

Beiträge zur Numerischen Mathematik 7

Beiträge zur Numerischen Mathematik 7

Herausgegeben von
Frieder Kuhnert und Jochen W. Schmidt



R. Oldenbourg Verlag München Wien 1979

CIP-Kurztitelaufnahme der Deutschen Bibliothek

Beiträge zur Numerischen Mathematik / hrsg. von
Frieder Kuhnert u. Jochen W. Schmidt. — München,

Wien: Oldenbourg.

NE: Kuhnert, Frieder [Hrsg.]

7. — 1. Aufl. — 1979.

ISBN 3-486-22411-5

© VEB Deutscher Verlag der Wissenschaften, Berlin 1979
Printed in the German Democratic Republic
Lizenz-Nr. 206 · 435/116/78
Gesamtherstellung: VEB Druckhaus „Maxim Gorki“, Altenburg
ISBN: 3-486-22411-5

Inhalt

J. ABAFFY and A. GALÁNTAI, Budapest Error estimations for conjugate direction methods	7
H. BIALY, Dresden Eine elementare Realisierung des lexikographischen Simplexverfahrens	13
R. FRANK und J. HERTLING, Wien Die Anwendung der Iterierten Defektkorrektur auf das Dirichletproblem	19
V. FRIEDRICH und A. UHLIG, Karl-Marx-Stadt Zur stochastischen Regularisierung linearer Gleichungen in Hilberträumen	33
B. HEINRICH, Karl-Marx-Stadt Monotone Differenzenapproximationen für lineare elliptische Differentialgleichungen mit gemischten Randbedingungen	49
J. HERZBERGER, Oldenburg Global konvergente Interpolationsmethoden zur Nullstelleneinschließung	65
B. HOFMANN, Karl-Marx-Stadt Über Quelldarstellungen bei einigen linearen Regularisierungsverfahren	75
G. PORATH, Güstrow, und E. TABBERT, Schwerin Das Tschebyscheffsche Iterationsverfahren für lineare Volterrasche Integralgleichungen zweiter Art	83
K. STREHMEL, Halle Eine Klasse A -stabiler Mehrschrittverfahren für Anfangswertaufgaben gewöhnlicher Diffe- rentialgleichungen	97
A. UHLIG, Karl-Marx-Stadt Numerische Vergleiche zwischen stochastischer Regularisierung und Projektionsverfahren am Beispiel der Rekonstruktion vertikaler Temperaturprofile der Erdatmosphäre auf der Grundlage von Satellitenmeßdaten	113
K. VETTERS, Dresden Asymptotisch symmetrische Verfahren für Nullstellen- und Extremwertaufgaben einer Ver- änderlichen	121

6 **Inhalt**

W. WEINELT und G. HELMERT, Karl-Marx-Stadt
Iterative Lösung spezieller nichtlinearer Differenzenschemata **139**

G. WINDISCH, Karl-Marx-Stadt
Ein Beispiel zur physikalischen Interpretation der Potenzmethode für spezielle dreidiagonale
Matrizen **159**

G. WINDISCH, Karl-Marx-Stadt
Zur numerischen Lösung eindimensionaler Wärmeleitprobleme mit nichtlinearen Randbedin-
gungen durch Differenzenmethoden **163**

G. ZIELKE, Halle
Motivation und Darstellung von verallgemeinerten Matrixinversen **177**

Error estimations for conjugate direction methods

JÓZSEF ABAFFY and AURÉL GALÁNTAI

This paper gives an error estimation for the class of conjugate direction methods when the objective function is quadratical.

1. Introduction

It is an important problem of the unconstrained function minimization to minimize quadratical functions of the form

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \quad (\mathbf{b}, \mathbf{x} \in R^n, c \in R^1) \quad (1.1)$$

where A is an $n \times n$ positive definite and symmetrical matrix.

The most important methods are based on the computation of conjugate directions ([1–5]). As it is known, the conjugate directions corresponding to the matrix A are defined as a linearly independent system of vectors $\{\mathbf{d}_i\}_{i=1}^n$ for which

$$\mathbf{d}_i^T A \mathbf{d}_j = \delta_{ij} \cdot c_{ij} \quad (i, j = 1, 2, \dots, n), \quad (1.2)$$

where δ_{ij} denotes the Kronecker symbol and $c_{ij} \neq 0$ are constant.

The conjugate direction methods then have the following general form:

Let $\mathbf{x}_1 \in R^n$ be arbitrary and

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \quad (i = 1, 2, \dots, n), \quad (1.3)$$

where α_i is defined by

$$\alpha_i = -\mathbf{g}(\mathbf{x}_i)^T \mathbf{d}_i / \mathbf{d}_i^T A \mathbf{d}_i \quad (1.4)$$

and $\mathbf{g}(\mathbf{x}_i) = \text{grad } f(\mathbf{x}_i)$.

In this scheme the methods differ from each others only in the computation way of conjugate directions $\{\mathbf{d}_i\}_{i=1}^n$.

Without loss of generality we can assume that

$$\|\mathbf{d}_i\| = 1 \quad (i = 1, 2, \dots, n). \quad (1.5)$$

Note that the number α_i for non-quadratical functions is defined by the relation

$$f(\mathbf{x}_{i+1}) = \min_{\alpha \in \mathbb{R}} f(\mathbf{x}_i + \alpha \mathbf{d}_i) \quad (1.6)$$

which is equivalent to (1.4) for the quadratical functions.

In two special cases the error of the above scheme was estimated by M. SACHET and S. KAHNE [4] for quadratical functions. The first case supposes error only in one step of the process and exact computations in the further steps. In the second case it is assumed that the errors in the computation of α_i and \mathbf{d}_i are independent of each others. In this way a simple estimation of the error was obtained.

In this paper we give a general error estimation (that means the errors of the different steps may depend on each others) which is obviously sharper than the result of SACHET and KAHNE.

2. A closed form of the error for quadratical functions

This section gives a closed form for the error propagation.

We denote the perturbed conjugate directions by \mathbf{p}_i ($i = 1, \dots, n$; $\mathbf{p}_i \neq \mathbf{0}$). The computed minimumpoints in the linear subspace of the vectors $\mathbf{d}_1, \dots, \mathbf{d}_i$ are denoted by \mathbf{y}_i ($i = 1, \dots, n + 1$). Similarly, β_i denotes the computed value of α_i ($i = 1, \dots, n$). Let $\delta_i = \mathbf{p}_i - \mathbf{d}_i$ be the error of direction \mathbf{d}_i ($i = 1, \dots, n$) and let us decompose β_i in the form

$$\beta_i = \alpha_i + \gamma_i + \varepsilon_i \quad (i = 1, \dots, n), \quad (2.1)$$

where

$$\gamma_i = -\mathbf{g}^\top(\mathbf{y}_i) \mathbf{p}_i / \mathbf{p}_i^\top A \mathbf{p}_i - \alpha_i \quad (i = 1, \dots, n) \quad (2.2)$$

is the accumulated error and ε_i is the round-off error of the computer actually used.

With the above notations we have

Theorem 1. *Consider the process (1.3)–(1.4) for the quadratical function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$ and assume that the first k steps are performed exactly i.e.*

$$\beta_i = \alpha_i, \quad \mathbf{p}_i = \mathbf{d}_i, \quad \mathbf{y}_i = \mathbf{x}_i \quad (i = 1, \dots, k; k < n). \quad (2.3)$$

Then

$$\mathbf{y}_j = \mathbf{x}_{k+1} + \sum_{i=k+1}^{j-1} \beta_i \mathbf{p}_i \quad (j = k + 1, \dots, n + 1) \quad (2.4)$$

and for the accumulated error

$$\gamma_i = \left[\alpha_i \delta_i^\top A \delta_i - 2\alpha_i \delta_i^\top A \mathbf{d}_i - \sum_{j=k+1}^{i-1} \beta_j \delta_j^\top A \mathbf{d}_j - \mathbf{g}^\top(\mathbf{y}_i) \delta_i \right] / \mathbf{p}_i^\top A \mathbf{p}_i \quad (2.5)$$

($i = k + 1, \dots, n$) holds.

Proof. Since the first relation is obvious we prove only (2.5) by induction corresponding to i . It is noted that (2.3) implies $\mathbf{y}_{k+1} = \mathbf{x}_{k+1}$. Using this fact and the relation

$$\mathbf{d}_i^\top A \mathbf{d}_i = \mathbf{p}_i^\top A \mathbf{p}_i - 2\delta_i^\top A \mathbf{p}_i + \delta_i^\top A \delta_i \tag{2.6}$$

we have

$$\begin{aligned} -\mathbf{g}^\top(\mathbf{y}_{k+1}) \mathbf{p}_{k+1} &= -\mathbf{g}^\top(\mathbf{x}_{k+1}) (\mathbf{d}_{k+1} + \delta_{k+1}) \\ &= \alpha_{k+1} \mathbf{d}_{k-1}^\top A \mathbf{d}_{k+1} - \mathbf{g}^\top(\mathbf{x}_{k+1}) \delta_{k+1} \\ &= \alpha_{k+1} \mathbf{p}_{k-1}^\top A \mathbf{p}_{k+1} + [\alpha_{k+1} \delta_{k-1}^\top A \delta_{k+1} - 2\alpha_{k+1} \delta_{k-1}^\top A \mathbf{p}_{k+1} \\ &\quad - \mathbf{g}^\top(\mathbf{y}_{k+1}) \delta_{k+1}] \\ &= (\alpha_{k+1} + \gamma_{k+1}) \mathbf{p}_{k+1}^\top A \mathbf{p}_{k+1} \end{aligned}$$

which gives the requested form of γ_{k+1} . Assume the truth of (2.5) for $i (1 \leq i < n)$. Then we have

$$\begin{aligned} \gamma_{i+1} &= -x_{i+1} - [\mathbf{g}^\top(\mathbf{x}_{k+1}) \mathbf{d}_{i+1} + \sum_{j=k+1}^i \beta_j (\mathbf{d}_j + \delta_j)^\top A \mathbf{d}_{i+1} \\ &\quad + \mathbf{g}^\top(\mathbf{y}_{i+1}) \delta_{i+1}] / \mathbf{p}_{i+1}^\top A \mathbf{p}_{i+1}. \end{aligned}$$

Now using the relation

$$\mathbf{g}^\top(\mathbf{x}_{i+1}) \mathbf{d}_{i+1} = [\mathbf{g}^\top(\mathbf{x}_{k+1}) + \sum_{j=k+1}^i \alpha_j \mathbf{d}_j^\top A] \mathbf{d}_{i+1} = \mathbf{g}^\top(\mathbf{x}_{k+1}) \mathbf{d}_{i+1}$$

and the relation (2.6) for $i + 1$ we have

$$\begin{aligned} \gamma_{i+1} &= -x_{i+1} + [\alpha_{i+1} \mathbf{p}_{i-1}^\top A \mathbf{p}_{i+1} + (\alpha_{i+1} \delta_{i-1}^\top A \delta_{i+1} \\ &\quad - 2\alpha_{i+1} \delta_{i-1}^\top A \mathbf{p}_{i+1} - \sum_{j=k+1}^i \beta_j \delta_j^\top A \mathbf{d}_{i+1} - \mathbf{g}^\top(\mathbf{y}_{i+1}) \delta_{i+1})] / \mathbf{p}_{i+1}^\top A \mathbf{p}_{i+1} \end{aligned}$$

which also implies (2.5) for γ_{i+1} , q.e.d.

3. A priori bound for the error propagation

In this section we give a bound for the error $\|\mathbf{y}_j - \mathbf{x}_j\|$ using the result of the previous section. The estimation needs a priori informations on the spectrum of A ($\lambda_{\max}, \lambda_{\min}$) and on the errors of conjugate directions actually used.

We need to prove the following lemma:

Lemma 1. *For the quadratical function (1.1) and the sequence (1.3) the inequality*

$$\|\mathbf{g}(\mathbf{x}_i)\| \leq \|\mathbf{g}(\mathbf{x}_1)\| \sqrt{\lambda_{\max}} / \sqrt{\lambda_{\min}} \quad (i = 1, \dots, n + 1) \tag{3.1}$$

holds.

Proof. The principal axis transformation keeps the norm of the gradient vector, so it is enough to prove (3.1) for the function of form

$$f^*(\mathbf{x}) = \sum_{i=1}^n \lambda_i x_i^2 \quad (\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0) \quad (3.2)$$

with the gradient

$$\mathbf{g}^*(\mathbf{x}) = (2\lambda_1 x_1, \dots, 2\lambda_n x_n)^\mathbf{T}. \quad (3.3)$$

consider the level set

$$S(\mathbf{z}) = \{\mathbf{x} \in R^n \mid f^*(\mathbf{x}) \leq f^*(\mathbf{z})\} \quad (\mathbf{z} \in R^n) \quad (3.4)$$

and look for $\max (\mathbf{g}^*(\mathbf{x}))^\mathbf{T} \mathbf{g}^*(\mathbf{x})$ in the set $S(\mathbf{z})$. Using the inequality $f^*(\mathbf{x}) \leq K$ ($K = f^*(\mathbf{z})$) we have

$$\begin{aligned} (\mathbf{g}^*(\mathbf{x}))^\mathbf{T} \mathbf{g}^*(\mathbf{x}) &= 4 \sum_{i=1}^n \lambda_i^2 x_i^2 \leq 4\lambda_1 \left(K - \sum_{i=2}^n \lambda_i x_i^2 \right) + 4 \sum_{i=2}^n \lambda_i^2 x_i^2 \\ &= 4\lambda_1 K + 4 \sum_{i=2}^n \lambda_i (\lambda_i - \lambda_1) x_i^2 \\ &\leq 4\lambda_1 K \leq 4 \frac{\lambda_1}{\lambda_n} \sum_{i=1}^n \lambda_i \lambda_n z_i^2 \leq \frac{\lambda_1}{\lambda_n} 4 \sum_{i=1}^n \lambda_i^2 z_i^2 = \frac{\lambda_1}{\lambda_n} \|\mathbf{g}^*(\mathbf{z})\|^2 \end{aligned}$$

which implies

$$\|\mathbf{g}^*(\mathbf{x})\| \leq \|\mathbf{g}^*(\mathbf{z})\| \sqrt{\lambda_{\max}} / \sqrt{\lambda_{\min}} \quad (\mathbf{x} \in S(\mathbf{z})).$$

For the conjugate direction method (1.3)–(1.4) holds $f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i)$ ($i = 1, \dots, n+1$), so the process cannot step out of the level set $S(\mathbf{x}_1)$. This obviously implies (3.1). q.e.d.

It is noted that Lemma 1 also implies

$$\|\mathbf{g}(\mathbf{y}_k)\| \leq \|\mathbf{g}(\mathbf{x}_1)\| \sqrt{\lambda_{\max}} / \sqrt{\lambda_{\min}} \quad (k > 1) \quad (3.5)$$

under the fairly formal assumption $f(\mathbf{y}_k) \leq f(\mathbf{x}_1)$ ($k > 1$), which holds for $\{\varepsilon_i\}$, $\{\delta_i\}$ small enough.

Another consequence of Lemma 1 is the inequality

$$|\alpha_i| \leq \alpha = \|\mathbf{g}(\mathbf{x}_1)\| \sqrt{\lambda_{\max}} / (\sqrt{\lambda_{\min}})^3 \quad (3.6)$$

for $i = 1, \dots, n$.

Let $\varepsilon = \max_i |\varepsilon_i|$ be the maximal round-off error and similarly, let $\delta = \max_i \|\delta_i\|$ be the maximal error of the conjugate directions. We assume that

$$\|\mathbf{p}_i\| = 1 \quad (i = 1, \dots, n) \quad (3.7)$$

also holds for the computed conjugate directions.

We then have

Theorem 2. Consider the process (1.3)–(1.4) for the quadratical function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$ and assume that the first k steps are performed exactly ($\beta_i = \alpha_i$, $\mathbf{p}_i = \mathbf{d}_i$, $\mathbf{y}_i = \mathbf{x}_i$; $i = 1, \dots, k$, $k < n$) and

$$f(\mathbf{y}_i) \leq f(\mathbf{x}_1) \quad (i = 1, \dots, n + 1) \quad (3.8)$$

is satisfied. Then

$$\|\mathbf{y}_j - \mathbf{x}_j\| \leq (j - k - 1) \left[\varepsilon + \frac{\alpha \delta}{\lambda_{\min}} \left(\delta \lambda_{\max} + 2\lambda_{\min} + \frac{j - k + 2}{2} \|\mathbf{A}\| \right) \right] \quad (3.9)$$

holds for $j \geq k + 1$.

Proof. From Lemma 1 and the relation $\mathbf{x}_j = \mathbf{x}_{k+1} + \sum_{i=k+1}^{j-1} \alpha_i \mathbf{d}_i$ it follows that

$$\|\mathbf{y}_j - \mathbf{x}_j\| \leq \left\| \sum_{i=k+1}^{j-1} [\alpha_i \delta \mathbf{d}_i + (\gamma_i + \varepsilon_i) \mathbf{p}_i] \right\| \leq \sum_{i=k+1}^{j-1} (\alpha \delta + \varepsilon + |\gamma_i|). \quad (3.10)$$

So we need to estimate $|\gamma_i|$ using the inequalities

$$|\gamma_i| \leq \left[\alpha \delta \mathbf{d}_i^T \mathbf{A} \delta \mathbf{d}_i + 2\alpha \delta \|\mathbf{A}\| + \|\mathbf{g}^T(\mathbf{y}_i)\| \delta + \sum_{j=k+1}^{i-1} |\beta_j| \delta \|\mathbf{A}\| \right] / \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i$$

and

$$|\beta_i| = |\mathbf{g}^T(\mathbf{y}_i) \mathbf{p}_i / \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i| \leq \|\mathbf{g}(\mathbf{x}_1)\| \sqrt{\lambda_{\max}} / (\sqrt{\lambda_{\min}})^3$$

as well as the inequality

$$\lambda_{\min} \mathbf{x}^T \mathbf{x} \leq \mathbf{x}^T \mathbf{A} \mathbf{x} \leq \lambda_{\max} \mathbf{x}^T \mathbf{x}. \quad (3.11)$$

Thus we get the bound

$$|\gamma_i| \leq \alpha \delta [\delta \lambda_{\max} + \lambda_{\min} + (i - k + 1) \|\mathbf{A}\|] / \lambda_{\min}. \quad (3.12)$$

Hence, a simple calculation yields the bound (3.9), q.e.d.

This result means that the error propagation is linear in the term of the maximal error of the conjugate directions. The bound can be decreased if it is possible to choose the conjugate directions such that

$$\lambda_{\min} \leq \mathbf{d}_i^T \mathbf{A} \mathbf{d}_i < \mathbf{d}_{i+1}^T \mathbf{A} \mathbf{d}_{i+1} \quad (i = 1, \dots, n - 1). \quad (3.13)$$

However a considerable improvement can be obtained only from the suitable choice of the initial value \mathbf{x}_1 according to the condition

$$\|\mathbf{g}(\mathbf{x}_1)\| = \|\mathbf{A} \mathbf{x}_1 + \mathbf{b}\| \leq \|\mathbf{b}\|. \quad (3.14)$$

At last we remark that using Gerschgorin type estimates for the spectrum of \mathbf{A} the dependence on λ_{\max} and λ_{\min} can be omitted in the bound (3.9). The simple estimate $\lambda_{\max} \leq \|\mathbf{A}\|$ yields the independence of λ_{\max} .

References

- [1] ABAFFY, J., A new three parameter class of Quasi-Newton methods, *Alkalmazott Matematikai Lapok* (to be published).
- [2] ABAFFY, J., and F. SLOBODA, General algorithms for generating conjugate directions and their applications, in: *Algoritmy vo vypoctovej technike*, Slovenska vedechotechnisa spolcnost (1975), p. 97—108.
- [3] FLETCHER, R., and C. M. REEVES, Function minimisation by conjugate gradients, *Computer J.* **7** (1964), 149.
- [4] SACHET, M., and S. KAHNE, Error analysis in conjugate direction methods (to be published).
- [5] WALSH, G. R., *Methods of optimization*, John Wiley, London 1975.
- [6] PSENICNY, B. N., and YU. N. DANILIN, Numerical methods of extremal problems (Russian), Nauka, Moscow 1975.

Manuskripteingang: 14. 2. 1977

VERFASSER:

Dr. JÓZSEF ABAFFY, Automation and Computer Institute of the Hungarian Academy of Sciences, Budapest

Dr. AURÉL GALÁNTAI, Department of Numerical Methods and Computer Science of the Eötvös Loránd University Budapest

Eine elementare Realisierung des lexikographischen Simplexverfahrens

HORST BIALY

1. Einleitung

Es gibt heute mehrere Zusatzmaßnahmen, die die Endlichkeit des Simplexverfahrens erzwingen. Am bekanntesten sind die Störungsmethode von CHARNES [1] und die Ausnutzung der Vektorlexikographie von DANTZIG, ORDEN und WOLFE [2]. Sinn beider Methoden ist die Verhinderung von Basiszyklen. Eine von JERKE in [3] für das duale Simplexverfahren vorgeschlagene und von SCHIEBEL, TERNO und UNGER in [4] noch vereinfachte lexikographische Vorgehensweise läßt sich in analoger Weise auch für das primale Simplexverfahren angeben. Das zugehörige Pivotwahlverfahren kommt ohne wesentliche Zusatzrechnung aus, ist numerisch einfacher als die Methode [5] und kann bedenkenlos auch in einer Einführungsvorlesung vorgetragen werden.

Wir bezeichnen Vektoren bzw. Matrizen mit halbfetten Klein- bzw. Großbuchstaben. Hochgestelltes \mathbf{T} ist das Transponierungszeichen.

Es sei

$$\begin{aligned} z &= d_0 + \mathbf{d}^T \mathbf{x}_N, \\ \mathbf{x}_B &= \mathbf{b} + \mathbf{B} \mathbf{x}_N \quad \text{mit } \mathbf{b} \geq \mathbf{0} \end{aligned} \quad (1)$$

ein Simplextableau eines in Normalform gebrachten linearen Optimierungsproblems. Die Komponenten der Vektoren \mathbf{x}_B bzw. \mathbf{x}_N sind die Basis- bzw. Nichtbasisvariablen. Die durch die Nichtbasisvariablen dargestellte Zielfunktion z sei zu minimieren. Für die numerische Rechnung benutzt man (1) in der Tabellenform

$$\begin{array}{c|cc} & 1 & \mathbf{x}_N^T \\ \hline z & d_0 & \mathbf{d}^T \\ \mathbf{x}_B & \mathbf{b} & \mathbf{B} \end{array} \quad (2)$$

Die Zeilen werden von 0 bis m , die Spalten von 0 bis n numeriert. $Z = \{0, 1, \dots, m\}$ ist die Menge der Zeilenindizes, $S = \{0, 1, \dots, n\}$ die Menge der Spaltenindizes. Im Sinne einheitlicher Transformationsformeln (3) ist die Verwendung der Tabellenmatrix

$$\mathbf{A} = \begin{pmatrix} d_0 & \mathbf{d}^T \\ \mathbf{b} & \mathbf{B} \end{pmatrix} = (a_{ik}), \quad i = 0(1)m, \quad k = 0(1)n,$$

vorteilhaft. Ein Basiswechsel mit dem Pivot $p := a_{st}$, also der Austausch der Basisvariablen in der Zeile $s \in Z \setminus \{0\}$ gegen die Nichtbasisvariable in der Spalte $t \in S \setminus \{0\}$,

erfolgt bekanntlich nach den Formeln

$$\begin{aligned}
 a_{st}^* &= 1/p, \\
 a_{sk}^* &= a_{sk}/(-p), & k \in S \setminus \{t\}, \\
 a_{it}^* &= a_{it}/p, & i \in Z \setminus \{s\}, \\
 a_{ik}^* &= a_{ik} + a_{it}a_{sk}^*, & i \in Z \setminus \{s\}, \quad k \in S \setminus \{t\}.
 \end{aligned}
 \tag{3}$$

Hier und im folgenden sind die Koeffizienten der neu entstehenden Tabelle durch einen Stern gekennzeichnet.

2. Das lexikographische Simplexverfahren

Definition 1. Eine Zeile der Tabelle (2) heißt *l-positiv* (*linkspositiv*), wenn das erste von Null verschiedene Element dieser Zeile positiv ist. Die Tabelle (2) heißt *l-positiv*, wenn die Zeilen mit den Nummern 1 bis m *l-positiv* sind.

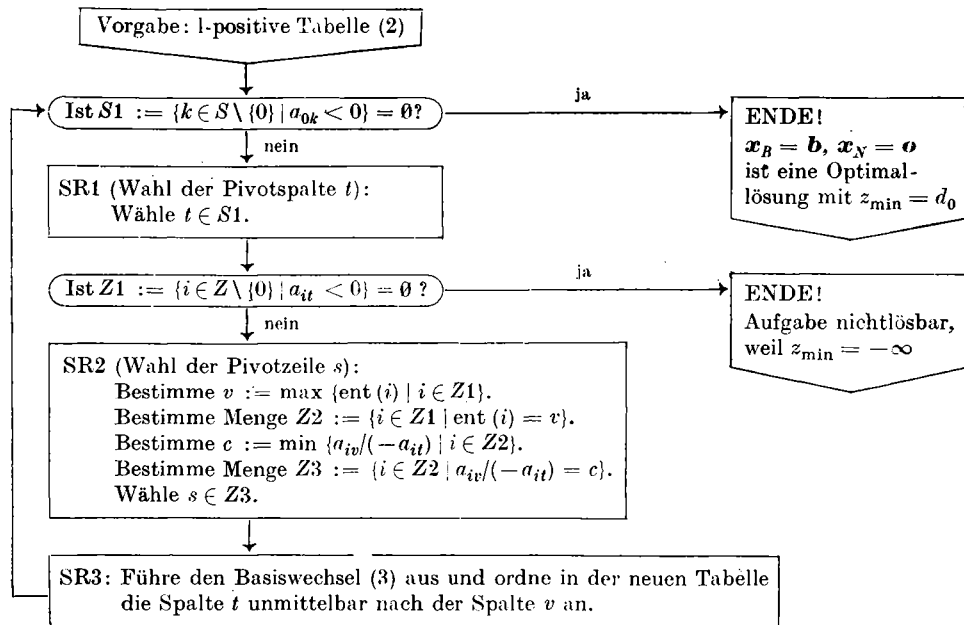
Definition 2. Tabelle (2) sei *l-positiv*. Dann heißt die Zahl

$$\text{ent}(i) := \min \{k \mid k \in S \wedge a_{ik} > 0\}$$

Entartungszahl der Zeile $i \in Z \setminus \{0\}$.

Die Entartungszahl gibt diejenige Stelle einer Zeile an, in der das erste von Null verschiedene Element steht.

Definition 3. Das *lexikographische Simplexverfahren* (LSV)



Die Zahl v heißt *Entartungsstufe des Basiswechsels*.

Satz 1. *Das LSV ist stets ausführbar. Nach jedem Basiswechsel entsteht wieder eine l-positive Tabelle, und es gilt $a_{0v}^* < a_{0v}$, wobei v die Entartungsstufe des Basiswechsels ist.*

Beweis. Die Ausführbarkeit der Simplexschritte SR1 bis SR3 ist klar, da Z2 und Z3 keine leeren Mengen sind, wenn Z1 nicht die leere Menge ist.

Wegen $Z3 \subset Z1$ gilt $p := a_{st} < 0$. Wegen $a_{sv} > 0$ ist sicher $v < t$. Die Spalten mit den Nummern $k = 0, 1, \dots, v - 1$ (das sind die ersten v Spalten!) ändern sich nicht.

Es gilt $a_{sv}^* = a_{sv}/(-p) > 0$. Außerdem ist $a_{sv}^* = c > 0$.

Für eine beliebige Zeile $i \in Z \setminus \{0, s\}$ mit der Entartungszahl $w := \text{ent}(i)$ sind folgende Fälle denkbar:

F 1: $w < v$. Es folgt $a_{iw}^* = a_{iw} > 0$.

F 2.1: $w = v \wedge a_{it} \geq 0$. Es gilt $a_{iw}^* = a_{iv}^* = a_{iv} + a_{it}a_{sv}^* \geq a_{iv} > 0$.

F 2.2: $w = v \wedge a_{it} < 0$. Wegen $i \in Z2$ gilt

$$a_{iw}^* = a_{iv}^* = a_{iv} + a_{it}(c) \geq a_{iv} + a_{it}(a_{iv}/(-a_{it})) = 0.$$

F 3.1: $w > v \wedge a_{it} > 0$. Dann ist $a_{iv}^* = 0 + c \cdot a_{it} > 0$ das erste von Null verschiedene Element in der neuen Zeile i .

F 3.2: $w > v \wedge a_{it} = 0$. Im Fall $t < w$ verändert sich die Zeile i nicht. Der Fall $t = w$ ist nicht möglich. Im Fall $t > w$ erfolgt lediglich die Umstellung einer Null aus Spalte t nach Spalte $v + 1$. In der neuen Tabelle gilt dann $\text{ent}(i) = w + 1$.

F 3.3: $w > v \wedge a_{it} < 0$. Dieser Fall ist wegen $i \in Z1$ und der Definition von v nicht möglich.

Die im Simplexschritt SR3 vorzunehmende Spaltenvertauschung ist ausschließlich wegen der in F 2.2 nicht vermeidbaren Beziehung $a_{iw}^* = 0$ erforderlich. Aus der Fallunterscheidung folgt die l-Positivität der neuen Tabelle.

Die Ungleichung $a_{0v}^* = a_{0v} + c \cdot a_{0t} < a_{0v}$ ist wegen $a_{0t} < 0$ trivial.

Satz 2. *Eine Wiederkehr der gleichen Basis ist im LSV unmöglich. Daher ist das LSV endlich, d. h., nach endlich vielen Ausführungen der Simplexregeln SR1 bis SR3 bricht das Verfahren ab.*

Beweis. Wir nehmen an, daß nach r aufeinanderfolgenden Basiswechseln wieder die gleiche Basis auftritt. v_1, v_2, \dots, v_r seien die zu den Basiswechseln gehörigen Entartungsstufen und $v := \min(v_1, v_2, \dots, v_r)$. Dann sind die ersten $v + 1$ Spalten während der Basiswechsel sicher nicht umgeordnet worden. Es sei x_u die zur Spalte v gehörende Nichtbasisvariable. Dann ist der z -Koeffizient von x_u nach Satz 1 kleiner geworden. Die Darstellung der Zielfunktion durch die Nichtbasisvariablen einer Basis ist aber eindeutig, unsere Annahme somit falsch. Im Verlauf des LSV ist also die Wiederkehr der gleichen Basis nicht möglich. Da jede Aufgabe nur endlich viele verschiedene Basen besitzt, muß das LSV nach endlich vielen Basiswechseln abbrechen.

3. Das Aufsuchen einer l-positiven Tabelle

Das LSV verlangt eine l-positive Ausgangstabelle, die stets auch ein Simplextableau ist. Aber die Umkehrung gilt nicht. Wir schlagen folgenden einfachen *Weg zur Aufstellung einer l-positiven Tabelle* vor.

1. *Bestimme ein Simplextableau (2)*. Hierfür sind zahlreiche Methoden bekannt, die an dieser Stelle nicht erörtert werden müssen.

2. Ist dieses Simplextableau nicht l-positiv, so gehe man zur *gestörten Tabelle*

$$\begin{array}{c|ccc} & 1 & \varepsilon & \mathbf{x}_N^T \\ \hline z & d_0 & M & \mathbf{d}^T \\ \mathbf{x}_B & \mathbf{b} & \mathbf{e} & \mathbf{B} \end{array} \quad \text{mit } \mathbf{e} = \begin{pmatrix} 1 \\ 1 \\ \dots \\ 1 \end{pmatrix} \in R^m \tag{4}$$

über, die sich durch eine zusätzliche positive Spalte mit der Nichtbasisvariablen ε auszeichnet. Dabei ist $M > 0$ eine hinreichend große Zahl, die garantiert, daß die Zusatzspalte niemals Pivotspalte gemäß der Simplexregel SR1 im LSV wird.

Da ε im Verlauf des LSV niemals Basisvariable wird, besitzen Originalproblem und gestörtes Problem die gleichen Lösungseigenschaften. Speziell ist die im LSV erzeugte Optimallösung des gestörten Problems wegen $\varepsilon = 0$ gleichzeitig Optimallösung des ungestörten Problems. Bei der numerischen Rechnung muß man M weder eintragen noch nach (3) umrechnen. Die Tabellen (2) und (4) sind äquivalent.

4. Ein Zahlenbeispiel

Das folgende Beispiel dient lediglich zur Demonstration der Regeln des LSV, beansprucht also kein Interesse als Optimierungsproblem. Die Tabellen sind links oben numeriert und die Pivots durch gekennzeichnet. Die erste Tabelle T1 ist bereits l-positiv. Die im LSV benötigten Daten sind ebenfalls angegeben.

T 1	1	x_1	x_2	x_3	x_4	x_5	x_6	x_7
z	2	1	0	4	-2	-2	1	1
x_8	0	0	0	2	-1	-1	0	1
x_9	3	2	1	1	0	-1	-1	-1
x_{10}	0	0	0	3	1	-1	-1	2
x_{11}	0	0	0	4	-3	-2	1	0
x_{12}	0	0	0	0	1	0	-2	3
x_{13}	0	0	0	0	0	0	1	-1

$S1 = \{4, 5\}$. Wir wählen $t = 5$. $Z1 = \{1, 2, 3, 4\}$, $v = 3$, $Z2 = \{1, 3, 4\}$, $c = 2$, $Z3 = \{1, 4\}$. Wir wählen $s = 1$. Nach dem Austausch von x_8 gegen x_5 ist die Spalte mit x_8 unmittelbar nach der Spalte mit x_3 anzuordnen. Andernfalls wäre die Zeile mit x_{11} in T 2 nicht mehr l-positiv (Beweisfall F 2.2 zu Satz 1). In den Zeilen mit x_{12} und x_{13} liegt der Fall 3.2 des Beweises zu Satz 1 vor.

T 2	1	x_1	x_2	x_3	x_8	x_4	x_6	x_7
z	2	1	0	0	2	0	1	-1
x_5	0	0	0	2	-1	-1	0	1
x_9	3	2	1	-1	1	1	-1	-2
x_{10}	0	0	0	1	1	2	-1	1
x_{11}	0	0	0	0	2	-1	1	-2
x_{12}	0	0	0	0	0	1	-2	3
x_{13}	0	0	0	0	0	0	1	-1

$S1 = \{7\}$. Also wird $t = 7$. $Z1 = \{2, 4, 6\}$, $v = 6$, $Z2 = Z3 = \{6\}$. Also wird $s = 6$. Da die Pivotspalte unmittelbar auf die Spalte mit der Nummer $v = 6$ folgt, unterbleibt ein Spaltenaustausch beim Übergang von T2 nach T3.

T 3	1	x_1	x_2	x_3	x_8	x_4	x_6	x_{13}
z	2	1	0	0	2	0	0	1
x_5	0	0	0	2	-1	-1	1	-1
x_9	3	2	1	-1	1	1	-3	2
x_{10}	0	0	0	1	1	2	0	-1
x_{11}	0	0	0	0	2	-1	-1	2
x_{12}	0	0	0	0	0	1	1	-3
x_7	0	0	0	0	0	0	1	-1

Wegen $S1 = \emptyset$ ist das LSV mit T3 beendet. Mit Ausnahme von $x_9 = 3$ sind sämtliche Komponenten der Optimallösung Null. Es gilt $z_{\min} = 2$.

Literatur

[1] CHARNES, A., Optimality and degeneracy in linear programming, *Econometrica* **20** (1952), 160–170.
 [2] DANTZIG, G. B., A. ORDEN and P. WOLFE, The generalized simplex method for minimizing a linear form under linear inequality restraints, *Pacific J. Math.* **5** (1955), 183–195.
 [3] JERKE, W., Duale Schnittverfahren der ganzzahligen linearen Optimierung unter besonderer Berücksichtigung der Entartung, Dissertation, TU Dresden 1974.
 [4] JERKE, W., W. SCHIEBEL, J. TERNO und G. UNGER, Ein Beitrag zur Behandlung der Entartung beim Simplexverfahren, *Beitr. Numer. Math.* **4** (1975), 105–114.
 [5] BIALY, H., Eine elementare Methode zur Behandlung des Entartungsfalles in der linearen Optimierung, *Unternehmensforschung* **10** (1966), 118–123.

Zusammenfassung

In der vorliegenden Arbeit wird eine numerische Realisierung des lexikographischen Simplexverfahrens angegeben, die gegenüber dem gewöhnlichen Simplexverfahren nur wenig Zusatzrechnung erfordert.

Summary

This paper presents a numerical realization of the lexicographic simplex method, which requires only a few supplement calculationwork compared to the usual simplex method.

Résumé

Dans l'article présent on donne une procédure numérique pour l'algorithme du simplexe lexicographique, quelle en comparaison de l'algorithme du simplex usuel demande peu de calcul supplémentaire.

Резюме

В данной работе указывают численную реализацию лексикографического симплекс-метода, требующая по сравнению с обычным симплекс-методом относительно мало добавочных вычислений.

Manuskripteingang: 15. 10. 1976

VERFASSER:

Dr. HORST BIALY, Sektion Mathematik, Naturwissenschaften und Rechentechnik der Hochschule für Verkehrswesen „Friedrich List“ Dresden

Die Anwendung der Iterierten Defektkorrektur auf das Dirichletproblem

REINHARD FRANK und JÖRG HERTLING

1. Darstellung der Methode

In [1] und [3] wurde die Methode der Iterierten Defektkorrektur entwickelt und deren Anwendung auf Anfangswertprobleme bei gewöhnlichen Differentialgleichungen und in [2] die Anwendung auf Zweipunkt-Randwertprobleme diskutiert. Die dort angegebenen Algorithmen, die sich in numerischen Experimenten als außerordentlich effektiv erwiesen haben [6], sollen nun auf partielle Differentialgleichungen vom elliptischen Typ erweitert werden. Während für Zweipunkt-Randwertprobleme in [2] die asymptotische Theorie ($h \rightarrow 0$) entwickelt wurde, beschränken wir uns in dieser Arbeit darauf, den Algorithmus der Iterierten Defektkorrektur für gewisse elliptische Randwertaufgaben (RWA) und numerische Experimente anzugeben. Wir betrachten

$$\begin{aligned} L[u] &\equiv a^{(11)}(x, y) \frac{\partial^2 u}{\partial x^2} + a^{(22)}(x, y) \frac{\partial^2 u}{\partial y^2} + a^{(12)}(x, y) \frac{\partial u}{\partial x} + a^{(21)}(x, y) \frac{\partial u}{\partial y} \\ &= f(x, y, u) \end{aligned} \quad (1.1)$$

auf dem Einheitsquadrat $R \equiv [0, 1] \times [0, 1]$ mit den Randbedingungen

$$\begin{aligned} u(0, y) &= r_1(y), & u(1, y) &= r_2(y), & 0 &\leq y \leq 1, \\ u(x, 0) &= \bar{r}_1(x), & u(x, 1) &= \bar{r}_2(x), & 0 &\leq x \leq 1. \end{aligned} \quad (1.2)$$

Natürlich setzen wir die Stetigkeit der Randfunktion in den Eckpunkten voraus, d. h., $r_1(0) = \bar{r}_1(0)$, $r_1(1) = \bar{r}_2(0)$, $r_2(0) = \bar{r}_1(1)$, $r_2(1) = \bar{r}_2(1)$.

Die exakte Lösung dieses Problems sei $z(x, y)$. In unserer Notation hat die Bedingung der gleichmäßigen Elliptizität die folgende Gestalt:

$$a^{(11)}(x, y) \xi^2 + a^{(22)}(x, y) \eta^2 \geq \alpha(\xi^2 + \eta^2), \quad \alpha > 0. \quad (1.3)$$

Zunächst möchten wir die Idee unserer Methode skizzieren. Dazu nehmen wir an, daß unser Problem mittels eines klassischen Differenzenschemas (Fünfpunkt-Sterne) gelöst wird. Dadurch erhält man auf dem Punktgitter

$$\mathbf{G}_h := \{(x_\nu, y_\mu); x_\nu = \nu h, y_\mu = \mu h, h = \frac{1}{l}, \nu, \mu = 0(1)l\} \quad (1.4)$$

die Näherungswerte $\eta_{\nu\mu}^{[0]}$ für die exakte Lösung $z(x_\nu, y_\mu)$. Nun interpolieren wir die Näherungswerte $\eta_{\nu\mu}^{[0]}$ auf $[0, 1] \times [0, 1]$ mit der Interpolationsfunktion $P_h^{[0]}(x, y)$:

$$P_h^{[0]}(x_\nu, y_\mu) = \eta_{\nu\mu}^{[0]}, \quad \nu, \mu = 0(1)m.$$

Der Index h deutet an, daß die von den $\eta_{\nu\mu}^{[0]}$ -Werten abhängige Interpolationsfunktion damit auch von der Schrittweite h abhängt. Nun konstruieren wir ein neues Randwertproblem

$$L[u] = f(x, y, u) + L[P_h^{[0]}] - f(x, y, P_h^{[0]}(x, y)), \quad (1.5)$$

$$u(0, y) = P_h^{[0]}(0, y), \quad u(1, y) = P_h^{[0]}(1, y), \quad (1.6)$$

$$u(x, 0) = P_h^{[0]}(x, 0), \quad u(x, 1) = P_h^{[0]}(x, 1),$$

wobei in L dieselben Koeffizientenfunktionen $a^{(11)}(x, y), \dots$ auftreten wie in (1.1). Offensichtlich hat dieses Randwertproblem $P_h^{[0]}(x, y)$ als exakte Lösung. Die Differentialgleichung (1.5) weicht von der ursprünglichen Gleichung (1.1) um die Störung

$$d_h^{[0]}(x, y) := L[P_h^{[0]}] - f(x, y, P_h^{[0]}(x, y)) \quad (1.7)$$

ab. Gilt $P_h^{[0]}(x, y) \approx z(x, y)$, $\frac{\partial^k}{\partial x^k} P_h^{[0]}(x, y) \approx \frac{\partial^k}{\partial x^k} z(x, y)$ und $\frac{\partial^k}{\partial y^k} P_h^{[0]}(x, y) \approx \frac{\partial^k}{\partial y^k} z(x, y)$ ($k = 1, 2$), was dadurch erreicht werden kann, daß erstens in (1.4) mit

hinreichend kleinem h gearbeitet wird, so daß die Größen $\eta_{\nu\mu}^{[0]}$ hinreichend gute Approximationen von $z(x_\nu, y_\mu)$ sind, und daß zweitens $P_h^{[0]}(x, y)$ aus einer geeigneten Klasse von Interpolationsfunktionen gewählt wird, dann ist die Störung $d_h^{[0]}(x, y)$ „klein“. Das bedeutet, daß das Problem (1.5), (1.6) „nahe“ unserem ursprünglichen Problem (1.1), (1.2) ist; daher nennen wir ab nun (1.5), (1.6) „Nachbarproblem“ (NP). Entsprechend bezeichnen wir nun (1.1), (1.2) als „Originalproblem“ (OP). Obwohl wir die exakte Lösung des NPs kennen, lösen wir es mit demselben Differenzenschema auf demselben Gitter. Das liefert die Näherungswerte $\pi_{\nu\mu}^{[0]}$ für $P_h^{[0]}(x_\nu, y_\mu)$. Natürlich kennen wir jetzt den globalen Diskretisierungsfehler $\pi_{\nu\mu}^{[0]} - P_h^{[0]}(x_\nu, y_\mu)$ von (1.5), (1.6). Da das NP nahe dem Originalproblem (1.1), (1.2) ist, können wir erwarten, daß $\pi_{\nu\mu}^{[0]} - P_h^{[0]}(x_\nu, y_\mu)$ eine gute Schätzung für den unbekanntem Diskretisierungsfehler $\eta_{\nu\mu}^{[0]} - z(x_\nu, y_\mu)$ von (1.1), (1.2) ist. Der eben beschriebene Gedanke geht auf ZADUNAIISKY zurück, der ihn für die Anwendung von Runge-Kutta-Methoden auf Anfangswertprobleme bei gewöhnlichen Differentialgleichungen eingeführt hat und ihn durch die obigen heuristischen Argumente motiviert hat [7]. Da wir eine echte Schätzung (mit richtigem Vorzeichen) und keine Abschätzung gewonnen haben, können wir diese Schätzung dazu benutzen, um unsere erste numerische Approximation auf die folgende Art zu verbessern: Wenn wir in der Identität

$$z(x_\nu, y_\mu) = \eta_{\nu\mu}^{[0]} - (\eta_{\nu\mu}^{[0]} - z(x_\nu, y_\mu))$$

den Ausdruck $\eta_{\nu\mu}^{[0]} - z(x_\nu, y_\mu)$ durch seine Schätzung $\pi_{\nu\mu}^{[0]} - P_h^{[0]}(x_\nu, y_\mu)$ ersetzen, erhalten wir voraussichtlich bessere Approximationen

$$\eta_{\nu\mu}^{[1]} := \eta_{\nu\mu}^{[0]} - (\pi_{\nu\mu}^{[0]} - P_h^{[0]}(x_\nu, y_\mu)) \quad (1.8)$$

für $z(x_\nu, y_\mu)$.

Es ist offensichtlich, daß diese Prozedur iteriert werden kann. Natürlich benutzt man im zweiten Schritt dieser Iteration eine neue Interpolationsfunktion $P_h^{[1]}(x, y)$, die die neuen Werte $\eta_{\nu\mu}^{[1]}$ interpoliert, und erhält

$$\eta_{\nu\mu}^{[2]} := \eta_{\nu\mu}^{[0]} - (\pi_{\nu\mu}^{[1]} - P_h^{[1]}(x_\nu, y_\mu)). \tag{1.9}$$

Weitere Iterationsschritte liefern

$$\eta_{\nu\mu}^{[j]} := \eta_{\nu\mu}^{[0]} - (\pi_{\nu\mu}^{[j-1]} - P_h^{[j-1]}(x_\nu, y_\mu)), \quad j = 3, 4, \dots \tag{1.10}$$

Um die Iteration (1.9) auf heuristische Art zu motivieren, beachten wir:

- a) Die Werte $\eta_{\nu\mu}^{[1]}$ sind (voraussichtlich) bessere Approximationen für $z(x_\nu, y_\mu)$ als $\eta_{\nu\mu}^{[0]}$.
- b) Daher ist die Funktion $P_h^{[1]}(x, y)$ (bzw. ihre Ableitungen), die die „besseren“ Werte $\eta_{\nu\mu}^{[1]}$ interpoliert, der exakten Lösung (bzw. ihren Ableitungen) näher.
- c) Daher läßt das neue NP (das $P_h^{[1]}(x, y)$ als exakte Lösung hat) eine noch kleinere Störung

$$d_h^{[1]}(x, y) := L[P_h^{[1]}] - f(x, y, P_h^{[1]}(x, y))$$

erwarten. (Das neue NP ist an das OP „nähergerückt“.)

- d) Das bedingt voraussichtlich, daß die neue Fehlerschätzung $\pi_{\nu\mu}^{[1]} - P_h^{[1]}(x_\nu, y_\mu)$ besser ist als die erste Fehlerschätzung $\pi_{\nu\mu}^{[0]} - P_h^{[0]}(x_\nu, y_\mu)$.

In verschiedenen Arbeiten über gewöhnliche Differentialgleichungen (Anfangs- und Randwertprobleme) ist diese Methode „Iterierte Defektkorrektur“ genannt und ihr asymptotisches Verhalten für $h \rightarrow 0$ analysiert worden. Es stellte sich heraus, daß für eine Basismethode der Konsistenzordnung p und für Interpolationsfunktionen aus einer geeigneten Funktionenklasse der j -te Schritt der Iteration die Ordnung $(j + 1)p$ hatte. Bei eindimensionalen Problemen beruhte der Beweis dieser Tatsache auf der Existenz einer asymptotischen Entwicklung des globalen Diskretisierungsfehlers (nach h -Potenzen). Dieser Beweis läßt sich nicht unmittelbar auf Probleme vom Typ (1.1), (1.2) übertragen, da bei physikalisch relevanten Problemen im allgemeinen die Existenz „hinreichend langer“ asymptotischer Entwicklungen nicht gesichert ist. Auch bei glatten Daten $\alpha^{(1)}(x, y)$, $\alpha^{(22)}(x, y)$, $\alpha^{(1)}(x, y)$, $\alpha^{(2)}(x, y)$, $f(x, y, u)$, $r_1(y)$, $r_2(y)$, $\bar{r}_1(x)$, $\bar{r}_2(x)$ können nämlich in den Eckpunkten Singularitäten der höheren Ableitungen von $z(x, y)$ auftreten. Die numerischen Experimente, die wir durchgeführt haben, zeigen jedoch, daß die erzielten Resultate ähnlich gut sind wie bei eindimensionalen Problemen.

Die Konstruktion der Interpolationsfunktion und die Berechnung der Störungen $d_h^{[j]}(x, y_\mu)$ erscheint aufwendig. In Abschnitt 2 werden wir jedoch sehen, daß diese Berechnung in sehr einfacher Weise möglich ist, falls mit stückweise polynomialer Interpolation gearbeitet wird.

Ein wesentlicher Vorteil unserer Methode ist der folgende: Um das Originalproblem zu lösen (Berechnung der Werte $\eta_{\nu\mu}^{[0]}$) und um das j -te NP zu lösen (Berechnung der Werte $\pi_{\nu\mu}^{[j]}$), muß man nichtlineare Gleichungssysteme lösen. Wenn die Werte $\eta_{\nu\mu}^{[j]}$ näher an $z(x_\nu, y_\mu)$ rücken und wenn gleichzeitig die Fehlerschätzungen $\pi_{\nu\mu}^{[j]} - P_h^{[j]}(x_\nu, y_\mu)$ besser mit dem tatsächlichen Diskretisierungsfehler $\eta_{\nu\mu}^{[0]} - z(x_\nu, y_\mu)$ übereinstimmen, dann rücken die Werte $\pi_{\nu\mu}^{[j]}$ näher an unsere erste numerische

Approximation $\eta_{\nu\mu}^{[0]}$ (Abb. 1). Für die numerische Lösung des j -ten NPs, d. h. für die Lösung des entsprechenden nichtlinearen Gleichungssystems, erhalten wir auf Grund unserer letzten Bemerkung die guten Startwerte $\eta_{\nu\mu}^{[0]}$.

Dies bedeutet, daß wir bei allen Schritten der iterierten Defektkorrektur mit demselben Startvektor arbeiten. Da die Störungen $d_h^{[j]}(x, y)$ nicht von u abhängen, verschwinden sie bei der Konstruktion der Jacobimatrix, die dem j -ten NP entspricht (Differentiation von $f(x, y, u) + d_h^{[j]}(x, y)$ nach u). Falls nun mit einem Pseudonewtonverfahren gearbeitet wird (Auswertung der Jacobimatrix nur an der Startstelle), so bedeutet dies, daß während der iterierten Defektkorrektur stets mit derselben Jacobimatrix gearbeitet wird. (Das heißt, Dreieckszerlegungen treten nur bei der Lösung des OPs auf, jedoch nicht bei der Lösung der NPs!)

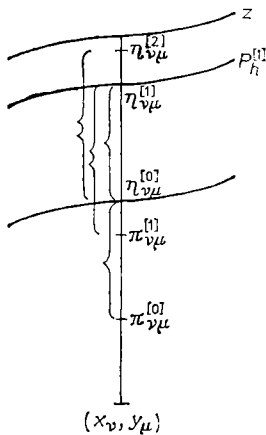


Abb. 1

Vergleichen wir unsere Methode mit einer Extrapolationsmethode, so bemerken wir, daß die Schrittweite fest bleibt und die Gleichungssysteme nicht explodieren. Das ist für mehrdimensionale Probleme ein sehr wichtiger Gesichtspunkt.

2. Algorithmische Details

Praktische Gründe, die noch später besprochen werden, legen nahe, daß für die Konstruktion unserer NPs stückweise polynomiale Interpolation in zwei Variablen sehr günstig ist. Da die NPs numerisch auf demselben Gitter gelöst werden müssen, müssen die Störungen $d_h^{[j]}(x, y)$ nur auf den Gitterpunkten berechnet werden. Daraus folgt, daß wir die Interpolationspolynome nicht explizit konstruieren müssen, sondern nur ihre ersten und zweiten partiellen Ableitungen an den Gitterpunkten.

Für die stückweise polynomiale Interpolation in zwei Variablen teilen wir das Einheitsquadrat in Teilquadrate (Abb. 2). Die Eckpunkte dieser Teilquadrate bilden das Gitter

$$\mathbf{G}_H := \left\{ (x^i, y^k); x^i = iH, y^k = kH, H = \frac{1}{n}, i, k = 0(1)n \right\} \quad (2.1)$$