.e. for any $\varepsilon > 0$ there exist $n_\varepsilon, r_\varepsilon$ and a splitting of A of type $A = \mathcal{A} + R + E$, where the matrices R and E satisfy the conditions $\text{rank}(R) < r_\varepsilon, \|E\| < \varepsilon$ for all $n > n_\varepsilon$. As a consequence the eigenvalues of $\mathcal{A}^{-1}A$ cluster around 1. On the contrary the eigenvalues of A are equally distributed (under the same assumption on t_k).

Theorem(clusterTC)

Let $t_k, k = 0, 1, \dots$, be real numbers such that $\sum_{k=0}^{+\infty} |t_k| < +\infty$. Set

$$t(\theta) = \sum_{k \in \mathbb{Z}} t_{|k|} e^{ik\theta} = t_0 + 2 \sum_{k=1}^{+\infty} t_k \cos(k\theta).$$

Note that the function $t(\theta)$ is a real, even, continuous (absolutely convergent sum of continuous functions), 2π -periodic, and

$$t_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} t(\theta) e^{-ik\theta} d\theta.$$

(In fact,

$$t(\theta) e^{-ij\theta} = \sum_{k \in \mathbb{Z}} t_{|k|} e^{i(k-j)\theta} \Rightarrow \int_{-\pi}^{\pi} t(\theta) e^{-ij\theta} d\theta = \sum_{k \in \mathbb{Z}} t_{|k|} \int_{-\pi}^{\pi} e^{i(k-j)\theta} d\theta \dots)$$

Set $t_{\min} = \min t(\theta), t_{\max} = \max t(\theta)$. Consider the real symmetric Toeplitz matrices $A = [t_{|i-j|}]_{i,j=1}^n, n = 1, 2, \dots$, and the corresponding circulant matrices $\mathcal{A} = C(\mathbf{z}_A^{(TC)})$. Then

i) $\lambda(A) \in [t_{\min}, t_{\max}]$ ($\lambda(A)$ are dense in the same interval), $\lambda(\mathcal{A}) \in [t_{\min}, t_{\max}]$ and $\lambda(\mathcal{A} - A)$ cluster around zero;

ii) if $t_{\min} > 0$ then A and \mathcal{A} are p.d. $\forall n$ and $\lambda(S^{-1}AS^{-T})$, where $SS^T = \mathcal{A}$, cluster around 1, i.e. $\forall \varepsilon > 0 \exists \nu_\varepsilon, k_\varepsilon \mid \forall n > \nu_\varepsilon$ at least $n - k_\varepsilon$ eigenvalues $\lambda(S^{-1}AS^{-T})$ are in $(1 - \varepsilon, 1 + \varepsilon)$.

[We will call (assu) the hypothesis $\sum_{k=0}^{+\infty} |t_k| < +\infty, t_{\min} > 0$]

Proof of i):

By the definition of A , if $\mathbf{z} \in \mathbb{C}^n$ then

$$\begin{aligned} \mathbf{z}^* \mathbf{A} \mathbf{z} &= \sum_{k,j} \bar{z}_k \left[\frac{1}{2\pi} \int_{-\pi}^{\pi} t(\theta) e^{-i(k-j)\theta} d\theta \right] z_j = \frac{1}{2\pi} \int_{-\pi}^{\pi} (\sum_k \bar{z}_k e^{-ik\theta}) t(\theta) (\sum_j z_j e^{ij\theta}) d\theta \\ &\leq t_{\max} \frac{1}{2\pi} \int_{-\pi}^{\pi} (\dots)(\dots) d\theta = t_{\max} \frac{1}{2\pi} \sum_{k,j} \bar{z}_k z_j \int_{-\pi}^{\pi} e^{i(j-k)\theta} d\theta = t_{\max} \mathbf{z}^* \mathbf{z}. \end{aligned}$$

Thus $t_{\min} \mathbf{z}^* \mathbf{z} \leq \mathbf{z}^* \mathbf{A} \mathbf{z} \leq t_{\max} \mathbf{z}^* \mathbf{z}$, $\forall \mathbf{z} \in \mathbb{C}^n$. In particular we have $\lambda(A) \in [t_{\min}, t_{\max}]$.

The fact that $\lambda(\mathcal{A})$ are in $[t_{\min}, t_{\max}]$ follows from the stronger result that for any hermitian matrix A its best circulant squares fit \mathcal{A} is hermitian and $\min \lambda(A) \leq \lambda(\mathcal{A}) \leq \max \lambda(A)$ (see below).

The fact that $\lambda(\mathcal{A} - A)$ cluster around 0 and assertion ii) are proved after the section 'An interesting problem'.

end(Theorem(clusterTC))

Note that the proposed definition of \mathcal{A} can be extended from real symmetric Toeplitz A to any arbitrary matrix A . Moreover, as a consequence of the following expression

$$\mathcal{A} = Fd(F^* A F)F^*, \quad d(M) = \text{diag}([M]_{ii}),$$

verified for any A , the matrix \mathcal{A} is p.d. whenever A is p.d..

If $-1 < t < 1$, then $t_k = t^k$, $k = 0, 1, \dots$, verifies the hypothesis of Theorem(clusterTC). In fact

$$\begin{aligned} t(\theta) &= \sum_{k \in \mathbb{Z}} t^{|k|} e^{ik\theta} = \sum_{k=0}^{+\infty} t^k e^{ik\theta} + \sum_{k=1}^{+\infty} t^k e^{-ik\theta} \\ &= \frac{1}{1-te^{i\theta}} + \frac{1}{1-te^{-i\theta}} - 1 = \frac{1-t^2}{1+t^2-2t \cos \theta}, \\ 0 &< \frac{1-|t|}{1+|t|} \leq t(\theta) \leq \frac{1+|t|}{1-|t|}. \end{aligned}$$

Thus $A = [t^{|i-j|}]_{i,j=1}^n$ is p.d. for all n (for a direct proof of this fact see the decomposition below), \mathcal{A} is p.d. $\forall n$, and the eigenvalues of $\mathcal{A}^{-1}A$ cluster around 1.

G.STRANG ($n = 2m$ even). Set $\mathcal{A} = C(\mathbf{z}_A^{(GS)})$ where, if $\mathbf{z}^{GS} = \mathbf{z}_A^{(GS)}$

$$z_{i+1}^{GS} = t_i, \quad i = 0, \dots, m, \quad z_{i+1}^{GS} = t_{n-i}, \quad i = m+1, \dots, n-1.$$

Note that \mathcal{A} is real symmetric like A and $\|A - \mathcal{A}\|_? = \min_{\mathbf{z}} \|A - C(\mathbf{z})\|_?$ [RC].

The proof of the following result shows that under the same assumptions (assu) of Theorem(clusterTC) on the real sequence $\{t_k\}$ such \mathcal{A} (besides being p.d. as A for large n) has exactly the required PROPERTY, i.e. for any $\varepsilon > 0$ there exist n_ε , r_ε and a splitting of A of type $A = \mathcal{A} + R + E$, where the matrices R and E satisfy the conditions $\text{rank}(R) < r_\varepsilon$, $\|E\| < \varepsilon$ for all $n > n_\varepsilon$. As a consequence the eigenvalues of $\mathcal{A}^{-1}A$ cluster around 1.

Theorem(clusterGS). Identical to Theorem(clusterTC) except that \mathcal{A} is p.d. only for large n and it is not known if $\lambda(\mathcal{A}) \subset [t_{\min}, t_{\max}]$ (but then how can point ii) be proved? check. . .)

Proof [SIAM Review 38 (1996), pp.427-482]

However, the G.Strang definition of \mathcal{A} cannot be easily extended to a matrix which is not a real symmetric Toeplitz matrix. Moreover, it can happen that \mathcal{A} does not inherit p.d. from A .

For example

$$A = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 2 \end{bmatrix} \Rightarrow \mathcal{A} = \begin{bmatrix} 2 & -1 & & -1 \\ -1 & 2 & \ddots & \\ & \ddots & \ddots & -1 \\ -1 & & -1 & 2 \end{bmatrix}$$

and A is p.d., whereas the eigenvalues of \mathcal{A} are $2 - 2 \cos \frac{2\pi(j-1)}{n}$, $j = 1, \dots, n$, so \mathcal{A} is only semi p.d.. In fact, not both the sufficient conditions (assu) in Theorem(clusterGS) are verified: $\sum_{k=0}^{+\infty} |t_k| = 3 < +\infty$, but $t(\theta) = 2 - 2 \cos \theta$ is not positive in $[-\pi, \pi]$.

If $-1 < t < 1$, then $t_k = t^k$, $k = 0, 1, \dots$, verifies the hypothesis of Theorem(clusterGS). Thus $A = [t^{|i-j|}]_{i,j=1}^n$ is p.d., \mathcal{A} is p.d. for large n and the eigenvalues of $\mathcal{A}^{-1}A$ cluster around 1.

Actually, in this particular case the eigenvalues of $\mathcal{A}^{-1}A$ are known in explicit form. In fact, it can be shown that

$$Ae_m = \mathcal{A}e_m, Ae_{m+1} = \mathcal{A}e_{m+1}, A \begin{bmatrix} t \\ \mathbf{0} \\ -t^m \\ \mp t^m \\ \mathbf{0} \\ \pm t \end{bmatrix} = \frac{1}{1 \pm t} \mathcal{A} \begin{bmatrix} t \\ \mathbf{0} \\ -t^m \\ \mp t^m \\ \mathbf{0} \\ \pm t \end{bmatrix},$$

$$A \begin{bmatrix} \mathbf{y}_i \\ \pm \mathbf{y}_i \end{bmatrix} = \frac{1}{1 \mp t^m} \mathcal{A} \begin{bmatrix} \mathbf{y}_i \\ \pm \mathbf{y}_i \end{bmatrix}, \begin{bmatrix} 1 & t & \dots & t^{m-1} \\ t^{m-1} & \dots & t & 1 \end{bmatrix} \mathbf{y}_i = \mathbf{0}$$

[ES,Tor Vergata] (Proof ...). So, they are: $\frac{1}{1-t}$, $\frac{1}{1+t}$, $\frac{1}{1-t^m}$ $m-2$ times, $\frac{1}{1+t^m}$ $m-2$ times, and 1 twice (Proof ...). There is a clustering around 1 (there are two outliers), as we expected from Theorem(clusterGS). Moreover, only five eigenvalues are distinct.

E.TYRTY. Set $\mathcal{A} = C(\mathbf{z}_A^{(ET)})$, where $\mathbf{z}^{ET} = \mathbf{z}_A^{(ET)}$ is such that $C(\mathbf{z}^{ET})$ is an optimal circulant solution of the $C + R$ approximation problem (of A), C circulant, rank(R) low [ZOT], [OT]. An approximation of \mathbf{z}^{ET} can be computed via the black dot algorithm [OT]. Experiments show that such approximation of \mathcal{A} inherits from A real symmetry but not p.d.. ... ET construction via black dot more expensive than TC, GS ? more recent references ? ET construction via black dot for the coefficient matrix of the google system ? Consider the $\mathcal{L} + R$ approximation problem, \mathcal{L} more general spaces ? extend black dot algorithm to solve it ? Algorithms alternative to black dot?

See, in the enclosed Figure, the eigenvalues of the three matrices $A = [(0.5)^{|i-j|}]_{i,j=1}^n$, $\mathcal{A}^{-1}A$ with $\mathcal{A} = C(\mathbf{z}_A^{(TC)})$ and $\mathcal{A}^{-1}A$ with $\mathcal{A} = C(\mathbf{z}_A^{(GS)})$ ($n = 16, 32$).

AN INTERESTING PROBLEM

An interesting problem is inverting $D + C$ where D is diagonal and C circulant (by Tyrty, Rome 2006).

Observe that there exist a unitary matrix Q and a unitary diagonal matrix Λ such that

$$F = Q^* \Lambda Q,$$

thus

$$Q(D+C)Q^* = Q(D+F\tilde{D}F^*)Q^* = Q(D+Q^*\Lambda Q\tilde{D}Q^*\Lambda^*Q)Q^* = QDQ^* + \Lambda Q\tilde{D}Q^*\Lambda^*.$$

Remarks. Set $W = \sqrt{2}F$ where F is the Fourier matrix of order 2, i.e.

$$W = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Then $\det(\lambda I - W) = (\lambda^2 - 1) - 1 = \lambda^2 - 2$. It follows that the eigenvalues of F are: 1, -1. The corresponding eigenvectors are

$$\begin{bmatrix} \frac{1}{2}\sqrt{2+\sqrt{2}} \\ \frac{1}{2}(\sqrt{2}-1)\sqrt{2+\sqrt{2}} \end{bmatrix}, \quad \begin{bmatrix} \frac{1}{2}\sqrt{2-\sqrt{2}} \\ -\frac{1}{2}(\sqrt{2}+1)\sqrt{2-\sqrt{2}} \end{bmatrix}.$$

Proof: ...

Set $W = \sqrt{3}F$ where F is the Fourier matrix of order 3, i.e.

$$W = \begin{bmatrix} 1 & 1 & 1 \\ 1 & a & \bar{a} \\ 1 & \bar{a} & a \end{bmatrix}, \quad a = -\frac{1}{2} + \mathbf{i}\frac{\sqrt{3}}{2}.$$

Then

$$\begin{aligned} \det(W - \lambda I) &= (1 - \lambda)[(a - \lambda)^2 - \bar{a}^2] - [(a - \lambda) - \bar{a}] + [\bar{a} - (a - \lambda)] \\ &= (1 - \lambda)(a^2 - 2a\lambda + \lambda^2 - \bar{a}^2) + 2[(\bar{a} - a) + \lambda] \\ &= (1 - \lambda)[\lambda - (a + \bar{a})][\lambda - (a - \bar{a})] + 2[\lambda + (\bar{a} - a)] \\ &= (1 - \lambda)(\lambda + 1)(\lambda - \mathbf{i}\sqrt{3}) + 2(\lambda - \mathbf{i}\sqrt{3}) \\ &= (\lambda - \mathbf{i}\sqrt{3})(3 - \lambda^2). \end{aligned}$$

So, the eigenvalues of F are: 1, -1, \mathbf{i} .

The eigenvalues of the 4×4 Fourier matrix

$$F = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \mathbf{i} & -1 & -\mathbf{i} \\ 1 & -1 & 1 & -1 \\ 1 & -\mathbf{i} & -1 & \mathbf{i} \end{bmatrix}$$

are: 1, -1, \mathbf{i} , 1.

By using MATHEMATICA, the eigenvalues of F are:

$n = 6$: 1, -1, \mathbf{i} , 1, $-\mathbf{i}$, -1;

$n = 8$: 1, -1, \mathbf{i} , 1, $-\mathbf{i}$, -1, \mathbf{i} , 1.

We conjecture that

$$F = Q^*\Lambda Q = Q^* \begin{bmatrix} I & & & \\ & \mathbf{i}I & & \\ & & -I & \\ & & & -\mathbf{i}I \end{bmatrix} Q.$$

It follows that

$$F^2 = Q^* \begin{bmatrix} I & & & \\ & -I & & \\ & & I & \\ & & & -I \end{bmatrix} Q, \quad (QQ^T)\Lambda(\overline{Q}Q^*) = \Lambda.$$

NOTE: $F^4 = I$ (it is in fact known that F^2 is a permutation matrix) \Rightarrow for the eigenvalues λ of F we have $\lambda^4 = 1$, so the conjecture is proved.

Find Q ; properties of Q ; find \mathbf{v} such that, for all i , $[Q^T \mathbf{v}]_i \neq 0$, so that for the algebra of matrices simultaneously diagonalized by Q , $sdQ = \{QDQ^* : D \text{ diagonal}\}$ we have the representation

$$sdQ = \{Qd(Q^T \mathbf{z})d(Q^T \mathbf{v})^{-1}Q^* : \mathbf{z} \in \mathbb{C}^n\}.$$

What structure has sdQ ?

$\mathbf{v} = \mathbf{e}_1$: $QF = \Lambda Q \Rightarrow FQ^T \mathbf{e}_1 = Q^T \mathbf{e}_1 \Rightarrow$ PROBLEM: find \mathbf{z} , $F\mathbf{z} = \mathbf{z}$, $\|\mathbf{z}\|_2 = 1$; $z_i \neq 0 \forall i$?

Proof of TheoremClusterTC

Let us show that $\lambda(\mathcal{A} - A)$ cluster around 0. Let $R^{(N)}$ and $E^{(N)}$ be $n \times n$ matrices such that $\mathcal{A} - A = R^{(N)} + E^{(N)}$ and $E^{(N)}$ is null except for its upper-left $n - N \times n - N$ submatrix which is equal to the corresponding upper-left submatrix of $\mathcal{A} - A$. Moreover, observe that

$$\begin{aligned} (\mathcal{A} - A)_{ij} &= \mathcal{A}_{ij} - t_{|i-j|} = \frac{1}{n} [(n - |i - j|)t_{|i-j|} + |i - j|t_{n-|i-j|}] - t_{|i-j|} \\ &= -\frac{1}{n}|i - j|(t_{|i-j|} - t_{n-|i-j|}). \end{aligned}$$

Thus

$$\begin{aligned} \rho(E^{(N)}) &\leq \|E^{(N)}\|_1 = \max_j \sum_i |e_{ij}^{(N)}| = \max_{j=1 \dots n-N} \sum_{i=1}^{n-N} |(\mathcal{A} - A)_{ij}| \\ &= \max_{j=1 \dots n-N} \frac{1}{n} \sum_{i=1}^{n-N} |i - j| |t_{|i-j|} - t_{n-|i-j|}| \leq \frac{2}{n} \sum_{j=1}^{n-N-1} j |t_j - t_{n-j}| \\ &\leq \frac{2}{n} \sum_{j=1}^{n-N-1} j |t_{n-j}| + \frac{2}{n} \sum_{j=1}^{n-N-1} j |t_j| = \frac{2}{n} \sum_{j=N+1}^{n-1} (n-j) |t_j| + \frac{2}{n} \sum_{j=1}^{n-N-1} j |t_j| \\ &= 2 \sum_{j=N+1}^{n-1} |t_j| - \frac{2}{n} \sum_{j=N+1}^{n-1} j |t_j| + \frac{2}{n} \sum_{j=1}^{n-N-1} j |t_j| \leq 2 \sum_{j=N+1}^{n-1} |t_j| + \frac{2}{n} \sum_{j=1}^N j |t_j|. \end{aligned}$$

(regarding the latter inequality, it is obvious if $n - N - 1 \leq N$; otherwise the terms with $j > N$ are liquidated by the first terms of the other sum).

Let $\varepsilon > 0$ be fixed. There exist $N_\varepsilon \mid 2 \sum_{j=N_\varepsilon+1}^{\infty} |t_j| < \varepsilon/2$, and $\bar{\nu}_\varepsilon \mid \forall n > \bar{\nu}_\varepsilon$ $\frac{2}{n} \sum_{j=1}^{N_\varepsilon} j |t_j| < \varepsilon/2$.

Thus there exist $N_\varepsilon, \bar{\nu}_\varepsilon \mid$ for all $n > \nu_\varepsilon = \max\{\bar{\nu}_\varepsilon, 2N_\varepsilon\}$ we have:

$$\rho(E^{(N_\varepsilon)}) < \varepsilon \text{ and } \text{rank}(R^{(N_\varepsilon)}) \leq 2N_\varepsilon.$$

The latter result implies that $R^{(N_\varepsilon)}$ has at least $n - 2N_\varepsilon$ null eigenvalues. So, if γ_i , $i = 1, \dots, n$, denote the eigenvalues of $\mathcal{A} - A$ in non-decreasing order, then, by the min-max eigenvalue representation theory (see below), we have that

$\forall n > \nu_\varepsilon$ at least $n - 2N_\varepsilon$ eigenvalues γ_i of $\mathcal{A} - A$ are such that

$$-\varepsilon < 0 + \min \lambda(E^{(N_\varepsilon)}) \leq \gamma_i \leq 0 + \max \lambda(E^{(N_\varepsilon)}) < \varepsilon$$

i.e. $\lambda(\mathcal{A} - A)$ cluster around 0.

Proof of ii):

By i), $t_{\min} > 0 \Rightarrow A$ p.d. $\Rightarrow \mathcal{A}$ p.d..

Now let $\tilde{\gamma}_i$, $i = 1, \dots, n$, be the eigenvalues of $I - \mathcal{A}^{-1}A$ or, equivalently, of $I - S^{-1}AS^{-T}$ (S real non singular s.t. $SS^T = A$) in non-decreasing order. We want to show that they cluster around 0, as the eigenvalues γ_i of $\mathcal{A} - A$.

By the min-max eigenvalue representation theory (see below), we have

$$\begin{aligned} \gamma_j &= \min_{V_j} \max_{\mathbf{x} \in V_j} f(\mathbf{x}), \\ \tilde{\gamma}_j &= \min_{V_j} \max_{\mathbf{x} \in V_j} \frac{\mathbf{x}^*(I - S^{-1}AS^{-T})\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \min_{V_j} \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} \end{aligned} \quad (1)$$

where $f(\mathbf{x}) = \mathbf{x}^*(\mathcal{A} - A)\mathbf{x}/\mathbf{x}^*\mathbf{x}$, $g(\mathbf{x}) = \mathbf{x}^*\mathcal{A}\mathbf{x}/\mathbf{x}^*\mathbf{x}$, $0 < t_{\min} \leq g(\mathbf{x}) \leq t_{\max}$.

It follows that, for $j = 1, \dots, n$,

$$\frac{1}{t_{\max}}|\gamma_j| \leq |\tilde{\gamma}_j| \leq \frac{1}{t_{\min}}|\gamma_j|. \quad (2)$$

This with what we know about the γ_i yields the assertion:

For all $n > \nu_\varepsilon$ at least $n - 2N_\varepsilon$ eigenvalues $\tilde{\gamma}_i$ of $I - \mathcal{A}^{-1}A$ satisfy the inequality $|\tilde{\gamma}_i| < \varepsilon/t_{\min}$

i.e. the thesis.

Proof of (1):

$$\begin{aligned} \min_{V_j} \max_{\mathbf{x} \in V_j} \frac{\mathbf{x}^*(I - S^{-1}AS^{-T})\mathbf{x}}{\mathbf{x}^*\mathbf{x}} &= \min_{V_j} \max_{\mathbf{x} \in V_j} \frac{\mathbf{x}^*S^{-1}(SS^T - A)S^{-T}\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \\ \min_{V_j} \max_{\mathbf{y} \in S^{-T}V_j} \frac{\mathbf{y}^*(\mathcal{A} - A)\mathbf{y}}{\mathbf{y}^*\mathcal{A}\mathbf{y}} &= \min_{S^{-T}V_j} \max_{\mathbf{y} \in S^{-T}V_j} \frac{\mathbf{y}^*(\mathcal{A} - A)\mathbf{y}}{\mathbf{y}^*\mathcal{A}\mathbf{y}} = \\ \min_{V_j} \max_{\mathbf{y} \in V_j} \frac{\mathbf{y}^*(\mathcal{A} - A)\mathbf{y}}{\mathbf{y}^*\mathcal{A}\mathbf{y}} &= \min_{V_j} \max_{\mathbf{y} \in V_j} \frac{\mathbf{y}^*(\mathcal{A} - A)\mathbf{y}/\mathbf{y}^*\mathbf{y}}{\mathbf{y}^*\mathcal{A}\mathbf{y}/\mathbf{y}^*\mathbf{y}} = \min_{V_j} \max_{\mathbf{y} \in V_j} \frac{f(\mathbf{y})}{g(\mathbf{y})} \end{aligned}$$

Proof of the right inequality in (2):

First note that

$\gamma_j \geq 0$ iff $\max_{\mathbf{x} \in V_j} f(\mathbf{x}) \geq 0$, $\forall V_j$ iff $\max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} \geq 0$, $\forall V_j$ iff $\tilde{\gamma}_j \geq 0$.

Then, for j such that $\gamma_j \geq 0$ we have:

$$\begin{aligned} 0 \leq \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} &= \max_{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0} \frac{f(\mathbf{x})}{g(\mathbf{x})} \leq \max_{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0} \frac{f(\mathbf{x})}{t_{\min}} \\ &= \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{t_{\min}} = \frac{1}{t_{\min}} \max_{\mathbf{x} \in V_j} f(\mathbf{x}), \quad \forall V_j \end{aligned}$$

$(\{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0\} \neq \emptyset)$ which imply $\tilde{\gamma}_j \leq \gamma_j/t_{\min}$; whereas, for j such that $\gamma_j < 0$,

$$\exists \hat{V}_j \mid 0 > \tilde{\gamma}_j = \max_{\mathbf{x} \in \hat{V}_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} \geq \max_{\mathbf{x} \in \hat{V}_j} \frac{f(\mathbf{x})}{t_{\min}} = \frac{1}{t_{\min}} \max_{\mathbf{x} \in \hat{V}_j} f(\mathbf{x}) \geq \frac{1}{t_{\min}} \min_{V_j} \max_{\mathbf{x} \in V_j} f(\mathbf{x}) = \frac{1}{t_{\min}} \gamma_j$$

$(f(\mathbf{x}) > 0, \text{ on the right of the first equality})$ which imply $-\tilde{\gamma}_j \leq -\gamma_j/t_{\min}$.

Proof of the left inequality in (2):

For j such that $\gamma_j \geq 0$ we have:

$$\begin{aligned} \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} &= \max_{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0} \frac{f(\mathbf{x})}{g(\mathbf{x})} \geq \max_{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0} \frac{f(\mathbf{x})}{t_{\max}} \\ &= \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{t_{\max}} = \frac{1}{t_{\max}} \max_{\mathbf{x} \in V_j} f(\mathbf{x}), \quad \forall V_j \end{aligned}$$

$(\{\mathbf{x} \in V_j, f(\mathbf{x}) \geq 0\} \neq \emptyset)$ which imply $\tilde{\gamma}_j \geq \gamma_j/t_{\max}$; whereas, for j such that $\gamma_j < 0$,

$$\exists \hat{V}_j \mid 0 > \gamma_j = \max_{\mathbf{x} \in \hat{V}_j} f(\mathbf{x}),$$

$$0 > \frac{1}{t_{\max}} \gamma_j = \max_{\mathbf{x} \in \hat{V}_j} \frac{f(\mathbf{x})}{t_{\max}} \geq \max_{\mathbf{x} \in \hat{V}_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} \geq \min_{V_j} \max_{\mathbf{x} \in V_j} \frac{f(\mathbf{x})}{g(\mathbf{x})} = \tilde{\gamma}_j$$

which imply $-\tilde{\gamma}_j \geq -\gamma_j/t_{\max}$.

LINEAR ALGEBRA et al. :

$A\mathbf{x} = \lambda\mathbf{x}$, $A\mathbf{y} = \mu\mathbf{y}$, $\lambda \neq \mu \Rightarrow \mathbf{x}$ and \mathbf{y} are linearly independent.

$\alpha \mathbf{x} + \beta \mathbf{y} = \mathbf{0} \Rightarrow p(A)(\alpha \mathbf{x} + \beta \mathbf{y}) = \mathbf{0}, \forall$ polynomials $p \Rightarrow \alpha p(\lambda) \mathbf{x} + \beta p(\mu) \mathbf{y} = \mathbf{0}$
 $\Rightarrow \alpha \mathbf{x} = \mathbf{0}$ ($p(x) = (x - \mu)/(\lambda - \mu)$) and $\beta \mathbf{y} = \mathbf{0}$ ($p(x) = (x - \lambda)/(\mu - \lambda)$) \Rightarrow
 $\alpha = \beta = 0$

$A \mathbf{x} = \lambda \mathbf{x}, A \mathbf{y} = \mu \mathbf{y}, \lambda \neq \mu, A = A^* \Rightarrow \mathbf{x}$ and \mathbf{y} are orthogonal, that is $\mathbf{y}^* \mathbf{x} = 0$.

$$\mathbf{y}^* A \mathbf{x} = \lambda \mathbf{y}^* \mathbf{x},$$

$$(\mathbf{y}^* A^* \mathbf{x})^* = \mathbf{x}^* A \mathbf{y} = \mu \mathbf{x}^* \mathbf{y} = \mu (\mathbf{y}^* \mathbf{x})^* = \overline{\mu \mathbf{y}^* \mathbf{x}}$$

$$\mathbf{y}^* A^* \mathbf{x} = \overline{\mu \mathbf{y}^* \mathbf{x}}$$

$$(\lambda - \overline{\mu}) \mathbf{y}^* \mathbf{x} = \mathbf{y}^* (A - A^*) \mathbf{x}$$

if A is normal?

$A \mathbf{x} = \lambda \mathbf{x}, A \mathbf{y} = \mu \mathbf{y}, A$ normal $\Rightarrow A(A^* \mathbf{x}) = A^* A \mathbf{x} = \lambda A^* \mathbf{x} \Rightarrow A^* \mathbf{x}$ is an
eigenvector of λ (besides \mathbf{x}); $A^* A \mathbf{y} = \mu A^* \mathbf{y} \Rightarrow \mathbf{x}^* A^* A \mathbf{y} = \mu \mathbf{x}^* A^* \mathbf{y}$

Polar decomposition - Marco Maddalena

Fourier series on Gantmacher

p.486 Wilkinson: Goldstine and Horwitz Jacobi for normal

Goldstine (Gauss and Fourier)

Isaacson Keller: a page on charact. meth. for 1st order equ.

and reference to Lax, Douglas

min-max eigenvalues representation theory

Theorem.

i) Let A be a hermitian matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. Then

$$\lambda_j = \min_{V_j \subset \mathbb{C}^n, \dim V_j = j} \max_{\mathbf{x} \in V_j, \mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}}$$

(in particular, $\lambda_1 = \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}}, \lambda_n = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}}, \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \in [\lambda_1, \lambda_n]$)

ii) Let A, B, C be hermitian matrices such that $C = A + B$, and let $\alpha_j, \beta_j, \gamma_j, j = 1, \dots, n$, be their eigenvalues in non-decreasing order. Then

$$\alpha_j + \beta_1 \leq \gamma_j \leq \alpha_j + \beta_n.$$

Proof ...

SVD

A $n \times n, a_{ij} \in \mathbb{C} \Rightarrow A = U \sigma V^*, U, V$ $n \times n$ unitary, $\sigma = \text{diag}(\sigma_i)$ with
 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0. \dots$

FFT: $F \cdot \mathbf{z}, \mathbf{z} \in \mathbb{C}^n$ can be computed in $O(n \log n)$ a.o.

$F = [\frac{1}{\sqrt{n}} \omega^{(i-1)(j-1)}]_{ij=1}^n, \omega = e^{-i2\pi/n}; W_n = \sqrt{n} F, n = 2m$ even. Then

$$W_n = \begin{bmatrix} I & D \\ I & -D \end{bmatrix} \begin{bmatrix} W_m & O \\ O & W_m \end{bmatrix} Q, Q = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & & & & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & & & 0 \\ 0 & 0 & 0 & 1 & & & 0 \\ 0 & & & & & 0 & 1 \end{bmatrix}$$

$D = \text{diag}(1, \omega, \dots, \omega^{m-1}),$

$$W_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Proof ...

decomposition of $A = [t^{|i-j|}]_{ij=1}^n$

$$A = [t^{|i-j|}]_{ij=1}^n = \begin{bmatrix} 1 & t & \dots & t^{n-1} \\ t & & & \\ \vdots & & & \\ t^{n-1} & & & \end{bmatrix}, \quad t \in \mathbb{R} \Rightarrow$$

$$A = \begin{bmatrix} 1 & & & \\ t & 1 & & \\ \vdots & & & \\ t^{n-1} & \dots & t & 1 \end{bmatrix} \begin{bmatrix} 1 & & & \\ & 1-t^2 & & \\ & & \ddots & \\ & & & 1-t^2 \end{bmatrix} \begin{bmatrix} 1 & t & \dots & t^{n-1} \\ & 1 & & \vdots \\ & & & t \\ & & & 1 \end{bmatrix}.$$

The formula follows easily from the equalities

$$\begin{bmatrix} 1 & t & \dots & t^{n-1} \\ t & & & \\ \vdots & & & \\ t^{n-1} & & & \end{bmatrix} = \begin{bmatrix} 1 & & & \\ t & 1 & & \\ \vdots & & & \\ t^{n-1} & \dots & t & 1 \end{bmatrix} + \begin{bmatrix} 1 & t & \dots & t^{n-1} \\ & 1 & & \vdots \\ & & & t \\ & & & 1 \end{bmatrix} - I,$$

$$\begin{bmatrix} 1 & -t & & \\ & 1 & \ddots & \\ & & \ddots & -t \\ & & & 1 \end{bmatrix} \begin{bmatrix} 1 & t & \dots & t^{n-1} \\ & 1 & & \vdots \\ & & & t \\ & & & 1 \end{bmatrix} = I.$$

If $t^2 = 1$ then

$$A = \begin{bmatrix} 1 & (\pm 1) & \dots & (\pm 1)^{n-1} \\ (\pm 1) & 1 & & \\ \vdots & & & \\ (\pm 1)^{n-1} & & & \end{bmatrix} = \begin{bmatrix} 1 \\ (\pm 1) \\ \vdots \\ (\pm 1)^{n-1} \end{bmatrix} [1 \ (\pm 1) \ \dots \ (\pm 1)^{n-1}]$$

is a rank 1 matrix.

If $t^2 \neq 1$ then A is non singular and A^{-1} is a tridiagonal (non Toeplitz) matrix (compute it!).

If $t^2 < 1$ then A is p.d.

REFERENCES

- [TC] Tony Chan
- [RC] Raymond Chan, CUHK, Hong Kong
- [L] Levinson
- [Tr] Trench
- [AG] Ammar, Gragg
- [Dick] B. Dickinson, rational Toeplitz, Math. Comput. 34 (1980), pp.227–233
- [ES] Oakland conference on PDE's. L.Bragg, J.Dettman, eds., Longmans, London 1987, G.Strang, A. Edelman, pp.109-117 (you can find this paper in Tor Vergata)

- [OT] I. Oseledets, E. Tyrtyshnikov, A unifying approach . . .
- [ZOT] N. Zamarashkin, I. Oseledets, E. Tyrtyshnikov
- [CDFFZ] C. Di Fiore, S. Fanelli, P.Zellini, On the best least squares fit to a matrix and its applications" (maiop2)