

Klausur zur Vordiplom-Prüfung

Numerische Verfahren

17. März 2008

Sie haben **90** Minuten Zeit zum Bearbeiten der Klausur.

Bitte kennzeichnen Sie jedes Blatt mit Ihrem Namen und Ihrer Matrikelnummer in DRUCKSCHRIFT.

Tragen Sie bitte zunächst Ihren Namen, Ihren Vornamen und Ihre Matrikelnummer in **Druckschrift** in die folgenden jeweils dafür vorgesehenen Felder ein.

Diese Eintragungen werden auf Datenträger gespeichert.

| | | | | | | | | | | | | | | |
|------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| Name: | <input type="text"/> |
| Vorname: | <input type="text"/> |
| Matr.-Nr. | <input type="text"/> |

Ich bin darüber belehrt worden, dass meine Ausarbeitung nur dann als Prüfungsleistung bewertet wird, wenn die Nachprüfung durch das Zentrale Prüfungsamt der TUHH meine offizielle Zulassung vor Beginn der Prüfung ergibt.

(Unterschrift)

Lösen Sie die folgenden 12 Aufgaben!

| | | | | | | | | | | | | |
|----------------|---|---|---|---|---|---|---|---|---|----|----|----|
| Aufgabe | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| Punkte | | | | | | | | | | | | |

| | |
|----------|--|
| Σ | |
|----------|--|

Aufgabe 1: (4 Punkte)

Zur Berechnung der Nullstellen eines reellen quadratischen Polynoms $z(x) = x^2 + px + q$ wird oft die sogenannte „pq-Formel“

$$x_{\pm} = -\frac{p}{2} \pm \sqrt{\frac{p^2}{4} - q}$$

für *beide* Nullstellen verwendet. Da

$$(x - x_+)(x - x_-) = x^2 - (x_+ + x_-)x + x_+x_- \quad \Rightarrow \quad q = x_+ \cdot x_-$$

gilt, könnte man aber auch x_+ aus der „pq-Formel“ und x_- aus

$$x_- = \frac{q}{x_+}$$

bestimmen. Welche dieser beiden mathematisch äquivalenten Varianten würden Sie im Fall $p = -10^8$ und $q = 1$ in Gleitkommaarithmetik verwenden? Warum?

Lösung zu Aufgabe 1: Die „pq-Formel“ basiert für x_+ auf der Addition zweier positiver Terme, also ist hier keine Auslöschung zu erwarten. Für die Berechnung von x_- werden allerdings zwei fast gleich grosse Terme voneinander abgezogen, und die Berechnung mittels des zweiten Weges ist bei den gegebenen Daten deutlich besser konditioniert. (4 Punkte)

Unter Matlab/Octave liefert die naive Auswertung der „pq-Formel“ approximativ

$$\text{fl}_1(x_+) = 10^8, \quad \text{fl}_1(x_-) \approx 7.45058059692383 \cdot 10^{-9},$$

wohingegen die Auswertung nach der zweiten Idee als Näherung für x_+ logischerweise denselben Wert und für x_- den approximativen Wert

$$\text{fl}_2(x_+) = 10^8, \quad \text{fl}_2(x_-) = 10^{-8}$$

liefert. Die echten Werte unterscheiden sich von denen der zweiten Berechnung nur um

$$\frac{|x_+ - \text{fl}_2(x_+)|}{x_+} \leq 1.0000000000000003 \cdot 10^{-16} = 10^{-16} + 3 \cdot 10^{-32},$$

$$\frac{|x_- - \text{fl}_2(x_-)|}{x_-} \leq 1.0000000000000001 \cdot 10^{-16} = 10^{-16} + 1 \cdot 10^{-32},$$

wohingegen

$$\frac{|x_- - \text{fl}_1(x_-)|}{x_-} \leq 0.2549419403076171 < \frac{1}{4} + \frac{5}{1000}$$

gilt. Der zweite Weg bestimmt also beide Nullstellen bis auf einen relativen Fehler von 10^{-16} , also bis auf (IEEE 754 double precision) Maschinengenauigkeit.

Aufgabe 2: (5+5 Punkte)

Geben Sie die Lagrangsche Form der Polynominterpolation wieder. Berechnen Sie das Interpolationspolynom der Funktion

$$f(x) = \sin(\pi x)$$

zu den Punkten

$$\begin{array}{c|ccccc} i & 0 & 1 & 2 & 3 & 4 \\ \hline x_i & -2 & -1 & 0 & 1 & 2 \end{array}.$$

Lösung zu Aufgabe 2: Die Lagrangsche Interpolation zu den Datenpaaren $\{(x_j, y_j)\}_{j=0}^n$ mit $x_i \neq x_j$ für alle $i \neq j$ ist eine Vorschrift um das eindeutige Polynom p vom Höchstgrad n zu berechnen, dass die Interpolationsbedingungen

$$p(x_j) = y_j, \quad j = 0, 1, \dots, n$$

erfüllt. Es läßt sich leicht explizit angeben als

$$p(x) = \sum_{j=0}^n y_j \ell_j(x),$$

wobei die Lagrangeschen Basispolynome gegeben sind durch die Interpolationspolynome zu den speziellen Aufgaben

$$\ell_j(x_i) = \delta_{ij}, \quad i = 0, 1, \dots, n, \quad j = 0, 1, \dots, n,$$

und explizit angegeben lauten

$$\ell_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^n (x - x_i) \bigg/ \prod_{\substack{i=0 \\ i \neq j}}^n (x_j - x_i), \quad j = 0, \dots, n. \quad (5 \text{ Punkte})$$

Die Interpolationsaufgabe läßt sich auffassen als die Interpolation der Datenpaare

$$\begin{array}{c|ccccc} x_i & -2 & -1 & 0 & 1 & 2 \\ \hline y_i & 0 & 0 & 0 & 0 & 0 \end{array}.$$

Das Interpolationspolynom p_4 hat den *Höchstgrad* vier und ist eindeutig bestimmt. Damit ist das Interpolationspolynom das Nullpolynom. (5 Punkte)

Natürlich konnte man p_4 auch mittels der Lagrangeschen oder der Newtonschen Form der Polynominterpolation berechnen und dabei zu dem selben Ergebnis kommen.

Aufgabe 3: (5+5 Punkte)

Beschreiben Sie die Grundidee der Newtonschen Form der Polynominterpolation und geben Sie diese in Formeln wieder.

Lösung zu Aufgabe 3: Die Grundidee der Newtonschen Form der Polynominterpolation ist es, bei einem bereits berechneten Interpolationspolynom p_n vom Höchstgrad n zu den $n + 1$ Datenpaaren $\{(x_i, y_i)\}_{i=0}^n$ mit paarweise verschiedenen Knoten $\{x_i\}_{i=1}^n$ nach Hinzufügen eines weiteren Datums (x_{n+1}, y_{n+1}) , wobei selbstverständlich $x_{n+1} \neq x_i$ für alle $i = 0, \dots, n$, schnell das zugehörige Interpolationspolynom zu den jetzt $n + 2$ Datenpaaren berechnen zu können. Dieses geschieht mittels der Addition des bereits berechneten Polynomes p_n und eines polynomialen Termes vom Grad n oder 0 , der allen vorherigen Stellen gleich Null ist und dessen Vorfaktor a so gewählt ist, dass

$$p_{n+1}(x_{n+1}) = p_n(x_{n+1}) + a \prod_{i=1}^n (x_{n+1} - x_i) = y_{n+1},$$

also

$$a = \frac{y_{n+1} - p_n(x_{n+1})}{\prod_{i=1}^n (x_{n+1} - x_i)}$$

gilt. (5 Punkte)

Die Newtonsche Interpolationsformel ist gegeben durch

$$p_n(x) = \sum_{j=0}^n [x_0, \dots, x_j] \prod_{k=0}^{j-1} (x - x_k),$$

wobei die sogenannten dividierten Differenzen $[x_i, \dots, x_j]$ rekursiv definiert sind durch

$$\begin{aligned} [x_j] &:= y_j, \\ [x_i, \dots, x_j] &:= \frac{[x_{i+1}, \dots, x_j] - [x_i, \dots, x_{k-1}]}{x_j - x_i}, \quad j > i \geq 0. \end{aligned} \quad (5 \text{ Punkte})$$

Aufgabe 4: (4+4 Punkte)

Was ist eine Quadraturformel? Berechnen Sie approximativ das Integral

$$\int_0^\pi \cos(x) dx = \sin(\pi) - \sin(0) = 0$$

mit einer Quadraturformel.

Lösung zu Aufgabe 4: Eine Quadraturformel ist eine Formel zur numerischen Annäherung eines Integrales, welche nur auf Funktionsauswertungen der Funktion unter dem Integral basiert und allgemein die Gestalt

$$\int_0^1 f(x) dx \approx \sum_{j=0}^n \omega_j f(x_j)$$

hat. Die Angabe der Auswertungsstellen, der Stützstellen $\{x_j\}_{j=0}^n$ und der sogenannten Gewichte $\{\omega_j\}_{j=0}^n$ bestimmt in eindeutiger Weise eine Quadraturformel. (4 Punkte)

Da in der Aufgabenstellung keine Rede von der Art der zu verwendenden Quadraturformel zur Approximation war, ist viel Freiheit gegeben. Wir verwenden exemplarisch erst die Mittelpunkregel, danach die Trapezregel und zu guter Letzt die optimale Gauß-Formel mit zwei Knoten.

Die Mittelpunkt- oder Rechteckregel ersetzt das Integral durch die Auswertung des Integranden in der Mitte des Intervalles und berechnet den Flächeninhalt des (vorzeichenbehafteten) Rechtecks, welches durch diese (vorzeichenbehaftete) Höhe mal der Länge, also, der Intervallbreite gegeben ist,

$$\begin{aligned} \int_0^\pi \cos(x) dx &\approx \cos\left(\frac{\pi}{2}\right) \cdot (\pi - 0) \\ &= 0 \cdot \pi = 0. \end{aligned} \quad (4 \text{ Punkte})$$

Die Trapezregel verwendet das von den Funktionswerten am Anfang und am Ende des Intervalles aufgespannte Trapez. Der (vorzeichenbehaftete) Flächeninhalt dieses Trapezes berechnet sich aus dem Mittelwert der beiden Funktionswerte mal der Länge des Intervalles, also gemäß

$$\begin{aligned} \int_0^\pi \cos(x) dx &\approx \frac{1}{2} (\cos(\pi) + \cos(0)) \cdot (\pi - 0) \\ &= (-1 + 1) \cdot \pi = 0. \end{aligned} \quad (4 \text{ Punkte})$$

Die optimale Gauß-Formel mit zwei Knoten basiert in Nachschlagewerken meist auf dem Intervall $[-1, 1]$, ist im Skript aber bereits für das Referenzintervall $[0, 1]$ gegeben und hat dort die Knoten

$$x_1 = \frac{1}{2} \left(1 - \frac{1}{\sqrt{3}}\right), \quad x_2 = \frac{1}{2} \left(1 + \frac{1}{\sqrt{3}}\right)$$

und die Gewichte

$$w_1 = w_2 = \frac{1}{2}.$$

Die lineare Transformation $y = \pi \cdot x$ bildet das Intervall $[0, 1]$ auf das Intervall $[0, \pi]$ ab und liefert die (neuen, aber auch mit x_1 und x_2 bezeichneten) Knoten

$$x_1 = \frac{\pi}{2} \left(1 - \frac{1}{\sqrt{3}}\right), \quad x_2 = \frac{\pi}{2} \left(1 + \frac{1}{\sqrt{3}}\right).$$

Damit ergibt die Gauß-Quadratur mit zwei Knoten unter Berücksichtigung des Additionstheoremes des Kosinus und der Feststellung, dass der Kosinus eine ge-

rade Funktion ist, den Wert

$$\begin{aligned}\int_0^\pi \cos(x) dx &\approx \frac{1}{2} \left(\cos \left(\frac{\pi}{2} \left(1 - \frac{1}{\sqrt{3}} \right) \right) + \cos \left(\frac{\pi}{2} \left(1 + \frac{1}{\sqrt{3}} \right) \right) \right) \\ &= \cos \left(\frac{\pi}{2} \right) \cos \left(\frac{\pi}{2\sqrt{3}} \right) = 0.\end{aligned}\quad (4 \text{ Punkte})$$

Man sieht also, dass bereits einfache Quadraturformeln gute Ergebnisse liefern können, aber das sollte nicht darüber hinwegtäuschen, dass der Fehler bei den genaueren Formeln (also solchen mit höherer Fehlerordnung und kleinerer Fehlerkonstante) *garantiert* für beliebige genügend oft stetig differenzierbare Funktionen kleiner beschränkt ist.

Aufgabe 5: (5+5 Punkte)

Welche Varianten der LR-Zerlegung kennen Sie? Welche würden Sie im Hinblick auf die numerische Stabilität und den nötigen Arbeitsaufwand verwenden?

Lösung zu Aufgabe 5: Es gibt drei Varianten der LR-Zerlegung, nämlich die LR-Zerlegung ohne Pivotisierung,

$$\mathbf{A} = \mathbf{LR},$$

wobei \mathbf{L} eine linke untere Dreiecksmatrix mit Einsen auf der Diagonalen und \mathbf{R} eine rechte obere Dreiecksmatrix ist, die LR-Zerlegung mit partieller Pivotisierung (auch Spaltenpivotsuche genannt),

$$\mathbf{PA} = \mathbf{LR},$$

wobei \mathbf{L} und \mathbf{R} wie vorher gewählt sind und \mathbf{P} eine Permutationsmatrix ist, welche das größte Element in der jeweiligen aktuell behandelten Spalte in die Diagonale tauscht und die LR-Zerlegung mit vollständiger Pivotisierung,

$$\mathbf{PAQ} = \mathbf{LR},$$

wobei \mathbf{L} , \mathbf{R} und \mathbf{P} wie oben sind und \mathbf{Q} zusammen mit \mathbf{P} dazu verwendet wird um das in der aktuell behandelten Restmatrix größte Element in die aktuell behandelte Position in der Diagonalen zu tauschen. (5 Punkte)

Am billigsten ist die einfache LR-Zerlegung zu implementieren, aber diese ist numerisch nicht stabil, die LR-Zerlegung existiert oft auch gar nicht, bekanntestes Beispiel ist

$$\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

Die LR-Zerlegung mit partieller Pivotisierung reicht theoretisch ohne Rundungsfehler aus, eine reguläre Matrix zu zerlegen, ist aber immer noch rundungsfehleranfällig, siehe das Beispiel im Skript. Die LR-Zerlegung mit vollständiger Pivotisierung ist am aufwändigsten, ist aber numerisch sehr stabil. Meist würde man,

sofern keine begründeten Ängste vorliegen, die LR-Zerlegung mit partieller Pivotisierung verwenden und bei unglaublichen Resultaten auf die vollständige Pivotisierung zurückgreifen. (5 Punkte)

Wenn die Ergebnisse dann immer noch nicht die Praxis reflektieren, würde man zu Regularisierung greifen, da dann (vorausgesetzt, es gibt keine groben Modellierungsfehler) meist die Kondition der Matrix zu hoch ist, um sinnvolle numerische Ergebnisse zu erhalten.

Aufgabe 6: (5 Punkte)

Lösen Sie das Gleichungssystem

$$\begin{pmatrix} 110 & 1 \\ 100 & 11 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 110 \\ 110 \end{pmatrix}$$

in zweistelliger dezimaler Fließkommazahlenarithmetik, also unter

$$\text{fl}(x \circ y) := \text{Rundung des Ergebnisses von } x \circ y \text{ auf zwei Stellen}$$

für jede Operation $\circ \in \{+, -, \cdot, /\}$, es gilt also z. B.

$$\text{fl}(100 + 10) = 110, \quad \text{fl}(120 + 3) = 120, \quad \text{fl}(120 + 5) = 130,$$

mittels der LR-Zerlegung.

Lösung zu Aufgabe 6: Die LR-Zerlegung kann hier ohne Pivotisierung erfolgen, das größte Element ist bereits in der Position (1, 1). Da die rechte Seite bereits gegeben ist, transformieren wir die Elemente gleich mit, arbeiten also mit der Matrix $(\mathbf{A} \ \mathbf{b})$,

$$\left(\begin{array}{cc|c} 110 & 1 & 110 \\ 100 & 11 & 110 \end{array} \right).$$

Zur Berechnung des ersten und einzigen interessanten Wertes in L wird 100 durch 110 in Fließkommazahlenarithmetik geteilt,

$$\text{fl}(100/110) = \text{fl}(0.\overline{90}) = 0.91.$$

Nun wird, wiederum in Fließkommazahlenarithmetik, die erste Zeile mit 0.91 malgenommen,

$$\begin{aligned} \text{fl}(0.91 \cdot (110 \ 1 \mid 110)) &= (\text{fl}(0.91 \cdot 110) \ \text{fl}(0.91 \cdot 1) \mid \text{fl}(0.91 \cdot 110)) \\ &= (\text{fl}(100.1) \ \text{fl}(0.91) \mid \text{fl}(100.1)) \\ &= (100 \ 0.91 \mid 100). \end{aligned}$$

Jetzt wird, selbstverständlich auch in Fließkommazahlenarithmetik, die so erhaltene Zeile von der zweiten Zeile abgezogen,

$$\begin{aligned} \text{fl}((100 \ 11 \mid 110) - (100 \ 0.91 \mid 100)) \\ &= (\text{fl}(100 - 100) \ \text{fl}(11 - 0.91) \mid \text{fl}(110 - 100)) \\ &= (\text{fl}(0) \ \text{fl}(10.01) \mid \text{fl}(10)) = (0 \ 10 \mid 10). \end{aligned}$$

Die resultierende in Fließkommazahlenarithmetik berechnete LR-Zerlegung hat demnach in Kompaktschreibweise die Form

$$\left(\begin{array}{cc|cc} 110 & 1 & 110 & \\ \hline 0.91 & 10 & 10 & \end{array} \right),$$

womit wir bereits Approximationen für \mathbf{L} , \mathbf{R} aus $\mathbf{A} = \mathbf{LR}$ und $\mathbf{y} = \mathbf{L}^{-1}\mathbf{b}$ haben, nämlich

$$\tilde{\mathbf{L}} = \begin{pmatrix} 1 & 0 \\ 0.91 & 1 \end{pmatrix}, \quad \tilde{\mathbf{R}} = \begin{pmatrix} 110 & 1 \\ 0 & 10 \end{pmatrix} \quad \text{und} \quad \tilde{\mathbf{y}} = \begin{pmatrix} 110 \\ 10 \end{pmatrix}.$$

Um das Gleichungssystem in Fließkommazahlenarithmetik zu lösen, lösen wir das Gleichungssystem

$$\tilde{\mathbf{R}}\mathbf{x} = \begin{pmatrix} 110 & 1 \\ 0 & 10 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 110 \\ 10 \end{pmatrix} = \tilde{\mathbf{y}}$$

mittels Rückwärtsauflösen in Fließkommazahlenarithmetik. Die letzte Komponente der genäherten Lösung ist durch

$$\tilde{x}_2 = \text{fl}(10/10) = \text{fl}(1) = 1$$

gegeben, die erste Komponente erhält man aus

$$\tilde{x}_1 = \text{fl}(\text{fl}(110 - 1)/110) = \text{fl}(\text{fl}(109)/110) = \text{fl}(110/110) = \text{fl}(1) = 1.$$

Die in Fließkommazahlenarithmetik berechnete (approximative) Lösung ist demnach

$$\tilde{\mathbf{x}} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}. \quad (5 \text{ Punkte})$$

Aufgabe 7: (5+7 Punkte)

Bestimmen Sie die Pseudonormallösungen der beiden Ausgleichsprobleme

$$\left\| \begin{pmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 5 \\ 1 & 5 \\ 1 & 3 \\ 1 & 1 \end{pmatrix} \mathbf{x}^1 - \begin{pmatrix} 1 \\ 2 \\ 3 \\ 3 \\ 2 \\ 1 \end{pmatrix} \right\|_2 = \min, \quad \left\| \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \mathbf{x}^2 - \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} \right\|_2 = \min.$$

Lösung zu Aufgabe 7: Das erste Ausgleichsproblem wird durch den Vektor

$$\mathbf{x}^1 = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

gelöst, wie man leicht sehen kann, da die erste rechte Seite im Spann der Spaltenvektoren der ersten Matrix liegt und die Matrix den vollen Rang Zwei hat. Wenn man diese Lösung nicht sieht, so kann man natürlich mittels Normalgleichungen, QR-Zerlegung oder SVD die Pseudonormallösung bestimmen. (5 Punkte)

Das zweite Ausgleichsproblem hat eine Rang-Eins-Matrix, da

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} (1 \quad 1 \quad 1) = \underbrace{\frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}}_{=U_1} \underbrace{\frac{\sqrt{12}}{\sqrt{3}} \frac{1}{\sqrt{3}} (1 \quad 1 \quad 1)}_{=V_1^T},$$

wobei die letzte Schreibung $U_1 \Sigma_1 V_1^T$ bereits die (ökonomische Variante der) SVD dieser Rang-Eins-Matrix bezeichnet. Die Pseudonormallösung \hat{x}^2 ist charakterisiert durch die Multiplikation der rechten Seite mit der Pseudoinversen,

$$\begin{aligned} \hat{x}^2 &= \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}^\dagger \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} = V_1 \Sigma_1^{-1} U_1^T \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} \\ &= \frac{1}{\sqrt{3}} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \frac{1}{\sqrt{12}} \cdot \frac{1}{2} (1 \quad 1 \quad 1 \quad 1) \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} \\ &= \frac{1}{12} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} (1 + 2 + 3 + 4) = \frac{5}{6} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}. \quad (7 \text{ Punkte}) \end{aligned}$$

Aufgabe 8: (8+4+4 Punkte)

Beschreiben Sie die Potenzmethode. Ist es möglich, dass die Potenzmethode den Eigenwert 2 der Matrix

$$G = \begin{pmatrix} 1 & 1 & -1 \\ 0 & 2 & -2 \\ 0 & 0 & -2 \end{pmatrix}$$

approximiert? Wenn nein, warum nicht? Wenn ja, ist dieses für einen zufällig gewählten Startvektor überhaupt wahrscheinlich?

Lösung zu Aufgabe 8: Die Potenzmethode geht aus von einer quadratischen Matrix A und einem gegebenen (bereits normalisierten) Startvektor u^0 . Für eine (gewünschte) Anzahl von Indizes $k = 1, 2, \dots$ wird dann rekursiv (z. B.) berechnet:

$$v = Au^{k-1}, \quad u^k = v/\ell^H v,$$

wobei ℓ ein fest vorgegebener Vektor ist. Unter geeigneten Voraussetzungen konvergiert dann die Folge $\{\mathbf{u}^k\}_{k=0}^{\infty}$ gegen einen (meist den zum betragsgrößten Eigenwert gehörigen) Eigenvektor und die Skalierungsgröße $\ell^H \mathbf{v}$ gegen den zugehörigen Eigenwert. (8 Punkte)

Es ist möglich, dass die Potenzmethode nicht den betragsgrößten Eigenwert und zugehörigen Eigenvektor zurückgibt. Für geeignete Startvektoren kann jedes Eigenpaar zurückgegeben werden, z. B. wenn der Startvektor gleich dem (gesuchten) Eigenvektor ist. (4 Punkte)

Konvergenz gegen den Eigenwert 2 ist der Fall, wenn der Startvektor keinen Anteil am Eigenvektor zum Eigenwert -2 enthält, aber einen nichttrivialen Anteil am Eigenvektor zum Eigenwert 2 hat. Ein zufällig gewählter Startvektor hat mit großer Wahrscheinlichkeit Anteile an jedem Eigenvektor, damit ist die Konvergenz gegen den Eigenwert 2 unter diesen Umständen nahezu (aber nicht völlig) ausgeschlossen. (4 Punkte)

Die Anteile des Startvektors lassen sich in diesem Fall leicht beschreiben. Ein Eigenvektor zum Eigenwert 1 ist gegeben durch

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix},$$

ein Eigenvektor zum Eigenwert 2 durch

$$\mathbf{v}_2 = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix},$$

und ein Eigenvektor zum Eigenwert -2 lässt sich herleiten zu

$$\mathbf{v}_3 = \begin{pmatrix} 1 \\ 3 \\ 6 \end{pmatrix}.$$

Damit ist eine Matrix von Eigenvektoren und ihre Inverse gegeben als

$$\mathbf{V} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 6 \end{pmatrix}, \quad \mathbf{V}^{-1} = \frac{1}{18} \begin{pmatrix} 18 & -18 & 6 \\ 0 & 18 & -9 \\ 0 & 0 & 3 \end{pmatrix}.$$

Die Anteile $\{\alpha_i\}_{i=1}^3$ von \mathbf{u}^0 an den einzelnen Eigenvektoren erhält man aus

$$\mathbf{u}^0 = \sum_{i=1}^3 \mathbf{v}_i \alpha_i = \mathbf{V} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix},$$

also gemäß

$$\begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} = \mathbf{V}^{-1} \mathbf{u}^0 = \frac{1}{18} \begin{pmatrix} 18 & -18 & 6 \\ 0 & 18 & -9 \\ 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} u_{10} \\ u_{20} \\ u_{30} \end{pmatrix}.$$

Um den Eigenwert 2 zu approximieren, muss $\alpha_3 = u_{30}/6$ gleich Null sein und (unter der Voraussetzung, dass bereits u_{30} gleich Null ist) $\alpha_2 = u_{20}$ ungleich Null sein, mit anderen Worten, jeder Startvektor der Form

$$\mathbf{u}^0 = \begin{pmatrix} \star \\ c \\ 0 \end{pmatrix}, \quad c \neq 0$$

mit beliebigem Eintrag in der ersten Komponente führt zur Konvergenz der Potenzmethode gegen den Eigenwert 2. Ein zufällig gewählter Startvektor wird mit an Sicherheit grenzender Wahrscheinlichkeit einen Eintrag ungleich Null in der letzten Komponente haben, womit bei einer Skalierung mittels des zweiten Einheitsvektors

$$\boldsymbol{\ell} = \mathbf{e}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

die Iterierten der Potenzmethode sich im Limes wie die Glieder einer alternierenden Folge, nämlich wie

$$\begin{aligned} \mathbf{u}^{2k-1} &\rightarrow \frac{\alpha_2 \mathbf{v}_2 - \alpha_3 \mathbf{v}_3}{\alpha_2 - 3\alpha_3}, \\ \mathbf{u}^{2k} &\rightarrow \frac{\alpha_2 \mathbf{v}_2 + \alpha_3 \mathbf{v}_3}{\alpha_2 + 3\alpha_3}, \end{aligned}$$

verhalten, wenn der Startvektor die Gestalt

$$\mathbf{u}^0 = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \alpha_3 \mathbf{v}_3$$

hat. Insbesondere versagt die Skalierung (spätestens bei Konvergenz) bei $\alpha_2 = \pm 3\alpha_3$, was man z. B. an

$$3\mathbf{v}_2 - \mathbf{v}_2 = \begin{pmatrix} 3 \\ 3 \\ 0 \end{pmatrix} - \begin{pmatrix} 1 \\ 3 \\ 6 \end{pmatrix} = \begin{pmatrix} 2 \\ 0 \\ -6 \end{pmatrix}$$

sofort ablesen kann.

Aufgabe 9: (6 Punkte)

Beschreiben Sie die Berechnung der QR-Zerlegung einer Matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ unter Verwendung von Householdermatrizen (Spiegelungen)

$$\mathbf{H}(\mathbf{x}) = \mathbf{E} - 2 \frac{\mathbf{x}\mathbf{x}^T}{\mathbf{x}^T \mathbf{x}}, \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{x} \neq \mathbf{o}_n.$$

Der Vektor $\mathbf{o}_n \in \mathbb{R}^n$ bezeichne hierbei den Nullvektor, \mathbf{E} die $n \times n$ -Einheitsmatrix.

Lösung zu Aufgabe 9: Die Idee der QR-Zerlegung unter Verwendung von Householder-Spiegelungen ist es, eine Matrix \mathbf{A} spaltenweise in eine obere Dreiecksmatrix zu transformieren. Dazu sucht man einen (Normalen-)Vektor, an dessen zugehöriger Hyperebene der Vektor in der ersten Spalte von \mathbf{A} auf ein Vielfaches des ersten Einheitsvektors abgebildet wird. Eine Spiegelung ist eine Kongruenztransformation, ändert also keine Längen, und das Vielfache des ersten Einheitsvektors ist bis auf das Vorzeichen bereits eindeutig bestimmt zu

$$\pm \|\mathbf{a}_1\|_2 \mathbf{e}_1$$

wobei \mathbf{a}_1 die erste Spalte von \mathbf{A} bezeichne. Die Differenz der Vektoren vorher, also \mathbf{a}_1 , und nachher, also $\pm \|\mathbf{a}_1\|_2 \mathbf{e}_1$, steht senkrecht auf der Hyperebene, welche wir für die Spiegelung nehmen wollen, also bildet

$$\mathbf{H}_1 := \mathbf{H}(\mathbf{z}_\pm) \quad \text{wobei} \quad \mathbf{z}_\pm := \mathbf{a}_1 \mp \|\mathbf{a}_1\|_2 \mathbf{e}_1$$

die erste Spalte auf $\pm \|\mathbf{a}_1\|_2 \mathbf{e}_1$ ab,

$$\mathbf{H}(\mathbf{z}_\pm) \mathbf{a}_1 = \pm \|\mathbf{a}_1\|_2 \mathbf{e}_1,$$

das Vorzeichen sucht man jetzt gemeinhin so aus, dass in der ersten Komponente von \mathbf{z}_\pm keine Auslöschung stattfindet. Jetzt hat man, da Householder-Matrizen orthogonal und symmetrisch sind, also die Inverse gleich der Matrix ist, \mathbf{A} bereits zerlegt in

$$\mathbf{A} = \mathbf{H}_1 \mathbf{A}_1,$$

wobei \mathbf{A}_1 in der ersten Spalte ein Vielfaches des ersten Einheitsvektors stehen hat. Nun spiegelt man in der Matrix, die man erhält, wenn man von \mathbf{A}_1 die erste Spalte und Zeile entfernt, die neue erste Spalte in den (jetzt um Eins kürzeren) Einheitsvektor mit einer Matrix $\mathbf{H}_2 \in \mathbb{R}^{(n-1) \times (n-1)}$ analog zu dem eben beschriebenen Vorgehen. Damit gilt

$$\mathbf{A}_1 = \begin{pmatrix} 1 & \mathbf{o}_n^T \\ \mathbf{o}_n & \mathbf{H}_2 \end{pmatrix} \mathbf{A}_2,$$

zusammen also bereits

$$\mathbf{A} = \mathbf{H}_1 \begin{pmatrix} 1 & \mathbf{o}_n^T \\ \mathbf{o}_n & \mathbf{H}_2 \end{pmatrix} \mathbf{A}_2.$$

Nun wird dieses Schema iterativ durchgeführt, bis die Matrix, die nach Multiplikation mit einer Householder-Matrix und Wegnehmen der ersten Zeile und Spalte keine Elemente mehr hat. Am Ende gilt (hier am Beispiel des für ein Ausgleichsproblem typischen Falles $m \leq n$)

$$\mathbf{A} = \mathbf{H}_1 \begin{pmatrix} 1 & \mathbf{o}_n^T \\ \mathbf{o}_n & \mathbf{H}_2 \end{pmatrix} \cdots \begin{pmatrix} \mathbf{E}_{m-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{H}_m \end{pmatrix} \mathbf{R},$$

wobei \mathbf{R} nach Konstruktion eine obere Dreiecksmatrix ist. Das Produkt der erweiterten Householder-Matrizen ist ein Produkt von orthogonalen Matrizen, also wiederum orthogonal. Damit ist die QR-Zerlegung von \mathbf{A} gegeben als

$$\mathbf{A} = \mathbf{QR}, \quad \text{wobei} \quad \mathbf{Q} := \mathbf{H}_1 \begin{pmatrix} 1 & \mathbf{o}_n^T \\ \mathbf{o}_n & \mathbf{H}_2 \end{pmatrix} \cdots \begin{pmatrix} \mathbf{E}_{m-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{H}_m \end{pmatrix}. \quad (6 \text{ Punkte})$$

Selbstverständlich multipliziert man \mathbf{Q} nie aus, sondern verwendet die Darstellung als Produkt von Householder-Matrizen, wobei auf die Reihenfolge zu achten ist, da Matrixmultiplikation bekanntlich nicht kommutativ ist. Die Householder-Matrizen werden auch nur in Form einer Skalierung und eines Vektors gespeichert und angewendet, was sowohl Speicherplatz als auch Arbeitsaufwand einspart.

Aufgabe 10: (6+3 Punkte)

Der zentrale Differenzenquotient

$$D_2 f(x; h) := \frac{f(x+h) - f(x-h)}{2h}$$

zur Approximation der ersten Ableitung an der Stelle x hat die Fehlerordnung 2. Beschreiben Sie in eigenen Worten, welches Verhalten Sie auf einem Rechner mit einer Maschinengenauigkeit von ca. 10^{-16} in der Genauigkeitsfunktion

$$G(h) := |f'(x) - D_2 f(x; h)|$$

für zweimal stetig differenzierbare Funktionen f , z. B. $f(x) = \cos(x)$, an einer Auswertungsstelle, z. B. $x = 1$ erwarten. In welcher Größenordnung wählen Sie h , damit die Formel trotz Rundungsfehlern „optimale“ Ergebnisse liefert?

Lösung zu Aufgabe 10: Die Genauigkeit bestimmt sich aus zwei Einflüssen. Wenn h klein gewählt wird, so beschreibt die Fehlerordnung, in welcher Potenz von h der Differenzenquotient gegen $f'(x)$ strebt, durch die Auslöschung ist aber immer ein Anteil gegeben, der für kleine h sich invers linear (siehe Skript), also wie

$$\frac{C_1}{h} 10^{-16}$$

verhält. Insgesamt kann man den Verlauf der Genauigkeit durch eine Funktion der Art

$$G(h) \approx V(h) := \frac{C_1}{h} 10^{-16} + C_2 h^2 + O(10^{-16}), \quad C_1, C_2 > 0$$

beschreiben. (6 Punkte)

Die Wahl des optimalen h orientiert sich an der Approximation V , man berechnet das optimale h als Minimum der gegebenen Approximation, also unter Verwendung der Ableitung

$$V'(h) = -\frac{C_1}{h^2} 10^{-16} + 2C_2 h = 0 \quad \Rightarrow \quad -C_1 10^{-16} + 2C_2 h^3 = 0$$

zu

$$h_{\text{optimal}} = \sqrt[3]{\frac{C_1}{2C_2}} \cdot \sqrt[3]{10^{-16}} = O(10^{-16/3}),$$

ist also bis auf eine Konstante gleich der dritten Wurzel aus der Maschinengenauigkeit.

Aufgabe 11: (5+4+4 Punkte)

Berechnen Sie die Pseudoinversen der Matrizen

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 0 & 1 \\ 1 & 2 \\ 2 & 0 \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}.$$

Lösung zu Aufgabe 11: Die Pseudoinverse ist die eindeutige lineare Abbildung, die einer gegebenen rechten Seite die Pseudonormallösung eines Ausgleichsproblems zuordnet. Da im Falle einer Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ mit $m \geq n$ mit vollem (Spalten-)Rang n die Pseudonormallösung die eindeutige Lösung der Normalgleichungen ist, gilt im Falle eines vollen Ranges

$$\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T.$$

Die Matrix \mathbf{A} hat nicht vollen Rang, sondern ist eine Rang-Eins-Matrix, welche die folgende (speicherplatzsparende Variante der) SVD hat,

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{pmatrix} = \frac{1}{\sqrt{3}} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \cdot \sqrt{42} \cdot \frac{1}{\sqrt{14}} (1 \ 2 \ 3) = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T.$$

Damit ist die Pseudoinverse gegeben als

$$\mathbf{A}^\dagger = \mathbf{V}_1 \mathbf{\Sigma}_1^{-1} \mathbf{U}_1^T = \frac{1}{42} \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} (1 \ 1 \ 1) = \frac{1}{42} \begin{pmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 3 & 3 & 3 \end{pmatrix} = \frac{1}{42} \mathbf{A}^T. \quad (5 \text{ Punkte})$$

Die Matrizen \mathbf{B} und \mathbf{C} haben vollen Spaltenrang, also können wir die obige Formel zur Berechnung der Pseudoinversen heranziehen, es gilt

$$\begin{aligned} \mathbf{B}^\dagger &= (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T = \begin{pmatrix} 5 & 2 \\ 2 & 5 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 1 & 2 \\ 1 & 2 & 0 \end{pmatrix} \\ &= \frac{1}{21} \begin{pmatrix} 5 & -2 \\ -2 & 5 \end{pmatrix} \begin{pmatrix} 0 & 1 & 2 \\ 1 & 2 & 0 \end{pmatrix} \\ &= \frac{1}{21} \begin{pmatrix} -2 & 1 & 10 \\ 5 & 8 & -4 \end{pmatrix} \end{aligned} \quad (4 \text{ Punkte})$$

und

$$\mathbf{C}^\dagger = (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T = \frac{1}{4} \begin{pmatrix} 1 & 1 & 1 & 1 \end{pmatrix}. \quad (4 \text{ Punkte})$$

Aufgabe 12: (4+4 Punkte)

Berechnen Sie die Singulärwertzerlegungen (SVD) der folgenden Matrizen:

$$\mathbf{Q} = \begin{pmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}.$$

Lösung zu Aufgabe 12: Die erste Matrix ist eine Vielfache einer orthogonalen Matrix, alle Zeilen und Spalten sind paarweise orthogonal und die Norm ist jeweils 3. Damit ist z. B.

$$\mathbf{Q} = \begin{pmatrix} \frac{1}{3} \mathbf{Q} \end{pmatrix} \cdot (3\mathbf{E}) \cdot (\mathbf{E}) = \mathbf{U}_Q \mathbf{\Sigma}_Q \mathbf{V}_Q^T,$$

wobei \mathbf{E} die 3×3 -Einheitsmatrix bezeichnet, eine (nicht eindeutig bestimmte) SVD, z. B. wäre auch

$$\mathbf{Q} = (\mathbf{E}) \cdot (3\mathbf{E}) \cdot \begin{pmatrix} \frac{1}{3} \mathbf{Q} \end{pmatrix} = \tilde{\mathbf{U}}_Q \tilde{\mathbf{\Sigma}}_Q \tilde{\mathbf{V}}_Q^T$$

eine SVD. (4 Punkte)

Da \mathbf{y} die erste Spalte von \mathbf{Q} ist, gilt

$$\mathbf{y} = \begin{pmatrix} \frac{1}{3} \mathbf{Q} \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 0 \\ 0 \end{pmatrix} \cdot 1 = \mathbf{U}_y \mathbf{\Sigma}_y \mathbf{V}_y^T$$

und dieses ist eine (vollständige) SVD. (4 Punkte)

Man hätte auch ohne Kenntnis der SVD von \mathbf{Q} die ökonomische Variante

$$\mathbf{y} = \begin{pmatrix} \frac{1}{3} \mathbf{y} \end{pmatrix} \cdot 3 \cdot 1 = \mathbf{U}_y^{(0)} \mathbf{\Sigma}_y^{(0)} (\mathbf{V}_y^{(0)})^T$$

der SVD von \mathbf{y} angeben können.