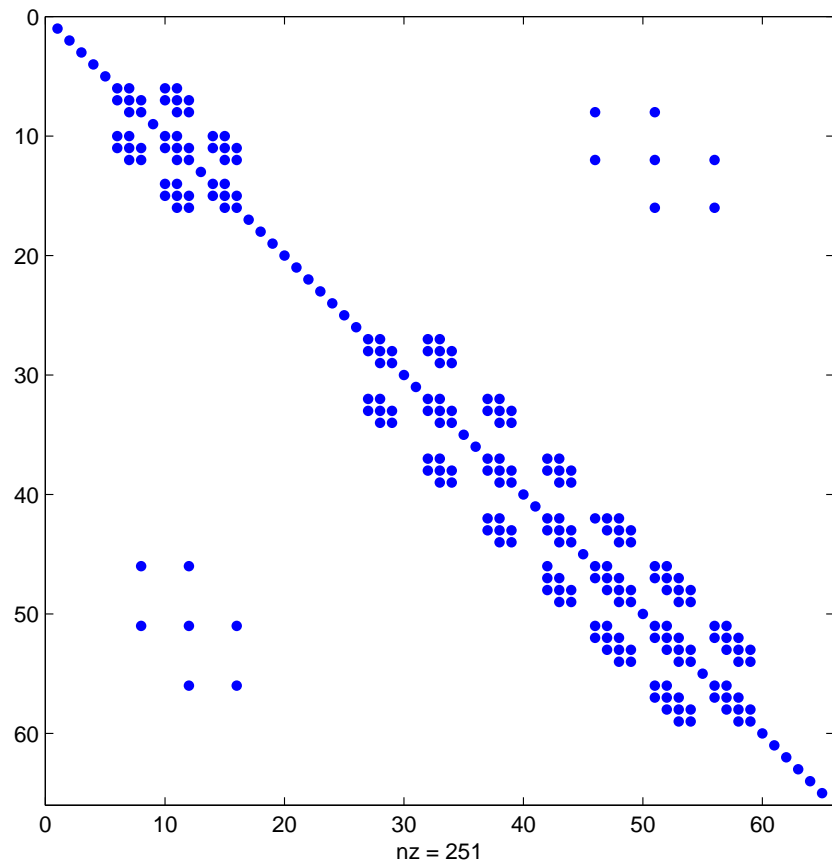


4 Die Lösung der diskreten Poisson-Gleichung

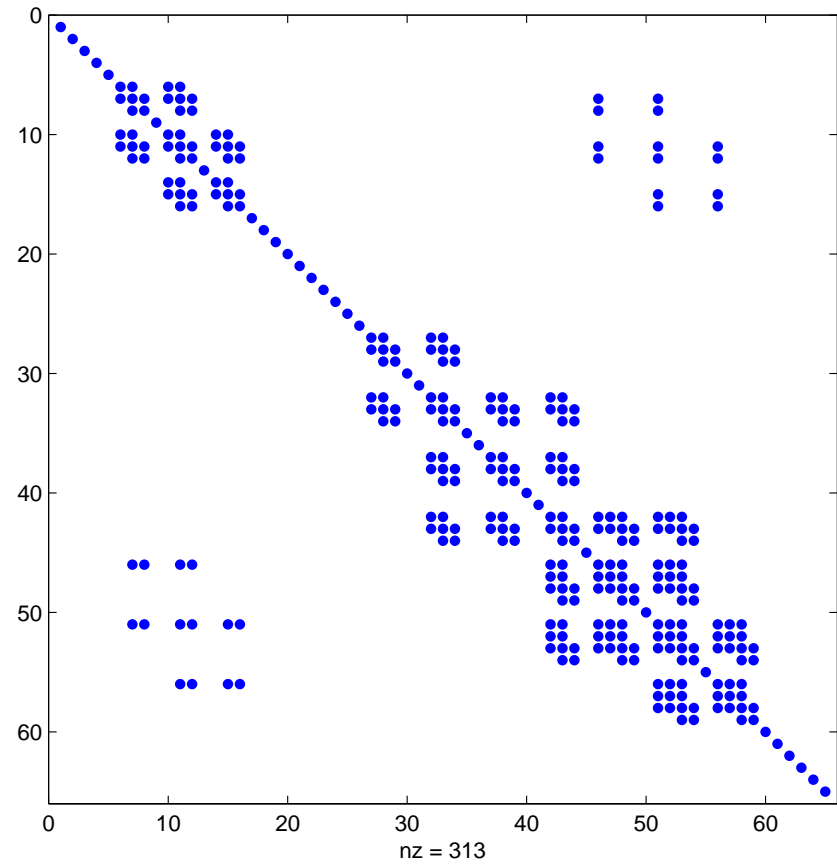
Die Galerkin-Matrix der FE-Diskretisierung der Poisson-RWA (3.1) ist symmetrisch und positiv-definit (gilt nicht notwendig für andere Dgln.)

Allen FE-Diskretisierungen gemeinsam: die Galerkin-Matrix ist **dünn besetzt** (sparse). Wir betrachten die Matrix aus Beispiel 2 (L-Gebiet).

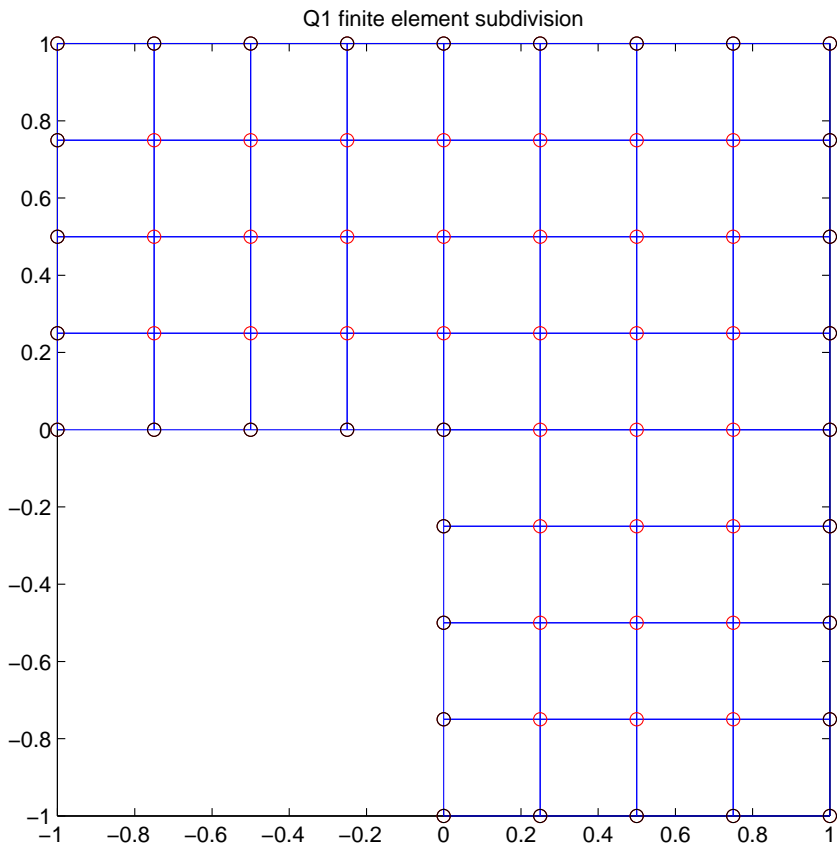
N	bilinear			biquadratisch		
	h	# Einträge $\neq 0$	Anteil	h	# Einträge $\neq 0$	Anteil
65	0.25	251	5.94%	0.5	313	7.4%
225	0.125	1339	2.64%	0.25	2073	4.1%
883	0.0625	6107	0.78%	0.125	10141	1.5%
3201	0.0312	26011	0.25%	0.0625	44767	0.44%
12545	0.0156	107291	0.068%	0.0312	187742	0.12%



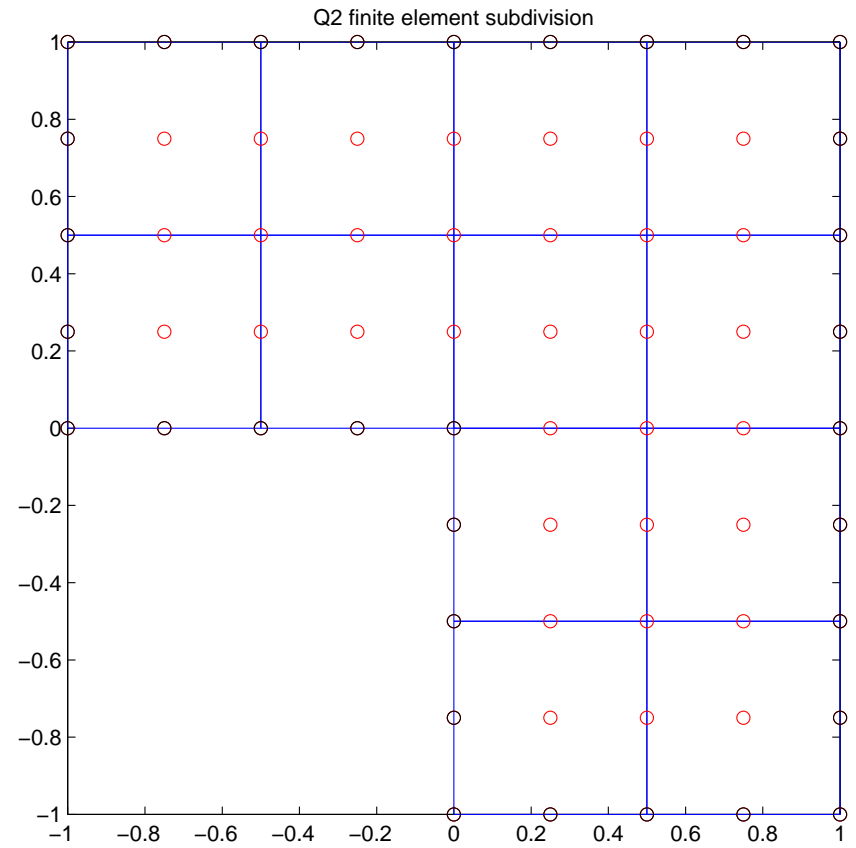
$h = 0.25$, bilineare Elemente



$h = 0.5$, biquadratische Elemente



$h = 0.25$, bilineare Elemente



$h = 0.5$, biquadratische Elemente

Direkte oder iterative Verfahren ?

- In der FE-Praxis wurde lange Zeit mit direkten Verfahren, also Varianten der Gauß-Elimination, gearbeitet.
- Stark verbreitet: **frontal solvers**. Hier wird die Faktorisierung der Matrix so früh wie möglich während der Assemblierung vorgenommen.
- Weiterentwicklungen auch heute noch relevant, auch als eigenständige Eliminationsverfahren, z.B. UMFPACK.
- Faustregel: in 2D reichen ausgereifte direkte, auf dünn-besetzte Matrizen spezialisierte Löser aus. Für 3D-Diskretisierungen muss auf **Iterationsverfahren** zurückgegriffen werden.
- Vorteil direkter Verfahren: typischerweise als „black box“ einsetzbar. Bei Iterationsverfahren ist deutlich mehr Expertise des Anwenders erforderlich.

4.1 Erste Iterationsverfahren

Zunächst sei allgemein gegeben

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{C}^{n \times n} \text{ invertierbar, } \mathbf{b} \in \mathbb{C}^n. \quad (4.1)$$

Grundidee:

Startnäherung $\mathbf{x}_0 \approx \mathbf{A}^{-1}\mathbf{b}$,

Erzeuge rekursiv Folge $\{\mathbf{x}_m\}_{m \in \mathbb{N}}$ mit $\lim_{m \rightarrow \infty} \mathbf{x}_m = \mathbf{A}^{-1}\mathbf{b}$.

Bezeichnungen:

$$\mathbf{r}_m := \mathbf{b} - \mathbf{A}\mathbf{x}_m$$

$$\mathbf{x}^* := \mathbf{A}^{-1}\mathbf{b}$$

$$\mathbf{e}_m := \mathbf{x}^* - \mathbf{x}_m$$

Residuum, Defekt,

(exakte) Lösung,

Fehler.

Die ersten Iterationsverfahren (**Handrechnung!**) beruhen auf Umstellung einzelner Gleichungen, Auflösen nach den entsprechenden Unbekannten und Wiedereinsetzen (**Relaxation**).

- Carl-Friedrich Gauß (~ 1832), Normalgleichungen aus Astronomie, Landesvermessung, 20–40 Unbekannte

„Ich empfehle Ihnen diesen Modus zur Nachahmung. Schwerlich werden Sie je wieder direkt eliminieren, wenigstens nicht, wenn Sie mehr als 2 Unbekannte haben. Das indirekte Verfahren lässt sich halb im Schläfe ausführen, oder man kann während desselben an andere Dinge denken.“

Brief von Gauß an Gerling, 1823

- Carl Gustav Jacobi (~ 1846)
- Ludwig Seidel (~ 1862)
- Christian August Nagel (Sächsische Triangulation, 1867-1886) 159 Unbekannte

Formale Beschreibung von Relaxationsverfahren durch **Zerlegungen** (splittings):

$$A = M - N, \quad M \text{ invertierbar}, \quad (4.2a)$$

überführt (4.1) in Fixpunktform $Mx = Nx + b$, zugehörige Fixpunktiteration

$$x_{m+1} = Tx_m + c, \quad T = M^{-1}N, \quad c = M^{-1}b. \quad (4.2b)$$

Typische algorithmische Umsetzung: (beachte: $Tx_m + c = x_m + M^{-1}r_m$)

Algorithmus 1 : Relaxationsverfahren.

Gegeben : A invertierbar, b , Startnäherung x_0

1 for $m = 0, 1, 2 \dots$ *bis Abbruchkriterium erfüllt* **do**

2 $r_m \leftarrow b - Ax_m$

3 Löse $Mh = r_m$

4 $x_{m+1} \leftarrow x_m + h$

Klassische Zerlegungen

- $M = I$
Richardson-Verfahren
- $M = D := \text{diag}(A)$
Jacobi-Verfahren, Einzelschrittverfahren
- $M = D - L, L = \text{tril}(A)$
Gauß-Seidel-Verfahren, Einzelschrittverfahren, Liebmann-Verfahren
- $M = D - \omega L, \omega \in \mathbb{R}$
SOR-Verfahren (successive-overrelaxation)
- $M = \frac{\omega}{2-\omega} \left(\frac{1}{\omega} D - L \right) D^{-1} \left(\frac{1}{\omega} D - U \right), U := \text{triu}(A), \omega \in \mathbb{R}$
SSOR-Verfahren (symmetric SOR)

Für den Fehler der m -ten Iterierten aus (4.2b) gilt

$$\mathbf{e}_{m+1} = \mathbf{x}^* - \mathbf{x}_{m+1} = (\mathbf{I} - \mathbf{M}^{-1}\mathbf{A})(\mathbf{x}^* - \mathbf{x}_m) = \mathbf{T}\mathbf{e}_m = \cdots = \mathbf{T}^m \mathbf{e}_0.$$

Lemma 4.1 Für eine beliebige Matrix $\mathbf{T} \in \mathbb{C}^{n \times n}$ gilt

$$\lim_{m \rightarrow \infty} \mathbf{T}^m = \mathbf{O} \quad \Leftrightarrow \quad \rho(\mathbf{T}) < 1.$$

Satz 4.2 Das durch die Zerlegung (4.2a) definierte Relaxationsverfahren (4.2b) konvergiert genau dann für alle Startnäherungen \mathbf{x}_0 , wenn für die Iterationsmatrix $\mathbf{T} = \mathbf{M}^{-1}\mathbf{N}$ gilt $\rho(\mathbf{T}) < 1$.

4.2 Das CG-Verfahren

Satz 4.3 *Die Iterierten, Residuen und Fehler des durch die Zerlegung (4.2a) definierten Relaxationsverfahrens (4.2b) sind gegeben durch*

$$\mathbf{x}_m = \mathbf{x}_0 + \sum_{j=0}^{m-1} (\mathbf{I} - \mathbf{M}^{-1} \mathbf{A})^j \mathbf{M}^{-1} \mathbf{r}_0,$$

$$\mathbf{M}^{-1} \mathbf{r}_m = (\mathbf{I} - \mathbf{M}^{-1} \mathbf{A})^m \mathbf{M}^{-1} \mathbf{r}_0,$$

$$\mathbf{e}_m = (\mathbf{I} - \mathbf{M}^{-1} \mathbf{A})^m \mathbf{e}_0 \quad m = 1, 2, \dots$$

also durch

$$\mathbf{x}_m = \mathbf{x}_0 + q_{m-1}(\mathbf{T}) \mathbf{M}^{-1} \mathbf{r}_0, \quad \mathbf{M}^{-1} \mathbf{r}_m = p_m(\mathbf{T}) \mathbf{M}^{-1} \mathbf{r}_0, \quad \mathbf{e}_m = p_m(\mathbf{T}) \mathbf{e}_0$$

mit Polynomen

$$q_{m-1}(\lambda) = 1 + \lambda + \dots + \lambda^{m-1} \in \mathcal{P}_{m-1} \quad \text{bzw.} \quad p_m(\lambda) = \lambda^m \in \mathcal{P}_m .$$

Das CG-Verfahren (conjugate gradients), oder allgemein **polynomiale Beschleunigungsverfahren**, auch **Krylov-Unterraumverfahren** genannt, können als Iterationsverfahren interpretiert werden, bei denen die Iterierten bzw. Residuen die Form

$$\mathbf{x}_m = \mathbf{x}_0 + q_{m-1}(\mathbf{A})\mathbf{r}_0, \quad \mathbf{r}_m = p_m(\mathbf{A})\mathbf{r}_0,$$

mit

$$q_{m-1} \in \mathcal{P}_{m-1}, \quad p_m(\lambda) = 1 - \lambda q_{m-1}(\lambda) \in \mathcal{P}_m$$

besitzen, wobei versucht wird, diese Polynome in jedem Schritt möglichst günstig zu wählen.

Entscheidend: in jedem Schritt der Iteration ist nur eine Matrix-Vektor-Multiplikation mit der Koeffizientenmatrix \mathbf{A} (bzw. der Iterationsmatrix \mathbf{T}) erforderlich.

4.2.1 Herleitung des CG-Verfahrens

Gegeben: $A \in \mathbb{R}^{n \times n}$ spd, $b \in \mathbb{R}^n$ (komplexer Fall analog).

Gesucht: Lösung x^* des LGS

$$Ax = b. \quad (4.3)$$

Sei (\cdot, \cdot) ein gegebenes Innenprodukt auf \mathbb{R}^n (mit zugehöriger Norm $\|\cdot\|$).

Wichtige Beobachtung: da A spd^a ist durch

$$(u, v)_A := (Au, v), \quad u, v \in \mathbb{R}^n,$$

ebenfalls ein Innenprodukt gegeben, das **A-Innenprodukt** oder auch **Energie-Innenprodukt** (mit zugehöriger Norm $\|\cdot\|_A$).

Da mit A auch A^{-1} spd, gilt dasselbe auch für $(\cdot, \cdot)_{A^{-1}}$.

^aeigentlich: A selbstadjungiert bez. (\cdot, \cdot)

Ausgangspunkt dieser Herleitung des CG-Verfahrens (vgl. Numerik-Vorlesung, Abschnitt über Variationsrechnung):

$$x^* \text{ löst (4.3)} \iff x^* \text{ minimiert } J(x) = \frac{1}{2}(x, x)_A - (b, x) \text{ über } \mathbb{R}^n.$$

Minimierung auf Unterräumen: Anstatt auf ganz \mathbb{R}^n minimieren wir J auf einer Folge geschachtelter affiner Unterräume der Gestalt

$$\mathcal{S}_m = x_0 + \mathcal{V}_m, \quad \mathcal{V}_m = \text{span}\{p_1, p_2, \dots, p_m\}.$$

Hierbei bilden p_1, \dots, p_m eine Folge linear unabhängiger **Richtungsvektoren**.

Als Näherungen $x_m \approx x^*$ erhalten wir die Lösungen der Unterraumaufgaben

$$\text{Bestimme } x_m \in \mathcal{S}_m \text{ mit } J(x_m) = \min_{x \in \mathcal{S}_m} J(x).$$

Wegen

$$J(\mathbf{x}) = \frac{1}{2} (\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{A^{-1}}^2 - \|\mathbf{b}\|_{A^{-1}}^2) = \frac{1}{2} \|\mathbf{x}^* - \mathbf{x}\|_A^2 + J(\mathbf{x}^*)$$

ist die Minimierung von J äquivalent zur Minimierung von $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{A^{-1}}$, d.h. gesucht ist ein $\mathbf{x} \in \mathcal{S}_m$ welches den Ausdruck

$$\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{A^{-1}} = \|\mathbf{b} - \mathbf{A}\mathbf{x}_0 - \mathbf{A}\mathbf{p}\|_{A^{-1}} = \|\mathbf{r}_0 - \mathbf{A}\mathbf{p}\|_{A^{-1}}$$

minimiert mit einem $\mathbf{p} \in \mathcal{V}_m$.

Äquivalent: bestimme Bestapproximation bez. $\|\cdot\|_{A^{-1}}$ an \mathbf{r}_0 aus dem Unterraum $\mathbf{A}\mathcal{V}_m$.

Äquivalent: \mathbf{A}^{-1} -orthogonale Projektion von \mathbf{r}_0 nach $\mathbf{A}\mathcal{V}_m$

Charakterisierung:

$$\mathbf{r}_0 - \mathbf{A}\mathbf{p} \perp_{A^{-1}} \mathbf{A}\mathbf{p}_j, \quad j = 1, 2, \dots, m,$$

oder

$$(\mathbf{p}, \mathbf{p}_j)_A = (\mathbf{r}_0, \mathbf{p}_j), \quad j = 1, 2, \dots, m,$$

oder

$$\mathbf{r}_m := \mathbf{b} - \mathbf{A}\mathbf{x}_m \perp \mathcal{V}_m.$$

Besitzt \mathbf{p} die Koordinatendarstellung

$$\mathbf{p} = y_1 \mathbf{p}_1 + \dots + y_m \mathbf{p}_m,$$

so erfüllt der Koordinatenvektor $\mathbf{y} = [y_1, y_2, \dots, y_m]^\top \in \mathbb{R}^m$ das LGS

$$[(\mathbf{p}_j, \mathbf{p}_i)_A]_{i,j=1}^m \mathbf{y} = [(\mathbf{r}_0, \mathbf{p}_i)]_{i=1}^m. \quad (4.4)$$

Bemerkung: Wegen $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{A^{-1}} = \|\mathbf{A}^{-1}\mathbf{b} - \mathbf{x}\|_A$ besitzt die Bestapproximation $\mathbf{p} \in \mathcal{V}_m$ auch die Eigenschaft, dass die zugehörige Approximation $\mathbf{x}_m = \mathbf{x}_0 + \mathbf{p}$ den Fehler $\mathbf{A}^{-1}\mathbf{b} - \mathbf{x}_m$ bezüglich der A -Norm minimiert.

Konjugierte Suchrichtungen

Sind die Richtungsvektoren p_j **A-konjugiert** (ein Synonym für *A*-orthogonal), so vereinfacht sich das LGS (4.4) für die Koordinaten:

- Die Koeffizientenmatrix wird diagonal,
- die Gleichungen sind entkoppelt, d.h. beim Übergang von \mathcal{V}_m nach \mathcal{V}_{m+1} bleiben die Koeffizienten y_1, \dots, y_m unverändert,
- es gilt

$$y_i = \frac{(\mathbf{r}_0, \mathbf{p}_i)}{(\mathbf{p}_i, \mathbf{p}_i)_A}, \quad i = 1, 2, \dots, m.$$

Wir erhalten so folgenden Algorithmus:

Algorithmus 2 : Verfahren der konjugierten Richtungen (CD-Verf.)

Gegeben : A spd, b

1 Bestimme Startnäherung x_0 , berechne $r_0 \leftarrow b - Ax_0$, $m \leftarrow 0$

2 **while** nicht konvergiert **do**

3 $m \leftarrow m + 1$

4 Bestimme nächsten, zu p_1, \dots, p_{m-1} konjugierten Richtungsvektor p_m

5 $\alpha_m \leftarrow \frac{(r_{m-1}, p_m)}{(p_m, p_m)_A}$

6 $x_m \leftarrow x_{m-1} + \alpha_m p_m$

7 $r_m \leftarrow r_{m-1} - \alpha_m A p_m$

Bemerkung: Eine Folge A -orthogonaler Suchrichtungen kann man erzeugen aus einer beliebigen Folge (linear unabhängiger) Vektoren $\tilde{p}_1, \tilde{p}_2, \dots$ mit Hilfe des **Gram-Schmidtschen Orthogonalisierungsverfahrens**:

$$p_m = \tilde{p}_m - \sum_{j=1}^{m-1} \frac{(\tilde{p}_m, p_j)_A}{(p_j, p_j)_A} p_j, \quad m = 1, 2, \dots \quad (4.5)$$

Lemma 4.4 *Im m -ten Schritt des CD-Verfahrens gilt:*

$$\mathbf{p}_m \perp_A \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{m-1}\}, \quad (4.6)$$

$$\mathbf{r}_m \perp \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}, \quad (4.7)$$

$$(\mathbf{p}_m, \mathbf{r}_{m-1}) = (\mathbf{p}_m, \mathbf{r}_{m-2}) = \dots = (\mathbf{p}_m, \mathbf{r}_0). \quad (4.8)$$

Lemma 4.5 *Das CD-Verfahren berechnet die Lösung x^* von (4.3) nach einer endlichen Anzahl $L \leq n$ Schritten.*

Bemerkungen:

1. Die endliche Abbrucheigenschaft geht i.A. verloren bei rundungsbehafteter Rechnung.
2. Obwohl der Fehler die interessierende Größe, ist oft nur das Residuum verfügbar. Der nächste Satz gibt prinzipielle Beziehungen zwischen diesen an.

Satz 4.6 Sei x eine Approximation an $A^{-1}b$. Dann gelten für den Fehler $e := A^{-1}b - x$ und das zugehörige Residuum $r = b - Ax$

$$\frac{\|r\|}{\|r\|_A} \leq \frac{\|e\|_A}{\|r\|} \leq \frac{\|e\|}{\|e\|_A}, \quad (4.9)$$

$$\frac{\|r\|^2}{\|r\|_A^2} \leq \frac{\|e\|}{\|r\|} \leq \frac{\|e\|^2}{\|e\|_A^2}. \quad (4.10)$$

Lemma 4.7 Im m -ten Schritt des CD-Verfahrens wird der Abstand in der A -Norm zur Lösung reduziert gemäß

$$\|e_{m-1}\|_A^2 - \|e_m\|_A^2 = \frac{(\mathbf{r}_{m-1}, \mathbf{p}_m)^2}{(\mathbf{p}_m, \mathbf{p}_m)_A}.$$

Konjugierte Gradienten

A -orthogonale Suchrichtungen machen einen Iterationsschritt zwar effizient, müssen aber per se nicht zu schneller Konvergenz führen. (Warum?)

Idee: Wähle in (4.5) den neuen Vektor $\tilde{\mathbf{p}}_{m+1}$ als Richtung stärkster Abnahme von J an der Stelle \mathbf{x}_m :

$$\begin{aligned}\tilde{\mathbf{p}}_{m+1} &= -\nabla J|_{\mathbf{x}=\mathbf{x}_m} = \mathbf{b} - \mathbf{A}\mathbf{x}_m = \mathbf{r}_m, & m > 1, \\ \tilde{\mathbf{p}}_1 &= \mathbf{p}_1 = \mathbf{r}_0.\end{aligned}$$

Danach A -Orthogonalisierung gegen $\{\mathbf{p}_1, \dots, \mathbf{p}_m\}$.

Vorteil: Liefert solange neue Suchrichtungen bis Lösung gefunden. (Warum?)

Aus dieser Wahl folgt sofort

$$\text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{m-1}\} = \text{span}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}, \quad m = 1, 2, \dots, L. \quad (4.11)$$

Wir erhalten so den

Algorithmus 3 : CG-Verfahren (vorläufige Version).

Gegeben : A spd, b

1 Bestimme Startnäherung x_0 , berechne $r_0 \leftarrow b - Ax_0$

2 Setze $m \leftarrow 0$, $p_1 \leftarrow r_0$

3 **while** *nicht konvergiert* **do**

4 $m \leftarrow m + 1$

5 $\alpha_m \leftarrow \frac{(r_{m-1}, p_m)}{(p_m, p_m)_A}$

6 $x_m \leftarrow x_{m-1} + \alpha_m p_m$

7 $r_m \leftarrow r_{m-1} - \alpha_m A p_m$

8 $p_{m+1} \leftarrow r_m - \sum_{j=1}^m \frac{(r_m, p_j)_A}{(p_j, p_j)_A} p_j$

Vereinfachungen:

(1) Solange $m < L$ ist $\alpha_m \neq 0$ und aus $\mathbf{r}_m = \mathbf{r}_{m-1} - \alpha_m \mathbf{A} \mathbf{p}_m$ folgt zunächst $\mathbf{A} \mathbf{p}_m \in \text{span}\{\mathbf{r}_{m-1}, \mathbf{r}_m\}$ und allgemein

$$\text{span}\{\mathbf{A} \mathbf{p}_1, \mathbf{A} \mathbf{p}_2, \dots, \mathbf{A} \mathbf{p}_{m-1}\} \subset \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{m-1}\} \\ \stackrel{(4.11)}{=} \text{span}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}.$$

Da aber $\mathbf{r}_m \perp \text{span}\{\mathbf{p}_1, \dots, \mathbf{p}_m\}$ nach (4.7) folgt hieraus $\mathbf{r}_m \perp_A \text{span}\{\mathbf{p}_1, \dots, \mathbf{p}_{m-1}\}$ sodass Zeile 8 sich vereinfacht zu

$$\mathbf{p}_{m+1} \leftarrow \mathbf{r}_m - \frac{(\mathbf{r}_m, \mathbf{p}_m)_A}{(\mathbf{p}_m, \mathbf{p}_m)_A} \mathbf{p}_m =: \mathbf{r}_m + \beta_{m+1} \mathbf{p}_m. \quad (4.12)$$

(2) Aus $\mathbf{p}_m = \mathbf{r}_{m-1} + \beta_m \mathbf{p}_{m-1}$ folgt wegen $\mathbf{r}_{m-1} \perp \mathbf{p}_{m-1}$ für Zeile 5

$$\alpha_m = \frac{(\mathbf{r}_{m-1}, \mathbf{p}_m)}{(\mathbf{p}_m, \mathbf{p}_m)_A} = \frac{(\mathbf{r}_{m-1}, \mathbf{r}_{m-1})}{(\mathbf{p}_m, \mathbf{p}_m)_A}. \quad (4.13)$$

(3) Aus $\mathbf{r}_m = \mathbf{r}_{m-1} - \alpha_m \mathbf{A} \mathbf{p}_m$ folgt

$$(\mathbf{r}_m, \mathbf{r}_m) = -\alpha_m (\mathbf{r}_m, \mathbf{A} \mathbf{p}_m) = -\alpha_m (\mathbf{r}_m, \mathbf{p}_m)_A$$

und daher wegen (4.13) für (4.12)

$$\beta_{m+1} = -\frac{(\mathbf{r}_m, \mathbf{p}_m)_A}{(\mathbf{p}_m, \mathbf{p}_m)_A} = \frac{(\mathbf{r}_m, \mathbf{r}_m)}{\alpha_m (\mathbf{p}_m, \mathbf{p}_m)_A} = \frac{(\mathbf{r}_m, \mathbf{r}_m)}{(\mathbf{r}_{m-1}, \mathbf{r}_{m-1})}.$$

Mit diesen Vereinfachungen erhalten wir die endgültige Formulierung des CG-Verfahrens:

Algorithmus 4 : CG-Verfahren.

Gegeben : A spd, b

1 Bestimme Startnäherung x_0 , berechne $r_0 \leftarrow b - Ax_0$

2 Setze $m \leftarrow 0$, $p_1 \leftarrow r_0$

3 **while** *nicht konvergiert* **do**

4 $m \leftarrow m + 1$

5 $\alpha_m \leftarrow \frac{(r_{m-1}, r_{m-1})}{(p_m, p_m)_A}$

6 $x_m \leftarrow x_{m-1} + \alpha_m p_m$

7 $r_m \leftarrow r_{m-1} - \alpha_m A p_m$

8 $\beta_{m+1} \leftarrow \frac{(r_m, r_m)}{(r_{m-1}, r_{m-1})}$

9 $p_{m+1} \leftarrow r_m + \beta_{m+1} p_m$

4.2.2 Eigenschaften des CG-Verfahrens

Aufwand: Operationen pro Schritt: 1 Matrix-Vektor Produkt,
2 Innenprodukte,
3 Vektoraufdatierungen.

Speicher: 4 Vektoren.

Lemma 4.8 *Für die Suchrichtungen des CG-Verfahrens gilt*

$$\mathbf{p}_m = \|\mathbf{r}_{m-1}\|^2 \sum_{j=0}^{m-1} \frac{\mathbf{r}_j}{\|\mathbf{r}_j\|^2}, \quad m = 1, 2, \dots, L. \quad (4.14)$$

Satz 4.9 *Beim CG-Verfahren gelten die Beziehungen*

$$(\mathbf{r}_m, \mathbf{p}_j) = 0, \quad 1 \leq j \leq m, \quad 1 \leq m \leq L - 1, \quad (4.15)$$

$$(\mathbf{r}_m, \mathbf{r}_j) = 0, \quad m \neq j, \quad 0 \leq m, j \leq L - 1, \quad (4.16)$$

$$(\mathbf{r}_m, \mathbf{p}_j) = \|\mathbf{r}_{j-1}\|^2, \quad m + 1 \leq j \leq L, \quad 0 \leq m \leq L - 1, \quad (4.17)$$

$$(\mathbf{p}_m, \mathbf{p}_j)_A = 0, \quad m \neq j, \quad 1 \leq m, j \leq L, \quad (4.18)$$

$$(\mathbf{p}_m, \mathbf{r}_{m-1})_A = \|\mathbf{p}_m\|_A^2, \quad 1 \leq m \leq L, \quad (4.19)$$

$$(\mathbf{p}_m, \mathbf{r}_j)_A = 0, \quad m + 1 \leq j \leq L - 1, \quad m = 1, 2, \dots, L, \quad (4.20)$$

Lemma 4.10 *Für die Residuen und Richtungsvektoren des CG-Verfahrens gelten*

$$(\mathbf{p}_m, \mathbf{p}_j) = \frac{\|\mathbf{r}_{j-1}\|^2 \|\mathbf{p}_m\|^2}{\|\mathbf{r}_{m-1}\|^2}, \quad j \geq m, \quad (4.21)$$

$$\|\mathbf{p}_m\|^2 = \|\mathbf{r}_{m-1}\|^2 + \beta_m^2 \|\mathbf{p}_{m-1}\|^2 = \|\mathbf{r}_{m-1}\|^4 \sum_{j=0}^{m-1} \frac{1}{\|\mathbf{r}_j\|^2}. \quad (4.22)$$

Lemma 4.11 *Im m -ten Schritt des CG-Verfahrens gilt*

$$\|\mathbf{e}_{m-1}\|_A^2 - \|\mathbf{e}_m\|_A^2 = \alpha_m \|\mathbf{r}_{m-1}\|^2 = \frac{\|\mathbf{p}_m\|_A^2}{\|\mathbf{p}_m\|^2} \|\mathbf{x}_m - \mathbf{x}_{m-1}\|^2. \quad (4.23)$$

Für $j < m$ gilt

$$\|\mathbf{e}_j\|_A^2 - \|\mathbf{e}_m\|_A^2 = \alpha_{j+1} \|\mathbf{r}_j\|^2 + \cdots + \alpha_m \|\mathbf{r}_{m-1}\|^2. \quad (4.24)$$

Die m -te Iterierte besitzt die Darstellung

$$\mathbf{x}_m = \mathbf{x}_0 + \sum_{j=1}^m \alpha_j \mathbf{p}_j = \mathbf{x}_0 + \sum_{j=1}^m \frac{\|\mathbf{e}_j\|_A^2 - \|\mathbf{e}_m\|_A^2}{\|\mathbf{r}_{j-1}\|^2} \mathbf{r}_{j-1}. \quad (4.25)$$

Satz 4.12 *Im m -ten Schritt des CG-Verfahrens gilt*

$$\|\mathbf{e}_{m-1}\|^2 - \|\mathbf{e}_m\|^2 = \frac{\|\mathbf{p}_m\|^2}{\|\mathbf{p}_m\|_A^2} [\|\mathbf{e}_m\|_A^2 + \|\mathbf{e}_{m-1}\|_A^2]. \quad (4.26)$$

Insbesondere ist der CG-Fehler auch bezüglich der Norm $\|\cdot\|$ monoton fallend.

CG als Krylov-Unterraumverfahren

Für $A \in \mathbb{R}^{n \times n}$ und $v \in \mathbb{R}^n$ heißt der Unterraum

$$\mathcal{K}_m(A, v) := \text{span}\{v, Av, \dots, A^{m-1}v\} \subset \mathbb{R}^n$$

der m -te Krylov-Unterraum von A und v .

Es ist

$$\mathcal{K}_m(A, v) = \{p(A)v : p \in \mathcal{P}_{m-1}\}.$$

Proposition 4.13

(a) Die m -te CG-Iterierte liegt im affinen Unterraum $x_0 + \mathcal{K}_m(A, r_0)$.

(b) Für den Abbruchindex L des CG-Verfahrens gilt

$$\begin{aligned} L &= \min\{m : \mathcal{K}_m(A, r_0) = \mathcal{K}_{m+1}(A, r_0)\} \\ &= \min\{\deg q : q \text{ monisch und } q(A)r_0 = \mathbf{0}\}. \end{aligned}$$

4.2.3 Konvergenz des CG-Verfahrens

Die Rolle der Konditionszahl

- Bei einer $n \times n$ -Matrix A liefert das CG -Verfahren die Lösung von $Ax = b$ nach $L \leq n$ Schritten (in exakter Arithmetik). In diesem Sinne „konvergiert“ das CG-Verfahren immer gegen $x^* = A^{-1}b$. Gesucht sind hier obere Schranken für den Fehler $x^* - x_m$ ($m < L$).
- Von entscheidender Bedeutung ist die Optimalitätseigenschaft

$$\|x^* - x_m\|_A = \min_{p \in \mathcal{P}_m, p(0)=1} \|p(A)(x^* - x_0)\|_A.$$

- Sei $\Omega \subset (0, \infty)$ eine kompakte Menge, welche die Eigenwerte von A enthält. Für jedes $p_m \in \mathcal{P}_m$ mit $p_m(0) = 1$ gilt dann

$$\frac{\|x^* - x_m\|_A}{\|x^* - x_0\|_A} \leq \|p_m\|_{\infty, \Omega} := \max_{\lambda \in \Omega} |p_m(\lambda)|.$$

- Alle echten Fehlerschranken für das CG-Verfahren basieren auf obiger (trivialer) Beobachtung: ausgehend von einer kompakten Obermenge $\Omega \supset \Lambda(\mathbf{A})$ (typischerweise $\Omega = [\lambda_{\min}(\mathbf{A}), \lambda_{\max}(\mathbf{A})]$) versuchen wir, Polynome p_m zu konstruieren mit $p_m(0) = 1$, die auf Ω ein möglichst kleines Maximum haben. Dies führt auf natürliche Weise zu einer Aufgabe der **polynomialen Bestapproximation**:
Bestimme $p_m \in \mathcal{P}_m$, $p_m(0) = 1$, sodass

$$\|p_m\|_{\infty, \Omega} = \min_{\substack{p \in \mathcal{P}_m, \\ p(0)=1}} \|p\|_{\infty, \Omega}.$$

- Für $\Omega = \Lambda(\mathbf{A})$ ist die Abschätzung

$$\frac{\|\mathbf{x}^* - \mathbf{x}_m\|_A}{\|\mathbf{x}^* - \mathbf{x}_0\|_A} \leq \min_{\substack{p \in \mathcal{P}_m, \\ p(0)=1}} \max_{\lambda \in \Lambda(\mathbf{A})} |p(\lambda)|$$

in folgendem Sinne scharf: Für jedes m existiert ein Startvektor $\mathbf{x}_0 = \mathbf{x}_0(m)$, sodass im m -ten Schritt Gleichheit gilt.

- Die Aufgabe, $p_m \in \mathcal{P}_m$, $p_m(0) = 1$, so zu bestimmen, dass

$$\|p_m\|_{\infty, \Omega} = \min_{\substack{p \in \mathcal{P}_m, \\ p(0)=1}} \|p\|_{\infty, \Omega}$$

ist lösbar für jedes $m \in \mathbb{N}_0$ und jede kompakte Teilmenge Ω der komplexen Ebene. Die Lösung ist eindeutig, sofern Ω aus mindestens m Punkten besteht. Eine geschlossene Darstellung dieser Lösung existiert jedoch nur in Spezialfällen, etwa wenn Ω ein Intervall ist.

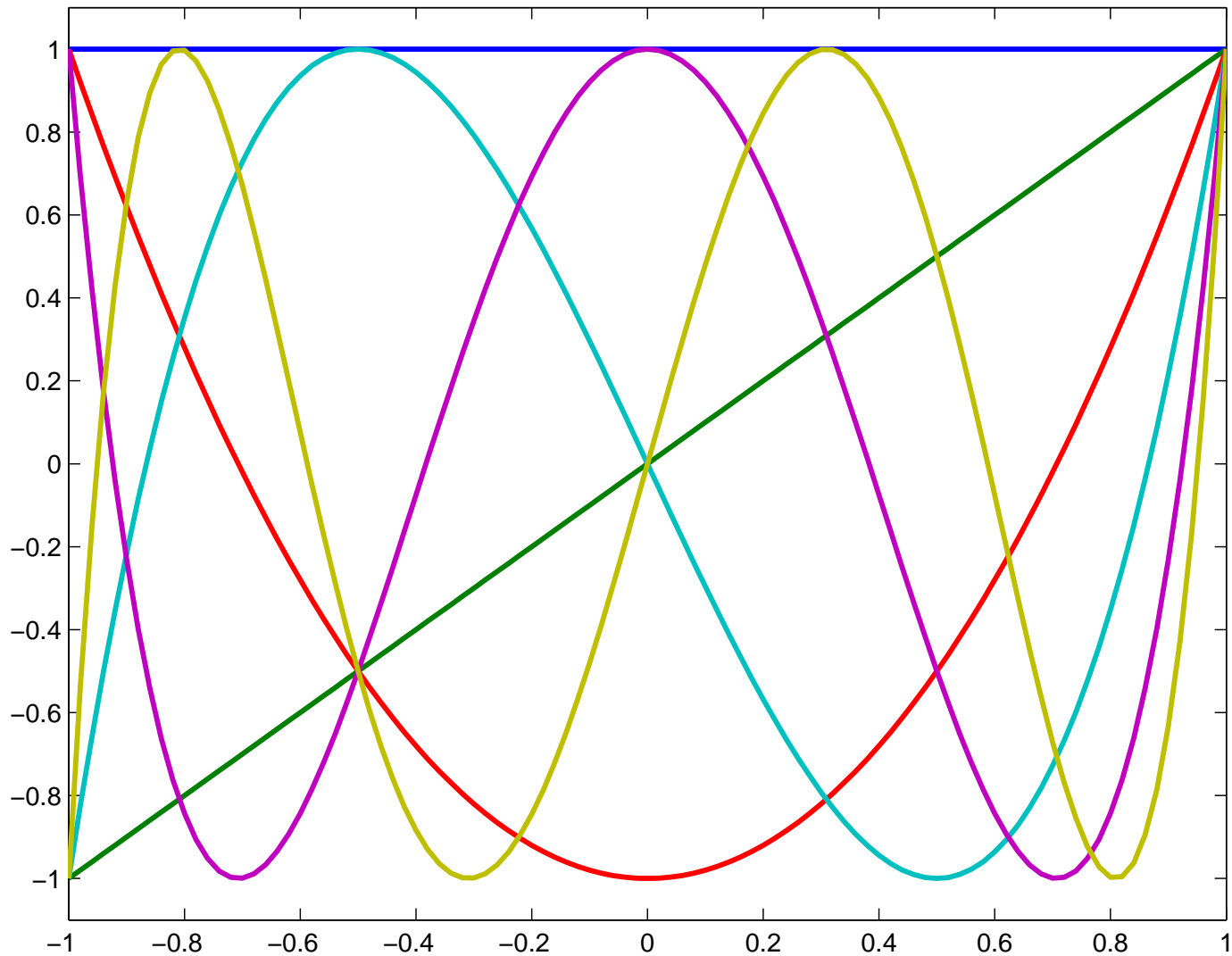
- Die **Chebyshev-Polynome** sind definiert durch

$$T_m(\xi) := \cos(m \arccos \xi) \in \mathcal{P}_m \quad (\xi \in [-1, 1], m = 0, 1, \dots) \quad (4.27)$$

- Sie genügen der dreistufigen Rekursion

$$T_m(\xi) = 2\xi T_{m-1}(\xi) - T_{m-2}(\xi) \quad (m = 2, 3, \dots)$$

mit $T_0(\xi) = 1$ and $T_1(\xi) = \xi$. Hieran erkennt man, dass T_m tatsächlich ein Polynom vom Grad (exakt) m ist.

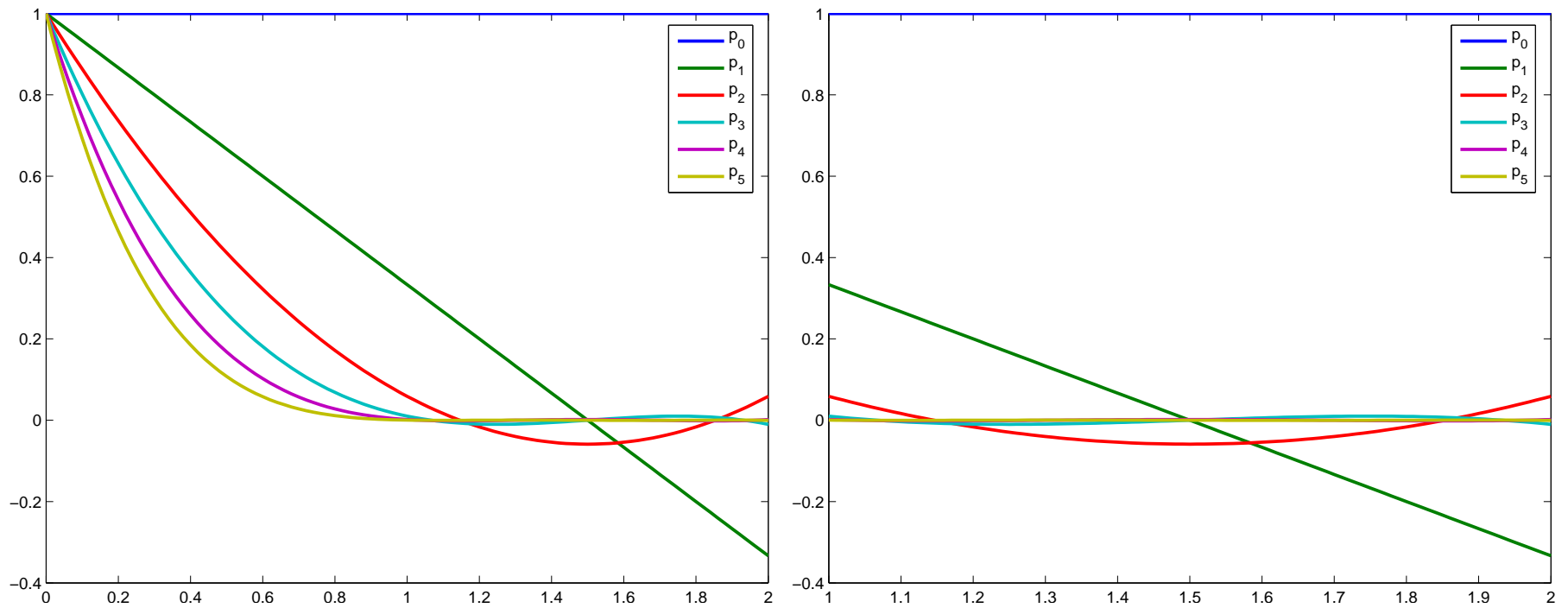


Graphen von T_0, T_1, \dots, T_5 auf $[-1, 1]$.

- T_m besitzt die m einfachen Nullstellen, $\xi_k^{(m)} = \cos((2k + 1)\pi/(2m))$, $k = 0, \dots, m - 1$, welche sämtlich in $(-1, 1)$ liegen.
- Es ist $|T_m(\xi)| \leq 1$ für alle $\xi \in [-1, 1]$. Gleichheit gilt für $\eta_k^{(m)} := \cos(k\pi/m)$, $k = 0, \dots, m$, genauer gesagt: $T_m(\eta_k^{(m)}) = (-1)^k$.
- Als nächstes transformieren wir die Polynome T_m auf ein beliebiges reelles Intervall $[\delta - \gamma, \delta + \gamma] \subset (0, +\infty)$, d.h., $0 < \gamma < \delta$.
 $\lambda \mapsto \ell(\lambda) := (\lambda - \delta)/\gamma$ ist eine Bijektion von $[\delta - \gamma, \delta + \gamma]$ auf $[-1, 1]$.
 Die Umkehrabbildung ist $\xi \mapsto \ell^{-1}(\xi) = \delta + \xi\gamma$.
 Wegen $-\delta/\gamma < -1$ gilt $T_m(-\delta/\gamma) \neq 0$ und

$$p_m(\lambda) := \frac{T_m(\ell(\lambda))}{T_m(\ell(0))} = \frac{T_m((\lambda - \delta)/\gamma)}{T_m(-\delta/\gamma)} \quad (m = 0, 1, \dots) \quad (4.28)$$

ist somit eine Folge von Residualpolynomen, d.h., $p_m(0) = 1$.



Graphen von p_0, p_1, \dots, p_5 auf $[0, 2]$ bzw. auf $[\delta - \gamma, \delta + \gamma] = [1, 2]$.

Lemma 4.14 Sei $\Omega := [\delta - \gamma, \delta + \gamma] \subset (0, \infty)$. Für die Polynome p_m aus (4.28) gilt

$$\|p_m\|_{\infty, \Omega} = \min_{\substack{p \in \mathcal{P}_m, \\ p(0)=1}} \|p\|_{\infty, \Omega}.$$

- Es verbleibt die Bestimmung von $T_m(-\delta/\gamma)$. Hierzu setzen wir

$$\theta := \frac{\delta - \sqrt{\delta^2 - \gamma^2}}{\gamma}.$$

- Es ist $0 < \theta < 1$ und

$$\max_{\lambda \in \Omega} |p_m(\lambda)| = \frac{1}{T_m(-\delta/\gamma)} = \frac{2\theta^m}{1 + \theta^{2m}} \leq 2\theta^m \quad (m = 0, 1, \dots).$$

Satz 4.15 (Fehlerschranke für CG) Für die Näherungen des CG-Verfahrens gilt die Fehlerabschätzung

$$\frac{\|\mathbf{x}^* - \mathbf{x}_m\|_A}{\|\mathbf{x}^* - \mathbf{x}_0\|_A} \leq \frac{1}{T_m\left(\frac{\kappa+1}{\kappa-1}\right)} = \frac{2\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^m}{1 + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2m}} \leq 2\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^m.$$

Hierbei bezeichnet $\kappa = \text{cond}_2(\mathbf{A}) = \lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$ die Konditionszahl von \mathbf{A} bezüglich der Euklid-Norm.

Diese Schranke ist (nahezu) scharf in folgendem Sinne: es gibt lineare Gleichungssysteme $\mathbf{A}\mathbf{x} = \mathbf{b}$ (mit $\mathbf{A} \in \mathbb{R}^{n \times n}$) sowie $\text{cond}_2(\mathbf{A}) = \kappa$ und Startnäherungen \mathbf{x}_0 sodass

$$\frac{\|\mathbf{x}^* - \mathbf{x}_{n-1}\|_A}{\|\mathbf{x}^* - \mathbf{x}_0\|_A} = \frac{1}{T_{n-1}\left(\frac{\kappa+1}{\kappa-1}\right)}$$

(Setze $\mathbf{A} = \text{diag}(\tilde{\eta}_1, \dots, \tilde{\eta}_n)$ mit $\tilde{\eta}_j$ die Extremalstellen des transformierten Chebyshev-Polynoms p_{n-1} .)

Abbruchbedingung

Typische Abbruchbedingung: iteriere bis $\|e_m\|_A / \|e_0\|_A \leq \epsilon$ mit vorgegebener Toleranz $\epsilon > 0$. Ausgehend von Satz 4.15 ist dies spätestens der Fall, wenn

$$\epsilon \geq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^m = 2 \left(1 - \frac{2/\sqrt{\kappa}}{1 + 1/\sqrt{\kappa}} \right)^m .$$

Logarithmieren und die Approximation $\log(1 + x) \approx x$ für x klein führt auf die Bedingung

$$m \geq \frac{1}{2} \left| \log \frac{\epsilon}{2} \right| \sqrt{\kappa} .$$

Beispiel: Poisson-Gleichung, uniformes Gitter, bilineare FE. Hier gilt $\kappa(\mathbf{A}) = O(h^{-2})$. Dies bedeutet: bei Halbierung der Gitterweite verdoppelt sich die Anzahl Iterationen.

Beispiel 1: die symmetrische 500×500 Matrix $\mathbf{A} = \mathbf{A}(\tau)$ definiert durch

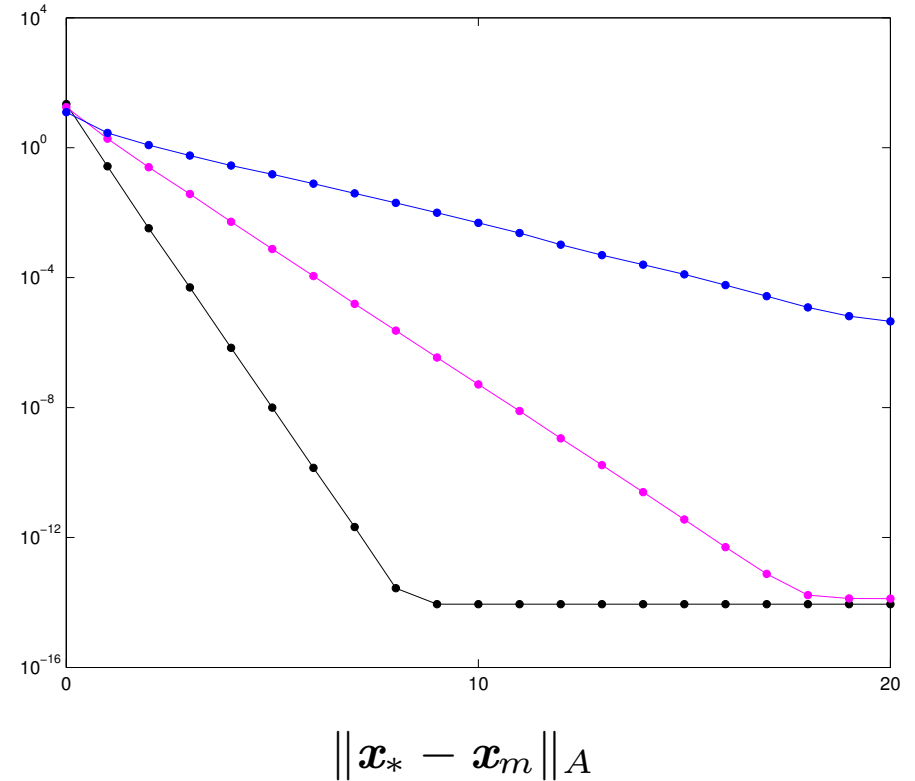
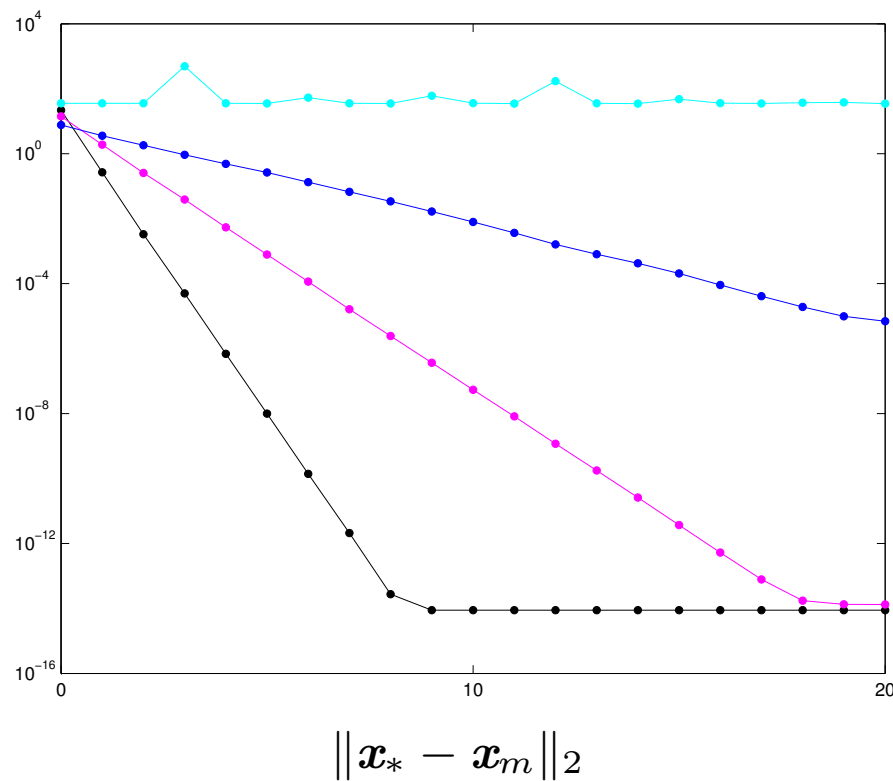
$$a_{i,j} = \begin{cases} 1 & \text{falls } i = j, \\ X_{i,j} & \text{falls } i \neq j \text{ und } X_{i,j} \leq \tau, \\ 0 & \text{sonst.} \end{cases}$$

wobei $X_{i,j}$ auf $[0, 1]$ gleichverteilte Zufallsvariable bezeichnet.

Man beachte: mit wachsendem τ nimmt die Anzahl von Null verschiedenen Einträgen in \mathbf{A} ebenso zu wie die Konditionszahl. Für $\tau = 0.2$ ist $\lambda_{\min}(\mathbf{A}) < 0$, \mathbf{A} also nicht mehr positiv definit. Wir erhalten so:

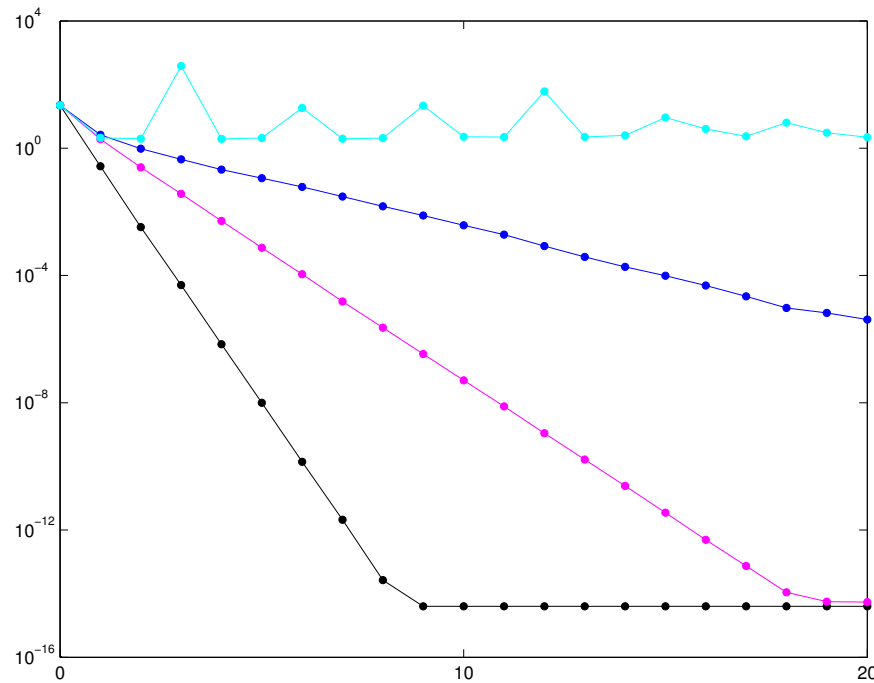
τ		$\text{nnz}(\mathbf{A})$	$\lambda_{\min}(\mathbf{A})$	$\lambda_{\max}(\mathbf{A})$	$\text{cond}_2(\mathbf{A})$	θ
0.01	(schwarz)	3016	0.9713	1.0325	1.0630	0.0153
0.05	(magenta)	12940	0.7150	1.6479	2.3049	0.2058
0.1	(blau)	25248	0.2097	3.5457	16.9059	0.6087
0.2	(cyan)	50510	-1.1677	11.1262	—	—

Wir wählen $\mathbf{b} = [1, 1, \dots, 1]^\top$, $\mathbf{x}_0 = \mathbf{0}$ und führen 20 Schritte des CG-Verfahrens durch.

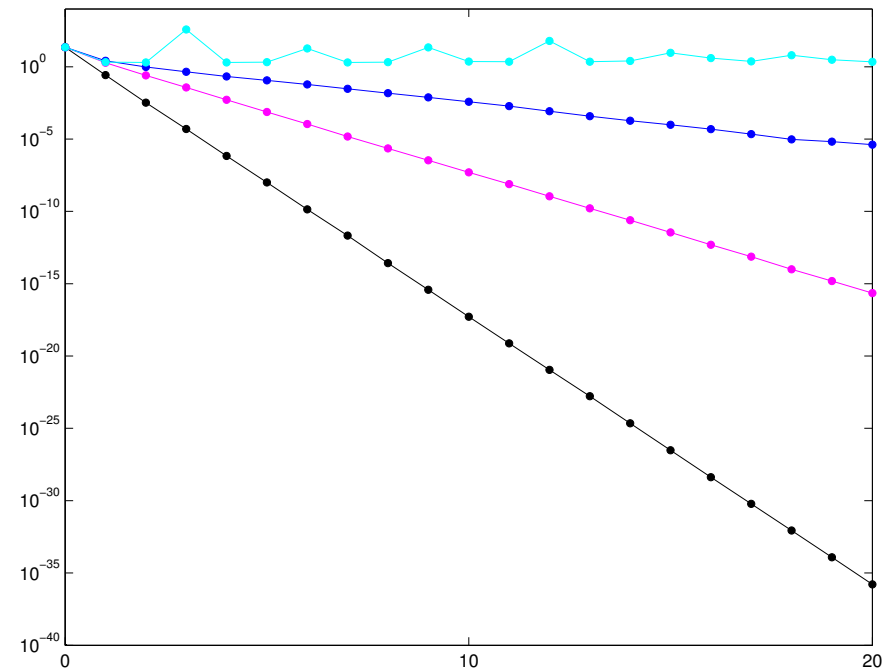


Beobachtungen:

- Divergenz im Fall $\tau = 0.2$.
- Im Fall $\tau = 0.1$ wird der Fehler nach 20 Schritten nur auf 10^{-5} reduziert, bei den übrigen Beispielen mit kleineren Konditionszahlen beobachten wir rasche Konvergenz.
- Für $\tau = 0.01$ erreicht der Fehler nach 9 Schritten den Bereich der Maschinengenauigkeit. In diesem Fall löst das CG-Verfahren $Ax = b$ in etwa 6×10^4 Flops und ist damit um den Faktor 700 schneller als die Cholesky-Zerlegung.



$$\|r_m\|_2, \text{ mit } r_m = b - Ax_m$$



$$\|r_m\|_2, \text{ mit } r_m = r_{m-1} - \alpha_m Ap_m$$

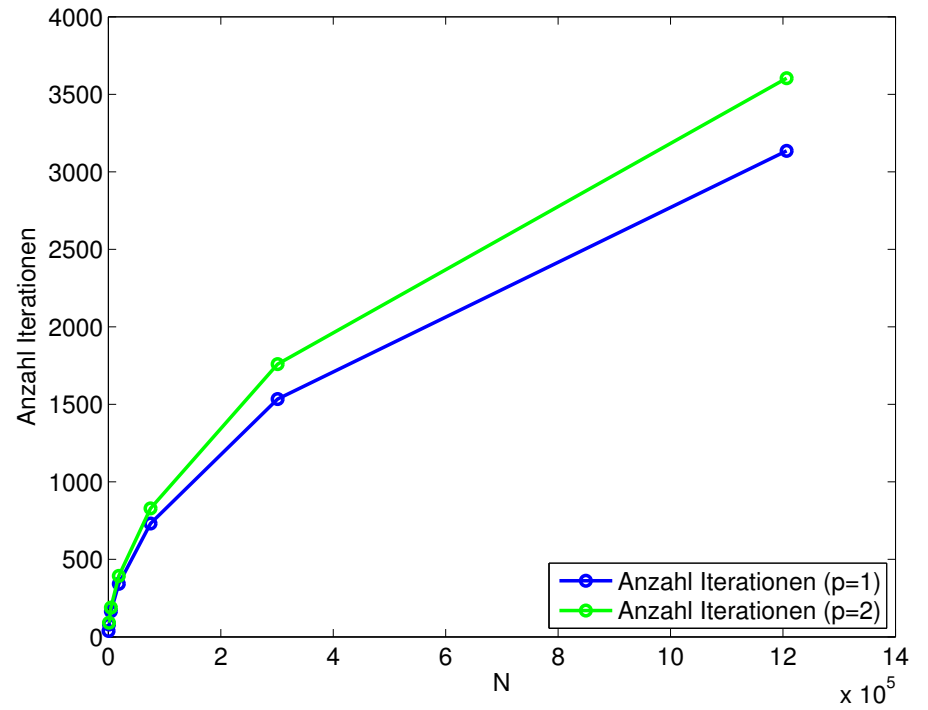
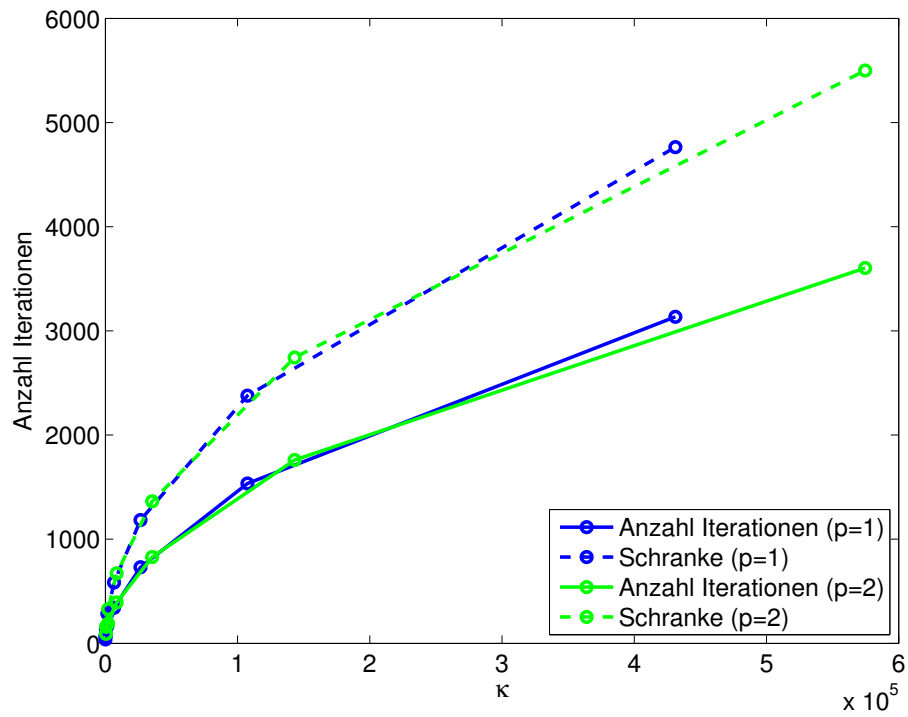
In der Praxis ist der Fehler nicht verfügbar und Abbruchkriterien basieren oft auf dem Residuum r_m . In exakter Arithmetik ist es einerlei, ob wir $r_m = b - Ax_m$ oder die rekursive Formel aus dem Algorithmus auswerten, wobei die erste Variante eine weitere Matrix-Vekto-Multiplikation erfordert. In Gleitpunktrechnung kann der Unterschied allerdings dramatisch sein.

Die obere Schranke aus Satz 4.15 ist in vielerlei Hinsicht unbefriedigend:

- Sie spiegelt nicht den Einfluß von r_0 (also b und x_0) wider.
- Sie hängt nur von den Extremaleigenwerten ab; die Verteilung der Eigenwerte innerhalb des Intervalls $[\lambda_{\min}(\mathbf{A}), \lambda_{\max}(\mathbf{A})]$ bleibt unberücksichtigt.
- In vielen für die Praxis relevanten Fällen ist die schranke aus Satz 4.15 viel zu pessimistisch.

In vielen Fällen wird bei CG ein „superlineares“ Konvergenzverhalten beobachtet. Dies ist erst vor kurzem mathematisch zufriedenstellend erklärt worden.

Beispiel 2: Wir betrachten die FE-Diskretisierung des ersten Referenzbeispiels aus Abschnitt 3.2 (Poisson-Gleichung auf Einheitsquadrat, homogene Dirichlet-Randbedingung, Quellterm $\equiv 1$) mit linearen und quadratischen Dreieckelementen.



4.2.4 Vorkonditionierung

Idee: Falls CG für gegebenes LGS zu langsam konvergiert, löse stattdessen ein **äquivalentes**, für das CG schneller konvergiert.

Ein solches ist gegeben durch

$$M^{-1}Ax = M^{-1}b, \quad M \text{ geeignete spd Matrix.}$$

Die resultierenden Modifikationen des Algorithmus erhöhen den Aufwand um

- Die Lösung eines LGS mit der Matrix M in jedem Schritt,
- Die Speicherung eines weiteren Vektors z_m , des „vorkonditionierten Residuums“.

Algorithmus 5 : Vorkonditioniertes CG-Verfahren.

Gegeben : A spd, b , M spd

- 1 Bestimme Startnäherung x_0 , berechne $r_0 \leftarrow b - Ax_0$
 - 2 Löse $Mz_0 = r_0$
 - 3 Setze $m \leftarrow 0$, $p_1 \leftarrow z_0$
 - 4 **while** nicht konvergiert **do**
 - 5 $m \leftarrow m + 1$
 - 6 $\alpha_m \leftarrow \frac{(r_{m-1}, z_{m-1})}{(p_m, p_m)_A}$
 - 7 $x_m \leftarrow x_{m-1} + \alpha_m p_m$
 - 8 $r_m \leftarrow r_{m-1} - \alpha_m A p_m$
 - 9 Löse $Mz_m = r_m$
 - 10 $\beta_{m+1} \leftarrow \frac{(r_m, z_m)}{(r_{m-1}, z_{m-1})}$
 - 11 $p_{m+1} \leftarrow z_m + \beta_{m+1} p_m$
-

Obwohl unverzichtbar in der Praxis ändert die Berücksichtigung von Vorkonditionierung wenig für die Konvergenzanalyse: es müssen lediglich folgende Ersetzungen vorgenommen werden:

$$\begin{aligned} \mathbf{A} &\leftarrow \mathbf{M}^{-1} \mathbf{A}, \\ \mathbf{b} &\leftarrow \mathbf{M}^{-1} \mathbf{b}, \\ (\mathbf{x}, \mathbf{y}) &\leftarrow (\mathbf{x}, \mathbf{y})_M := (\mathbf{M}\mathbf{x}, \mathbf{y}). \end{aligned}$$

Beachte: die $\mathbf{M}^{-1}\mathbf{A}$ -Norm bezüglich $(\cdot, \cdot)_M$ ist nichts als die \mathbf{A} -Norm bezüglich (\cdot, \cdot) . Das vorkonditionierte CG-Verfahren minimiert also weiterhin die \mathbf{A} -Norm der Fehlers.

Mit unvollständiger Cholesky-Zerlegung als Vorkonditionierer ergeben sich in Beispiel 2 für $p = 1$ folgende Iterationszahlen.

