

Sitzungsberichte

der

mathematisch-naturwissenschaftlichen

Klasse

der

Bayerischen Akademie der Wissenschaften

zu München

Jahrgang 1951

München 1952

Verlag der Bayerischen Akademie der Wissenschaften

In Kommission bei der C. H. Beck'schen Verlagsbuchhandlung München

Zur numerischen Auflösung algebraischer Gleichungen

Von Fritz Wenzl in München

Vorgelegt von Herrn Josef Lense am 6. Juli 1951

I.

Ph. Furtwängler verwendete zur numerischen Auflösung von Gleichungen 4. Grades mit zwei Paaren konjugiert komplexer Wurzeln ein offenbar wenig bekannt gewordenes Verfahren, das an folgendem Beispiel erläutert sei:¹ Gegeben sei die Gleichung

$$(1.1) \quad F(x) = x^4 + 2,5x^3 + 10x^2 + 4x + 1.$$

Sie hat, wie man leicht erkennt, ein Paar „großer“ konjugiert komplexer Wurzeln (vom Betrag $R > 1$), zusammengefaßt in einer quadratischen Gleichung $G(x) = 0$, und ein Paar „kleiner“ konjugiert komplexer Wurzeln (vom Betrag $\frac{1}{R} < 1$), zusammengefaßt in einer quadratischen Gleichung $K(x) = 0$. Die Gleichung (1.1) ist im wesentlichen gelöst, wenn die Aufspaltung

$$(1.2) \quad F(x) = G(x) \cdot K(x)$$

in die Faktoren G und K hinreichend genau bekannt ist². Dazu kann man folgendermaßen vorgehen:

¹ Der Verfasser verdankt die Kenntnis des Verfahrens einer Vorlesung von Herrn Prof. Dr. Lense. Erst nachträglich wurde der Verfasser der vorliegenden Arbeit auf eine Abhandlung von B. Friedman aufmerksam (Comm. on Pure and Applied Math. 1949, S. 195), in der im wesentlichen derselbe Gegenstand wie hier behandelt ist. Da sich jedoch weder die Fragestellung noch der eingeschlagene Weg und die erzielten Ergebnisse decken, mag das Folgende trotzdem von Interesse sein. (Friedman beschränkt sich von vornherein auf den Fall hinreichend guter erster Näherungen und erzielt dabei auch sehr allgemeine Ergebnisse. Für Gleichungen 4. Grades mit 2 Paaren konjugiert komplexer Wurzeln werden jedoch auch unter der Voraussetzung hinreichend guter erster Näherungen die Friedmanschen Ergebnisse hier beträchtlich verschärft; vgl. dazu auch Anm. S. 88.)

² Der Koeffizient von x^2 sei in beiden Faktoren gleich 1.

1. Näherung K_1 für die „kleinen“ Wurzeln (worunter hier wie auch weiterhin immer die Wurzeln mit dem kleineren Betrag verstanden seien; entsprechend ist die Bezeichnung „große“ Wurzeln zu verstehen):

$$(1.3) \quad K_1 = x^2 + 0,4x + 0,1$$

(entstanden durch Vernachlässigung der Glieder 3. und 4. Grades in (1.1)). Daraus 1. Näherung G_1 für die „großen“ Wurzeln mittels Division unter Vernachlässigung des Restes:

$$(1.4) \quad F : K_1 \sim G_1 = x^2 + 2,1x + 9,06.$$

Daraus 2. Näherung K_2 für die „kleinen“ Wurzeln mittels angenäherter Division nach steigenden Potenzen von x , so daß

$$(1.5) \quad \underbrace{1 + 4x + \dots x^4}_F \sim \underbrace{(9,06 + 2,1x + x^2)}_{G_1} \cdot \underbrace{(0,1104 + 0,4159x + x^2)}_{K_2}.$$

(Übereinstimmung des konstanten Gliedes, des linearen Gliedes und des Gliedes mit x^4 auf beiden Seiten.)

Nun wieder mittels gewöhnlicher Division unter Vernachlässigung des Restes, ausgehend von K_2 , eine zweite Näherung G_2 für die „großen“ Wurzeln aus

$$(1.6) \quad F : K_2 \sim G_2 = x^2 + 2,084x + 9,023$$

usw. Ein Vergleich mit $G_1 = x^2 + 2,10x + 9,06$ zeigt, daß entsprechende Koeffizienten von G_1 und G_2 nur mehr um weniger als 1% voneinander abweichen.

II.

Das in I. dargestellte Verfahren weist eine gewisse Unsymmetrie zwischen den Näherungsdivisionen $F : K_v$ und den Näherungsdivisionen $F : G_v$ auf. Während nämlich z. B. bei der Division (1.4) der Quotient G_1 wie üblich so bestimmt wird, daß F und $K_1 \cdot G_1$ in den drei höchsten auftretenden Potenzen von x (d. h. in x^4 ; x^3 und x^2) übereinstimmen, wird in (1.5) der Quotient K_2 so berechnet, daß F und $G_1 \cdot K_2$ in den Koeffizienten von x^0 ,

x^1 und x^4 übereinstimmen. Es liegt nahe, den Schritt (1.5) und die entsprechenden späteren Näherungsdivisionen $F:G_v$ so abzuändern, daß F und $G_v \cdot K_{v+1}^*$ in den drei niedrigsten Potenzen von x (also in x^0 ; x^1 und x^2) übereinstimmen. Dieser Forderung entspricht das folgende symmetrisch gebaute Rechenchema (die Indizes sind hier abweichend von I. so gewählt, daß K_v ; G_{v+1} ; K_{v+2} usw. aufeinander folgen):

Gegeben sei:

$$(2.1) \quad F = x^4 + Ax^3 + Bx^2 + Cx + 1$$

und eine irgendwie gewonnene Näherung

$$(2.2) \quad K_v = B_v x^2 + C_v x + 1$$

für die „kleinen“ Wurzeln von $F = 0$. Dann erhält man eine Näherung G_{v+1} für die „großen“ Wurzeln aus

(2.2 a)

$$\underbrace{x^4 + Ax^3 + Bx^2 + \dots}_{F} \sim \frac{1}{B_v} \underbrace{(B_v x^2 + C_v x + 1)}_{K_v} \underbrace{(x^2 + A_{v+1} x + B_{v+1})}_{G_{v+1}}$$

mit

$$(2.3 a) \quad A_{v+1} = A - \frac{C_v}{B_v}; \quad B_{v+1} = B - A_{v+1} \cdot \frac{C_v}{B_v} - \frac{1}{B_v}$$

und daraus eine Näherung K_{v+2} für die „kleinen“ Wurzeln mittels

(2.2 b)

$$\underbrace{1 + Cx + Bx^2 + \dots}_{F} \sim \frac{1}{B_{v+1}} \underbrace{(B_{v+1} + A_{v+1}x + x^2)}_{G_{v+1}} \underbrace{(1 + C_{v+2}x + B_{v+2}x^2)}_{K_{v+2}}$$

mit

$$(2.3 b) \quad C_{v+2} = C - \frac{A_{v+1}}{B_{v+1}}; \quad B_{v+2} = B - C_{v+2} \cdot \frac{A_{v+1}}{B_{v+1}} - \frac{1}{B_{v+1}}$$

Die in den Rekursionsformeln (2.3 a) bzw. (2.3 b) auftretenden Quotienten und Produkte lassen sich jeweils mit einer Rechenschiebereinstellung ermitteln.

* Bis auf einen Proportionalitätsfaktor, der für die Nullstellen nichts ausmacht.

Rechnet man das eben dargestellte Verfahren für das Beispiel (1.1) durch, so ergibt sich, ausgehend von

$$(2.4) \left\{ \begin{array}{l} K_1 = 1 + 4x \quad + 10x^2 \\ \quad \quad \quad \text{zunächst} \quad G_2 = x^2 + 2,1x + 9,06 \quad \text{und dann} \\ K_3 = 1 + 3,768x \quad + 9,016x^2; \\ \quad \quad \quad \quad \quad \quad G_4 = x^2 + 2,08206x \quad + 9,0189; \\ K_5 = 1 + 3,769145x + 9,018996x^2; \\ \quad \quad \quad \quad \quad \quad G_6 = x^2 + 2,082088x + 9,018994. \end{array} \right.$$

Der Vergleich von G_4 mit G_6 läßt vermuten, daß bereits nach 3 Divisionen, d. h. bei G_4 , die Koeffizienten in den ersten 4 Dezimalen festliegen. Einen Anhaltspunkt für die Güte der Näherung gibt die Übereinstimmung bzw. Differenz des Koeffizienten von x^2 in K_v und des konstanten Gliedes in G_{v+1} bzw. G_{v-1} (falls die sämtlichen Näherungspolynome entsprechend den Formeln (2.2a) und (2.2b) normiert sind).

Die in (2.1) enthaltene Voraussetzung, daß F das konstante Glied 1 besitzt, stellt keine wesentliche Einschränkung der Allgemeinheit dar, denn jede Gleichung $\xi^4 + a\xi^3 + b\xi^2 + c\xi + d = 0$ mit $d > 0$ läßt sich durch $\xi = \sqrt[4]{d} \cdot x$ auf die Form (2.1) bringen (für den zunächst betrachteten Fall zweier Paare konjugiert komplexer Wurzeln ist immer $d > 0$). Natürlich kann das durch (2.3a) und (2.3b) dargestellte Verfahren sukzessiver Divisionen sinngemäß auch auf den Fall $d \neq 1$ unmittelbar (ohne vorherige Transformation der Unbekannten) angewendet werden. Für die folgenden Konvergenzbetrachtungen genügt es jedoch wegen der Möglichkeit der Transformation, wenn wir uns von vornherein auf den rechnerisch besonders bequemen Fall $d = 1$, d. h. auf den Fall der normierten Gleichung (2.1) beschränken.

III. (Konvergenzbetrachtung)

Das Polynom (2.1) enthalte die quadratischen Faktoren

$$(3.1) \quad G(x) = x^2 + \mu \cdot Rx + R^2 \quad \text{und} \quad K(x) = R^2 x^2 + xRx + 1.$$

Nach (2.1) ergeben sich dann die Koeffizienten von

$$F(x) = G(x) \cdot K(x) \cdot R^{-2} \text{ aus}$$

$$(3.2) \quad A = \mu R + \varkappa R^{-1}; B = R^2 + \mu \varkappa + R^{-2}; C = \varkappa R + \mu R^{-1}.$$

Für die Koeffizienten B_v und C_v der Näherung (2.2) gelte unter Verwendung zweier Korrekturgrößen β_v und γ_v die Darstellung

$$(3.3) \quad B_v = R^2(1 + \beta_v); \quad C_v = R(\varkappa + \gamma_v).$$

Dann folgt nach (2.3 a) unter Berücksichtigung von (3.2) und (3.3) zunächst

$$A_{v+1} = \mu R + \varkappa R^{-1} - \frac{\varkappa + \gamma_v}{1 + \beta_v} \cdot R^{-1} \quad \text{oder}$$

$$(3.4a) \quad A_{v+1} = R(\mu + \alpha_{v+1}) \text{ mit } \alpha_{v+1} = \frac{\varkappa \beta_v - \gamma_v}{1 + \beta_v} \cdot R^{-2}.$$

Ebenso folgt nach (2.3 a) unter Berücksichtigung von (3.2), (3.3) und (3.4 a) die Gleichung

$$B_{v+1} = R^2 + \mu \varkappa + R^{-2} - R(\mu + \alpha_{v+1}) \cdot \frac{\varkappa + \gamma_v}{1 + \beta_v} \cdot R^{-1} - \frac{1}{1 + \beta_v} R^{-2}$$

oder

(3.4b)

$$B_{v+1} = R^2(1 + \beta_{v+1}) \text{ mit } \beta_{v+1} = \alpha_{v+1} \left(\mu - \frac{\varkappa + \gamma_v}{1 + \beta_v} R^{-2} \right) + \frac{\beta_v R^{-4}}{1 + \beta_v}.$$

Bei zwei Paaren konjugiert komplexer Wurzeln der Gleichung $F(x) = 0$ stellt der Faktor

$$(3.5) \quad \rho = R^{-2}, \quad |\rho| < 1 \quad \text{vorausgesetzt,}$$

wegen (3.1) das Verhältnis des Betrags einer „kleinen“ Wurzel zum Betrag einer „großen“ Wurzel dar. Benutzt man (3.5) zur Abkürzung, so erhält man aus (3.4a) und (3.4b) die Rekursionsformeln

$$(3.6) \quad \alpha_{v+1} = \rho \cdot \frac{\varkappa \beta_v - \gamma_v}{1 + \beta_v}; \quad \beta_{v+1} = \alpha_{v+1} \left(\mu - \rho \cdot \frac{\varkappa + \gamma_v}{1 + \beta_v} \right) + \rho^2 \cdot \frac{\beta_v}{1 + \beta_v}$$

und daher für γ_{v+2} entsprechend der Gleichung für α_{v+1} (Ver-

tauschung von α mit γ und x mit μ ; v durch $v+1$ ersetzt) die Rekursionsformel

$$(3.7) \quad \gamma_{v+2} = \rho \cdot \frac{\mu \beta_{v+1} - \alpha_{v+1}}{1 + \beta_{v+1}} = \frac{\rho \cdot \alpha_{v+1} \left(\mu^2 - 1 - \rho \mu \cdot \frac{x + \gamma_v}{1 + \beta_v} \right) + \rho^3 \mu \frac{\beta_v}{1 + \beta_v}}{1 + \beta_{v+1}}$$

(Bedeutung der Korrekturgröße γ_{v+2} entsprechend derjenigen von γ_v in (3.3); im Zähler der rechten Seite wurde β_{v+1} mit Hilfe von (3.6) eliminiert).

Folgerungen: a) Unter Vernachlässigung höherer Potenzen der Korrekturgrößen β_v und γ_v folgt aus (3.6) bzw. (3.7) zunächst

$$(3.8) \quad \beta_{v+1} = \alpha_{v+1} (\mu - \rho x) + \rho^2 \beta_v; \quad \gamma_{v+2} = \rho \alpha_{v+1} (\mu^2 - 1 - \rho \mu x) + \rho^3 \mu \beta_v$$

und ganz entsprechend

$$(3.9) \quad \beta_{v+2} = \gamma_{v+2} (x - \rho \mu) + \rho^2 \beta_{v+1}; \quad \alpha_{v+3} = \rho \gamma_{v+2} (x^2 - 1 - \rho \mu x) + \rho^3 x \beta_{v+1}.$$

Durch Einsetzen von (3.8) in (3.9) könnte man nun β_{v+2} und α_{v+3} unmittelbar aus β_v und α_{v+1} erhalten in der Gestalt

$$(3.9a) \quad \begin{aligned} \beta_{v+2} &= c_{11} \beta_v + c_{12} \alpha_{v+1}; \\ \alpha_{v+3} &= c_{21} \beta_v + c_{22} \alpha_{v+1}, \end{aligned}$$

wobei die c_{ik} nur noch von μ ; x und ρ abhängen. Daraus folgt sicher $\beta_v \rightarrow 0$ und $\alpha_v \rightarrow 0$ mit $v \rightarrow \infty$, wenn sich für ein geeignetes positives λ aus $|\beta_v| < \varepsilon$; $|\alpha_{v+1}| < \lambda \varepsilon$ immer $|\beta_{v+2}| < q \varepsilon$ und $|\alpha_{v+3}| < q \cdot \lambda \varepsilon$ mit einem festen Faktor $q < 1$ ergibt. Dies ist sicher der Fall, wenn zugleich

$$\begin{aligned} |c_{11}| + \lambda |c_{12}| &< q, \text{ also } \lambda < \frac{q - |c_{11}|}{|c_{12}|} \text{ und} \\ |c_{21}| + \lambda |c_{22}| &< q \lambda, \text{ also } \lambda > \frac{|c_{21}|}{q - |c_{22}|} \text{ mit } |c_{22}| < q \end{aligned}$$

für ein geeignetes positives λ gilt, d. h. für denjenigen Bereich der Größen μ ; x und ρ , für den die Ungleichung

$$(3.10) \quad \frac{1 - |c_{11}|}{|c_{12}|} > \frac{|c_{21}|}{1 - |c_{22}|} \text{ mit } |c_{22}| < 1$$

gilt (der Grenzfall, daß $c_{12} = 0$; $|c_{11}| < 1$; $|c_{22}| < 1$, sei dabei sinngemäß eingeschlossen).

In diesem Bereich konvergiert bei hinreichend kleinen Beträgen von β_1 und α_2 das angegebene Verfahren immer (vgl. hierzu auch die Bemerkung vor Satz 1). Wir beschränken uns jedoch für das Weitere auf die folgenden Abschätzungen: Es sei

$$(3.11) \quad |\varkappa| < 2; \quad |\mu| < 2$$

(was bei reellen \varkappa ; μ und R nach (3.1) nichts anderes bedeutet als zwei Paare konjugiert komplexer – nicht reeller – Nullstellen von F). Ferner sei

$$(3.12) \quad |\beta_v| < 2\varepsilon; \quad |\alpha_{v+1}| < 3\rho \cdot \varepsilon.$$

Aus (3.8) folgt dann (\varkappa ; μ und ρ weiterhin reell vorausgesetzt)

$$(3.13) \quad |\beta_{v+1}| < 2\varepsilon(3\rho + 4\rho^2); \quad |\gamma_{v+2}| < 3\rho\varepsilon \cdot \left(3\rho + \frac{16}{3}\rho^2\right).$$

Beide Klammern sind kleiner als 0,99 für $\rho < \frac{1}{4,3}$. Durch Vergleich mit (3.12) ergeben sich dann, indem man die entsprechende Rechnung unter Berücksichtigung von (3.9) statt (3.8) nochmal durchführt, die Ungleichungen

$$|\beta_{v+2}| < 2\varepsilon \cdot 0,99^2; \quad |\alpha_{v+3}| < 3\rho\varepsilon \cdot 0,99^2.$$

Bei Verwendung der exakten Rekursionsformeln (3.6) und (3.7) hätte man statt dessen

(3.13 a)

$$|\beta_{v+2}| < 2\varepsilon \cdot 0,99^2 + M\varepsilon^2; \quad |\alpha_{v+3}| < 3\rho\varepsilon \cdot 0,99^2 + M\rho\varepsilon^2$$

mit einer von ε unabhängigen Schranke M (wenn nur ε klein genug ist, daß man innerhalb des Konvergenzbereichs der Reihenentwicklungen bleibt, von denen (3.8) und (3.9) die ersten Näherungen geben). Aus (3.13 a) folgt nun aber für jedes $\rho < \frac{1}{4,3}$ bei hinreichend kleinem ε sofort

$$(3.13 b) \quad |\beta_{v+2}| \leq q \cdot 2\varepsilon \text{ und } |\alpha_{v+3}| \leq q \cdot 3\rho\varepsilon \text{ mit } (0 <) q < 1.$$

Vergleicht man dies wieder mit (3.12), so erhält man folgenden

Satz 1: Das in II. zusammengestellte Rechenschema konvergiert bei hinreichend guten Anfangswerten für die Koeffizienten des Faktorpolynoms $K(x)$ immer dann, wenn das gegebene Polynom 4. Grades $F(x)$ zwei Paare konjugiert komplexer Wurzeln besitzt, deren Beträge sich mindestens wie 4,3 : 1 verhalten.¹

Der Vergleich von (3.13) mit (3.12) gibt besonders für kleinere ρ zugleich auch eine Abschätzung für die Güte der Konvergenz.

Unter der zusätzlichen Annahme $\mu\kappa > 0$ (was bei reellen μ ; κ und R nach (3.1) zwei Paare konjugiert komplexer Wurzeln bedeutet, deren Realteile dasselbe Vorzeichen haben) genügen wesentlich schwächere Voraussetzungen für ρ , wie sich folgendermaßen ergibt:

Man kann ohne Einschränkung der Allgemeinheit $|\mu| < |\kappa|$ annehmen. Ferner läßt sich der Fall negativer μ und κ nach (3.8) und (3.9) auf denjenigen positiver μ und κ zurückführen, indem man die β_v durch $(-\beta_v)$ ersetzt. Dementsprechend sei zunächst (1. Fall)

$$(3.14) \quad 1,6 < \mu < \kappa < 2 \quad \text{und}$$

$$(3.14a) \quad (0 <) \rho < 0,5.$$

Dann ist

$$\mu - \rho\kappa \begin{cases} < \kappa(1 - \rho) < 2 \cdot (1 - \rho) \\ > 1,6 - 2\rho > 0,6 \quad \text{nach (3.14a),} \end{cases}$$

also

$$(3.14b) \quad 0,6 < \mu - \rho\kappa < 2 \cdot (1 - \rho)$$

und daher

$$\mu^2 - 1 - \rho\kappa\mu = \mu \cdot (\mu - \rho\kappa) - 1 \begin{cases} < 2 \cdot 2 \cdot (1 - \rho) - 1 = 3 - 4\rho \\ > 1,6 \cdot 0,6 - 1 = -0,04, \end{cases}$$

¹ Dieses Ergebnis ist bereits schärfer als die entsprechende Aussage bei Friedman, wo (S. 205) $\rho < 1:5$ und dazu noch gleiches Vorzeichen des Realteils bei allen Wurzeln verlangt wird (allerdings mit der Zusatzbemerkung, daß diese hinreichende Bedingung zweifellos noch durch eine bedeutend schwächere ersetzt werden könnte).

also, weil $3 - 4\rho > 1$ ($> 0,04$) nach (3.14a),

$$(3.14c) \quad |\mu^2 - 1 - \rho\mu x| < 3 - 4\rho.$$

Aus (3.8), (3.12) und (3.14b) folgt dann

$$(3.15a) \quad |\beta_{v+1}| < 3\rho\varepsilon \cdot 2(1-\rho) + \rho^2 \cdot 2\varepsilon = 2\varepsilon(3\rho - 2\rho^2) < 2\varepsilon,$$

$$\text{weil } \frac{d}{d\rho}(3\rho - 2\rho^2) = 3 - 4\rho > 0$$

nach (3.14a) und weil daher $3\rho - 2\rho^2$ im abgeschlossenen Bereich (3.14a)¹ seine obere Grenze für $\rho = 0,5$ annimmt.

Aus (3.8), (3.12), (3.14) und (3.14c) folgt ferner

$$(3.15b) \quad \begin{aligned} |\gamma_{v+2}| &< 3\rho^2\varepsilon \cdot (3 - 4\rho) + 2\rho^3 \cdot 2\varepsilon \\ &= 3\rho\varepsilon \left(3\rho - \frac{8}{3}\rho^2 \right) < \frac{5\rho\varepsilon}{2}, \end{aligned}$$

$$\text{weil } \frac{d}{d\rho} \left(3\rho - \frac{8}{3}\rho^2 \right) = 3 - \frac{16}{3}\rho > 0$$

nach (3.14a) und weil daher $\left(3\rho - \frac{8}{3}\rho^2 \right)$ im abgeschlossenen Bereich (3.14a) seine obere Grenze für $\rho = 0,5$ erreicht.

Nun folgt aber aus (3.14) und (3.14a) auch noch

(3.15c)

$$x - \rho\mu \begin{cases} < 2 - 1,6\rho \\ > x - \rho x > 1,6(1 - \rho) > 0,8 \end{cases}$$

und

$$x^2 - 1 - \rho\mu x \begin{cases} > x^2 - 1 - \rho x^2 = x^2(1 - \rho) - 1 > 1,6^2 \cdot 0,5 - 1 > 0 \\ < x^2 - 1 - \rho \cdot 1,6x < 3 - 3,2\rho, \end{cases}$$

weil $x^2 - 1 - \rho \cdot 1,6x$ nach (3.14a) für $1,6 < x < 2$ monoton mit x wächst und dabei seine obere Grenze für $x = 2$ annimmt.

Also ist

$$(3.15d) \quad 0 < x^2 - 1 - \rho\mu x < 3 - 3,2\rho.$$

Aus (3.9), (3.15b), (3.15c) und (3.15a) folgt dann

$$(3.15e) \quad |\beta_{v+2}| < \frac{5\rho\varepsilon}{2} (2 - 1,6\rho) + \rho^2 \cdot 2\varepsilon = 2\varepsilon \left(\frac{5\rho}{2} - \rho^2 \right) < 2\varepsilon$$

¹ d. h. im Bereich $0 \leq \rho \leq 0,5$.

(nach (3.14a), ähnlich wie in der Begründung zu (3.15a) und (3.15b)); ferner ebenso unter Berücksichtigung von (3.9), (3.15a), (3.15b), (3.15d) und (3.14)

$$(3.15f) \quad |\alpha_{v+3}| < \frac{5\rho^2\varepsilon}{2} \cdot (3 - 3,2\rho) + \rho^3 \cdot 4\varepsilon = \rho\varepsilon \left(\frac{15}{2} \rho - 4\rho^2 \right) < \frac{11}{4} \rho\varepsilon.$$

Der Vergleich von (3.15e) und (3.15f) mit (3.12) gibt wieder die Konvergenz (hinreichend kleine Beträge von β_1 und α_2 vorausgesetzt). Nun sei (2. Fall)

$$(3.16) \quad 0 < \mu < 1,6; \quad \mu < \varkappa < 2 \quad \text{und}$$

$$(3.16a) \quad (0 <) \rho < 0,44.$$

Dann ist (α) $\rho(\mu - \rho\varkappa) < \rho(\mu - \rho\mu) < \rho(1 - \rho) \cdot 1,6 < 0,4$, weil für $0 < \rho < 0,5$, also erst recht für (3.16a) immer

$$\frac{d}{d\rho} (\rho - \rho^2) = 1 - 2\rho > 0, \text{ also } \rho(1 - \rho) < \frac{1}{4}.$$

Außerdem ist dann aber

$$(\beta) \quad \rho(\rho\varkappa - \mu) < \rho^2\varkappa < \rho^2 \cdot 2 < 0,44^2 \cdot 2 < 0,4,$$

also, (α) und (β) zusammengefaßt,

$$(3.16b) \quad \rho|\mu - \rho\varkappa| < 0,4.$$

Schließlich folgt aus (3.16) und (3.16a) auch noch

$$(3.16c) \quad \mu^2 - 1 - \rho\mu\varkappa < \mu^2 - 1 - \rho\mu^2 = \mu^2(1 - \rho) - 1 < 1,56 - 2,56\rho;$$

$$\text{und} \quad 1 + \rho\mu\varkappa - \mu^2 < 1 + \rho\mu \cdot 2 - \mu^2.$$

Aus (3.8), (3.12) und (3.16b) folgt nun

$$(3.17a) \quad |\beta_{v+1}| < 3\varepsilon \cdot 0,4 + \rho^2 \cdot 2\varepsilon < 1,6\varepsilon \quad \text{nach (3.16a).}$$

Aus (3.8), (3.12), (3.16c) und (3.16) folgt im Fall $0 < \mu^2 - 1 - \rho\mu\varkappa$

$$\begin{aligned} |\gamma_{v+2}| &< 3\rho^2\varepsilon \cdot (1,56 - 2,56\rho) + \rho^3 \cdot 1,6 \cdot 2\varepsilon \\ &< 3\rho\varepsilon \cdot (1,56\rho - 1,49\rho^2) \\ &< 3\rho\varepsilon \cdot (1,56 \cdot 0,44 - 1,49 \cdot 0,44^2)^1 < 3\rho\varepsilon \cdot 0,4, \end{aligned}$$

¹ Nach (3.16a) wegen $\frac{d}{d\rho} (1,56\rho - 1,49\rho^2) > 0$ für $\rho < 0,44$.

dagegen im Fall $0 > \mu^2 - 1 - \rho\mu\kappa$

$$\begin{aligned} |\gamma_{v+2}| &< 3\rho^2\varepsilon \cdot (1 + 2\rho\mu - \mu^2) + \rho^3 \cdot 2\mu\varepsilon \\ &< 3\rho^2\varepsilon \cdot \left[1 - \mu^2 + 2\mu \left(1 + \frac{1}{3} \right) \cdot 0,44 \right] \\ &< 3\rho^2\varepsilon \cdot \left[1 + \frac{1}{4} \cdot \left(\frac{8}{3} \cdot 0,44 \right)^2 \right], \end{aligned}$$

weil Max. $[1 - \mu^2 + c\mu] = 1 + \frac{c^2}{4}$ für festes c ; damit hat man in diesem Fall (mit $c = \frac{8}{3} \cdot 0,44 < 1,2$) die Ungleichung

$$(3.17b) \quad |\gamma_{v+2}| < 3\rho \cdot 0,44\varepsilon \cdot \left[1 + \frac{1,2^2}{4} \right] = 3\rho \cdot 0,44\varepsilon \cdot 1,36 < 1,8\rho\varepsilon.$$

Nun folgt aus (3.16) und (3.16a) aber auch noch

$$(3.17c) \quad 0 < \kappa - \rho\mu < 2 \quad \text{und}$$

$$(3.17d) \quad \kappa^2 - 1 - \rho\kappa\mu \begin{cases} < \kappa^2 - 1 < 3 \\ > \kappa^2 - 1 - \rho\kappa^2 > -1. \end{cases}$$

Aus (3.9), (3.17b), (3.17c) und (3.17a) folgt dann

$$(3.17e) \quad |\beta_{v+2}| < 1,8\rho\varepsilon \cdot 2 + \rho^2 \cdot 1,6\varepsilon < 1,9\varepsilon \quad \text{nach (3.16a)}$$

und unter Berücksichtigung von (3.17d) und (3.16)

$$(3.17f) \quad |\alpha_{v+3}| < 1,8\rho^2\varepsilon \cdot 3 + \rho^3 \cdot 2 \cdot 1,6\varepsilon < 3\rho\varepsilon \quad \text{nach (3.16a)}.$$

Der Vergleich mit (3.12) ergibt wieder (unter der Voraussetzung hinreichend kleiner Beträge von β_1 und α_2) die Konvergenz.

Auf Grund der Bemerkungen vor (3.14) lassen sich nun aber die Ergebnisse der Abschätzungen zu (3.14) und (3.16) zusammenfassen in dem folgenden

Satz 2: Das in II. zusammengestellte Rechenschema konvergiert bei hinreichend guten Anfangswerten für die Koeffizienten des Faktorpolynoms $K(x)$ immer dann, wenn das gegebene Polynom 4. Grades $F(x)$ zwei Paare konjugiert komplexer Wurzeln besitzt, deren Beträge sich mindestens wie $(1 : 0,44 \approx) 2,3 : 1$ verhalten und deren Realteile dasselbe Vorzeichen haben (Vgl. auch hierzu Anm. S. 88, d. h. die Friedmansche Bedingung $\frac{1}{\rho} > 1 : 0,2 = 5 : 1$).

Aus (3.2) und (3.5) folgt $\rho + \frac{1}{\rho} = B - \mu\kappa$ mit $\frac{1}{\rho} > 1$; $0 < \rho < 1$ (reelles R wie bisher vorausgesetzt). Aus der Voraussetzung

$$(3.18) \quad B > 6,5$$

folgt daher

$$\text{für } \mu\kappa < 4 \quad \text{sofort} \quad \rho + \frac{1}{\rho} > 2,5; \quad \rho < 0,5 \quad (\alpha),$$

$$\text{für } \mu\kappa < 3,2 \quad \text{dagegen} \quad \rho + \frac{1}{\rho} > 3,3; \quad \rho < 0,4 \quad (\beta),$$

$$\text{für } \mu\kappa < 0 \quad \text{sogar} \quad \rho + \frac{1}{\rho} > 6,5; \quad \rho < 0,2 \quad (\gamma).$$

Wir unterscheiden nun wieder die Fälle

$$(3.14) \quad 1,6 < \mu < \kappa < 2 \quad \text{und}$$

$$(3.16) \quad 0 < \mu < 1,6; \quad \mu < \kappa < 2.$$

Aus (3.14) und (3.18) folgt dann nach (α) sofort $\rho < 0,5$, d. h. (3.14a), was („1. Fall“) zusammen mit (3.14) für den Konvergenzbeweis genügt. Aus (3.16) und (3.18) folgt nach (β) dagegen sogar $\rho < 0,4$, also erst recht (3.16a), d. h. $\rho < 0,44$, was wieder („2. Fall“) zusammen mit (3.16) die Konvergenz ergab.

Die Fälle (3.14) und (3.16) zusammengefaßt besagen aber nichts anderes als die Voraussetzung

$$0 < \mu < \kappa < 2.$$

Dies führt also zusammen mit (3.18) in jedem Fall zur Konvergenz. Nach der Bemerkung vor (3.14), S. 88, genügt dazu nun aber auch (3.18) und

$$0 < |\mu| < 2; \quad 0 < |\kappa| < 2, \quad \text{solange nur } \mu\kappa > 0.^1$$

(3.18) und $\mu\kappa < 0$ gibt jedoch nach (γ) zunächst $\rho < 0,2$ und daher erst recht $\rho < \frac{1}{4,3}$, d. h. die Voraussetzung von Satz 1, wenn man an 2 Paaren konjugiert komplexer Wurzeln festhält. Daher genügt (3.18) unter dieser Voraussetzung (2 Paare konj. kompl. Wurzeln) immer und es gilt als zusammenfassende Abschätzung der

¹ (μ sowie κ reell vorausgesetzt).

Satz 3: Das in II. zusammengestellte Rechenschema konvergiert bei hinreichend guten Anfangswerten und 2 Paaren konjugiert komplexer Wurzeln immer, wenn in (2.1) $B > 6,5$ ist.

Geht man statt von (2.1) von der allgemeinen Gleichung 4. Grades $\xi^4 + a\xi^3 + b\xi^2 + c\xi + d = 0$ aus, so lautet die entsprechende (hinreichende) Konvergenzbedingung $b : \sqrt{d} > 6,5$ (vgl. letzten Abschnitt vor III.).

Ebenso wie Satz 3 bloß eine hinreichende Konvergenzbedingung gibt, bedeutet z. B. auch $\rho < 0,44$ nur eine hinreichende, aber keine notwendige Bedingung für die Gültigkeit des Satzes 2 bei sonst unveränderten Voraussetzungen. Daß jedoch das in II. besprochene Verfahren nicht immer brauchbar ist, zeigt das Beispiel

$$\mu = 0; \quad \kappa = 2.$$

Aus (3.8) bzw. (3.9) folgt dann nämlich

- (a) $\beta_{v+1} = -2\rho\alpha_{v+1} + \rho^2\beta_v; \quad \gamma_{v+2} = -\rho\alpha_{v+1}$ und daher
 (b) $\beta_{v+2} = -2\rho\alpha_{v+1}(1 + \rho^2) + \rho^4\beta_v; \quad \alpha_{v+3} = -\rho^2\alpha_{v+1}(3 + 4\rho^2) + 2\rho^5\beta_v.$

Nun sei

$$(c) \quad \beta_v \alpha_{v+1} > 0 \quad \text{und}$$

$$(d) \quad |\beta_v| < 2\rho(1 + \rho^2) \cdot |\alpha_{v+1}|.$$

Dann ist nach (3.5), d. h. wegen $(0 <)\rho < 1$, erst recht

$$(e) \quad \rho^4 |\beta_v| < 2\rho(1 + \rho^2) \cdot |\alpha_{v+1}|,$$

außerdem aber (wieder wegen (d))

$$(f) \quad 2\rho^5 |\beta_v| < \rho^2(4\rho^4 + 4\rho^6) \cdot |\alpha_{v+1}| < \rho^2(3 + 4\rho^2) \cdot |\alpha_{v+1}|$$

$$\text{für} \quad \rho^2 \leq 0,9$$

(es ist nämlich $4\rho^4 + 4\rho^6 < 3 + 4\rho^2$, falls $\rho^2 + \rho^4 < \frac{3}{4\rho^2} + 1$, was für $\rho^2 = 0,9$ und damit erst recht für $\rho^2 < 0,9$ erfüllt ist).

Nach (e) und (f) ist für die Vorzeichen von β_{v+2} und α_{v+3} wegen (b) das Vorzeichen von $(-\alpha_{v+1})$ maßgebend. Es gilt also

$$(g) \quad \beta_{v+2} \alpha_{v+3} > 0 \quad \text{sowie}$$

(wegen (b) und (c))

$$(h) \quad |\beta_{v+2}| < 2\rho(1 + \rho^2) \cdot |\alpha_{v+1}| < 2\rho(1 + \rho^2) \cdot |\alpha_{v+3}|, \\ \text{wenn nur } |\alpha_{v+1}| < |\alpha_{v+3}|.$$

Nun ist aber nach (b), (c) und (f) sicher

$$(i) \quad |\alpha_{v+3}| > |\alpha_{v+1}| \cdot [(3\rho^2 + 4\rho^4) - 2\rho^5(2\rho + 2\rho^3)] > |\alpha_{v+1}| \\ \text{für} \quad 0,27 < \rho^2 < 0,8$$

(es ist nämlich $3\rho^2 + 4\rho^4 - 4\rho^6 - 4\rho^8 > 1$, falls $3\rho^2 + 4\rho^4 - 1 > 4\rho^6(1 + \rho^2)$, d. h. (Division durch $1 + \rho^2$), falls $4\rho^2 - 1 > 4\rho^6$ (erfüllt für $0,27 < \rho^2 < 0,8$)).

Unter den für (f) und (i) benötigten (miteinander verträglichen) Voraussetzungen für ρ^2 ist also nach (i) der Betrag von α_{v+3} größer als derjenige von α_{v+1} . Entsprechendes folgt durch Vergleich von (g) und (h) mit (c) und (d) mittels vollständiger Induktion dann aber auch für alle größeren v , solange man die für Korrekturgrößen geringen Betrags gültigen Näherungen (3.8) und (3.9) benutzen kann. Dabei konnten β_v und α_{v+1} dem Betrag nach beliebig klein sein, wenn nur die Ungleichungen (c) und (d) galten.

b) Setzt man als erste Näherung (indem man x^4 und x^3 in (2.1) vernachlässigt)

$$(3.19) \quad K_1 = Bx^2 + Cx + 1,$$

so folgt aus (3.2) und (3.3) $\gamma_1 = \mu R^{-2}$ und $\beta_1 = \mu \varkappa R^{-2} + R^{-4}$. Nach (3.5) bedeutet dies

$$(3.20) \quad \gamma_1 = \mu\rho; \beta_1 = \varkappa\mu\rho + \rho^2 = \varkappa\gamma_1 + \rho^2, \text{ also } \alpha_2 = \frac{(\varkappa^2 - 1) \cdot \mu\rho^2 + \varkappa\rho^3}{1 + \varkappa\mu\rho + \rho^2} \\ \text{nach (3.6).}$$

Wegen der Symmetrie der gestellten Aufgabe und der in II. zusammengestellten Rekursionsformeln bezüglich der Faktorpolynome G und K kann man sich auch hier wieder auf die Annahme $|\mu| < |\varkappa|$ beschränken; im Fall $|\varkappa| < |\mu|$ würde man nämlich zweckmäßigerweise sowieso statt mit einer ersten Näherung für

K mit einer ersten Näherung für G beginnen, indem man $Cx + 1$ statt $x^4 + Ax^3$ in $F(x)$ wegläßt, da dann (bei unverändertem β_1) an Stelle der Korrekturgröße $\gamma_1 = \mu\rho$ die absolut genommen kleinere Korrekturgröße $\alpha_1 = \kappa\rho$ tritt.

Darüber hinaus kann man sich auch hier wieder auf $\mu > 0$ beschränken; denn ersetzt man in (3.20) zugleich μ durch $(-\mu)$ und κ durch $(-\kappa)$, so wechseln γ_1 und α_2 ihr Vorzeichen, während β_1 unverändert bleibt; andererseits bleiben aber die Rekursionsformeln (3.6) und (3.7) erhalten, wenn man bei μ und κ sowie bei den Korrekturgrößen γ und α das Vorzeichen wechselt, während man die β unverändert läßt.

Zunächst sei nun (mit später genauer zu bestimmenden μ_0 und ρ_0)

$$(3.21) \quad 1 < \mu_0 < \mu < \kappa < 2 \quad \text{und} \quad (0 <) \rho < \rho_0 (< 1).$$

Nach (3.20) ist dann einerseits

$$(3.22) \quad 2\beta_1 > \gamma_1 > 0 \quad \text{und} \quad \alpha_2 > 0.$$

Andererseits folgt aus der Annahme, für ein bestimmtes v sei

$$(3.23) \quad 2\beta_v > \gamma_v > 0 \quad \text{und} \quad \alpha_{v+1} > 0,$$

wegen $\kappa < 2$ zunächst

$$(3.23a) \quad \frac{\kappa + \gamma_v}{1 + \beta_v} < \frac{2 + 2\beta_v}{1 + \beta_v} = 2.$$

Nach (3.6), (3.21) und (3.23) folgt daraus aber

$$(3.23b) \quad \beta_{v+1} > \alpha_{v+1} \cdot (\mu_0 - \rho_0 \cdot 2) > \frac{\alpha_{v+1}}{2}, \quad \text{falls nur } \mu_0 - 2\rho_0 > \frac{1}{2}.$$

Ferner ergibt sich aus (3.6) und (3.23) die Ungleichung

$$\alpha_{v+1} < \rho \cdot \frac{\kappa \beta_v}{1 + \beta_v}, \quad \text{also}$$

$$\frac{\rho \cdot \beta_v}{1 + \beta_v} > \frac{\alpha_{v+1}}{\kappa} > \frac{\alpha_{v+1}}{2} \quad \text{nach (3.21)}$$

und damit

$$(3.23c) \quad \frac{\rho^3 \mu \beta_v}{1 + \beta_v} > \frac{\rho^2 \cdot \mu \alpha_{v+1}}{2}.$$

Aus (3.23 a) und (3.21) folgt andererseits

$$(3.23 d) \quad \mu^2 - 1 - \rho\mu \cdot \frac{\kappa + \gamma_v}{1 + \beta_v} > \mu^2 - 1 - \rho\mu \cdot 2.$$

Nach (3.23) und (3.23 b) ist nun $\alpha_{v+1} > 0$ und $\beta_{v+1} > 0$.¹ Aus (3.7) ergibt sich dann wegen (3.23 d) und (3.23 c) schließlich

$$\begin{aligned} \gamma_{v+2} &> \frac{1}{1 + \beta_{v+1}} \cdot \left[\rho \alpha_{v+1} (\mu^2 - 1 - \rho\mu \cdot 2) + \rho^2 \cdot \mu \frac{\alpha_{v+1}}{2} \right] \\ &= \frac{\rho \alpha_{v+1}}{1 + \beta_{v+1}} \cdot \left(\mu^2 - 1 - \rho\mu \cdot \frac{3}{2} \right) \end{aligned}$$

und damit nach (3.21) und (3.23) sicher

$$(3.23 e) \quad \gamma_{v+2} > 0, \text{ wenn nur } \left(\mu_0^2 - 1 - \rho_0 \mu_0 \cdot \frac{3}{2} \right) > 0.^2$$

Die in (3.23 b) und (3.23 e) ausgesprochenen Bedingungen für μ_0 und ρ_0 sind z. B. erfüllt, wenn

$$(3.24) \quad \mu_0 = 1,6; \quad \rho_0 = 0,5$$

oder auch, wenn

$$(3.25) \quad \mu_0 = 1,281; \quad \rho_0 = \frac{1}{3}.$$

In allen Fällen, in denen ρ , κ und μ den Ungleichungen (3.21) mit μ_0 und ρ_0 aus (3.24) oder aus (3.25) genügen, folgt dann durch Vergleich von (3.23 b) und (3.23 e) mit (3.23) mittels vollständiger Induktion wegen (3.22) die Gültigkeit von (3.23) für alle v (Daß beim nächsten Schritt α und γ sowie κ und μ ihre Rollen vertauschen, macht nichts aus, da bisher die in (3.21) steckende Voraussetzung $\mu < \kappa$ nicht benutzt wurde, sondern nur daß $\mu_0 < \mu < 2$ und $\mu_0 < \kappa < 2$). Im besonderen sind in diesen Fällen also die sämtlichen Korrekturgrößen α_v , β_v und γ_v positiv.

Es gilt also dann allgemein (3.23) und, indem man v in (3.23 b) und (3.23 e) durch $(v-2)$ ersetzt, ebenfalls ganz allgemein (für $v > 2$)

$$(3.26) \quad 2\beta_{v-1} > \alpha_{v-1} > 0; \quad \gamma_v > 0.$$

¹ Wenn nur $\mu_0 - 2\rho_0 > \frac{1}{2}$ wie in (3.23 b).

² Nach (3.21) ist $\mu^2 - 1 - \rho\mu \cdot \frac{3}{2} > \mu^2 - 1 - \rho_0 \cdot \mu \cdot \frac{3}{2}$ und $\frac{d}{d\mu} (\mu^2 - 1 - \rho_0 \cdot \mu \cdot \frac{3}{2}) = 2\mu - \rho_0 \cdot \frac{3}{2} > 2 - \frac{3}{2} > 0$, also $\mu^2 - 1 - \rho \cdot \mu \cdot \frac{3}{2} > \mu_0^2 - 1 - \rho_0 \mu_0 \cdot \frac{3}{2}$.

Daraus und aus (3.21) folgt zunächst entsprechend (3.23 a) auch

$$(3.26 a) \quad \frac{\mu + \alpha_{v-1}}{1 + \beta_{v-1}} < 2,$$

weiterhin aber

$$(3.26 b) \quad x - \rho \cdot \frac{\mu + \alpha_{v-1}}{1 + \beta_{v-1}} < 2 - \rho \cdot \frac{\mu_0}{1 + \beta_{v-1}} \quad \text{und}$$

$$(3.26 c) \quad x^2 - 1 - \rho x \cdot \frac{\mu + \alpha_{v-1}}{1 + \beta_{v-1}} < 3 - \rho \cdot 2 \cdot \frac{\mu_0}{1 + \beta_{v-1}},$$

weil die Ableitung der linken Seite nach x wegen (3.26 a) und (3.21) sicher positiv ist.

In Anlehnung an (3.15 a) und (3.15 b) sei nun (für ein bestimmtes v)

$$(3.27) \quad 0 < \beta_{v-1} < 2\varepsilon; \quad 0 < \gamma_v < \frac{5}{2} \cdot \frac{\rho \cdot \varepsilon}{1 + \beta_{v-1}}.$$

Durch sinngemäße Anwendung¹ von (3.6) bzw. (3.7) folgen dann wegen (3.26 b) bzw. (3.26 c) die Ungleichungen

$$(3.27 a) \quad \beta_v < \left[\frac{5}{2} \frac{\rho \cdot \varepsilon}{1 + \beta_{v-1}} \cdot \left(2 - \frac{\rho \mu_0}{1 + \beta_{v-1}} \right) + \frac{2\varepsilon \rho^2}{1 + \beta_{v-1}} \right] \quad \text{bzw.}$$

$$(3.27 b) \quad \alpha_{v+1} < \frac{1}{1 + \beta_v} \cdot \left[\frac{5}{2} \cdot \frac{\rho^2 \varepsilon}{1 + \beta_{v-1}} \left(3 - \frac{2\rho \mu_0}{1 + \beta_{v-1}} \right) + \frac{4\varepsilon \rho^3}{1 + \beta_{v-1}} \right].$$

Dabei gelte weiterhin entweder (3.24) oder (3.25). In den beiden letzten Zeilen wird dann die eckige Klammer (für feste Werte von $\rho; \varepsilon; \mu_0$) am größten, wenn der Faktor $u = \frac{1}{1 + \beta_{v-1}}$ am größten, d. h. wegen (3.27), wenn $u = 1$ ist; denn die Ableitung von

$$(5 + 2\rho)u - \frac{5}{2} \rho \mu_0 u^2 \quad (\text{vgl. (3.27 a)})$$

$$\text{und von} \quad \left(\frac{15}{2} + 4\rho \right) u - 5 \rho \mu_0 u^2 \quad (\text{vgl. (3.27 b)})$$

nach u ist für $0 < u < 1$ durchweg positiv, wenn (3.24), d. h. $\mu_0 = 1,6; \rho < 0,5$, oder auch wenn (3.25), d. h. $\mu_0 = 1,281; \rho < \frac{1}{3}$

¹ μ mit x , α mit γ vertauscht und v durch $(v-1)$ ersetzt.

vorausgesetzt ist. Aus (3.27a) wird damit (indem man $\beta_{v-1} = 0$ einsetzt)

$$(3.27c) \quad (0 <) \beta_v < \frac{5}{2} \rho \varepsilon \cdot (2 - \rho \mu_0) + 2 \rho^2 \varepsilon < 2 \varepsilon$$

(wie bei (3.15e)) zunächst im Fall (3.24) und erst recht im Fall (3.25). Aus (3.27b) wird auf entsprechende Weise

(3.27d)

$$(0 <) \alpha_{v+1} < \frac{1}{1 + \beta_v} \cdot \left[\frac{5}{2} \rho^2 \varepsilon (3 - 2 \rho \mu_0) + 4 \rho^3 \varepsilon \right] < \frac{11 \rho \varepsilon}{4(1 + \beta_v)}$$

(wie bei (3.15f)) zunächst im Fall (3.24) und erst recht im Fall (3.25). Nach (3.21) und (3.23) ist weiter

$$(3.28a) \quad \mu - \rho \cdot \frac{x + \gamma_v}{1 + \beta_v} < x - \rho \cdot \frac{x}{1 + \beta_v} < 2 - \frac{2\rho}{1 + \beta_v}$$

und

$$(3.28b) \quad \mu^2 - 1 - \rho \mu \cdot \frac{x + \gamma_v}{1 + \beta_v} < x^2 - 1 - \rho x \cdot \frac{x}{1 + \beta_v} < 3 - \frac{4\rho}{1 + \beta_v},$$

weil wegen (3.23a) und (3.21) die Ableitung der linken Seite nach μ und die Ableitung der Mitte nach x immer positiv ist.

Aus (3.6) bzw. (3.7), (3.27c) und (3.27d) folgt damit

$$(3.28c) \quad (0 <) \beta_{v+1} < \frac{11 \rho \varepsilon}{4(1 + \beta_v)} \cdot \left(2 - \frac{2\rho}{1 + \beta_v} \right) + \frac{\rho^2 \cdot 2\varepsilon}{1 + \beta_v} \quad \text{bzw.}$$

$$(3.28d) \quad (0 <) \gamma_{v+2} < \frac{\rho}{1 + \beta_{v+1}} \cdot \left[\frac{11 \rho \varepsilon}{4(1 + \beta_v)} \cdot \left(3 - \frac{4\rho}{1 + \beta_v} \right) + \frac{\rho^2 \cdot 4\varepsilon}{1 + \beta_v} \right].$$

Denkt man sich ρ und ε fest, so ist die rechte Seite von (3.28c) wieder am größten für $v = \frac{1}{1 + \beta_v} = 1$, weil die Ableitung von $\rho \varepsilon \left\{ \left(\frac{11}{2} + 2\rho \right) v - \frac{11}{2} \rho v^2 \right\}$ nach v wegen $\beta_v > 0$ (d.h. $0 < v < 1$) und wegen $(0 <) \rho < \frac{1}{2}$ im Fall (3.24) und im Fall (3.25) durchweg positiv ist. Aus (3.28c) wird daher

$$(3.28e) \quad (0 <) \beta_{v+1} < \frac{11 \rho \varepsilon}{4} (2 - 2\rho) + \rho^2 \cdot 2\varepsilon \\ = \varepsilon \cdot \left(\frac{11}{2} \rho - \frac{7}{2} \rho^2 \right) < \frac{15}{8} \varepsilon \quad \text{für } \rho < \frac{1}{2}.$$

Die eckige Klammer in (3.28d) ist bei festem ε und β_v proportional zu $g(\rho) = \frac{33}{4} \rho - \frac{11\rho^2}{1+\beta_v} + 4\rho^2$ und damit wegen

$$\frac{dg}{d\rho} = \frac{33}{4} - \left(\frac{22}{1+\beta_v} - 8 \right) \rho > \frac{33}{4} - 14\rho \quad \left(> 0 \text{ für } \rho < \frac{1}{2} \right)$$

im Bereich $0 \leq \rho \leq \frac{1}{2}$ am größten für $\rho = \frac{1}{2}$. Die eckige Klammer selbst ist daher kleiner als

$$\frac{\varepsilon}{2} \left(\frac{33}{4} v - \frac{11}{2} v^2 + 2v \right) = \frac{\varepsilon}{2} \left(\frac{41}{4} v - \frac{11}{2} v^2 \right) \leq \frac{\varepsilon}{2} \cdot \frac{41^2}{8 \cdot 44} < \frac{5}{2} \varepsilon$$

(wobei wieder $\frac{1}{1+\beta_v} = v$ gesetzt wurde; Max. $(av - bv^2) = \frac{a^2}{4b}$).

Aus (3.28d) folgt damit

$$(3.28f) \quad (0 <) \gamma_{v+2} < \frac{5\rho\varepsilon}{2(1+\beta_{v+1})} \quad \left(\text{für } \rho < \frac{1}{2} \right).$$

Durch Vergleich von (3.28e) und (3.28f) mit (3.27) folgt dann unmittelbar die Konvergenz des Verfahrens, wenn man zu den Voraussetzungen (3.20) und (3.21) noch entweder (3.24) oder (3.25) hinzunimmt.

Es bleibt noch der Fall

$$(3.29) \quad 0 < \mu < 1,29; \quad \mu < \varkappa < 2,$$

wobei wieder $\rho < \frac{1}{3}$ vorausgesetzt sei.

In diesem Fall folgt aus (3.20) unmittelbar bzw. durch einfache Abschätzungen

$$(3.29a) \quad 0 < \gamma_1 < 0,43; \quad 0 < \beta_1 < -1; \quad -0,1 < \alpha_2 < 0,3; \quad -0,2 < \beta_2 < 0,4$$

(für β_2 unter Berücksichtigung von (3.6)).

Der Konvergenzbeweis unter den Voraussetzungen (3.29) läßt sich wegen der Möglichkeit negativer β_v hier nicht ohne weiteres in Anlehnung an (3.16) führen. Man kann aber trotzdem zum Ziel kommen, indem man z. B. noch $\varkappa \geq 1,5$ unterscheidet. Im einzelnen sei jedoch der Konvergenzbeweis für (3.29) hier nicht mehr wiedergegeben, da er etwas mehr Rechnung als die bisherigen Abschätzungen erfordert, ohne wesentlich Neues zu enthalten.

Das Ergebnis läßt sich unter Berücksichtigung der vorausgegangenen Abschätzungen zusammenfassen in

Satz 4: Benutzt man (3.19) oder $G_1 = x^2 + Ax + B$ als erste Näherung, je nachdem $|\mu| \leq |\varkappa|$, so konvergiert das in II. angegebene Verfahren bei zwei Paaren konjugiert komplexer Wurzeln mit einheitlichem Vorzeichen des Realteils immer, wenn das Verhältnis ρ des Betrags der „kleinen“ zum Betrag der „großen“ Wurzeln unter der Schranke $\rho = \frac{1}{3}$ bleibt.

Außerdem lassen sich nun aber auch die Überlegungen nach (3.18), die zu Satz 3 führten, zunächst bei einheitlichem Vorzeichen der Realteile aller 4 Wurzeln der gegebenen Gleichung auf den Fall übertragen, daß man von $K_1 = Bx^2 + Cx + 1$ oder von $G_1 = x^2 + Ax + B$ ausgeht

(Aus $B > 6,5$ und $0 < \mu < 1,6$; $\mu < \varkappa < 2$ folgt zwar zunächst nur $\rho + \frac{1}{\rho} > 3,3$ und damit $\rho < \frac{1}{2,95}$, also noch nicht genau das, was in (3.25) von ρ vorausgesetzt war; man erkennt jedoch aus den durchgeführten Abschätzungen sofort, daß für $\mu\varkappa \approx 3,2$ die Ungleichung $\rho < \frac{1}{2,95}$ völlig ausreicht).

Im folgenden soll jetzt aber gezeigt werden, daß die Voraussetzung einheitlicher Realteile hier ebenso unnötig ist wie in Satz 3. Es sei nämlich

$$(3.30) \quad -4 < \mu\varkappa < 0 \text{ und } B = \rho + \mu\varkappa + \frac{1}{\rho} > 6,5.$$

Dann ist

$$(3.30a) \quad \frac{1}{B} = \frac{\rho}{1 + \mu\rho + \rho^2} < \frac{1}{6,5} \text{ und } \rho + \frac{1}{\rho} > 6,5 - \mu\varkappa (> 6,5),$$

also, da nach (3.5) notwendig $|\rho| < 1$, jedenfalls $\frac{1}{\rho} > 6,34$; $\rho < 0,16$ und damit

$$(3.31) \quad \frac{1}{\rho} > 6,34 - \mu\varkappa; \quad \rho < \frac{1}{6,34 - \mu\varkappa} \left(< \frac{1}{6,34} \right).$$

Nach (3.20) und (3.30a) folgt

$$(3.32) \quad |\alpha_2| < \frac{1}{6,5} |(\varkappa^2 - 1) \cdot \mu\rho + \varkappa\rho^2|.$$

Nun haben nach (3.30) die Glieder $(-\mu\rho)$ und $x\rho^2$ das entgegengesetzte Vorzeichen von $x^2\mu\rho$. Daher ist nach (3.32), (3.31) und (3.30), je nachdem $|x^2\mu\rho| \geq |-\mu\rho + x\rho^2|$, entweder

$$(3.32a) \quad |\alpha_2| < \frac{1}{6,5} \cdot |x^2\mu\rho| < \frac{2}{6,5} |x\mu\rho| < \frac{2}{6,5} \left| \frac{\mu x}{6,34 - \mu x} \right| \\ < \frac{2}{6,5} \cdot \frac{4^1}{10,34} < 0,12 \quad \text{oder}$$

$$(3.32b) \quad |\alpha_2| < \frac{1}{6,5} \cdot \left| \frac{2}{6,34} + \frac{2}{6,34^2} \right|,$$

also jedenfalls

$$(3.32c) \quad |\alpha_2| < 0,12.$$

Außerdem folgen aber aus (3.20), (3.30) und (3.31) die Ungleichungen

$$(3.33) \quad |\gamma_1| < \frac{2}{6,34} < 0,32 \quad \text{und} \quad |\beta_1| < |x\mu\rho| < \frac{4^1}{10,34} \quad \text{oder} \quad |\beta_1| < \rho^2,$$

also in beiden Fällen sicher $|\beta_1| < 0,39$.

Nun sei allgemein vorausgesetzt, für ein bestimmtes v sei

$$(3.34) \quad |\gamma_v| < 0,32; \quad |\beta_v| < 3\varepsilon \quad \text{und} \quad |\alpha_{v+1}| < \varepsilon \quad \text{mit} \quad \varepsilon < 0,13 \\ \text{(also } |\beta_v| < 0,39\text{)}.$$

Nach (3.31) ist dann

$$\left| \rho \cdot \frac{x + \gamma_v}{1 + \beta_v} \right| < \frac{1}{6,34 + |\mu x|} \cdot \frac{|x| + 0,32}{1 - 0,39} < \frac{2,32}{6,34 + 2|\mu|} \cdot \frac{1}{0,61},$$

weil $\frac{|x| + 0,32}{6,34 + |\mu x|}$ wegen $|\mu| < 2$ mit $|x|$ zunimmt;

daher ist

$$|\mu - \rho \cdot \frac{x + \gamma_v}{1 + \beta_v}| < M(\mu) = |\mu| + \frac{2,32}{6,34 + 2|\mu|} \cdot \frac{1}{0,61} < 2 + \frac{2,32}{10,34 \cdot 0,61},$$

weil $M(\mu)$ mit $|\mu|$ zunimmt.

Nach (3.6) wird damit

$$|\beta_{v+1}| < \varepsilon \left(2 + \frac{2,32}{10,34} \cdot \frac{1}{0,61} \right) + \frac{3\varepsilon}{6,34^2} \cdot \frac{1}{0,61} < \varepsilon \cdot 2,5 (< 0,13 \cdot 2,5 < 0,33).$$

$\frac{1}{6,34 - y}$ nimmt monoton zu, also $\left| \frac{y}{6,34 - y} \right|$ für $-4 < y < 0$ monoton ab, wenn y wächst.

Entsprechend bekommt man jetzt aus (3.7)

$$|\gamma_{v+2}| < \frac{\varepsilon}{6,34 \cdot 0,67} \left[\left(3 + \frac{4,64}{10,34 \cdot 0,61} \right) + \frac{6}{6,34^2 \cdot 0,61} \right] < \varepsilon \cdot 0,95 (< 0,13).$$

Der Vergleich mit (3.34) gibt zusammen mit (3.33) und (3.32c) durch vollständige Induktion wieder den Konvergenzbeweis. Daher gilt allgemein der

Satz 5: Benutzt man (3.19) oder $G_1 = x^2 + Ax + B$ als erste Näherung, je nachdem $|\mu| \leq |x|$, so konvergiert das in II. angegebene Verfahren bei zwei Paaren konjugiert komplexer Wurzeln immer, wenn nur $B > 6,5$, d. h. bei nicht „normierten“ Gleichungen, wenn nur $b: \sqrt{d} > 6,5$.

IV. Bemerkungen zur Anwendbarkeit des Verfahrens

Für die praktische Verwendung des in II. angegebenen Verfahrens wird man im allgemeinen eine 1. Näherung für $K(x)$ durch Vernachlässigung von x^3 und x^4 oder eine 1. Näherung für $G(x)$ durch Vernachlässigung von x^1 und x^0 gewinnen. Für die Brauchbarkeit des Verfahrens ist nicht nur die Konvergenz selbst, sondern vor allem auch die Güte der Konvergenz maßgebend. Hierfür kann man unter Vernachlässigung höherer Potenzen von ρ die folgenden Näherungen gewinnen: Nach (3.20) gilt für die Korrekturgrößen β_v und γ_v zunächst

$$\beta_1 \sim x\gamma_1, \text{ also } x\beta_1 - \gamma_1 \sim (x^2 - 1)\gamma_1$$

und daher nach (3.6), (3.8) und (3.9) allgemein (für beschränktes v)

$$(4.1) \quad \begin{aligned} \gamma_{v+2} &\sim \rho^2 \cdot (\mu^2 - 1) \cdot (x^2 - 1) \cdot \gamma_v = q\gamma_v^1 \text{ und} \\ \beta_{v+2} &\sim \rho^2 \cdot (\mu^2 - 1) \cdot (x^2 - 1) \cdot \beta_v = q\beta_v. \end{aligned}$$

Nach (3.1) und (3.5) bedeutet dabei

ρ das Verhältnis des Betrags der „kleinen“ Wurzeln zu demjenigen der „großen“ Wurzeln,

¹ D. h. $|\gamma_{v+2} - q\gamma_v| < M\rho^3$ mit beschränktem M für beschränktes v .

— $\frac{1}{2}\mu$ das Verhältnis des Realteils zum Betrag bei den „großen“
Wurzeln und

— $\frac{1}{2}\kappa$ das Verhältnis des Realteils zum Betrag bei den „kleinen“
Wurzeln

(2 Paare konjugiert komplexer Wurzeln vorausgesetzt).

Einen Überblick über das Konvergenzverhalten bekommt man unmittelbar aus $F(x)$, wenn man (3.2) berücksichtigt. Aus (3.2) und (3.5) folgt nämlich (unter Vernachlässigung höherer Potenzen von ρ)

$$(4.2) \quad \rho \sim \frac{1}{B\left(1 - \frac{AC}{B^2}\right)}; \quad \mu \sim \sqrt{\rho(A - \rho C)}; \quad \kappa \sim \sqrt{\rho(C - \rho A)}.$$

Für das Beispiel (1.1) bedeutet dies

$$\rho \sim \frac{1}{10\left(1 - \frac{10}{100}\right)} = \frac{1}{9}; \quad \mu \sim \frac{2,5 - 0,4}{3} = 0,7; \quad \kappa \sim \frac{4 - 0,3}{3} \sim 1,2$$

und daher nach (4.1) für den nach 2 Divisionen erzielten Konvergenzfaktor q jeweils näherungsweise den Faktor

$$q \sim \frac{-0,51 \cdot 0,44}{81} \sim -0,003,$$

wodurch die Güte der Konvergenz schon recht genau erfaßt ist, wie ein Vergleich mit (2.4) zeigt.

Die Näherungsregel (4.2) läßt sich natürlich auch übertragen auf Gleichungen, die nicht so „normiert“ sind, daß der Faktor von x^4 und derjenige von x^0 gleich 1 ist. Ist

$$(4.3) \quad \xi^4 + a\xi^3 + b\xi^2 + c\xi + d$$

gegeben, so ergibt sich mittels $\xi = x \sqrt[4]{d}$ und Division durch d sofort

$$(4.3a) \quad B = \frac{b}{\sqrt[4]{d}}; \quad A = \frac{a}{\sqrt[4]{d}}; \quad C = \frac{c}{\sqrt[4]{d^3}}.$$

Nach Satz 5 genügt $B > 6,5$ für die Konvergenz des Verfahrens.

Das geschilderte Näherungsverfahren hat den Nachteil, daß es sich nicht für alle Fälle von Gleichungen 4. Grades mit zwei Paaren konjugiert komplexer Wurzeln eignet; es hat außerdem den Nachteil, daß es nach (4.1) bloß „von der 1. Ordnung“ konvergiert, d. h. daß der Fehler beim nächsten Schritt jeweils näherungsweise proportional dem Fehler der vorausgegangenen Näherung und nicht etwa (wie z. B. beim Newtonschen Verfahren) proportional dem Quadrat des Fehlers der vorausgegangenen Näherung ist. Das Verfahren hat aber den Vorteil einer sehr einfachen Rechenvorschrift und den Vorteil, daß man mit Hilfe von (4.3 a), (4.2) und (4.1), wenn eine Gleichung 4. Grades vorliegt, die Brauchbarkeit des Verfahrens und die Güte der Konvergenz in vielen Fällen durch eine einfache Überschlagsrechnung sofort feststellen kann. Es wird vor allem dann dem Verfahren von Graeffe (das wie das Newtonsche im wesentlichen „quadratisches“ Konvergenzverhalten zeigt) vorzuziehen sein, wenn die gewünschte Genauigkeit bereits durch 3 oder 4 Divisionen erreichbar ist (da sich der Vorteil des quadratischen Verhaltens bei den ersten Schritten noch nicht so auswirkt). Besonders im Fall $\mu x > 0$ (was z. B. bei Schwingungen um stabile Gleichgewichtslagen gilt, wenn es sich um die charakteristische Gleichung einer linearen Differentialgleichung handelt) wird das Verfahren häufig auch dann gut konvergieren, wenn die Beträge beider Wurzelpaare wenig voneinander abweichen.

Aus den Überlegungen in Abschnitt III. folgt übrigens zunächst für Gleichungen 4. Grades mit zwei Paaren konjugiert komplexer Wurzeln auch die Begründung für die folgende Rechenregel:

Ist von der linken Seite der Gleichung $F = x^4 + Ax^3 + \dots = 0$ näherungsweise ein quadratischer Faktor K für die Wurzeln mit dem kleineren Betrag bekannt, so ist das geeignete Verfahren zur Berechnung eines Näherungspolynoms für den anderen quadratischen Faktor (für die Wurzeln mit dem größeren Betrag) die gewöhnliche Division $F:K$; ist dagegen näherungsweise ein quadratischer Faktor G für die Wurzeln mit dem größeren Betrag bekannt, so verwendet man im Interesse größerer Genauig-

keit des anderen Faktors (für die Wurzeln mit dem kleineren Betrag) zweckmäßigerweise nicht die gewöhnliche, sondern die „umgekehrte“ Division (d. h. die Division nach steigenden Potenzen von x im Sinne der Gleichung (2.2 b)).

Noch ein Wort zu dem unter I. dargelegten Verfahren: die einzelnen Schritte erfordern dort etwas weniger Rechenarbeit als bei dem unter II. behandelten Verfahren; eine einfache Konvergenzbetrachtung für das Verfahren I. (für $\rho \ll 1$ und hinreichend gute erste Näherungen) ergibt jedoch, daß diese einfachere Rechenarbeit bei den einzelnen Schritten durch den Nachteil erkauft werden muß, daß man etwa mit doppelt so vielen Schritten rechnen muß, um zu ebenso guten Näherungen wie beim Verfahren II. zu kommen.

V. Gleichungen beliebigen Grades

Das in den vorausgegangenen Abschnitten behandelte Verfahren läßt sich unter geeigneten Voraussetzungen auf Gleichungen höheren Grades verallgemeinern; um dies zu zeigen, sei die gegebene Gleichung $F = G \cdot K = 0$ (falls es nicht von vorneherein so ist) mit Hilfe einer geeigneten Transformation $\xi = \lambda \cdot x$ so normiert, daß der größte unter ihren Wurzeln auftretende Betrag gleich 1 ist. Die Gleichung läßt sich dann, falls sie m Wurzeln vom Betrag 1 (zusammengefaßt in G) und k Wurzeln von kleinerem Betrag (zusammengefaßt in K) hat, darstellen in der Gestalt

$$(5.1) \quad 0 = F = K \cdot G =$$

$$= \underbrace{[x^k + \varepsilon_1 x^{k-1} + \varepsilon_2 x^{k-2} + \dots + \varepsilon_k]}_{K'} \cdot \underbrace{[x^m + c_1 x^{m-1} + c_2 x^{m-2} + \dots + c_{m-1} x + 1]}_{G''}$$

$$= \underbrace{[x^k + \varepsilon'_1 x^{k-1} + \varepsilon'_2 x^{k-2} + \dots + \varepsilon'_k]}_{K'} \cdot \underbrace{[x^m + (c_1 + \delta_1) x^{m-1} + \dots + (1 + \delta_m)]}_{G''} + \text{Rest,}$$

wobei die Größen ε'_h irgendwelche Näherungen für die Koeffizienten ε_h des Polynoms K bedeuten, während die Größen δ_μ die Korrekturglieder für die Koeffizienten des Polynoms G bedeuten, wenn man G durch $G'' = F:K'$ (unter Vernachlässigung des bei

der Division auftretenden Restes) ersetzt. Durch Koeffizientenvergleich folgt dann

$$\begin{aligned}c_1 + \varepsilon_1 &= (c_1 + \delta_1) + \varepsilon'_1; \\c_2 + \varepsilon_1 c_1 + \varepsilon_2 &= (c_2 + \delta_2) + \varepsilon'_1(c_1 + \delta_1) + \varepsilon'_2; \\c_3 + \varepsilon_1 c_2 + \varepsilon_2 c_1 + \varepsilon_3 &= (c_3 + \delta_3) + \varepsilon'_1(c_2 + \delta_2) + \varepsilon'_2(c_1 + \delta_1) + \varepsilon'_3 \text{ usw.},\end{aligned}$$

also unter Vernachlässigung von Gliedern der Gestalt $\varepsilon'_\lambda \delta_\mu$

$$\begin{aligned}\delta_1 &= \varepsilon_1 - \varepsilon'_1; \\\delta_2 &= (\varepsilon_1 - \varepsilon'_1) \cdot c_1 + (\varepsilon_2 - \varepsilon'_2); \\\delta_3 &= (\varepsilon_1 - \varepsilon'_1) \cdot c_2 + (\varepsilon_2 - \varepsilon'_2) \cdot c_1 + (\varepsilon_3 - \varepsilon'_3) \text{ und allgemein} \\(5.2) \quad &\vdots \\&\vdots \\&\vdots \\\delta_\mu &= (\varepsilon_1 - \varepsilon'_1) \cdot c_{\mu-1} + (\varepsilon_2 - \varepsilon'_2) \cdot c_{\mu-2} + \dots = \sum (\varepsilon_\lambda - \varepsilon'_\lambda) \cdot c_{\mu-\lambda} \\&\text{mit } c_0 = 1.\end{aligned}$$

Andererseits ergibt sich („umgekehrte“ Division durch $\frac{G''}{1+\delta_m}$) aus

(5.3)

$$\begin{aligned}F &= [1 + c_{m-1}x + c_{m-2}x^2 + \dots + c_1x^{m-1} + x^m] \cdot [\varepsilon_k + \varepsilon_{k-1}x + \varepsilon_{k-2}x^2 + \dots + x^k] \\&= \left[1 + \frac{c_{m-1} + \delta_{m-1}}{1 + \delta_m} x + \frac{c_{m-2} + \delta_{m-2}}{1 + \delta_m} x^2 + \dots \right] \\&\cdot [\varepsilon_k + (\varepsilon_{k-1} + \varepsilon_{k-1}^*)x + (\varepsilon_{k-2} + \varepsilon_{k-2}^*)x^2 + \dots] + \text{Rest}\end{aligned}$$

durch Koeffizientenvergleich für die Korrekturgrößen ε_k^* zunächst die Gleichung

$$\begin{aligned}c_{m-1} \varepsilon_k + \varepsilon_{k-1} &= \frac{c_{m-1} + \delta_{m-1}}{1 + \delta_m} \cdot \varepsilon_k + (\varepsilon_{k-1} + \varepsilon_{k-1}^*), \text{ also} \\(5.4) \quad \varepsilon_{k-1}^* &= \varepsilon_k \cdot \left[c_{m-1} - \frac{c_{m-1} + \delta_{m-1}}{1 + \delta_m} \right] = \varepsilon_k \cdot \frac{\delta_m c_{m-1} - \delta_{m-1}}{1 + \delta_m}.\end{aligned}$$

Ferner folgt aus (5.3) (durch Vergleich der Koeffizienten von x^2)

$$\begin{aligned}&c_{m-2} \varepsilon_k + c_{m-1} \varepsilon_{k-1} + \varepsilon_{k-2} \\&= \frac{c_{m-2} + \delta_{m-2}}{1 + \delta_m} \varepsilon_k + \frac{c_{m-1} + \delta_{m-1}}{1 + \delta_m} (\varepsilon_{k-1} + \varepsilon_{k-1}^*) + (\varepsilon_{k-2} + \varepsilon_{k-2}^*),\end{aligned}$$

also (wegen $c_{m-\lambda} - \frac{c_{m-\lambda} + \delta_{m-\lambda}}{1 + \delta_m} = \frac{c_{m-\lambda} \delta_m - \delta_{m-\lambda}}{1 + \delta_m}$)

$$\varepsilon_{k-2}^* = \varepsilon_k \cdot \frac{c_{m-2} \delta_m - \delta_{m-2}}{1 + \delta_m} + \varepsilon_{k-1} \cdot \frac{c_{m-1} \delta_m - \delta_{m-1}}{1 + \delta_m} - \varepsilon_{k-1}^* \cdot \frac{c_{m-1} + \delta_{m-1}}{1 + \delta_m}$$

und damit wegen (5.4) unter Vernachlässigung höherer Potenzen der δ_μ schließlich

$$(5.4a) \quad \begin{aligned} \varepsilon_{k-2}^* &= \varepsilon_k \cdot (c_{m-2} \delta_m - \delta_{m-2} - c_{m-1}^2 \delta_m + c_{m-1} \delta_{m-1}) + \\ &\quad + \varepsilon_{k-1} (c_{m-1} \delta_m - \delta_{m-1}) \end{aligned}$$

oder kurz

$$(5.4b) \quad \varepsilon_{k-2}^* = \sum \varepsilon_\rho \delta_\sigma C_{\rho\sigma}, \quad \text{wobei} \quad |C_{\rho\sigma}| < M$$

mit einer für alle Gleichungen gegebenen Grades gemeinsamen Schranke M , da die $|c_\mu|$ nach Voraussetzung (G soll lauter Wurzeln von Betrag 1 haben) höchstens gleich den Binomialkoeffizienten $\binom{m}{\mu}$ sein können. Ganz entsprechend folgt dann allgemein

$$(5.5) \quad \varepsilon_{k-r}^* = \sum \varepsilon_\rho \delta_\sigma C_{\rho\sigma r} \quad \text{mit} \quad |C_{\rho\sigma r}| < M.$$

Daraus und aus (5.2) folgt die Konvergenz des Verfahrens, falls nur die ε_ρ und die $(\varepsilon_\lambda - \varepsilon'_\lambda)$ hinreichend klein sind. Dies gilt sicher, wenn die Wurzeln von K alle hinreichend kleinen Betrag haben und wenn man für die erste Näherung K' dasjenige Polynom k . Grades verwendet, das entsteht, wenn man in F alle höheren Potenzen von x als x^k wegläßt; vgl. dazu die Bemerkungen nach (5.7). Übrigens kann man jede Gleichung $(m+k)$. Grades, die nicht lauter Wurzeln desselben Betrags hat, z. B. durch vorherige Anwendung des Graeffeschen Verfahrens¹ so transformieren, daß ihre Wurzeln die eben ausgesprochene Konvergenzbedingung erfüllen.

Aus dem Gedankengang des Konvergenzbeweises ergibt sich natürlich ohne weiteres, daß die Bedingung, G solle lauter Wurzeln gleichen Betrags haben, keineswegs notwendig ist, solange sich die Beträge der Wurzeln von G nicht allzustark voneinander unterscheiden.

¹ Und Normierung, so daß die „größten“ Wurzeln den Betrag 1 haben.

Als Beispiel diene hierzu die folgende Gleichung 8. Grades (nicht im Sinn von (5.1) „normiert“, was aber nichts Wesentliches ausmacht):

$$x^8 + 12x^7 + 110x^6 + 60x^5 + 400x^4 + 70x^3 + 75x^2 + 8x + 1 = 0.$$

Aus $G_0 = x^4 + 12x^3 + 110x^2 + 60x + 400$ folgt

$$K_1 = 1 + 7,85x + 73,55x^2 + 56,8x^3 + 371x^4$$

und dann $G_2 = x^4 + 11,85x^3 + 108,0x^2 + 41,1x + 372,05$;

$$K_3 = 1 + 7,89_0x + 73,84x^2 + 59,5x^3 + 371,7x^4$$

$$G_4 = x^4 + 11,84_0x^3 + 107,91x^2 + 40,4x + 371,85$$

$$K_5 = 1 + 7,89_{15}x + 73,85_{35}x^2 + 59,60x^3 + 371,84x^4.$$

Zu einer Abschätzung des Konvergenzbereichs und der Güte der Konvergenz kann man folgendermaßen kommen: Der größte unter den Wurzeln von K auftretende Betrag sei ρ ; der Betrag der Wurzeln von G sei wie in (5.1) durchwegs 1. Wir halten m und k fest und beschränken uns auf den Fall hinreichend kleiner ρ , so daß es genügt, jeweils nur die niedrigsten Potenzen von ρ , die vorkommen, zu betrachten.¹ Dann ist $|\varepsilon_1| \leq k\rho$ und

$$(5.6) \quad |\varepsilon_\nu| \leq \binom{k}{\nu} \cdot \rho^\nu \ll k\rho \quad \text{für } \nu > 1, \quad \text{ferner } |c_\mu| \leq \binom{m}{\mu}.$$

Außerdem sei

$$(5.7) \quad (\varepsilon_\nu - \varepsilon'_\nu) = \alpha_\nu \cdot \rho^{\nu+1},$$

wobei $|\alpha_\nu| \leq M_0$ mit einer nur vom Grad $m+k$ abhängigen Schranke M_0 . Dies ist sicher der Fall, wenn man die erste Näherung K' für K durch Weglassen der höheren Potenzen von x aus F gewinnt. Nach (5.1) ist nämlich

$$\begin{aligned} F = & x^{m+h} + \dots + x^h \cdot (1 + \varepsilon_1 c_{m-1} + \varepsilon_2 c_{m-2} + \dots) + \\ & + x^{h-1} (\varepsilon_1 + \varepsilon_2 c_{m-1} + \varepsilon_3 c_{m-2} + \dots) \\ & + x^{h-2} (\varepsilon_2 + \varepsilon_3 c_{m-1} + \dots) \\ & \vdots \\ & + x^0 \cdot \varepsilon_h \end{aligned}$$

¹ Wir setzen dabei voraus, daß die angeschriebenen Koeffizienten $\neq 0$ sind, sonst bekäme man noch bessere Konvergenzverhältnisse, weil z. B. die ε_ν^* in (5.9) und (5.10) eine höhere Potenz von ρ als Faktor hätten.

und daher im vorliegenden Fall (mit der Abkürzung $1 + \varepsilon_1 c_{m-1} + \varepsilon_2 c_{m-2} + \dots = 1 + \varepsilon$)

$$K' = x^k + x^{k-1} \cdot \frac{\varepsilon_1 + \varepsilon_2 c_{m-1} + \dots}{1 + \varepsilon} + x^{k-2} \cdot \frac{\varepsilon_2 + \varepsilon_3 c_{m-1} + \dots}{1 + \varepsilon} + \dots$$

$$= x^k + x^{k-1} \cdot \varepsilon'_1 + x^{k-2} \cdot \varepsilon'_2 + \dots \quad \text{nach (5.1).}$$

Daraus folgt

$$\varepsilon_v - \varepsilon'_v = \varepsilon_v - \frac{\varepsilon_v + \varepsilon_{v+1} c_{m-1} + \dots}{1 + \varepsilon} = \frac{\varepsilon_v \cdot \varepsilon - \varepsilon_{v+1} c_{m-1} + \dots}{1 + \varepsilon}.$$

Nun ist bis auf höhere Potenzen von ρ immer

$$|\varepsilon| = |\varepsilon_1 c_{m-1} + \varepsilon_2 c_{m-2} + \dots| \leq k \rho \cdot m \quad \text{nach (5.6)}$$

und damit wieder nach (5.6) bis auf höhere Potenzen von ρ

(5.7a)

$$|\varepsilon_v - \varepsilon'_v| \leq \binom{k}{v} \rho^v \cdot k \rho m + \binom{k}{v+1} \rho^{v+1} \cdot m = m \cdot \underbrace{\left[\binom{k}{v} \cdot k + \binom{k}{v+1} \right]}_{\text{beschränkt, wenn } m+k \text{ gegeben.}} \cdot \rho^{v+1}.$$

beschränkt, wenn $m + k$ gegeben.

Nach (5.2) und (5.7) folgt dann

$$(5.8) \quad \delta_\mu \sim c_{\mu-1} \cdot \alpha_1 \rho^2, \text{ speziell } \delta_m \sim c_{m-1} \alpha_1 \rho^2; \delta_{m-1} \sim c_{m-2} \alpha_1 \rho^2.$$

Unter Benutzung von (5.4) folgt daraus

$$(5.9) \quad \varepsilon_{k-1}^* \sim \varepsilon_k (c_{m-1} \delta_m - \delta_{m-1}) \sim \varepsilon_k (c_{m-1}^2 - c_{m-2}) \cdot \alpha_1 \rho^2.$$

Entsprechend folgt aus (5.4a) und (5.8)

$$\varepsilon_{k-2}^* \sim \varepsilon_{k-1} (c_{m-1}^2 - c_{m-2}) \alpha_1 \rho^2$$

und entsprechend aus (5.3) allgemein

$$(5.10) \quad \varepsilon_{k-\lambda}^* \sim \varepsilon_{k-\lambda+1} \cdot (c_{m-1}^2 - c_{m-2}) \cdot \alpha_1 \rho^2,$$

also speziell, weil nach (5.7) $\alpha_1 \rho^2 = \varepsilon_1 - \varepsilon'_1$,

$$\varepsilon_1^* \sim \varepsilon_2 \cdot (c_{m-1}^2 - c_{m-2}) \cdot (\varepsilon_1 - \varepsilon'_1);$$

$$\varepsilon_0^* \sim \varepsilon_1 \cdot (c_{m-1}^2 - c_{m-2}) \cdot (\varepsilon_1 - \varepsilon'_1).$$

Man erhält daher als zweite Näherung für $K=0$ auf Grund der in (5.3) angegebenen Bedeutung der ε_v^* die Gleichung

(5.10a)

$$0 = (1 + \varepsilon_0^*) x^k + (\varepsilon_1 + \varepsilon_1^*) x^{k-1} + \dots = (1 + \varepsilon_0^*) [x^k + \varepsilon_1'' x^{k-1} + \dots]$$

mit $\varepsilon_1'' = \frac{\varepsilon_1 + \varepsilon_1^*}{1 + \varepsilon_0^*}$, also $\varepsilon_1 - \varepsilon_1'' = \frac{\varepsilon_1 \varepsilon_0^* - \varepsilon_1^*}{1 + \varepsilon_0^*}$.

Nach (5.10) folgt daraus

$$(5.11) \quad |\varepsilon_1 - \varepsilon_1''| \sim |\varepsilon_1^2 - \varepsilon_2| \cdot |c_{m-1}^2 - c_{m-2}| \cdot |\varepsilon_1 - \varepsilon_1'|.$$

Bei hinreichend kleinem ρ ist also nach (5.6) eine Schranke für den Konvergenzfaktor q , mit dem sich der ursprüngliche Fehler nach 2 Divisionen multipliziert hat, zunächst für die Fehler bei ε_1 näherungsweise gegeben durch

$$(5.12) \quad q_0 = \rho^2 \cdot [k^2 + \binom{k}{2}] \cdot [m^2 + \binom{m}{2}].$$

Bei hinreichend kleinem ρ sind aber nach (5.9) und (5.10) die Fehler $|\varepsilon_v - \varepsilon_v''|$ alle entscheidend durch $\alpha_1 \rho^2 = (\varepsilon_1 - \varepsilon_1')$ bestimmt; nach (5.10a) ist nämlich

$$\varepsilon_v'' = \frac{\varepsilon_v + \varepsilon_v^*}{1 + \varepsilon_0^*}, \quad \text{also} \quad \varepsilon_v - \varepsilon_v'' = \frac{\varepsilon_v \varepsilon_0^* - \varepsilon_v^*}{1 + \varepsilon_0^*};$$

nach (5.10) folgt daraus

$$(5.11a) \quad |\varepsilon_v - \varepsilon_v''| \sim |\varepsilon_v \varepsilon_1 - \varepsilon_{v+1}| \cdot |c_{m-1}^2 - c_{m-2}| \cdot |\varepsilon_1 - \varepsilon_1'| \\ < M \cdot \rho^{v+1} \cdot |\varepsilon_1 - \varepsilon_1'| \quad \text{nach (5.6)}.$$

Dies entspricht der Annahme (5.7) für die $(\varepsilon_v - \varepsilon_v')$ und erlaubt damit die entsprechenden Schlüsse für weitere Näherungen, im besonderen z. B.

$$|\varepsilon_v - \varepsilon_v'''| \sim |\varepsilon_v \varepsilon_1 - \varepsilon_{v+1}| \cdot |c_{m-1}^2 - c_{m-2}| \cdot |\varepsilon_1 - \varepsilon_1''|.$$

Daher gibt q_0 ganz allgemein die gesuchte Schranke für den „Konvergenzfaktor“ bei Wiederholung des Verfahrens.

Das Ergebnis läßt sich zusammenfassen im folgenden

Satz 6: Bei gegebenem Grad $k + m$ der Gleichung $F = 0$ konvergiert das Verfahren sukzessiver Divisionen im Sinn von (5.1) und (5.3) (wobei G lauter Wurzeln gleichen Betrags hat) immer, wenn nur das Verhältnis ρ des größten Wurzelbetrags von K zum Betrag der Wurzeln von G klein genug ist. Unter Vernachlässigung höherer Potenzen von ρ ist bei wieder-

holter Anwendung des Verfahrens eine obere Schranke q_0 für den „Konvergenzfaktor“ q gegeben durch (5.12), d. h. durch

$$q_0 = \rho^2 \cdot [k^2 + \binom{k}{2}] \cdot [m^2 + \binom{m}{2}].$$

Dies bedeutet (vgl. die entsprechenden Überlegungen S. 87): Es ist $q \leq q_0 \cdot (1 + \varepsilon)$, wobei ε eine beliebig kleine positive Zahl bedeutet, wenn nur ρ unterhalb einer geeigneten (positiven) Schranke $\rho(\varepsilon)$ bleibt.

Im Fall reeller Polynome G vom 2. Grad mit zwei konjugiert komplexen Wurzeln ist nach (5.1) $c_{m-2} = 1$; ($4 \geq$) $c_{m-1}^2 > 0$, also $|c_{m-1}^2 - c_{m-2}| \leq 4 - 1 = 3$. Statt (5.12) folgt aus (5.11) dann

(5.13)

$$q_0 = 3 \cdot (k^2 + \binom{k}{2}) \cdot \rho^2, \text{ d. h. z. B. } q_0 = \begin{cases} 66 \rho^2 \approx (8,1 \rho)^2 \text{ für } k = 4 (m + k = 6), \\ 153 \rho^2 \approx (12,4 \rho)^2 \text{ für } k = 6 (m + k = 8). \end{cases}$$

Im Fall $m = 1$ ist in (5.8) sinngemäß $\delta_{m-1} = 0$ zu setzen, während $c_{m-1} = 1$. Statt (5.12) folgt dann aus (5.11) in diesem Fall

(5.14)

$$q_0 = (k^2 + \binom{k}{2}) \cdot \rho^2, \text{ d. h. z. B. } q_0 = \begin{cases} 22 \rho^2 < (5 \rho)^2 \text{ für } k = 4 (m + k = 5), \\ 35 \rho^2 < (6 \rho)^2 \text{ für } k = 5 (m + k = 6). \end{cases}$$

Von besonderem Interesse ist noch der Fall, daß F aus p quadratischen Faktoren mit je einem Paar konjugiert komplexer Wurzeln besteht, deren Beträge sich verhalten wie $1 : \rho_1 : \rho_2 : \dots : \rho_{p-1}$. Ist dann G der quadratische Faktor mit den Wurzeln größten Betrags (1), so folgt aus (5.11) unter der Voraussetzung

$$\rho_{v+1} : \rho_v \leq \rho \ll 1 \text{ für alle } v,$$

daß nicht nur $|c_{m-1}^2 - c_{m-2}| \lesssim 3$, sondern auch $|\varepsilon_1^2 - \varepsilon_2| \lesssim 3 \rho^2$ und damit

(5.15)

$$q_0 \lesssim (3 \rho)^2$$

für hinreichend kleine ρ in Übereinstimmung mit den Ergebnissen bei den Gleichungen 4. Grades.