

THE PROBABILITY THAT A RATIONAL RANDOM
MATRIX HAS MULTIPLE EIGENVALUES

Zhang Zhinan

College of Mathematics and System Science

Xinjiang University

Urumqi, 830046, P.R. CHINA

e-mail: znz01@xju.edu.cn

Abstract: In this paper, we show that the Probability of a random floating point matrix to be diagonalizable nears 1. Thus no one numerical algorithm is effective for determining the Jordan block structure of a random floating point matrix having nonlinear elementary divisor.

AMS Subject Classification: 15A52

Key Words: Jordan canonical form, random matrix, Lebesgue measure

*

It is well known that ascertaining the Jordan Block Structure of a matrix is an important task in matrix computations [2]. Golub pointed out that the set of n -by- n diagonalizable matrixes is dense in $C^{n \times n}$, see [1]. On the other hand the computed result of the Jordan block structure of matrix A is

$$fl(JBS(A)) = JBS(A + \Delta A),$$

where ΔA is the perturbation as a result from rounding errors, $\|\Delta A\| \leq \delta, \delta > 0$ is the unique restriction on ΔA , thus ΔA in set $\{\Delta A \mid \|\Delta A\| \leq \delta, \delta > 0\}$ is random. So it is possible that $JBS(A + \Delta A)$ is diagonalizable in spite of matrix A having nonlinear elementary divisor, in this case the numerical computation has failed. Hence, the probability of that random matrix $A + \Delta A$ being diagonalizable plays a key role in the failure mentioned above.

As $A \in R^{n \times n}$,

$$\det(\lambda I - A) = \lambda^n + p_1 \lambda^{n-1} + \cdots + p_n.$$

This is a monic polynomial which corresponds to $P = (p_1, \cdots, p_{n-1}, p_n) \in R^n$ as 1-1. Therefore, the problem mentioned above can be summed up as follows:

$$\Omega(\delta) = \{(p_1, \cdots, p_n) \mid |p_i - a_i| \leq \delta, \delta > 0, p_i \in R, i = 1, 2, \cdots, n\},$$

$a_i, i = 1, \cdots, n$ are real constants.

Subset

$$W = \{(p_1, \cdots, p_n) \mid (p_1, \cdots, p_n) \in \Omega(\delta)$$

and $\lambda^n + p_1 \lambda^{n-1} + \cdots + p_n$ having multiple roots $\}$.

If $\mu_n \Omega(\delta)$ and $\mu_n W$ represent the Lebesgue measure of set $\Omega(\delta)$ and W respectively, thus

$$p = \frac{\mu_n W}{\mu_n \Omega(\delta)} \quad (1)$$

denotes the probability of a monic polynomial having multiple roots in the neighborhood of a given monic polynomial $f(\lambda) = \lambda^n + a_1 \lambda^{n-1} + \cdots + a_n$. When considering the question in the real field, we have already proved the following conclusion: $\mu_n W = 0$ [4] and since $\mu_n \Omega(\delta) = (2\delta)^n > 0$, so

$$p = \frac{\mu_n W}{\mu_n \Omega(\delta)} = 0, \quad (2)$$

i.e. the probability of a real random monic polynomial having multiple roots equal to zero.

The question we are going to discuss next is whether the probability of a rational random monic polynomial also have similar properties. Suppose the complete set

$$\widehat{\Omega}(\delta) = \{(p_1, \cdots, p_n) \mid |p_i - a_i| \leq \delta, \delta > 0, p_i \in Q, i = 1, 2, \cdots, n\}.$$

Here Q denotes the rational number field. Subset

$$\widehat{W} = \{(p_1, \cdots, p_n) \mid (p_1, \cdots, p_n) \in \widehat{\Omega}(\delta)$$

and $\lambda^n + p_1 \lambda^{n-1} + \cdots + p_n$ having multiple roots $\}$.

At first, the question on the rational number field cannot be simply regarded as the special situations of the similar questions on the real number field. For

example, $\mu_n \widehat{W} = 0$ and $\widehat{\Omega}(\delta) = 0$, too, because the elements of \widehat{W} and $\widehat{\Omega}(\delta)$ are rational points (i.e. that the points which coordinate are all rational numbers). So they are countable sets. If according to the definition of the probability in real number field, it will be

$$p = \frac{\mu_n \widehat{W}}{\mu_n \widehat{\Omega}(\delta)} = 0/0.$$

Obviously this is indefinite, for this we have to redefine the probability in rational number field as $p = \theta/\Sigma$, where θ is number of points in \widehat{W} , Σ is number of the points in $\widehat{\Omega}(\delta)$. Even though \widehat{W} and $\widehat{\Omega}(\delta)$ both contain many infinite elements in the usual situation, but according to the concrete meaning of elements of \widehat{W} and $\widehat{\Omega}(\delta)$ in this paper, as a fact, when $\Omega(\delta) \supseteq W$, thus $\widehat{\Omega}(\delta) \supseteq \widehat{W}$ and must hold $\Sigma \geq \theta$. Finally, without losing the generality, this paper uses the decimal base to express the rational number.

Definition. The set of the point $P = (p_1, \dots, p_{n-1}, p_n) \in R^n$ bounded by inequality

$$a_i < p_i < b_i, \quad b_i - a_i = 10^{-t}, \quad i = 1, \dots, n,$$

where a_i and b_i are real constant called n -dimensional t -order unit open-interval. Here we do not have to produce the general definition of n -dimensional t -order unit open-interval, because this paper only uses t -order unit open-interval whose edge and coordinating axis are parallel respectively.

Lemma 1. W is a bounded closed set.

The proof process is omitted.

Finite t -order open covering of W .

Lemma 2. W can be covered by finite n -dimensional t -order unit open-interval G_1, G_2, \dots, G_r and $G_i, (i = 1, \dots, n)$ overlapping each other at most 2^n times.

Because W is the bounded closed set, and each point of W has at least one n -dimensional t -order unit open-interval covering this point, therefore according to the finite covering theorem, there exists limited n -dimensional t -order unit open-interval G_1, \dots, G_r covering W , here $r = r(t)$. We further prove $G_j (j = 1, \dots, r)$ overlap each other at most 2^n times. In fact, we do not lose the generality. Suppose $\omega = (\omega_1, \omega_2, \dots, \omega_n) \in W$, simultaneously belongs to $G_1, G_2, \dots, G_l, l > 2^n$, namely, the following inequalities hold simultaneously

$$\begin{aligned} a_{11} < \omega_1 < b_{11}, & a_{12} < \omega_2 < b_{12}, & \dots & a_{1n} < \omega_n < b_{1n}, \\ a_{21} < \omega_1 < b_{21}, & a_{22} < \omega_2 < b_{22}, & \dots & a_{2n} < \omega_n < b_{2n}, \\ & \dots & \dots & \\ a_{l1} < \omega_1 < b_{l1}, & a_{l2} < \omega_2 < b_{l2}, & \dots & a_{ln} < \omega_n < b_{ln}, \end{aligned}$$

where $G_j = \{(p_1, \dots, p_n) | a_{ji} < p_i < b_{ji}, b_{ji} - a_{ji} = 10^{-t}, i = 1, \dots, n.\}$, $j = 1, 2, \dots, l$. Then there exist $2n$ constants $a'_1, b'_1; a'_2, b'_2; \dots, a'_n, b'_n$ such that

$$\begin{aligned} \max\{a_{11}, a_{21}, \dots, a_{l1}\} &< a'_1 < \omega_1 < b'_1 < \min\{b_{11}, b_{21}, \dots, b_{l1}\}, \\ \max\{a_{12}, a_{22}, \dots, a_{l2}\} &< a'_2 < \omega_2 < b'_2 < \min\{b_{12}, b_{22}, \dots, b_{l2}\}, \\ &\vdots \\ \max\{a_{1n}, a_{2n}, \dots, a_{ln}\} &< a'_n < \omega_n < b'_n < \min\{b_{1n}, b_{2n}, \dots, b_{ln}\}. \end{aligned} \tag{3}$$

Obviously

$$0 < b'_i - a'_i < \min\{b_{1i}, b_{2i}, \dots, b_{li}\} - \max\{a_{1i}, a_{2i}, \dots, a_{li}\} \leq b_{li} - a_{li} = 10^{-l}.$$

From this we may further find the constants $b''_1, b''_2, \dots, b''_n; a''_1, a''_2, \dots, a''_n$, so that the following n pairs of equalities hold

$$\begin{cases} b''_1 - a'_1 = 10^{-t}. \\ b'_1 - a''_1 = 10^{-t}; \end{cases} \quad \begin{cases} b''_2 - a'_2 = 10^{-t}, \\ b'_2 - a''_2 = 10^{-t}; \end{cases} \quad \dots \quad \begin{cases} b''_n - a'_n = 10^{-t}, \\ b'_n - a''_n = 10^{-t}. \end{cases}$$

From this we can see that $b''_1 > b'_1, a''_1 < a'_1; b''_2 > b'_2, a''_2 < a'_2; \dots; b''_n > b'_n, a''_n < a'_n$.

Respectively, take (a'_i, b'_i) and (a''_i, b''_i) as open-intervals in the i -th coordinate axis, establish 2^n n -dimensional t -order unit open-intervals G_k^* in $R^n, k = 1, 2, \dots, 2^n$. Notice that

$$(a'_i, b'_i) \bigcup (a''_i, b''_i) = (a''_i, b''_i), \quad i = 1, 2, \dots, n.$$

Therefore

$$\bigcup_{k=1}^{2^n} G_k^* = \{(p_1, \dots, p_n) | a''_i < p_i < b''_i, i = 1, \dots, n\}.$$

We can see that

$$G_j \subset \bigcup_{k=1}^{2^n} G_k^*,$$

thus

$$\bigcup_{k=1}^l G_k \subseteq \bigcup_{k=1}^{2^n} G_k^*. \tag{4}$$

From this, replacing $\bigcup_{k=1}^{2^n} G_k^*$ with $\bigcup_{k=1}^l G_k$, we still obtain limited n -dimensional t -order unit open-interval to cover W . But passing through this

kind of replacement, causes reducing $q = l - 2^n > 0$ of the number of n -dimensional unit open-intervals which cover W . Notice that r is limited, therefore such replacement stops, inevitably, after limited times, namely, W can be covered by limited n -dimensional t -order unit open-intervals, also for each point of W is covered at most 2^n n -dimensional t -order unit open-intervals. Without loss of the generality, supposing n -dimensional t -order unit open-intervals which satisfy this condition is G_1, G_2, \dots, G_r , thus

$$W \subseteq \bigcup_{i=1}^r G_i$$

and G_i overlap each other at most 2^n times.

We agree that, in R^n , the points whose coordinates all are rational numbers are called rational points. Specially, the point which coordinates are all as $K_i \star 10^{-t}$ where t (non-negative integer), K_i (integer) are called t -order rational mesh point.

Theorem. *The probability that the monic polynomial with random rational coefficients has multiple roots is zero.*

Proof. Because G_i is t -order unit open-interval, therefore G_i at most contains one t -order rational mesh point and

$$H(t) = \bigcup_{i=1}^r G_i$$

at most contains $r = r(t)$ t -order rational mesh points. $\Gamma(t)$ represents the number of t -order rational mesh points in W , $\theta(t)$ number of t -order rational mesh points in \widehat{W} , due to $\widehat{W} \subset W \subset H(t)$, so

$$\theta(t) \leq \Gamma(t) \leq r(t). \tag{5}$$

When calculating $\sum_{i=1}^{r(t)} \mu_n G_i$, the non-overlap section in

$$H(t) = \bigcup_{i=1}^r G_i$$

is measured only one time, but overlapping sections must be measured as many times as they overlap. When calculating $\mu_n(\sum_{i=1}^{r(t)} G_i)$, regardless of how many times they overlap in

$$H(t) = \bigcup_{i=1}^r G_i,$$

it is measured only one time. Therefore,

$$\sum_{i=1}^{r(t)} \mu_n G_i \leq 2^n \mu_n \left(\sum_{i=1}^{r(t)} G_i \right) = 2^n \mu_n H(t). \tag{6}$$

Suppose $\widehat{\Omega}(\delta)$ have $m_i(t)$ numbers t -order rational mesh points in the i -th coordinate axis, then the number of the t -order rational mesh points contained in $\widehat{\Omega}(\delta)$ is $\prod_{i=1}^n m_i(t)$. If $\Sigma(t)$ denotes the number of the t -order rational mesh points containing in $\Omega(t)$, thus $\prod_{i=1}^n m_i(t) = \Sigma(t)$. If $V(t)$ denotes the lebesgue measure (volume) of the t -order unit open-intervals, then (cf. (5))

$$\begin{aligned} \frac{\theta(t)}{\Sigma(t)} &\leq \frac{V(t)\Gamma(t)}{V(t) \prod_{i=1}^n (m_i(t) - 1)} \leq \frac{V(t)r(t)}{V(t) \prod_{i=1}^n (m_i(t) - 1)} \\ &= \frac{\sum_{i=1}^{r(t)} \mu_n G_i}{V(t) \prod_{i=1}^n (m_i(t) - 1)} \leq \frac{2^n \mu_n \left(\sum_{i=1}^{r(t)} G_i \right)}{V(t) \prod_{i=1}^n (m_i(t) - 1)}, \end{aligned} \tag{7}$$

thus

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\theta(t)}{\Sigma(t)} &\leq \lim_{t \rightarrow \infty} \frac{2^n \mu_n \left(\sum_{i=1}^{r(t)} G_i \right)}{V(t) \prod_{i=1}^n (m_i(t) - 1)} \\ &= 2^n \frac{\lim_{t \rightarrow \infty} \mu_n \left(\sum_{i=1}^{r(t)} G_i \right)}{\lim_{t \rightarrow \infty} (V(t) \prod_{i=1}^n (m_i(t) - 1))} = 2^n \frac{\mu_n W}{\mu_n \Omega(\delta)} = 0. \end{aligned} \tag{8}$$

On the other hand

$$\lim_{t \rightarrow \infty} \frac{\theta(t)}{\Sigma(t)} = \frac{\theta}{\Sigma}. \tag{9}$$

According to (8) and (9), we know that

$$\frac{\theta}{\Sigma} = 0. \tag{10}$$

Corollary 1. *The probability that a random rational matrix has no multiple eigenvalues is 1.*

Corollary 2. *The probability that a random rational matrix being diagonalizable equals to 1.*

In fact, if the matrix has no multiple eigenvalues then it must be diagonalizable.

Corollary 3. *Probability that polynomial whose coefficient is random floating point number has multiple roots approaching zero.*

In fact, suppose the set of floating point numbers is F (finite set), M is the maximal number in F , thus $\Omega'_f = \{(f_1, f_2, \dots, f_m) \mid f_i \leq M, f_i \in F\}$ may be regarded as a sample taken from $\widehat{\Omega} = \{(p_1, p_2, \dots, p_m) \mid p_i \leq M, p_i \in Q\}$, generally, this sample is quite big, therefore basic properties of Ω'_f are uniformed with that of $\widehat{\Omega}$. From the theorem it follows that Corollary 3 holds.

Corollary 4. *Probability that random floating point matrix to be diagonalizable approaches 1.*

In fact, this corollary is result of Corollary 1 and Corollary 2.

Corollaries 1 to 4 have to provide theory background for Wilkinson's experience theory "Matrix having exact non-linear devisors are almost non-existent in practical work"; "Rounding errors will usually lead to a matrix which no longer has non-linear elementary devisors", see [4].

It is not effective in practice to determinate the Jordan block structure of a matrix having non-linear elementary divisor by the numerical computation. Because according to the above theory the numerical computed result of the Jordan block structure of a matrix is almost, always diagonalizable, so it is almost, always defeated.

Specially, it is impossible to overcome this difficulty through increasing the machine precision, because according to the proof of Corollary 3 the higher the precision goes, the bigger the possibility the numerical computation has to treat the non-linear elementary divisor as a linear elementary divisor.

Acknowledgments

This work is supported by State key Laboratory of Scientific and Engineering Computing, Chinese Academy Sciences.

References

- [1] Gene H. Golub, Charles F. Van Loan, *Matrix Computations*, The Johns Hopkins University Press (1983).
- [2] Gene H. Golub, Charles F. Van Loan, *Matrix Computations*, The Johns Hopkins University Press (1996).
- [3] Gene H. Golub, J.H. Wilkinson, Conditional eigensystems and the computation of the Jordan canonical form, *SIAM Review*, **18** (1976), 578-619.

- [4] J.H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford University Press (1965).
- [5] Zhang Zhi-nan, The Jordan canonical form of a real random matrix, *Numerical Mathematics. A Journal of Chinese Universities*, **23**, No.4 (2001), 363-367.