

Numerische Mathematik 1

Wintersemester 2020, Übungsblatt 5

Ausarbeitung über Moodle bis 6. November 2020. Nach diesem Datum erscheinen die nachträglichen Kommentare und die [Lösungen](#) der Teilnehmer.

1. Für $n \in \mathbb{N}$, $n > 1$, sei die invertierbare Matrix A durch $A = D + L + U$ mit einer invertierbaren diagonalen Matrix D , einer streng unterdreieckigen Matrix L und einer streng oberdreieckigen Matrix U zerlegt. Für $\mathbf{x}, \mathbf{b} \in \mathbb{R}^n$ wird die Jacobi-Methode $\mathbf{x}^{(k)} = T_J \mathbf{x}^{(k-1)} + M_J^{-1} \mathbf{b}$, $k \in \mathbb{N}$, auf Seite 79 im Skriptum verwendet, um das lineare Gleichungssystem $A\mathbf{x} = \mathbf{b}$ iterativ zu lösen. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Wenn A streng diagonal dominant ist, folgt $\|T_J\|_\infty < 1$.
- (b) Wenn A streng diagonal dominant ist, folgt $\|T_J\|_1 < 1$.
- (c) Die Jacobi Iterierten $\{\mathbf{x}^{(k)}\}_{k=0}^\infty$ können divergieren, wenn die folgende Iteration konvergiert

$$\mathbf{y}^{(k)} = -(L + U)D^{-1}\mathbf{y}^{(k-1)} + \mathbf{b}, \quad k \in \mathbb{N}.$$

- (d) Wenn A^\top streng diagonal dominant ist, müssen die Jacobi Iterierten $\{\mathbf{x}^{(k)}\}_{k=0}^\infty$ zu $\mathbf{x} = A^{-1}\mathbf{b}$ konvergieren.

Kommentare: Die Kombination von (c) und (d) ist der Beweis, dass die Jacobi Iteration konvergiert, nicht nur wenn $\|T_J\|_\infty < 1$ gilt, sondern auch wenn $\|T_J\|_1 < 1$ gilt.

2. Für $n \in \mathbb{N}$, $n > 1$, sei die streng diagonal dominante Matrix $A = \{a_{i,j}\}_{i,j=1}^n \in \mathbb{R}^{n \times n}$ durch $A = D + L + U$ mit einer diagonalen Matrix D , einer streng unterdreieckigen Matrix L und einer streng oberdreieckigen Matrix U zerlegt. Mit $\mathbf{x}, \mathbf{b} \in \mathbb{R}^n$ wird die Gauß-Seidel-Methode $\mathbf{x}^{(k)} = T_{GS} \mathbf{x}^{(k-1)} + M_{GS}^{-1} \mathbf{b}$, $k \in \mathbb{N}$, auf Seite 82 im Skriptum verwendet, um das lineare Gleichungssystem $A\mathbf{x} = \mathbf{b}$ iterativ zu lösen. Sei (λ, \mathbf{x}) ein Eigenwert-Eigenvektor-Paar für die Iterationsmatrix T_{GS} mit $|x_k| = \|\mathbf{x}\|_\infty = 1$. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Es gilt $-U\mathbf{x} = \lambda(A - U)\mathbf{x}$.
- (b) Es gilt

$$-\sum_{j=i}^n a_{i,j}x_j = \lambda \sum_{j=1}^i a_{i,j}x_j, \quad 1 \leq i \leq n.$$

- (c) Es gilt

$$\lambda a_{k,k}x_k = -\sum_{j=k}^n a_{k,j}x_j - \lambda \sum_{j=1}^{k-1} a_{k,j}x_j$$

- (d) Es gilt

$$|\lambda| \leq \sum_{j=k+1}^n |a_{k,j}| \Big/ \left[|a_{k,k}| - \sum_{j=1}^{k-1} |a_{k,j}| \right] < 1.$$

Kommentare: Seien Teile (c) und (d) mit

$$-\sum_{j=i+1}^n a_{i,j}x_j = \lambda \sum_{j=1}^i a_{i,j}x_j, \quad 1 \leq i \leq n.$$

bzw.

$$\lambda a_{k,k}x_k = -\sum_{j=k+1}^n a_{k,j}x_j - \lambda \sum_{j=1}^{k-1} a_{k,j}x_j$$

korrigiert. Dann aus der letzten Gleichung folgt

$$|\lambda||a_{k,k}| \leq \sum_{j=k+1}^n |a_{k,j}| + |\lambda| \sum_{j=1}^{k-1} |a_{k,j}|$$

was zu (d) führt. So sind die Schritte (a) – (d) der Beweis, dass die Gauß-Seidel Iteration konvergiert, wenn A streng diagonal dominant ist.

3. Basierend auf Beispiel 8 auf dem 2. Übungsblatt seien die Matrix A und die Vektoren \mathbf{v} und \mathbf{u}^* gegeben durch den folgenden Matlab Code.

```
N = 50; N1 = N-1; h = 1/N;
mu = 0.01;
Q = spdiags([-ones(N1,1),ones(N1,1)], [0,1], N1, N); i = speye(N);
Dx = kron(i,Q); Dy = kron(Q,i); D = Dx'*Dx + Dy'*Dy; I = speye(N*N);
v = sin(1:N*N)'; A = I + mu*D/h^2; ustar = A \ v;
```

Das System $A\mathbf{u} = \mathbf{v}$ soll nun mit der SSOR Methode $\mathbf{u}^{(k+1)} = T_{\text{SSOR}}^{(\omega)} \mathbf{u}^{(k)} + M_{\text{SSOR}}^{(\omega)} \mathbf{v}$, $k \in \mathbb{N}$, iterativ gelöst werden. Seien $\{\mathbf{u}^{(k)}\}_{k=0}^M$ die Iterierten der SSOR Methode mit $\mathbf{u}^{(0)} = \mathbf{v}$. Der Dämpfungsparameter ω ist noch auszuwählen, aber der optimale Wert ist

$$\omega^* = \underset{\omega \in (0,2)}{\operatorname{argmin}} \rho(T_{\text{SSOR}}^{(\omega)}).$$

Kreuzen Sie bei den wahren Behauptungen an.

- (a) Sei $\rho(T_J)$ der Spektralradius der Iterationsmatrix T_J der Jacobi Methode. Dann ist ω^* gegeben durch

$$\frac{2}{1 + \sqrt{1 - [\rho(T_J)]^2}}.$$

- (b) Wenn die Vektoriteration auf Seite 100 im Skriptum für die Matrix $T_{\text{SSOR}}^{(\omega^*)}$ und mit dem Startvektor `ones(N*N,1)` angewendet wird, ergibt sich zu 4 signifikanten Ziffern der Grenzwert $\rho(T_{\text{SSOR}}^{(\omega^*)}) = 0.8594$ als Abschätzung des dominanten Eigenwerts.
- (c) Mit $\omega = 0.5$ ist $M = 465$ die kleinste Anzahl der SSOR Iterationen, mit der der relative Fehler $\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2 / \|\mathbf{u}^*\|_2$ unterhalb der Schwelle von einfacher Genauigkeit liegt.
- (d) Mit $\omega = 1.5$ ist $M = 119$ die kleinste Anzahl der SSOR Iterationen, mit der der relative Fehler $\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2 / \|\mathbf{u}^*\|_2$ unterhalb der Schwelle von einfacher Genauigkeit liegt.

Kommentare: Sehen Sie diesen [Code](#). Wie auf Seite 466 im Lehrbuch gesehen, ist die Formel im Teil (a) bekannt für $T_{\text{SOR}}^{(\omega)}$ (nicht SSOR) in dem Fall, dass A SPD und tridiagonal ist. Dies lässt sich für das reduzierte Beispiel $A = \text{speye}(N) + \text{mu}*Q'*Q/h^2$ bestätigen. Weitere Details und eine grafische Darstellung der Wirkung des Dämpfungsparameters befinden sich hier. Für das obige 2D Beispiel gelten $\omega^* = 1.679$ und $\rho(T_{\text{SSOR}}^{(\omega^*)}) = 0.8594$ zu 4 signifikanten Ziffern. Die richtig Anzahl der Iterationen für Teil (c) ist $M = 865$.

4. Basierend auf Beispiel 8 auf dem 2. Übungsblatt seien die Matrix A und die Vektoren \mathbf{v} und \mathbf{u}^* gegeben durch den folgenden Matlab Code.

```
N = 50; N1 = N-1; h = 1/N;
mu = 0.01;
Q = spdiags([-ones(N1,1),ones(N1,1)], [0,1], N1, N); i = speye(N);
Dx = kron(i,Q); Dy = kron(Q,i); D = Dx'*Dx + Dy'*Dy; I = speye(N*N);
v = sin(1:N*N)'; A = I + mu*D/h^2; ustar = A \ v;
```

Das System $A\mathbf{u} = \mathbf{v}$ soll nun mit der Methode der konjugierten Gradienten iterativ gelöst werden. Seien $\{\mathbf{u}^{(k)}\}_{k=0}^{\infty}$ die Iterierten des Verfahrens mit $\mathbf{u}^{(0)} = \mathbf{v}$. Seien das A -Skalarprodukt und die A -Norm durch $(\mathbf{x}, \mathbf{y})_A = \mathbf{y}^\top A \mathbf{x}$ bzw. $\|\mathbf{x}\|_A = (\mathbf{x}, \mathbf{x})_A^{1/2}$ definiert. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Seien die Matrizen Q und D im Code definiert. Seien $\{(\lambda_j, \mathbf{v}_j)\}_{i=1}^N$ die Eigenwert-Eigenvektor-Paare für $Q^\top Q$, wie in den Kommentaren für Beispiel 5 auf Übungsblatt 3 detailliert. Dann sind die Eigenwerte der Matrix D gegeben durch $\{\lambda_i + \lambda_j\}_{i,j=1}^{N^2}$.
- (b) Es gilt $\kappa_2(A) = 1 + 8\mu/h^2$.
- (c) Es gelten

$$\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_A \leq 2\|\mathbf{u}^{(0)} - \mathbf{u}^*\|_A \left[\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right]^k, \quad k \in \mathbb{N}.$$

- (d) Es gelten

$$\|\mathbf{u}^{(k)} - \mathbf{u}^*\|_A \leq 2\|\mathbf{u}^{(0)} - \mathbf{u}^*\|_A \left[\frac{\sqrt{\kappa_2(A)} - 1}{\sqrt{\kappa_2(A)} + 1} \right]^{2k}, \quad k \in \mathbb{N}.$$

Kommentare: Sehen Sie diesen [Code](#). Anhand vom Teil (a) ist $\{1 + \mu(\lambda_i + \lambda_j)/h^2\}_{i,j=1}^{N^2}$ das Spektrum von A . Mit $\lambda_i = 4\sin^2(\pi(i-1)/(2N))$ liegt das Spektrum im Intervall $[1, 1+8\mu/h^2]$. Jedoch gilt $\lambda_i + \lambda_j = 8$ für kein Paar (i, j) , und daher gilt $\kappa_2(A) < 1+8\mu/h^2$. Weitere Details über die Methode der konjugierten Gradienten befinden sich hier, und die Fehlerabschätzung im Teil (c) erscheint auf Seite 17.

5. Basierend auf Beispiel 8 auf dem 2. Übungsblatt seien die Matrix A und die Vektoren \mathbf{v} und \mathbf{u}^* gegeben durch den folgenden Matlab Code.

```
N = 50; N1 = N-1; h = 1/N;
mu = 0.01;
Q = spdiags([-ones(N1,1),ones(N1,1)], [0,1], N1, N); i = speye(N);
Dx = kron(i,Q); Dy = kron(Q,i); D = Dx'*Dx + Dy'*Dy; I = speye(N*N);
v = sin(1:N*N)'; A = I + mu*D/h^2; ustar = A \ v;
```

Das System $A\mathbf{u} = \mathbf{v}$ soll nun mit der präkonditionierten Methode der konjugierten Gradienten iterativ gelöst werden. Wie auf Seite 94 im Skriptum angedeutet ist eine SPD Matrix C mit $A \approx C^2$ auszuwählen. Seien $\{\mathbf{u}^{(k)}\}_{k=0}^{\infty}$ die Iterierten des Verfahrens mit $C = I$ und $\mathbf{u}^{(0)} = \mathbf{v}$. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Mit $B = C^2$ lässt sich die präkonditionierte Methode der konjugierten Gradienten zur Lösung des Systems $A\mathbf{u} = \mathbf{v}$ mit \mathbf{kmax} Iterationen folgendermaßen implementieren.

```

u  = v;
r  = v - A*u;
z  = B \ r;
r1 = z'*r; r2 = r1; a = 0*r;
for k=1:kmax
    s  = r1/r2;
    a  = z + s*a;
    w  = A*a;
    t  = r1/(a'*w);
    u  = u + t*a;
    r  = r - t*w;
    z  = B \ r;
    r2 = r1;
    r1 = z'*r;
end

```

- (b) Die approximierten Inversen für die Jacobi Methode und die Gauß-Seidel Methode seien durch M_J^{-1} bzw. M_{SGS}^{-1} bezeichnet. Es gelten

$$\kappa_2(M_J^{-1/2}AM_J^{-1/2}) < \kappa_2(M_{SGS}^{-1/2}AM_{SGS}^{-1/2}) < \kappa_2(A).$$

- (c) Seien $\{\mathbf{u}_J^{(k)}\}_{k=0}^{\infty}$ die Iterierten des Verfahrens mit $C^2 = M_J$ und $\mathbf{u}_J^{(0)} = \mathbf{v}$. Es gelten

$$\|\mathbf{u}_J^{(k)} - \mathbf{u}^*\|_2 < \|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2, \quad k = 1, \dots, N.$$

- (d) Seien $\{\mathbf{u}_{SGS}^{(k)}\}_{k=0}^{\infty}$ die Iterierten des Verfahrens mit $C^2 = M_{SGS}$ und $\mathbf{u}_{SGS}^{(0)} = \mathbf{v}$. Es gelten

$$\|\mathbf{u}_{SGS}^{(k)} - \mathbf{u}^*\|_2 < \|\mathbf{u}^{(k)} - \mathbf{u}^*\|_2, \quad k = 1, \dots, N.$$

Kommentare: Sehen Sie diesen [Code](#). Die Methode der konjugierten Gradienten auf Seite 93 im Skriptum hat die folgende Form für das präkonditionierte System $(C^{-1}AC^{-1})(C\mathbf{x}) = C^{-1}\mathbf{b}$ auf Seite 94,

$$\begin{aligned}
& C\mathbf{x}^{(0)} \text{ sei gegeben} \\
\mathbf{v}^{(1)} &= (C^{-1}\mathbf{r}^{(0)}) = (C^{-1}\mathbf{b}) - (C^{-1}AC^{-1})(C\mathbf{x}^{(0)}) \\
& \text{Für } k = 1, 2, \dots \\
t^{(k)} &= (C^{-1}\mathbf{r}^{(k-1)}) \cdot (C^{-1}\mathbf{r}^{(k-1)}) / \mathbf{v}^{(k)} \cdot (C^{-1}AC^{-1})\mathbf{v}^{(k)} \\
(C\mathbf{x}^{(k)}) &= (C\mathbf{x}^{(k-1)}) + t^{(k)}\mathbf{v}^{(k)} \\
(C^{-1}\mathbf{r}^{(k)}) &= (C^{-1}\mathbf{r}^{(k-1)}) - t^{(k)}(C^{-1}AC^{-1})\mathbf{v}^{(k)} \\
s^{(k)} &= (C^{-1}\mathbf{r}^{(k)}) \cdot (C^{-1}\mathbf{r}^{(k)}) / (C^{-1}\mathbf{r}^{(k-1)}) \cdot (C^{-1}\mathbf{r}^{(k-1)}) \\
\mathbf{v}^{(k+1)} &= (C^{-1}\mathbf{r}^{(k)}) + s^{(k)}\mathbf{v}^{(k)}
\end{aligned}$$

Mit $\tilde{\mathbf{v}}^{(k)} = C^{-1}\mathbf{v}^{(k)}$ folgt der Algorithmus, der im Teil (a) implementiert ist,

$$\begin{aligned}
 & \mathbf{x}^{(0)} \text{ sei gegeben} \\
 \tilde{\mathbf{v}}^{(1)} &= C^{-2}\mathbf{r}^{(0)}, \quad \mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)} \\
 \text{Für } k &= 1, 2, \dots \\
 t^{(k)} &= \mathbf{r}^{(k-1)} \cdot C^{-2}\mathbf{r}^{(k-1)} / \tilde{\mathbf{v}}^{(k)} \cdot A\tilde{\mathbf{v}}^{(k)} \\
 \mathbf{x}^{(k)} &= \mathbf{x}^{(k-1)} + t^{(k)}\tilde{\mathbf{v}}^{(k)} \\
 \mathbf{r}^{(k)} &= \mathbf{r}^{(k-1)} - t^{(k)}A\tilde{\mathbf{v}}^{(k)} \\
 s^{(k)} &= \mathbf{r}^{(k)} \cdot C^{-2}\mathbf{r}^{(k)} / \mathbf{r}^{(k-1)} \cdot C^{-2}\mathbf{r}^{(k-1)} \\
 \tilde{\mathbf{v}}^{(k+1)} &= \mathbf{r}^{(k)} + s^{(k)}\tilde{\mathbf{v}}^{(k)}
 \end{aligned}$$

mit dem nur die Matrix C^{-2} erscheint und nicht ihre Wurzel. Dieser Algorithmus befindet sich auf Seite 23 hier.

Für eine SPD Matrix A ist die Konditionszahl $\kappa_2(A)$ gegeben durch den Quotienten $\max\{\sigma(A)\} / \min\{\sigma(A)\}$, wobei $\sigma(A)$ das Spektrum der Matrix A bezeichnet. Für eine SPD Matrix $M = V\Lambda V^\top$, Λ diagonal und V orthogonal, ist die pte Potenz der Matrix durch $M^p = V\Lambda^p V^\top$ definiert. Jedoch wird das Spektrum der Matrix $M^{-\frac{1}{2}}AM^{-\frac{1}{2}}$ bequem durch die verallgemeinerten Eigenwerte λ des Systems $A\mathbf{x} = \lambda M\mathbf{x}$ bestimmt. Folglich sind Konditionszahlen für Teil (b) gegeben durch

```

E_A = eig(full(A));
K_A = max(E_J)/min(E_J); % = 200.8027
                           % ist auch cond(full(A))

L = tril(A,-1);
U = triu(A,+1);
D = A - L - U;
M_J = D;
E_J = eig(full(A),full(M_J));
K_J = max(E_J)/min(E_J); % = 197.2977
                           % ist auch cond(full(Mh*A*Mh)) mit
                           % Mh = D^(-0.5)

M_SGS = (D + U);
M_SGS = D \ MSGS;
M_SGS = (D + L)*M_SGS;
E_SGS = eig(full(A),full(M_SGS)); %
K_SGS = max(E_SGS)/min(E_SGS); % = 25.4398
                           % ist auch cond(full(Mh*A*Mh)) mit
                           % Mh = V*Md^(-0.5)*V' wobei
                           % [V,Md] = eig(full(M_SGS))

```

Da $\kappa_2(M_J^{-1/2}AM_J^{-1/2}) < \kappa_2(A)$ gilt, ist zu erwarten, dass (c) richtig ist. Dies ist aber nicht der Fall, wie die Rechnungen zeigen. Der Unterschied zwischen den Konditionszahlen 200.8027 und 197.2977 ist nicht ausreichend. Jedoch ist $\kappa_2(M_{SGS}^{-1/2}AM_{SGS}^{-1/2}) = 25.4398$ viel kleiner als $\kappa_2(A) = 200.8027$, und die Rechnungen zeigen, dass (d) richtig ist.

6. Basierend auf Beispiel 7 auf dem 3. Übungsblatt seien $\{\mathbf{W}^k\}_{k=0}^N \subset \mathbb{R}^N$ gegeben durch die folgende Diskretisierung der Wärmegleichung.

```

N = 50; N1 = N-1; h = 1/N; ta = h; ka = 1;
x = linspace(0,1,N+1); x = x(1:N)' + h/2;
Q = spdiags([-ones(N1,1),ones(N1,1)], [0,1], N1, N);
D = Q'*Q; B = -ka*D/h^2; I = speye(N); A = I - ta*B;
W = sign(x-0.5);
for k=1:N
    W = A\W;
end

```

Die Systeme $A\mathbf{W}^k = \mathbf{W}^{k-1}$, $k \in \mathbb{N}$, sollen nun mit der Methode der konjugierten Gradienten iterativ gelöst werden, und die iterativ berechneten Lösungen seien mit $\{\tilde{\mathbf{W}}^k\}_{k=0}^N$ bezeichnet. Anfänglich gilt $\tilde{\mathbf{W}}^0 = \mathbf{W}^0$. Weiters seien die Iterierten beim k ten Zeitschritt mit $\{\tilde{\mathbf{W}}^{k,l}\}_{l=0}^M$ bezeichnet. Diese Iteration endet mit $\tilde{\mathbf{W}}^k = \tilde{\mathbf{W}}^{k,M}$ und startet mit $\tilde{\mathbf{W}}^{k,0}$, das noch auszuwählen ist. Kreuzen Sie bei den wahren Behauptungen an.

(a) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ gelten für jede Anzahl $M = 1, \dots, 10$ der Iterationen,

$$|(\mathbf{W}^k)_i| < |(\tilde{\mathbf{W}}^k)_i|, \quad i = 1, \dots, N, \quad k = 1, \dots, N.$$

(b) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ gelten für jede Anzahl $M = 1, \dots, 10$ der Iterationen,

$$(\mathbf{W}^k)_i < (\mathbf{W}^k)_{i+1}, \quad (\tilde{\mathbf{W}}^k)_i < (\tilde{\mathbf{W}}^k)_{i+1}, \quad i = 1, \dots, N-1, \quad k = 0, \dots, N.$$

(c) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ ist $M = 25$ die kleinste Anzahl der Iterationen, mit der der relative Fehler $\|\{\tilde{\mathbf{W}}^k - \mathbf{W}^k\}_{k=0}^N\|_\infty / \|\{\mathbf{W}^k\}_{k=0}^N\|_\infty$ unterhalb der Schwelle von einfacher Genauigkeit liegt.

(d) Mit $\tilde{\mathbf{W}}^{k,0} = \text{ones}(N, 1)$ ist $M = 25$ die kleinste Anzahl der Iterationen, mit der der relative Fehler $\|\{\tilde{\mathbf{W}}^k - \mathbf{W}^k\}_{k=0}^N\|_\infty / \|\{\mathbf{W}^k\}_{k=0}^N\|_\infty$ unterhalb der Schwelle von einfacher Genauigkeit liegt.

Kommentare: Sehen Sie diesen [Code](#). Der Diffusionsprozess in der Wärmegleichung ist dämpfend, und daher ist die Abweichung des Zustandes \mathbf{W}^k von seinem Fliessgleichgewicht eine monoton fallende Funktion von k . Weiters bleiben die Komponenten $\mathbf{W}^k = \{W_i^k\}_{i=1}^N$ eine steigende Funktion von i für jedes k . Insofern die iterativ berechneten Zustände $\tilde{\mathbf{W}}^k$ akkurat sind, widerspiegeln sie diese qualitativen Eigenschaften der Lösung der Wärmegleichung. Obwohl gewisse iterative Verfahren ihre eigenen dämpfenden Eigenschaften genießen, ist die Methode der konjugierten Gradienten nicht derartig.

Mit $M = 25$ Iterationen gilt die Ungleichung im Teil (a) nicht mehr. Mit $M = 5, 6, 7$ gelten die Ungleichungen im Teil (b) nicht. Da die Methode der konjugierten Gradienten ein System in einer endlichen Anzahl von Iterationen mit exakter Arithmetik lösen soll, sinkt der relative Fehler schlagartig von $5.84081\text{e-}02$ für $M = 24$ auf $3.88578\text{e-}15$ für $M = 25$ im Teil (c). Da die Zustände \mathbf{W}^k und \mathbf{W}^{k-1} nicht sehr unterschiedlich sind, ist \mathbf{W}^{k-1} ein sinnvoller Startvektor für die Iteration im k ten Zeitschritt. Wenn dies nicht verwendet wird, sind $M = 26$ Iterationen im Teil (d) notwendig.

7. Basierend auf Beispiel 8 auf dem 3. Übungsblatt seien $\{\mathbf{W}^k\}_{k=0}^N \subset \mathbb{R}^N$ gegeben durch die folgende Diskretisierung der Wellengleichung.

```

N = 50; N1 = N-1; h = 1/N; ta = h; nu = 1;
x = linspace(0,1,N+1); x = x(1:N)' + h/2;
Q = spdiags([-ones(N1,1),ones(N1,1)], [0,1], N1, N);
B = (nu/h)*[sparse(N1,N1), Q;-Q', sparse(N,N)];
I = speye(2*N-1); A = I - (ta/2)*B; C = I + (ta/2)*B;
a = cos(pi*(2*x-1));
W = [(nu/h)*Q*a; 0*a];
for k=1:N
    Ws = W;
    W = A\ (C*Ws);
    a = a + (ta/2)*(W(N:(2*N-1)+Ws(N:2*N-1)));
end

```

Die Systeme $A\mathbf{W}^k = C\mathbf{W}^{k-1}$, $k \in \mathbb{N}$, sollen nun mit der symmetrischen Gauß-Seidel Methode gelöst werden, und die iterativ berechneten Lösungen seien mit $\{\tilde{\mathbf{W}}^k\}_{k=0}^N$ bezeichnet. Anfänglich gilt $\tilde{\mathbf{W}}^0 = \mathbf{W}^0$. Die Auslenkungen seien durch

$$\mathbf{a}^k = \mathbf{a}^{k-1} + (\tau/2)\{(\mathbf{W}^k + \mathbf{W}^{k-1})_i\}_{i=N}^{2N-1}, \quad k = 1, \dots, N$$

und

$$\tilde{\mathbf{a}}^k = \tilde{\mathbf{a}}^{k-1} + (\tau/2)\{(\tilde{\mathbf{W}}^k + \tilde{\mathbf{W}}^{k-1})_i\}_{i=N}^{2N-1}, \quad k = 1, \dots, N$$

gegeben. Weiters seien die Iterierten beim k ten Zeitschritt durch $\{\tilde{\mathbf{W}}^{k,l}\}_{l=0}^M$ bezeichnet. Diese Iteration endet mit $\tilde{\mathbf{W}}^k = \tilde{\mathbf{W}}^{k,M}$ und startet mit $\tilde{\mathbf{W}}^{k,0}$, das noch auszuwählen ist. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ gelten für jede Anzahl $M = 1, \dots, 10$ der Iterationen zu 3 signifikanten Ziffern,

$$\|\mathbf{W}^k\|_2 = \|\tilde{\mathbf{W}}^k\|_2 = 31.4, \quad k = 0, \dots, N.$$

- (b) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ gelten für jede Anzahl $M = 1, \dots, 10$ der Iterationen,

$$|(\mathbf{a}^k)_i| < |(\tilde{\mathbf{a}}^k)_i|, \quad i = 1, \dots, N, \quad k = 1, \dots, N.$$

- (c) Mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ ist $M = 4$ die kleinste Anzahl der Iterationen, mit der der relative Fehler $\|\{\tilde{\mathbf{W}}^k - \mathbf{W}^k\}_{k=0}^N\|_\infty / \|\{\mathbf{W}^k\}_{k=0}^N\|_\infty$ unterhalb der Schwelle von einfacher Genauigkeit liegt.
- (d) Mit $\tilde{\mathbf{W}}^{k,0} = 10^{-5}\text{sign}(\tilde{\mathbf{W}}^{k-1})$ ist $M = 780$ die kleinste Anzahl der Iterationen, mit der der relative Fehler $\|\{\tilde{\mathbf{W}}^k - \mathbf{W}^k\}_{k=0}^N\|_\infty / \|\{\mathbf{W}^k\}_{k=0}^N\|_\infty$ unterhalb der Schwelle von einfacher Genauigkeit liegt.

Kommentare: Sehen Sie diesen [Code](#). Die Methode auf Seiten 228-229 im Skriptum für die Wellengleichung ist im obigen Code mit “\” implementiert. Wie im Skriptum angedeutet, entspricht $\{W_i^k\}_{i=1}^{N-1} = \nu \mathbf{a}_x^k$ der potentiellen Energie der Schwingungen und $\{W_i^k\}_{i=N}^{2N-1} \approx \mathbf{a}_t^k$ entspricht der kinetischen Energie. Also ist $\|\mathbf{W}^k\|_2$ die gesamte Energie zur Zeit $t^k = k\tau$, die im Lauf der Zeit erhalten bleiben soll. Diese Energieerhaltung gelingt der Methode im obigen Code, und deswegen gilt $\|\mathbf{W}^k\|_2 = 31.4$, $k = 0, \dots, N$. Jedoch wenn nur $M = 1$ Iteration der symmetrischen Gauß-Seidel Methode anstatt “\” verwendet wird, gilt $\|\tilde{\mathbf{W}}^k\|_2 = 31.4$, $k = 0, \dots, N$, nicht mehr. Obwohl die symmetrische Gauß-Seidel

Methode eine dämpfende Wirkung für Teil (b) im Beispiel 5 genießen könnten, gibt es kein Gegenstück für Teil (b) hier.

Da die Zustände \mathbf{W}^k und \mathbf{W}^{k-1} nicht sehr unterschiedlich sind, ist \mathbf{W}^{k-1} ein sinnvoller Startvektor für die Iteration im k ten Zeitschritt. Wenn dies nicht verwendet wird, sind $M = 780$ Iterationen im Teil (d) notwendig. So viele Iterationen sind auch notwendig, weil der Startvektor $\tilde{\mathbf{W}}^{k,0} = \{\tilde{W}_i^{k,0}\}_{i=1}^{2N-1}$ keine glatte Funktion von i ist. Mit einem glatten Startvektor $\tilde{\mathbf{W}}^{k,0} = \mathbf{0}$ sind die relativen Fehler weniger mit $1.54786\text{e-}06$ für $M = 3$ und $6.05741\text{e-}09$ für $M = 4$. Jedoch mit $\tilde{\mathbf{W}}^{k,0} = \tilde{\mathbf{W}}^{k-1}$ sind die relativen Fehler noch besser mit $1.92030\text{e-}07$ für $M = 3$ und $7.57108\text{e-}10$ für $M = 4$.

8. Die folgenden Matrizen seien definiert:

$$A = \begin{bmatrix} +1 & +1 & +1 \\ -1 & +1 & +1 \\ -1 & -1 & +1 \end{bmatrix}, \quad B = \begin{bmatrix} +8 & -2 & +4 \\ -1 & +7 & -2 \\ -3 & +3 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} +2 & -1 & -1 \\ -1 & +2 & +1 \\ +1 & +1 & +2 \end{bmatrix}.$$

Seien $\{\mu^{(k)}\}_{k=1}^{\infty}$ und $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ durch die Vektoriteration auf Seite 100 oder die inverse Vektoriteration auf Seite 102 im Skriptum gegeben, wobei $\mathbf{x}^{(0)}$ durch

```
x = randn(3,1);
x = x/norm(x,inf);
```

bestimmt wird. Kreuzen Sie bei den wahren Behauptungen an.

- (a) Für A führt die Vektoriteration zum Ergebnis $\mu^{(k)} \xrightarrow{k \rightarrow \infty} \rho(A)$.
- (b) Für B führt die Vektoriteration zum Ergebnis, dass $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ zu einem Vielfachen von $(1, 1, 0)$ konvergiert.
- (c) Für C führt die Vektoriteration zum Ergebnis $|\mu^{(k)} - 3|/(2/3)^k \xrightarrow{k \rightarrow \infty} c$, $c \in (0, \infty)$.
- (d) Für C führt die inverse Vektoriteration mit $q = 0$ zum Ergebnis $\mu^{(k)} \xrightarrow{k \rightarrow \infty} 1$.

Kommentare: Sehen Sie diesen [Code](#). Die Matrix A hat Eigenwerte $\{1, 1 \pm \sqrt{3}i\}$ und daher wird $\rho(A) = |1 \pm \sqrt{3}i|$ gegeben durch 2 betragsmäßig gleiche Eigenwerte. Daher sind die Voraussetzungen der Vektoriteration nicht erfüllt, und die Rechnungen mit den Anfangsbedingungen zeigen, dass die Werte $\{\mu^{(k)}\}$ zum Spektralradius von A nicht konvergieren.

Die Matrix B hat die Eigenwerte $\{3, 6\}$ und einen 2-dimensionalen Eigenraum für den größten Eigenwert. Daher sind die Voraussetzungen der Vektoriteration nicht erfüllt, und die Rechnungen mit den Anfangsbedingungen zeigen, dass die Vektoren $\{\mathbf{x}^{(k)}\}$ nicht eindeutig konvergieren.

Die Matrix C hat die Eigenwerte $\{1, 2, 3\}$, die Voraussetzungen der Vektoriteration sind erfüllt und die Konvergenzrate der Werte $\{\mu^{(k)}\}$ hängt vom Quotienten $2/3$ der zwei größten Eigenwerte ab. Die inverse Vektoriteration mit $q = 0$ ist die Vektoriteration für die Matrix C^{-1} mit Eigenwerten $\{1/3, 1/2, 1\}$. Die Voraussetzungen der Vektoriteration sind erfüllt, und die Werte $\{\mu^{(k)}\}$ konvergieren zum Wert $\rho(C^{-1}) = 1$.