

Institut für Angewandte und Numerische Mathematik  
TU Wien

WS 2002/03

## Numerische Mathematik für Informatiker

Prüfung am 28. Januar 2003

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)	c)		
	d)		
	e)		
	f)		

Der gesamte Rechengang ist auf den beiliegenden  
Blättern zu dokumentieren.  
Zusätzlich beigefügte Zetteln werden bei der  
Korrektur nicht berücksichtigt.

1) (10 Punkte)

Der Ausdruck

$$\varphi(x) = \frac{e^x - 1}{x^2}$$

sei für eine gegebene Maschinenzahl  $x > 0$  in Gleitpunktarithmetik (mit relativer Genauigkeit  $eps$ ) auszuwerten. Dabei sei angenommen, dass der Exponentialterm  $e^x$  mit Hilfe einer vordefinierten Prozedur `exp` “sauber” ausgewertet wird, d.h. mit einem relativen Auswertefehler  $\leq eps$ .

- a) Geben Sie eine möglichst scharfe Schranke für den sich dabei ergebenden relativen Auswertefehler an, d.h. bestimmen Sie  $\varepsilon$  so, dass<sup>1</sup>  $\tilde{\varphi}(x) = \varphi(x)(1 + \varepsilon)$ ,  $|\varepsilon| \leq ??? eps$ .

Antwort zu a):

<sup>1</sup> $\tilde{\varphi}$  symbolisiert die Auswertung von  $\varphi$  in Gleitpunktarithmetik.

- b) Für welche Werte von  $x$  ergibt sich hier ein großer relativer Fehler? Kommentieren Sie Ihr Resultat ausführlich!

Antwort zu b):

- c) Betrachten Sie nun die unter b) diagnostizierte numerisch instabile Situation. Modifizieren Sie für diesen Fall den Ausdruck  $\varphi(x)$  durch Linearisierung<sup>2</sup> des Terms  $e^x$  und schreiben Sie die sich ergebende Approximation für  $\varphi(x)$  so an, dass Sie in numerisch stabiler Weise ausgewertet werden kann. (Mit ausführlicher Begründung: worin besteht der entscheidende Unterschied zur direkten Auswertung von  $\varphi(x)$ ?)

Antwort zu c):

---

<sup>2</sup>mittels nach dem linearen Glied abgebrochener Taylor-Entwicklung von  $e^x$

2) (10 Punkte)

a) Wie löst man mehrere lineare Systeme

$$Ax_1 = b_1, \quad Ax_2 = b_2, \quad \dots, \quad Ax_k = b_k$$

mit möglichst geringem Aufwand?

(Anmerkung: Die Operatoren **factorize** für eine LU-Zerlegung und **solve** für eine Rücksubstitution sollen verwendet werden.)

Antwort zu 3a):

b) Welcher Aufwand ist erforderlich, um ein lineares Gleichungssystem  $Tx = b$  mit einer  $n \times n$  Tridiagonalmatrix effizient zu lösen?

Anzahl der arithmetischen Operationen  $O(n^p)$ :

$p =$

c) Welchen Aufwand kann man reduzieren, wenn man weiß, daß die Matrix  $A$  des linearen Systems  $Ax = b$  symmetrisch und positiv definit ist?

(i)

(ii)

(iii)

d) Was kann man mit Hilfe der Konditionszahl  $\kappa(A) = \|A\| \|A^{-1}\|$  abschätzen?

$$\dots \leq \kappa(A) \cdot \dots$$

Verbale Beschreibung:

e) Wann ist ein lineares Gleichungssystem numerisch singulär?

Antwort zu 3e):

f) Kann man durch Einsatz geeigneter Pivot-Strategien numerisch singuläre Systeme mit zufriedenstellender Genauigkeit lösen?

ja

nein

Begründung:

3) (10 Punkte)

a) Man bestimme die Koeffizienten  $a$  und  $b$  in der Funktion

$$y(t) = a \cdot t + b \cdot e^{-t} \quad (1)$$

so, dass sie dem folgenden Datensatz möglichst gut im Sinn der kleinsten quadratischen Ordinatenabweichung angepasst sind:

$t_i$	0	1	2	3
$y_i$	-2.10	1.30	7.70	17.90

(Rechnung in Taschenrechnergenauigkeit.)

$a =$

$b =$

b) Machen Sie eine Skizze der soeben errechneten Kurve und vergleichen Sie ihren Verlauf für wachsendes  $t$  mit dem Verhalten der Datenpunkte. Geben Sie einen zu (1) alternativen Ansatz, der Ihrer Meinung nach das Verhalten der Daten besser wiedergibt.

Alternativer Ansatz:

10 Punkte)

Bei Umrechnung von Währungen in Bezug auf den Euro gelten folgende Regeln!:

- Um Pfund auf Euro zu konvertieren wird der Pfund-Wert durch den Kurs dividiert.
- Um Euro auf Pfund umzurechnen wird der Euro-Wert mit dem Kurs multipliziert.
- Bei jeder Umrechnung ist optimale Rundung auf den nächsten Cent (1/100 Euro) oder Penny (1/100 Pfund) anzuwenden. (*Round away from zero* im Fall des Mittelpunktes zwischen zwei Werten.)

Den folgenden Rechnungen soll ein Kurs von 1 Euro = 0.654 325 Pfund zugrundegelegt werden.

a) Man konvertiere

Euro	Pfund
	500
	1000
	7000
1000	

b) Man berechne

	Pfund	Euro
Item 1	1000	
Item 2	1000	
Item 3	1000	
Item 4	1000	
Item 5	1000	
Item 6	1000	
Item 7	1000	
Summe		

Die Differenz zur Umrechnung von 7000 Pfund in Euro nennt man „vertikalen“ Rundungsfehler.

Vertikaler Rundungsfehler:	Euro
----------------------------	------

Institut für Angewandte und Numerische Mathematik  
TU Wien

WS 2001/02

### Numerische Mathematik für Informatiker

Prüfung am 29. Jänner 2002

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)	c)	c)	
d)			
e)			

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.

Zusätzlich beigefügte Zettel werden bei der Korrektur nicht berücksichtigt.

e) Man konvertiere die folgenden Bereichsgrenzen:

Pfund	$\leq 1000$	1000.01 – 4000	4000.01 – 8000	$\geq 8000.01$
Euro	$\leq$	–	–	$\geq$

Welches Problem tritt hier auf?

Antwort:

c) Man berechne (nach obigen Regeln) die folgenden Konversionen Euro  $\rightarrow$  Pfund  $\rightarrow$  Euro:

Euro	Pfund	Euro	Fehler
1.00			
1.01			
1.02			
1.03			
1.04			
1.05			
1.06			
1.07			
1.08			

Die bei der Konversion Euro  $\rightarrow$  Pfund  $\rightarrow$  Euro auftretenden Fehler nennt man „horizontale“ Rundungsfehler.

d) Wieso tritt bei der Konversion Pfund  $\rightarrow$  Euro  $\rightarrow$  Pfund *kein* horizontaler Rundungsfehler auf?

Antwort:



b) Lösen Sie die Aufgabe für folgenden Datensatz:

$i$	$x_i$	$y_i$	$\varphi_i$
1	0.0	0.0	0.0
2	1.1	1.1	9.9
3	1.1	2.1	14.9
4	2.1	1.1	13.9

Antwort zu b):

$$a =$$

$$b =$$

$$c =$$

c) Kommentieren Sie den Spezialfall, bei dem die  $(x_i, y_i)$  auf einer gemeinsamen Geraden liegen.

Antwort zu c):

2) (10 Punkte)

Bei linearen Ausgleichsproblemen, bei denen eine „Ausgleichsebene“  $f(x, y) = ax + by + c$  gesucht ist, sind die Parameter  $a, b, c \in \mathbb{R}$  so zu bestimmen, dass

$$\sum_{i=1}^n [f(x_i, y_i) - \varphi_i]^2 \rightarrow \text{minimal} \quad (\text{LS})$$

für einen gegebenen Satz von  $n$  Datenpunkten  $\{(x_i, y_i; \varphi_i), i = 1, 2, \dots, n\}$  gilt.

a) Wie groß muss die Anzahl  $n$  der zur Verfügung stehenden Datenpunkte mindestens sein, damit das Problem (LS) eine eindeutige Lösung  $(a, b, c)$  besitzen kann? Ist diese Bedingung auch hinreichend? (Begründung!)

Antwort zu a):

$$n \geq$$

Diese Bedingung ist  hinreichend

nicht hinreichend

Begründung:

b) Will man den Wert  $f(x)$  für ein  $x$  mit  $x_i < x < x_{i+1}$  berechnen, so bezieht man die Tabellenwerte  $(x_i, f_i = f(x_i))$ ,  $(x_{i+1}, f_{i+1} = f(x_{i+1}))$  ein und wertet die folgende Formel mit  $h = x_{i+1} - x_i$  aus:

$$f(x) = \frac{x_{i+1} - x}{h} f_i - \frac{x_i - x}{h} f_{i+1}.$$

Berechnen Sie die Schranke für den absoluten Rechenfehlereffekt, der bei der Auswertung des obigen Ausdruckes in einer Gleitpunktarithmetik entsteht. Nehmen Sie dabei an, dass bei den Subtraktionen  $x_{i+1} - x$  und  $x_i - x$  keine Rechenfehler entstehen und alle involvierten Größen Maschinenzahlen sind. Die Rechnungen werden in der IEC/IEEE-Arithmetik, im einfachen Format mit optimaler Rundung durchgeführt.

Hinweis: Im Intervall  $[0, 2]$  ist  $f(x)$  positiv und monoton fallend.

Antwort zu b):

3) (10 Punkte)

Für die Funktion

$$f(x) = \frac{1}{\cosh(x)} = \frac{2}{e^x + e^{-x}}, \quad x \in [0, 2]$$

wird eine Wertetabelle mit der Eigenschaft benötigt, dass bei einer linearen Interpolation zwischen zwei Tabellenwerten für alle  $x$  aus dem gegebenen Intervall maximal ein absoluter Gesamtfehler (Dateifehler + Verfahrensfehler + Rechenfehler) von  $10^{-6}$  auftritt.

a) Wie groß darf der Abstand  $\Delta x$  der Tabellenwerte höchstens sein, damit der Verfahrensfehler auf dem gesamten Intervall  $5 \times 10^{-7}$  nicht übersteigt? Der Abstand sollte eine Zehnerpotenz sein (oder die Form  $5 \times$  Zehnerpotenz haben).

Als Vorüberlegung zeigen Sie, dass folgendes gilt:

$$f''(x) = \frac{2 \sinh^2 x - \cosh^2 x}{\cosh^3 x}$$

$$|f''(x)| \leq \frac{3}{\cosh x}$$

Um letztere Ungleichung zu zeigen, benutzen Sie die Abschätzung  $\sinh^2 x \leq \cosh^2 x$ .

Antwort zu a):

$\Delta x \leq$

- c) Kann man garantieren, dass die Schranke für den Gesamtfehler bei der obigen Interpolation gleich  $10^{-6}$  ist, falls für die Datenfehler aller Tabellenwerte  $|\tilde{f}_i - f_i| \leq 3 \cdot 10^{-7}$  gilt? Wenn nicht, welche Maßnahme(n) würden Sie ergreifen, um die Toleranzforderung für den absoluten Gesamtfehler von  $10^{-6}$  zu erfüllen?

Antwort zu c):

**Numerische Mathematik für Informatiker**

Prüfung am 26. Juni 2001

Vorname	Kennzahl / Matrikelnummer
---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)	c)	c)	
d)			
e)			

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.

Zusätzlich beigefügte Zettel werden bei der Korrektur nicht berücksichtigt.

1) (10 Punkte)

Das Gleitpunktzahlensystem der Rechnerfamilie IBM System/390 ist hexadezimal:  $\mathbb{F} = \mathbb{F}(16, 6, -64, 63)$  (einfache Genauigkeit, 4 Byte Format). Gleitpunktzahlen  $x = \pm m \cdot 16^e \in \mathbb{F}$  sind codiert als

CCMMMMM,

wobei die Hexadezimal-Ziffern<sup>1</sup>  $C, M \in \{0, \dots, 9, A, \dots, F\}$  wie folgt zu interpretieren sind:

- Im führenden Bit des ersten Bytes CC ('Charakteristik') ist das Vorzeichen von  $x$  codiert gemäß
- Die restlichen 7 Bits der Charakteristik CC entsprechen einer Binärzahl  $c$  zwischen 0 und  $127 = 7F_{\text{HEX}}$ , die den Exponenten  $e$  der Zahl  $x$  gemäß

$$e := c - 64 = c - 40_{\text{HEX}}$$

festlegt.

- Die restlichen 3 Bytes bilden die 6-stellige hexadezimale Mantisse  $m$ . Wir verwenden hier die Schreibweise mit  $m \in [0, 1)$ , d.h.  $m$  steht für die hexadezimale Festpunktzahl

.MMMMMM<sub>HEX</sub>.

Beispiel: Die Dezimalzahlen 10, -10 besitzen in diesem Format die Darstellung 41A0000, C1A00000.

Antworten zu folgenden Fragen bitte in Form der hexadezimalen Codierung CMMMMMM:

- a) Wie lauten die kleinste und die größte positive normalisierte Zahl  $x_{\min}, x_{\max} \in \mathbb{F}$ ?

Antwort:

- b) Was passiert bei der Gleitpunkt-Division  $1. \oslash x_{\min}$ ?

Antwort:

- c) Wie lautet die Gleitpunktzahl  $eps$  (relative Maschinengenauigkeit) bezüglich der Rundung gegen 0 (Rundung durch Abschneiden)?

Antwort:

- d) Wie lautet die kleinste positive Maschinenzahl  $x$ , für die die Gleitpunkt-Subtraktion  $1. \ominus x$  einen Wert  $< 1$  ergibt?

Antwort:

- e) Welche Zahl  $x \in \mathbb{F}$  liegt der Dezimalzahl  $-0.1$  am nächsten?

Antwort:

<sup>1</sup>mit der naheliegenden 4-Bit-Codierung als Binärzahl:  $0 = 0000, 1 = 0001, \dots, F = 1111$

2) (10 Punkte)

a) Gegeben sei das folgende, sehr speziell strukturierte lineare Gleichungssystem:

$$\begin{pmatrix} a & a & a & a & \dots & a \\ 1 & a & a & a & \dots & a \\ 0 & 1 & a & a & \dots & a \\ 0 & 0 & 1 & a & \dots & a \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & 1 & a \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \quad (1)$$

Man betrachte für diesen einfachen Fall die Gauß-Elimination (ohne Pivotstrategie), die zur Faktorisierung  $LUx = b$  führt und gebe die Matrizen  $L$  und  $U$  an:

$$L = \begin{pmatrix} \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \end{pmatrix}$$

$$U = \begin{pmatrix} \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \end{pmatrix}$$

Weiters gebe man  $x_n$  als Formel Ausdruck in  $a$  an:

$$x_n =$$

b) Für den Spezialfall  $a = 2$  berechne man  $U^{-1}$  und  $L^{-1}$ .

$$L^{-1} = \begin{pmatrix} \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \end{pmatrix}$$

$$U^{-1} = \begin{pmatrix} \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \end{pmatrix}$$

Neben (1) (im Spezialfall  $a = 2$ ) betrachte man nun das gestörte System

$$\begin{pmatrix} 2 + \epsilon_{11} & 2 + \epsilon_{12} & 2 + \epsilon_{13} & \dots & 2 + \epsilon_{1n} \\ 1 + \epsilon_{21} & 2 + \epsilon_{22} & 2 + \epsilon_{23} & \dots & 2 + \epsilon_{2n} \\ \epsilon_{31} & 1 + \epsilon_{32} & 2 + \epsilon_{33} & \dots & 2 + \epsilon_{3n} \\ \epsilon_{41} & \epsilon_{42} & 1 + \epsilon_{43} & \dots & 2 + \epsilon_{4n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \epsilon_{n1} & \epsilon_{n2} & \dots & 1 + \epsilon_{n,n-1} & 2 + \epsilon_{nn} \end{pmatrix} \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \tilde{x}_3 \\ \tilde{x}_4 \\ \vdots \\ \tilde{x}_n \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$

Unter der Annahme  $|\epsilon_{ij}| \leq 10^{-8}$  schätze man  $\frac{\|\tilde{x} - x\|_\infty}{\|\tilde{x}\|_\infty}$  ab.

$$\frac{\|\tilde{x} - x\|_\infty}{\|\tilde{x}\|_\infty} \leq$$

Hinweis: Für das Produkt  $M \cdot N$  zweier Matrizen  $M$  und  $N$  gilt

$$\|M \cdot N\| \leq \|M\| \cdot \|N\|.$$

3) (10 Punkte)

Die CO-Monatsmittelwerte der Luft in Wien sollen durch eine Funktion der Form

$$\text{CO}(t) = a_0 + a_1 \sin \frac{2\pi t}{12} + a_2 \cos \frac{2\pi t}{12}, \quad t \text{ in Monaten}$$

nach der Methode der kleinsten Quadrate approximiert werden. Folgende Daten (aus dem Jahr 1993) bilden hierfür die Grundlage:

Monat	1	2	3	4	5	6	7	8	9	10	11	12
CO [mg/m <sup>3</sup> ]	1.3	1.4	1.0	0.8	0.7	0.6	0.5	0.6	0.7	1.1	1.2	1.3

- Man stelle die Normalgleichungen auf.
- Man berechne  $a_0$ ,  $a_1$  und  $a_2$ .
- Man skizziere den Funktionsverlauf und trage die Datenpunkte ein. Was läßt sich über die Approximationsqualität sagen?

**Numerische Mathematik für Informatiker**

Prüfung am 30. Januar 2001

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)		a)	
b)		b)	
c)		c)	
d)		d)	

1) (10 Punkte)

Es sei  $\mathbb{F} = \mathbb{F}(10, 9, -98, 100)$  (eine typische Taschenrechner-Arithmetik mit 9 signifikanten Dezimalstellen), mit optimaler Rundung.  $p$  bezeichne denjenigen Wert der durch Rundung von  $\pi = 3.141592653\dots$  nach  $\mathbb{F}$  entsteht.

Im Folgenden bedeutet das Symbol ' $*$ ' Gleitpunkt-Multiplikation in  $\mathbb{F}$ . Alle Antworten sind genau zu begründen.

- a) Wie lautet der größte Wert  $K_1 \in \mathbb{N}$  mit der Eigenschaft dass der Wert  $k * p \in \mathbb{F}$  für alle  $k \leq K_1$  ( $k \in \mathbb{N}$ ) mit dem exakten Wert  $k \cdot p$  übereinstimmt?

Antwort zu a):

- b) Wie lautet der größte Wert  $K_2 \in \mathbb{N}$  so dass

$$k * p \in I(k) := [k \cdot \pi - 10^{-5}, k \cdot \pi + 10^{-5}]$$

für alle  $k \leq K_2$  ( $k \in \mathbb{N}$ ) garantiert zutrifft?

Antwort zu b):

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.

Zusätzlich beigefügte Zetteln werden bei der Korrektur nicht berücksichtigt.

c) Bestimmen Sie die maximale relative Konditionszahl für die beiden Funktionen

(i)  $f(x) = \sin(x)$  und (ii)  $g(x) = \cos(x)$

für  $x \in I(k)$  (das Intervall  $I(k)$  ist wie unter b) definiert). Antwort erbeten als Ausdruck in  $k$  (für beliebiges  $k \in \mathbb{N}_0$ ).

(Tipp des Tages: Man verwende die Approximation  $|\tan \delta| \approx |\delta|$  für kleine Werte  $|\delta|$ .)

Antwort zu c):

d) Kommentieren Sie Ihr unter c) erhaltenes Ergebnis.

Für beide Funktionen stellt  $k = 0$  einen Sonderfall dar. In welchem Sinn jeweils?

Antwort zu d):

2) (10 Punkte)

In dem Gleichungssystem  $Ax = b$  hat die Matrix  $A \in \mathbb{R}^{n \times n}$  die Gestalt

$$A = \begin{pmatrix} 1 + \delta_{11} & \delta_{12} & \delta_{13} & \dots & \delta_{1,n-1} & 1 + \delta_{1,n} \\ \delta_{21} & 1 + \delta_{22} & \delta_{23} & \dots & \delta_{2,n-1} & 1 + \delta_{2,n} \\ \delta_{31} & \delta_{32} & 1 + \delta_{33} & \dots & \delta_{3,n-1} & 1 + \delta_{3,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \delta_{n-1,1} & \delta_{n-1,2} & \delta_{n-1,3} & \dots & 1 + \delta_{n-1,n-1} & 1 + \delta_{n-1,n} \\ 1 + \delta_{n,1} & 1 + \delta_{n,2} & 1 + \delta_{n,3} & \dots & 1 + \delta_{n,n-1} & n + \delta_{n,n} \end{pmatrix}$$

mit  $|\delta_{ij}| \leq 10^{-8}$ . Es sei  $\bar{x}$  die Lösung des vereinfachten Systems  $\bar{A}\bar{x} = b$  mit

$$\bar{A} = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 1 \\ 0 & 1 & 0 & \dots & 0 & 1 \\ 0 & 0 & 1 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 1 \\ 1 & 1 & 1 & \dots & 1 & n \end{pmatrix}$$

Man schätze  $\frac{\|\bar{x} - x\|}{\|x\|}$  bezüglich der Maximumnorm  $\|\cdot\|_\infty$  möglichst scharf ab!

$$\frac{\|\bar{x} - x\|}{\|x\|} \leq$$



3) (10 Punkte)

Für die Funktion

$$f(x) = \sinh(x) = \frac{e^x - e^{-x}}{2}, \quad x \in [0, 2]$$

wird eine Wertetabelle benötigt, mit der Eigenschaft, daß bei einer linearen Interpolation zwischen zwei Tabellenwerten für alle  $x$  aus dem gegebenen Intervall ungünstigstenfalls ein absoluter Fehler von von  $10^{-6}$  auftritt.

a) Wie groß darf der Abstand der Tabellenwerte höchstens sein, damit der Verfahrensfehler auf dem gesamten Intervall  $3 * 10^{-7}$  nicht übersteigt? Der Abstand sollte eine Zehnerpotenz sein oder die Form  $5 * \text{Zehnerpotenz}$  haben.

Antwort zu a):

c) Will man den Wert  $f(x)$  für ein  $x$  mit  $x_i < x < x_{i+1}$  berechnen, so bezieht man die Tabellenwerte  $(x_i, f_i = f(x_i))$ ,  $(x_{i+1}, f_{i+1} = f(x_{i+1}))$  ein und wertet die folgende Formel mit  $h = x_{i+1} - x_i$  aus:

$$f(x) = \frac{x_{i+1} - x}{h} f_i - \frac{x_i - x}{h} f_{i+1}.$$

Berechnen Sie die Schranke für den absoluten Rechenfehler, der bei der Auswertung des obigen Ausdruckes in einer Gleitpunktarithmetik entsteht. Nehmen Sie dabei an, daß bei den Subtraktionen  $x_{i+1} - x$  und  $x_i - x$  keine Rechenfehler entstehen und alle involvierten Größen Maschinenzahlen sind. Die Rechnungen werden in der IEC/IEEE-Arithmetik, im einfachen Format mit optimaler Rundung durchgeführt.

Hinweis: Im Intervall  $[0, 2]$  ist  $f(x)$  positiv und monoton wachsend.

Antwort zu c):

b) Wie groß darf die Datengenauigkeit (in den Funktionswerten) sein, damit der dadurch entstehende Datenfehlereffekt kleiner als  $10^{-7}$  ist?

Antwort zu b):

d) Kann man garantieren, daß die Schranke für den Gesamtfehler bei der obigen Interpolation kleiner als  $10^{-6}$  ist? Wenn nicht, welche Maßnahme(n) würden Sie ergreifen, um die Toleranzforderung für den absoluten Gesamtfehler von  $10^{-6}$  zu erfüllen?

Antwort zu d):

Numerische Mathematik für Informatiker

Prüfung am 1. Dezember 1998

Name	Kennzahl / Matrikelnummer
Vorname	

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)		c)	
		d)	

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.  
Zusätzlich beigefügte Zettelchen werden bei der Korrektur nicht berücksichtigt.

1) (10 Punkte)

Die Nullstellen  $x_1$  und  $x_2$  der quadratischen Gleichung  $x^2 + bx + c = 0$  seien bekannt. Jemand möchte daraus die Koeffizienten  $b = -(x_1 + x_2)$  und  $c = x_1 \cdot x_2$  berechnen und anschließend die Größen  $d = c/b$  und  $e = b/c$ . Bei allen Überlegungen dieses Beispiels werde angenommen, daß  $x_1$  und  $x_2$  reelle Zahlen sind, für welche  $\frac{1}{10} \leq |x_i| \leq 1$  ( $i = 1, 2$ ) gilt und daß die Berechnung von  $b, c, d$  und  $e$  auf einem Taschenrechner mit der Arithmetik  $M(10, 10, +99, -99)$  durchgeführt wird (sodaß sowohl bei der Eingabe von  $x_1$  und  $x_2$  als auch bei der weiteren Rechnung entsprechende Rundungsfehler auftreten). Weiters nehme man an, daß nach jedem Rechenschritt abgeschnitten wird und daß keine Schutzstellen vorhanden sind.

- a) Man gebe ein Nullstellenpaar  $x_1$  und  $x_2$  an, sodaß sich bei der Berechnung von  $b$  nur vier richtige Stellen ergeben.
- b) Nun betrachte man die absoluten Rundungsfehler, die bei der Berechnung von  $b, c, d$  und  $e$  entstehen. Konkret untersuche man, für welche der Größen  $b, c, d$  oder  $e$  es Nullstellenpaare  $x_1$  und  $x_2$  gibt, sodaß das absolute Fehlerniveau
  - (i) in der Größenordnung von  $10^{-8}$  liegt
  - (ii) in der Größenordnung von  $10^0$  liegt
  - (iii) unendlich groß ist.
 (Kurze Begründung)
- c) Man gebe (mittels '(1+ε)-Technik') eine möglichst scharfe Rundungsfehlerschranke für den absoluten Rundungsfehler bei der Berechnung von  $e$  an, die für sämtliche Nullstellenpaare  $x_1, x_2$  mit  $\frac{1}{10} \leq |x_i| \leq 1$  gilt.

2) (10 Punkte)

In der graphischen Datenverarbeitung spielen Spline-Funktionen eine wichtige Rolle bei der Darstellung punktförmig gegebener Daten. Bei der Berechnung einer periodischen Spline-Interpolierenden treten lineare Gleichungssysteme der folgenden Bauart auf:

- Welche Bauart hat die obere Dreiecksmatrix  $U$  bei der Dreieckszerlegung? Was folgt daraus für deren Speicherung?
- Man formuliere einen Algorithmus, der die Lösung des obigen Gleichungssystems – ohne Pivottisierung – unter Berücksichtigung der speziellen Bauart der Matrix liefert.

3) (10 Punkte)

Um den Fehler bei der doppelt genauen Gleitpunktdivision in der Steinzeitversion des Pentium-Prozessors zu umgehen, faßt Fred folgende softwaremäßige Lösung ins Auge:

Um eine vorgegebene doppelt genaue Gleitpunktzahl  $x$  zu invertieren, wählt Fred einen geeigneten Näherungswert  $y_0$  für  $y := 1/x$  (also z.B. den durch einfach genaue Rechnung erhaltenen Wert) und will ausgehend von  $y_0$  die zu  $y = 1/x$  äquivalente Gleichung  $f(y) = xy - 1 = 0$  iterativ in doppelt genauer Arithmetik lösen. Er meint, daß sich leicht eine Iterationsvorschrift finden lassen müßte, die sehr rasch konvergiert aber keinerlei explizite Divisionen erfordert.

- Fred denkt an das *Newton-Verfahren* für  $f(y) = 0$ . Kommentieren Sie diese Idee.
- Wilma denkt über Freds Vorschlag nach und weiß es besser: Sie meint, es gibt ja auch vereinfachte Varianten des Newton-Verfahrens, bei denen man die in der Iterationsvorschrift auftretende Ableitung nicht exakt bildet sondern dafür einen geeigneten (von  $y_0$  abhängigen) Näherungswert verwendet.  
Zeigen Sie, in welcher Weise Wilmas Idee zu einer Iterationsvorschrift führt, die tatsächlich divisionsfrei abläuft und eine sehr gute Konvergenz gegen  $y = 1/x$  erwarten läßt (Begründung!).  
Schreiben Sie die so entstehende Iteration so an, daß sie (nach einer Initialisierungsphase) in jedem Iterationsschritt nur eine Addition und eine Multiplikation erfordert.
- Geben Sie für die unter b) eingeführte Iteration eine Abschätzung für die Konvergenzgeschwindigkeit an, d.h. bestimmen Sie  $\rho$  so, daß

$$|y_{i+1} - y| \leq \rho |y_i - y|$$

- gilt. Wodurch ist die Konvergenzgeschwindigkeit, d.h. die Kleinheit von  $\rho$  bestimmt?
- Approximieren Sie den Wert  $\varphi(x) := 1/x$  mittels Taylor-Entwicklung von  $\varphi$  bis zum linearen Term (Entwicklung um  $x_0 := 1/y_0$ ). Wie hängt die so entstehende Näherung für  $1/x$  mit obiger Iteration zusammen?

<sup>1</sup>Hier bezeichnet  $y = 1/x$  das exakte Ergebnis; Rundungsfehlereffekte bei der Iteration bleiben außer Betracht.

# Numerische Mathematik f. Informatiker

Prüfung am 6. Mai 1998

- (1.) Ausdruck  $x^2 + 3x + 2$   
(16) 2 Auswertungsvarianten  
result 1 =  $(2 + 3x) + x \cdot x$   
result 2 =  $2 + (x \cdot (3 + x))$

! (Achtung! Beispiel war ungefähr so aus dem Gedächtnis!)

a.) Für result 2 (Maschinenparenzigkeit eps) soll der relative Rundungsfehler berechnet werden.  
Begründung für den Fehler!

b.) ~~Wie~~ Ist result 1 besser? Ja/Nein? Warum! → nur Begründung keine Berechnung

c.) Gibt es eine bessere Variante → dafür Rundungsfehler berechnen.

(2.) Matrizenbeispiel Nr. 2 von der VO-Prüfung am 27. 1. 97  
(20)

(3.) Polygonzug - Fehlerberechnung - Beispiel Nr. 3 von ——— " ———  
(14)

Numerische Mathematik für Informatiker

Prüfung am 12. März 1998

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)	c)	c)	

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.  
Zusätzlich beigefügte Zetteln werden bei der Korrektur nicht berücksichtigt.

1) (10 Punkte)  
Der Ausdruck

$$\varphi(x) = \frac{\sqrt{x+1} - \sqrt{x}}{10} = \frac{1}{10(\sqrt{x+1} + \sqrt{x})}, \quad x > 1,$$

soll für große Werte von  $x$  berechnet werden.

a) Man gebe die relative Konditionszahl  $K_{\varphi(x)-x}$  als Funktion von  $x$  an.

Man gebe für  $K_{\varphi(x)-x}$  eine möglichst gute, von  $x$  unabhängige Schranke an, d.h. man bestimme ein möglichst kleines  $C > 0$ , sodaß

$$|K_{\varphi(x)-x}| \leq C, \quad \forall x > 1$$

gilt. Welche Aussage kann man demgemäß für die Kondition des Problems machen?

b) Man gebe für beide Auswertungsvarianten möglichst gute Fehlerschranken für den relativen Rundungsfehler als Ausdrücke in  $x$  und in der Maschinengenauigkeit  $\epsilon_{ps}$  an, d.h. man bestimme für beide Fälle ein  $\epsilon$ , mit

$$\text{rnd}(\varphi(x)) = \varphi(x)(1 + \epsilon)$$

und schätze es in Abhängigkeit von  $x$  und  $\epsilon_{ps}$  ab.

*Hinweis:* Der Einfachheit halber nehme man an, daß für die Auswertung der Wurzelfunktionen

$$\text{rnd}(\sqrt{1+x}) = \sqrt{1+x}(1 + \epsilon_1), \quad \text{rnd}(\sqrt{x}) = \sqrt{x}(1 + \epsilon_2)$$

mit  $|\epsilon_1| \leq 2\epsilon_{ps}$  und  $|\epsilon_2| \leq \epsilon_{ps}$  gilt und vernachlässige den Rundungsfehler bei der Division durch 10.

c) Welche der beiden Auswertungsvarianten ist für große Werte von  $x$  vorzuziehen? Ausführliche Begründung!

*Problem ist vorzeichenlos  $\sqrt{x}$  ist groß!*

<sup>1</sup>  $\text{rnd}(\varphi(x))$  symbolisiert die jeweilige Auswertung von  $\varphi$  in Gleitpunktarithmetik.

Tragen Sie die Ergebnisse ein und kreuzen Sie die zutreffenden Aussagen an !

$$K_{\varphi(x)-x} =$$

Für  $x > 1$  ist  $C =$

und deshalb ist die Auswertung  gut  schlecht konditioniert.

Fehlerschranke für den relativen Rundungsfehler von  $\varphi_1(x) := \frac{\sqrt{x+1}-\sqrt{x}}{10}$ :

$$|\varepsilon| \leq$$

Fehlerschranke für den relativen Rundungsfehler von  $\varphi_2(x) := \frac{1}{10(\sqrt{x+1}+\sqrt{x})}$ :

$$|\varepsilon| \leq$$

Die Formel   $\varphi_1(x)$    $\varphi_2(x)$  ist bei der Auswertung von  $\varphi(x)$  für große  $x$ -Werte vorzuziehen.

Begründung:

2) (10 Punkte)

Gegeben sei ein 4-dimensionales lineares Gleichungssystem  $Ax = b$ . Die konkrete Gestalt von  $A$  und  $b$  tut hier nichts zur Sache; es sei nur bekannt, daß gilt

$$\|A\|_2 = \frac{77}{60}, \quad \|A^{-1}\|_2 = 63420.$$

a) Wie groß ist die relative Konditionszahl von  $A$  bezüglich der Norm  $\|\cdot\|_2$ ?

Antwort zu a):

81588

b) Jemand verwendet ein Iterationsverfahren (z.B. Jacobi-Verfahren) und erhält nach einigen Schritten eine Näherungslösung  $\tilde{x}$ . Diese sieht vernünftig aus, denn ihr Residuum  $A\tilde{x} - b$  ist klein, und zwar (auf zwei Stellen genau angegeben):

$$A\tilde{x} - b = (3.0E-4, 3.0E-4, 5.0E-5, 2.0E-4)^T \quad \| \Delta b \| = \| b - \tilde{b} \|$$

Dieser Jemand gibt sich mit der Näherungslösung  $\tilde{x}$  zufrieden, erkennt aber später, daß sie extrem ungenau ist: In keinem der  $\tilde{x}_i$  ist auch nur eine Dezimalstelle richtig.

Erklären Sie die Ursache für diesen Effekt mittels ausführlicher schriftlicher Begründung (ohne Rechnung). Warum ist hier der Fehler noch so groß, obwohl das Residuum schon so klein ist?

Antwort zu b):

$$Y = Ax$$

$\|A\|$  groß: kleines  $\Delta x \rightarrow$  großes  $\Delta Y$

$\|A^{-1}\|$  groß: kleines  $\Delta Y \rightarrow$  großes  $\Delta x$

$$Ax = b \quad (3.0E-4, \dots) = \tilde{b} - b$$

$$A\tilde{x} = \tilde{b}$$

$$\frac{\|\tilde{b} - b\|_2}{\|b\|_2} \leq \|A\|_2 \cdot \frac{\|\tilde{x} - x\|_2}{\|x\|_2}$$

$$\frac{\|\tilde{x} - x\|_2}{\|x\|_2} \leq \underbrace{\|A^{-1}\|_2}_{\text{riesig}} \cdot \underbrace{\frac{\|\tilde{b} - b\|_2}{\|b\|_2}}_{\text{klein}}$$

relativ

$\|A\|$  groß: Ur-Bildbereich groß

$\|A^{-1}\|$  groß: Bild  $\rightarrow$  Urbild groß (sing. lösbar!)

$$Ax = b \rightarrow \tilde{x} \rightarrow A\tilde{x} - b = 0$$

Residuum

$\|A\|$  groß  $\rightarrow$  Residuum ist (relativ) immer groß

c) Geben Sie - unter Zuhilfenahme der über das Problem vorliegenden Information - an, wie groß dieser Fehler maximal gewesen sein kann: d.h. man schätze die Norm  $\|\tilde{x} - x^*\|_2$  nach oben ab? (Rechnung inklusive Begründung).

Antwort zu c):

$$\|\tilde{x} - x\| \leq \|A^{-1}\|_2 \cdot \|\tilde{b} - b\|_2$$

aus Angabe b.)  
ausrechnen  
(quadratisches, Summe, Wurzel, ...)

$$\|\tilde{x} - x\| \leq 63420$$

Begründung:

\*  $x^*$  bezeichnet die exakte Lösung des Gleichungssystems.

5/14/16

3) (10 Punkte)

Für die graphische Ausgabe von Funktionen sollen diese durch Polygonzüge (stückweise lineare Funktionen) ersetzt werden.

a) Die Funktion  $f(x) = \cos(x)$  soll am Intervall  $[0, 1.5]$  durch einen Polygonzug  $p$  dargestellt werden, der äquidistante (gleichabständige) Knoten besitzt und

$$|p(x) - f(x)| \leq 10^{-3}, \quad x \in [0, 1]$$

erfüllen soll. Wie groß ist der Abstand  $h$  der Knoten (Datenpunkte) zu wählen?

b) Wie kann man den Abstand  $h$  bestimmen, wenn anstelle des absoluten Fehlers der relative Fehler unter einer vorgegebenen Schranke liegen soll?

c) Für den Fall, daß man von der darzustellenden Funktion  $f$  keine Information (Ableitungen etc.) besitzt, die sich zur Fehlerabschätzung eignet, sind Vorschläge zu machen, wie man die algorithmische Bestimmung eines geeigneten Knotenabstandes durchführen könnte.

Polygone  
1. Grades  
(d=1)

S. 119-

Antwort zu a):  $x_k \leq x \leq x_{k+1}$

$$G_n(x) = \frac{\cos(\xi)}{2!} (x-x_k)(x-x_{k+1})$$

$$|e_n(x)| \leq \frac{1}{2} \cdot |\cos \xi| \cdot |x-x_k| \cdot |x-x_{k+1}|$$

$$\leq \frac{1}{2} \cdot 1 \cdot \frac{h}{2} \cdot \frac{h}{2} = \frac{h^2}{8}$$

$$10^{-3} = \frac{h^2}{8} \Rightarrow h = \sqrt{8 \cdot 10^{-3}} = 0,0894$$

genaue  
Lösung!

Antwort zu b): relative Fehler:

$$\left| \frac{e_n(x)}{f(x)} \right| = \left| \frac{e_n(x)}{\cos x} \right| = \frac{1}{2} \cdot \frac{|\cos \xi|}{|\cos x|} \cdot |x-x_k| \cdot |x-x_{k+1}|$$

Bei kleinem  $h$  wegen  $x_k \leq x, \xi \leq x_{k+1}$ :  $\cos x \approx \cos \xi$

$$\frac{1}{2} \cdot |x-x_k| \cdot |x-x_{k+1}| = \frac{h^2}{8}$$

$\rightarrow$  ändert sich nichts!

Antwort zu c):

1) in  $h$  quadrate

2. in die Halbwertspunkte od. Intervalls mittels

3. wenn an 1 Stelle pro  $10^{-3}$

$$\Rightarrow \frac{h}{2}$$

$\rightarrow$  so lange, bis die Kreis  $10^{-3}$

Halbe zur de  
Stabilität!

- wie (Nicht)stetig oder als "Pfad"?
- irgendwas mit "L-Feld" (ist vollendet?)
- (nicht) mit "Ausgewählter" oder "Abstände"?



Institut für Angewandte und Numerische Mathematik  
TU Wien

WS 1997/98

Numerische Mathematik für Informatiker

Prüfung am 14. Oktober 1997

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)	c)	c)	

Der gesamte Rechengang ist auf den beteiligten  
Blättern zu dokumentieren.  
Zusätzlich beigefügte Zetteln werden bei der  
Korrektur nicht berücksichtigt.

1) (10 Punkte)  
Der Ausdruck

$$\varphi(x) = 1 - \frac{1}{1+x} = \frac{x}{1+x}, \quad x \in I := (-\infty, -1.5] \cup [-0.5, +\infty)$$

sei für eine gegebene Maschinenzahl  $x$  in Gleitpunktarithmetik, mit einer relativen Maschinengenauigkeit  $\epsilon_{ps}$ , zu berechnen.

- a) Man gebe die relative Konditionszahl  $K_{\varphi(x)-x}$  als Funktion von  $x$  an.  
Man gebe für  $K_{\varphi(x)-x}$  eine möglichst gute, von  $x$  unabhängige Schranke an, d.h. man bestimme ein möglichst kleines  $C > 0$ , sodass

$$|K_{\varphi(x)-x}| \leq C \quad \forall x \in I$$

gilt. Welche Aussage kann man demgemäß für die Kondition des Problems machen?

- b) Man gebe für beide Auswertungsvarianten möglichst gute Fehlerschranken für den relativen Rundungsfehler als Ausdrücke in  $x$  und in der Maschinengenauigkeit  $\epsilon_{ps}$  an, d.h. man bestimme für beide Fälle ein  $\epsilon$ , mit<sup>1</sup>

$$\varphi(x) \approx \varphi(x)(1 + \epsilon)$$

und schätze es in Abhängigkeit von  $x$  und  $\epsilon_{ps}$  ab.

- c) Für welchen Bereich der  $x$ -Werte ist eine der beiden Auswertungsvarianten der anderen stark vorzuziehen? Ausführliche Begründung!

<sup>1</sup> $\varphi(x)$  symbolisiert die jeweilige Auswertung von  $\varphi$  in Gleitpunktarithmetik

Tragen Sie die Ergebnisse ein und kreuzen Sie die zutreffenden Aussagen an!

$$K_{\varphi(x) \rightarrow x} =$$

Für  $x \in I$  ist  $C =$

und deshalb ist die Auswertung  gut  schlecht konditioniert.

Fehlerschranke für den relativen Rundungsfehler von  $\varphi_1(x) := 1 - \frac{1}{1+x}$ :

$$|\varepsilon| \leq$$

Fehlerschranke für den relativen Rundungsfehler von  $\varphi_2(x) := \frac{x}{1+x}$ :

$$|\varepsilon| \leq$$

Die Formel   $\varphi_1(x)$    $\varphi_2(x)$  ist bei der Auswertung von  $\varphi(x)$

für  $x \approx$   
Begründung:

vorzuziehen.

2) (10 Punkte)

Gegeben sei die Matrix

$$M = \begin{pmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & 1 \end{pmatrix}$$

a) Berechnen Sie die  $LU$ -Zerlegung von  $M$  (ohne Pivotsuche):

$$L =$$

$$U =$$

b) Verwenden Sie diese Zerlegung, um das Gleichungssystem  $M^T x = y$  zu lösen, wobei  $y = (-12, -9, -3, 14)^T$ . ( $M^T \dots$  die transponierte Matrix.)

$$x =$$

Dokumentieren und erläutern Sie Ihre Rechnung.

3) (16 Punkte)

Die Datenpunkte

$i$	1	2	3	4	5	6	7	8
$x_i$	-2	-1	0	1	2	3	4	5
$y_i$	0	1	2	3	5	6	7	8

sind mit Hilfe der Akima-Methode zu interpolieren.

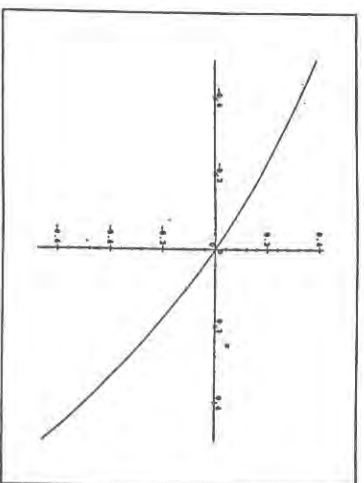
- Man fertige eine Skizze der Datenpunkte an.
- Man überlege (auf Grund der Definition der Akima-Funktion), wie die Interpolationsfunktion für  $x \in [0, 3]$  aussieht und trage deren Verlauf in die Skizze ein.
- Man berechne den Wert der Interpolationsfunktion an der Stelle  $x = 1.5$ .

Numerische Mathematik für Informatiker  
Prüfung am 26. Juni 1997

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
	c)	c)	
	d)		

Der gesamte Rechengang ist auf den beiliegenden Blättern zu dokumentieren.  
Zusätzlich beigefügte Zetteln werden bei der Korrektur nicht berücksichtigt.



1) (16 Punkte)  
Die Funktion

$$f(x) := \ln(1 - \sin x)$$

(vgl. obige Skizze) soll in der Nähe der Nullstelle  $x = 0$  genau ausgewertet werden. Für  $x = 0.1234567890E-8$  erhalte Auswertung in 10-stelliger dezimaler Gleitpunktarithmetik den Wert

$$x = -0.119999989E-8$$

mit nur einer einzigen richtigen Stelle im Vergleich zu dem (auf 10 Stellen) exakten Wert

$$x = -0.1234567891..E-8.$$

a) Liefere eine ausführliche Erklärung für die Ursache dieser ausgeprägten Ungenauigkeit.

b) Jemand verwendet die Identität  $\sin^2 x + \cos^2 x = 1$ , um  $f(x)$  mittels

$$\ln\left(\frac{\cos^2 x}{1 + \sin x}\right)$$

auszuwerten. Erhält dies ein genaueres Resultat und wenn ja, warum nicht? (Ausführliche Argumentation!)

<sup>1</sup>Korrekte 10-stellige Implementierung der auftretenden Elementarfunktionen wird vorausgesetzt.

2) (18 Punkte)

Man betrachte die beiden linearen Gleichungssysteme

$$\begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ 1 & 1 & 1 & \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{pmatrix} \quad (1)$$

und

$$\begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ 1 & 1 & 1 & \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} m \\ m-1 \\ m-2 \\ \vdots \\ 1 \end{pmatrix} \quad (2)$$

a) Man gebe für den Fall einer allgemein rechte Seite, also für das Gleichungssystem

$$\begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ 1 & 1 & 1 & \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} k_1 \\ k_2 \\ k_3 \\ \vdots \\ k_n \end{pmatrix} \quad (3)$$

die Lösung an:

$x_1 =$
$x_2 =$
$x_3 =$
$\vdots$
$x_n =$

b) Man gebe die Inverse  $A^{-1}$  der Matrix

$$A = \begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ 1 & 1 & 1 & \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix}$$

an:

$A^{-1} =$
$\left[ \right]$

c) Nach Folgend werde als Vektornorm die Maximumnorm  $\|\cdot\|_\infty$  zugrunde gelegt. Man bestimme (aufgrund von  $B$ ) die relative Konditionszahl  $K(A)$ :

$$K(A) =$$

d) Man betrachte man (1), (2) und (3) mit positiver reeller Seite  $B$ , d.h. in Fall (1) die reelle Seite im Fall (2) die reelle Seite  $\begin{pmatrix} m \\ \vdots \\ 1 \end{pmatrix} + \begin{pmatrix} \Delta B_1 \\ \vdots \\ \Delta B_n \end{pmatrix}$  und im allgemeinen Fall (3) die reelle Seite  $\begin{pmatrix} B_1 \\ \vdots \\ B_n \end{pmatrix} + \begin{pmatrix} \Delta B_1 \\ \vdots \\ \Delta B_n \end{pmatrix}$ . Unter der Annahme  $\left\| \begin{pmatrix} \Delta B_1 \\ \vdots \\ \Delta B_n \end{pmatrix} \right\|_\infty = \|\Delta B\|_\infty \leq \Delta$  ordnet man entsprechend wie der

unter Punkt a) getriebenen Lösungswahlungs Modifikation über die Größe  $F$  ist

$$\frac{\| \Delta \times \|_\infty}{\| \times \|_\infty} \leq F \frac{\| \Delta B \|_\infty}{\| B \|_\infty}$$

Res. (1) und (2) ab:

$$\begin{array}{l} F = \text{für (1)} \\ F = \text{für (2)} \end{array}$$

sind verglichen jeweils mit  $K(A)$ :

$$\begin{array}{l} \frac{F}{K(A)} = \text{für (2)} \\ \frac{F}{K(A)} = \text{für (2)} \end{array}$$

(10 Punkte)

In einem Programm zur Nullstellenbestimmung skalarer (eindimensionaler) nicht-linearer Gleichungen soll das Konvergenzverhalten der vom Programm erzeugten Folge automatisch untersucht werden.

- a) Es soll zwischen konvergentem und nicht-konvergentem Verhalten der Folge von Näherungswerten unterschieden werden. Man mache einen Entwurf für einen Programmabschnitt, der diese Unterscheidung trifft.
- b) Im Fall konvergenter Näherungswerte soll zwischen Konvergenzphase und Endphase (wo der Abstand zur Lösung in der Größe der Rechenfehler pro Schritt liegt) unterschieden werden. Man mache einen Entwurf für einen Programmteil, der dies leistet.
- c) Für die Konvergenzphase soll die Konvergenzordnung festgestellt werden. Man mache einen Entwurf für einen Teil des Programms, mit dem diese Ermittlung möglich ist.

*Hinmerkung:* Es ist zu beachten, daß die exakte Lösung  $x^*$  nicht bekannt ist.

**PRÜFUNGSORDNER - ein Service Deiner Fachschaft Informatik!**

LVA: NUMERISCHE MATHEMATIK -VO

Preis: ~~20,-~~

Institut für Angewandte und Numerische Mathematik  
TU Wien

WS 1996/97

**Numerische Mathematik für Informatiker**

Prüfung am 28. Januar 1997

Name	Vorname	Kennzahl / Matrikelnummer
------	---------	---------------------------

Beispiel 1	Beispiel 2	Beispiel 3	Gesamt
a)	a)	a)	
b)	b)	b)	
c)			



1) (16 Punkte)

Zur Auswertung des Ausdrucks

$$a(x) = \frac{1 - \cos x}{x}$$

für  $x$ -Werte in der Nähe von Null stehen zwei mathematisch äquivalente Formeln zur Verfügung:

$$f(x) = \frac{1 - \cos x}{x} \quad \text{bzw.} \quad g(x) = \frac{\tan(x/2) \sin x}{x}$$

- a) Man gebe die relative Konditionszahl  $K_{a(x) \leftarrow x}$  als Funktion von  $x$  an. Man überprüfe die Kondition der Aufgabe für kleine  $x$ -Werte,  $x \approx 0$ .

*Hinweis:* Es gilt  $\lim_{x \rightarrow 0} \frac{x}{\sin x} = 1$ .

- b) Man gebe für  $f(x)$  und  $g(x)$  Rundungsfehlerschranken für den relativen Rundungsfehler als Ausdrücke in  $x$  und in der Maschinengenauigkeit  $eps$  an. Man nehme dabei an, daß die relativen Fehler bei der Berechnung von trigonometrischen Funktionen durch  $eps$  beschränkt sind. *Hinweis:* Ist  $y\delta$  klein, so gilt in der ersten Näherung:

$$\tan(y(1 + \delta)) = \tan y + \frac{1}{\cos^2 y} y \delta = \tan y \left( 1 + \frac{2y}{\sin(2y)} \delta \right)$$

- c) Welche der obigen Formeln ist bezüglich der Rundungsfehler bei der Auswertung von  $a(x)$  in der Nähe von Null vorzuziehen? Ausführliche Begründung!

Tragen Sie die Ergebnisse ein  
und kreuzen Sie die zutreffenden Aussagen an !

$$K_{a(x)} =$$

Für  $x \approx 0$  ist  $K_{a(x) \leftarrow x} \approx$

und deshalb ist die Auswertung  gut  schlecht konditioniert.

Fehlerschranke für den relativen Rundungsfehler von  $f(x)$  :

Fehlerschranke für den relativen Rundungsfehler von  $g(x)$  :

Die Formel  $\boxed{f(x)}$   $\boxed{g(x)}$  ist bei der Auswertung von  $a(x)$   
für  $x \approx 0$  vorzuziehen.  
Ausführliche Begründung!

2) (20 Punkte)

a) Gegeben sei das folgende, sehr speziell strukturierte lineare Gleichungssystem:

$$\begin{pmatrix} a & a & a & a & \dots & a \\ 1 & a & a & a & \dots & a \\ 0 & 1 & a & a & \dots & a \\ 0 & 0 & 1 & a & \dots & a \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & 1 & a \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \quad (1)$$

Man betrachte für diesen einfachen Fall die Gauß-Elimination (ohne Pivotstrategie), die zur Faktorisierung  $LUx = b$  führt und gebe die Matrizen  $L$  und  $U$  an:

$$L = \begin{pmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix}$$

$$U = \begin{pmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix}$$

Weiters gebe man  $x_n$  als Formelausdruck in  $a$  an:

$$x_n =$$