

1 Die Methode der Konjugierten Gradienten

1.1 Problem

Lösung des linearen Gleichungssystems

$$Ax = b \quad (1.1)$$

mit $(N \times N)$ -Matrix A und N -komponentigen Spaltenvektoren b und x .

N ist so gross, dass nur Vektoren, aber keine Matrizen gespeichert werden können. A wird implementiert durch eine *Vorschrift*, wie es auf einen Vektor wirkt, d.h. wie $v = Au$ zu bilden ist.

Die bekannten Eliminationsverfahren (Gauss, Householder) zur Lösung des linearen Systems sind nicht anwendbar. Stattdessen werden Wege gesucht, die Lösung aus einer Basis aufzubauen, die durch Anwendung der Matrix A auf *einen* Vektor u entsteht:

$$\{u, Au, A^2u, A^3u, \dots, A^{n-1}u\}$$

Das nennt man einen **Krylov**-Raum.

Die Lösung x entsteht dann iterativ: $\{x_0 \rightarrow x_1 \rightarrow x_2 \dots\}$, ihre Qualität beurteilt man anhand des Restvektors

$$r_k = b - Ax_k \quad (1.2)$$

Von jetzt an nehmen wir der Einfachheit halber an, dass die Matrix A **symmetrisch und positiv** ist. (Ansonsten könnte man statt $Ax = b$ einfach $A^T Ax = A^T b$ lösen.)

Das folgende Verfahren der **Konjugierten Gradienten** geht von einem Startvektor x_0 aus (womöglich $x_0 = 0$) und baut den Krylov-Raum über dem anfänglichen Restvektor $u = r_0$ auf.

1.2 Algorithmus

Initialisierung:

```

if      flag
       $x_0 = 0$ 
       $r_0 = b$ 
else
       $r_0 = b - Ax_0$ 
      if       $r_0^T r_0 < tol$ 
              exit
       $p_0 = r_0$ 

```

Iteration:

```

 $s = Ap_k$ 
 $\alpha_k = \frac{p_k^T r_k}{p_k^T s}$ 
 $x_{k+1} = x_k + \alpha_k p_k$ 
 $r_{k+1} = r_k - \alpha_k s$ 
      if       $r_{k+1}^T r_{k+1} < tol$ 
              exit
       $\beta_k = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$ 
       $p_{k+1} = r_{k+1} + \beta_k p_k$ 

```

Bemerkungen zur Implementierung:

- Die Iteration leistet $\{x_k r_k p_k\} \rightarrow \{x_{k+1} r_{k+1} p_{k+1}\}$, die Vektoren kann man gleich **überschreiben**. Man sollte sich auch noch $r_k^T r_k$ für die nächste Iteration merken, alle anderen Größen haben nur temporäre Bedeutung.
- Es wird ein **Hilfsvektor** s eingeführt, damit man A nur einmal pro Iteration anwenden muss.
- Für α_k und β_k gibt es eine Reihe von **Varianten**, die mathematisch (aber nicht numerisch) äquivalent sind. Brauchbar ist vor allem $\alpha_k = (r_k^T r_k)/(p_k^T s)$, damit wird ein Skalarprodukt eingespart.
- Als **Abbruchbedingung** verlangt man sinnvollerweise, dass die Norm von r_k klein genug ist im Vergleich zu b , etwa $\|r_k\|/\|b\| < relerr$ mit einer *relativen* Genauigkeit $relerr = 10^{-8} \dots 10^{-12}$.
- An **Speicherplatz** braucht man insgesamt 5 Vektoren, wenn man b behalten will, aber nur 4, wenn man b mit r überschreibt.

1.3 Eigenschaften

Die Verschiebungen p_k sind A -konjugiert (daher der Name):

$$p_k^T A p_l = 0 \quad \text{für } k \neq l \quad (1.3)$$

Die Restvektoren sind orthogonal (aber nicht normiert):

$$r_k^T r_l = 0 \quad \text{für } k \neq l \quad (1.4)$$

Daraus folgt *mathematisch*, dass *spätestens* $r_N = 0$, d.h. die CG-Methode ist ein *N-Schritt-Verfahren*. Man beobachtet aber, dass wegen Rundungsfehlern oft erst $r_k = 0$ (Maschinengenauigkeit) für ein $k > N$ erreicht wird.

Die Konvergenz ist monoton in der Norm

$$r_k^T A^{-1} r_k \quad (1.5)$$

(Minimal-Eigenschaft).

Bei grossem N konvergiert das Verfahren exponentiell in dem Sinne, dass z.B.

$$|r_k| \sim e^{-\gamma k} \quad (1.6)$$

Die *Konvergenzrate* γ hängt von der **Konditionszahl** der Matrix ab:

$$\operatorname{cond}(A) \equiv \frac{\lambda_{\max}}{\lambda_{\min}} \quad (1.7)$$

$$\tanh \gamma/2 = \operatorname{cond}(A)^{-1/2} \quad (1.8)$$

$$(1.9)$$

Die für eine gewisse Genauigkeit nötige Zahl von Iterationen ist also proportional zu $\operatorname{cond}(A)^{1/2}$.

1.4 Anwendung: Laplace–Operator auf einem kubischen Gitter

- Rechteckiges Gitter in d Dimensionen mit Gitterkonstante a
- Skalarfeld: $\Phi(x)$

Laplace–Operator:

$$-\Delta\Phi(x) \rightarrow a^{-2} \sum_{k=1}^d [2\Phi(x) - \Phi(x + ae_k) - \Phi(x - ae_k)] \quad (1.10)$$

$$= 2da^{-2}\Phi(x) - a^{-2} \sum_{k=1}^d [\Phi(x + ae_k) + \Phi(x - ae_k)] \quad (1.11)$$

Diagonalisierung im Impulsraum:

$$-\Delta e^{ip \cdot x} = a^{-2} \sum_{k=1}^d 2(1 - \cos p_k a) e^{ip \cdot x} \quad (1.12)$$

$$= \sum_{k=1}^d 4a^{-2} \sin^2(p_k a/2) e^{ip \cdot x} \quad (1.13)$$

1.5 “Massiver” Fall

Einführung einer “Masse” m oder Abschirmlänge $\xi = m^{-1}$ führt auf den Operator

$$A = (-\Delta + m^2) \quad (1.14)$$

Bei periodischen Randbedingungen hat man $N_1 \times N_2 \times \dots \times N_d$ Gitterpunkte und Impulse

$$p_k = \frac{2\pi n_k}{N_k a}, \quad n_k = 0 \dots N_k - 1$$

Konditionszahl und Konvergenzrate:

$$\begin{aligned} \Rightarrow \lambda_{min}(A) &= m^2 \\ \lambda_{max}(A) &= 4da^{-2} + m^2 \\ \Rightarrow \text{cond}(A) &= \frac{4d}{m^2 a^2} + 1 \\ \Rightarrow \sinh(\gamma/2) &= \frac{ma}{2\sqrt{d}} \\ \Rightarrow \gamma &= \frac{ma}{\sqrt{d}} + \mathcal{O}((ma)^3) \end{aligned} \quad (1.15)$$

1.6 Laplace–Problem mit Dirichlet–Randbedingungen

Z.B. Gittergrösse $L_1 \times \dots \times L_d$ mit $L_k = N_k a$.

Das Feld sei vorgegeben, wenn ein $x_k = 0$ oder $x_k = L_k$, es gibt also insgesamt $(N_1 + 1) \times \dots \times (N_d + 1)$ Punkte, davon aber nur $(N_1 - 1) \times \dots \times (N_d - 1)$ innere Punkte, auf denen das Feld bestimmt werden muss.

Diagonlisierung von $A = -\Delta$:

$$\begin{aligned}\psi_p(x) &= \prod_{k=1}^d \sin p_k x_k \\ p_k &= \frac{\pi n_k}{L_k}, \quad n_k = 1 \dots (N_k - 1) \\ \lambda(p) &= \sum_{k=1}^d 4a^{-2} \sin^2(p_k a/2)\end{aligned}\tag{1.16}$$

Konvergenzrate:

$$\begin{aligned}\lambda_{min}(A) &= \sum_{k=1}^d 4a^{-2} \sin^2 \frac{\pi}{2N_k} \\ \lambda_{max}(A) &= 4da^{-2} \\ \Rightarrow \text{cond}(A) &= d / \sum_{k=1}^d \sin^2 \frac{\pi}{2N_k} \\ \Rightarrow \tanh \gamma/2 &= \left[d^{-1} \sum_{k=1}^d \sin^2 \frac{\pi}{2N_k} \right]^{1/2}\end{aligned}\tag{1.17}$$

Speziell für $N_k = N$:

$$\begin{aligned}\tanh \frac{\gamma}{2} &= \sin \frac{\pi}{2N} \\ &= \sin \frac{\pi a}{2L}\end{aligned}\tag{1.18}$$

$$\Rightarrow \gamma = \frac{\pi a}{L} + \mathcal{O}((a/L)^3)\tag{1.19}$$

1.7 Implementierung der Gittergeometrie

Auf einem d -dimensionalen Gitter (Gitterkonstante $a = 1$) der Größe $L_1 \times L_2 \times \dots \times L_d$ soll ein Feld $\Phi(x)$ als *eindimensionaler* Array implementiert werden. Man muss die Gitterpunkte also irgendwie durchnummernieren und dann dafür sorgen, dass man die Nachbarpunkte wiederfindet, ggf. unter Berücksichtigung der Randbedingungen.

Dazu folgende Idee: man durchläuft das Gitter systematisch, x_1 läuft am schnellsten, dann x_2 usw. Wenn man die Koordinaten als $x_k = 0..(L_k - 1)$ ansetzt, dann lässt sich die Nummerierung schreiben als

$$\begin{aligned} i &= x_1 + L_1 x_2 + L_1 L_2 x_3 + \dots \\ i &= 0..(V - 1) \end{aligned}$$

Das hat eine interessante Interpretation: die x_k erscheinen als *Ziffern* in einem Zahlensystem mit (positionsabhängiger) *Basis*

$$\begin{aligned} b_1 &= 1 \\ b_2 &= L_1 \\ b_3 &= L_1 L_2 \\ &\dots \\ \Rightarrow i &= \sum_{k=1}^d b_k x_k \\ b_k &= \prod_{l=1}^{k-1} L_l \end{aligned} \tag{1.20}$$

und der sich so ergebende *Zahlenwert* ist der gesuchte Index. Natürlich braucht man Koordinaten und Index nicht von 0 an zu zählen, aber damit ist es am einfachsten zu erklären.

1.7.1 Periodische Randbedingungen

Nun muss man noch das Indexfeld konstruieren, mit dem man Nachbarn findet. Für periodische Randbedingungen wird das durch das folgende Geometrieprogramm erledigt, das einmal das gesamte Gitter durchläuft, dabei sowohl die x_k als auch den Index i im Auge behält und so für jeden Punkt die Indizes der Nachbarn bestimmen kann:

```

1 #include <stdlib.h>
2 #include "global.h"
3 /*
4     ueber global.h werden folgende globale Variablen eingebunden:
5     int ndim, nvol, *lsize, **nn;
6 */
7
8 void geom_pbc(){

```

```

9  /*
10 Folgende Groessen muessen gesetzt sein:           B Bunk 12/2005
11     Dimension      ndim                         rev    4/2013
12     Gittergroesse lsize[k], k=1..ndim
13
14 Angelegt und berechnet wird
15     Volumen      nvol
16     NN-Indexfeld nn[k][i], k=0..2*ndim, i=0..(nvol-1)
17
18 nn[k][i] gibt den Index des Nachbarn von i in Richtung +k,
19 unter Beruecksichtigung periodischer Randbedingungen.
20 Fuer einen Schritt in Richtung -k setze man den Index auf (ndim+k).
21 nn[0][i] ist reserviert.
22 */
23 int i, k;
24 int *ibase, *ix;
25
26 ibase = (int *) malloc((ndim+2) * sizeof(int));
27 ix = (int *) malloc((ndim+1) * sizeof(int));
28
29 /* Basis fuer Punktindizes */
30 ibase[1] = 1;
31 for (k=1; k<=ndim; k++) ibase[k+1] = ibase[k]*lsize[k];
32 nvol = ibase[ndim+1];
33
34 if (nn) free(nn[0]);
35 free(nn);
36 nn = (int **) malloc((2*ndim+1) * sizeof(int *));
37 nn[0] = (int *) malloc(nvol*(2*ndim+1) * sizeof(int));
38 for (k=1; k<=2*ndim; k++) nn[k] = nn[0] + nvol*k;
39
40 for (k=1; k<=ndim; k++) ix[k] = 0; /* Koord. des Anfangspunkts */
41
42 for (i=0; i<nvol; i++){          /* Schleife ueber Punkte */
43     for (k=1; k<=ndim; k++){
44         nn[k][i] = i + ibase[k];      /* Nachbar x + e_k */
45         if (ix[k] == (lsize[k]-1)) nn[k][i] -= ibase[k+1];
46
47         nn[ndim+k][i] = i - ibase[k]; /* Nachbar x - e_k */
48         if (ix[k] == 0) nn[ndim+k][i] += ibase[k+1];
49     }
50
51     for (k=1; k<=ndim; k++){      /* Koord. des naechsten Punkts */
52         ix[k]++;
53         if (ix[k] < lsize[k]) break;
54         ix[k] = 0;
55     }
56 }
57 free(ibase); free(ix);

```

58 }

Wenn dann das Feld wirklich bearbeitet wird, beispielsweise im CG, werden nur noch die Indizes i durchlaufen, von den $\{x_k\}$ ist nicht mehr die Rede.

1.7.2 Dirichlet–Randbedingungen

Hier sieht man für die festgehaltenen Werte von $\Phi(x)$ auf dem Rand ebenfalls Feldkomponenten vor, die anfangs auf die vorgegebenen Werte gesetzt werden und dann nicht mehr verändert werden dürfen, sondern nur als Nachbarn der inneren (aktiven) Punkte ins Spiel kommen. Der Einfachheit halber werden auch alle anderen Felder entsprechend erweitert. Man implementiert nun den Laplace–Operator so, dass er (bei Anwendung auf irgendein Feld) am Rand Nullen zurückgibt. Dann muss bei einem Problem wie $(-\Delta + m^2)\Phi(x) = \eta(x)$ das Quellfeld ergänzt werden um $\eta(x) = m^2\Phi(x)$ auf dem Rand. Im CG–Verfahren werden damit alle Shift– und Restvektoren auf den Randpunkten ebenfalls verschwinden, und die Gleichungen für $\Phi(x)$ auf den inneren Punkten werden konsistent mit den Randbedingungen gelöst.

Genauer: der anfängliche Restvektor

$$r_0 = \eta - (-\Delta + m^2)\phi_0 \quad (1.21)$$

muss auf dem Rand verschwinden. Dann kann man anhand der Rekursion zeigen, dass *alle* Rest- und Shiftvektoren am Rand verschwinden, wenn eine der beiden Varianten implementiert wird:

$$\Delta \rightarrow 0 \quad \eta = m^2\phi_0 \quad \text{auf dem Rand} \quad (1.22)$$

oder

$$(-\Delta + m^2) \rightarrow 0 \quad \eta = 0 \quad \text{auf dem Rand} \quad (1.23)$$

Bei $m = 0$ fällt beides zusammen.

Um die Unterscheidung zwischen inneren und Randpunkten schnell anhand des Index treffen zu können, führt man ein Feld von Flags ein. Dazu kann das Geometrieprogramm z.B. die Komponente $nn(i, 0)$ verwenden – sie wird ja ansonsten nicht benutzt. Im Übrigen wird $nn(i, k)$ für $k \neq 0$ wie bei periodischen Randbedingungen gesetzt.