

# Numerische Mathematik - Projektteil 2

Richard Weiss

Florian Schager

Christian Sallinger

Fabian Zehetgruber

Paul Winkler

Christian Göth

Random code snippet, damit alle checken, wie man code displayt:

```
1 print("Hello World!")
```

## 1 Dünn besetzte Matrizen

Eine häufige Problemstellung in der Numerischen Mathematik lautet lineare Gleichungssysteme mit großen, dünn besetzten Matrizen zu lösen. Dabei kommen meist iterative Verfahren zum Einsatz, die in diesem Projekt effizient implementiert werden.

a)

Generieren Sie eine symmetrisch positiv definite Zufallsmatrix  $A \in \mathbb{R}^{n \times n}$ , wobei pro Zeile eine fixe Anzahl an Einträgen ungleich Null sind. Testen Sie für eine zufällige rechte Seite  $b \in \mathbb{R}^n$  bis zu welcher Größe  $n$  das lineare Gleichungssystem

$$Ax = b$$

mit einem direkten Löser (`numpy.linalg.solve`) in akzeptabler Zeit gelöst werden kann. Welchen Aufwand erwarten Sie abhängig von  $n$ ? Plotten Sie die Rechenzeit in Abhängigkeit von der Problemgröße.

```
1
2 def zufallsmatrix(n, nonzeros):
3     A = np.concatenate((np.zeros((n,nonzeros)), np.random.rand(n, n-nonzeros)), axis = 1)
4     for i in range(n):
5         np.random.shuffle(A[i])
6     return A + A.T + np.diag(np.random.rand(n)*10)
7
8 A_base = zufallsmatrix(5000,100)
```

Wie in Abbildung 1 ersichtlich verhält sich die Rechenzeit kubisch in Relation zur Problemgröße. Zum Testen wurde eine Zufallsmatrix mit 100 Nicht-Null-Einträgen pro Zeile (nicht exakt, da die Symmetriesierung den Wert pro Zeile verzerrt) und einer Gesamtgröße von 5000 erstellt. Der direkte Löser wurde schließlich auf die  $(200k \times 200k)$ -dimensionalen Ausschnitte der oberen rechten Ecke angewandt ( $2 \leq k \leq 25$ ). Wie man sieht erreichen wir damit schon langsam die Grenze des akzeptabel Berechenbaren, für die volle  $5000 \times 5000$ -Matrix braucht der Algorithmus schon über 10 Sekunden.

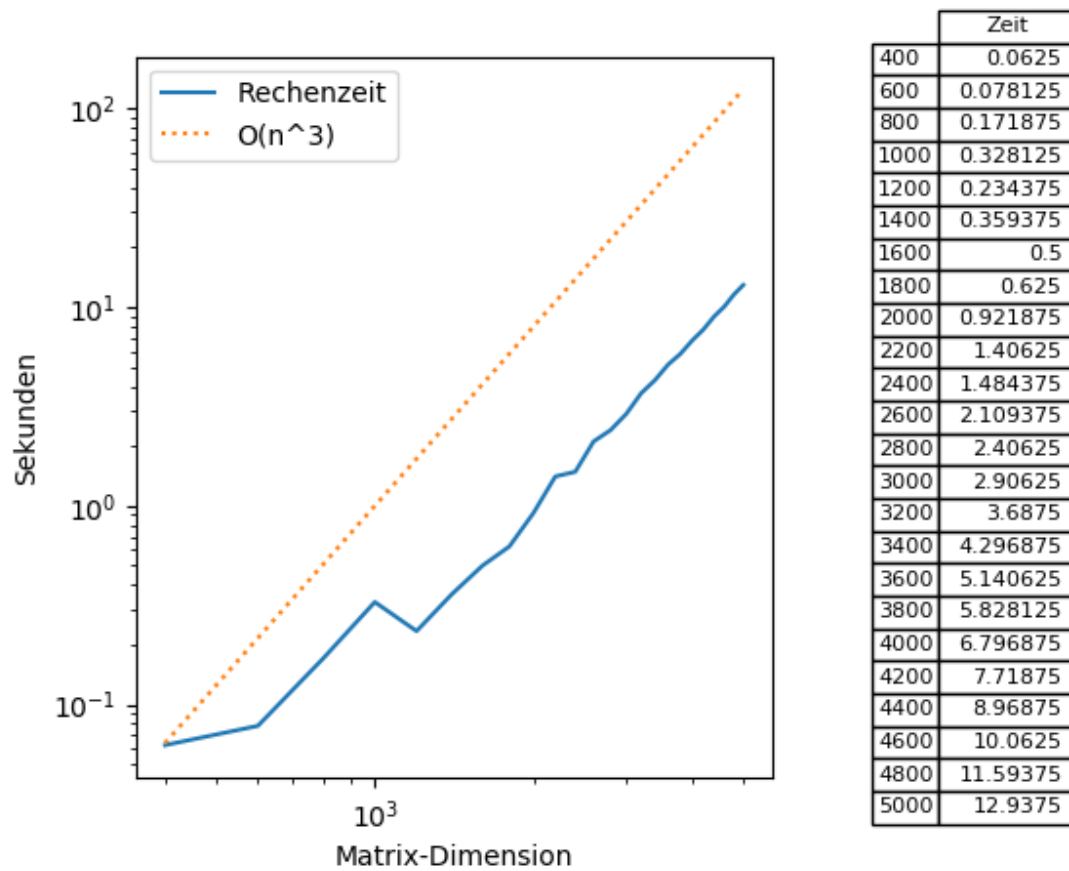


Abbildung 1: Rechenzeit abhängig von der Problemgröße

b)

Implementierung des CG-Verfahrens:

```

1 def cg(A,b,x0,tol):
2     xt = x0
3     r0 = b - np.dot(A,xt)
4     d = r0
5     while(np.linalg.norm(r0) > tol):
6         prod = np.dot(np.transpose(r0),r0)
7         prod2 = np.dot(A,d)
8         alpha = prod/np.dot(np.transpose(d),prod2)
9         xt = xt + alpha*d
10        r0 = r0 - alpha*prod2
11        beta = np.dot(np.transpose(r0),r0)/prod
12        d = r0 + beta*d
13    return xt

```

Beweis der Äquivalenz obigen Algorithmus zu Algorithmus 8.10 im Numerik-Skript:

Wir führen den Beweis mittels Induktion:

Dabei bezeichnen wir mit \* die Variablen aus unserem Algorithmus und ohne \* die Variablen des Algorithmus aus dem Skript.

Induktionsanfang:

$$\begin{aligned}
 A &= A^*, b = b^*, x_0 = x_0^* \\
 r_0 &= b - Ax_0 = b^* - A^*x_0^* = r_0^* \\
 d_0 &= r_0 = r_0^* = d_0^* \\
 \alpha_0 &= \frac{r_0^T d_0}{d_0^T A d_0} = \frac{r_0^{*T} d_0^*}{d_0^{*T} A d_0^*} = \frac{r_0^{*T} r_0^*}{d_0^{*T} A d_0^*} = \alpha_0^* \\
 x_1 &= x_0 + \alpha_0 d_0 = x_0^* + \alpha_0^* d_0^* = x_1^* \\
 r_1 &= b - Ax_1 = b^* - A^*x_1^* = b^* - A^*x_0^* - \alpha_0^* A d_0^* = r_0^* - \alpha_0^* A^* d_0^* = r_1^* \\
 \beta_0 &= -\frac{r_1^T A d_0}{d_0^T A d_0} = \frac{-r_1^{*T} A d_0^*}{d_0^{*T} A d_0^*} = \frac{-r_1^{*T} A r_0^*}{r_0^{*T} A r_0^*}
 \end{aligned}$$

Unter Ausnutzung der Orthogonalität der Residuen erhalten wir:  $r_1^{*T} r_0^* = 0$  und somit können wir den Bruch folgendermaßen erweitern:

$$\frac{-r_1^{*T} A r_0^*}{r_0^{*T} A r_0^*} = \frac{r_1^{*T} r_0^* - \alpha_0^* r_1^{*T} A r_0^*}{\alpha_0^* r_0^{*T} A r_0^*}$$

Nun berechnen wir

$$r_1^{*T} r_1^* = r_1^{*T} (r_0^* - \alpha_0^* A d_0^*) = r_1^{*T} r_0^* - \alpha_0^* r_1^{*T} A r_0^*$$

und setzen  $\alpha_0^* = \frac{r_0^{*T} r_0^*}{d_0^{*T} A d_0^*}$  ein:

$$\begin{aligned}
 \frac{r_1^{*T} r_0^* - \alpha_0^* r_1^{*T} A r_0^*}{\alpha_0^* r_0^{*T} A r_0^*} &= \frac{r_1^{*T} r_1^*}{\frac{r_0^{*T} r_0^*}{r_0^{*T} A r_0^*} r_0^{*T} A r_0^*} = \frac{r_1^{*T} r_1^*}{r_0^{*T} r_0^*} = \beta_0^* \\
 d_1 &= r_1 + \beta_0 d_0 = r_1^* + \beta_0^* d_0^* = d_1^*
 \end{aligned}$$

Damit haben wir die Gleichheit der Variablen nach dem ersten Schleifendurchlauf gezeigt.

Sei nun nach  $n$  Schleifendurchläufen die Gleichheit aller vorhergehenden Variablen vorausgesetzt:

$$\alpha_{n-1} = \alpha_{n-1}^*, x_n = x_n^*, r_n = r_n^*, \beta_{n-1} = \beta_{n-1}^*, d_n = d_n^*$$

$$\alpha_n = \frac{r_n^T d_n}{d_n^T A d_n} = \frac{r_n^{*T} d_n^*}{d_n^{*T} A^* d_n^*} = \frac{r_n^{*T} (r_n^* + \beta_{n-1}^* d_{n-1}^*)}{d_n^{*T} A^* d_n^*}$$

Jetzt nutzen wir die Eigenschaft:  $\forall 0 \leq j < m : r_m^T d_j = 0$  und erhalten:

$$\frac{r_n^{*T} (r_n^* + \beta_{n-1}^* d_{n-1}^*)}{d_n^{*T} A^* d_n^*} = \frac{r_n^{*T} r_n^*}{d_n^{*T} A^* d_n^*} = \alpha_n^*$$

$$x_{n+1} = x_n + \alpha_n d_n = x_n^* + \alpha_n^* d_n^* = x_{n+1}^*$$

$$r_{n+1} = b - A x_{n+1} = b - A(x_n + \alpha_n d_n) = b - A x_n - \alpha_n A d_n =$$

$$= r_n - \alpha_n A d_n = r_n^* - \alpha_n^* A^* d_n^* = r_{n+1}^*$$

$$\beta_n = -\frac{r_{n+1}^T A d_n}{d_n^T A d_n} = \frac{r_{n+1}^{*T} r_n^* - \alpha_n^* r_{n+1}^{*T} A^* d_n^*}{d_n^{*T} A^* d_n^*} = \frac{r_{n+1}^{*T} r_{n+1}^*}{r_n^{*T} r_n^*} = \beta_n^*$$

$$d_{n+1} = r_{n+1} + \beta_n d_n = r_{n+1}^* + \beta_n^* d_n^* = d_{n+1}^*$$

Und der Beweis ist vollständig.

Rechenschritte pro Schleifendurchlauf:

Am aufwändigsten ist natürlich die Matrix-Vektor-Multiplikation, wovon wir in unserem Algorithmus im Gegensatz zu jenem aus dem Skript nur eine pro Schleifendurchlauf benötigen. Zusätzlich dazu brauchen wir 3 Vektor-Vektor-Multiplikationen und 4 Vektor-Skalar-Operationen.

Damit erhalten wir insgesamt  $n^2 + 7n$  Flops pro Durchlauf.

Im Vergleich dazu verwendet der Algorithmus aus dem Numerik-Skript pro Iteration zwei Matrix-Vektor-Multiplikationen und ist daher aus Effizienzgründen unterlegen.

Die Theorie besagt, dass das CG-Verfahren spätestens nach  $n$  Durchläufen die exakte Lösung liefert und für die Iterierten folgende Fehlerabschätzung gilt:

$$\|x^{(t)} - A^{-1}b\|_A \leq 2 \left( \frac{1 - 1/\sqrt{\kappa}}{1 + 1/\sqrt{\kappa}} \right)^t \|x^{(0)} - A^{-1}b\|_A, \quad t \in \mathbb{N},$$

mit der spektralen Konditionszahl  $\kappa = \text{cond}_2(A)$ .

Also sollte das Verfahren exponentiell konvergieren ( $\mathcal{O}(AB^t)$ ) mit Konstanten

$$A = 2\|x^{(0)} - A^{-1}b\|_A, \quad B = \frac{1-1/\sqrt{\kappa}}{1+1/\sqrt{\kappa}}$$

Wie in untenstehender Grafik ersichtlich scheint bei unseren Zufallsmatrizen der Wert von  $B$  sich in etwa bei 0.92 einzupendeln und unsere vorgegebene Toleranz von  $10^{-8}$  wird bereits nach etwa 350 Iterationen erreicht, noch weit vor dem theoretisch (bis auf Rechenfehler) garantierten exakten Resultat nach  $n = 5000$  Durchläufen. Ansonsten verhält sich der Fehler wie wir ihn erwarten würden.

Unsere Testwerte:

```
1 n = 5000
2 A = np.random.rand(n,n)
3 A = np.dot(A,np.transpose(A))
4 A += 10*np.diag(abs(np.random.random(n)))
5 b = np.random.rand(n)
6 tol = 10**(-8)
7 x0 = np.random.rand(n)*10
```

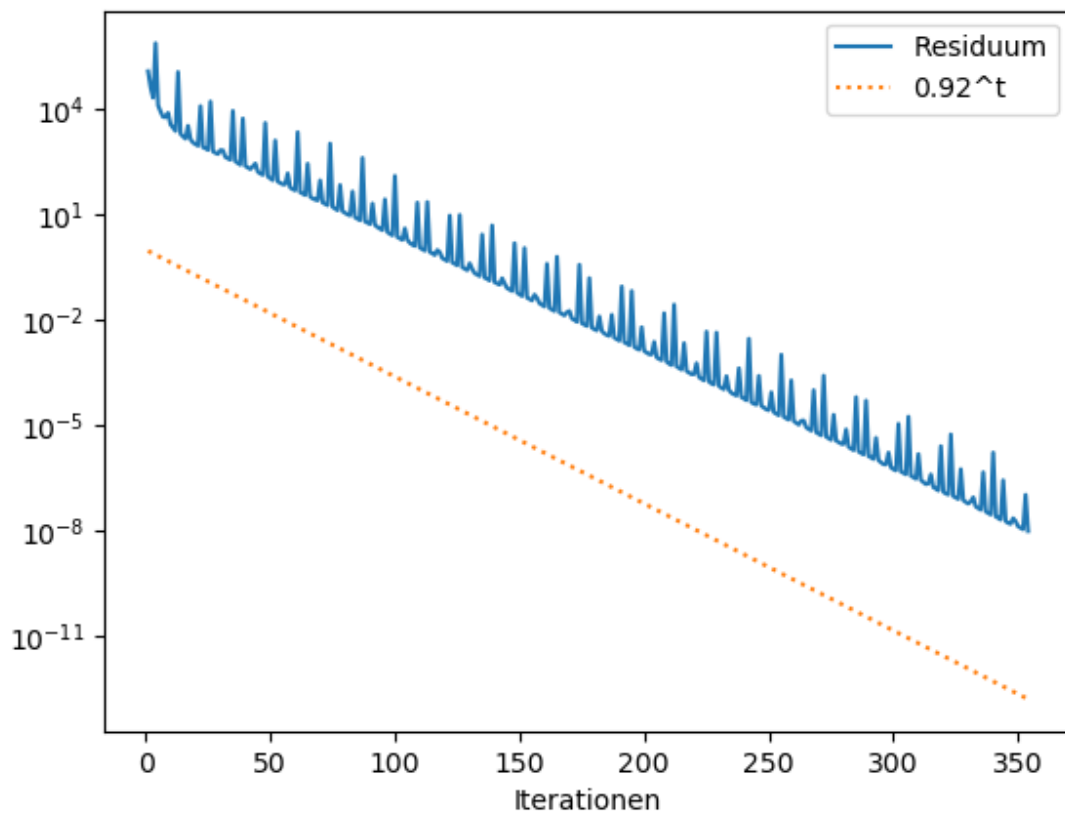


Abbildung 2: Residuum abhängig von der Anzahl an Iterationen

c)

Dünn besetzte Matrizen erlauben effizientere Implementierungen als voll besetzte, indem beim Speichern und Rechnen nur Einträge die ungleich Null sind berücksichtigt werden. Eine Möglichkeit einer solchen Implementierung ist das sogenannte *compressed sparse row* Format. Anstelle aller Einträge  $A_{ij}$ ,  $i, j = 1, \dots, n$  einer Matrix  $A \in \mathbb{R}^{n \times n}$  werden ein Vektor  $v \in \mathbb{R}^{n \times n}$  aller Einträge ungleich Null, ein Vektor  $J \in \mathbb{N}_0^m$  von Spaltenindizes und ein Vektor  $I \in \mathbb{N}_0^{n+1}$  von Zeigern gespeichert. Die  $i$ -te Zeile von  $A$  ist dann gegeben durch

$$A_{ij} = \begin{cases} v_{k(j)}, & \text{falls } j \in \{J_{I_i}, J_{I_i} + 1, \dots, J_{I_{i+1}} - 1\} \\ 0, & \text{sonst} \end{cases}$$

Implementierung in Python:

```

1 class Sparse:
2
3     def __init__(self, b, v, J = np.zeros(0), I = np.zeros(0)):
4         if b:
5             self.v = np.array(v)
6             self.J = np.array(J)
7             self.I = np.array(I)
8             self.n = len(self.I) - 1
9         else:
10            self.v, self.J, self.I = self.fromdense(v)
11            self.n = len(self.I) - 1
12
13     def __matmult__(self, b):
14         d = np.zeros(self.n)
15         for i in range(self.n):
16             x = np.array(self.J[self.I[i]:self.I[i+1]]).astype(int)
17             d[i] = self.v[self.I[i]:self.I[i+1]] @ b[x]
18         return d
19
20
21     def todense(self):
22         A = np.zeros([self.n, self.n])
23         for i in range(self.n):
24             for j in range(self.I[i], self.I[i+1]):
25                 A[i][self.J[j]] = self.v[j]
26         return A
27
28     def fromdense(self, A):
29         v, J = np.zeros(0), np.zeros(0)
30         I = np.array([0])
31         c = 0
32         for i in range(np.shape(A)[0]):
33             for j in range((np.shape(A))[0]):
34                 if A[i][j] != 0:
35                     v = np.append(v, A[i][j])
36                     J = np.append(J, j)
37                     c += 1
38             I = np.append(I, c)
39         return v, J, I

```

Da wir bei diesem Projekt uns nur auf das cg-Verfahren konzentrieren, haben wir nur die Matrix-Vektor Multiplikation für die Klasse Sparse implementiert.

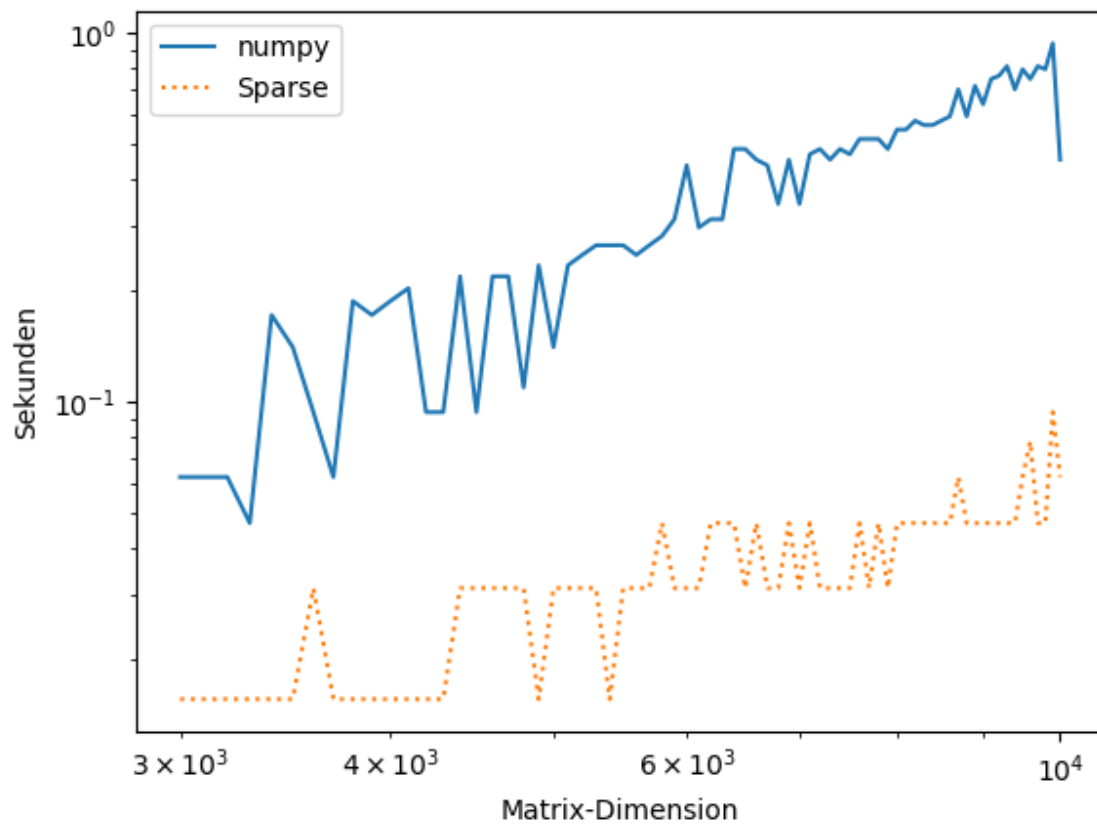


Abbildung 3: Vergleich Numpy Matrixmultiplikation vs. Sparse Matrixmultiplikation

Hier sehen wir, dass erst ab einer gewissen Größe die Sparse-Multiplikation effizienter gegenüber der Numpy-Multiplikation ist.

d)

Kombinieren Sie Ihre CG-Implementierung mit Ihrer **sparse** Klasse und testen Sie die Effizienz: Implementierung:

```
1 def Scg(A,b,x0,tol):
2     xt = x0
3     r0 = b - A.__matmult__(xt)
4     d = r0
5     while(np.linalg.norm(r0) > tol):
6         prod = np.dot(np.transpose(r0),r0)
7         prod2 = A.__matmult__(d)
8         alpha = prod/np.dot(np.transpose(d),prod2)
9         xt = xt + alpha*d
10        r0 = r0 - alpha*prod2
11        beta = np.dot(np.transpose(r0),r0)/prod
12        d = r0 + beta*d
13    return xt
```

Bei dieser CG-Implementierung verwenden wir anstelle der Numpy-Matrix-Vektor-Multiplikation die Implementierung von der Sparse-Klasse. Zu beachten ist, dass die Matrix A ein Objekt der Klasse Sparse sein muss, damit die Funktion durchgeführt werden kann. Ansonsten ist die Implementierung ident zum vorherigen cg-Verfahren.



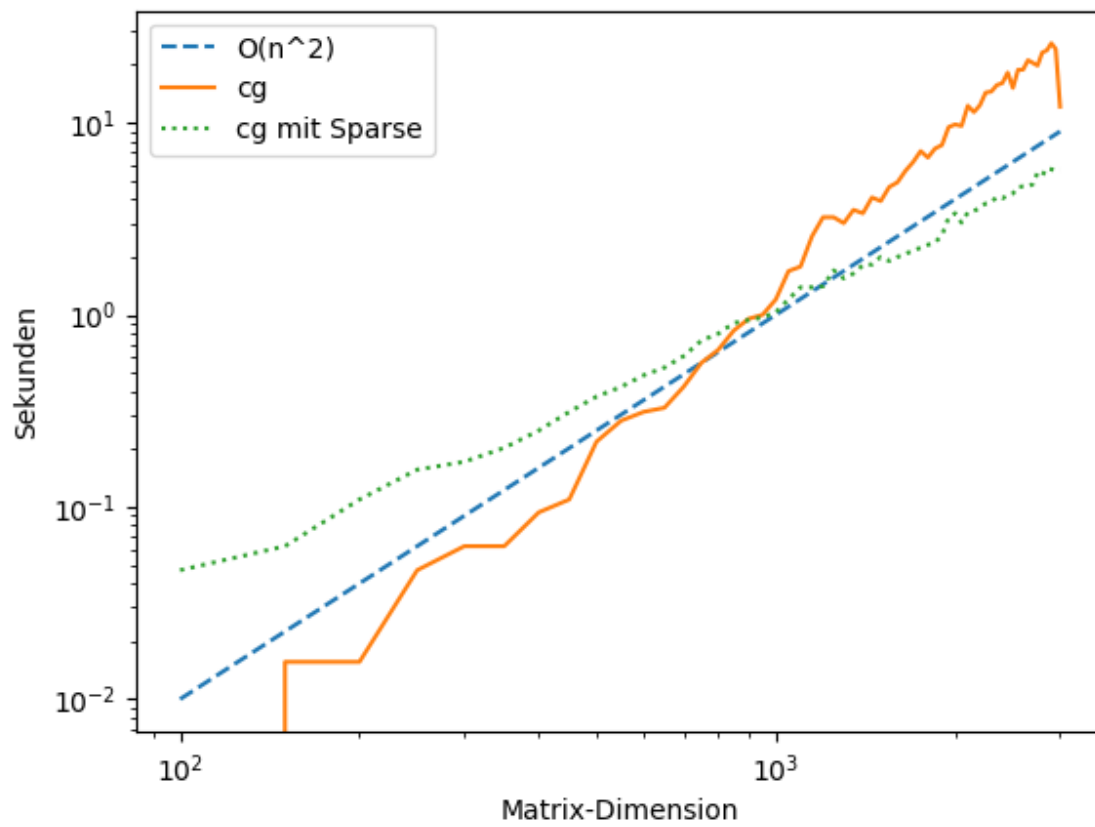


Abbildung 4: cg-Verfahren vs. cg-Verfahren mit Klasse Sparse

Ab einer Matrixgröße von ca. 1000 ist das Scg-Verfahren effizienter als die übliche Implementierung. Man sieht außerdem, dass beide Algorithmen eine Laufzeit von etwa  $\mathcal{O}(n^2)$

e)

Die Konvergenzgeschwindigkeit des CG-Verfahrens ist durch die spektrale Konditionszahl  $\text{cond}(A)$  der Matrix  $A$  bestimmt. Um die Konvergenzgeschwindigkeit zu erhöhen löst man das vorkonditionierte System

$$D^{-1}AD^{-1T} = D^{-1}b$$

und gewinnt die Lösung  $x$  dann durch  $x = D^{-1T}y$ . Die Matrix  $D$  wird dabei so gewählt, dass

- für beliebige Vektoren  $z \in \mathbb{R}^n$  der Vektor  $D^{-1T}D^{-1}z$  einfach zu berechnen ist und
- $\text{cond}(D^{-1}AD^{-1T}) < \text{cond}(A)$ .

Implementierung des vorkonditionierten CG-Verfahrens:

```

1  def vcg(A,b,x0,P,tol):
2      r0 = b - A.__matmult__(x0)
3      P_inv = np.linalg.inv(P)
4      z0 = P_inv*r0
5      d = z0
6      while(np.linalg.norm(r0) > tol):
7          prod = np.dot(np.transpose(z0),r0)
8          prod2 = A.__matmult__(d)
9          alpha = np.dot(np.transpose(r0),z0)/np.dot(np.transpose(d),prod2)
10         x0 = x0 + alpha*d
11         r0 = r0 - alpha*prod2
12         z0 = P_inv*r0
13         beta = np.dot(np.transpose(z0),r0)/prod
14         d = z0 + beta*d
15     return x0

```

Wieder schlechte Nachrichten: Der Scheiß konvergiert!

## 2 Eigenschwingungen

### 2.1 Problembeschreibung

Das Projekt beschäftigt sich mit den Eigenschwingungen einer fest eingespannten Saite. Sei dazu  $u(t, x)$  die vertikale Auslenkung der Saite an der Position  $x \in [0, 1]$  zur Zeit  $t$ .  $u$  wird näherungsweise durch die sogenannte Wellengleichung

$$\frac{\partial^2 u}{\partial x^2}(t, x) = \frac{1}{c^2} \frac{\partial^2 u}{\partial t^2}(t, x) \quad (1)$$

für alle  $x \in (0, 1)$  und  $t \in \mathbb{R}$  beschrieben, wobei  $c$  die Ausbreitungsgeschwindigkeit der Welle ist. Wenn die Saite an beiden Enden fest eingespannt ist, so gelten die Randbedingungen

$$u(t, 0) = u(t, 1) = 0$$

für alle  $t \in \mathbb{R}$ .

Zur Berechnung der Eigenschwingungen suchen wir nach Lösungen  $u$ , die in der Zeit harmonisch schwingen. Solche erfüllen folgenden Ansatz

$$u(x, t) = \Re(v(x)e^{-i\omega t})$$

mit einer festen, aber unbekannten Kreisfrequenz  $\omega > 0$  und einer Funktion  $v$ , welche nur noch vom Ort  $x$  abhängt. Durch Einsetzen erhalten wir für  $v$  die sogenannte Helmholtz-Gleichung

$$-v''(x) = \kappa^2 v(x), \quad x \in (0, 1), \quad (2)$$

mit der unbekannten Wellenzahl  $\kappa := \frac{\omega}{c}$  und den Randbedingungen

$$v(0) = v(1) = 0. \quad (3)$$

## 2.2 Analytische Lösung

$$v_\kappa(x) = C_1 \cos(\kappa x) + C_2 \sin(\kappa x), \quad x \in [0, 1], \quad (4)$$

mit beliebigen Konstanten  $C_1, C_2$  löst die Helmholtz-Gleichung (2). Das erkennt man durch stumpfes Einsetzen.

$$\begin{aligned} -v_\kappa''(x) &= -\frac{\partial^2}{\partial x^2}(C_1 \cos(\kappa x) + C_2 \sin(\kappa x)) = -\frac{\partial}{\partial x}(-C_1 \kappa \sin(\kappa x) + C_2 \kappa \cos(\kappa x)) \\ &= -(-C_1 \kappa^2 \cos(\kappa x) - C_2 \kappa^2 \sin(\kappa x)) = \kappa(C_1 \cos(\kappa x) + C_2 \sin(\kappa x)) = \kappa^2 v_\kappa(x) \end{aligned}$$

Wir fragen uns, für welche  $\kappa > 0$ , Konstanten  $C_1$  und  $C_2$  existieren, sodass  $v_\kappa$  auch die Randbedingungen (3) erfüllt.

$$0 \stackrel{!}{=} \begin{cases} v_\kappa(0) = C_1 \cos 0 + C_2 \sin 0 = C_1 \\ v_\kappa(1) = C_1 \cos \kappa + C_2 \sin \kappa = C_2 \sin \kappa \end{cases}$$

Nachdem  $\cos 0 = 1$  und  $\sin 0 = 0$ , erhält man, aus der oberen Gleichung,  $C_1 = 0$ . Mit der unteren Gleichung folgt aber auch  $C_2 \sin \kappa = 0$ . Wenn nun auch  $C_2 = 0$ , dann erhielte man die triviale Lösung  $v_\kappa = 0$ . Für eine realistischere Modellierung, d.h.  $v_\kappa \neq 0$ , müsste  $\sin \kappa = 0$ , also  $\kappa \in \pi\mathbb{Z}$ .

Das sind die gesuchten  $\kappa > 0$ . Sei nun eines dieser  $\kappa$  fest. Offensichtlich ist  $C_1 = 0$  eindeutig,  $C_2 \in \mathbb{R}$  jedoch beliebig.

## 2.3 Numerische Approximation

Häufig lassen sich solche Probleme nicht analytisch lösen, sodass auf numerische Verfahren zurückgegriffen wird, welche möglichst gute Näherungen an die exakten Lösungen berechnen sollen. Als einfachstes Mittel dienen sogenannte Differenzenverfahren. Sei dazu  $x_j := jh$ ,  $j = 0, \dots, n$  eine Zerlegung des Intervalls  $[0, 1]$  mit äquidistanter Schrittweite  $h = 1/n$ . Die zweite Ableitung in (2) wird approximiert durch den Differenzenquotienten

$$v''(x_j) \approx D_h v(x_j) := \frac{1}{h^2}(v(x_{j-1}) - 2v(x_j) + v(x_{j+1})), \quad j = 1, \dots, n-1. \quad (5)$$

Es sei angemerkt, dass (5) tatsächlich einen Differenzenquotienten beschreibt. Um das einzusehen, verwenden wir den links- und rechts-seitigen Differenzenquotient erster Ordnung, sowie  $x_{j-1} = x_j - h$ ,  $x_{j+1} = x_j + h$ . Wir erhalten  $\forall j = 1, \dots, n-1$ :

$$\begin{aligned} v''(x_j) &= \lim_{h \rightarrow 0} \frac{1}{h} (v'(x_j + h) - v'(x_j)) \\ &= \lim_{h \rightarrow 0} \frac{1}{h} \left( \frac{1}{h} (v(x_j + h) - v(x_j)) - \frac{1}{h} (v(x_j) - v(x_j - h)) \right) \\ &= \lim_{h \rightarrow 0} \frac{1}{h^2} (v(x_j + h) - 2v(x_j) + v(x_j - h)) \\ &= \lim_{h \rightarrow 0} D_h v(x_j) \end{aligned}$$

### 2.3.1 Approximationsfehler

Für hinreichend glatte Funktionen  $v$  mit einer geeigneten Konstanten  $C > 0$  wird der Approximationsfehler quadratisch in  $h$  klein, d.h. dass

$$|v''(x_j) - D_h v(x_j)| \leq Ch^2. \quad (6)$$

Nachdem  $v$  hinreichend glatt ist, gilt nach dem Satz von Taylor, dass  $\forall j = 1, \dots, n-1$ :

$$\begin{aligned} v(x_j + h) &= \sum_{\ell=0}^{n+2} \frac{h^\ell}{\ell!} v^{(\ell)}(x_j) + \mathcal{O}(h^{n+3}), \\ v(x_j - h) &= \sum_{\ell=0}^{n+2} \frac{(-h)^\ell}{\ell!} v^{(\ell)}(x_j) + \mathcal{O}(h^{n+3}). \end{aligned}$$

Man beachte, dass sich die ungeraden Summanden, der oberen Taylor-Polynome, gegenseitig aufheben. Damit erhalten wir für den Differenzenquotient  $D_h v(x_j)$ ,  $j = 1, \dots, n-1$  eine asymptotische Entwicklung.

$$\begin{aligned} D_h v(x_j) &= \frac{1}{h^2} (v(x_j - h) + v(x_j + h) - 2v(x_j)) \\ &= \frac{1}{h^2} \left( 2v(x_j) + h^2 v''(x_j) + \sum_{\substack{\ell=4 \\ \ell \in 2\mathbb{N}}}^{n+2} \frac{h^\ell}{\ell!} v^{(\ell)}(x_j) (1 + (-1)^\ell) - 2v(x_j) \right) + \mathcal{O}(h^{n+3}) \\ &= v''(x_j) + 2 \sum_{\ell=1}^{\lfloor \frac{n}{2} \rfloor} \frac{h^{2\ell}}{(2\ell+2)!} v^{(2\ell)}(x_j) + \mathcal{O}(h^{n+1}) \end{aligned}$$

Daraus folgt unmittelbar die quadratische Konvergenz (6), d.h.  $\forall j = 1, \dots, n-1$ :

$$D_h v(x_j) - v''(x_j) = \mathcal{O}(h^2), \quad h \rightarrow 0.$$

### 2.3.2 Eigenwertproblem

Wir wollen nun den Differenzenquotienten  $D_h v(x_j)$  verwenden, um ein Eigenwertproblem der Form  $A\mathbf{v} = \lambda\mathbf{v}$  mit einer Matrix  $A \in \mathbb{R}^{(n-1) \times (n-1)}$  zu dem Eigenvektor  $\mathbf{v} := (v(x_1), \dots, v(x_{n-1}))^T$  und dem Eigenwert

$\lambda := \kappa^2$  herzuleiten. Das soll die Helmholtz-Gleichung (2), mit Tupeln und Eigenwerten für die Funktionen  $v_\kappa$  bzw. Vorfaktoren  $\kappa$ , approximieren.

Es wird eine Matrix  $-A_n$ ,  $n \geq 2$  gesucht, die den Differenzenquotienten  $D_h v(x_j)$  auf den oberen Vektor  $\mathbf{v}^{(n)}$  komponentenweise anwendet. Wir rufen in Erinnerung, dass  $h = 1/n$  und definieren die naheliegende Matrix

$$-A_n := \frac{1}{h^2} \begin{pmatrix} -2 & 1 & & & \\ 1 & \ddots & \ddots & & \\ & \ddots & \ddots & 1 & \\ & & & 1 & -2 \end{pmatrix} \in \mathbb{R}^{(n-1) \times (n-1)}.$$

Wenn nun die Randbedingungen (3) gelten sollen, d.h.  $v(x_0), v(x_n) = 0$ , leistet diese Matrix  $A_n$  tatsächlich das Gewünschte. Sie approximiert die linke Seite der Helmholtz-Gleichung (2).

$$A_n \mathbf{v}^{(n)} = -\frac{1}{h^2} \begin{pmatrix} v(x_0) - 2v(x_1) + v(x_2) \\ v(x_1) - 2v(x_2) + v(x_3) \\ \vdots \\ v(x_{n-3}) - 2v(x_{n-2}) + v(x_{n-1}) \\ v(x_{n-2}) - 2v(x_{n-1}) + v(x_n) \end{pmatrix} = - \begin{pmatrix} D_h v(x_1) \\ \vdots \\ D_h v(x_{n-1}) \end{pmatrix} \xrightarrow{n \rightarrow \infty} -v''$$

Die Matrix  $A_n$  wurde in der Funktion `my_numpy_matrix` implementiert.

```
1 def my_numpy_matrix(n):
2
3     assert n >= 2
4
5     h = 1/n
6
7     a = -2 * np.ones(n-1)
8     b = np.ones(n-2)
9
10    A = np.diag(b, -1) + np.diag(a, 0) + np.diag(b, 1)
11    A = -A/h**2
12
13    return A
```

Das Eigenwertproblem wurde mit `np.linalg.eig`, für beliebige  $n \geq 2$ , gelöst. `np.linalg.eig` retourniert die Eigenwerte und Eigenvektoren nicht zwangsläufig, der gröÙe der Eigenwerte nach, sortiert. Nachdem der Zusammenhang zwischen Eigenwert und Eigenvektor nicht verloren gehen soll, erfordert dies einigen logistischen Aufwand. Das ist aber nicht wesentlich und wird daher nicht weiter erklärt. Wir vergleichen lieber die Eigenwerte und Eigenvektoren mit den analytischen Ergebnissen.

In der folgenden Tabelle, sind die 3 Eigenwerte (sollten diese bereits existieren), der Matrizen  $A_2, \dots, A_{10}$ , aufgelistet. Diese Tabelle ist analog zu (7) zu verstehen.

	2	3	4	5	6	7	8	9	10
1	8.0	9.0	9.372583	9.54915	9.646171	9.705051	9.743420	9.769795	9.788697
2	NaN	27.0	32.000000	34.54915	36.000000	36.897999	37.490332	37.900800	38.196601
3	NaN	NaN	54.627417	65.45085	72.000000	76.192948	79.016521	81.000000	82.442950

Die Matrix  $A_n$  besitzt also scheinbar  $n - 1$  paarweise verschiedene Eigenwerte  $\lambda_{1,n} < \dots < \lambda_{n-1,n}$ , welche jeweils gegen  $\lambda_i := (i\pi)^2$ ,  $i \in \mathbb{N}$  konvergieren.

$$\begin{array}{ccccccc}
\lambda_{1,2} & \rightarrow & \lambda_{1,3} & \rightarrow & \dots & \rightarrow & \lambda_{1,n} \xrightarrow{n \rightarrow \infty} \lambda_1 = (1\pi)^2 \\
& & \lambda_{2,3} & \rightarrow & \dots & \rightarrow & \lambda_{2,n} \xrightarrow{n \rightarrow \infty} \lambda_2 = (2\pi)^2 \\
& & & & \ddots & & \vdots \\
& & & & & & \lambda_{i,n} \xrightarrow{n \rightarrow \infty} \lambda_i = (i\pi)^2
\end{array} \tag{7}$$

Nachdem  $\{\sqrt{\lambda_i} : i \in \mathbb{N}_0\} = \pi\mathbb{N}_0 \subsetneq \pi\mathbb{Z}$ , könnte man sich nun fragen, wo die andere Hälfte der analytischen Ergebnisse steckt. Die Erklärung ist ganz unspektakulär  $\lambda := (\pm\kappa)^2$ .

Wir bezeichnen mit  $\epsilon_i(n) := |\lambda_i - \lambda_{i,n}|$ ,  $i = 1, \dots, n - 1$  den absoluten Konvergenz-Fehler des  $i$ -ten Eigenwertes. In der folgenden Abbildung 5 wurde  $\epsilon_i$  mit der Vergleich-Geraden  $\text{id}^2$ , doppelt logarithmisch, geplottet.

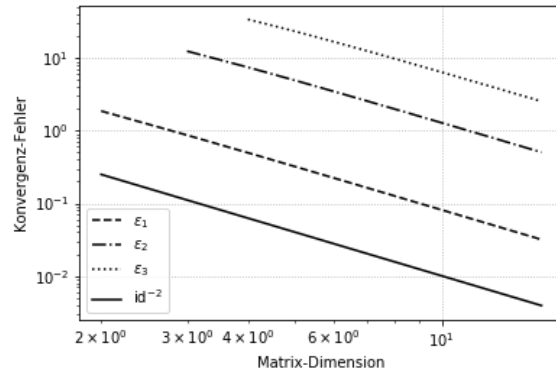
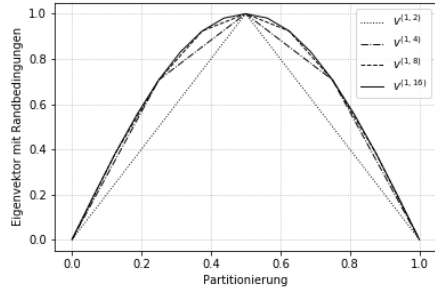


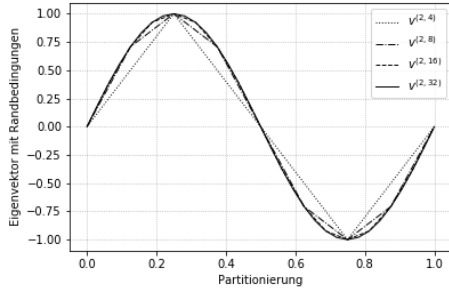
Abbildung 5: Konvergenz-Fehler der Eigenwerte von  $A_n$

Allem Anschein nach, verschwindet  $\epsilon_i$  quadratisch. Das korreliert mit dem Ergebnis (6). Man beachte, dass der  $i$ -te Eigenwert erst ab einer Matrix  $A_n$ ,  $n > i$  existiert. Daher fangen die Plots von  $\epsilon_i$  desto später an, je größer  $i$  ist. Für größeres  $i$  ist auch der initiale Fehler größer. Obwohl dieser ebenfalls quadratisch konvergiert, werden mehr Rechenoperationen für ein genaues Ergebnis benötigt.

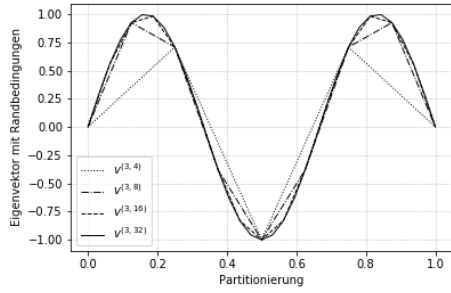
Seien  $\mathbf{v}^{(1,n)}, \dots, \mathbf{v}^{(n-1,n)}$  die Eigenvektoren (modulo Konstante), zu den Eigenwerten  $\lambda_{1,n} < \dots < \lambda_{n-1,n}$ , der Matrix  $A_n$ . Diese sollten nun gegen die Funktionen  $v_{\kappa_i}$ ,  $\kappa_i = \sqrt{\lambda_i}$ , vielleicht sogar quadratisch, konvergieren. Folgende Abbildungen sollen dies veranschaulichen.



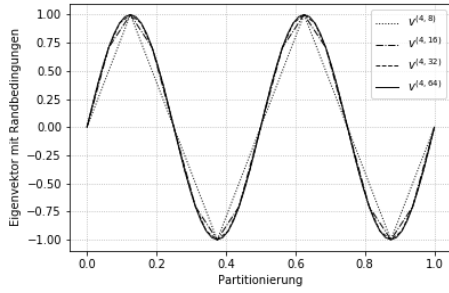
(a) Eigenvektoren  $v^{(1,n)}$ ,  $n = 2, 4, 8, 32$



(b) Eigenvektoren  $v^{(2,n)}$ ,  $n = 4, 8, 16, 64$



(c) Eigenvektoren  $v^{(3,n)}$ ,  $n = 8, 16, 32, 64$



(d) Eigenvektoren  $v^{(4,n)}$ ,  $n = 8, 16, 32, 64$

Abbildung 6: Eigenvektoren  $v^{(i,n)}$ ,  $i = 1, \dots, 4$  der Matrizen  $A_n$

Erstaunlicherweise, gibt es scheinbar keinen Konvergenz-Fehler, da die Eigenvektoren direkt an den Grenzfunktionen liegen. Mit anderen Worten,  $\forall n \in \mathbb{N}, \forall i, j = 1, \dots, n-1$ :

$$\mathbf{v}_j^{(i,n)} = v_{\kappa_i}(x_j).$$

Anschaulich, erhält man ein zu (7) analoges Schema (8).

$$\begin{array}{ccccccc}
 \mathbf{v}^{(1,2)} & \rightarrow & \mathbf{v}^{(1,3)} & \rightarrow & \dots & \rightarrow & \mathbf{v}^{(1,n)} & \xrightarrow{n \rightarrow \infty} & v_{\kappa_1} \\
 & & \mathbf{v}^{(2,3)} & \rightarrow & \dots & \rightarrow & \mathbf{v}^{(2,n)} & \xrightarrow{n \rightarrow \infty} & v_{\kappa_2} \\
 & & & & \ddots & & \vdots & & \vdots \\
 & & & & & & \mathbf{v}^{(i,n)} & \xrightarrow{n \rightarrow \infty} & v_{\kappa_i}
 \end{array} \tag{8}$$

## 2.4 Verallgemeinerte Problembeschreibung

Die Ausbreitungsgeschwindigkeit  $c$  in (1) hängt vom Material der Saite ab. Bisher haben wir sie als konstant angenommen, d.h. die Saite bestand aus einem Material. Sei nun für  $c_0, c_1 \in \mathbb{R}$

$$c(x) := \begin{cases} c_0, & x \in (0, 1/2) \\ c_1, & x \in (1/2, 1) \end{cases} . \tag{9}$$

## 2.5 Verallgemeinerte Analytische Lösung

Zuerst leiten wir eine zur Helmholtz-Gleichung (2) ähnliche Gleichung her und geben einen (4) entsprechenden Lösungsansatz an, wenn die Lösung  $v$  auf  $(0, 1)$  stetig differenzierbar sein soll. Dabei betrachten wir eine angepasste Version der Wellengleichung (1).

$$\frac{\partial^2 u}{\partial x^2}(t, x) = \frac{1}{c^2(x)} \frac{\partial^2 u}{\partial t^2}(t, x), \quad x \in (0, 1), \quad t \in \mathbb{R}$$

Wir verwenden jedoch den selben Ansatz, wie Vorher. Das war  $u(x, t) = \Re(v(x)e^{-i\omega t})$ , mit einer festen, aber unbekannten Kreisfrequenz  $\omega > 0$  und einer Funktion  $v$ , welche nur noch vom Ort  $x$  abhängt.

Einsetzen und analoges Nachrechnen, gibt, mit der unbekannten Wellenzahl  $\kappa(x) := \frac{\omega}{c(x)}$ , die Randbedingungen (3) und

$$-v''(x) = \kappa^2(x)v(x), \quad x \in (0, 1).$$

Um Probleme mit der Differenzierbarkeit von  $\kappa$  zu vermeiden, führen wir die Abkürzungen  $\kappa_0 := \frac{\omega}{c_0}$ ,  $\kappa_1 := \frac{\omega}{c_1}$  ein, und definieren den Lösungsansatz durch Fallunterscheidung und mit (vorerst) beliebigen Konstanten  $C_{01}, C_{02}, C_{11}, C_{12}$ .

$$v(x) := \begin{cases} C_{01} \cos(\kappa_0 x) + C_{02} \sin(\kappa_0 x), & x \in (0, 1/2) \\ C_{11} \cos(\kappa_1 x) + C_{12} \sin(\kappa_1 x), & x \in (1/2, 1) \end{cases}$$

Durch Berücksichtigung der Randbedingungen (3), erhält man (fast analog zu Vorher)

$$C_{01} = 0, \quad C_{11} \cos \kappa_1 + C_{12} \sin \kappa_1 = 0.$$

Soll  $v$  auf  $1/2$  stetig fortgesetzt werden, so müssen dessen links- und rechts-seitiger Grenzwert übereinstimmen.

$$C_{02} \sin(\kappa_0/2) = \lim_{x \rightarrow 1/2-} v(x) = \lim_{x \rightarrow 1/2+} v(x) = C_{11} \cos(\kappa_1/2) + C_{12} \sin(\kappa_1/2)$$

Um stetige Differenzierbarkeit zu erhalten, muss auch die Ableitung

$$v'(x) = \begin{cases} C_{02} \kappa_0 \cos(\kappa_0 x), & x \in (0, 1/2) \\ -C_{11} \kappa_1 \sin(\kappa_1 x) + C_{12} \kappa_1 \cos(\kappa_1 x), & x \in (1/2, 1) \end{cases}$$

auf  $1/2$  stetig fortgesetzt werden.

$$C_{02} \kappa_0 \cos(\kappa_0/2) = \lim_{x \rightarrow 1/2-} v'(x) = \lim_{x \rightarrow 1/2+} v'(x) = -C_{11} \kappa_1 \sin(\kappa_1/2) + C_{12} \kappa_1 \cos(\kappa_1/2)$$

Aus den Randbedingungen und stetigen Fortsetzungen, ergibt sich also das homogene lineare Gleichungssystem  $R\mathbf{C} = 0$ , mit

$$R := \begin{pmatrix} \sin(\kappa_0/2) & -\cos(\kappa_1/2) & -\sin(\kappa_1/2) \\ \kappa_0 \cos(\kappa_0/2) & \kappa_1 \sin(\kappa_1/2) & -\kappa_1 \cos(\kappa_1/2) \\ 0 & \cos \kappa_1 & \sin \kappa_1 \end{pmatrix} \in \mathbb{R}^{3 \times 3}, \quad \mathbf{C} := \begin{pmatrix} C_{02} \\ C_{11} \\ C_{12} \end{pmatrix} \in \mathbb{R}^{3 \times 1}.$$



Sei  $R \in \text{GL}_3(\mathbb{R})$  regulär, so ist deren Kern trivial, d.h.  $\ker R = \{0\}$ , und somit auch die Lösung  $\mathbf{C} = 0$ . Dieser Trivialfall wurde jedoch vorhin bereits ausgeschlossen. Darum betrachten wir  $\det R = 0$ . Mit SymPy berechnet man

$$\begin{aligned}\det R &= \sin\left(\frac{\kappa_0}{2}\right) \cos\left(\frac{\kappa_1}{2}\right) \kappa_1 + \sin\left(\frac{\kappa_1}{2}\right) \cos\left(\frac{\kappa_0}{2}\right) \kappa_0 \\ &= \sin\left(\frac{\omega}{2c_0}\right) \cos\left(\frac{\omega}{2c_1}\right) \frac{\omega}{c_1} + \sin\left(\frac{\omega}{2c_1}\right) \cos\left(\frac{\omega}{2c_0}\right) \frac{\omega}{c_0} =: f_c(\omega)\end{aligned}\quad (10)$$

Für feste  $c_0, c_1$ , lässt sich das gewünschte  $\omega$ , als (nicht eindeutige) Nullstelle dieser Funktion  $f_c$ , charakterisieren.

## 2.6 Verallgemeinerte Semi-Analytische Lösung

Das Ergebnis aus (10) wurde, in Form von folgender Funktion, implementiert.

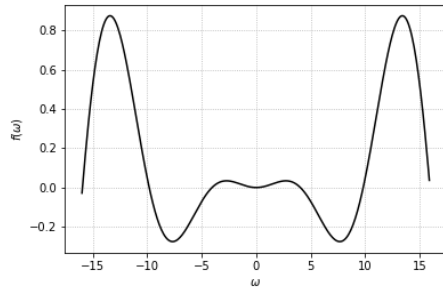
```

1 # c ... pair of propagation speeds
2
3 def get_zero_function(c):
4
5     # allocate some sympy symbols:
6     omega = sp.Symbol('\omega')
7     kappa = sp.IndexedBase('\kappa')
8
9     # implement the matrix R (properly):
10    R = sp.Matrix([[ sp.sin(kappa[0]/2),  kappa[0]*sp.cos(kappa[0]/2),  0],
11                  [-sp.cos(kappa[1]/2),  kappa[1]*sp.sin(kappa[1]/2),  sp.cos(kappa[1])],
12                  [-sp.sin(kappa[1]/2),  -kappa[1]*sp.cos(kappa[1]/2),  sp.sin(kappa[1])])
13    R = R.T
14
15    # calculate R's determinant (properly):
16    det = sp.det(R)
17    det = sp.simplify(det)
18
19    # substitute for kappa_0 and kappa_1:
20    kappa_0 = omega/c[0]
21    kappa_1 = omega/c[1]
22    substitution = {kappa[0]: kappa_0, kappa[1]: kappa_1}
23    det = det.subs(substitution)
24
25    # transform expression det into proper numpy function:
26    zero_function = sp.lambdify(omega, det, 'numpy')
27
28    return zero_function

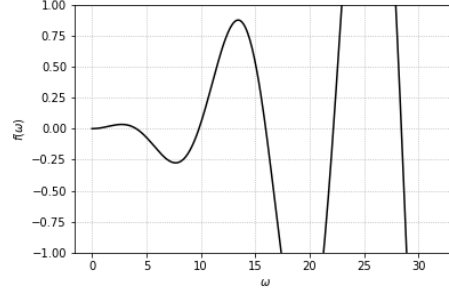
```

Nun können wir  $f_c$  für fixe  $c_0, c_1$  plotten lassen. Dann bekommen wir ein besseres Verständnis dafür, welchen Startwert wir wählen sollen, um die Gleichung  $f_c(\omega) = 0$ , mit `scipy.optimize.fsolve` lösen zu lassen.

Für die folgenden Plots in Abbildung 7, wurden die arbiträren Werte  $c_0 = 100$ ,  $c_1 = 1$  gewählt. Die davon abhängige Funktion  $f_c$  ist scheinbar gerade. Das liegt an (10), sowie dass  $\cos$  gerade und  $\sin$  ungerade ist.



(a) auf dem Intervall  $(-16, 16)$



(b) auf dem Intervall  $(0, 64)$  und herangezoomt

Abbildung 7: Plots von  $f_c$  für  $c_0 = 100$ ,  $c_1 = 1$

Dementsprechend, können passende Startwerte  $\tilde{\omega}$  für iterative Verfahren gewählt werden. Die jeweils ersten Ergebnisse  $\omega$  vom, bereits erwähnten, `scipy.optimize.fsolve` sind in der folgenden Tabelle. Die Quadrate  $\omega^2$  dieser Ergebnisse sind Approximationen der Grenzwerte der Eigenwerte des nächsten Unterkapitels.

	1	2	3	4	5	6
$\tilde{\omega}$	5.000000	10.000000	15.000000	20.000000	25.000000	30.000000
$\omega$	4.057425	9.826058	15.956815	22.170349	15.956815	28.413934
$\omega^2$	16.462695	96.551423	254.619961	491.524375	254.619961	807.351663

## 2.7 Verallgemeinertes Eigenwertproblem

Wir wollen nun den Differenzenquotienten  $D_h v(x_j)$  verwenden, um ein verallgemeinertes Eigenwertproblem der Form  $A\mathbf{v} = \lambda B\mathbf{v}$  mit Matrizen  $A, B \in \mathbb{R}^{(n-1) \times (n-1)}$  herzuleiten.

Sei abermals  $x_j := jh$ ,  $j = 0, \dots, n$  unsere Zerlegung des Intervalls  $[0, 1]$  mit äquidistanter Schrittweite  $h = 1/n$ . Die Matrix  $-A_n$ , für den Differenzenquotienten  $D_h v(x_j)$ , und der Vektor  $\mathbf{v}^{(n)} := (v(x_1), \dots, v(x_{n-1}))^T$ , bleiben ebenfalls nach wie vor so, wie sie waren.

$B_n \lambda$  soll nun, analog zu Vorher,  $\kappa^2$  repräsentieren. Diesmal, ist  $\kappa$  jedoch als (stückweise konstante) Funktion zu verstehen. Also wird die Matrix  $B_n$  deren Fallunterscheidungen übernehmen und  $\lambda$  konstant bleiben. Es läuft darauf hinaus, dass

$$B_n := \begin{cases} \text{diag}^{-2}(c_0, \dots, c_0, c_1, \dots, c_1), & n-1 \in 2\mathbb{N} \\ \text{diag}^{-2}(c_0, \dots, c_0, \frac{c_0+c_1}{2}, c_1, \dots, c_1), & n-1 \in 2\mathbb{N}+1 \end{cases}, \quad \lambda := \omega^2,$$

wobei  $c_0, c_1$  in  $B_n \in \text{GL}_{n-1}(\mathbb{R})$  jeweils  $\lfloor \frac{n-1}{2} \rfloor$ -mal vorkommen. Dabei sei vorausgesetzt, dass  $c_0, c_1 \neq 0$  und  $c_0 + c_1 \neq 0$ . Die Wahl von  $B_n$  lässt sich wie folgt begründen.

Seien  $\mathbf{a}, \mathbf{b}$  Vektoren mit gleich vielen Komponenten. Dann ist die Matrix-Vektor-Multiplikation  $\cdot$ , mit einer erzeugten Diagonalmatrix, äquivalent zur komponentenweisen Multiplikation  $\odot$ .

$$\text{diag}(\mathbf{a}) \cdot \mathbf{b} = \mathbf{a} \odot \mathbf{b} = \mathbf{b} \odot \mathbf{a} = \text{diag}(\mathbf{b}) \cdot \mathbf{a} \quad (11)$$

Bei dem vorherigen Eigenwertproblem wäre  $B_n$  als Einheits-Matrix  $I_n$  zu interpretieren. Der Eigenwert  $\lambda$  konnte gleich ganz  $\kappa^2$  approximieren, weil dieser Wert konstant war. Man hätte aber freilich auch mit der Skalarmatrix  $(I_n c)^{-2}$  und  $\omega^2$  anstelle von  $\kappa^2$  arbeiten können. Da  $\kappa^2$  nun aber, als Funktion, zwei unterschiedliche Werte

$$\left(\frac{\omega}{c_0}\right)^2, \left(\frac{\omega}{c_1}\right)^2$$

annehmen kann, müssen wir die obere Eigenschaft (11) von Diagonalmatrizen ausnutzen. Damit realisieren wir die Fallunterscheidung zwischen  $x_j < 1/2$  und  $x_j > 1/2$ . Für  $x_j = 1/2$ , was genau bei  $n - 1 \in 2\mathbb{N}$  auftritt, wird gemittelt.

Nachdem die Inverse einer Diagonalmatrix genau die Matrix selbst mit komponentenweise Kehrwerten ist, lassen sich gleich  $B_n$  und  $B_n^{-1}$  leicht implementieren. Die zuständigen Funktionen besitzen die kreativen Namen `my_other_numpy_matrix` bzw. `my_other_numpy_matrix_inverse`.

```

1 def my_other_numpy_matrix_inverse(n, c):
2
3     times = np.floor((n-1)/2)
4     times = int(times)
5
6     lower = [c[0]]*times
7     upper = [c[1]]*times
8
9     middle = [(c[0] + c[1])/2]
10
11     if n%2 != 0:
12         B_inverse = np.diag(lower + upper)**2
13         return B_inverse
14     else:
15         B_inverse = np.diag(lower + middle + upper)**2
16         return B_inverse
17
18 def my_other_numpy_matrix(n, c):
19     return 1/my_other_numpy_matrix_inverse(n, c)

```

Das verallgemeinerte Eigenwertproblem  $A_n \mathbf{v} = \lambda B_n \mathbf{v}$  werden wir zunächst auf  $B_n^{-1} A_n \mathbf{v} = \lambda \mathbf{v}$  umformulieren. Dieses kann nun ebenfalls mit `np.linalg.eig`, für beliebige  $n \geq 2$ , gelöst werden.

Die Matrix  $B_n^{-1} A_n$  besitzt hoffentlich wieder  $n - 1$  paarweise verschiedene Eigenwerte  $\lambda_{1,n}^c < \dots < \lambda_{n-1,n}^c$ , die jeweils konvergieren.

Um einen ersten Eindruck des möglichen Konvergenz-Verhaltens zu bekommen, vergleichen wir die Eigenwerte mit den oberen semi-analytischen Ergebnissen mit  $c_0 = 100$ ,  $c_1 = 1$ . Es folgt eine zur oberen analoge Tabelle mit Eigenwerten. Die Ergebnisse lassen zwar zu wünschen übrig, aber immerhin pendelt sich die Größenordnung rasch ein.

	2	3	4	5	6
1	20402.0	13.499662	21.329378	15.313793	20.048572
2	NaN	180004.500338	56813.374144	68.017688	96.940091
3	NaN	NaN	344805.296478	250012.501563	102552.962302
4	NaN	NaN	NaN	750004.166956	422576.319931

Wir bezeichnen (optimistischerweise) mit  $\epsilon_i^c(n) := |\lambda_i^c - \lambda_{i,n}^c|$ ,  $i = 1, \dots, n - 1$  den absoluten Konvergenz-

Fehler des  $i$ -ten Eigenwertes.  $\lambda_i^c$  erhalten wir durch  $\omega^2$  von `scipy.optimize.fsolve`, wobei  $\tilde{\omega} := \sqrt{\lambda_{i,n}^c}$  als Startwert für das iterative Verfahren gewählt wird. Theoretisch hängt  $\lambda_i^c$  also noch von  $n$  ab. In der folgenden Abbildung 8 wurde  $\epsilon_i^c$  mit der Vergleich-Geraden  $\text{id}^2$ , doppelt logarithmisch, geplottet.

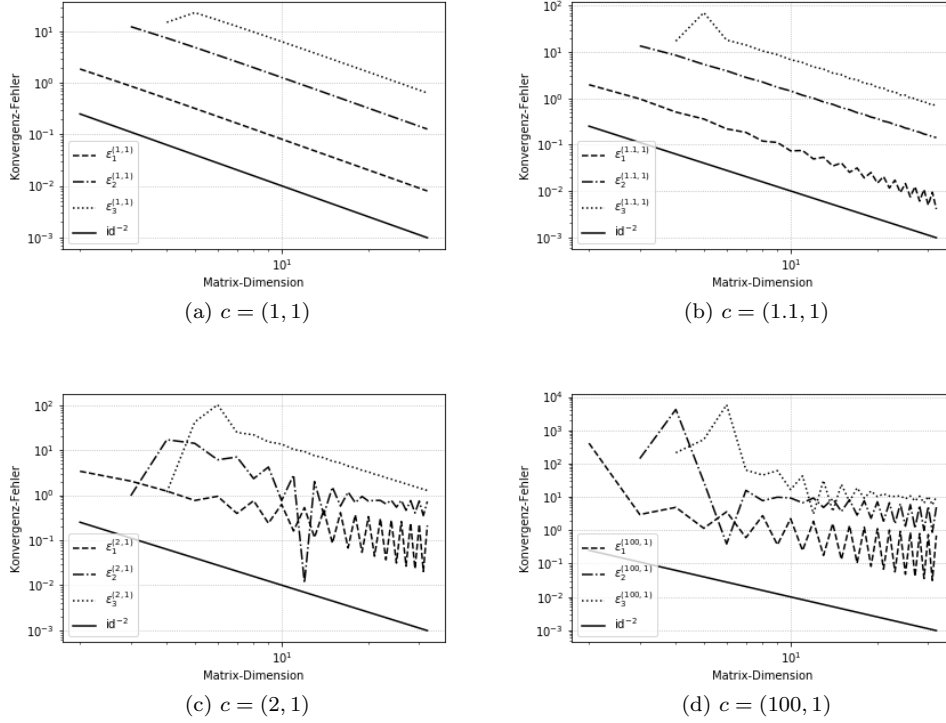


Abbildung 8: Konvergenz-Fehler der Eigenwerte von  $B_n^{-1}A_n$ , für  $n = 2, \dots, 32$  und

Es macht Sinn, dass Abbildung 8a mit Abbildung 5 korreliert. Diese „Bergsteiger-Konvergenz“ steigt anscheinend mit dem Verhältnis  $c_0/c_1$  (no pun intended). Etwas Aufschlussreicher werden 8c und 8d, wenn man gerade und ungerade  $n$  unterscheidet. Diese Abbildung 8 bleibt übrigens unverändert, wenn man die Komponenten von  $c$  vertauscht.

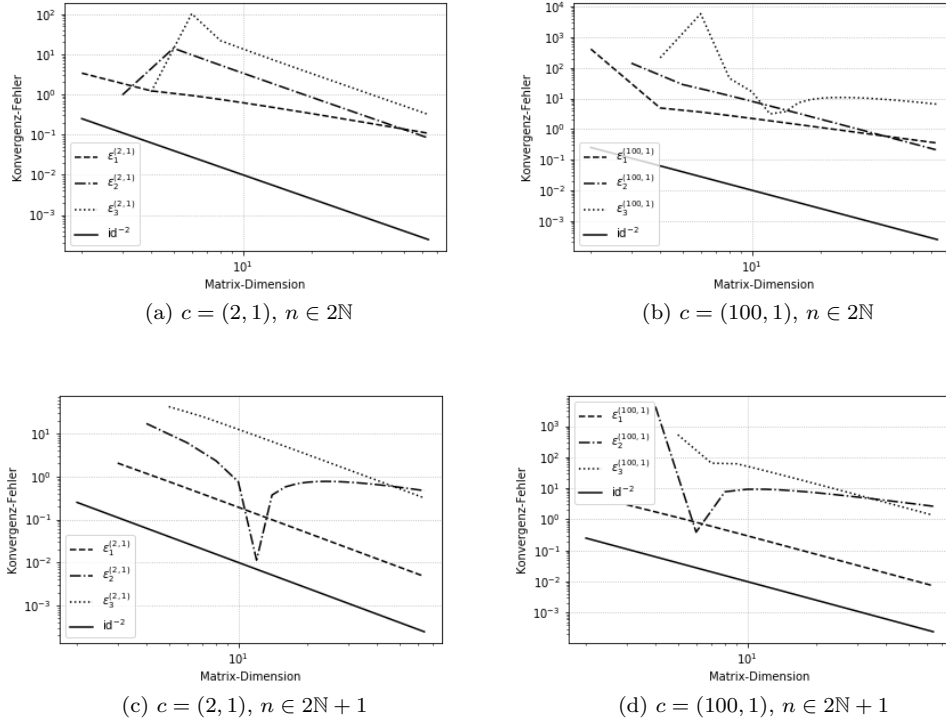


Abbildung 9: Konvergenz-Fehler der Eigenwerte von  $B_n^{-1}A_n$ , für  $n = 2, \dots, 64$  und

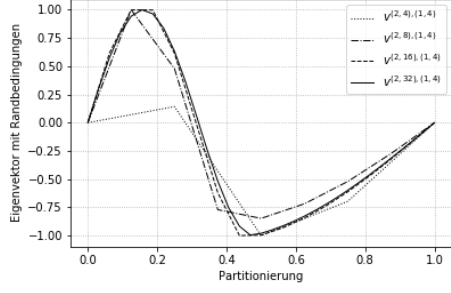
Wenn man für  $c = (100, 1)$  gerade und ungerade betrachtet, so bemerkt man, dass  $\epsilon_2^{(100,1)}$  bei Abbildung 9b konvergiert und  $\epsilon_1^{(100,1)}, \epsilon_3^{(100,1)}$  bei Abbildung 9d. Dies lässt vermuten, dass im Allgemeinen die besten Approximationsstrategie ist, Teilfolgen zu betrachten.

$$\text{Benutze } \begin{cases} (\lambda_{i,n}^c)_{n \in 2\mathbb{N}} & \text{für } i \in 2\mathbb{N}, \\ (\lambda_{i,n}^c)_{n \in 2\mathbb{N}+1} & \text{für } i \in 2\mathbb{N} + 1. \end{cases}$$

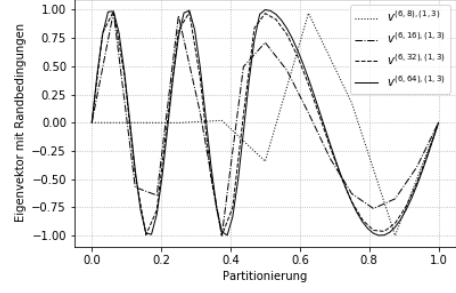
Seien  $\mathbf{v}^{(1,n),c}, \dots, \mathbf{v}^{(n-1,n),c}$  die Eigenvektoren (modulo Konstante), zu den Eigenwerten  $\lambda_{1,n}^c < \dots < \lambda_{n-1,n}^c$ , der Matrix  $B_n^{-1}A_n$ . Wir antizipieren ein weniger schönes Konvergenz-Verhalten, als das der Eigenvektoren der Matrizen  $(A_n)_{n \in \mathbb{N}}$ . Nichts desto trotz, wurden die Eigenvektoren  $\mathbf{v}^{(i,n_i),c}$ , normiert bzgl.  $\|\cdot\|_\infty$ , geplottet, wobei

$$\begin{aligned} c_{\max} &= 4, & c &\in \{(c_0, c_1) : c_0, c_1 = 1, \dots, c_{\max}\}, \\ i &= 1, \dots, 2c_{\max}, & n_i &\in \{2^{p_{\min}} + p_{\text{add}} : p_{\min} = \lceil \log_2(i+1) \rceil, p_{\text{add}} = 0, \dots, 3\} =: N_i. \end{aligned}$$

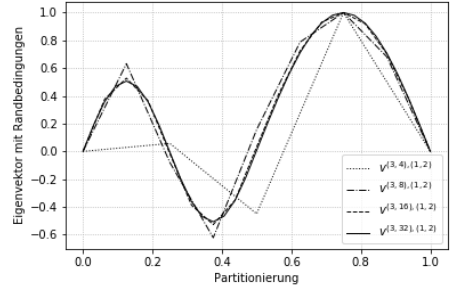
Dabei wurden die Vektoren  $\{\mathbf{v}^{(i,n),c} : n \in N_i\}$  jeweils zu einem Bild zusammengefasst. Nachdem wir dadurch auf 128 Bilder kommen, werden wir nicht alle herzeigen. Außerdem, sieht nur ein Bruchteil davon „schön“ aus.



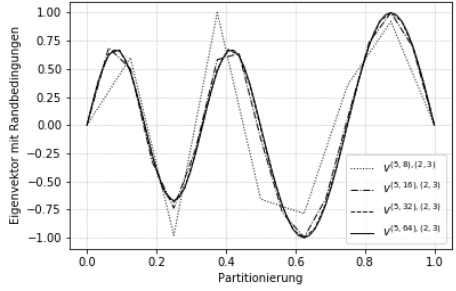
(a) mit  $i = 2$ ,  $n = 4, 8, 16, 32$ , und  $c = (1, 4)$



(b) mit  $i = 6$ ,  $n = 8, 16, 32, 64$ , und  $c = (1, 3)$



(c) mit  $i = 3$ ,  $n = 4, 8, 16, 32$ , und  $c = (1, 2)$



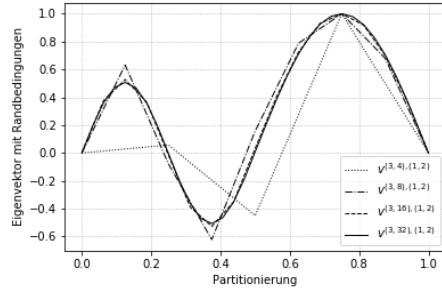
(d) mit  $i = 5$ ,  $n = 8, 16, 32, 64$ , und  $c = (2, 3)$

Abbildung 10: Eigenvektoren  $v^{(i,n),c}$  der Matrizen  $B_n^{-1}A_n$

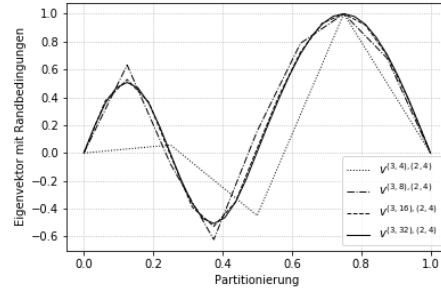
Mit „schön“ ist gemeint, dass die Grenzfunktionen bzgl.  $n$  der Eigenvektoren  $\mathbf{v}^{(i,n),c}$  seinen letzte (halbe) Schwingungs-Periode vollenden kann, bevor die nächste Ausbreitungsgeschwindigkeit übernimmt. Mit anderen Worten, die beiden Funktionshälften treffen sich an der  $x$ -Achse. Man fragt sich nun vielleicht, für welche Ausbreitungsgeschwindigkeiten  $c$  und Eigenpaar-Nummerierung  $i$  diese Grenzfunktionen „schön“ aussehen.

Dazu bemerken wir zuerst, dass die selben Plots herauskommen, wenn  $c^{(1)}, c^{(2)}$  bis auf eine Konstante übereinstimmen. Wegen der Normierung bzgl.  $\|\cdot\|_\infty$ , sieht man das durch hinschauen (auf  $A\mathbf{v} = \lambda B\mathbf{v}$ ), also eigentlich auch bereits a priori.

$$v^{(i,n),c^{(1)}} = v^{(i,n),c^{(2)}}, \text{ für } c^{(1)} \equiv c^{(2)} \text{ mod Konstante.}$$



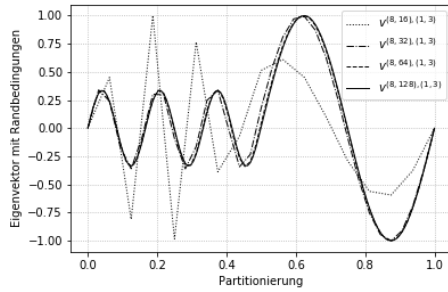
(a)  $c = (1, 2)$



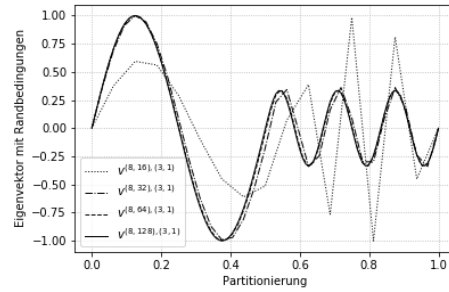
(b)  $c = (2, 4)$

Abbildung 11: Eigenvektoren  $v^{(3,n),c}$  der Matrizen  $B_n^{-1}A_n$ , mit  $n = 8, 16, 32, 64$  und

A posteriori hingegen, bemerken wir, dass zwei Plots mit  $(c_0, c_1)$  bzw.  $(c_1, c_0)$  auch graphisch zusammenhängen. Die erste und zweite Hälfte der Grenzfunktionen tauschen, wenn man zwischen den Plots wechselt.



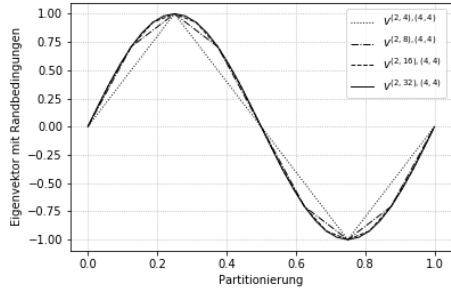
(a)  $c = (1, 3)$



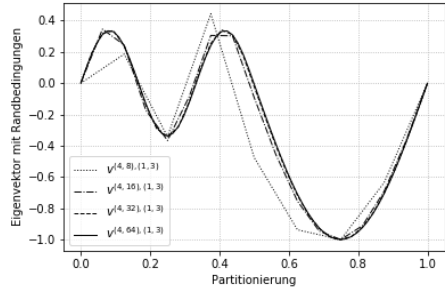
(b)  $c = (3, 1)$

Abbildung 12: Eigenvektoren  $v^{(8,n),c}$  der Matrizen  $B_n^{-1}A_n$ , mit  $n = 16, 32, 64, 128$  und

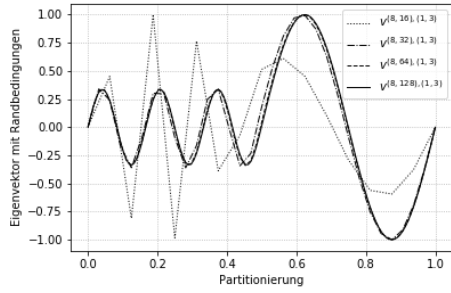
Es sticht aber noch etwas Anderes ins Auge. Wenn man irgendeinen dieser Plots betrachtet, dann geht die Grenzfunktion genau dann durch den ausgezeichneten Punkt  $(1/2, 0)$ , wenn  $\tilde{c}_0 + \tilde{c}_1 \mid i$ , wobei  $\tilde{c}_0 : \tilde{c}_1 = c_0 : c_1$  und  $\tilde{c}$  vollständig gekürzt ist.



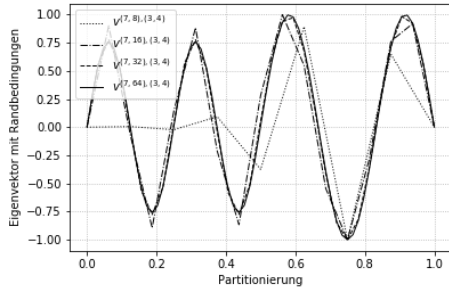
(a) mit  $i = 2$ ,  $n = 4, 8, 16, 32$ , und  $c = (4, 4)$



(b) mit  $i = 4$ ,  $n = 8, 16, 32, 64$ , und  $c = (1, 3)$



(c) mit  $i = 8$ ,  $n = 8, 16, 32, 64$ , und  $c = (1, 3)$



(d) mit  $i = 7$ ,  $n = 8, 16, 32, 64$ , und  $c = (3, 4)$

Abbildung 13: Eigenvektoren  $v^{(i,n),c}$  der Matrizen  $B_n^{-1}A_n$

Das kann man heuristisch so begründen, dass die Ausbreitungsgeschwindigkeiten  $c$ , gemeinsam mit  $i$ , für die relativen Anzahlen der Schwingungsperioden verantwortlich sind; d.h. diese stehen im reziproken Verhältnis zu einander. Für  $i = 4$ ,  $c = (1, 3)$  aus Abbildung 13c zum Beispiel, schwingt die erste Hälfte doppelt so schnell, wie die zweite. Dasselbe geschieht für  $i = 8$ , wobei die Grenzfunktion als Ganzes doppelt so schnell schwingt, wie bei  $i = 4$ . Durch dieses Zwischenspiel von  $c$  und  $i$ , finden wir die meisten „schönen“ Funktionen.

Was ist aber mit dem „künstlich verallgemeinerten“ Fall  $c = (1, 1) = \dots = (c_{\max}, c_{\max})$ ,  $i \in \mathbb{N} - 1$ ? Wieso sieht der so „schön“ aus? Um eine Erklärung zu finden, verallgemeinern wir fröhlich weiter und betrachten die Material-Funktion  $c = (c_0, \dots, c_m) \in (\mathbb{Q}^+)^{m+1}$ , wobei diese, analog zu (9), bzgl.  $1/(m+1)$  äquidistant zu verstehen ist.

$$c(x) := \begin{cases} c_0, & x \in (0, 1/(m+1)) \\ \vdots & \vdots \\ c_m, & x \in (m/(m+1), 1) \end{cases}$$

Wir stellen fest, dass sich die naheliegende Verallgemeinerung der oberen Regel für diese Fälle gilt. Die Grenzfunktion geht genau dann durch die ausgezeichneten äquidistanten Punkte  $\{k/(m+1)\}_{k=0}^{m+1}$ , wenn  $|\tilde{c}| := \tilde{c}_0 + \dots + \tilde{c}_m \mid i$ , wobei  $\tilde{c}_0 : \dots : \tilde{c}_m = c_0 : \dots : c_m$  und  $\tilde{c}$  als Ganzes vollständig gekürzt ist.

Tatsächlich gibt es im „künstlich verallgemeinerten“ Fall mit beliebigem  $i$  ein  $m$ , sodass für  $c \in (\mathbb{Q}^+)^{m+1}$  gilt  $|\tilde{c}| = m+1 \mid i$ . Dieser Fall liefert übrigens, wegen der bereits angesprochenen Skalarmatrix, immer dasselbe Resultat.



## 2.8 Vektor-Iteration

Sei  $\rho > 0$ . Wir verwenden die Vektor-Iteration angewendet auf das Eigenwertproblem

$$(A - \rho B)^{-1} B \tilde{\mathbf{v}} = \mu \tilde{\mathbf{v}} \quad (12)$$

mit den Matrizen  $A$  und  $B$  aus dem vorigen Unterkapitel.

Das macht man, weil die Vektor-Iteration in ihrer einfachsten Form, `vector_iteration_simple`, nur den Eigenvektor mit dem betragsgrößten Eigenwert liefert.

```
1 # m ... number of iterations
2
3 def vector_iteration_simple(M, m):
4
5     # some random (non zero) vector
6     randy = np.random.rand(M.shape[0])
7     # normalize
8     randy = randy/np.linalg.norm(randy)
9
10    for _ in range(m):
11
12        # apply matrix
13        randy = M @ randy
14        # normalize
15        randy = randy/np.linalg.norm(randy)
16
17    return randy
```

Es besteht ein Zusammenhang zwischen den Eigenpaaren  $(\mu, \tilde{\mathbf{v}})$  und den Eigenpaaren  $(\lambda, \mathbf{v})$  aus der vorigen Aufgabe. Für  $\mu \neq 0$ , ist dieser  $\lambda = \rho + \frac{1}{\mu}$ , weil dann

$$A\mathbf{v} = \lambda B\mathbf{v} \Leftrightarrow A\mathbf{v} = \left(\rho + \frac{1}{\mu}\right)B\mathbf{v} \Leftrightarrow B\mathbf{v} = \mu(A - \rho B)\mathbf{v} \Leftrightarrow (A - \rho B)^{-1}B\mathbf{v} = \mu\mathbf{v}.$$

Bei der Vektor-Iteration wird der potentielle Eigenvektor  $\mathbf{v}$  ständig normiert. Damit kann man wegen  $A\mathbf{v} = \lambda\mathbf{v}$  sehr rasch auf den zugehörigen Eigenwert  $\lambda$  kommen.

$$\langle A\mathbf{v}, \mathbf{v} \rangle = \langle \lambda\mathbf{v}, \mathbf{v} \rangle = \lambda \|\mathbf{v}\|^2 = \lambda$$

Jetzt können wir auch ein schlauereres Abbruch-Kriterium wählen. Die Vektor-Iteration soll terminieren, wenn die Änderung der Eigenwerte hinreichend klein ist. Die Implementierung, die diese Überlegungen beherzigt, folgt mit `vector_iteration_unshifted`.

```
1 def vector_iteration_unshifted(M, tol):
2
3     # some random (non zero) vector
4     randy = np.random.rand(M.shape[0])
5     # normalize
6     randy = randy/np.linalg.norm(randy)
7
8     # stop iteration when differences are small enough
9     eigen_randy_old = -tol
10    eigen_randy_new = tol
11
```

```

12 while abs(eigen_randy_old - eigen_randy_new) > tol:
13
14     # get eigen values of randy
15     eigen_randy_old = np.dot(M @ randy, randy)
16
17     # apply matrix
18     randy = M @ randy
19     # normalize
20     randy = randy/np.linalg.norm(randy)
21
22     # get eigen values of randy
23     eigen_randy_new = np.dot(M @ randy, randy)
24
25     # use best approximation of eigen value
26     eigen_randy = eigen_randy_new
27
28     # get eigen pair
29     eigen_pair = (eigen_randy, randy)
30
31 return eigen_pair

```

Die untere Implementierung, `vector_iteration_shifted`, der Vektor-Iteration des geshifteten Eigenwertproblems (12), ist ein gutes Beispiel eines sinnvollen Einsatzes der *LU*-Zerlegung. Diese ist zwar aufwändig, aber nachdem sie in jedem Iterationsschritt verwendet werden kann, zahlt sich dieser einmalige Aufwand aus.

```

1 from scipy.linalg import lu, solve_triangular
2
3 def vector_iteration_shifted(n, c, rho, tol):
4
5     # some random (non zero) vector
6     randy = np.random.rand(n-1)
7     # normalize
8     randy = randy/np.linalg.norm(randy)
9
10    # get matrices
11    A = my_numpy_matrix(n)
12    B = my_other_numpy_matrix(n, c)
13    B_inverse = my_other_numpy_matrix_inverse(n, c)
14    B_inverse_A = B_inverse @ A
15
16    # calculate lu-decomposition and apply permutation matrix
17    M = A - rho*B
18    P, L, U = lu(M)
19
20    # used to get eigen value of randy
21    M = np.linalg.inv(A - rho*B) @ B
22
23    # stop iteration when differences are small enough
24    eigen_randy_old = -tol
25    eigen_randy_new = tol
26
27    while abs(eigen_randy_old - eigen_randy_new) > tol:
28
29        # get eigen values of randy
30        eigen_randy_old = np.dot(M @ randy, randy)
31
32        # apply first part of matrix to randy
33        randy = P @ B @ randy
34
35        # solve L @ forwards = randy via forwards substitution
36        forwards = solve_triangular(L, randy, lower = True)
37        # solve U @ backwards = forwards per backwards substitution
38        backwards = solve_triangular(U, forwards, lower = False)

```

```

39         # apply second part of matrix to randy
40         randy = backwards
41         # normalize
42         randy = randy/np.linalg.norm(randy)
43
44         # get eigen values of randy
45         eigen_randy_new = np.dot(M @ randy, randy)
46
47     # use best approximation of eigen value
48     eigen_randy = eigen_randy_new
49
50     # get eigen pair of shifted problem
51     eigen_pair_shifted = (eigen_randy, randy)
52
53     # get eigen pair of unshifted problem
54     eigen_randy_unshifted = 1/eigen_randy + rho
55     eigen_pair_unshifted = (eigen_randy_unshifted, randy)
56
57     return eigen_pair_shifted, eigen_pair_unshifted
58

```

Angenommen, wir wüssten bereits, dass all unsere Eigenwerte  $\lambda_{1,n}^c < \dots < \lambda_{n-1,n}^c \in \mathbb{R}^+$  reell und positiv sind. `vector_iteration_simple` liefert uns den Eigenvektor mit dem betragsgrößten Eigenwert und `vector_iteration_shifted` den Eigenvektor mit dem Eigenwert, der am nächsten bei  $\rho$  liegt. Nun können wir mit einer binären Suche, ähnlich zum Bisektionsverfahren, alle Eigenwerte finden und ein eigenes `np.linalg.eig` programmieren.

```

1 from timeit import default_timer as timer
2
3 def recursion(n, c, tol, eigen_pairs, bound_lower, bound_upper):
4
5     bound_middle = rho = (bound_lower + bound_upper)/2
6
7     print("searching for eigen pair with eigen value near", rho, "...")
8     start = timer()
9     eigen_pair_middle = vector_iteration_shifted(n, c, rho, tol)[1]
10    end = timer()
11    bound_middle = eigen_pair_middle[0]
12    print("found eigen value", bound_middle, "in", end-start, "seconds", "\n")
13
14    if abs(bound_middle - bound_lower) < 10 or bound_middle <= bound_lower:
15        print("bound_lower ~ bound_middle")
16        print(bound_lower, "~", bound_middle)
17        print("or")
18        print("bound_middle <= bound_lower")
19        print(bound_middle, "<=", bound_lower)
20        print("stopping pursuit", "\n")
21        return False
22    if abs(bound_upper - bound_middle) < 10 or bound_upper <= bound_middle:
23        print("bound_middle ~ bound_upper")
24        print(bound_middle, "~", bound_upper)
25        print("or")
26        print("bound_upper <= bound_middle")
27        print(bound_upper, "<=", bound_middle)
28        print("stopping pursuit", "\n")
29        return False
30
31    eigen_pairs.update([eigen_pair_middle])
32    print("eigen_pair_middle (new one) =", eigen_pair_middle, "\n")
33
34    if len(eigen_pairs) == n-1:
35        print("len(eigen_pairs) == n-1")

```

```

36     print(len(eigen_pairs), "=", n-1)
37     print("All eigen pairs found", "\n")
38     return True
39
40     print("repeating recursion with:")
41     print("bound_lower =", bound_lower)
42     print("bound_middle =", bound_middle, "\n")
43     recursion(n, c, tol, eigen_pairs, bound_lower, bound_middle)
44
45     if len(eigen_pairs) == n-1:
46         print("len(eigen_pairs) == n-1")
47         print(len(eigen_pairs), "=", n-1)
48         print("All eigen pairs found", "\n")
49         return True
50
51     print("repeating recursion with:")
52     print("bound_middle =", bound_middle)
53     print("bound_upper =", bound_upper, "\n")
54     recursion(n, c, tol, eigen_pairs, bound_middle, bound_upper)
55
56     return True
57
58 def get_my_eigen_pairs_with_vector_iteration(n, c, tol):
59
60     print("Starting search with n =", n, "...", "\n")
61
62     B_inverse_A = my_general_numpy_matrix(n, c)
63     eigen_pairs = {}
64
65     rho = 0
66
67     eigen_pair_upper = vector_iteration_unshifted(B_inverse_A, tol)
68     bound_upper = eigen_pair_upper[0]
69
70     eigen_pairs.update([eigen_pair_upper])
71     print("eigen_pair_upper =", eigen_pair_upper, "\n")
72
73     if len(eigen_pairs) == n-1:
74         print("len(eigen_pairs) == n-1")
75         print(len(eigen_pairs), "=", n-1)
76         print("All eigen pairs found", "\n")
77         return eigen_pairs
78
79     eigen_pair_lower = vector_iteration_shifted(n, c, rho, tol)[1]
80     bound_lower = eigen_pair_lower[0]
81
82     eigen_pairs.update([eigen_pair_lower])
83     print("eigen_pair_lower =", eigen_pair_lower, "\n")
84
85     if len(eigen_pairs) == n-1:
86         print("len(eigen_pairs) == n-1")
87         print(len(eigen_pairs), "=", n-1)
88         print("All eigen pairs found", "\n")
89         return eigen_pairs
90
91
92     print("bound_lower =", bound_lower)
93     print("bound_upper =", bound_upper, "\n")
94     recursion(n, c, tol, eigen_pairs, bound_lower, bound_upper)
95
96     return eigen_pairs

```

Das Abbruch-Kriterium ist aber anscheinend noch immer suboptimal, weil die damit berechneten Eigenwerte, sogar für  $\text{tol} = 1\text{e-}15$ , nicht ganz denen von `np.linalg.eig` entsprechen.

	2	3	4	5	6
1	20402.0	13.499662	21.329378	15.313793	20.048572
2	NaN	180004.500338	56813.374137	68.003831	96.940090
3	NaN	NaN	344805.296478	180595.546136	102552.962297
4	NaN	NaN	NaN	750004.166956	392219.488916

Zum Vergleich, zeigen wir nochmal die Eigenwerte via `np.linalg.eig`. Betrachtet man die Eigenwerte  $\lambda_{3,5}^{(100,1)}, \lambda_{4,6}^{(100,1)}$ , so merkt man, dass Ungenauigkeiten beim vorletzten, d.h. zweitgrößten, Eigenpaar auftreten, wenn  $n$  groß ist. Offensichtlich, lässt sich diese Eigenwertsuche noch optimieren, das würde aber den Rahmen dieses Projektes sprengen.

	2	3	4	5	6
1	20402.0	13.499662	21.329378	15.313793	20.048572
2	NaN	180004.500338	56813.374144	68.017688	96.940091
3	NaN	NaN	344805.296478	250012.501563	102552.962302
4	NaN	NaN	NaN	750004.166956	422576.319931

Zuletzt, wollen wir noch wissen, ob die Eigenwerte unserer Matrix  $B_n^{-1}A_n$  tatsächlich alle positiv und reell sind. Sonst wäre ja der Algorithmus, aus mathematische Sicht, wertlos. Eine Möglichkeit bietet die explizite Darstellung von Eigenpaaren von Tridiagonalmatrizen. Wir wissen aber, dass das Gewünschte genau dann gilt, wenn  $B_n^{-1}A_n$  positiv definit ist. Das können wir mit dem Hauptminoren-Kriterium locker überprüfen.

$$A'_n := \begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix}$$

Zuerst berechnen wir die Determinante der Matrix  $A'_n$ . Dazu induzieren wir  $\det(A'_n) = n$ . Der Induktionsanfang,  $n = 2, 3$ , ist trivial. Für den Induktionsschritt muss man einmal nach der letzten Spalte und dann Zeile entwickeln (oder umgekehrt).

Nun gilt aber  $A_n = \frac{1}{h^2}A'_n$  und  $B_n^{-1} > 0$  ist eine Diagonalmatrix mit positiven Einträgen. Weil die Determinante solcher Matrizen bloß das Produkt ihrer Komponenten ist, sind alle Hauptminoren dieser Matrizen positiv. Unsere Matrizen sind also tatsächlich positiv definit. Deren Produkt und Vielfaches  $A_n$  ist es also auch, weil

$$0 < x^*Ax, x^*Bx \Rightarrow 0 < x^*Axx^*Bx = x^*A\langle x, x \rangle Bx = x^*A\|x\|^2 Bx \Rightarrow 0 < x^*ABx.$$

### 3 Titel