A RAPID ITERATIVE DIAGONALIZATION METHOD FOR LARGE MATRICES

MIKLÓS. I. BÁN

Institute of Physical Chemistry, Attila József University P.O.Box 105., H-6701 Szeged, Hungary (Received November 7, 1989)

AN ACCELERATED ITERATIVE MATRIX DIAGONALIZATION TECHNIQUE REDUCING THE COMPUTATIONAL WORK IS DESCRIBED: THE METHOD IS USEFUL WHEN THE RATE OF CONVERGENCE IS NOTOBIOUSLY SLOW AND THE ITERATIVE PROCEDURE CAN NOT BE AVOIDED.

1. Introduction

Iterative procedures used to determine the eigenvectors (and the spectrum) of matrices yield the eigenvector belonging to the eigenvalue largest in absolute sense. The overall rate of convergence depends on the structure of the spectrum very strongly and the usual procedures require in most cases too many iterations to be of practical value. In order to determine the rest of eigenvectors (and eigenvalues) several methods reducing the rank of the matrix, or projective elimination techniques are generally used.

The diagonalization of large matrices is an acute problem extending to many fields of computational sciences (molecular physics, quantum chemistry, etc.). A typical computational task in recent molecular physics and quantum chemistry is the determination of some (e.g. generally the largest) eigenvalues and eigenvectors of the Hamiltonian at the post-SCF level.

In this paper an effective acceleration procedure reducing substantially the number of iteration steps necessary to reach convergence will be described. The procedure is especially useful in cases where the separation of eigenvalues is small.

2. Hethod

Let us consider a positive (semi-)definite matrix A with real elements and denote the ordered set of its eigenvalues by λ_i :

$$\lambda_{j} := (\lambda_{1} > \lambda_{2} \ge \lambda_{3} \ge \dots \ge \lambda_{n}) \in \mathbb{R}, \quad j = 1, \dots, n$$
(1)

and its eigenvectors by e;, where R is the assembly of real numbers.

As is well known, the existence of the following limit values! can be verified²[1-4]:

$$\lim_{\mathbf{k}\to\infty}\frac{1}{\lambda_1^{\mathbf{k}}}\mathbf{A}^{\mathbf{k}}{}^{\mathbf{0}}\mathbf{u} = \alpha_1 \mathbf{e}_1, \quad {}^{\mathbf{0}}\mathbf{u} = {}^{\mathbf{0}}\alpha_1 \mathbf{e}_1 + {}^{\mathbf{0}}\alpha_2 \mathbf{e}_2 + \dots + {}^{\mathbf{0}}\alpha_n \mathbf{e}_n, \quad (2)$$

$$\min_{(u)} Q = \frac{\overline{u} A u}{\overline{u} u} = \lambda_n, \qquad (3)$$

and

$$\max_{(u)} Q = \frac{\overline{u} A u}{\overline{u} u} = \lambda_1, \qquad (4)$$

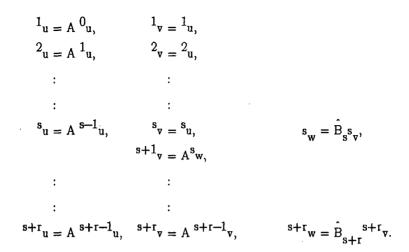
where k is the sequential number of the iteration step, ${}^{0}\alpha_{1}$ is the 0th coefficient vector and the vector³ u in the formulas (3) and (4) run over the whole configuration space, while Q is the so called Rayleigh quotient [5]. The classic iterative procedure, the "power method"[5] consists of consecutive multiplications by a trial vector. In the following section a modification of this method will be described.

Let us examine the following schemes (in the left and right columns the steps of the power method and the modified procedure respectively, are compared):

The superscripts on the left and right of symbols refer to the sequential number of iteration and the exponents, *resp.*

²Theorem of von MISES.

 $^{3\}overline{u}$ is the transpose of u.



The operator B_s performs a non-linear transformation of the trial vector. This step is embedded into the basic procedure and it will repeatedly be applied by using the actual vectors stored in the given cycle.

Now we try to construct the operator B_s in such a manner that the vector $s^{+1}v = A^s w$ approximates a vector k^u of the basic procedure where k>s+1 (or even more favourably: k>s+1).

3. Details of the procedure

To define the operator \hat{B}_s we express ${}^{0}u$, ${}^{k}u$, *etc.* by the eigenvectors $\{e_i\}$ in the space used for the computations⁴:

$${}^{\mathbf{0}}\mathbf{u} = \alpha_1 \mathbf{e}_1 + \alpha_2 \mathbf{e}_2 + \dots + \alpha_n \mathbf{e}_n, \tag{5}$$

$${}^{0}\mathbf{u}_{i} = \alpha_{1}\mathbf{x}_{i1} + \alpha_{2}\mathbf{x}_{i2} + \dots + \alpha_{n}\mathbf{x}_{in}, \tag{6}$$

$${}^{k}_{k} u = \alpha_{1} \lambda_{1}^{k} e_{1} + \alpha_{2} \lambda_{2}^{k} e_{2} + \dots + \alpha_{n} \lambda_{n}^{k} e_{n},$$
(7)

$${}^{s}u_{i} = \alpha_{1}\lambda_{1}^{k}x_{i1} + \alpha_{2}\lambda_{2}^{k}x_{i2} + \dots + \alpha_{n}\lambda_{n}^{k}x_{in}.$$
 (8)

⁴The coordinates of the eigenvectors in this system are: $\{e_i\} = \{x_{ii}\}$.

By constructing the sequence (2) we obtain:

$${}^{k}\widetilde{u}_{i} = \alpha_{1}x_{i1} + \alpha_{2}\left[\frac{\lambda_{2}}{\lambda_{1}}\right]^{k}x_{i2} + \dots + \alpha_{n}\left[\frac{\lambda_{n}}{\lambda_{1}}\right]^{k}x_{in} \approx \alpha_{1}x_{i1} + M\left[\frac{\lambda_{2}}{\lambda_{1}}\right]^{k},$$
(9)

where M is constant.

Each component vanishes separately according to a function

$$f_{ij}(k) = \gamma_{ij} a_j^k, \qquad (10)$$

where $a_j = (\lambda_j/\lambda_1)$ and the meaning of γ_{ij} follows from eqn.(9). Regarding (10) and (1) the rate of convergence is apparently determined by the second largest eigenvalue. Now combining (9) and(10)

$${}^{k}\widetilde{u}_{i} = \alpha_{1}x_{i1} + \sum_{j=2}^{n} \gamma_{ij}a_{j}^{k}$$
(11)

is obtained.

The next goal is to determine a single function

$$\Theta_{i}(\mathfrak{t}) = \xi_{i} e^{-\nu_{i} \mathfrak{t}}, \quad 0 \leq \mathfrak{t} < \infty$$
(12)

which approximates satisfactorily the function⁵

$$F_{i}(k) = \sum_{j=2}^{n} f_{ij}(k).$$
 (13)

⁵Function (13) having an integer argument is not a continuous function in contrast to (12) which has a real argument.

Substituting $f_{ij}(k)$ by

$$\varphi_{ij}(t) = \gamma_{ij} e^{-\beta_{ij}t}, \qquad (14)$$

then for integer values of t the equality

$${}^{k}\widetilde{u}_{i} = \alpha_{1}x_{i1} + \sum_{j=2}^{n}f_{ij}(k) = \alpha_{1}x_{i1} + \sum_{j=2}^{n}\varphi_{ij}(t)\Big|_{t=k}$$
(15)

will hold. It can readily be shown that as an exact substitute for (13) no function of type (12) is expected to be found. To realize this, let us consider the difference function $\omega_{i}(t)$ and its derivatives $\omega_{i}^{(r)}(t)$:

$$\omega_{i}(t) = \Theta_{i}(t) - \sum_{j=2}^{n} \varphi_{ij}(t) = \xi_{i} e^{-\nu_{i}t} - \sum_{j=2}^{n} \gamma_{ij} e^{-\beta_{ij}t} =$$

$$= \sum_{j=2}^{n} \left[\frac{1}{n-1} \xi_{i} e^{-\nu_{i}t} - \gamma_{ij} e^{-\beta_{ij}t} \right], \qquad (16)$$

$$\omega_{i}^{(r)}(t) = (-1)^{r} \xi_{i} \nu_{i}^{r} e^{-\nu_{i}t} - \sum_{j=2}^{n} (-1)^{r} \gamma_{ij} \beta_{ij}^{r} e^{-\beta_{ij}t} =$$

$$\sum_{j=2}^{n} (-1)^{r} \gamma_{i}^{r} \left\{ \frac{1}{n-1} \xi_{i} e^{-\nu_{i}t} - \gamma_{ij} \left[\frac{\beta_{ij}}{\nu_{i}} \right]^{r} e^{-\beta_{ij}t} \right\}. \qquad (17)$$

From (16) and (17) follows that at a given point, in general, $\omega_i(t)$ and its derivatives can not simultaneously be zero unless one of the following conditions is

fulfilled⁶:

$$\frac{\beta_{ij}}{\nu_i} = 1, \quad j = 2,...,n,$$

$$\beta_{il} \neq 0$$
, but $\beta_{ij} = 0$ if $j \neq l$,

$$\gamma_{ij} \neq 0$$
, but $\gamma_{ij} = 0$ if $j \neq l$.

Now for judging the possibility of an approximation, let us examine the behaviour of (15) describing the changes of the components of the iterated vectors in the course of iteration. We represent first (15) in the form

$${}^{k}\widetilde{u}_{i} = \alpha_{1}x_{i1} + \sum_{j=2}^{n} \gamma_{ij}e^{-\beta_{ij}t} = \alpha_{1}x_{i1} + \left\{\sum_{j=2}^{n} \gamma_{ij}e^{-z_{ij}t}\right\}e^{-\xi_{i}t}.$$
 (18)

The expression in the parentheses can be transcribed as follows:

$$\sum_{j=2}^{n} \gamma_{ij} e^{-z_{ij}t} = \sum_{j=2}^{n} \gamma_{ij} e^{-\delta_i(t)t}, \qquad (19)$$

where

$$\delta_{i}(t) = -\frac{\ln \left[\sum_{j=2}^{n} \gamma_{ij} e^{-z_{ij}t} \sum_{j=2}^{n} \gamma_{ij}\right]}{t} =$$

⁶All these cases lie outside the region of practical importance.

$$= -\frac{\ln \sum_{j=2}^{n} d_{ij} e^{-z_{ij}t}}{t},$$
 (20)

$$d_{ij} = \frac{\gamma_{ij}}{\sum_{l=2}^{n} \gamma_{ll}}$$
(21)

and

$$\sum_{j=2}^{n} d_{ij} = 1.$$
 (22)

The sum in the argument of the logarithm is increasing or decreasing strictly monotonously depending upon the signs of z_{ij} . In the first case ($z_{ij} < 0$) the range of the argument function is between 1 and ∞ if $0 \le t < \infty$ and the relations⁷

$$\ln\left[\sum_{j=2}^{n} d_{ij} e^{-z_{ij}t}\right] < \ln\left[\sum_{j=2}^{n} d_{ij} e^{-z_{i,\min}t}\right] \equiv$$
$$\equiv \ln(e^{-z_{i,\min}t}) = |z_{i,\max}|^{t}$$
(23)

hold.

If the second case $(z_{ij} > 0)$ occurs, the argument function ranges between 1 and 0 if $0 \le t < \infty$ and the relations

$$\ln\left[\sum_{j=2}^{n} d_{ij} e^{-z_{ij}t}\right] < \ln\left[\sum_{j=2}^{n} d_{ij} e^{-z_{i\min}t}\right] \equiv$$

 $r_{z_{i,\min}}$ and $z_{i,\max}$ denote the minimum and maximum elements in the set $z_{i,i}$

$$\equiv \ln(e^{-z_{i,\min}t}) = |z_{i,\min}|^{t}$$
(24)

will be satisfied.

For $z_{ij} < 0$

$$\lim_{t\to\infty} \delta_{i}(t) = -\frac{\ln |z_{i}|_{\max}}{t} = -|z_{i,\max}|$$
(25)

and for $z_{ij} > 0$

$$\lim_{t \to \infty} \delta_{i}(t) = -\frac{\ln |\mathbf{e}| \mathbf{z}_{i,\min}|^{t}}{t} = -|\mathbf{z}_{i,\min}|$$
(26)

will be valid.

It seems that in any case $\delta_i(t)$ becomes constant in the course of iteration. By using these last results we get an approximation for the exact formula⁸:

$${}^{k}\widetilde{u}_{i} = \alpha_{1}x_{i1} + \sum_{j=2}^{n}\gamma_{ij}e^{-\delta_{i}(t)t}e^{-\zeta_{i}t} \approx \alpha_{1}x_{i1} + \sum_{j=2}^{n}\gamma_{ij}e^{-(z_{i},\max/\min+\zeta_{i})t}.$$
 (27)

We expect this approximating formula (27) to be satisfactory for practical purposes only if it is used after having reached the required accuracy.

The computations proceed as follows:

i) Because of the possibility of the presence of negative eigenvalues (in contrast

⁸The subscript max/min means either maximum or minimum in conformity with the former remarks.

ii)

to our assumption about the positive semidefinite nature of A), if the accuracy is below a predetermined limit, each second vectors start to be stored. For computational efficiency it is practical and enough to use but three vectors. The formula (27) will be used in the form

$${}^{k}\widetilde{u}_{i} \approx S_{i}(t) = \Theta_{i}(t) + K_{i}.$$
⁽²⁸⁾

To have a descending sequence, if necessary, the ith components of the vectors used are to be transformed. By supposing that the following equation system can be set up for each component of the trial vector separately, it will be solved for all the unknown parameters⁹:

$$\xi_{i} e^{-\nu_{i} k_{1}} + K_{i} = {}^{k_{1}} b_{i}, \qquad (29)$$

$$\xi_{i}e^{-\nu_{i}k_{2}} + K_{i} = {}^{k_{2}}b_{i}, \qquad (30)$$

$$\xi_{i} e^{-\nu_{i} k_{3}} + K_{i} = {}^{k_{3}} b_{i}.$$
(31)

 k_1 ; $k_2 = k_1 + 2$; $k_3 = k_2 + 2$ are the serial numbers of iteration, ${}^{k_j}b_i$ are the ith components of the iterated vectors, and ξ_i , ν_i and K_i are the unknown parameters to be determined. To get ν_i one has to solve a second order equation, however, ξ_i and K_i are in linear dependence.

⁹It is justifiable to use an equation system, instead of using a least-square type procedure, only if we assume that (27) is exactly valid. Computational experience shows this procedure to be admissible without any significant deterioration concerning the accuracy of the results.

iii) The ith component of the last vector will be replaced by the limit value¹⁰ of (28):

$$\lim_{t \to \infty} S_i(t) = K_i.$$
(32)

 iv) The original iterative scheme will be continued by repeating the same procedure, until the accuracy is above a predetermined limit.

If there are only positive eigenvaluee, to enchance the efficiency of the procedure three consecutive vectors are to be used in *eqns*. (29-31). In special cases, for increasing the separation of eigenvalues any usual method (*e.g.* raising the matrix to power, *etc.*) can be appealed to.

Examples

The results obtained in the course of testing the procedure are summarized in Table I. In the columns 1, 2 and 3 the dimensions of matrices (MD)¹¹, the ratios of the total number of iterations (NIR) in the accelerated procedure related to that in the basic procedure and the ratios of the total iteration times (TIR) formed similarly as NIR are displayed. For testing purposes real symmetrical matrices having pseudo random number elements were used.

All the testing examples -quite independently from the dimensions of the matrices used-show the savings in the total numbers and times of iterations mainly to be between 0.4 and 0.6, *i.e.* the acceleration reduces both quantities roughly to the halves of their original values, thus giving hope to expect similar gains at larger matrices. When using a backstorage device for the high dimension of the matrix to be diagonalized, the gain in

¹⁰If necessary, it has to be retransformed.

¹¹The dimensions of matrices were confined to 55 regarding that -for practical reasons- the method was tested on a microcomputer.

computer time is expected to be even larger.

Data comparing results		
MD	NIR	TIR
3	0.52	0.60
3	0.56	0.69
4	0.59	0.68
5	0.60	0.70
6	0.52	0.58
6	0.48	0.52
7	0.60	0.68
8	0.51	0.58
8	0.52	0.55
10	0.51	0.54
10	0.32	0.44
12	0.48	0.51
20	0.51	0.53
30	0.56	0.57
55	0.56	0.57

Table I

Acknowledgement

In writing up this paper thanks are due to Dr. Imre Bálint for the consent of using many results of his work, consultations and advice during the period working in my group between January 1985 and July 1987. The support of this work is acknowledged to the Hungarian Ministry of Higher Education (grant Nos. P-437/84. and 746/86).

Leferences

- Varga, L.S.: Matrix Iterative Analysis. Prentice Hall, Englewood Cliffs, 1962.
- [2] Householder, A.S.: Principles of Numerical Analysis. McGraw-Hill Book Co., New York, 1953.
- [3] Householder, A.S.: The Theory of Matrices in Numerical Analysis. Blaisdell, New York, 1964.
- [4] Lánczos, C.: Applied Analysis. Prentice Hall, Englewood Cliffs, 1961.
- [5] Axelsson, O.: "Solution of Linear System of Equations: Iterative Methods" in Sparse Matrix Techniques (V.A. Barker, ed.), Springer, Berlin-Heidelberg-New York, 1977.

МЕТОД БЫСТРОЙ ИТЕРАТИВНОЙ ДИАГОНАЛИЗАЦИИ ДЛЯ БОЛЬШИХ МАТРИКСОВ

М.И. БАН

Описана ускоренная техника итеративной диагонализации матриксов, позволяющая сократить время необходимое для проведения расчетов метод особенно полезень, когда скорость конвергенции чрезвычайно мала и итерационный процесс не может быть завершен.