

II. NUMERISCHE VERFAHREN

§5 Konstruktion von Näherungen

Es sei I das Intervall [a,b] subseteq IR und f in C(I); N in N. Es werden nun Verfahren zur Konstruktion von Näherungen beschrieben; das entscheidende Problem hierbei ist die Behandlung eines nichtlinearen Gleichungssystems.

5.1 Näherungen nach Meinardus

Gegeben sei eine Punktmenge X={x_i | 0 <= i <= 2N-1} subseteq I; gesucht ist

eine Funktion E(x) = sum_{i=1}^N a_i e^{t_i x} in V_N^0 mit

(5.1) E(x_i) = f(x_i) 0 <= i <= 2N-1

Nach [13] wird eine zumindest formale Lösung dieses Interpolationsproblems angegeben; hierzu wird X als äquidistante Punktmenge angenommen:

x_i = a + ih 0 <= i <= 2N-1, h > 0.

Setzt man für 1 <= j <= N A_j := a_j e^{t_j a}, E_j := e^{t_j h} und f_j := f(x_j) für 0 <= j <= 2N-1,

so erhält man aus (5.1)

(5.2) sum_{j=1}^N A_j E_j^i = f_i 0 <= i <= 2N-1

Zur Lösung dieses nichtlinearen Gleichungssystems mit den 2N Unbekannten A_j und E_j wird ein Verfahren angewandt, das auf Srinvasa Ramanujan, [14], zurückgeht:

Die rationale Funktion R(x) := sum_{j=1}^N A_j / (1 - E_j x) besitzt für x in U,

U = {x in IR | max_{1 <= j <= N} |E_j x| < 1}, die Darstellung

R(x) = sum_{j=1}^N A_j sum_{i=0}^inf E_j^i x^i = sum_{i=0}^inf x^i sum_{j=1}^N A_j E_j^i

$$F_N = \begin{bmatrix} f_{N-1} & f_{N-2} & \dots & f_0 \\ f_N & f_{N-1} & \dots & f_1 \\ \vdots & \vdots & & \vdots \\ f_{2N-2} & f_{2N-3} & \dots & f_{N-1} \end{bmatrix}$$

Zur Lösung von (5.2) und damit des Ausgangsproblems geht man vor, wie folgt:

1. Ermittlung der Frequenzen

Zuerst bestimmt man eine Lösung q des linearen Gleichungs-

systems $F_N q = \begin{bmatrix} f_N \\ f_{N+1} \\ \vdots \\ f_{2N-1} \end{bmatrix}$ mit $q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_N \end{bmatrix}$

Existiert eine Lösung q , so werden die Nullstellen des

Polynoms $Q(x) = \sum_{i=0}^N q_i x^i$ mit $q_0=1$ ermittelt.

Besitzt Q die N Nullstellen z_i , $1 \leq i \leq N$, so sind diese wegen $q_0=1$ von Null verschieden und nach (5.4) sei o.B.d.A.

$$(5.7) \quad E_i = z_i^{-1} \quad 1 \leq i \leq N.$$

Besitzt Q nur einfache, reelle Nullstellen und ist weiter $z_i > 0$ für $1 \leq i \leq N$ erfüllt, so erhält man die reellen

Frequenzen $t_i = h^{-1} \ln E_i = -h^{-1} \ln z_i$, $1 \leq i \leq N$.

2. Die Koeffizienten A_i und a_i

Die Koeffizienten p_i , $1 \leq i \leq N-1$, erhält man aus

$$L_{N-1} \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_{N-1} \end{bmatrix} + \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{N-1} \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_{N-1} \end{bmatrix}.$$

Mit $p_0=f_0$ ergeben sich nunmehr aus (5.3) und (5.7) durch Koeffizientenvergleich die Größen A_i , $1 \leq i \leq N$; es genügt dazu,

daß Q genau N reelle Nullstellen besitzt und so $E_i, 1 \leq i \leq N$, nach (5.7) bestimmt ist. Für den Fall, daß Q nur einfache, reelle Nullstellen besitzt, kann man nach Satz 5.1 A_i explizit angeben.

Sind N Frequenzen $t_i \in \mathbb{R}, 1 \leq i \leq N$, ermittelt, so gilt

$$a_i = A_i e^{-t_i a} \in \mathbb{R}, 1 \leq i \leq N.$$

Bemerkung:

1. Sind die Frequenzen $t_i, 1 \leq i \leq N$, bekannt, dann können die Koeffizienten $a_i, 1 \leq i \leq N$, auch durch Lösung des (nunmehr linearen) Gleichungssystems (5.1) bestimmt werden.
2. Ein nichtlineares Gleichungssystem wie (5.2) erhält man, wenn man an Stelle des Interpolationsproblems (5.1) ein Approximationsproblem auf einer Punktmenge $\{x_i \mid 0 \leq i \leq 2N, x_i = x_0 + ih\}$ zu lösen versucht; man vergleiche hierzu Abschnitt 5.2.
3. Kelly beschreibt in [9] zur Lösung von (5.1) ein Verfahren von Prony ("Prony's method"): Es wird benutzt, daß die $f_i, 0 \leq i \leq 2N-1$, falls $f_i = E(x_i)$ mit $E \in V_N^0$ erfüllt ist, einer Differenzgleichung N -ter Ordnung genügen. Die Bestimmung der Koeffizienten dieser Differenzgleichung führt auf das Gleichungssystem (5.6) und die Frequenzen ergeben sich wie oben beschrieben.

Zur Theorie des Verfahrens:

Die Bezeichnungen von oben werden beibehalten.

Formeln zur direkten Bestimmung der A_i gibt Satz 5.1 an:

Satz 5.1

Q besitze die N einfachen, reellen Nullstellen $E_j^{-1}, 1 \leq j \leq N$; E_j ist für $1 \leq j \leq N$ wegen $q_0 = 1$ definiert.

Behauptung:

$$A_k = \frac{\sum_{j=0}^{N-1} p_j E_k^{N-1-j}}{N \prod_{\substack{j=1 \\ j \neq k}} (E_k - E_j)} \quad 1 \leq k \leq N.$$

Beweis:

Wegen $E_j^{-1} \neq 0$ für $1 \leq j \leq N$ erhält man aus (5.3) für $1 \leq k \leq N$:

$$\sum_{j=0}^{N-1} p_j E_k^{-j} = \sum_{i=1}^N A_i \prod_{\substack{j=1 \\ j \neq i}}^N \left(1 - \frac{E_j}{E_k}\right) = A_k \prod_{\substack{j=1 \\ j \neq k}}^N \left(1 - \frac{E_j}{E_k}\right)$$

$$A_k = \sum_{j=0}^{N-1} p_j E_k^{-j} \prod_{\substack{j=1 \\ j \neq k}}^N E_k (E_k - E_j)^{-1} = \sum_{j=0}^{N-1} p_j E_k^{N-1-j} \prod_{\substack{j=1 \\ j \neq k}}^N (E_k - E_j)^{-1}$$

Damit ist die Behauptung gezeigt.

Es folgen nun Aussagen zur Existenz von Lösungen für das Gleichungssystem (5.6) und deren Eindeutigkeit.

Hilfssatz 5.1

Es sei $n \in \mathbb{N}$; besitzt das Polynom $\sum_{i=0}^n c_i x^i$ eine Nullstelle $z \neq 0$, dann ist z^{-1} eine Nullstelle des Polynoms $\sum_{i=0}^n c_{n-i} x^i$.

Beweis:

Nach Voraussetzung gilt

$$z^n \sum_{i=0}^n c_{n-i} (z^{-1})^i = \sum_{i=0}^n c_{n-i} z^{n-i} = 0;$$

wegen $z^n \neq 0$ folgt daraus die Behauptung.

Satz 5.2

Es gelte $f_i = E(x_i) = \sum_{j=1}^n a_j e^{t_j x_i}$ für $0 \leq i \leq 2N-1$ mit $\text{grad}(E) = n \leq N$;

A_i und E_i sind für $1 \leq i \leq n$ definiert wie oben. Es sei $z_i \in \mathbb{R}$ mit $z_i \neq 0$ für $n+1 \leq i \leq N$ beliebig gewählt.

Die Koeffizienten q_i , $0 \leq i \leq N$, seien gegeben durch

$$\sum_{i=0}^N q_i x^i := \prod_{i=1}^n (1 - E_i x) \prod_{i=n+1}^N (1 - z_i x).$$

Behauptung:

(5.8) Es gilt $F_N \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_N \end{bmatrix} = - \begin{bmatrix} f_N \\ f_{N+1} \\ \vdots \\ f_{2N-1} \end{bmatrix}$ und

für $n=N$ bilden die Koeffizienten q_i , $1 \leq i \leq N$, die eindeutig bestimmte Lösung von (5.8).

Beweis:

Das Polynom $\sum_{i=0}^N q_i x^i$ besitzt nach Voraussetzung die n Nullstellen E_i^{-1} , $1 \leq i \leq n$, (E_i^{-1} ist wegen $E_i = e^{t_i h}$ stets definiert).

Nach Hilfssatz 5.1 gilt für $1 \leq k \leq n$:

$$0 = \sum_{i=0}^N q_{N-i} E_k^i = \sum_{i=1}^N q_i E_k^{N-i} + q_0 E_k^N; \text{ wegen } q_0 = 1 \text{ folgt für } 1 \leq k \leq n: \\ \sum_{i=1}^N q_i E_k^{N-i} = -E_k^N$$

Für $1 \leq j \leq N$ gilt hiermit $\sum_{i=1}^N f_{N+j-1-i} q_i =$

$$= \sum_{i=1}^N \left(\sum_{k=1}^n A_k E_k^{N+j-1-i} \right) q_i = \sum_{k=1}^n \left(\sum_{i=1}^N q_i E_k^{N-i} \right) A_k E_k^{j-1} = \\ = \sum_{k=1}^n A_k E_k^{j-1} (-E_k^N) = -f_{N+j-1};$$

damit ist (5.8) gezeigt.

Es sei nun $n=N$; es gilt also $\sum_{i=0}^N q_i x^i = \prod_{i=1}^N (1 - E_i x)$ und

es sei $\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_N \end{bmatrix} = v$ eine zweite Lösung von (5.8); für $1 \leq j \leq N$ folgt

also

$$(5.9) \quad \sum_{i=1}^N f_{N+j-1-i} v_i = \sum_{k=1}^N A_k E_k^{j-1} \left(\sum_{i=1}^N v_i E_k^{N-i} \right) = -f_{N+j-1}.$$

Die $(N-N)$ -Matrix dieses Gleichungssystems besteht aus den Elementen $A_k E_k^{j-1}$, $1 \leq k, j \leq N$. Wegen $E \in V_N^0 - V_{N-1}$ gilt

$$(5.10) \quad A_k \neq 0, \quad 1 \leq k \leq N, \quad \text{und } 0 < E_k < E_{k+1}, \quad 1 \leq k \leq N-1.$$

Somit erhält man:

$$\det \begin{bmatrix} A_1 E_1^0 & A_2 E_2^0 & \dots & A_N E_N^0 \\ A_1 E_1^1 & A_2 E_2^1 & \dots & A_N E_N^1 \\ \vdots & \vdots & & \vdots \\ A_1 E_1^{N-1} & A_2 E_2^{N-1} & \dots & A_N E_N^{N-1} \end{bmatrix} =$$

$$= \left(\prod_{i=1}^N A_i \right) \det \begin{bmatrix} 1 & 1 & \dots & 1 \\ E_1 & E_2 & \dots & E_N \\ \vdots & \vdots & & \vdots \\ E_1^{N-1} & E_2^{N-1} & \dots & E_N^{N-1} \end{bmatrix} \neq 0,$$

da eine Vandermondesche Determinante vorliegt und die Ungleichungen (5.10) bestehen. Das Gleichungssystem (5.9) ist also eindeutig lösbar und es gilt daher für $1 \leq k \leq N$

$$\sum_{i=1}^N v_i E_k^{N-i} = -E_k^N,$$

da $\begin{bmatrix} E_1^N \\ \vdots \\ E_N^N \end{bmatrix}$ eine Lösung von (5.9) ist.

Das Polynom $\sum_{i=0}^N v_i x^{N-i} = \sum_{i=0}^N v_{N-i} x^i$ mit $v_0=1$ besitzt damit die

N Nullstellen E_k , $1 \leq k \leq N$, und nach Hilfssatz 5.1 gilt

für $x=E_k^{-1}$:

$$\sum_{i=0}^N v_i x^i = 0 \quad 1 \leq k \leq N.$$

Nach Voraussetzung stimmen die Nullstellen von $\sum_{i=0}^N v_i x^i$

und $\sum_{i=0}^N q_i x^i$ überein, womit die Eindeutigkeit gezeigt ist:

$$q_i = v_i, \quad 1 \leq i \leq N.$$

Bemerkung:

Nimmt man an, daß $f_i = E(x_i)$ für $0 \leq i \leq 2N-1$ mit $E \in V_N^0 - V_{N-1}$

erfüllt ist - (5.1) besitze also die eindeutig bestimmte

Lösung E -, so folgt mit Satz 5.2, daß das beschriebene

Verfahren die Lösung des Interpolationsproblems (5.1) ergibt.

Es soll nun noch gezeigt werden, daß die Aussage von Satz 5.2 zur Existenz einer Lösung nicht nur für $E \in V_N^0$ richtig ist.

Hilfssatz 5.2

Für $n \in \mathbb{N}$ und $t \in \mathbb{R}$ gilt mit $0 \leq m < n$:

$$\sum_{k=0}^n (-1)^k \binom{n}{k} (t-k)^m = 0$$

Beweis (Meinardus, [13]):

Es sei $K_r = \{z \in \mathbb{C} \mid |z| = r\}$, $r \in \mathbb{R}$. Unter Verwendung des Residuenkalküls erhält man für $r > n$:

$$\begin{aligned} \frac{1}{2\pi i} \int_{K_r} (t-z)^m \prod_{j=0}^n (z-j)^{-1} dz &= \sum_{k=0}^n \operatorname{Res}_{z=k} (t-z)^m \prod_{j=0}^n (z-j)^{-1} = \\ &= \sum_{k=0}^n \lim_{\substack{z \rightarrow k \\ z \neq k}} (z-k) (t-z)^m \prod_{j=0}^n (z-j)^{-1} = \sum_{k=0}^n (t-k)^m \prod_{\substack{j=0 \\ j \neq k}}^n (k-j)^{-1} = \\ &= \sum_{k=0}^n \frac{(t-k)^m}{k!(n-k)!} (-1)^{n-k} = (-1)^n \frac{1}{n!} \sum_{k=0}^n (-1)^k (t-k)^m \binom{n}{k} \end{aligned}$$

Es sei $A(z) := \sum_{j=0}^m a_j z^j := (t-z)^m$ und $B(z) := \sum_{j=0}^{n+1} b_j z^j := \prod_{j=0}^n (z-j)$.

Es gilt $a_m = (-1)^m$, $b_{n+1} = 1$ und man erhält für $z \in \mathbb{C}$, $z \neq 0$:

$$|A(z)| \leq |z|^m \left(1 + \sum_{j=0}^{m-1} |a_j| |z|^{j-m}\right)$$

$$|B(z)| \geq |z|^{n+1} \left(1 - \sum_{i=0}^n |b_i| |z|^{i-n-1}\right)$$

Wählt man daher $r_0 \in \mathbb{R}$, $r_0 > n$, hinreichend groß, so gilt

$$|A(z)| \leq 2|z|^m \quad \text{und} \quad |B(z)| \geq \frac{1}{2} |z|^{n+1} \quad \text{für } z \in \mathbb{C} \text{ mit } |z| \geq r_0.$$

Daraus folgt mit $r \geq r_0$: $\left| \int_{K_r} \frac{A(z)}{B(z)} dz \right| \leq 8\pi r^{m-n}$.

Wegen $m < n$ erhält man somit für alle $r \in \mathbb{R}$ mit $r \geq r_0$

$$\frac{1}{n!} \left| \sum_{k=0}^n (-1)^k (t-k)^m \binom{n}{k} \right| \leq 4r^{-1},$$

woraus die Behauptung folgt.

Bemerkung:

Diesen Hilfssatz erhält man auch aus der Beziehung

$$\Delta^n(0,1,\dots,n)g = \frac{1}{n!} \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} g(k) \quad (\text{Willers, [23], S.73})$$

mit $g(x) = (t-x)^m$; denn aus $m < n$ folgt nach (1.9)

$$\Delta^n(0,1,\dots,n)g = 0 \text{ und damit } \sum_{k=0}^n (-1)^k \binom{n}{k} (t-k)^m = 0$$

Satz 5.3

Es sei $N \geq 2$ und $f_i = E(x_i)$ mit $E(x) = \left(\sum_{i=0}^{N-1} a_i x^i \right) e^{tx} \in V_N - V_{N-1}$

für $0 \leq i \leq 2N-1$; ferner sei mit $E_1 = e^{tx}$ für $0 \leq i \leq N$ q_i gegeben

$$\text{durch } \sum_{i=0}^N q_i x^i = (1 - E_1 x)^N.$$

Behauptung:

$$F_N \cdot \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_N \end{bmatrix} = - \begin{bmatrix} f_N \\ f_{N+1} \\ \vdots \\ f_{2N-1} \end{bmatrix}$$

Beweis:

Mit dem Binomischen Lehrsatz erhält man

$$\sum_{i=0}^N q_i x^i = \sum_{i=0}^N (-1)^i \binom{N}{i} E_1^i x^i$$

und daraus durch Koeffizientenvergleich

$$q_i = (-1)^i \binom{N}{i} E_1^i, \quad 0 \leq i \leq N;$$

für $1 \leq j \leq N$ gilt daher mit $A = e^{ta}$:

$$\sum_{i=1}^N f_{N+j-1-i} q_i = A \sum_{i=1}^N \left(\sum_{k=0}^{N-1} a_k x_{N+j-1-i}^k \right) E_1^{N+j-1-i} q_i =$$

$$= A E_1^{N+j-1} \sum_{k=0}^{N-1} a_k \left(\sum_{i=1}^N E_1^{-i} q_i x_{N+j-1-i}^k \right) =$$

$$= A E_1^{N+j-1} \sum_{k=0}^{N-1} a_k \left(\sum_{i=1}^N (-1)^i \binom{N}{i} x_{N+j-1-i}^k \right)$$

Mit $0 \leq k < N$ erhält man durch Anwendung von Hilfssatz 5.2 :

$$\begin{aligned} \sum_{i=1}^N (-1)^i \binom{N}{i} x_{N+j-1-i}^k &= \sum_{i=1}^N (-1)^i \binom{N}{i} \sum_{m=0}^k \binom{k}{m} a^{k-m} (N+j-1-i)^m h^m = \\ &= \sum_{m=0}^k h^m \binom{k}{m} a^{k-m} \sum_{i=1}^N (-1)^i \binom{N}{i} ((N+j-1)-i)^m = \\ &= - \sum_{m=0}^k \binom{k}{m} a^{k-m} (N+j-1)^m h^m = -x_{N+j-1}^k \end{aligned}$$

Damit gilt $\sum_{i=1}^N f_{N+j-1-i} a_i = A E_1^{N+j-1} \sum_{k=0}^{N-1} a_k (-x_{N+j-1}^k) = -f_{N+j-1}$,

was zu zeigen war.

Bemerkung:

Versucht man (5.2) mit $A_i > 0$ für $1 \leq i \leq N$ zu lösen ($0 < E_i < E_{i+1}$ gilt nach Definition), dann liegt ein endliches Momentenproblem vor (Gantmacher, [25], S. 211); nach [25] besitzt (5.2) in diesem Fall genau dann eine Lösung, wenn die quadratischen

Formen $\sum_{j,k=0}^{N-1} f_{j+k} y_j y_k$, $\sum_{j,k=0}^{N-1} f_{j+k+1} y_j y_k$ positiv definit sind.

Läßt sich f als Dirichletsche Reihe mit positiven Koeffizienten darstellen, gilt also

$$f(x) = \sum_{n=1}^{\infty} c_n e^{s_n x} \text{ mit } c_n, s_n \in \mathbb{R}, c_i > 0 \text{ für } 1 \leq i \leq N \text{ und } c_n \geq 0, n \in \mathbb{N},$$

dann ist (5.2) lösbar, wie in [11] gezeigt ist:

Es gilt $f_i > 0$ für $0 \leq i \leq 2N-1$; mit $f_{j+k} = \sum_{n=1}^{\infty} c_n e^{s_n a} e^{(j+k) h s_n}$ gilt

$$\begin{aligned} \sum_{j,k=0}^{N-1} f_{j+k} y_j y_k &= \sum_{n=1}^{\infty} c_n e^{s_n a} \sum_{j,k=0}^{N-1} e^{j h s_n} y_j e^{k h s_n} y_k = \\ &= \sum_{n=1}^{\infty} c_n e^{s_n a} \left(\sum_{j=0}^{N-1} e^{j h s_n} y_j \right)^2, \text{ so daß } \sum_{j,k=0}^{N-1} f_{j+k} y_j y_k > 0 \text{ gilt,} \end{aligned}$$

falls nicht $y_i = 0$, $0 \leq j \leq N-1$, erfüllt ist (wegen $c_i > 0$ für $1 \leq i \leq N$). Ebenso folgt die positive Definitheit der zweiten quadratischen Form. Man beachte hierzu die Beispiele von §8.

Es wird nun gezeigt, daß die Bestimmung von $\det F_N$ keine Aussage über die Durchführbarkeit des Verfahrens zuläßt:

Beispiel 5.1

Es sei $N=2$, $X=\{0,1,2,3\}$.

1. $f_0=1$, $f_1=e$, $f_2=e^2$, $f_3=e^3$

Es gilt $\det F_2 = \det \begin{bmatrix} e & 1 \\ e^2 & e \end{bmatrix} = 0$, jedoch lassen sich

Lösungen von $F_2 \begin{bmatrix} q_1 \\ q_2 \end{bmatrix} = \begin{bmatrix} e^2 \\ e^3 \end{bmatrix}$ angeben:

Die Lösungsmenge des Systems ist gegeben durch

$$\left\{ \begin{bmatrix} -r \\ -e^2 + er \end{bmatrix} \mid r \in \mathbb{R} \right\}.$$

Es sei nun $r \in \mathbb{R}$, $r \neq e$, fest gewählt:

Damit folgt $Q(x) = (-e^2 + er)x^2 - rx + 1$; die Nullstellen von Q

sind $z_1 = e^{-1}$ und $z_2 = (r-e)^{-1}$. Damit folgt $E_1 = e$ und $E_2 = r-e$.

Es gilt $L_{N-1} = L_1 = 1$ und daher $p_1 = -r+e$ und $p_0 = 1$.

Die Gleichung (5.3) lautet hier:

$$1 + (e-r)x = A_1 + A_2 + x(e-r)A_1 - A_2 e x.$$

Dies ergibt: $A_2 = 1 - A_1$ und $e - r = (e-r)A_1 - (1 - A_1)e$

$$2e - r = A_1(2e - r) \implies A_1 = 1, A_2 = 0$$

Von Bedeutung ist also nur die Bestimmung von t_1 :

$$t_1 = \ln E_1 = 1.$$

Man erhält $E(x) = e^x$; dies ist die Lösung des Interpolationsproblems.

2. $f_0=1$, $f_1=-1$, $f_2=1$, $f_3=1$.

Hierfür gibt es keine Lösung $E \in V_2$ mit $E(i) = f_i$, $0 \leq i \leq 3$,

da $E \neq 0$ höchstens eine reelle Nullstelle besitzt.

Die Gleichungen

$$\begin{aligned} -q_1 + q_2 &= -1 \\ q_1 - q_2 &= -1 \end{aligned}$$

sind unverträglich und es gilt

$$\det F_2 = \det \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} = 0$$

5.2 Konstruktion einer Näherung nach Rice

=====

Es sei in $I=[a,b]$ die äquidistante Punktmenge

$$X=\{x_i \mid x_i \in I, x_i=a+ih, 1 \leq i \leq 2N+1, h>0\}$$

gegeben.

Unter der Annahme, daß es eine beste Approximation $E \in V_N^0 - V_{N-1}$ an $f \in C(I)$ bzgl. V_N^0 auf X gibt, wird ein Verfahren zur Bestimmung von E nach Rice, [16], hergeleitet.

Man kann hoffen, auf diesem Wege Näherungen für die beste Approximation an f auf I bzgl. V_N^0 zu erhalten, da diese, sofern sie existiert und die Länge N besitzt, in I eine Alternante der Länge $2N+1$ hat.

Es sei also $E(x) = \sum_{i=1}^N a_i e^{t_i x} \in V_N^0 - V_{N-1}$ die beste Approximation an f auf X ; es gilt daher

$$(5.11) \quad E(x_i) = f(x_i) + (-1)^i r \quad 1 \leq i \leq 2N+1$$

$|r|$ ist also die Minimalabweichung von f auf X ; eliminiert man r durch Addition aufeinanderfolgender Gleichungen, so erhält man das Gleichungssystem

$$(5.12) \quad \sum_{i=1}^N a_i e^{t_i a} (e^{t_i h} + 1) e^{t_i j h} = f(x_j) + f(x_{j+1}) \quad 1 \leq j \leq 2N.$$

Die Bestimmung der Koeffizienten und Frequenzen von E aus (5.12) wird nun nach Cesàro, [4], zurückgeführt auf die Ermittlung von Nullstellen eines Polynoms und die anschließende Lösung eines linearen Gleichungssystems. Hierzu werden die folgenden Bezeichnungen eingeführt:

$$f_j := f(x_j) + f(x_{j+1}) \quad 1 \leq j \leq 2N$$

$$A_i := a_i e^{t_i a} (e^{t_i h} + 1), \quad E_i := e^{t_i h} \quad 1 \leq i \leq N$$

Damit erhält man aus (5.12)

$$\sum_{i=1}^N A_i E_i^j = f_j \quad 1 \leq j \leq 2N$$

und da nach Voraussetzung $A_i \neq 0$ für $1 \leq i \leq N$ erfüllt ist, gilt

für $1 \leq k \leq N$ (wegen linearer Abhängigkeit der ersten Spalte):

$$\det \begin{bmatrix} f_k & E_1^k & \dots & E_N^k \\ f_{k+1} & E_1^{k+1} & \dots & E_N^{k+1} \\ \vdots & \vdots & \dots & \vdots \\ f_{k+N} & E_1^{k+N} & \dots & E_N^{k+N} \end{bmatrix} = 0$$

Wegen $E_i \neq 0$, $1 \leq i \leq N$, folgt daher

$$(5.13) \quad \det \begin{bmatrix} f_k & 1 & \dots & 1 \\ f_{k+1} & E_1^1 & \dots & E_N^1 \\ \vdots & \vdots & \dots & \vdots \\ f_{k+N} & E_1^N & \dots & E_N^N \end{bmatrix} = 0 \quad 1 \leq k \leq N .$$

Entwickelt man die Determinanten in (5.13) nach der ersten Spalte und ist D_i für $0 \leq i \leq N$ die Adjunkte von f_{k+i} , so gilt (da D_i unabhängig ist von k):

$$(5.14) \quad \sum_{i=0}^N f_{k+i} D_i = 0 \quad 1 \leq k \leq N .$$

Ebenso erhält man

$$(5.15) \quad \sum_{i=0}^N E_j^i D_i = 0 \quad 1 \leq j \leq N ,$$

da in den entsprechenden Matrizen zwei Spalten übereinstimmen. Aus $E_i < E_{i+1}$ folgt

$$(5.16) \quad D_0 \neq 0, D_N \neq 0$$

Je eine Gleichung aus dem System (5.15) ergibt zusammen mit (5.14) ein lineares, homogenes Gleichungssystem mit den Unbekannten D_i , $0 \leq i \leq N$:

Es sei $j \in \{1, \dots, N\}$:

$$(5.17) \quad \begin{aligned} \sum_{i=0}^N E_j^i D_i &= 0 \\ \sum_{i=0}^N f_{k+i} D_i &= 0 \quad 1 \leq k \leq N \end{aligned}$$

Wegen (5.16) folgt hieraus für $1 \leq j \leq N$:

$$\det \begin{bmatrix} 1 & E_j^1 & \dots & E_j^N \\ f_1 & f_2 & \dots & f_{N+1} \\ f_2 & f_3 & \dots & f_{N+2} \\ \vdots & \vdots & & \vdots \\ f_N & f_{N+1} & \dots & f_{2N} \end{bmatrix} = 0$$

Dies bedeutet, daß das Polynom

$$P(x) := \det \begin{bmatrix} x^0 & x^1 & \dots & x^N \\ f_1 & f_2 & \dots & f_{N+1} \\ \vdots & \vdots & & \vdots \\ f_N & f_{N+1} & \dots & f_{2N} \end{bmatrix}$$

die N Nullstellen $E_i = e^{t_i h}$, $1 \leq i \leq N$, besitzt.

Dieses Ergebnis legt folgendes Vorgehen zur Ermittlung der Parameter einer Näherung E nahe:

Man bestimmt die Nullstellen von P ; besitzt P die N Nullstellen E_i , $1 \leq i \leq N$, und sind diese reell und weiterhin einfach und positiv, dann erhält man die N reellen Frequenzen

$$t_i = h^{-1} \ln E_i, \quad 1 \leq i \leq N.$$

Mit diesen Werten löst man das lineare Gleichungssystem (5.11) oder (5.12), so daß damit die restlichen Parameter bestimmt sind.

Bemerkung:

Die Lösung von (5.11) ist vorteilhaft, falls $\|f-E\|_X$ von Bedeutung ist, etwa zur Bestimmung einer unteren Schranke für $\|f-E\|_I$.

Da die Koeffizienten von P sich aus der Berechnung von Determinanten ergeben, können hier Rundungsfehler das Ergebnis besonders leicht beeinflussen.

5.3 Konstruktion einer Näherung nach Willers

=====

In [23] behandelt Willers ein Verfahren zur "Annäherung" von Funktionen durch Elemente von V_N^0 , falls m Funktionswerte (Messwerte) mit $m \geq 2N$ gegeben sind.

Für die Konstruktion einer Näherung aus V_N^0 für eine gegebene Funktion $f \in C(I)$, $I=[a,b]$, bietet sich also mit diesem Verfahren die Möglichkeit, die Zahl der Funktionswerte, die in die Rechnung eingehen, zu variieren; mit $m=2N$ erhält man das in 5.1 beschriebene Verfahren.

Es sei also $m \geq 2N$ und $X=\{x_i | x_i=a+ih, h>0, 1 \leq i \leq m\} \subseteq I$ gegeben.

Nimmt man an, daß

$$E(x_j) = f(x_j) \quad 1 \leq j \leq m$$

mit $E(x) = \sum_{i=1}^N a_i e^{t_i x} \in V_N^0 - V_{N-1}^0$ erfüllt ist, so erhält man

mit $A_i := a_i e^{t_i a}$, $E_i := e^{t_i h}$ für $1 \leq i \leq N$ und $f_j := f(x_j)$ für $1 \leq j \leq m$:

$$(5.18) \quad \sum_{i=1}^N A_i E_i^j = f_j \quad 1 \leq j \leq m$$

Durch $Q(x) := \sum_{j=0}^N q_j x^j := \prod_{j=1}^N (x - E_j)$ sei q_i für $0 \leq i \leq N$ gegeben;

Es gilt also

$$(5.19) \quad Q(E_j) = 0 \quad 1 \leq j \leq N$$

Es sei nun $k \in \{1, \dots, m-N\}$ fest und man betrachte in (5.18) die $N+1$ Gleichungen

$$(5.20) \quad \sum_{i=1}^N A_i E_i^{j+k} = f_{j+k} \quad 0 \leq j \leq N$$

Multipliziert man für $0 \leq j \leq N$ die $(j+1)$ -te Gleichung von (5.20) auf beiden Seiten mit q_j , so erhält man mit (5.19) durch anschließende Addition der $N+1$ Gleichungen:

$$\sum_{j=0}^N f_{j+k} q_j = \sum_{j=0}^N (q_j \sum_{i=1}^N (A_i E_i^{k+j})) = \sum_{i=1}^N (A_i E_i^k \sum_{j=0}^N q_j E_i^j) = 0$$

Da diese Beziehung für $1 \leq k \leq m-N$ gilt, erhält man auf diese Weise ein lineares Gleichungssystem mit $m-N$ Gleichungen und N Unbekannten q_i , $0 \leq i \leq N-1$, wegen $q_N=1$:

$$(5.21) \quad \sum_{j=0}^{N-1} f_{j+k} q_j = -f_{N+k} \quad 1 \leq k \leq m-N$$

Es wird daher folgendes Verfahren vorgeschlagen:

Für $m=2N$ bestimme man die Unbekannten q_i , $0 \leq i \leq N-1$, durch Lösen des Gleichungssystems (5.21), falls dieses lösbar ist. Für $m > 2N$ bestimme man q_i , $0 \leq i \leq N-1$, durch Anwendung der Methode der kleinsten Quadrate.

Hat man so Koeffizienten q_i , $0 \leq i \leq N-1$, gefunden, so bestimmt man die Nullstellen E_i , $1 \leq i \leq N$, von $\sum_{i=0}^N q_i x^i$ mit $q_N=1$.

Gilt $E_i \in \mathbb{R}$ und $E_i > 0$ für $1 \leq i \leq N$, dann erhält man die

N Frequenzen $t_i = h^{-1} \ln E_i$, $1 \leq i \leq N$.

Sind die Nullstellen des Polynoms einfach, gilt also $t_i \neq t_j$ für $i \neq j$, dann bestimmt man a_i für $1 \leq i \leq N$ durch Lösung eines linearen Gleichungssystems, das aus N Gleichungen von (5.18) besteht.

5.4 Numerische Beispiele

=====

Zur numerischen Erprobung der drei beschriebenen Verfahren sind mit der CD 3300 des Rechenzentrums der Universität Erlangen-Nürnberg Berechnungen durchgeführt worden; die Programme sind in FORTRAN IV geschrieben und es ist stets mit doppelter Genauigkeit gerechnet worden. Alle Zahlenangaben sind daher gerundete Größen. Die Zeichnungen sind mit dem Plotter des Rechenzentrums und den dort vorhandenen Plotterroutinen angefertigt worden.

In den folgenden Beispielen werden die Bezeichnungen benutzt, wie sie bei der Herleitung der Verfahren eingeführt worden sind.

Für $1 \leq N \leq 4$ treten mit $f(x) = \frac{1}{x+1}$, $f(x) = \sqrt{x}$ in $I = [0,1]$ oder der Riemannschen Zetafunktion in $I = [2,3]$, $I = [2,4]$ keine Schwierigkeiten auf: Es wurde keine Punktmenge X gefunden, so daß hier ein Verfahren gescheitert wäre. Die Güte der jeweiligen Näherung E , d.h. $\|f-E\|_I$, hängt wesentlich von der Menge X ab und man wird daher im allgemeinen, um gute Näherungen zu erhalten, verschiedene Punktmenge X verwenden. In den Zeichnungen sind stets die jeweiligen Fehlerfunktionen $f-E$ dargestellt.

Beispiel 5.2

Es sei $f(x) = \sqrt{x}$, $I = [0,1]$ und $N=2$.

Zu bestimmen ist also eine Näherung $E(x) = a_1 e^{t_1 x} + a_2 e^{t_2 x}$; zur Beurteilung der folgenden Ergebnisse beachte man die Resultate von §7.

Eine günstige Näherung liefert das Verfahren von 5.1 mit den Punkten $x_1 = 0.01 + ih$, $0 \leq i \leq 3$, wobei $x_3 = 0.800$ und damit $h = 0.2633$ ist. Man erhält die Polynomkoeffizienten $q_2 = 0.345\ 112\ 153$, $q_1 = -1.467\ 229\ 623$, $q_0 = 1.000\ 000\ 000$, $p_1 = 0.376\ 089\ 942$, $p_0 = 1.000\ 000\ 000$.

Die Nullstellen von Q sind $3.398\ 958\ 380$ und $0.852\ 499\ 044$.

Es ergeben sich die Parameter

$$a_1 = -0.483\ 373\ 717$$

$$t_1 = -4.646\ 084\ 908$$

$$a_2 = 0.558\ 037\ 443$$

$$t_2 = 0.606\ 012\ 117$$

(Zeichnung 1)

Die Norm der Fehlerfunktion ist 0.07466 .

Mit $x_0=0.1$, $x_3=0.9$, $h=0.2666$ erhält man

$$\begin{aligned} a_1 &= -0.496\ 967\ 379 & t_1 &= -2.961\ 059\ 561 \\ a_2 &= 0.655\ 627\ 088 & t_2 &= 0.450\ 327\ 955 \end{aligned} \quad (\text{Zeichnung 1A})$$

Die Norm der Fehlerfunktion auf I ist hier 0.15866, womit die Abhängigkeit von X demonstriert sei.

Das Verfahren von Rice ergibt mit $x_i = \frac{i-1}{4}$, $1 \leq i \leq 5$, das

Polynom $P(x) = p_2 x^2 + p_1 x + p_0$ mit

$$p_2 = -0.670\ 540\ 689 \quad , \quad p_1 = 0.965\ 925\ 826 \quad , \quad p_0 = -0.222\ 252\ 952$$

und den Nullstellen 1.153 063 231 , 0.287 454 610 .

Man erhält so die Parameter

$$\begin{aligned} a_1 &= -0.565\ 826\ 602 & t_1 &= -4.986\ 761\ 225 \\ a_2 &= 0.570\ 571\ 292 & t_2 &= 0.569\ 688\ 323 \end{aligned} \quad (\text{Zeichnung 2})$$

Für $1 \leq i \leq 5$ gilt:

$$|\sqrt{x_i} - a_1 e^{t_1 x_i} - a_2 e^{t_2 x_i}| = 0.004\ 744\ 688$$

Damit ist gezeigt, daß, obwohl $x_1=0.0$ und $x_5=1.0$ Alternantenpunkte der Fehlerfunktion für die beste Approximation an f bzgl. V_2 auf I sind, die Näherung nach Rice, die ja die beste Approximation auf $\{x_i | 1 \leq i \leq 5\}$ für f ist, in x_1 und x_5 wesentlich von der Minimallösung für f bzgl. V_2 verschieden sein kann

Mit dem Verfahren von Willers kann man oft bessere Näherungen für die Frequenzen erhalten; es werden m Punkte $x_i = 0.0 + (i-1)h$ mit $1 \leq i \leq m$ und $x_m = 1.0$ verwendet:

$m=10$: Es ist also $h = \frac{1}{9}$; man errechnet die Koeffizienten

$$q_2 = 1.000\ 000\ 000 \quad , \quad q_1 = -1.456\ 251\ 707 \quad , \quad q_0 = 0.406\ 468\ 037$$

und damit die Nullstellen 1.079 834 854 und 0.376 416 853 .

Die Frequenzen sind also:

$$\begin{aligned} t_1 &= -8.793\ 522\ 892 \\ t_2 &= 0.691\ 273\ 052 \end{aligned}$$

Durch Lösung des Systems

$$\sum_{i=1}^2 a_i e^{t_i x_k} = \sqrt{x_k} \quad k=2,3$$

folgt

$$\begin{aligned} a_1 &= -0.420\ 954\ 177 \\ a_2 &= 0.455\ 428\ 511 \end{aligned} \quad (\text{Zeichnung 3})$$

Durch Vergrößern von m erhält man nicht unbedingt bessere Werte:

$$\begin{array}{lll}
 m=22: & t_1 = -15.928\ 838 & t_2 = 0.799\ 146 \\
 m=32: & t_1 = -21.251\ 001 & t_2 = 0.838\ 908
 \end{array}$$

Schwierigkeiten ergeben sich mit $x \in [-1,1]$ für $f(x) = |x|$ bei $N \geq 3$

$$\text{und } f(x) = \begin{cases} -5 & : x \leq -0.5 \\ 10x & : x \in (-0.5, 0.5) \\ +5 & : x \geq 0.5 \end{cases} \quad \text{oder } f(x) = \begin{cases} 2+x & : x \leq 0 \\ 2-x & : x \geq 0 \end{cases}$$

für $N \geq 2$:

Beispiel 5.3

$$I = [-1,1], N=2, f(x) = \begin{cases} 2+x & : x \leq 0 \\ 2-x & : x \geq 0 \end{cases}$$

Bei der Berechnung nach 5.1 erhält man für Q meist komplexe Nullstellen, wie die folgende Tabelle zeigt:

(es genügt die Angabe von x_0 und x_3)

| x_0, x_3 | Die Nullstellen von Q | |
|------------|-------------------------|--------------|
| | Realteil | Imaginärteil |
| -1.0 1.0 | 0.8000 | ± 0.6000 |
| -0.9 0.9 | 0.8235 | ± 0.5673 |
| -0.8 0.8 | 0.8462 | ± 0.5329 |
| -0.7 0.7 | 0.8679 | ± 0.4967 |
| -0.5 0.5 | 0.9091 | ± 0.4166 |
| -0.8 0.7 | 0.7282 | ± 0.4225 |
| -0.7 0.8 | 1.0274 | ± 0.5962 |
| -0.5 0.4 | 0.6986 | ± 0.2149 |

Mit $x_0 = -0.5$ und $x_3 = 0.35$ jedoch ist das Verfahren durchführbar:

Man erhält die Polynomkoeffizienten

$$q_2 = 2.837\ 462\ 834 \quad q_1 = -3.470\ 763\ 132 \quad q_0 = 1.000\ 000\ 000$$

$$p_1 = -3.422\ 811\ 364 \quad p_0 = 1.500\ 000\ 000 \quad \text{und damit}$$

$$\text{die Nullstellen } 0.758\ 641\ 874 \quad \text{und } 0.464\ 550\ 581$$

Es folgt damit

$$a_1 = 2.820\ 549\ 263$$

$$t_1 = 0.974\ 913\ 359$$

$$a_2 = -0.898\ 915\ 427$$

$$t_2 = 2.705\ 946\ 470 \quad (\text{Zeichnung 4})$$

Die Minimalabweichung für f in I bzgl. V_1 ist 0.5 (die beste Approximation ist die konstante Funktion $E(x) \equiv 1.5 \in V_1$); die gegebene "Näherung" approximiert also schlechter als die beste Approximation bzgl. V_1 .

Weitere Ergebnisse:

Mit $x_0 = -0.5$, $x_3 = 0.3$ erhält man

$$a_1 = 2.138\ 865\ 180 \qquad t_1 = 0.692\ 894\ 768$$

$$a_2 = -0.185\ 691\ 049 \qquad t_2 = 5.381\ 257\ 439$$

Die Norm der Fehlerfunktion ist größer als 37.0 .

Mit $x_0 = -0.5$, $x_3 = 0.25$ erhält man

$$a_1 = 2.050\ 428\ 834 \qquad t_1 = 0.624\ 759\ 329$$

$$a_2 = -0.050\ 428\ 834 \qquad t_2 = 10.207\ 441\ 475$$

Die Norm der Fehlerfunktion ist größer als 1365.0 ;

für diese drei Fehlerfunktionen ist $x = 1.0$ der Extrempunkt.

Bei der Ermittlung einer Näherung nach Rice treten die gleichen Schwierigkeiten auf, da auch P meist komplexe Nullstellen besitzt:

| x_1, x_5 | Die Nullstellen von P | |
|------------|-------------------------|--------------|
| | Realteil | Imaginärteil |
| -1.0 1.0 | 0.8571 | ± 0.5151 |
| -0.9 0.9 | 0.8732 | ± 0.4873 |
| -0.8 0.8 | 0.8888 | ± 0.4581 |
| -0.5 0.5 | 0.9333 | ± 0.3590 |

Für $x_1 = -0.5$ und $x_5 = 0.25$ erhält man das Polynom $P(x) = \sum_{i=0}^2 p_i x^i$

mit $p_2 = -0.539\ 062\ 500$, $p_1 = 1.828\ 125\ 000$, $p_0 = -1.398\ 437\ 500$

und den Nullstellen $2.225\ 778\ 005$, $1.165\ 526\ 343$.

Damit folgt

$$a_1 = 2.314\ 520\ 833 \qquad t_1 = 0.816\ 921\ 503$$

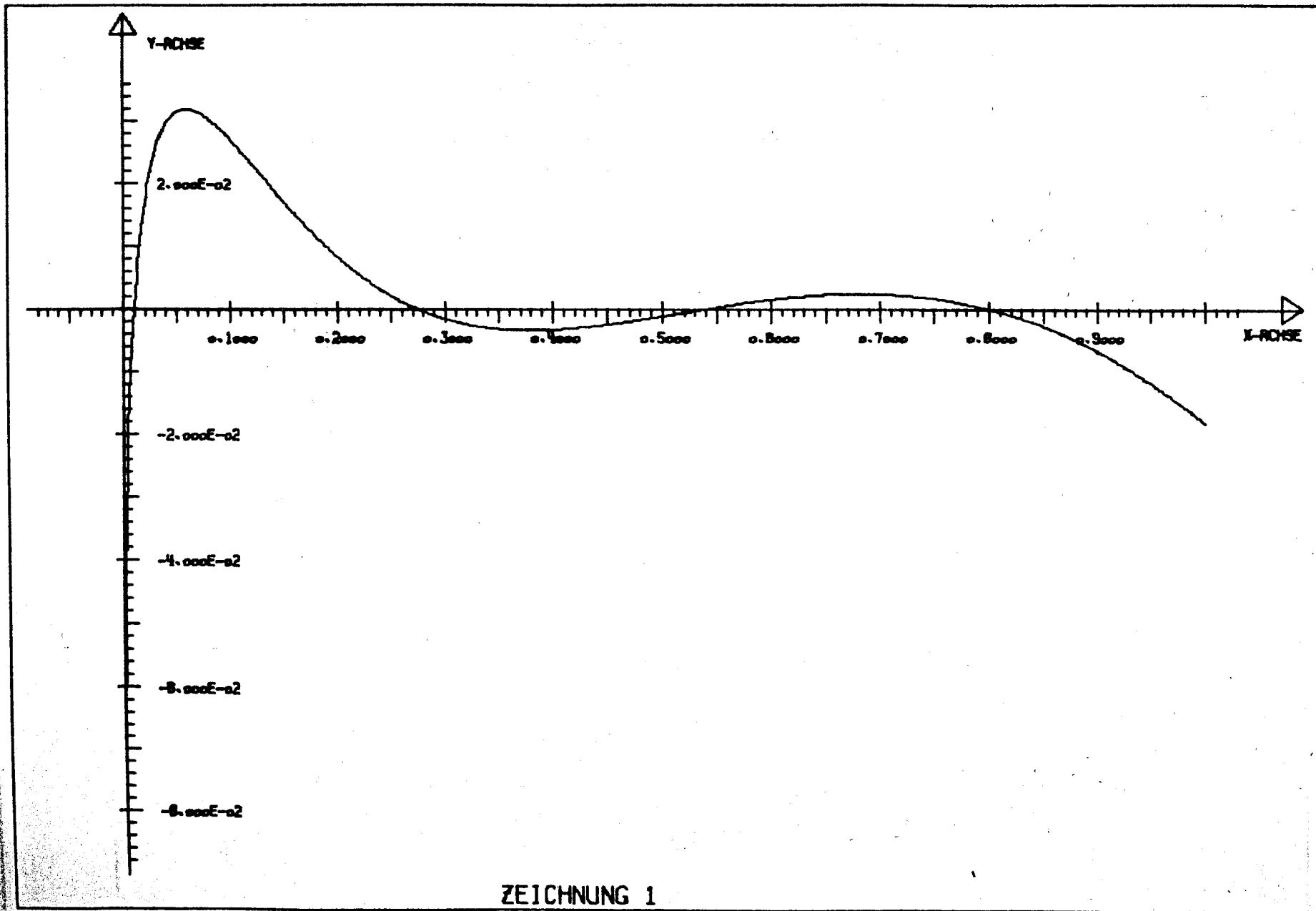
$$a_2 = -0.376\ 853\ 596 \qquad t_2 = 4.267\ 234\ 771 \qquad (\text{Zeichnung 5})$$

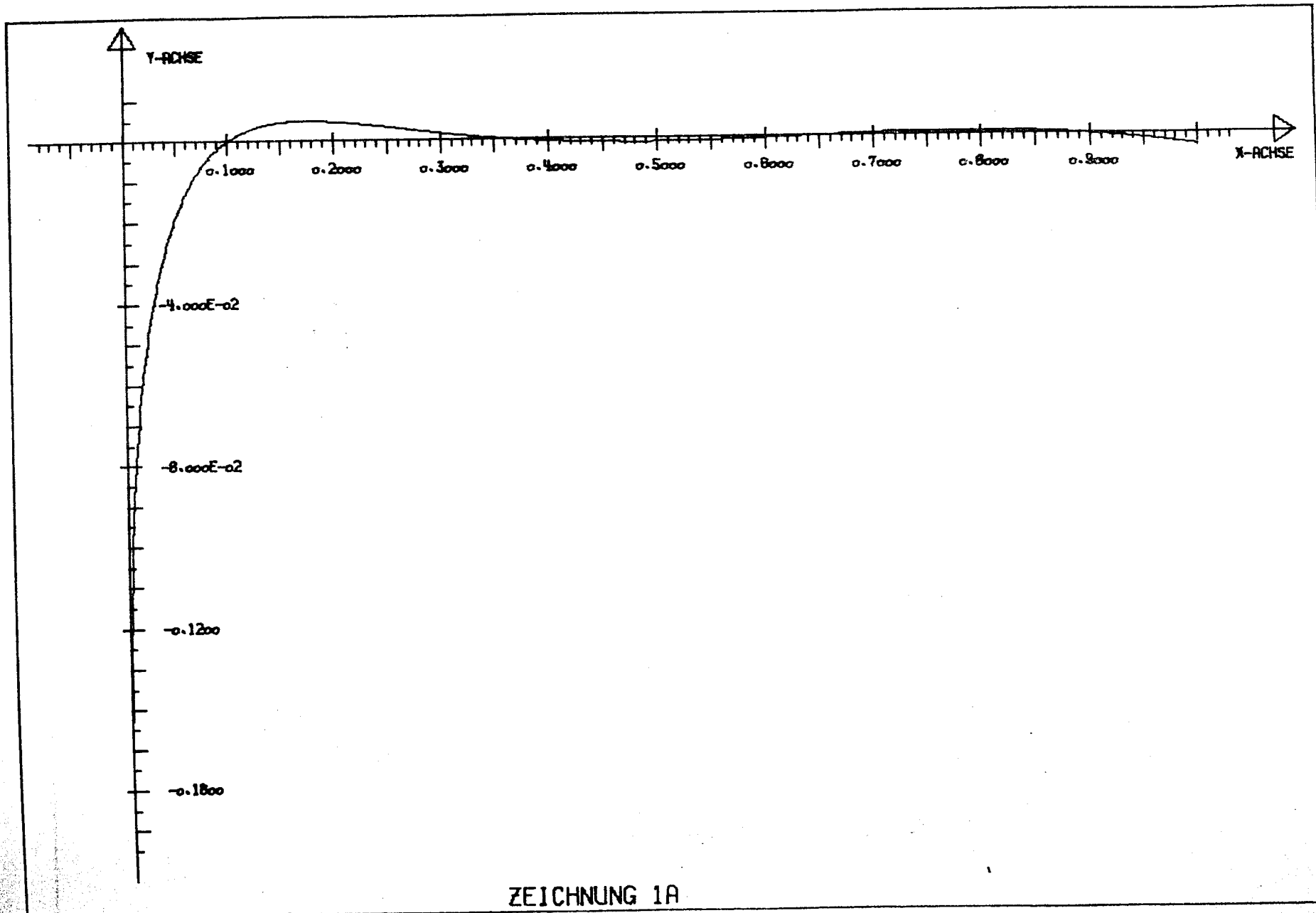
Auch die Berechnung einer Näherung nach 5.3 mit verschiedenen Werten für m führt auf komplexe Nullstellen; es sei $x_1 = -1.0$ und $x_m = 1.0$:

| m | Die Nullstellen von Q | |
|----|-----------------------|--------------|
| | Realteil | Imaginärteil |
| 10 | 0.9329 | ± 0.1892 |
| 22 | 0.9740 | ± 0.0831 |
| 32 | 0.9828 | ± 0.0567 |

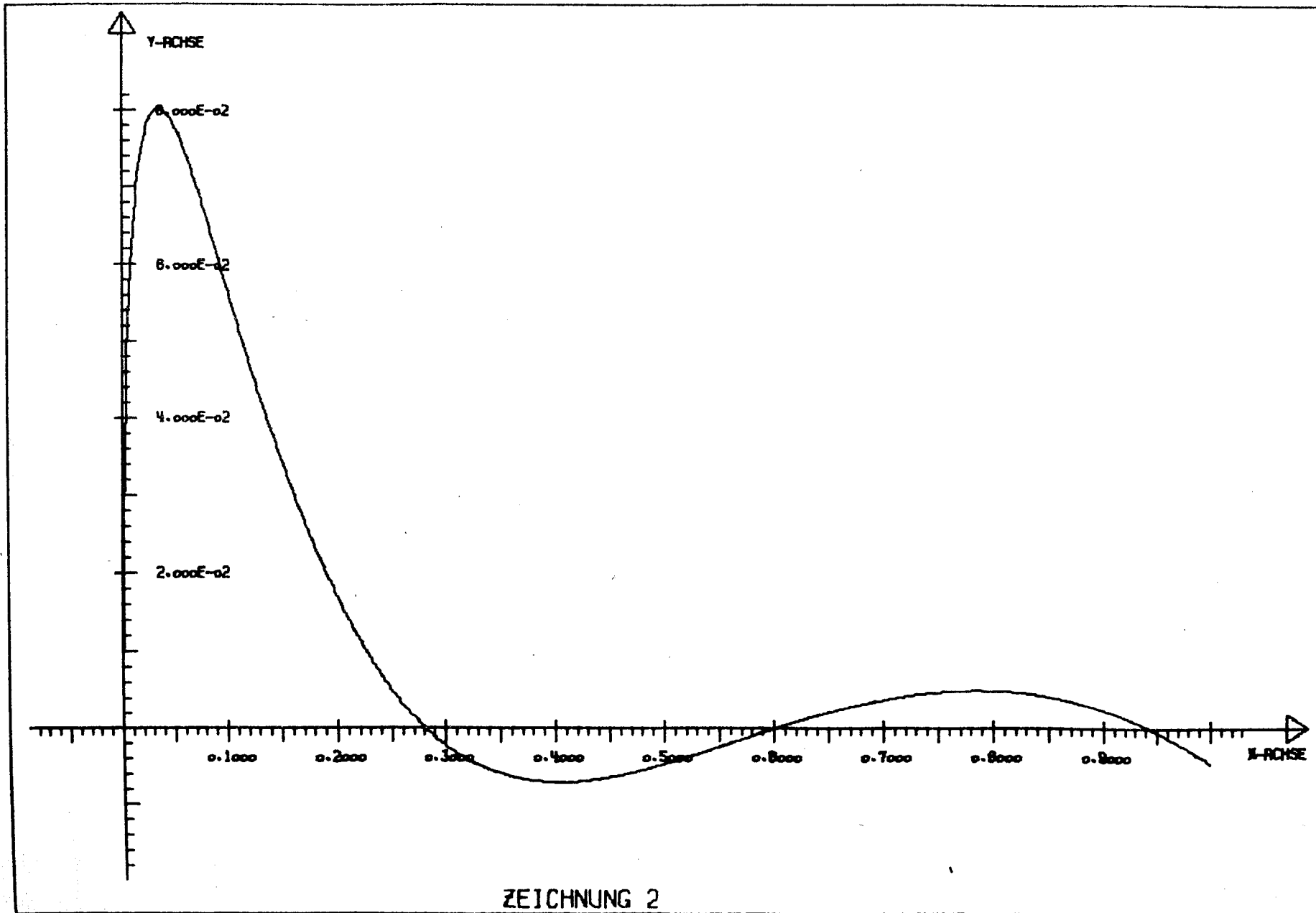
Die hier auftretenden Schwierigkeiten dürften darin begründet sein, daß es für f in $[-1,1]$ nach Satz 3.3 zwei Minimallösungen E_1, E_2 bzgl. V_2 gibt und daher keine beste Approximation bzgl. V_2^0 existiert; es liegt also nahe zu vermuten, daß $f - E_i$, $i=1,2$, in I keine Alternante der Länge 5 besitzt, was sich besonders auf die Verfahren von 5.1 und 5.2 auswirkt.

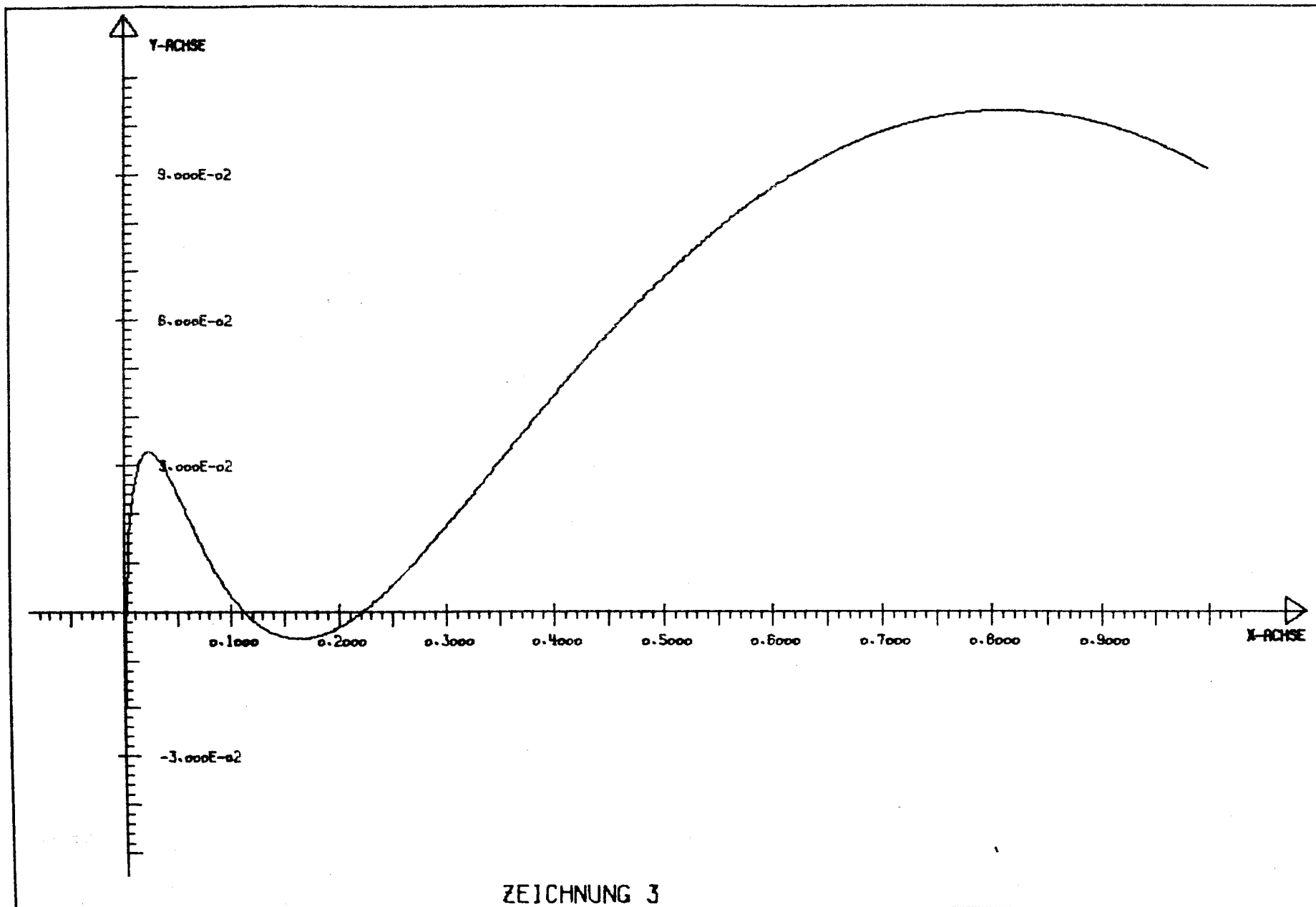
Zur Anwendung des Verfahrens von 5.1 auf $f(x) = \frac{1}{1+x}$ bzw. die Riemannsche Zetafunktion sei auf §7 bzw. §8 verwiesen.

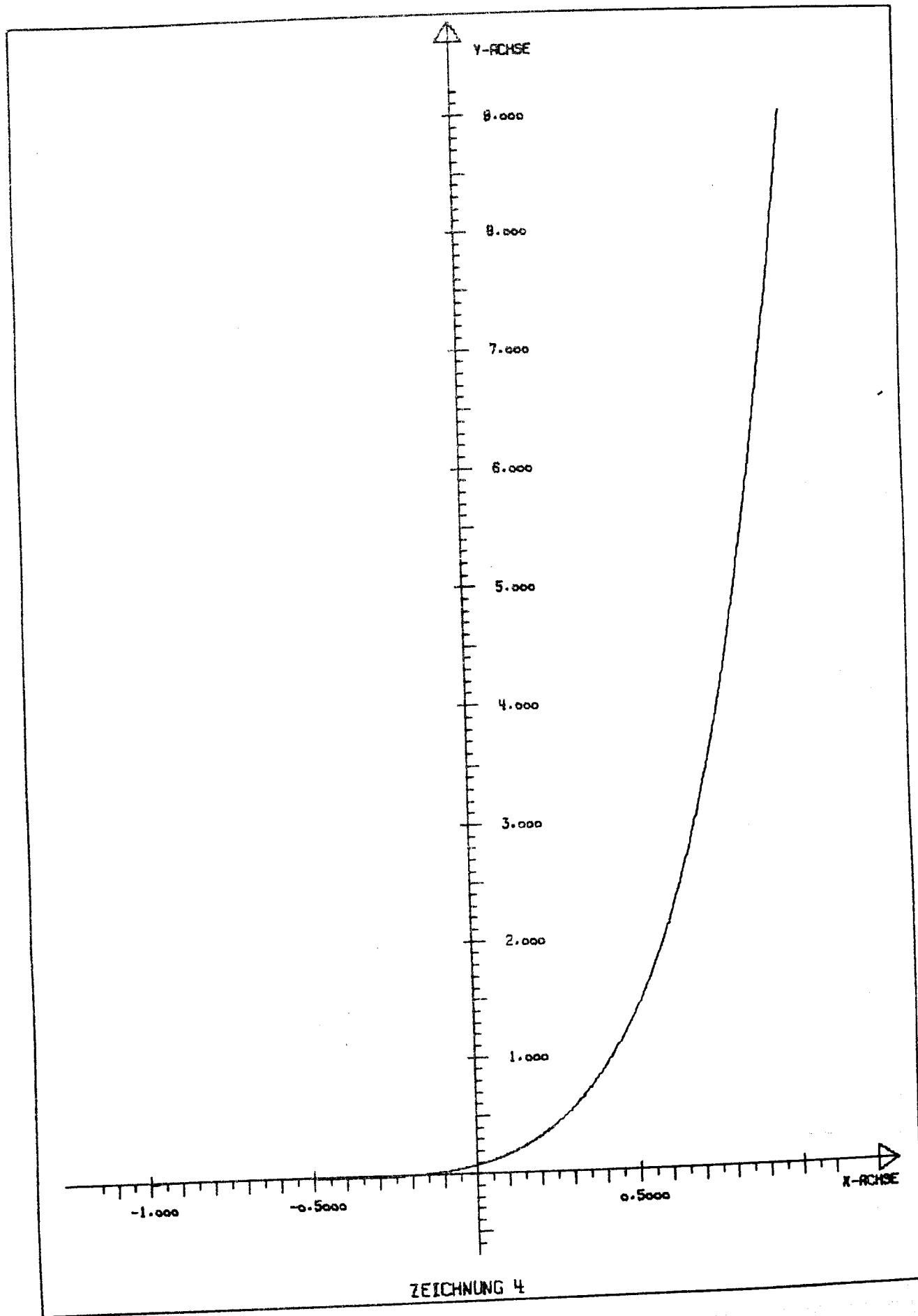


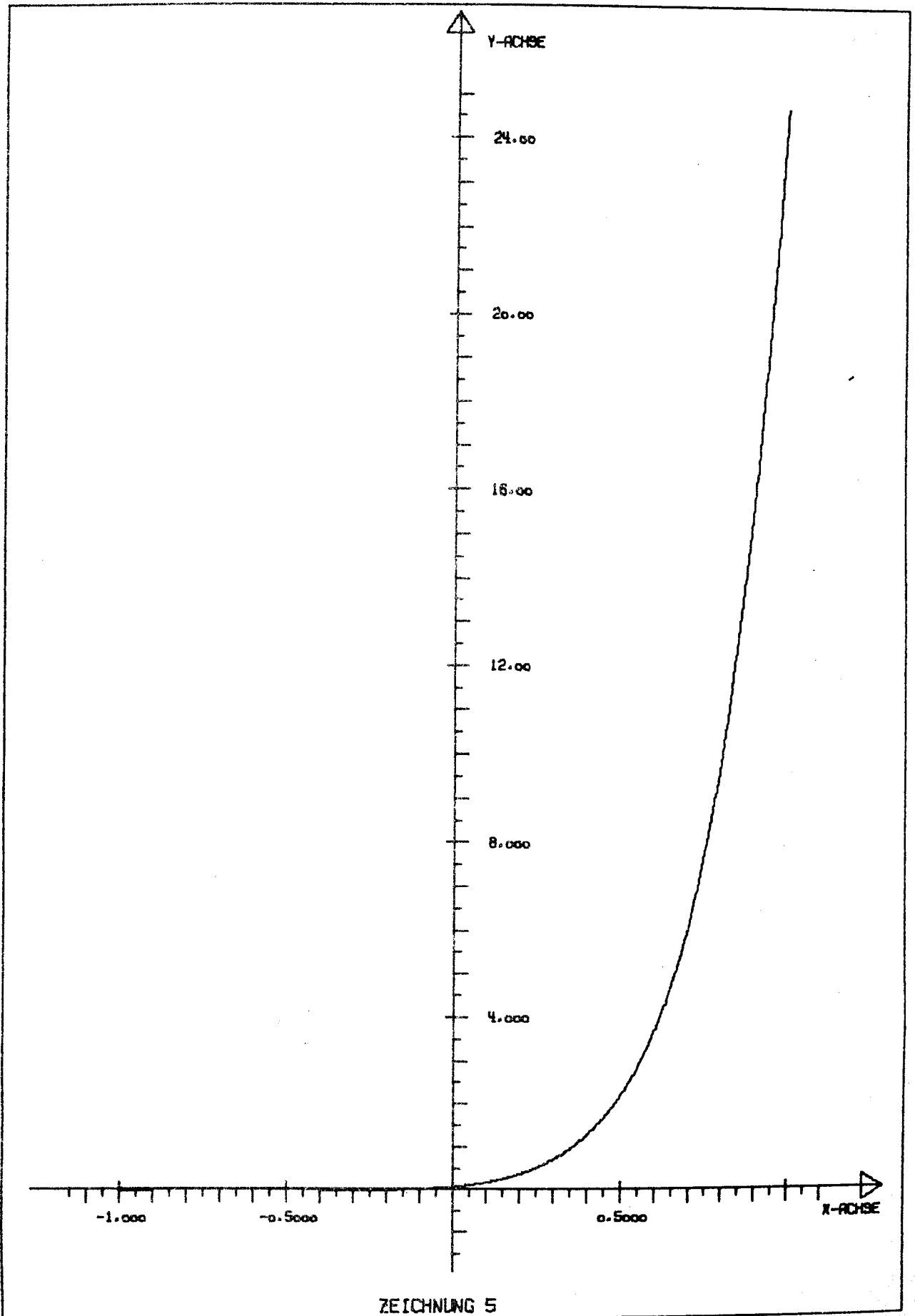


ZEICHNUNG 1A









ZEICHNUNG 5