

目次

第 1 章	数の表現と数値解析の概観	3
1.1	実数の浮動小数点表示	3
1.2	演算	5
1.3	連立一次方程式	6
1.4	固有値問題	6
第 2 章	行列とその分解と変換	7
2.1	連立一次方程式と Cramer の公式から見た線型計算の世界	7
2.2	行列の標準形	7
2.3	自己共役行列の性質	7
2.3.1	定義と特徴付け	8
2.3.2	一般化	8
2.4	優対角行列と既約行列	8
2.4.1	定義	8
2.4.2	M-行列	9
2.5	特異値分解	9
2.6	LU 分解と Cholesky 分解	10
2.7	基本直交変換	10
2.7.1	Householder 変換	10
2.7.2	Givens 変換	10
2.8	三角化 / QR 分解	10
2.9	Hessenberg 化 / 3 重対角化	10
2.10	Schur 分解	10
2.11	Moore-Penrose の一般化逆行列	10
2.12	行列のノルム	10
2.12.1	p-ノルム	10
2.12.2	Frobenius ノルム	11
2.12.3	スペクトル半径	11
2.12.4	条件数	11
第 3 章	代数：連立一次方程式の消去法	12
3.1	Gauss の消去法	12
第 4 章	解析：連立一次方程式の反復法	13
4.1	ベクトルのノルム	13
4.1.1	p-ノルム	13
4.1.2	ノルムと位相	13
4.2	行列のノルム	14
4.2.1	p-行列ノルム	14
4.2.2	複素行列の行列ノルム	14

4.3	定常反復法	15
4.4	安定性と条件数	15
4.4.1	後退誤差解析	15
第 5 章	連立一次方程式と共役勾配法	17
第 6 章	非線形方程式	18
6.1	C^k 級関数と Taylor の定理	18
6.2	二分法	18
6.3	反復法と不動点定理	19
6.4	ベクトル値関数の微分	20
6.5	多変数の反復法	20
第 7 章	固有値問題	21
第 8 章	関数近似	22
第 9 章	補間と積分	23
第 10 章	常微分方程式の初期値問題	24
参考文献		25

第 1 章

数の表現と数値解析の概観

コンピュータの性能が向上しようと、その数学的な限界（有限性や数の表現など）を議論できるのは、その母体である数学のみである。

1.1 実数の浮動小数点表示

指数表現によって可能な十分に広い絶対値の範囲内において、仮数部の桁数に依って常に一定の範囲内の相対誤差で任意の実数を近似できる。基本的には m -進小数展開を途中で切ったものであり、表現 $[\beta^m, \beta^{m+1}]$ を指数部により移動させる。結果、最大正数は $1.7976 \cdots \times 10^{308}$ ある。倍精度の計算機イプシロンは $\epsilon_M = \beta^{-52} = 2.2204 \cdots \times 10^{-16}$ なので誤差は 16 桁目から生じる。

記法 1.1.1. $\beta \geq 2$ を偶数とし、これを基数と呼ぶ。

補題 1.1.2 (m 進小数展開は連続). $m \geq 2$ を自然数とする。 m を離散空間とし、無限積空間 $m^{\mathbb{N}}$ からの写像 $e_m : m^{\mathbb{N}} \rightarrow [0, 1]$ を m 進小数展開

$$e_m((x_n)) := \frac{1}{m} \sum_{n=0}^{\infty} \frac{x_n}{m^n}$$

で定める。この時、 e_m は連続である。

[証明].

(1) 自然数 $n \in \mathbb{N}$ に対し、 n 番目以降を切断する写像

$$\begin{array}{ccc} q_n : m^{\mathbb{N}} & \longrightarrow & m^n \\ \downarrow & & \downarrow \\ (x_l) & \longmapsto & (x_0, \dots, x_{n-1}) \end{array}$$

は連続である。実際、任意の $(x_0, \dots, x_{n-1}) \in m^n$ に対して、逆像は $q_n^{-1}(x_0, \dots, x_{n-1}) = \{x_0\} \times \{x_1\} \times \cdots \times \{x_{n-1}\} \times m \times m \times \cdots \in \mathcal{U}$ は $m^{\mathbb{N}}$ の開集合系の基底である。よって、これを用いて $V_n(x) := q_n^{-1}(q_n(x))$ とおくと、これは x の開近傍である。

(2)

$$e_m(V_n(x)) = \left[\frac{1}{m} \sum_{l=0}^{n-1} \frac{x_l}{m^l}, \frac{1}{m} \sum_{l=0}^{n-1} \frac{x_l}{m^l} + \frac{1}{m^n} \right]$$

である。左端は $x(i) = 0$ ($i \geq l$) となる数列、右端は $x(i) = m - 1$ ($i \geq l$) となる数列で、全体として閉区間で、直径は $\frac{1}{m^n}$ である。

(3) 任意の実数 $r > 0$ に対し、 $m^n r > 1$ を満たすように n を取ると、

$$\forall x \in m^{\mathbb{N}}, e_m(V_n(x)) \subset U_r(e_m(x))$$

である。よって、任意の $e_m(x) \in [0, 1]$ の任意の開近傍の基本系の逆像は、その下に開集合が見つかる訳だから、命題??より、 $e_m : m^{\mathbb{N}} \rightarrow \mathbb{R}$ は連続。

要諦 1.1.3. 書籍の方では $\{V_n(x)\}_{n \in \mathbb{N}}$ が開近傍の基本系であることを確認せずに証明しているの、実はやばくないか？それにしても、射影 pr_i が連続写像になるだけでなく、その拡張(?) q_n も連続写像になり、 $n \in \mathbb{N}$ の調節によって自由にピントを合わせることができる、という道具立てを使う。

補題 1.1.4 (β -adic decimal expansion). $\beta \geq 2$ を偶数とする。

- (1) 任意の実数 $x \in \mathbb{R}_{\neq 0}$ に対して、整数 $m \in \mathbb{Z}$ と列 $(d_i)_{i \in \mathbb{N}} : \mathbb{N} \rightarrow \beta$, $b_0 \neq 0$ が存在し、

$$x = \pm \left(\frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \frac{d_2}{\beta^2} + \cdots \right) \times \beta^m$$

を満たす。

- (2) 各 $m \in \mathbb{N}$ について、こうして定めた写像

$$\begin{array}{ccc} E_m : \beta^{\mathbb{N}} & \longrightarrow & [\beta^m, \beta^{m+1}] \\ \downarrow & & \downarrow \\ (d_i)_{i \in \mathbb{N}} & \longmapsto & x = \frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \frac{d_2}{\beta^2} + \cdots \end{array}$$

は全射であるが、 \mathbb{Z} との共通部分において単射ではない。

[証明]。

- (1) 算譜を構成する。まず x の整数部分 $\lfloor x \rfloor$ を β -進数で表し、これを (d_0, d_1, \dots, d_r) とする。続いて全射 $e : \beta^{\mathbb{N}} \rightarrow [0, 1]$ を用いて、小数部分を繋げれば良い。
- (2) $x \in \mathbb{Z}$ の場合、上の構成法の他に、 $x - 1$ を整数部分として、小数部分を用いて 1 を作ることで、対応する $(d_i)_{i \in \mathbb{N}}$ が得られる。 B_m の非単射性は本質的にはこれのみにより、 B_m の任意の長さ 1 の閉区間への制限は全単射である。

要諦 1.1.5. 要は、 B_m の非単射性は、無限和であることにより、有限和の範囲では各 $1/\beta^i$ は基底だから、表示は一意的になる。

定義 1.1.6 (floating point number). 合成

$$\begin{array}{ccc} E_{m,n} := E_m \circ q_n : \beta^{n+1} & \longrightarrow & [\beta^m, \beta^m(\beta - \beta^{-n})] \\ \downarrow & & \downarrow \\ (d_i)_{i \in n+1} & \longmapsto & \left(\underbrace{\frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \cdots + \frac{d_n}{\beta^n}}_{=: \alpha} \right) \times \beta^m \end{array}$$

は単射だが、もはや全射ではない。 m が大きいほど間隙は大きくなる。 $\alpha \in [1, \beta]$ を仮数部、 m を指数部という。さらに、正数 $L, U \in \mathbb{Z}_{>0}$ に対して、

$$\mathcal{F}(\beta, n, L, U) = \{0\} \cup \bigcup_{i=-L}^U \pm \text{Im } E_{i,n}$$

を β 進 $n+1$ 桁の浮動小数点数系と呼ぶ。IEEE754-1985 規格に拠ると、 $(\beta, n, L, U) = (2, 23, 126, 127)$ を単精度^{†1}、 $(\beta, n, L, U) = (2, 52, 1022, 1023)$ を倍精度^{†2}という。^{†3}

注 1.1.7 (ケチ表現). 2 進法で正規化をすると、最上位ビット d_0 は常に 1 になるので、これを表さず常に 1 があるものとみなす省略が可能で、省略した表現をケチ表現などと言う。この省略を使うと、仮数部に割り当てたビット数が n であれば、有効桁数は $n+1$ となる。

^{†1} 符号 1bit, 指数 L, U で 8bit, 仮数部 n が 23bit.

^{†2} 指数部には 11bit 使っている。符号 1bit を除いて仮数部 52bit.

^{†3} IEEE754-2008 では binary32, binary64 と改名されている。

定義 1.1.8 (round, rounding error, underflow, overflow).

(1) 浮動小数点数系 \mathcal{F} に対して, 写像

$$\begin{array}{ccc} \text{round} : \mathbb{R} & \longrightarrow & \mathcal{F} \\ \psi & & \psi \\ \tilde{x} & \longmapsto & x \end{array}$$

を1つ定め (丸めの方法と呼ぶ), それに沿って値を対応させることを丸めるといふ. その時に生じる誤差 $d(\tilde{x}, x)$ を丸め誤差という.

(2) 代表的な丸めの方法には, 最近点への丸めと切り捨てがある.

(3) $0 \in \mathcal{F}$ に丸められてしまうことをアンダーフロー, $\max \mathcal{F}$ を超えてしまうことをオーバーフローという.

例 1.1.9 (最近偶数への丸め). デフォルトの丸め関数であるから, この関数を指して「最近数への丸め」ともいう. 基本的には $\tilde{x} := \max_{y \in \mathcal{F}} |y - \tilde{x}|$ と定めるが, \max 関数が定まらない場合は,

$$x := \begin{cases} \sigma \cdot \left(\frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \cdots + \frac{d_n}{\beta^n} \right) \cdot \beta^m & d_n \text{ が偶数のとき} \\ \sigma \cdot \left(\frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \cdots + \frac{d_n + 1}{\beta^n} \right) \cdot \beta^m & d_n \text{ が奇数のとき} \end{cases}$$

命題 1.1.10 (相対誤差: machine epsilon). 任意の $\tilde{x} \in \mathbb{R}$ と $x := \text{round}(\tilde{x})$ について,

$$\left| \frac{\tilde{x} - x}{\tilde{x}} \right| \leq \begin{cases} \frac{1}{2} \beta^{-n}, & \text{最近点への丸め} \\ \beta^{-n}, & \text{切り捨て} \end{cases}$$

が成り立つ. 特徴量 $\epsilon_M := \beta^{-n} = \min(\mathcal{F} \setminus (-\infty, 1])$ を計算機イブシロンという.

例 1.1.11 (倍精度の計算機イブシロン). 倍精度の計算機イブシロンは $\epsilon_M = \beta^{-52} = 2.2204 \cdots \times 10^{-16}$ となる.

実験事実 1.1.12 (IEEE754 の実装: denormalized number, NaN).

$$x_{\min} := 2^{-1022} = \min_{y \in \mathcal{F}(2, 52, 1022, 1023)} |y|$$

とする. 正規化数の仮数部は必ず $1 \leq \alpha \leq \beta$ をみたしたが, 非正規化数

$$\mathbb{Y} := \left\{ \pm \left(\frac{d_0}{2^0} + \frac{d_1}{2^1} + \cdots + \frac{d_{51}}{2^{51}} \right) \cdot 2^{-1023} \mid d_i \in 2 \ (i \in 52) \right\} \subset [0, x_{\min})$$

が定義されている. よって, $|\tilde{x}| \leq \frac{1}{2} y_{\min}$ はアンダーフローの可能性がある.

また, $0/0, \sqrt{-1}$ として NaN (Not a Number) が定義されており, ∞ として Inf が定義されている.

歴史 1.1.13 (IBM 方式). IBM 方式は, IBM 社が System/360 で導入し, 以後同社の標準として System/370 などのメインフレームで使った方式である. 指数部が2の幕ではなく16の幕を表すという特徴がある. この方式は, より大きな範囲を少ないビット数の指数部で表すことができ, そのぶんビットを仮数部の桁数に使うことで精度も確保できるように一見思える. しかし, 仮数部にケチ表現を使うことができず, さらに指数部の変化の前後で, 仮数部の LSB が表現する値の刻み幅が16倍変化するため, 2べきの場合に比べて最悪の場合には2進で3ビット分の精度が損なわれるため, 一般には大成功であったと評された System/360 の設計において良くなかった点の一つとして挙げられる.

1.2 演算

遠すぎる数の和, 近すぎる数の差には注意.

整数演算と同じ操作で処理が済む固定小数点と違い, 通常の整数演算命令を使って実装すると, 多くの命令と時間が必要になる. 処理の軽減のため, 演算にはハードウェアで実装した FPU などのコプロセッサを用い, 現在のマイクロプロセッサなどの多くでは内蔵されていることが多い.

例 1.2.1 (情報落ち / loss of trailing digit). 絶対値の大きさが極端に違う 2 数を足す際に吸収律が成り立ってしまうこと. Basel 級数 $\sum_{n=1}^{\infty} \frac{1}{n^2}$ の計算は末尾から足すと精度が出るが, 頭から足すと途中の時点で停止する.

例 1.2.2 (桁落ち / loss of significance). 値の近い 2 数の減算をすると, 有効桁数が大きく減る. 二次方程式の解の小さい方を求める $-b - \sqrt{b^2 - 4ac}$ ときに影響が大きいので, もう一方は解と係数の関係 $x_1 x_2 = c$ で計算する.

1.3 連立一次方程式

- (1) 消去法: 数値計算誤差がない状況においては必ず有限停止する算譜である.
- (2) 反復法: 漸化式によって近似解の列を生成する手法の総称であり, 本質的に有限停止はしない.
- (3) 共役勾配法: 二つの性質を併せ持つ.

関連する問題として最小二乗問題がある. これは, A が縦長の行列の場合には, 連立一次方程式 $Ax = b$ は一般には解を持たず, $\|Ax - b\|_2$ を最小にする x を求めることを考えることとなる.

1.4 固有値問題

- (1) Hermite 行列: 固有値が全て実数となることを利用した専用の算譜がある.
- (2)

関連する問題として特異値問題があり, 特異値は $A^T A, AA^T$ の固有値の平方根に等しいから理論的には問題にならないが, 数値計算手法としては特別な手法が考えられる.

第 2 章

行列とその分解と変換

2.1 連立一次方程式と Cramer の公式から見た線型計算の世界

計算量の世界：三角行列への注目

Cramer の公式は行列式計算を含むため、乗除算の階数は階乗のオーダーで増えていく。これは最悪な種類の計算量である。もはや計算速度があと 10^n 倍速くなろうと意味をなさない算譜である。また $\det(\epsilon A) = \epsilon^n \det A$ でもあるから、さらに意味を持たない。よって、 \det を用いた正則性の判断は実は難しい。こうして、三角行列への注目が誘導される。

また、逆行列が転置によって瞬時に求まることと、数値的に安定である直交行列が定める変換も重用される。

2.2 行列の標準形

なるべく素な行列で、また対角成分付近に密集していることが好ましい。

定義 2.2.1. 対角行列を一般化したクラスが次の通り。

- (1) 3 重対角行列： $|i - j| > 1 \Rightarrow a_{ij} = 0$ 。副対角成分が非零のとき、既約であるという。
- (2) w -帯行列： $|i - j| > w \Rightarrow a_{ij} = 0$ 。 $2w + 1$ を帯幅という。

$|i - j| = 1$ を満たす成分 a_{ij} を副対角成分という。

非対称性を考慮したクラスは次の通り。

- (3) 下三角行列： $i < j \Rightarrow a_{ij} = 0$ 。対角成分が全て 1 であるとき、単位的であるという。対角成分が全て 0 であるとき狭義であるという。
- (4) Hessenberg 行列： $i > j + 1 \Rightarrow a_{ij} = 0$ 。

2 つ合わせたクラスは次の通り。

- (5) $i < j \vee i > j + 1 \Rightarrow a_{ij} = 0$ 。

2.3 自己共役行列の性質

- Hermite 行列は、内積 $\langle \cdot, \cdot \rangle : V \times W \rightarrow K$ について定まる随伴関手 $*$: $\text{Hom}(V, W) \rightarrow \text{Hom}(W^*, V^*)$ への対応であり、Hermite 性は $\langle Ax, x \rangle = \overline{\langle x, Ax \rangle} = \langle x, Ax \rangle$ より、 $\langle Ax, x \rangle = \exists \lambda \langle x, x \rangle \in \mathbb{R}$ より、固有値は実数になる。
- 正則行列による共役変換を相似変換とよび、ある対角行列と相似であることを対角化可能という。Hermite 行列は、固有ベクトルの中から正規直交基底を選び出せるので、ユニタリー行列による相似変換で対角化する。

2.3.1 定義と特徴付け

定義 2.3.1 ((semi-)positive definite).

- (1) Hermite 行列が正定値であるとは, $\forall x \in \mathbb{C}^n \setminus \{0\} \langle Ax, x \rangle > 0$. 固有値が全て正であることと同値.
- (2) Hermite 行列が半正定値であるとは, $\forall x \in \mathbb{C}^n \setminus \{0\} \langle Ax, x \rangle \geq 0$. 固有値が全て非負であることと同値.

補題 2.3.2. Hermite 行列 A に対して, $B := AA^*$ と定める.

- (1) B は半正定値である.
- (2) B が正定値であることと, A が正則であることは同値.

定義 2.3.3 (submatrix, principal minor, leading principal minor). 行列 $(a_{ij})_{(i,j) \in I \times J}$ について,

- (1) ある集合 $\Lambda := \{i_1 < \dots < i_k\} \subset I \cap J$ を用いて, $(a_{i_\lambda j_{\lambda'}})_{\lambda, \lambda' \in \Lambda}$ と表せる行列を Λ -部分行列または小行列と呼ぶことにしよう. この行列式を k 次主小行列式という.
- (2) k 次の首座小行列式とは, $\Lambda = [k]$ である場合をいう.

補題 2.3.4 (半正定値性の特徴付け). 正方行列 $A \in M_n(\mathbb{C})$ について, 次の3条件は同値.

- (1) A は半正定値である: $\forall x \in \mathbb{C}^n \setminus \{0\} \langle Ax, x \rangle \geq 0$.
- (2) 任意の主小行列式が非負である.
- (3) 任意の固有値が全て非負である.

2.3.2 一般化

定義 2.3.5 (P-matrix (positive)). 一般の正方行列 $P \in M_n(\mathbb{C})$ について, 任意の主小行列式が非負である行列を P_0 -行列という. 任意の任意の主小行列式が正である行列を P -行列といい, P_0 -行列全体の集合は P -行列全体の集合の閉包となる.

補題 2.3.6. P -行列かつ Z -行列ならば, M -行列である.

2.4 優対角行列と既約行列

対角成分に非零要素が集積した疎行列を捉えるための概念を導入する. 優対角行列が可逆であるためには, 狭義性と既約性が必要.

2.4.1 定義

定義 2.4.1 ((strictly) diagonally dominant matrix, reducible).

- (1) 各 i 行目について, 対角要素の絶対値が非対角要素の絶対値の和よりも大きい行列 A を, 行方向の優対角行列という:

$$|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}|.$$
- (2) 真の不等号が成り立つ場合, 行列 A を狭義優対角行列という.
- (3) 正の対角成分をもつ対角行列 $D = \text{diag}(d_i) \ d_i > 0$ で, AD が狭義優対角行列となるものが存在するとき, A を一般化狭義優対角行列という: $|a_{ii}|d_i \geq \sum_{j=1, j \neq i}^n |a_{ij}|d_j$.
- (4) A^\top が (一般化・狭義) 優対角行列であるとき, A は列方向に (一般化・狭義) 優対角であるという.
- (5) $n \geq 2$ のとき, 空でない真部分集合 $J \subset \{1, \dots, n\}$ であって, $\forall i \in J \ \forall j \notin J \ a_{ij} = 0$ が成り立つとき, 行列 A を可約という.
- (6) 可約でない行列を既約という. $n = 1$ の時は, 零行列のみを可約とする.

例 2.4.2. 行列 $\begin{pmatrix} 2 & -1 & 0 \\ 0 & 2 & -2 \\ 0 & 2 & -2 \end{pmatrix}$ は優対角行列だが，一般化狭義優対角行列ではなく，非正則である．

2.4.2 M -行列

M -行列は，(とくにラプラシアンのような最小・最大原理の存在する) 微分作用素の離散化の結果として自然に得られる行列 A を捉えたものである．これの逆 A^{-1} は積分の離散化になっていると考えられることが inverse positivity に対応する．

定義 2.4.3.

- (1) 非対角要素が全て非正である行列を Z -行列という．
- (2) 正則行列 $A = (a_{ij}) \in \text{GL}_n(\mathbb{R})$ が M -行列であるとは，非対角要素が全て非正であり (Z -行列)，逆行列 A^{-1} が非負であるもの (Inverse-positivity) をいう．
- (3) 行列 $A = (a_{ij}) \in M_n(\mathbb{C})$ が定める行列 $(\alpha_{ij}) \in M_n(\mathbb{R})$, $\alpha_{ij} := \begin{cases} |a_{ij}|, & i = j, \\ -|a_{ij}|, & i \neq j. \end{cases}$ を比較行列という．
- (4) 比較行列が M -行列になるような行列 $A = (a_{ij}) \in M_n(\mathbb{C})$ を H -行列という．

補題 2.4.4. M -行列の対角成分は正になる．

補題 2.4.5. 次の2条件は同値．

- (1) A は H 行列である．
- (2) A は一般化狭義優対角行列である．

系 2.4.6. 行方向に一般化狭義優対角であることと，列方向に一般化狭義優対角であることは同値．

2.5 特異値分解

Autonne (1915) が導入した．スペクトル定理の一般化とみれる．

定理 2.5.1 (singular value decomposition (実の場合)). 任意の $A \in M_{m,n}(\mathbb{R})$ について，直交行列 $U \in O_m(\mathbb{R}), V \in O_n(\mathbb{R})$ が存在して， $r := \text{rank}(A)$ とすると

$$A = U\Sigma V^{\top}, \quad \Sigma = \begin{bmatrix} D & O_{r,n-r} \\ O_{m-r,r} & O_{m-r,n-r} \end{bmatrix}, \quad D = \text{diag}(\sigma_1, \dots, \sigma_r), \quad (\sigma_1 \geq \dots \geq \sigma_r > 0)$$

と表せる．

要諦 2.5.2. A の定める線型変換 $\mathbb{R}^n \rightarrow \mathbb{R}^m$ は， U, V の列ベクトル u_1, \dots, u_m と v_1, \dots, v_n について簡単な形で表せる：
 $Av_i = \sigma_i u_i$ ．

系 2.5.3.

- (1) $A^{\top}A = V\Sigma^{\top}\Sigma V$, $AA^{\top} = U\Sigma\Sigma^{\top}U^{\top}$ ．
- (2) A の特異値は $A^{\top}A$ の非零固有値の正の平方根に等しい．
- (3) V の列ベクトルは $A^{\top}A$ の固有ベクトルであり， U の列ベクトルは AA^{\top} の固有ベクトルである．

定理 2.5.4 (一般化逆行列の公式). 任意の $A \in M_{m,n}(\mathbb{R})$ について，特異値分解を $A = U\Sigma V^{\top}$ とする．

2.6 LU 分解と Cholesky 分解

枢軸選択なしの Gauss 消去法は、係数行列 A を LU 分解していることに相当する。次の数学的事実が、Gauss 消去法の背後にある。

定義 2.6.1.

- (1) $A \in M_n(\mathbb{R})$ を、単位下三角行列 L と上三角行列 U の積 $A = LU$ に分解することを LU 分解という。
- (2) A が自己共役かつ正定値であるとき、下三角行列 C を用いて $A = CC^T$ と分解できる。この場合を特に **Cholesky** 分解という。
- (3) A が自己共役かつ正定値でないとき、単位下三角行列 L と対角行列 D を用いて $A = LDL^T$ と分解できる。この場合を特に LDL^T 分解という。

2.7 基本直交変換

直交変換を繰り返し適用して、特定の良い形の疎行列に収束させる技法が基本となる。

2.7.1 Householder 変換

2.7.2 Givens 変換

2.8 三角化 / QR 分解

2.9 Hessenberg 化 / 3 重対角化

2.10 Schur 分解

固有値問題の数値解法では、ユニタリ変換を反復適用して上三角行列に収束させることが多いが、その背景には次の数学的事実がある。

定理 2.10.1 (Schur decomposition). 任意の複素行列 $A \in M_n(\mathbb{C})$ は、ユニタリ行列 U と上三角行列 S を用いて $A = USU^*$ と分解できる。なお、一意性は保証されない。

2.11 Moore-Penrose の一般化逆行列

2.12 行列のノルム

2.12.1 p -ノルム

定義 2.12.1 (p -norm). $M_{m,n}(\mathbb{R})$ 上の p -ノルムとは、 $\mathbb{R}^n, \mathbb{R}^m$ の p -ノルムが定める作用素ノルムをいう：

$$\|A\|_p := \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \inf \{c > 0 \mid \forall x \in \mathbb{R}^n \quad \|Ax\|_p \leq c\|x\|_p\}.$$

補題 2.12.2. A の最大の特異値を $\sigma_1(A)$ で表す。

- (1) $\|A\|_1 = \max_{j \in [n]} \sum_{i=1}^m |a_{ij}|$. 列ベクトルの 1-ノルムの ∞ -ノルムになる。

- (2) $\|A\|_\infty = \max_{i \in [m]} \sum_{j=1}^n |a_{ij}|$. 行ベクトルの1-ノルムの ∞ -ノルムになる .
- (3) $\|A\|_2 = \sigma_1(A)$. これをスペクトルノルムともいう .

2.12.2 Frobenius ノルム

定義 2.12.3. $A \in \mathbb{R}^{m \times n}$ とみたときの2-ノルムを **Frobenius ノルム** という :

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2} = \left(\sum_{i=1}^r \sigma_i(A)^2 \right)^{1/2} .$$

2.12.3 スペクトル半径

命題 2.12.4 (スペクトル半径公式). $\forall p \in [1, \infty] \quad \lim_{k \rightarrow \infty} (\|A^k\|_p)^{1/k} = \rho(A) \leq \|A\|_p$.

補題 2.12.5.

- (1) 三角不等式 : 各要素の絶対値を要素とする行列を $|A| := (|a_{ij}|)$ とすると , $\rho(A) \leq \rho(|A|)$.
- (2) $\forall A, B \in M_n(\mathbb{R}) \quad \rho(AB) = \rho(BA)$.

2.12.4 条件数

定義 2.12.6 (condition number). 正則行列 $A \in \text{GL}_n(\mathbb{R})$ の p -ノルムに関する条件数とは ,

$$\kappa_p := \|A\|_p \|A^{-1}\|_p$$

をいう .

補題 2.12.7. 正則行列 $A \in \text{GL}_n(\mathbb{R})$ について ,

- (1) $\kappa_2 = \frac{\sigma_1(A)}{\sigma_n(A)}$.
- (2) $\kappa_2^{-1} = \inf_{S \in M_n(\mathbb{R}) \setminus \text{GL}_n(\mathbb{R})} \frac{\|S - A\|_2}{\|A\|_2}$.

第 3 章

代数：連立一次方程式の消去法

現実問題の数値モデル化においては、偏微分方程式の近似を経て、問題が連立一次方程式に帰着される場合が多い。素行列でも密行列でも対応できる解析的方法をまずは考えたい。このための靈性を湛えた数学の御堂が代数である。しかし、密行列の場合は、 $O(N^2)$ の記憶容量、 $O(N^3)$ の計算量が大きな欠点となる。

こうして、数値手法・算譜の数理的な扱いに集中するのがどうやらこの本の哲学である。ソフトウェアの適切な選択ということである。

3.1 Gauss の消去法

第 4 章

解析：連立一次方程式の反復法

未知数の多い連立一次方程式系で、係数行列はよく疎行列となる。このような連立一次方程式系は、Gauss の消去法などの直接法と共に、反復法も有効である。つまり、収束を狙って構成する技法である。そのための収束の議論のためには、ノルムが大事になる。

4.1 ベクトルのノルム

収束を距離ではなくノルムの時点から議論することで初めて解析が可能になる問題も多いということである。

ノルムを定めると、Hausdorff 位相が自然に定まる。位相に正確に対応するのが seminorm で、これが norm である（正値性を満たす）ことが分離可能性 = Hausdorff 性に同値。

4.1.1 p -ノルム

定義 4.1.1 (norm). k を絶対値を備えた体とする。 k^n 上の実数値関数 $\|\cdot\| : k^n \rightarrow \mathbb{R}$ が次の 3 条件を満たす時、ノルムであるという。

- (1) (positivity) $\forall x \in k^n \ x \neq 0 \Rightarrow \|x\| > 0$ ^{†1}
- (2) (linearity) $\forall \alpha \in k \ \forall x \in k^n \ \|\alpha x\| = |\alpha| \|x\|$.
- (3) (triangle inequality) $\forall x, y \in k^n \ \|x + y\| \leq \|x\| + \|y\|$.

命題 4.1.2 (p-norm).

- (1) $p \in [1, \infty]$ について定まる次の関数 $k^n \rightarrow \mathbb{R}$ はノルムである：

$$\|x\|_p := \begin{cases} \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, & p \in [1, \infty), \\ \max_{1 \leq i \leq n} |x_i|, & p = \infty. \end{cases}$$

- (2) $\forall p < q \in [1, \infty] \ \forall x \in k^n \ \|x\|_p \leq \|x\|_q \leq n^{\frac{1}{p} - \frac{1}{q}} \|x\|_p$.

4.1.2 ノルムと位相

有限次元ノルム空間において、任意のノルムは同値になる。セミノルムが同値ならば、同じ位相を定める。したがって、普通の p -ノルムについての開球は形は違えど本質的には変わらない。関数空間では任意の p -ノルムは違い、無数の位相の定め方がある。

命題 4.1.3 (continuousness of norm). $\|\cdot\| : k^n \rightarrow \mathbb{R}$ をノルムとする。

^{†1} この条件が落ちると seminorm という。全ての seminorm は位相を定め、その位相が Hausdorff であることとノルムであることが同値になる。

$$(1) \forall x, y \in k^n \quad \|x\| - \|y\| \leq \|x - y\|.$$

(2) ノルムは連続である。

命題 4.1.4 (equivalence of norm). k^n の任意の2つのノルム $\|\cdot\|, \|\cdot\|'$ について, $\exists C, C' > 0 \quad \forall x \in k^n \quad C\|x\|' \leq \|x\| \leq C'\|x\|'$.

4.2 行列のノルム

一般には行列のノルムは、整合性条件 $\forall x \in \mathbb{R}^n \quad \|Ax\| \leq \|A\|\|x\|$ を満たすものとして定める。ここでは有限次元を考えているので、どう構成しようと整合性条件を満たす行列ノルムは一意であるから、今回はベクトルへの作用を解して定める。 $A \in M_n(k)$ について、 $f_A: x \mapsto Ax \mapsto \|\cdot\| \|Ax\|$ は連続写像であるから、有界閉集合 $\partial\Delta \subset k^n$ 上で最大値が定まる。これを行列のノルムとすれば良い。これは、 A の「方向」にクリティカルヒットした時の拡大率である。^a

^a Jacobian よりも一般的な概念なのか！？

定義 4.2.1. k^n のノルム $\|\cdot\|$ に従属する $k^{n \times n}$ の行列ノルムを、

$$\|A\| := \max_{x \in k^n \setminus \{0\}} \frac{\|Ax\|}{\|x\|}$$

で定める。

補題 4.2.2 (行列ノルムの性質). $A, B \in M_n(k)$ を任意の行列とする。

$$(1) \forall x \in \mathbb{R}^n \quad \|Ax\| \leq \|A\|\|x\|.$$

$$(2) \|AB\| \leq \|A\|\|B\|.$$

$$(3) \|I\| = 1.$$

(4) 2つの同値なベクトルノルムに対して、それぞれに従属する行列ノルムも同値になる。

[証明]。

(1) $\|ABx\| \leq \|A\|\|Bx\| \leq \|A\|\|B\|\|x\|$ であるため、ノルムの定義から。

(2) 任意の $x \in k^n$ について、 $\frac{\|Ix\|}{\|x\|} = 1$ であるため。

(3) 任意の $x \in k^n$ について、 C'

■

4.2.1 p-行列ノルム

定義 4.2.3 (spectral radius). A の固有値の最大値 $\rho(A)$ をスペクトル半径という。

命題 4.2.4 (p-行列ノルムの特徴付け).

$$(1) \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|. \text{ 列ベクトルの } l^1 \text{ ノルムの最大値.}$$

$$(2) \|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|. \text{ 行ベクトルの } l^1 \text{ ノルムの最大値.}$$

$$(3) \|A\|_2 = \sqrt{\rho(A^*A)} = \sqrt{\rho AA^*}.$$

4.2.2 複素行列の行列ノルム

命題 4.2.5. $A \in \mathbb{C}^{n \times n}$ とする。

(1) 任意の行列ノルムについて、 $\rho(A) \leq \|A\|$ 。

(2) 任意の行列 $A \in \mathbb{C}^{n \times n}$ と $\epsilon > 0$ に対して, $\|A\| \leq \rho(A) + \epsilon$ を満たす行列ノルム $\|\cdot\| = \|\cdot\|_{A,\epsilon}$ が存在する.

4.3 定常反復法

Stationary Iterative Method

stationary とは, 反復計算中に解ベクトル以外の変数を変化させない方法をいい, この手法が一番計算が速くなる. が, 収束性がアプリケーションや境界条件の影響を受けやすいため, 前処理 (preconditioning) が必要になる. 非定常的な方法は Krylov 部分空間への写像を基底として使用するので, その数理手法から Krylov 部分空間法ともいう.

定義 4.3.1 (stationary iterative method). 連立一次方程式系 $Ax = b$ について, 係数行列を正則部分 M を用いて $A = M - N$ と分解し,

$$x = \underbrace{M^{-1}N}_{=:H}x + \underbrace{M^{-1}b}_{=:c}$$

と変形し, 解ベクトルの列 $(x^{(k)})_{k \in \mathbb{N}}$ を $x^{(k+1)} := Hx^{(k)} + c$ と更新していく手法を, H, c がループ数 k に依らないために定常反復法という. 以下, A の対角成分を D , 下・上三角成分を E, F とする

Jacobi method $M := D, N := -(E + F)$ とする.

Gauss-Seidel method $M := D + E, N := -F$ とする.^{†2}

Successive Over-Relaxation method 緩和係数 $\omega > 0$ を用いて, $M := \frac{1}{\omega}(D + \omega E), N := \frac{1}{\omega}((1 - \omega)D - \omega F)$ と定める. $\omega = 1$ のときが Gauss-Seidel 法である.

定理 4.3.2 (定常反復法の収束条件). 任意の $x^{(0)} \in \mathbb{R}^n, b \in \mathbb{R}^n$ に対して定常反復法による反復列 $(x^{(k)})_{k \in \mathbb{N}}$ が $Ax = b$ の解 x に収束するための必要十分条件は, $\rho(H) < 1$ である.

系 4.3.3 (SOR 法の収束の必要条件). SOR 法は, $\omega \leq 0, \omega \geq 2$ のとき, 収束しない.

4.4 安定性と条件数

Gauss の消去法の数値的安定性の厄介さ

2次元においては連立方程式系は直線の交点だが, 2直線がほぼ直交しているならば, 解は b の摂動に対して極めて安定的である.

定義 4.4.1 (condition number).

$$\text{cond}(A) = \begin{cases} \|A\| \cdot \|A^{-1}\|, & A \text{ が正則のとき,} \\ \infty, & A \text{ が正則でないとき.} \end{cases}$$

特に行列ノルムとして p -ノルムを用いた場合, $\text{cond}_p(A)$ と表す.

4.4.1 後退誤差解析

視点を逆にする

丸め誤差の拡大を順方向で追うのは難しく, 極端な評価しか得られない.

von Neumann and Goldstein (1947) – 素朴な (前進) 誤差解析で Gauss の消去法を解析した. これを信じると 100 元の方程式を解くのは絶望的と考えられるが, 実際には楽々解けてしまう.

^{†2} Jacobi 法の大体 2 倍の速さで収束すること.

しかし逆に考えて，得られた解が，どのような摂動問題の解であるかを見ることで，精度保証が出来る．^a (E, c) を後退誤差という．結論は，(残念ながら) 残差が小さくても，誤差も小さいとは保証されない！ 係数行列の条件数が小さければ，誤差が小さくなることが保証される。

^a J.H.Wilkinson "Rounding errors in algebraic process", 1963

定理 4.4.2 (後退誤差の存在). Gauss の消去法の数値解を \tilde{x} とし，残差を $r = b - A\tilde{x}$ とする． $|r|$ が十分に小さいとき，

$$(A + E)\tilde{x} = b + c, \|E\| \leq \epsilon_M \alpha \|A\|, \|c\| \leq \epsilon_E \alpha \|b\|$$

を満たす $E \in M_n(\mathbb{R}), c \in \mathbb{R}^n, \alpha \in \mathbb{R}_{>0}$ が存在する．

定理 4.4.3. Gauss の消去法の数値解を \tilde{x} とし， $E \in M_n(\mathbb{R}), c \in \mathbb{R}^n$ を後退誤差とする．このとき， $\|E\| < \|A^{-1}\|^{-1}$ ならば，

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \frac{2\epsilon_M \alpha \cdot \text{cond}(A)}{1 - \epsilon_M \alpha \cdot \text{cond}(A)}.$$

^{†3} いくら ϵ_M が小さい優秀な浮動小数点系を使っても，係数行列の条件数が多いと，Gauss の消去法では良い数値解が得られる保証はない．

第 5 章

連立一次方程式と共役勾配法

Krylov 部分空間はロシアの応用数学者で海軍技術者であったアレクセイ・クリロフにちなんで名づけられた。これは数値線形代数において最も成功した手法の一つである。

第 6 章

非線形方程式

例えば $n \geq 2$ 次多項式など，非線形な方程式は，例 1.2.2 など，代数的な解析解は必ずしも数値計算の文脈で有用な手段とはならない．ここが形式科学としての純粋数学の大きな違いである．計算機のための形式と，人類のための形式とは，違う．

ということで議論は簡単に解析学に移る．代数が無効というわけではないが，有限な算譜の構成を諦めるべきである．目指すは，実数値連続関数で，一次関数の形で表せないものに対する，縮小列の構成法を考える理論である．近似解の構成算譜と解の存在の証明とは表裏一体である．

例 6.0.1 (単独非線型方程式).

- (1) $n \geq 2$ 次の代数方程式．代数的方法 = 指数演算と冪根の有限合成では $n \geq 5$ のとき算譜がない．
- (2) Kepler 方程式: $\beta = x - \alpha \sin x$.
- (3) 期間中一定とみなしたときの利回り方程式 (yield equation) .

6.1 C^k 級関数と Taylor の定理

定義 6.1.1 (境界での微分係数は連続補完する). $f \in C^1[a, b]$ について, $f'(a) := \lim_{x \rightarrow a+0} f'(x)$ と定める. $f(a)$ が定義できない場合もあるので, その場合も関係なく簡単に定義できる.

定理 6.1.2 (Taylor の定理: Bernoulli の剰余項の変形). $f \in C^k[a, b]$ と $x, y \in [a, b]$ について,

$$f(x) = \sum_{m=0}^{k-1} \frac{f^{(m)}(y)}{m!} (x-y)^m + \frac{1}{(k-1)!} \int_0^1 (1-s)^{k-1} (x-y)^k f^{(k)}(y+s(x-y)) ds.$$

6.2 二分法

問題 6.2.1. 連続関数 $f: \mathbb{R} \supset I \rightarrow \mathbb{R}$ の方程式 $f(a) = 0$ の解 $a \in I$ を求める.

議論 6.2.2 (中間値の定理の証明抽出). 任意の $f(\alpha_0)f(\beta_0) < 0$ を満たす区間 $[\alpha_0, \beta_0]$ で, f はただ一つの解を持つならば, 各 $x_k := \frac{1}{2}(\alpha_k + \beta_k)$ について

$$[\alpha_{k+1}, \beta_{k+1}] := \begin{cases} [\alpha_k, x_k], & f(x_k)f(\beta_k) \geq 0 \text{ のとき,} \\ [x_k, \beta_k], & f(x_k)f(\beta_k) < 0 \text{ のとき.} \end{cases}$$

と更新すれば,

$$|x_k - a| \leq \beta_{k+1} - \alpha_{k+1} = \left(\frac{1}{2}\right)^{k+1} (\beta_0 - \alpha_0)$$

が保証される.

6.3 反復法と不動点定理

反復法の本質としては、縮小写像を構成し、そのただ一つの不動点として解を構成する。

定義 6.3.1 (iteration method). 開区間上の関数 $f: I \rightarrow \mathbb{R}$ が $f(x_0) > 0$ を満たすとする。

Newton's method

$f \in C^1(I)$ のとき, $x_{k+1} := x_k - \frac{f(x_k)}{f'(x_k)}$ と定めると近似列 (x_k) を得る。

simplified Newton's method

f が連続で $f'(x_0)$ の値が定まっているとき, $x_{k+1} := x_k - \frac{f(x_k)}{f'(x_0)}$ と定めると近似列 (x_k) を得る。^{†1}

secant method

$f(x_1) > 0$ も満たすとする。 $x_{k+1} := x_k - \frac{x_k - x_{k+1}}{f(x_k) - f(x_{k+1})} f(x_k)$ と定めると近似列 (x_k) を得る。

inverse linear interpolation method

$f(x_1) > 0$ も満たすとする。 $x_{k+1} := x_k - \frac{x_k - x_0}{f(x_k) - f(x_0)} f(x_k)$ と定めると近似列 (x_k) を得る。これを線型逆補間法という。

relaxation iteration method

β を任意の定数として, $x_{k+1} := x_k - \beta f(x_k)$ とする。これを緩和反復法という。^{†2}

一般に、傾きを $\varphi(x_k)$ とすると, $g := \text{id} - \varphi f$ とおけば、求める f の解は $g(a) = a$ を満たす不動点である。 g の条件を考える。

定義 6.3.2 (Lipschitz continuity, contraction mapping). 関数 $g: I \rightarrow \mathbb{R}$ について,

(1) 閉区間 J に於て、次が成り立つとき J 上リプシッツ連続であるという：

$$\exists \lambda > 0 \quad \forall x, x' \in J \quad |g(x) - g(x')| \leq \lambda |x - x'|.$$

(2) Lipschitz 定数 λ が $\lambda \in (0, 1)$ を満たすとき, g を J 上縮小写像という。^{†3}

定理 6.3.3 (縮小写像の定理). 区間上の実数値連続関数 $g: I \rightarrow \mathbb{R}$ が、ある閉区間 J 上に $g(J) \subset J$ を満たし、この上で縮小写像 $g: J \rightarrow J$ を定めるとき (すなわち、Lipschitz 定数が $\lambda \in (0, 1)$ を満たすとき), 次が成り立つ：

(1) $\exists! a \in J \quad g(a) = a$.

(2) 任意の $x_0 \in J$ について、これが定める列 $(x_{k+1} := g(x_k))_{k \in \mathbb{N}}$ は、もちろん $\{x_k\} \subset J$ を満たし、

$$|x_k - a| \leq \frac{1}{1 - \lambda} \lambda^k |x_1 - x_0|$$

を満たす。特に、 $\lim_{k \rightarrow \infty} x_k = a$.

[証明].

Cauchy 列性と収束の評価 任意の $k > m$ について、

$$\begin{aligned} |x_k - x_m| &\leq |x_k - x_{k-1}| + |x_{k-1} - x_{k-2}| + \cdots + |x_{m-1} - x_m| \\ &\leq (\lambda + 1)|x_{k-1} - x_{k-2}| + \cdots + |x_{m-1} - x_m| \\ &\leq (\lambda^2 + \lambda + 1)|x_{k-2} - x_{k-3}| + \cdots + |x_{m-1} - x_m| \\ &\vdots \\ &\leq (\lambda^{k-m-1} + \cdots + 1)|x_{m-1} - x_m| = \frac{1 - \lambda^{k-m}}{1 - \lambda} |x_{m-1} - x_m| \\ &\leq \frac{1 - \lambda^{k-m}}{1 - \lambda} \lambda^m |x_1 - x_0| \left(\xrightarrow{k \rightarrow \infty} \frac{1}{1 - \lambda} \lambda^m |x_1 - x_0| \text{ より評価を得る} \right) \end{aligned}$$

^{†1} 1 回目だけ微分を計算した後、あとはこれで誤魔化せないか。

^{†2} うまくいくためには、微分係数である必要はなく、それを模倣する必要もない。ということで、大海へ飛び出す。

^{†3} $\lambda = 0$ の場合は定値写像である。

$$\leq \frac{1}{1-\lambda} \lambda^m |x_1 - x_0| \xrightarrow{m,k \rightarrow \infty} 0.$$

より, (x_k) は Cauchy 列だから, 極限 $a := \lim_{k \rightarrow \infty} x_k$ が存在する.

収束先が値域に入っていること $\forall_{k \in \mathbb{N}} x_k \in J =: [\alpha, \beta]$ なのであるから, どうやっても背理法から導ける. が, 特に, 例えば $a > \beta$ とすると, $a - \beta > 0$ より, $\exists_{N>0} \forall_{n \geq N} |x_n - a| = a - x_n < a - \beta$ であるが, これは $\beta < x_n$ を意味し, $\forall_{k \in \mathbb{N}} x_k \in J$ に矛盾. つまり, 収束先が閉区間 J から出ているのなら, 数列の十分先も J から出ている必要が生じてしまうため, 矛盾.

不動点について g は連続関数としたから, $g(a) = g(\lim_{k \rightarrow \infty} x_k) = \lim_{k \rightarrow \infty} g(x_k)$. 一意性については, $g(b) = b$ を不動点とすると, J 上での縮小写像性より $|g(a) - g(b)| = |a - b| \leq \lambda |a - b|$ であるが, $\lambda \in (0, 1)$ より, $a = b$ が必要.

■

命題 6.3.4. ただ一つの不動点 $g(a) = a$ を持つ関数 g は, a の近傍で C^1 級かつ $|g'(a)| < 1$ を満たすとする. このとき, ある $\delta > 0$ が存在して, $J := [a - \delta, a + \delta]$ に対し, $g|_J$ は連続な縮小写像である. 特に, Lipschitz 定数を $\lambda := \frac{1}{2}(1 + |g'(a)|)$ と取れる.

6.4 ベクトル値関数の微分

6.5 多変数の反復法

第 7 章

固有値問題

第 8 章

関数近似

第 9 章

補間と積分

第 10 章

常微分方程式の初期値問題

参考文献

- [1] [J. C. A. Barata, M. S. Hussein, The Moore-Penrose Pseudoinverse. A Tutorial Review of the Theory](#)
- [2] 齊藤宣一『数値解析』
- [3] 森正武『数値解析』
- [4] 杉原正顯・室田一雄『数値計算法の数理』
- [5] 杉原正顯・室田一雄『線型計算の数理』