

COLLOQUES INTERNATIONAUX
DU
CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE

N° 165

PROGRAMMATION
EN MATHÉMATIQUES
NUMÉRIQUES

Besançon
7-14 Septembre 1966

EXTRAIT

ÉDITIONS DU CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE
15, quai Anatole-France — Paris - VII
1968

NEWTON-LAGUERRE ITERATION

by T.J. DEKKER

Mathematisch Centrum, Amsterdam

RESUME

Dans cette communication une généralisation de la formule de Laguerre est définie.

Appliquant cette formule on obtient un processus de convergence localement cubique, quelle que soit la multiplicité des zéros. Quelques expériences numériques de calcul des valeurs propres (complexes) de matrices réelles montrent que, pour des matrices ayant des valeurs propres multiples, la formule est préférable à la formule classique de Laguerre.

SUMMARY

This paper deals with a generalisation of Laguerre's formula for finding zeros of polynomials.

Applying this formula iteratively one obtains a (locally) cubically convergent process, irrespective the multiplicity of the zeros. Some numerical experiments with an ALGOL 60 program for finding the (complex) eigenvalues of real matrices show that, for matrices with multiple eigenvalues, the formula is preferable above the classical Laguerre formula.

*
* *
*

1. INTRODUCTION

The formula of Laguerre (see (2.4) below) is very useful for finding the zeros of analytic functions, especially polynomials.

Suppose we are given an analytic function f .

Let z be a complex number in a certain neighbourhood of a zero r of f . Replacing successively z by $L(z)$ (cf. 2.4), we obtain a series converging to r , if the initial value z is sufficiently close to r . The order of convergence depends on the choice of the two parameters p and q . In particular we get cubic convergence if p is equal to the multiplicity of r and q is a positive constant.

In practice, however, we do not know the multiplicity in advance. B. Parlett [2] assumes in his procedure for calculating the complex eigen-

values of a real matrix that the multiplicity of the eigenvalues is either one or two, the latter in the case that convergence is slow. The purpose of this paper is to show how in each step of the iteration the multiplicity may be estimated by means of a Newton formula. If this estimate is used as the value of p in Laguerre's formula, we obtain cubic convergence for multiple zeros also. The author has done some numerical experiments with this formula and has the impression that the method is favourable if many zeros are multiple or in clusters.

In sections 2, 3 and 4, some formulas are derived and some theorems about the order of convergence are proved. It is assumed that the given function f is analytic in a neighbourhood of a zero r of f , so that it makes sense to speak about the multiplicity of that zero and about Taylor expansion. As to convergence, only the local convergence, i.e. convergence in a sufficiently small neighbourhood of r is considered (although Laguerre's formula, in fact, is attractive especially because of its good global convergence property).

Section 5 discusses the choice of the two parameters, p and q , in Laguerre's formula and section 6 gives some practical upper bounds for the error in the case that the function is a polynomial.

Section 7 describes some features and results of an ALGOL 60 program for calculating the eigenvalues of a real matrix.

Acknowledgements.

The author is indebted to the staff members of the Computation department of the Mathematical Centre for their valuable suggestions and to R.J. De Vogelaere for his suggestions on iteration and iteration control.

2. THE FORMULA OF LAGUERRE

To derive Laguerre's formula, we start from the interpolating function :

$$2.1 \quad f^*(z) = c(z - a)^p (z - b)^q$$

and choose, for given p and q , the parameters a , b and c such that the values of f^* and its first and second derivative are equal to those of the given function f . Let :

$$2.2 \quad s_1 = f'(z)/f(z) , \quad s_2 = s_1^2 - f''(z)/f(z) .$$

We then find :

$$2.3 \quad s_1 = \frac{p}{z-a} + \frac{q}{z-b}, \quad s_2 = -\frac{ds_1}{dz} = \frac{p}{(z-a)^2} + \frac{q}{(z-b)^2}.$$

Let the zero a of f^* be denoted by $L(z)$.

Putting $s_0 = p + q$ and eliminating b , we then obtain Laguerre's formula :

$$2.4 \quad L(z) = z - s_0 / (s_1 + \sqrt{(q/p)(s_0 s_2 - s_1^2)})$$

(cf. [4] and [1] formula (7)).

Here in, we choose the sign of the square root such that the absolute value of $L(z) - z$ is minimal.

We then have

2.5. Theorem.

The series $z, L(z), L(L(z)), \dots$ converges cubically in a neighbourhood of a zero r of f having multiplicity m , if p and q satisfy

$$2.5.1 \quad p = m + O(z-r)^2, \quad q > \text{constant} > 0.$$

Proof.— For convenience, we carry out a translation such that the root r coincides with the origin. Then $s_1 = \frac{m}{z} + \sigma$, where σ tends to a finite value and $s_2 = (m + O(z^2))/z^2$.

The condition on p now reads $p = m + O(z^2)$ and we have :

$$\begin{aligned} z \sqrt{(q/p)(s_0 s_2 - s_1^2)} &= \sqrt{(q/m)(s_0 m - m^2 - 2m\sigma z) + O(z^2)} \\ &= \sqrt{q(q - 2\sigma z) + O(z^2)} = q - \sigma z + O(z^2). \end{aligned}$$

So we find for the Laguerre iterate (we have to take the + sign, since m and q are both positive) :

$$\begin{aligned} L(z) &= z - s_0 z / (m + \sigma z \pm (q - \sigma z) + O(z^2)) \\ &= z - s_0 z / (s_0 + O(z^2)) = O(z^3). \end{aligned}$$

If the sign of the square root is chosen the other way, the same theorem

holds but for the condition on q which has to be replaced by " $q < \text{constant} < 0$ ".

If, however, the condition on p is replaced by $p = m + O(z - r)$ the convergence is only quadratic.

3. THE FORMULA OF NEWTON

Let :

$$3.1 \quad f^*(z) = c(z - a)^p ,$$

in which a and c are chosen such that f^* and f have equal function value and derivative at z . Then we have :

$$3.2 \quad s_1 = f'(z)/f(z) = p/(z - a) .$$

Denoting the zero a of f^* by $N(z)$ we obtain Newton's formula :

$$3.3 \quad N(z) = z - p/s_1 .$$

3.4. Theorem.

The series $z, N(z), N(N(z)), \dots$ converges quadratically in a neighbourhood of an m -fold zero r of f if $p = m + O(z - r)$. Moreover, if the Taylor series of f is :

$$3.4.1 \quad f(z) = a_m(z - r)^m + a_{m+1}(z - r)^{m+1} + O(z - r)^{m+2}$$

and p satisfies :

$$3.4.2 \quad p = m + \frac{a_{m+1}}{a_m}(z - r) + O(z - r)^2 ,$$

then the convergence is cubic.

Proof.— We again shift r to the origin. Then we have :

$$1/s_1 = f(z)/f'(z) = \frac{z}{m} \left(1 - \frac{a_{m+1}}{a_m} \frac{z}{m} + O(z^2) \right) .$$

Let $p = m + \mu z + O(z^2)$, then :

$$\begin{aligned} N(z) &= z - z \left(1 + \mu \frac{z}{m} \right) \left(1 - \frac{a_{m+1}}{a_m} \frac{z}{m} \right) + O(z^3) \\ &= \left(\frac{a_{m+1}}{a_m} - \mu \right) \frac{z^2}{m} + O(z^3) . \end{aligned}$$

so the convergence is cubic, if $\mu = \frac{a_{m+1}}{a_m}$, and otherwise quadratic.

4. FORMULA FOR THE MULTIPLICITY

We again start from the interpolating function (3.1), but we now choose the parameters a , c and p such that the values of f^* and its first and second derivative are equal to those of f . Then :

$$4.1 \quad s_1 = p/(z - a) , \quad s_2 = p/(z - a)^2 .$$

Denoting the zero a of f^* by $M(z)$ we have for p and $M(z)$ (cf. [4]) :

$$4.2 \quad p = s_1^2/s_2$$

$$4.3 \quad M(z) = z - p/s_1 = z - s_1/s_2 .$$

This is, in fact, the ordinary Newton formula $N(z)$ with $p = 1$ for the function $f(z)/f'(z)$.

4.4. Theorem.

The series $z, M(z), M(M(z)), \dots$ converges quadratically in a neighbourhood of a zero r of f and the values p converge linearly to the multiplicity of r .

Proof.— Let f have the Taylor expansion (3.4.1). Then it easily follows that :

$$4.4.1 \quad s_1^2/s_2 = m + 2 \frac{a_{m+1}}{a_m} (z - r) + O(z - r)^2$$

Thus, p converges linearly to m and, according to theorem (3.4), the convergence of the iteration M is quadratic.

5. CHOICE OF p AND q IN LAGUERRE'S FORMULA

We may choose p according to (4.2) and use this value in Laguerre's formula (2.4). Since condition (2.5.1) is not satisfied (cf. 4.4.1), we obtain not cubic, but only quadratic, convergence. In fact, the formula thus obtained is equivalent to Newton's formula (4.3) provided the sign of the square root in (2.4) is suitably chosen. (If z is sufficiently near the limit, this choice coincides with the choice mentioned at (2.4)). In order to obtain a cubically convergent process, we must therefore use a better estimate for the multiplicity. As the multiplicity is an integer, we may simply choose p as the integer which is nearest to s_1^2/s_2 . If we are sufficiently near the limit, this yields the correct value and we get cubic convergence.

If, on the other hand, we are not near the limit, it may very well happen that this value for p is useless, either because it is negative or 0 (which happens often if f has non-real zeros) or because, in the case that f is a polynomial, it exceeds the degree. In these cases, the most obvious choices for p seem to be $p = 1$ or $p = \text{degree}$, respectively. In this way, we may obtain a negative argument for the square root and thus a real start leads in a natural way to non-real iterates. Of course, then s_1^2/s_2 will become non-real also, in which case we simply disregard its imaginary part for the determination of p .

As to the choice of q , reasonable values seem p and ∞ or, in case of polynomials, $\text{degree}-p$. In the latter case we should avoid the situation $p = \text{degree}$, $q = 0$, since this formula would not converge cubically; so it is better to take $p \leq \text{degree} - 1$ and thus $q \geq 1$.

Summarizing, we obtain the following choices for p and q :

- 5.1 if f is not a polynomial :
- $p =$ the positive integer nearest to s_1^2/s_2 ,
- $q = \infty$ or $q = p$;
- 5.2 if f is a polynomial of degree $n \geq 2$:

$p =$ the positive integer smaller than n nearest to s_1^2/s_2 ,
 $q = n - p$ or $q = \min(n - p, p)$.

6. UPPER BOUNDS FOR THE ERROR OF ZEROS OF POLYNOMIALS

If the given function f is a polynomial of degree n , we have the following upper bounds for the error $z - r$, where r is the nearest zero of f .

a) expressed in the Newton step $\Delta z = N(z) - z$ with $p = 1$ (cf. 3.3) :

$$6.1 \quad |z - r| \leq n/|s_1| = n|\Delta z| ;$$

b) expressed in the Laguerre step $\Delta z = L(z) - z$ (cf. 2.4), where p satisfies $1 \leq p < n$:

$$6.2 \quad q = n - p \longrightarrow |z - r| \leq (1 + \sqrt{2(n-1)})|\Delta z| ,$$

$$6.3 \quad q = \infty \longrightarrow |z - r| \leq \sqrt{n} |\Delta z| ,$$

$$6.4 \quad q = p \longrightarrow |z - r| \leq n |\Delta z| .$$

We prove only (6.2).

$$\begin{aligned} |z - r| &= \frac{|\Delta z|}{n} |s_1 \pm \sqrt{\frac{q}{p} (ns_2 - s_1^2)}| |z - r| \\ &\leq \frac{|\Delta z|}{n} (n + \sqrt{(n-1)(n^2 + n^2)}) = (1 + \sqrt{2(n-1)})|\Delta z|. \end{aligned}$$

These upper bounds are certainly not all best possible.

In practice, however, they are good enough, especially (6.2 & 3). So the criterion “ $|\Delta z|$ smaller than a desired tolerance” seems to be a good acceptance test for zeros of polynomials. If z is near the limit, we have $|z - r| \approx |\Delta z|$ and otherwise the error is at worst only a modest factor times $|\Delta z|$.

It should be borne in mind, however, that no rounding errors are considered here. Cancellation of figures may cause a much smaller $|\Delta z|$ and thus a far too optimistic error estimate. This may especially happen near already accepted, and removed, zeros.

Therefore, it is important to avoid circles around the accepted zeros during the iteration. This difficulty around the already accepted zeros is, in fact, the most serious drawback of any nondeflating method.

7. NUMERICAL EXPERIMENTS

The author did some experiments with an ALGOL 60 program for calculating the eigenvalues of a real matrix. Many ideas for this program have been taken from [2]. The main features of the program are :

a) The matrix is first transformed to Hessenberg form by means of Householder's transformation (see [5] and [7] p. 347).

b) For calculating f , f' and f'' Hyman's method is used (see [6] p. 327 and [7] p. 426).

c) The iteration formula used is Laguerre's formula (2.4), where p and q are chosen according to (5.2).

d) The iteration is continued until either $|\Delta z| < \text{norm} \times \text{eps}$, where eps is a given parameter and norm is the infinity norm of the matrix, or the number of iterations exceeds a given number.

Here, for the number of iterations, the program allows a higher maximum in case of convergence, i.e. in case $|\Delta z|$ is decreasing, then otherwise. (This strategy is inspired by J.W. Garwick's procedure "converge" [3] and an unpublished procedure for iteration control by R.J. De Vogelaere).

e) Iterates outside the circle around the origin with radius the infinity norm of the matrix are rejected and replaced by a suitable number on the edge of the circle. (This often saves an iteration from divergence).

f) Accepted zeros r_i are removed by subtracting $\sum (z - r_i)^{-k}$ from s_k ($k = 1, 2$). Moreover circles around these zeros with radius norm times a given parameter, eta , are avoided during the iteration, as long as $|\Delta z|$ is greater than 4 times this radius.

g) After accepting a zero (in fact either a real zero or a pair of complex conjugates), the program calculates a start for the next iteration by means of a Newton step (cf. Parlett [2] p. 473). The first start is $L(\infty)$ with $p = 1$, or in an earlier version, $L(0)$.

h) After accepting a non-real zero, the program accepts its complex conjugate without any checking.

The problem here is how to define non-reality. On the one hand, one has to prevent an accepted conjugate pair from causing a too-high multiplicity in a cluster of zeros, and on the other hand, one wants to deliver complex conjugates pair-wise. In this program the non-reality criterion is " $|\text{Im}(z)| > \text{norm} \times \text{eta}$ ".

i) In an earlier version the final value of each iteration was accepted as a p -tuple zero, where p is the last value used in Laguerre's formula. This, however, yields difficulties, because, especially in the last step, cancellation of figures may cause a useless value of s_1^2/s_2 and thus yield a wrong multiplicity. If this multiplicity turns out too high, it ruins the whole subsequent calculation. A more fundamental objection is the following. Because of the cubic convergence one may use a rather modest tolerance and expect a much higher precision for the last iterate. Thus one may accept the zero as a cluster of multiplicity p in the prescribed tolerance, but not in the higher precision expected. The newer version accepts each zero as a single one. In case of a multiple zero r the Newton step (cf. (g)) usually vanishes and thus the next iterate starts outside a circle around r with radius $\text{norm} \times \text{eta}$ (cf. (f)). In case r is multiple, the next iteration will again converge to r (this may be considered a numerical definition of multiple zeros).

The program was run on the Electrologica computers X1 and X8 by means of the ALGOL 60 systems of the Mathematical Centre, Amsterdam, the X1-system written by Dijkstra and Zonneveld and the X8-system by Kruseman Aretz. The author tried the following three versions of the program : the choice $q = n - p$, the choice $q = \min(n - p, p)$ (see 5.2) and the classical Laguerre formula, i.e. formula (2.4) with $p = 1$ and $q = n - 1$ (cf. [1] formula (7)). The test matrices used were, among others, Rosser's matrix of order 8, Frank's matrix of order 12, Eberlein's matrices of order 16 and 5 (see [2] and [8]), and 5 tenth order matrices (with linear elementary divisors) in Frobenius canonical form.

The experiments show that the choices $q = n - p$ and $q = \min(n - p, p)$ yield no significant difference in the reached precision and the required number of iterations. On the other hand these two versions required much less iterations for obtaining the same or higher precision for matrices having multiple eigenvalues, while the classical Laguerre formula was favourable for matrices having well-separated eigenvalues.

The text of the program together with some results of the version $q = n - p$ have been published in [9].

BIBLIOGRAPHY

- [1] E.N. LAGUERRE.— Oeuvres de Laguerre, Gauthier-Villars, Paris, Vol I, p. 87-103.
- [2] B. PARLETT.— Laguerre's Method applied to the Matrix Eigenvalue Problem, *MTAC*, 1964, 18, 464-485.
- [3] J.V. GARWICK.— Algorithm 1; *BIT*, 1961, 1, p. 64.
- [4] H.J. MAEHLY.— Zur iterativen Auflösung algebraischer Gleichungen *ZAMP*, 1954, 5, 261-263.
- [5] J.H. WILKINSON.— Householder's method for the Solution of the Algebraic Eigenproblem, *Comp. J.*, 1960, 3, 23-27.
- [6] J.H. WILKINSON.— Error Analysis of Floating point Computation, *Num. Mat.*, 1960, 2, 319-340.
- [6] J.H. WILKINSON.— *The Algebraic Eigenvalue Problem*, 1965, Oxford.
- [8] P.J. EBERLEIN.— A Jacobi-like method for the automatic computation of eigenvalues and eigenvectors of an arbitrary matrix, *J. SIAM*, 1962, 10, 74-88.
- [9] T.J. DEKKER.— *Newton-Laguerre iteration*, 1966, Report MR 82, Mathematical Centre, Amsterdam.

NOTE SUR DES BORNES POUR LES PLUS GRANDS ZEROS DE POLYNOMES PAR T.J. DEKKER

Les valeurs absolues des zéros d'un polynôme f de la forme :

$$f(z) = z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n$$

sont inférieures ou égales à ρv , où :

$$v = \max |a_k|^{1/k}$$

et ρ est la seule racine positive de l'équation :

$$\rho^n - \rho^{n-1} - \dots - \rho - 1 = 0$$

(v. [4] p. 19 et [5]).

Ceci résulte d'un Théorème de Cauchy :

Les valeurs absolues des zéros de f sont inférieures ou égales à la seule racine positive ξ de l'équation :

$$\xi^n - |a_1| \xi^{n-1} - \dots - |a_{n-1}| \xi - |a_n| = 0$$

(v. [1] Appendix 3).

De l'autre côté la plus grande valeur absolue $|r|_{\max}$ des zéros de f satisfait à :

$$|r|_{\max} \geq (2^{1/n} - 1) \xi .$$

M. Specht ([4] p. 18) attribue ce résultat à M. Birkhoff [2], qui cependant dérive seulement du résultat suivant lié :

$$\alpha \leq |r|_{\max} \leq \alpha / (2^{1/n} - 1) ,$$

où :

$$\alpha = \max_{k=1, \dots, n} \left(|a_k| / \binom{n}{k} \right)^{1/k} ,$$

remarquant que la borne inférieure α est déjà donnée en ([3] p. 20). Le nombre ρ est inférieur, et pour $n \gg 1$ presque égal, à deux ; le nombre α satisfait à $\alpha \geq v/n$. Alors en pratique on peut bien utiliser la relation simple :

$$v/n \leq |r|_{\max} \leq 2v ,$$

qui montre que la borne supérieure $2v$ n'est pas très pessimiste. Il ne semble pas facile de trouver de telles bornes assez voisines pour la plus grande valeur absolue des valeurs propres d'une matrice, ni pour le cas spécial d'une matrice presque triangulaire.

BIBLIOGRAPHIE

- [1] A. CAUCHY.— Cours d'Analyse de l'école royale polytechnique, 1ère partie : Analyse algébrique, 1821.
- [2] G.D. BIRKHOFF.— An elementary double inequality for the roots of an algebraic equation having greatest absolute value, *Bull. Amer. Math. Soc.*, 1914, 21, 494-495.

- [3] R.D. CARMICHAEL & T.E. MASON.— Note on the roots of algebraic equations
Bull. Amer. Math. Soc., 1914, 21, 14-22.
- [4] W. SPECHT.— Algebraische Gleichungen mit reellen oder komplexen Koeffizienten, *Enz. der Math. Wissenschaften*, 1958, Band I 1, Heft 3, Teil II 1-76.
- [5] J. ALBRECHT.— Zur simultanen Einkreisung sämtlicher Nullstellen eines polynoms, *ZAMM*, 1963, Band 43, Heft 7/8 377-379.