

航空宇宙技術研究所資料

TECHNICAL MEMORANDUM OF NATIONAL AEROSPACE LABORATORY

TM-315

1 次方程式系の解法 III

——長方形列を対象とするもの——

福田正大・末松俊二

1976年10月

航空宇宙技術研究所
NATIONAL AEROSPACE LABORATORY

目 次

1. 緒 言	1
2. 数値テスト	1
3. ハウスホルダー法による最小自乗問題の解法	2
方法の基礎	2
計算方法とサブルーチンの説明	3
用 途	3
サブルーチンの使用方法	4
数値テスト及び結果	6
4. 直交化法による最小自乗問題の解法	7
方法の基礎	7
計算方法とサブルーチンの説明	8
用 途	8
サブルーチンの使用方法	9
数値テスト及び結果	12
5. Singular Value Decomposition	15
方法の基礎	15
計算方法とサブルーチンの説明	15
用 途	16
サブルーチンの使用方法	17
数値テスト及び結果	19
6. 結 語	22

1 次 方 程 式 系 の 解 法 III

—— 長方形列を対象とするもの ——*

福 田 正 大**、 末 松 俊 二**

概 要

本報告は、一連の一次方程式系を取り扱ったものうち、長方形列を対象としたものである。そのうち主として最小自乗問題の解法を取り扱い、特殊な場合のものとして、実正方形列の逆行列を計算するものがある。方法としては、ハウスホルダー変換を利用したもの、消去法型の直交化法を利用したもの、及び singular value decompositionを利用したものの3つがある。

1 章 緒 言

本報告は、「1次方程式系の解法 I」⁽¹⁾にも述べてあるように、科学技術計算用サブルーチンライブラリーの標準化ならびに充実化という目標に沿ったものの一つである。これは、計算センターにとって各種解法の知識を集積していくことが必要と考えるからである。このような研究は地味ではあるが計算センターの質を高める上で甚だ重要なことと思われる。

今回報告するものは、長方形列 A を係数行列にもつ overdetermined system

$$Ax = b$$

を最小自乗法によって解くものを主としている。行列 A としては稠密なものを想定しており、帯状をした行列に対する特殊化は行っていない。なお、各サブルーチンの仕様で、「1次方程式系の解法 II」で採用した「1次方程式系の解法 I」からの変更をここでも採用した。それは次の二点である。

- 1) 1次元のワーキングエリアも引き数に加える。こうすることによって引き数の並びは若干煩雑になるが、その代わりに無駄なエリアを取らなくてすむ。
- 2) 整合配列の引き数を各2次元アレイごとにつけるのではなく、行列として同じ行数をもつものについては同じ整合配列引き数を使用する。

(例) A が $m \times n$, B が $m \times r$, X が $n \times r$ という行列のとき、 A と B には同じ整合配列引き数 KA が X には異なる引き数 KX が使用されている。従って、サブルーチンの中での dimension 宣言は

$$\text{DIMENSION } A(KAN), B(KAR), \\ X(KX, R)$$

となっている。

「1次方程式系の解法 I, II, III」全体のサブルーチンの用いられ方については、本報告末尾の綴じ込みに図版 A としてフローチャートを載せてある。

2 章 数値テスト

ここでいう最小自乗問題とは、 m 行 n 列 ($m > n$) の行列 A と右辺のベクトル b に対して、

$$\|b - Ax\| = \min_x \|b - Ax\|$$

となる x を求めることである。複数個のベクトル b_i ($i=1, \dots, r$) について解くときには、右辺はそれらのベクトルより作られる $m \times r$ 行列 B として扱う。この解を求めるルーチンのテストとして次の4つのものを使用した。なお、 r は右辺にあるベクトルの個数である。又、各長方形列の singular value の最大値 σ_1 と最小値 σ_n の値を記しておく。ここで、singular value というのは $A^T A$ の固有値の正の平方根のところである。正值対称行列 $A^T A$ のスペクトルノルムに関する条件数は $(\sigma_1 / \sigma_n)^2$ になる。

(1) $m=6$ $n=5$ $r=2$

$$A = \begin{pmatrix} 36 & -630 & 3360 & -7560 & 7560 \\ -630 & 14700 & -88200 & 211680 & -220500 \\ 3360 & -88200 & 564480 & -1411200 & 1512000 \\ -7560 & 211680 & -1411200 & 3628800 & -3969000 \\ 7560 & -220500 & 1512000 & -3969000 & 4410000 \\ -2772 & 83160 & -582120 & 1552320 & -1746360 \end{pmatrix}$$

$$\begin{pmatrix} b_1 & b_2 \\ 463 & -4157 \\ -13860 & -17820 \end{pmatrix}$$

* 昭和51年8月11日 受付
** 計算センター

$$\begin{pmatrix} 97020 & 93555 \\ -258720 & -261800 \\ 291060 & 288288 \\ -116424 & -118944 \end{pmatrix}$$

$$\sigma_1 = 8.888 \times 10^6 \quad \sigma_5 = 1.892 \quad \sigma_1/\sigma_5 = 4.699 \times 10^6$$

この行列 A は、6 次ヒルベルト行列の逆行列の最初の 5 列より成っているものである。右辺のベクトル b_1 は、 $\|b_1 - Ax\| = 0$ となるものであり、 b_2 は A の各列に直交するベクトルを b_1 に加えたものである。従って最小自乗問題の解はどちらも同じである。

(2) $m=6 \quad n=5 \quad r=3$

$$\begin{pmatrix} -74 & 80 & 18 & -11 & -4 \\ 14 & -69 & 21 & 28 & 0 \\ 66 & -72 & -5 & 7 & 1 \\ -12 & 66 & -30 & -23 & 3 \\ 3 & 8 & -7 & -4 & 1 \\ 4 & -12 & 4 & 4 & 0 \end{pmatrix} \begin{pmatrix} b_1 & b_2 & b_3 \\ 51 & -56 & -5 \\ -61 & 52 & -9 \\ -56 & 764 & 708 \\ 69 & 4096 & 4165 \\ 10 & -13276 & -13266 \\ -12 & 8421 & 8409 \end{pmatrix}$$

$$\sigma_1 = 173.5 \quad \sigma_6 = 0.1599 \quad \sigma_1/\sigma_6 = 1.085 \times 10^8$$

b_1 は $\|b_1 - Ax\| = 0$ となるものであり、 b_2 は A の各列に直交しており ($A^T b_2 = 0$)、 b_3 は b_1 と b_2 を加えたものである。従って、 b_1 と b_3 に対する最小自乗問題の解は同じであり、 b_2 に対するそれは 0 ベクトルである。

(3) $m=8 \quad n=5 \quad r=3$

$$\begin{pmatrix} 22 & 10 & 2 & 3 & 7 \\ 14 & 7 & 10 & 0 & 8 \\ -1 & 13 & -1 & -11 & 3 \\ -3 & -2 & 13 & -2 & 4 \\ 9 & 8 & 1 & -2 & 4 \\ 9 & 1 & -7 & 5 & -1 \\ 2 & -6 & 6 & 5 & 1 \\ 4 & 5 & 0 & -2 & 2 \end{pmatrix} \begin{pmatrix} b_1 & b_2 & b_3 \\ -1 & 1 & 0 \\ 2 & -1 & 1 \\ 1 & 10 & 11 \\ 4 & 0 & 4 \\ 0 & -6 & -6 \\ -3 & 6 & 3 \\ 1 & 11 & 12 \\ 0 & -5 & -5 \end{pmatrix}$$

$$\sigma_1 = 35.33 \quad \sigma_5 = 0.0$$

この行列 A のランクは 3 なので、通常最小自乗問題としては解が一意的には定まらない。 $\|x\|$ を最小にするという条件を付ければ一意に定まる。 b_1, b_2, b_3 の構成は (2) と同じである。

(4) $m=7 \quad n=5 \quad r=3$

$$\begin{pmatrix} 2 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 2 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 2 & 4 & 1 \end{pmatrix} \begin{pmatrix} b_1 & b_2 & b_3 \\ 0 & 7 & 21 \\ 47 & 40 & 32 \\ 22 & 22 & 25 \\ 69 & 55 & 36 \\ -4 & 3 & 17 \\ 15 & 8 & -21 \\ 8 & 15 & 26 \end{pmatrix}$$

$$\sigma_1 = 8.242 \quad \sigma_5 = 1.0 \quad \sigma_1/\sigma_5 = 8.242$$

この最小自乗問題の解は 3 つとも同じである。

以上 4 つのテストをそれぞれテスト 1, テスト 2, テスト 3, テスト 4 と呼ぶことにする。これら 4 つの問題の正確な解と、残差の自乗 $\|b - Ax\|^2$ は次の通りである。

- (1) $x_1 = x_2$; (1, 1/2, 1/3, 1/4, 1/5)
 $\|b_2 - Ax_2\|^2 = 7.255 \times 10^7$
- (2) $x_1 = x_3$; (1, 2, -1, 3, -4)
 $\|b_2 - Ax_2\|^2 = \|b_3 - Ax_3\|^2 = 2.645 \times 10^8$
- (3) $x_1 = x_3$; (-1/12, 0, 1/4, -1/12, 1/12)
 $\|b_2 - Ax_2\|^2 = \|b_3 - Ax_3\|^2 = 320$
- (4) $x_1 = x_2 = x_3$; (5, 4, 3, 2, 1)
 $\|b_1 - Ax_1\|^2 = 4880 \quad \|b_2 - Ax_2\|^2 = 2577$
 $\|b_3 - Ax_3\|^2 = 1913$

結果の整理には、計算機より出力された値 \bar{x} をそのまま記し、さらに正確な解を x として

$$y = x - \bar{x}$$

$$z_i = (x_i - \bar{x}_i) / x_i$$

を計算して、 $a_m = \max |y_i|$, $a_f = \|y\|$, $s_m = \max |z_i|$, $s_f = \|z\|$ の値を求めた。さらに、

$$T_S = 26 \log 2 - \log \|b - Ax\|$$

$$T_D = 61 \log 2 - \log \|b - Ax\|$$

の 2 つの値も併せ記してある。なお、 z_i (相対誤差) の計算で、 $x_i = 0$ なるときには、分子の値を z_i として採用した。

なおテスト計算はすべて、FACOM 230-75 を使用して、FORTRAN-H, V02, L05, OPT 2 のもとに行った。

3章 ハウスホルダー法による 最小自乗問題の解法

3-1 方法の基礎

行列 A は m 行 n 列 ($m \geq n$) で、そのランクは n であるとする。ベクトル b が与えられたとき、 $b - Ax$ のユークリッドノルム $\|b - Ax\| = \left\{ \sum_{i=1}^m (b_i - \sum_{j=1}^n a_{ij} x_j)^2 \right\}^{1/2}$ を最小にするベクトル x を求める問題を考える。

ユークリッドノルムは直交変換のもとで不変であるから、

$$\|b - Ax\| = \|c - QAx\| \quad c = Qb, \quad Q^T Q = I$$

となる。そこで Q を

$$QA = R = \begin{pmatrix} \tilde{R} \\ 0 \end{pmatrix} \begin{matrix} n \times n \\ (m-n) \times n \end{matrix} \quad (3.1)$$

\tilde{R} : 上三角行列

となるように選ぶ。そうすると、 \tilde{c} を c の始めの n コの要素から成るベクトルとして、求めを解 x は

$$x = \tilde{R}^{-1} \tilde{c}$$

で与えられる。

3-2 計算方法とサブルーチンの説明

(3.1) のような分解を行うための効果的な方法の一つとして、ハウスホルダー変換がある。ハウスホルダー変換は基本対称直交変換とも呼ばれ、それは長さ1のベクトル u を用いて、

$$P = I - 2uu^T$$

と表わされる。容易にわかるように、 P は対称な直交行列である。今、与えられたベクトル x と y の長さが等しければ ($x^T x = y^T y$)、 x を y に変換するような行列 P が常に存在する。そのためには u として、

$$u = (y - x) / \{(y - x)^T (y - x)\}^{1/2}$$

をとればよい。 $x = y$ のときには $u = 0$ とする。また u の第 i 要素が 0 ならば、 Pv の第 i 要素は変化しない。即ち

$$u_i = 0 \quad \text{ならば} \quad (Pv)_i = v_i$$

となる。この変換は数値計算の手法上甚だ応用の広いものである。

ここでは次のようにしてこの変換を利用する。 $A^{(1)} = A$ から始めて、

$$P^{(k)} = I - \beta_k u^{(k)} u^{(k)T}$$

という対称直交行列 $P^{(k)}$ を用いて、 $A^{(2)}$ 、 $A^{(3)}$ 、 $A^{(n+1)}$ を順次

$$A^{(k+1)} = P^{(k)} A^{(k)} \quad k = 1, 2, \dots, n$$

に従って計算していく。ここで $P^{(k)}$ の要素は

$$a_{ik}^{(k+1)} = 0 \quad i = k+1, \dots, n$$

となるように定める。そのためには $u^{(k)}$ を次のようにして求めればよい。

$$\sigma_k = \left(\sum_{i=k}^m (a_{ik}^{(k)})^2 \right)^{1/2}$$

$$\beta_k = (\sigma_k (\sigma_k + |a_{kk}^{(k)}|))^{-1}$$

$$u_i^{(k)} = 0 \quad i < k$$

$$u_k^{(k)} = \text{sgn}(a_{kk}^{(k)}) (\sigma_k + |a_{kk}^{(k)}|)$$

$$u_i^{(k)} = a_{ik}^{(k)} \quad i > k$$

$A^{(k+1)} = P^{(k)} A^{(k)}$ を計算するには、 $P^{(k)}$ を求めてから計算するということはせず、

$$y_k^T = \beta_k u^{(k)T} A^{(k)}$$

$$A^{(k+1)} = A^{(k)} - u^{(k)} y_k^T$$

という順に計算する。 y_k と $A^{(k+1)}$ の計算には $u^{(k)}$ の始めの $k-1$ 個の要素が 0 であるという事実を十分に利用する。

さらに、できるだけ精度を確保するために、次のようなビョッティングを行う。すなわち、第 k ステップで $|a_{k,k}^{(k+1)}|$ が最大になるような $A^{(k)}$ の列を選んで変換を行う。そのためには、

$$s_j^{(k)} = \sum_{i=k}^n (a_{ij}^{(k)})^2 \quad j = k, \dots, n$$

としたとき、 $|a_{k,k}^{(k+1)}| = \sigma_k$ だから、 $s_j^{(k)}$ が最大になる列を選べばよい。 $s_j^{(k)}$ の計算には、列の長さが直交変換のもとで不変であるという性質を用いて、 $A^{(k)}$ が計算された後に

$$s_j^{(k)} = s_j^{(k-1)} - (a_{k-1,j}^{(k)})^2$$

に従って求める。サブルーチン DECOMS はここまでのことをする。

サブルーチン SOLVES は、 $c = Qb$ の計算と $x = \tilde{R}^{-1} \tilde{c}$ の計算を行う。

サブルーチン LSQHDS ではイタレイションによって解の精度を上げる試みをしている。今 x を初期解、修正解を $x' = x + e$ とし、剰余 r を

$$r = b - Ax$$

とすれば、

$$\|b - Ax'\| = \|r - Ae\|$$

となる。即ち、修正ベクトル e それ自身も又同じ最小自乗問題の解となっている。 A を一度分解し、変換を保存しておけば、 r を計算して e を求めることは容易である。なお、 r の計算は倍精度で行う。イタレイションは収束するまで続けるが、収束解の全部の桁が正しいという保証はない。(3-5「数値テスト及び結果」を参照) なお、このルーチンでは右辺のベクトルが複数個の場合にも使用できるようになっているので、右辺はベクトルではなく行列である。

3-3 用途

長方形行列 A を取り扱っているとしばしば $(A^T A)^{-1}$ や $\det(A^T A)$ の計算が必要になる。このとき A を (3.1) のように分解すれば、 \tilde{R} が上三角行列だからその逆行列は容易に求まり、前者は

$$(A^T A)^{-1} = \tilde{R}^T \tilde{R}^{-1} = \tilde{R}^{-1} (\tilde{R}^{-1})^T$$

から計算できる。又後者は

$$\det(A^T A) = \det(\tilde{R})^2 = r_{11}^2 r_{22}^2 \dots r_{nn}^2$$

より求まる。三角行列 \tilde{R} の対角要素 r_{ii} の値は、DECOMS 或は LSQHDS のアウトプット D にストアされている。

A が正方形行列である場合には、最小自乗解の特殊な場合として、一次方程式系の解が得られる。この場合さらに右辺が単位行列であれば逆行列が求まることになる。

3-4 サブルーチンの使用方法

(1) DECOMS, DECOMD

このサブルーチンは、長方形行列 A を $A = Q^T R$ と分解する。変換に関する情報は A の全下三角部分に、上三角行列 R の非対角要素は A の真上三角部分にそれぞれストアされる。従って元の A の値は失われる。

呼び出し形式

DECOMS ($M, N, A, KA, D, IP, W1, W2, ILL$)

インプット

M, N : 整数型変数名又は整数

M は方程式の数を, N は未知数の数を表わす。

$M \geq N \geq 2$

A : 大きさ KA の 2 次元アレイ。

$M \times N$ の長方形行列 A の要素をストアする。

KA : 整数型変数名又は整数

アレイ A の大きさを指定, $KA \geq M$

アウトプット

A : 3 角行列 \tilde{R} の非対角要素が真上三角部分に, ハウスホルダー変換で使用するベクトル u の非零要素が全下三角部分にそれぞれストアされている。図 3-1 を参照。

D : 大きさ N の 1 次元アレイ

3 角行列 \tilde{R} の対角要素がストアされている。

IP : 大きさ N の整数型 1 次元アレイ

ピボットの値がストアされている。

$W1, W2$: 大きさ N の 1 次元アレイ

ワーキングエリアとして使用。

$W2$ には $s_j^{(k)}$ の値がストアされている。

ILL : 整数型変数名

= 0 正常に終了したとき

= 30000 $M < N$ 又は $KA < M$ 又は $N < 2$ のとき

= $K, 1 \leq K \leq N$ $RANK(A) = K - 1$

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ のときには ILL 値をセットし, さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

K 番目の変換を行なうとき, 対象となる列の自乗和が正確に 0 となれば, $RANK(A) = K - 1$ と判定して ILL 値をセットする。さらに改頁して, "SUB DECOMS IS FAIL" と印字する。

詳細説明

インプット A とアウトプット A, D は次のようになっている。

列に関するピボットングを行なっているため, インプットとアウトプットとでは列の番号が違っている。そ

の対応の仕方はアウトプット IP によってわかる。即ち, アウトプットの第 i 列はインプットの第 $IP(i)$ 列に対応している。

DECOMD を用いるときには, $A, D, W1, W2$ を倍精度実数型として指定する。

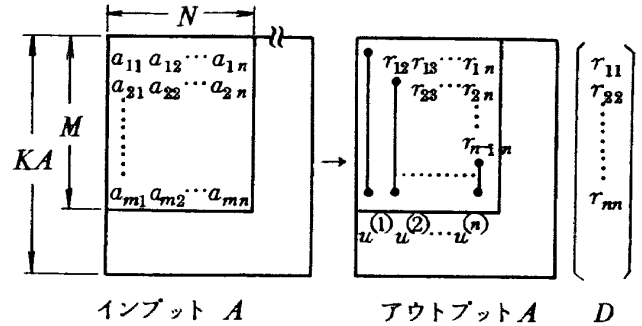


図 3-1

(2) SOLVES, SOLVED

このサブルーチンは, 与えられたベクトル b に対して, $c = Qb$ の計算と $x = \tilde{R}^{-1} c$ の計算を行なう。 b の値がストアされていたエリアに新しく c をストアするので, 元の b の値は失われる。本サブルーチンを使用する前に, DECOMS を用いて係数行列 A を分解しておく。

呼び出し形式

SOLVES ($M, N, A, KA, D, IP, B, X, W1$)

インプット

M, N : 整数型変数名又は整数

M は方程式の数, N は未知数の数を表わす。

$M \geq N \geq 2$

A : 大きさ KA の 2 次元アレイ

DECOMS からのアウトプット A

KA : 整数型変数名又は整数

アレイ A の大きさを指定, $KA \geq M$

D : 大きさ N の 1 次元アレイ

DECOMS からのアウトプット D

IP : 大きさ N の整数型 1 次元アレイ

DECOMS からのアウトプット IP

B : 大きさ M の 1 次元アレイ

方程式の右辺の要素をストアする。

アウトプット

B : $c = Qb$ の値がストアされている。

X : 大きさ N の 1 次元アレイ

最小自乗解が元の未知数の順にストアされている。

$W1$: 大きさ N の 1 次元アレイ

ワーキングエリアとして使用

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ であればなにもせずに

そのままリターンする。*ILL*値のセットやエラーメッセージの印字等も行なわない。

詳細説明

A, *D*, *IP* の値は変化しないが, *B* の値は変化する。*X* と *W1* の実引数を同じものにしてはならない。*B* と *X* 或は *B* と *W1* の実引数は同じのものであってもよい。その場合前者では *B(X)* に解がストアされ, 後者では *B* はワーキングエリアとして使用される。

SOLVED を用いるときには *A*, *D*, *B*, *X*, *W1* を倍精度実数型として指定する。

③ LSQHDS, LSQHDD

このサブルーチンは, 長方形行列 *A* を $A = Q^T R$ と分解し (DECOMS を使用), 後述するイタレイションによって最小自乗解を得る (SOLVES を使用)。右辺が複数個あってもよい。逐次式からもわかるように, 元の *A* と *B* の値は保存される。

呼び出し形式

LSQHDS (*M*, *N*, *NR*, *A*, *KA*, *B*, *X*, *KX*, *QR*, *D*, *IP*,
W1, *W2*, *W3*, *W4*, *ILL*)

インプット

M, *N*, *NR*: 整数型変数名又は整数定数

M は方程式の数, *N* は未知数の数を表わす。

NR は方程式の右辺の列の数

$M \geq N \geq 2$, $NR \geq 1$

A: 大きさ *KA* の 2 次元アレイ

$M \times N$ の長方形行列 *A* の要素をストアする。

KA: 整数型変数名又は整数定数

アレイ *A*, *B*, *QR* の大きさを指定, $KA \geq M$

B: 大きさ *KA* の 2 次元アレイ

方程式の右辺である $M \times NR$ 行列 *B* の要素をストアする。

KX: 整数型変数名又は整数定数

アレイ *X* の大きさを指定, $KX \geq N$

アウトプット

X: 大きさ *KX* の 2 次元アレイ

最小自乗問題の解が元の未知数の順にストアされている。

QR: 大きさ *KA* の 2 次元アレイ

3 角行列 \tilde{R} の非対角要素が真上 3 角部分に, ハウスホルダー変換で使用するベクトル *u* の非零要素が全下 3 角部分に, それぞれストアされている。DECOMS のアウトプット *A* を参照。

D: 大きさ *N* の 1 次元アレイ

3 角行列 \tilde{R} の対角成分がストアされている。

IP: 大きさ *N* の整数型 1 次元アレイ

ピボットの値がストアされている。

W1, *W2*, *W3*: 大きさ *N* の 1 次元アレイ

ワーキングエリアとして使用

W4: 大きさ *M* の 1 次元アレイ

ワーキングエリアとして使用

ILL: 整数型変数名

= 0 正常に終了したとき

= 30000 $M < N$ 又は $KA < M$ 又は $KX < N$
又は $N < 2$ のとき

= *K* $1 \leq K \leq N$ RANK(*A*) = *K* - 1

= -*L* $1 \leq L \leq NR$ 右辺の *L* 番目のベクトルに対してイタレイションが収束しなかった。

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ 又は $KX < N$ のときには, *ILL* 値をセットしさらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

DECOMS をコールして RANK(*A*) < *N* と判定されたときには, *ILL* 値をセットしてリターンする。このとき, 改頁して "SUB DECOMS IS FAIL" と印字される。

右辺が複数個あるとき ($NR \geq 2$), 各ベクトルごとイタレイションを行なうので, *L* 番目のベクトルに対して収束しなければ *ILL* 値をセットする。さらに改頁して, "SUB LSQHDS IS FAIL" と印字し, リターンする。

収束判定

イタレイションのスキームは

$$\begin{aligned} x^{(0)} &= 0 \\ r^{(k)} &= b - Ax^{(k)}, \quad e^{(k)}: \min \|r^{(k)} - Ae^{(k)}\| \\ x^{(k+1)} &= x^{(k)} + e^{(k)} \quad k = 0, 1, \dots \end{aligned}$$

である。収束の判定は

$$\|e^{(1)}\| / \|x^{(1)}\| > 1/4$$

の場合には, イタレイションループには入らずにエラーとして終了し,

$$\|e^{(k+1)}\| / \|e^{(k)}\| > 1/4 \quad \text{又は} \\ \|e^{(k)}\| / \|x^{(1)}\| < \varepsilon \quad \varepsilon = 2^{1-t} \quad t: \text{仮数部桁数}$$

となったときに収束したとして正常に終了する。

詳細説明

B は $M \times NR$ 行列, *X* は $N \times NR$ 行列なので *NR* が 1 であっても, それぞれ 2 次元アレイとして指定する。

ILL の値が負のときには, $|ILL| - 1$ までの右辺のベクトルに対してはイタレイションが収束しているので, その収束解が *X* にストアされている。

LSQHDD を用いるときには, *A*, *B*, *X*, *QR*, *D*, *W1*, *W2*, *W3*, *W4* を倍精度実数型として指定する。

3-5 数値テスト及び結果

(1)~(4)の最小自乗問題について行った。結果は表1-1

テストは、LSQHDSとLSQHDDを用いて、2章の ~表1-3である。

表 1-1

表 1-1

単精度

テスト 1 倍精度

	x_{1S}	x_{2S}
	1.0000000	5.3354787
	0.50000000	1.9476525
	0.33333334	0.95377918
	0.25000000	0.52142189
	0.20000000	0.29650499
a_m	0.0	4.335
a_f	0.0	4.622
s_m	0.0	4.335
s_f	0.0	5.662
r	6.199×10^{-7}	7.255×10^7
T_S		3.896
$-\log s_m$		-0.6370
$-\log s_f$		-0.7529

	x_{1D}	x_{2D}
	1.000000000000000000	0.999999999945136063
	0.500000000000000000	0.499999999982012801
	0.333333333333333333	0.333333333325702088
	0.250000000000000000	0.24999999996683421
	0.200000000000000000	0.19999999998826207
a_m	0.0	5.486×10^{-11}
a_f	0.0	5.835×10^{-11}
s_m	0.0	5.486×10^{-11}
s_f	0.0	7.098×10^{-11}
r	1.352×10^{-26}	7.255×10^7
T_D		14.43
$-\log s_m$		10.26
$-\log s_f$		10.15

単精度

表 1-2

テスト 2

	x_{1S}	x_{2S}	x_{3S}
	1.0000000	-8	0.99921529
	2.0000000	-9	1.9991392
	-1.0000000	2×10^{-4}	-0.99987446
	3.0000000	-20	2.9981695
	-4.0000000	34	-3.9970396
a_m	0.0	3.474×10^{-3}	2.960×10^{-3}
a_f	0.0	4.196×10^{-3}	3.673×10^{-3}
s_m	0.0		7.847×10^{-4}
s_f	0.0		1.318×10^{-3}
r	0.0	2.645×10^{-8}	2.645×10^8
T_S			3.616
$-\log s_m$			3.105
$-\log s_f$			2.880

表 1-2

倍精度

	x_{1D}	x_{2D}	x_{3D}
	1.000000000000000000	-5	0.999999999999954085
	2.000000000000000000	-5	1.99999999999994916
	-1.000000000000000000	1×10^{-14}	-0.999999999999989330
	3.000000000000000000	-12	2.99999999999988919
	-4.000000000000000000	20	-3.9999999999981238
a_m	0.0	1.959×10^{-13}	1.876×10^{-13}
a_f	0.0	2.401×10^{-13}	2.287×10^{-13}
s_m	0.0		4.691×10^{-14}
s_f	0.0		8.020×10^{-14}
r	0.0	2.645×10^8	2.645×10^8
T_D			14.15
$-\log s_m$			13.33
$-\log s_f$			13.10

表 1-3

テスト 4

単精度

	x_{1S}	x_{2S}	x_{3S}
	4.99999999 4.00000019 2.99999989 1.99999991 1.00000001	4.99999995 4.00000004 2.99999992 2.00000001 1.00000007	5.00000000 4.00000002 2.99999999 1.99999999 0.99999999
a_m	1.907×10^{-6}	7.749×10^{-7}	2.384×10^{-7}
a_f	2.361×10^{-6}	1.200×10^{-6}	2.815×10^{-7}
s_m	4.768×10^{-7}	6.855×10^{-7}	5.961×10^{-8}
s_f	7.515×10^{-7}	7.465×10^{-7}	8.574×10^{-8}
r	4880	2577	1913
T_S	5.983	6.121	6.186
$-\log s_m$	6.322	6.164	7.225
$-\log s_f$	6.124	6.127	7.067

表を見てわかるように、どの場合にも残差はほぼ正確な解のそれと一致している。又、テスト 3 を除いていつも“収束解”が得られてはいるが、計算機精度内の全桁があっているとは限らない。(3-2「計算方法とサブルーチンの説明」を参照) どの程度の有効数字が失われるのか量的関係は定かでないが、残差の大きさが影響するようである。テスト 1, 2 の残差が 0 になるベクトルでは、全桁が正しい。又、テスト 2, 4 では残差のオーダーほどの桁が失われている。(T_S , T_D と $-\log s_m$ を比べる) テスト 1 とテスト 2 では残差のオーダーはほぼ等しいが、失われる有効桁数はテスト 1 の方が多い。これは、最大 singular value と最小 singular value の比の違いによるのかもしれない。(この比の値は、前者で 10^6 のオーダー、後者では 10^3 のオーダーである)

テスト 3 については、 $\text{rank}(A) < n$ のため、本来ならこのルーチンでは計算不可能である。但し、DECOMS での $\text{rank}(A) < n$ の判定が「=0」で行われているため、そこでは $\text{rank}(A) < n$ とは判定されずに通過してきた。そのために、LSQHDS で通常どおり“解”が計算され、 $ILL = -2$ がセットされて戻ってきた。これは、右辺の第一ベクトルに対してはイタレーションが“収束”したことを意味している。しかし、結果の値は有効数字を 1 桁ももっていない。ちなみに、このときの 3 角行列 \tilde{R} の対角要素の値を記しておく、(上が単精度、下が倍精度)

$$(-29.53, -18.64, -16.17, 1.347 \times 10^{-7}, -1.242 \times 10^{-7})$$

$$(-29.53, -18.64, -16.17, 7.361 \times 10^{-18}, 1.387 \times 10^{-18})$$

である。これを見ると、 A のランクは 3 であると判断してさしつかえない。

4 章 直交化法による最小自乗問題の解法

4-1 方法の基礎

行列 A は m 行 n 列 ($m \geq n$) であるとする。このとき A の各列が 1 次独立であるとき (即ち A のランクは n) に限り、任意のベクトル b に対して $b - Ax$ のユークリッドノルム $\|b - Ax\|$ を最小にするようなベクトル x が一意に存在する。このとき次の 3 つの事柄が成り立っている。

(i) 残差ベクトル $b - Ax$ は A の各列に直交している。即ち

$$A^T(b - Ax) = 0$$

(ii) $A^+ = (A^T A)^{-1} A^T$ を A の擬逆行列とすると x は $x = A^+ b$

で与えられる。ここで、 A の擬逆行列とは

$$AXA = A, XAX = X$$

$$(AX)^T = AX, (XA)^T = XA$$

の 4 つの条件を満足する n 行 m 列の行列 X のことである。

(iii) $b - Ax$ は、 A を消去するような b 上への適当な射影によって得られる。即ち実エルミート行列 H で

$$H^2 = H, HA = 0$$

となる或る H によって

$$b - Ax = Hb$$

と表わされる。このような H としては

$$H = I - A(A^T A)^{-1} A^T$$

をとればよい。

普通の計算方法は、(i) から直接に得られる正規方程式

$$A^T A x = A^T b \quad (4.1)$$

を解く。或は上式の代わりに、正則な $n \times n$ 行列 S を用いて、方程式

$$(AS)^T Ax = (AS)^T b$$

を解いてもよい。ここで S を適当に選ぶと、 $AS=U$ は A の各列の正規直交基底となるようにできる。その結果得られる方程式

$$U^T Ax = U^T b$$

の係数行列 $U^T A$ の条件数は (4.1) の係数行列 $A^T A$ のそれよりも、一般に良い。

4-2 計算方法とサブルーチンの説明

ここで用いているアルゴリズムは行列 A を、0 ではない n 個の互いに直交する列より成る行列 U と n 次の正則な上三角行列 R と $D = (U^T U)^{-1}$ とによって、 $A = UDR$ と分解することに基づいている。

$$U^T (Ax - b) = 0$$

という条件から

$$Rx = U^T b$$

という系が得られる。これは正則な上三角行列を係数としているので逆進過程によって容易に解ける。ここで $R = U^T A$ と右辺の $U^T b$ は A と b に同じ変換 U^T を施すことによって得られる。

$Rx = U^T b$ へ変形する計算過程は n 段階で終了する。第 i 段階では、中間行列 $A^{(i-1)}$ から新しい中間行列 $A^{(i)}$ が、第 i 列より後の各列を第 i 列に直交させることによって得られる。これは、本質的にはシュミットの直交化法である。ただ後者では、第 i 列のベクトルを第 $i-1$ 列までの全ベクトルに直交させていく過程をとっている。

$A^{(0)} = A$ から始めて、 $A^{(i-1)}$ から $A^{(i)}$ を導くには、 $A^{(i)}$ の第 k 列を $a_k^{(i)}$ で表わすことにして、

$$a_k^{(i)} = a_k^{(i-1)} - a_i^{(i-1)} \frac{(a_i^{(i-1)})^T a_k^{(i-1)}}{(a_i^{(i-1)})^T a_i^{(i-1)}} \quad i < k$$

に従って各列を計算する。上式からわかるように、計算としては本質的にはガウス消去法型である。こうして n 段階後に得られる行列 $A^{(n)}$ は、直交基底の行列 U となっており、対角成分が 1 の上三角行列 S によって $U = AS$ と表わされる。又 $U^T A$ と $U^T b$ の計算も同時に進めることができ、消去過程が終わると基底行列 U 、正則な上三角行列 $R = U^T A$ 、系の右辺 $U^T b$ がすべて求まっている。

なお正規方程式 $A^T Ax = A^T b$ をガウス消去法で解いた場合にも同じ系に導かれる。というのも

$$A^T A = R^T DR$$

となり、 $R^T D$ は対角成分が 1 の下三角行列である。又、 U を正規化しておけば、 R は $A^T A$ をコレスキー分解して得られる上三角行列に一致する。

最終的に解を得るにあたり、次のようなイタレイシヨ

ンによって解の精度を上げる試みをしている。初期近似解を x_0 として順次 x_n を

$$R e_n = U^T (b - Ax_n)$$

$$x_{n+1} = x_n + e_n \quad n = 0, 1, \dots$$

に従って計算する。このとき $b - Ax_n$ の計算は倍精度演算によって行なう。このために行列 A 、 U の要素を保存しておく必要がある。イタレイシヨンは収束するまで続けるが、収束解の全部の桁が正しいという保証はない。

(4-5 「数値テスト及び結果」を参照)

ここには 4 種のサブルーチンがあって、そのうちの 2 つは最小自乗解を求めるもの、残りの 2 つは正方行列の逆行列を求めるものである。前者は、右辺が単一ベクトルの場合 (LSQE1S) と複数個ある場合 (LSQE2S) とからなっており、それぞれ上記のイタレイシヨンを行っている。後者は、イタレイシヨンを行なうもの (OTIV1S) と行なわないもの (OTIVNS) とからなっている。

4-3 用途

4 つのサブルーチン各々について自明である。なお $A^T A = R^T DR$ で $R^T D$ は対角成分が 1 の三角行列だから、

$$\det(A^T A) = \det(R)$$

となる。上三角行列 R の要素は、アウトプットの 1 次元アレイ R にストアされているので、次のようにすれば $\det(R)$ の計算ができる。

$$K = N * (N + 1) / 2$$

$$DET = R(K)$$

$$DO 10 I = 2, N$$

$$K = K - I$$

$$DET = DET * R(K)$$

10 CONTINUE

又、 $(A^T A)^{-1} = R^{-1} D^{-1} (R^{-1})^T$ から $(A^T A)^{-1}$ の計算も可能である。

4-4 サブルーチンの使用方法

(1) LSQE1S, LSQE1D

このサブルーチンは、右辺が1つしかない場合、即ち $Ax=b$ という形 (b, x はベクトル) の最小自乗問題を解くものである。イタレーションによって解の精度を上げる試みをしているので、インプットである A と b の値は保存される。

呼び出し形式

LSQE1S ($M, N, A, KA, B, X, U, R, W1,$
 $W2, ILL$)

インプット

M, N : 整数型変数名又は整数

M は方程式の数を、 N は未知数の数を表わす

$$2 \leq N \leq M$$

A : 大きさ KA の2次元アレイ

$M \times N$ の長方形アレイ A の要素をストアする。

KA : 整数型変数名又は整数

アレイ A と U の大きさを指定、 $KA \geq M$

B : 大きさ M の1次元アレイ

方程式の右辺であるベクトル b の要素をストアする。

アウトプット

X : 大きさ N の1次元アレイ

最小自乗問題の解がストアされている。

U : 大きさ KA の2次元アレイ

$M \times N$ 行列 U の要素がストアされている。

R : 大きさ $N(N+1)/2$ の1次元アレイ

3角行列 R の要素がストアされている。

図4-1を参照。

$W1$: 大きさ M の1次元アレイ

ワーキングエリアとして使用。

$W2$: 大きさ N の1次元アレイ

ワーキングエリアとして使用。

ILL : 整数型変数名

= 0 正常に終了したとき

= 30000 $M < N$ 又は $KA < M$ 又は $N < 2$ のとき

= 1 イタレーションが収束しなかったとき

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ のときには ILL 値をセットし、さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

イタレーションが収束しなかったときには、 ILL 値をセットし、さらに改頁して "SUB LSQE1S IS FAIL" と印字する。このとき X には最終近似解がストアされて

いる。

収束判定

イタレーションは、近似解を x_n 修正ベクトルを e_n として、

$$R x_0 = U^T b, \quad x_{n+1} = x_n + e_n$$

$$R e_n = U^T (b - A x_n)$$

に従って行い。収束の判定は

$$\|e_n\| > \|x_n\| / 2$$

となったとき、収束しないとして $ILL=1$ をセットして終了する。

$$\|e_n\| / \|x_n\| < \epsilon \text{ 又は } \|e_n\| > \|e_{n-1}\| / 10$$

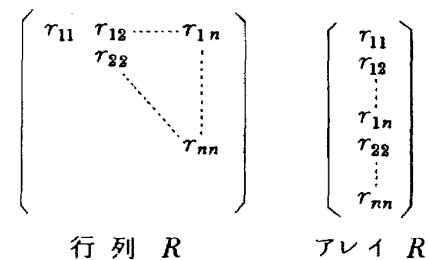
$$\epsilon = 2^{1-t}$$

t : 仮数部桁数

となったときに収束したとして正常に終了する。

詳細説明

3角行列 R の要素は下図のようにストアされている。



行列 R アレイ R

図4-1

LSQE1D を用いるときには $A, B, X, U, R, W1, W2$ を倍精度実数型として指定する。

(2) LSQE2S, LSQE2D

このサブルーチンは、右辺が複数個ある場合、即ち $AX=B$ という形 (B, X は行列) の最小自乗問題を解くものである。イタレーションによって解の精度を上げる試みをしているので、インプットである A と B の値は保存される。

呼び出し形式

LSQE2S ($M, N, NR, A, KA, B, X, KX, U,$
 $R, W1, W2, W3, ILL$)

インプット

M, N, NR : 整数型変数名又は整数

M は方程式の数を、 N は未知数の数を表わす。

NR は方程式の右辺の列の数

$$2 \leq N \leq M, \quad 2 \leq NR$$

A : 大きさ KA の2次元アレイ

$M \times N$ の長方形アレイ A の要素をストアする。

KA : 整数型変数名又は整数

アレイ A, B, U の大きさを指定, $KA \geq M$

B : 大きさ KA の 2次元アレイ

方程式の右辺である $M \times NR$ 行列 B の要素をストアする。

KX : 整数型変数名又は整定数

アレイの大きさを指定, $KX \geq N$

アウトプット

X : 大きさ KX の 2次元アレイ

$N \times NR$ 行列で最小自乗問題の解がストアされている。

U : 大きさ KA の 2次元アレイ

$M \times N$ 行列 U の要素がストアされている。

R : 大きさ $N(N+1)/2$ の 1次元アレイ

3角行列 R の要素がストアされている。

図 4-1 を参照。

$W1$: 大きさ M の 1次元アレイ

ワーキングエリアとして使用。

$W2, W3$: 大きさ N の 1次元アレイ

ワーキングエリアとして使用。

ILL : 整数型変数名

= 0 正常に終了したとき。

= 30000 $M < N$ 又は $KA < M$ 又は $KX < N$ 又は $N < 2$ 又は $NR < 2$ のとき。

= K $1 \leq K \leq NR$ 右辺の第 K 番目のベクトルに対してイタレーションが収束しなかったとき。

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ 又は $KX < N$ 又は $NR < 2$ のときには ILL 値をセットし, さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

右辺の第 K 番目のベクトルに対してイタレーションが収束しなかったときには, ILL 値をセットし, さらに改頁して "SUB LSQE2S IS FAIL" と印字する。このとき X の第 $K-1$ 列までには正常に終了した収束解がストアされている。

収束判定

LSQE1S と同じイタレーションを右辺の各ベクトルごとに行い, 収束判定も各ベクトルごとと同様の基準が適用される。

詳細説明

3角行列 R の要素の配列の仕方については図 4-1 を参照。

LSQE2D を用いるときには $A, B, X, U, R, W1, W2, W3$ を倍精度実数型として指定する。

(3) OTIVNS, OTIVND

このサブルーチンは, 正方行列 A の逆行列 AI を $AI = R^{-1}U^T$ として求めるものである。インプットである A の値は失われる。

呼び出し形式

OTIVNS ($N, A, KA, R, W1, ILL$)

インプット

N : 整数型変数名又は整定数

正方行列 A の次数, $2 \leq N$

A : 大きさ KA の 2次元アレイ

正方行列 A の要素をストアする。

KA : 整数型変数名又は整定数

アレイ A の大きさを指定, $KA \geq N$

アウトプット

A : 行列 A の逆行列の要素がストアされている。

R : 大きさ $N(N+1)/2$ の 1次元アレイ

3角行列 R の要素がストアされている。

$W1$: 大きさ N の 1次元アレイ

ワーキングエリアとして使用

ILL : 整数型変数名

= 0 正常に終了したとき。

= 30000 $KA < N$ 又は $N < 2$ のとき。

エラー処理

$N < 2$ 又は $KA < N$ のときには ILL 値をセットし, さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

詳細説明

このサブルーチンではイタレーションを行わないので, 行列 A の値を保存しておく必要がない。そこで A から計算される U の要素をアレイ A にストアし, ついでこれから計算される A の逆行列の要素を又アレイ A にストアしている。

3角行列 R の要素の配列の仕方については図 4-1 を参照。

OTIVND を用いるときには $A, R, W1$ を倍精度実数型として指定する。

(4) OTIVIS, OTIVID

このサブルーチンは, 正方行列 A の逆行列 AI を $AI = R^{-1}U^T$ で求め, その後イタレーションによって精度を上げている。インプットである A の値は保存される。

呼び出し形式

OTIVIS ($N, A, KA, AI, U, R, W1, W2, W3, ILL$)

インプット

N : 整数型変数名又は整数

正方行列 A の次数, $2 \leq N$

A : 大きさ KA の 2 次元アレイ

正方行列 A の要素をストアする。

KA : 整数型変数名又は整数

アレイ A, AI, U の大きさを指定, $KA \geq N$

アウトプット

AI : 大きさ KA の 2 次元アレイ

行列 A の逆行列の要素がストアされている。

U : 大きさ KA の 2 次元アレイ

N 次正方行列 U の要素がストアされている。

R : 大きさ $N(N+1)/2$ の 1 次元アレイ

3 角行列 R の要素がストアされている。

$W1, W2, W3$: 大きさ N の 1 次元アレイ

ワーキングエリアとして使用。

ILL : 整数型変数名

= 0 正常に終了したとき。

= 30000 $KA < N$ 又は $N < 2$ のとき。

= 1 イタレーションが収束しなかったとき。

エラー処理

$N < 2$ 又は $KA < N$ のときには ILL 値をセットし、さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

イタレーションが収束しなかったときには、 ILL 値をセットし、さらに改頁して "SUB OTIVIS IS FAIL" と印字する。このとき AI には最終近似逆行列の要素がストアされている。

収束判定

イタレーションは

$$A_{n+1} = A_n + E_n$$

E_n の各列ベクトルは LSQE2S と同じ方法で計算

に従って行なり。収束の判定は

$$\|A\|_E = \left(\sum_{i,j=1}^N a_{ij}^2 \right)^{1/2}$$

としたとき

$$\|E_n\|_E > \|A_n\|_E / 2$$

となると、収束しないとして $ILL=1$ をセットして終了する。

$$\|E_n\|_E / \|A_n\|_E < \epsilon \quad \epsilon = 2^{1-t} \quad t: \text{仮数部桁数}$$

$$\text{又は } \|E_n\|_E > \|E_{n-1}\|_E / 10$$

となったときに収束したとして正常に終了する。

詳細説明

このサブルーチンではイタレーションによって解の精

度を上げているので、行列 A 及び U の要素を保存しておく必要がある。そのために、OTIVNS に比べてメモリーを多く必要とする。

3 角行列 R の要素の配列の仕方については図 4-1 を参照。

OTIVID を用いるときには $A, AI, U, R, W1, W2, W3$ を倍精度実数型として指定する。

4-5 数値テスト及び結果

最小自乗問題を解くルーチンのテストは、LSQE2S と LSQE2D を用いて、2章の(1)~(4)の最小自乗問題について行った。結果は表2-1~表2-3である。テスト1の単精度、テスト3については収束せずに、ILL=1が出力されてきた。

ここでも表を見てわかるように、残差は正解のそれにほぼ一致しており、残差が0となるベクトルに対しては、“収束解”の全桁が正しい。しかし、残差が0でないものについては、“収束解”でも全桁が正しいとは限らない。テスト1の倍精度、テスト4についてはLSQHDD、LSQHDSよりも1桁ずつ精度が良いが、テスト2の倍精度では逆にLSQE2Dの方が1桁精度が悪い。

テスト3については、ここでも本来このルーチンでは計算不可能である。このときの3角行列Rの対角成分の値は、(上が単精度、下が倍精度)

$$(872, 264.5, 343.8, 4.741 \times 10^{-14}, 1.160 \times 10^{-14})$$

$$(872, 264.5, 343.8, 2.581 \times 10^{-35}, 1.866 \times 10^{-36})$$

である。これを見ると、Aのランクは3であると判断してさしつかえない。

正方行列の逆行列を求めるルーチンについては、表2-4と表2-5の2つの行列について行った。表には相対誤差の行列を掲げてある。X_SはOTIVNS, X_DはOTIVND, Y_SはOTIVIS, Y_DはOTIVIDに対応するものである。aとfは、X_Sについて書けば

表2-1

テスト1

倍精度

	x_{1D}	x_{2D}
	1.000000000000000000	1.00000000000106584
	0.500000000000000000	0.500000000000361341
	0.333333333333333333	0.333333333333489477
	0.250000000000000000	0.250000000000068665
	0.200000000000000000	0.200000000000024503
a_m	0.0	1.066×10^{-12}
a_f	0.0	1.139×10^{-12}
s_m	0.0	1.066×10^{-12}
s_f	0.0	1.403×10^{-12}
r	1.352×10^{-26}	7.255×10^7
T_D		14.43
$-\log s_m$		11.97
$-\log s_f$		11.85

表2-2

テスト2

単精度

	x_{1S}	x_{2S}	x_{3S}
	1.0000000	86	0.99911812
	2.0000000	-93	1.9990406
	-1.0000000	3×10^{-5}	-0.99995470
	3.0000000	-190	2.9980310
	-4.0000000	268	-3.9971867
a_m	0.0	2.678×10^{-3}	2.813×10^{-3}
a_f	0.0	3.519×10^{-3}	3.673×10^{-3}
s_m	0.0		8.819×10^{-4}
s_f	0.0		1.391×10^{-3}
r	0.0	2.645×10^8	2.645×10^8
T_S			3.616
$-\log s_m$			3.055
$-\log s_f$			2.857

$$a = \max_{i,j} |(X_S)_{ij}|$$

$$f = \sqrt{\sum_{i,j=1}^n (X_S)_{ij}^2} = \|X_S\|_E \quad n: \text{行列の次数}$$

の値である。第2のテスト行列に対する Y_S, Y_D は O 行列である。

第1の行列のスペクトルノルムに関する条件数 χ_1 は 3.664×10^6 であり、行列の条件としては甚だ悪い。第2の行列のそれは、 $\chi_2 = 2.984 \times 10^8$ である。 χ と計算機精度、解の精度の間の関係として、「1次方程式系の

解法 I」⁽¹⁾ で考えたものと同じのを採ると、 X_S, X_D の a の解の精度の代表として、

$26 \log 2 - \log \chi_1 = 1.263$	$\log a = 0.3947$
$61 \log 2 - \log \chi_1 = 11.80$	$\log a = 11.21$
$26 \log 2 - \log \chi_2 = 4.352$	$\log a = 4.361$
$61 \log 2 - \log \chi_2 = 14.89$	$\log a = 15.07$

となる。ここでも、この2つの行列に関する限り、解の精度は、計算機の精度よりも条件数のオーダー程度悪くなっているといえる。

表 2-2

倍精度

	x_{1D}	x_{2D}	x_{3D}
	1.0000000000000000	-19	0.999999999999807834
	2.0000000000000000	-21	1.99999999999978771
	-1.0000000000000000	4×10^{-14}	-0.999999999999959278
	3.0000000000000000	-46	2.99999999999954037
	-4.0000000000000000	77	-3.99999999999923328
a_m	0.0	7.730×10^{-13}	7.667×10^{-13}
a_f	0.0	9.469×10^{-13}	9.396×10^{-13}
s_m	0.0		1.922×10^{-13}
s_f	0.0		3.318×10^{-13}
r	0.0	2.645×10^8	2.645×10^8
T_D			14.15
$-\log s_m$			12.72
$-\log s_f$			12.48

表 2-3

テスト 4

単精度

	x_{1S}	x_{2S}	x_{3S}
	5.0000000	5.0000000	5.0000000
	4.0000001	4.0000001	4.0000000
	2.9999999	2.9999999	2.9999999
	1.9999999	1.9999999	2.0000000
	0.99999994	0.99999994	1.0000000
a_m	1.192×10^{-7}	1.192×10^{-7}	5.960×10^{-8}
a_f	1.712×10^{-7}	1.577×10^{-7}	6.664×10^{-8}
s_m	5.960×10^{-8}	5.960×10^{-8}	2.981×10^{-8}
s_f	8.266×10^{-8}	7.565×10^{-8}	3.583×10^{-8}
r	4880	2577	1913
T_S	5.983	6.121	6.186
$-\log s_m$	7.225	7.225	7.526
$-\log s_f$	7.083	7.121	7.446

表 2-4

$$\begin{array}{c}
 A \\
 \begin{pmatrix} -74 & 80 & 18 & -11 & -4 & -8 \\ 14 & -69 & 21 & 28 & 0 & 7 \\ 66 & -72 & -5 & 7 & 1 & 4 \\ -12 & 66 & -30 & -23 & 3 & -3 \\ 3 & 8 & -7 & -4 & 1 & 0 \\ 4 & -12 & 4 & 4 & 0 & 1 \end{pmatrix}
 \end{array}
 \begin{array}{c}
 A^{-1} \\
 \begin{pmatrix} 1 & 0 & -7 & -40 & 131 & -84 \\ 0 & 1 & 7 & 35 & -112 & 70 \\ -2 & 2 & 29 & 155 & -502 & 319 \\ 15 & -12 & -192 & -1034 & 3354 & -2130 \\ 43 & -42 & -600 & -3211 & 10406 & -6595 \\ -56 & 52 & 764 & 4096 & -13276 & 8421 \end{pmatrix}
 \end{array}$$

$$\begin{array}{c}
 X_S \\
 10^{-3} \times \begin{pmatrix} -210 & -145 & -46 & -1 & 0.7 & 0.3 \\ 179 & 123 & -40 & -1 & 0.7 & 0.3 \\ -403 & 278 & -43 & -1 & 0.7 & 0.3 \\ -359 & 309 & -43 & -1 & 0.7 & 0.3 \\ -388 & 274 & -43 & -1 & 0.7 & 0.3 \\ -380 & 282 & -43 & -1 & 0.7 & 0.3 \end{pmatrix}
 \end{array}
 \begin{array}{c}
 Y_S \\
 10^{-7} \times \begin{pmatrix} -3 & 10 & -5 & 4 & -3 & -5 \\ 2 & -9 & -4 & 3 & -3 & -5 \\ -5 & -20 & -5 & 3 & -3 & -5 \\ -4 & -22 & -5 & 4 & -3 & -5 \\ -5 & -20 & -5 & 3 & -3 & -5 \\ -5 & -20 & -5 & 3 & -3 & -5 \end{pmatrix}
 \end{array}$$

$$\begin{array}{c}
 a = 4.030 \times 10^{-1} \\
 f = 1.019
 \end{array}
 \begin{array}{c}
 a = 2.223 \times 10^{-6} \\
 f = 4.870 \times 10^{-6}
 \end{array}$$

$$\begin{array}{c}
 X_D \\
 10^{-14} \times \begin{pmatrix} 8 & -288 & 27 & -8 & 2 & 2 \\ -7 & 245 & 23 & -7 & 2 & 2 \\ 16 & 550 & 25 & -8 & 2 & 2 \\ 14 & 613 & 25 & -8 & 2 & 2 \\ 15 & 543 & 25 & -8 & 2 & 2 \\ 15 & 559 & 25 & -8 & 2 & 2 \end{pmatrix}
 \end{array}
 \begin{array}{c}
 Y_D \\
 \begin{pmatrix} 0 & 6 \times 10^{-36} & 0 & 0 & 0 & 0 \\ 8 \times 10^{-37} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}
 \end{array}$$

$$\begin{array}{c}
 a = 6.125 \times 10^{-12} \\
 f = 1.197 \times 10^{-11}
 \end{array}$$

表 2-5

$$\begin{array}{c}
 A \\
 \begin{pmatrix} 5 & 7 & 6 & 5 \\ 7 & 10 & 8 & 7 \\ 6 & 8 & 10 & 9 \\ 5 & 7 & 9 & 10 \end{pmatrix}
 \end{array}
 \begin{array}{c}
 A^{-1} \\
 \begin{pmatrix} 68 & -41 & -17 & 10 \\ -41 & 25 & 10 & -6 \\ -17 & 10 & 5 & -3 \\ 10 & -6 & -3 & 2 \end{pmatrix}
 \end{array}$$

$$\begin{array}{c}
 X_S \\
 10^{-6} \times \begin{pmatrix} 5 & -15 & -27 & 43 \\ 5 & -14 & -28 & 44 \\ 5 & -15 & -23 & 36 \\ 6 & -15 & -23 & 32 \end{pmatrix}
 \end{array}
 \begin{array}{c}
 X_D \\
 10^{-17} \times \begin{pmatrix} -8 & 22 & 41 & -85 \\ -9 & 22 & 42 & -85 \\ -9 & 23 & 35 & -71 \\ -9 & 22 & 34 & -63 \end{pmatrix}
 \end{array}$$

$$\begin{array}{c}
 a = 4.351 \times 10^{-6} \\
 f = 9.869 \times 10^{-6}
 \end{array}
 \begin{array}{c}
 a = 8.546 \times 10^{-16} \\
 f = 1.781 \times 10^{-16}
 \end{array}$$

固有値が分離されたかどうかの判定は次のようにして行なう。

$$\delta = \varepsilon \cdot \max_i (|q_i| + |e_i|) \quad \varepsilon : \text{計算機精度}$$

として,

$$|e_n| \leq \delta$$

となれば, $|q_n|$ を singular value とし, 行列の次数を 1 つ小さくする。

サブルーチン SVDS は, この方法に基づいて A を $U\Sigma V^T$ と分解する。

サブルーチン MINFIS も同様に, A を $\tilde{U}\tilde{\Sigma}V^T$ 分解し, 同時に行列積 $\tilde{U}B$ も計算する。このサブルーチンでは $m < n$ であってもよい。

なお, どちらのサブルーチンでも, 求めた singular values は大きさの順には並んでいない。

5-3 用途

A を $U\Sigma V^T$ と分解すれば, A の擬逆行列 A^\dagger が次のようにして計算できる。(擬逆行列については 4-1 を参照)

$$A^\dagger = V \Sigma^\dagger U^T$$

ここで $\Sigma^\dagger = (\sigma_i^\dagger)$ は

$$\sigma_i^\dagger = \begin{cases} \sigma_i^{-1} & \sigma_i > 0 \\ 0 & \sigma_i = 0 \end{cases}$$

として定まる対角行列である。

行列 A のランク r が $r < n$ のとき, 同次方程式

$$Ax = 0$$

の解が, 次のようにして求まる。このときには,

$$\sigma_i = 0 \quad i = r+1, \dots, n$$

であり, $A = U\Sigma V^T$ であることより, 行列 U, V の列ベクトルをそれぞれ u_i, v_i とすれば

$$Av_i = \sigma_i u_i \quad i = 1, \dots, n$$

と書ける。従って, $i = r+1, \dots, n$ に対しては

$$Av_i = 0$$

となり, これが同次方程式の解となる。

b を与えられたベクトルとして, $\|b - Ax\|$ を最小にするという最小自乗問題の解 x が次の式に従って求まる。

$$x = A^\dagger b = V \Sigma^\dagger U^T b = V \Sigma^\dagger c$$

このとき, A のランクが n より小さいときには, このままでは解が一意には定まらないが, さらに $\|x\|$ を最小にするという条件をつければ解は一意的に定まる。その解も同じ式によって計算される。

5-4 サブルーチンの使用方法

(1) SVDS, SVDD

このサブルーチンは長方形行列 A を $A=U\Sigma V^T$ と singular value 分解するものである。行列 U と V の要素は必要に応じて計算されるようになっている。

呼び出し形式

SVDS ($M, N, A, KA, U, V, KV, S, LU,$
 $LV, W1, ILL$)

インプット

M, N : 整数型変数名又は整数

M は行列 A の行数を, N は列数を表わす。 $M \geq N \geq 2$

A : 大きさ KA の2次元アレイ

$M \times N$ 行列 A の要素をストアする。

KA : 整数型変数名又は整数

アレイ A, U の大きさを指定, $KA \geq M$

KV : 整数型変数名又は整数

アレイ V の大きさを指定, $KV \geq N$

LU, LV : 論理型変数名又は論理型定数

アウトプットの行列 U が必要なときには, $LU = .TRUE.$, 必要でないときには, $LU = .FALSE.$

アウトプットの行列 V が必要なときには, $LV = .TRUE.$, 必要でないときには, $LV = .FALSE.$

アウトプット

U : 大きさ KA の2次元アレイ

$LU = .TRUE.$ のとき, 正規直交化された列をもつ $M \times N$ 行列 U の要素がストアされている。 $LU = .FALSE.$ のときには, ワーキングエリアとして使用される。

V : 大きさ KV の2次元アレイ

$LV = .TRUE.$ のとき, N 次直交行列 V の要素がストアされている。 $LV = .FALSE.$ のときには, V は使用されない。

S : 大きさ N の1次元アレイ

行列 A の singular value がストアされている。

$W1$: 大きさ N の1次元アレイ

ワーキングエリアとして使用

ILL : 整数型変数名

= 0 正常に終了したとき

= 30000 $M < N$ 又は $N < 2$ 又は $KA < M$
又は $KV < N$

エラー処理

$N < 2$ 又は $M < N$ 又は $KA < M$ 又は $KV < N$ のときには, ILL 値をセットし, さらに改頁してその旨の印字をする。この場合それ以上の計算はせずにそのままリターンする。

詳細説明

$LV = .FALSE.$ のときアレイ V は使用されない。このときには仮引数 V, KV の実引数は型が合っていて, $KV \geq N$ を満足していれば何であってもよい。特に, V がアレイである必要もない。一方 U については, $LU = .FALSE.$ であってもワーキングエリアとして使用されるので, 実アレイが対応していなければならない。仮引数 A と U の実引数が異なれば, 元の A の値は保存されるが, 同じであれば $LU = .TRUE.$ のときアレイ A に行列 U の要素がストアされ, $LU = .FALSE.$ のときアレイ A はワーキングエリアとして使用される。従って, 行列 A の singular values が欲しいだけのときには,

CALL SVDS ($M, N, A, KA, A, A, KA,$
 $S, .FALSE., .FALSE., W1, ILL$)

でよい。

SVDDを用いるときには, $A, U, V, S, W1$ を倍精度実数型として指定する。

(2) MINFIS, MINFID

このサブルーチンは, $m \times n$ 行列 A を $A = \tilde{U} \tilde{\Sigma} \tilde{V}^T$ と singular value 分解する。 A の要素がストアされていた場所に V の要素をストアするので, 元の A の値は失われる。これと同時に, $m \times r$ 行列 B に対して, $C = \tilde{U}^T B$ の値が計算され, B がストアされていた場所に C をストアする。

呼び出し形式

MINFIS ($M, N, NR, AB, KA, Q, W1, ILL$)

インプット

M, N, NR : 整数型変数名又は整数

M は行列 A の行数を, N は列数を表わす。

NR は, 最小自乗問題を解く場合の右辺の列数。

$M, N > 2, NR > 0$

AB : 大きさ KA の2次元アレイ

$\max(M, N) \times (N + NR)$ 行列であって, 行列 A と行列 B の要素を次のようにストアする。

$$AB(I, J) = a_{ij} \quad 1 < I < M, 1 < J < N$$

$$AB(I, N+J) = b_{ij} \quad 1 < I < M, 1 < J < NR$$

図5-2を参照。

KA: 整数型変数名又は整数

アレイ A の大きさを指定, $KA > \max(M, N)$

アウトプット

AB: V の要素と C の要素が次のようにストアされている。

$$AB(I, J) = v_{ij} \quad 1 \leq I, J \leq N$$

$$AB(I, N+J) = c_{ij} \quad 1 \leq I \leq \max(M, N), \\ 1 \leq J \leq NP$$

図 5-2 を参照。

Q: 大きさ N の 1 次元アレイ

行列 A の singular values がストアされている。

W1: 大きさ N の 1 次元アレイ

ワーキングエリアとして使用

ILL: 整数型変数名

= 0 正常に終了したとき。

= 30000 $N < 2$ 又は $M < 2$ 又は $KA < \max(M, N)$

エラー処理

$N < 2$ 又は $M < 2$ 又は $KA < \max(M, N)$ のときには,

ILL 値をセットし, さらに改頁してその旨の印字をする。

この場合それ以上の計算はせずにそのままリターンする。

詳細説明

このサブルーチンでは $M < N$ であってもよい。このときには, $N - M$ 個の自明な零 singular values がある。

インプットとアウトプットのアレイ AB は下図のようになっている。

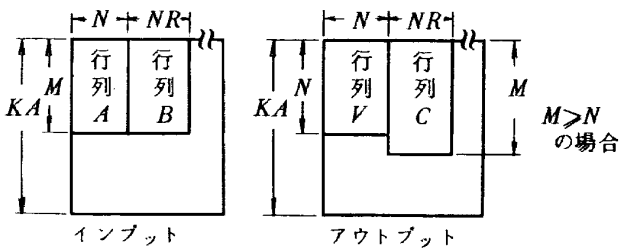


図 5-2

行列 \tilde{U} の要素が必要なきには, 右辺 B に M 次の単位行列をセットしておけば ($M \geq N$ のとき), C の所に \tilde{U}^T が計算されてくる。 $M < N$ のときには,

$$\begin{pmatrix} 1 & 0 \\ & \ddots \\ & & 1 \\ & & & 0 \end{pmatrix}$$

という形の $N \times M$ 行列をセットしておき, $N - M$ 個の自明な零 singular values に対応する行を取り去ってできる $M \times M$ 行列が \tilde{U}^T を与える。

MINFID を用いるときには, AB, Q, W1 を倍精度実数型として指定する。

5-5 数値テスト及び結果

テスト1~テスト4の行列Aを、サブルーチンSVDS, SVDDを用いてsingular value分解した結果は、表3-1のように特徴づけられる。表中の記号は下記の通りである。

$$A = (a_{ij}), Y = (y_{ij})$$

$$m(A) = \max |a_{ij}|, E(A) = \sqrt{\sum a_{ij}^2}$$

$$\bar{A} = U\Sigma V^T, X = A - \bar{A}, y_{ij} = (a_{ij} - \bar{a}_{ij})/a_{ij}$$

この表から、行列U, Vの各列は殆んど正規直交化されていると見てよい。又、 $A = U\Sigma V^T$ という分解も、テスト1の行列だけが若干悪いが、まずまずの精度であると判断できる。次に、これらの行列のsingular valueを単精度のものについてだけ、表3-2に掲げておく。そこでは、正確な値として倍精度で計算されたものの10桁を採用した。なお、相対誤差の欄で*印のある所は絶対誤差の値である。

表を見ると、絶対値の大きいsingular valueほど精度よく求まっているのがわかる。これは、固有値を求めるルーチンにとって普通のことである。テスト3の行列には0というsingular valueが2つあるわけだが、それぞれSVDSでは 2×10^{-7} と 9×10^{-7} 、SVDDでは 6×10^{-18} と 2×10^{-18} として求まっている。どちらの場合にも十分0と見なせる値である。これより、同次方程式 $Ax = 0$ は、テスト3の行列を除いては自明な解 $x = 0$ しかもたないことがわかる。(rank(A)=n) テスト3の行列では、SVDSのアウトプットVの第2列と第5列が自明でない1次独立な解を与えてくれる。この解を単精度のものについて表にしたのが表3-3である。正規化されているため値は小数になっているが、0でない絶

対値最小の要素で割り算して、適当整数倍することによって整数化すると、正確な解として、

$$(-137, 144, -17, 221, 135)$$

$$(0, 12, 10, 7, -23)$$

が得られる。

最小自乗問題を解くには、サブルーチンMINFIS或はMINFIDを用いて、

$$x = A^+ b = V \Sigma^+ U^T b = V \Sigma^+ c$$

という行列の掛け算によって求める。ここで、VとcはMINFISからのアウトプットABにセットされている。もちろん、中央の式に従ってSVDSによっても計算できる。ここで Σ^+ の計算にあたって、単精度では 10^{-5} より小さいsingular valueを、倍精度では 10^{-15} より小さいそれを0で置き代えた。その結果が表3-4~表3-7である。ここでも、解の精度が悪いときでも残差の値はほぼ正しいものになっている。又、テスト1, テスト2については残差が0になるベクトルに対しての方が精度はよくなっている。このルーチンでは残差の影響の他に、 $A = U\Sigma V^T$ と分解する精度の影響も考慮されねばならないだろう。

なお、このルーチンが前2者と最も異なる点は、テスト3のようにrank(A) < nとなっている問題に対しても、 $\|x\|$ を最小にするという付帯条件の解を求められることにある。但し、このときにはn - rank(A)個の自明な零-singular valueがあるが、計算機では丁度0という値になって出力されることは殆んどないので、何らかの判定基準に従ってそれらの値を0に置き代えてやる必要がある。

表3-1

単精度

テストNo	m(X)	E(X)	m(Y)	E(Y)	m(I-U ^T U)	m(I-V ^T V)
1	1.875×10^{-1}	3.206×10^{-1}	8.186×10^{-6}	1.429×10^{-5}	1.192×10^{-7}	8.941×10^{-8}
2	5.722×10^{-6}	1.317×10^{-6}	1.192×10^{-6}	1.782×10^{-6}	7.451×10^{-8}	8.941×10^{-8}
3	1.907×10^{-6}	3.088×10^{-6}	7.122×10^{-7}	1.039×10^{-6}	8.941×10^{-8}	5.960×10^{-8}
4	1.788×10^{-7}	4.059×10^{-7}	1.192×10^{-7}	2.882×10^{-7}	5.960×10^{-8}	5.960×10^{-8}

倍精度

テストNo	m(X)	E(X)	m(Y)	E(Y)	m(I-U ^T U)	m(I-V ^T V)
1	3.688×10^{-12}	7.262×10^{-12}	4.305×10^{-18}	6.931×10^{-16}	1.735×10^{-18}	3.469×10^{-18}
2	3.331×10^{-16}	5.105×10^{-16}	7.932×10^{-18}	1.656×10^{-17}	2.168×10^{-18}	3.036×10^{-18}
3	4.857×10^{-17}	1.073×10^{-18}	1.214×10^{-17}	2.630×10^{-17}	3.469×10^{-18}	3.469×10^{-18}
4	5.204×10^{-18}	9.422×10^{-18}	1.741×10^{-18}	6.651×10^{-18}	1.301×10^{-18}	1.735×10^{-18}

表 3-2

テスト 1		テスト 2	
$\bar{\sigma}_k$	$(\sigma_k - \bar{\sigma}_k) / \sigma_k$	$\bar{\sigma}_k$	$(\sigma_k - \bar{\sigma}_k) / \sigma_k$
8.8881582×10^6	2.197	1.7347876×10^9	0.8531
6.9916152×10^4	-5.753	6.4858294×10^1	3.352
1.2492552×10^2	45.24×10^{-8}	1.0034119×10^1	-0.1694×10^{-8}
1.8916602	38660	9.0089682×10^{-1}	75.01
3.8969962×10^1	-703.0	1.5990873×10^{-1}	137.6

テスト 3		テスト 4	
$\bar{\sigma}_k$	$(\sigma_k - \bar{\sigma}_k) / \sigma_k$	$\bar{\sigma}_k$	$(\sigma_k - \bar{\sigma}_k) / \sigma_k$
3.5327044×10^1	-1.514	8.2423568	-1.329×10^{-8}
1.7087926×10^{-7}	-17.09 *	1.0000000	0.0
2.0000001×10^1	-5.000×10^{-8}	2.2502348	-4.684×10^{-8}
1.9595918×10^1	-0.2960	1.0000000	0.0
9.3931860×10^{-7}	-93.93 *	1.0000000	0.0

表 3-3

\bar{v}_2	\bar{v}_5	$v_2 - \bar{v}_2$	$v_5 - \bar{v}_5$
-0.41909547	0.0	-1.5	0.0 (Def)
0.44050911	0.41854808	1.3	-1.6
-0.052004562	0.34879006	1.3	-0.7×10^{-8}
0.67605914	0.24415303	0.0	0.7
0.41297729	-0.80221714	1.3	1.8

表 3-4

単精度

テスト 1

	x_{1S}	x_{2S}
	0.97162955	0.067334449
	0.49080089	0.19346116
	0.32945643	0.20308967
	0.24832276	0.19333959
	0.19940837	0.17993266
a_m	2.837×10^{-2}	0.9327
a_f	3.013×10^{-2}	0.9922
s_m	2.837×10^{-2}	0.9327
s_f	3.650×10^{-2}	1.208
r	1.207×10^{-2}	7.255×10^7
T_S		3.896
$-\log s_m$	1.547	3.027×10^{-2}
$-\log s_f$	1.438	-8.215×10^{-2}

倍精度

	x_{1D}	x_{2D}
	1.000000000000001649	1.00000000003793734
	0.500000000000004691	0.500000000012505591
	0.333333333333335179	0.333333333338655164
	0.250000000000000764	0.25000000002317488
	0.200000000000000262	0.20000000000821350
a_m	1.649×10^{-14}	3.794×10^{-11}
a_f	1.726×10^{-14}	4.037×10^{-11}
s_m	1.649×10^{-14}	3.794×10^{-11}
s_f	2.004×10^{-14}	4.922×10^{-11}
r	3.116×10^{-23}	7.255×10^7
T_S		14.43
$-\log s_m$	13.78	10.42
$-\log s_f$	13.70	10.31

表3-5

単精度

テスト2

	x_{1S}	x_{2S}	x_{3S}
	1.0000035	-6	0.99386187
	2.0000039	-7	1.9932397
	-0.99999948	1×10^{-8}	-0.99882910
	3.0000077	-15	2.9854727
	-4.0000067	24	-3.9760442
a_m	7.689×10^{-6}	2.410×10^{-2}	2.396×10^{-2}
a_f	1.147×10^{-6}	2.972×10^{-2}	2.949×10^{-2}
s_m	3.487×10^{-6}		6.138×10^{-3}
s_f	5.065×10^{-6}		1.048×10^{-2}
τ	1.824×10^{-10}	2.645×10^8	2.645×10^8
T_S			3.616
$-\log s_m$	5.458		2.212
$-\log s_f$	5.295		1.980

倍精度

	x_{1D}	x_{2D}	x_{3D}
	0.9999999999999999	5	1.00000000000004893
	2.0000000000000000	5	2.00000000000005445
	-0.9999999999999999	-1×10^{-14}	-1.00000000000001384
	2.9999999999999998	12	3.00000000000012056
	-4.0000000000000000	-22	-4.00000000000021198
a_m	1.908×10^{-17}	2.155×10^{-13}	2.120×10^{-13}
a_f	1.922×10^{-17}	2.582×10^{-13}	2.550×10^{-13}
s_m	6.366×10^{-18}		5.300×10^{-14}
s_f	6.613×10^{-18}		8.804×10^{-14}
τ	2.829×10^{-31}	2.645×10^8	2.645×10^8
T_D			14.15
$-\log s_m$			13.28
$-\log s_f$			13.06

表3-6

テスト3

単精度

	x_{1S}	x_{2S}	x_{3S}
	-0.083333334	-3	-0.083333334
	-2.346×10^{-9}	-11	-3.668×10^{-9}
	0.25000001	22×10^{-10}	0.25000001
	-0.083333349	4	-0.083333347
	0.083333336	5	0.083333336
a_m	1.490×10^{-8}	2.193×10^{-9}	1.490×10^{-8}
a_f	2.129×10^{-8}	2.550×10^{-9}	2.022×10^{-8}
s_m	1.788×10^{-7}		1.565×10^{-7}
s_f	1.898×10^{-7}		1.690×10^{-7}
τ	1.372×10^{-18}	320	320
T_S			6.574
$-\log s_m$	6.748		6.806
$-\log s_f$	6.722		6.772

倍精度

	x_{1D}	x_{2D}	x_{3D}
	-0.083333333333333334 2.169×10^{-20} 0.250000000000000001 -0.083333333333333333 0.083333333333333336	1 -6 8×10^{-19} 4 2	-0.083333333333333334 -4.522×10^{-19} 0.250000000000000002 -0.083333333333333330 0.083333333333333338
a_m a_f	6.505×10^{-19} 7.154×10^{-19}	8.315×10^{-19} 1.153×10^{-18}	1.518×10^{-18} 1.690×10^{-18}
s_m s_f	3.249×10^{-18} 4.407×10^{-18}		6.068×10^{-18} 9.324×10^{-18}
r	1.945×10^{-34}	320	320
T_D -log s_m -log s_f	17.49 17.36		17.11 17.22 17.03

表 3-7

テスト 4

単精度

	x_{1S}	x_{2S}	x_{3S}
	5.0000004 3.9999996 2.9999997 2.0000008 1.0000005	5.0000004 3.9999999 2.9999999 2.0000006 1.0000002	5.0000006 3.9999998 2.9999999 2.0000004 1.0000002
a_m a_f	8.345×10^{-7} 1.126×10^{-6}	5.960×10^{-7} 7.468×10^{-7}	5.960×10^{-7} 8.205×10^{-7}
s_m s_f	4.768×10^{-7} 6.515×10^{-7}	2.980×10^{-7} 3.900×10^{-7}	2.086×10^{-7} 3.262×10^{-7}
r	4880	2577	1913
T_S -log s_m -log s_f	5.983 6.322 6.186	6.121 6.526 6.409	6.186 6.681 6.487

6章 結 語

4章でも述べたように, overotetermined system

$$Ax = b$$

を最小自乗法で解く最も普通の計算方法は,

$$A^T Ax = A^T b \tag{6.1}$$

を解くものである。こうすると, $A^T A$ は正値対称行列になるので, これを解くには「1次方程式系の解法 I」⁽¹⁾で述べた方法が採用できる。このときに失われる精度は, $A^T A$ の条件数-行列が正定値の場合の条件数は(最大固有値) ÷ (最小固有値) - のオーダー程度である。(cf. 「1次方程式系の解法 I」⁽¹⁾) ここで用いたテスト行列(1), (2), (4)に対する条件数はそれぞれ, 2.2×10^{13} , 1.2×10^6 , 6.8×10^1 であるから, (6.1) という方程式を「1次方程式系の解法 I」のサブルーチンで解いたと

きに失われる精度は, それぞれ大略13桁, 6桁, 2桁と予想される。(実際の結果は表5-2を参照) 上記の数字を本報告のサブルーチンによるものと比較すれば下表のようになる。(数字は失われた精度の桁数)

表 4-1

	テスト1	テスト2		テスト4
	倍精度	単精度	倍精度	単精度
LSQHDS	8.1~8.2	4.7~4.9	5.0~5.3	0.6~0.8
LSQE2S	6.4~6.5	4.8~5.0	5.6~5.9	0.3~0.4
MINFIS	7.9~8.1	5.6~5.8	5.1~5.3	1.1~1.3

これを見ると, 精度についてならテスト1の行列では, (6.1)を解くのは圧倒的に不利であるのがわかる。もう

一つ(6.1)を解くことの欠点は、行列の掛け算を計算する必要があるため計算量が多くなることである。このことはとりわけ行列の次数が大きくなったときに影響が大きくなる。($m \times n$ 行列 A について $A^T A$ を計算するには約 $2n^2 m$ 回の積と和の計算が必要である)

しかし、計算量を不問にし、 $A^T A$ 、 $A^T b$ が精度よく計算できれば、有効数字2桁程度の近似解が得られると、(6.1)の解はイタレーションによってかなりの精度まで求められる。ところが、本報告のルーチンではイタレーションした結果が表4-1である。このサブルーチンでは、残差が0に近い場合には殆んど正しい値が得られているが、残差が大きくなると有効数字の桁数が減っている。この点について改良しなければ、次数の小さい問題については(6.1)をイタレーションで解く方が無難だったということになってしまいます。そのためは、それがプログラム技術的な改良によって可能なのか、或は方法を改良する必要があるのか、それとも又別な方法を開発しなければならないのか、といった検討等が今後の課題である。

なお、「1次方程式系の解法 I」のサブルーチン『CLDT2S』と『CLSL2S』を用いて(6.1)を解いた結果は下表の通りである。又、イタレーションを行っているサブルーチン『ACLSLS』を用いた場合には、下表にあるテスト全部について正しい解が得られた。

表4-2

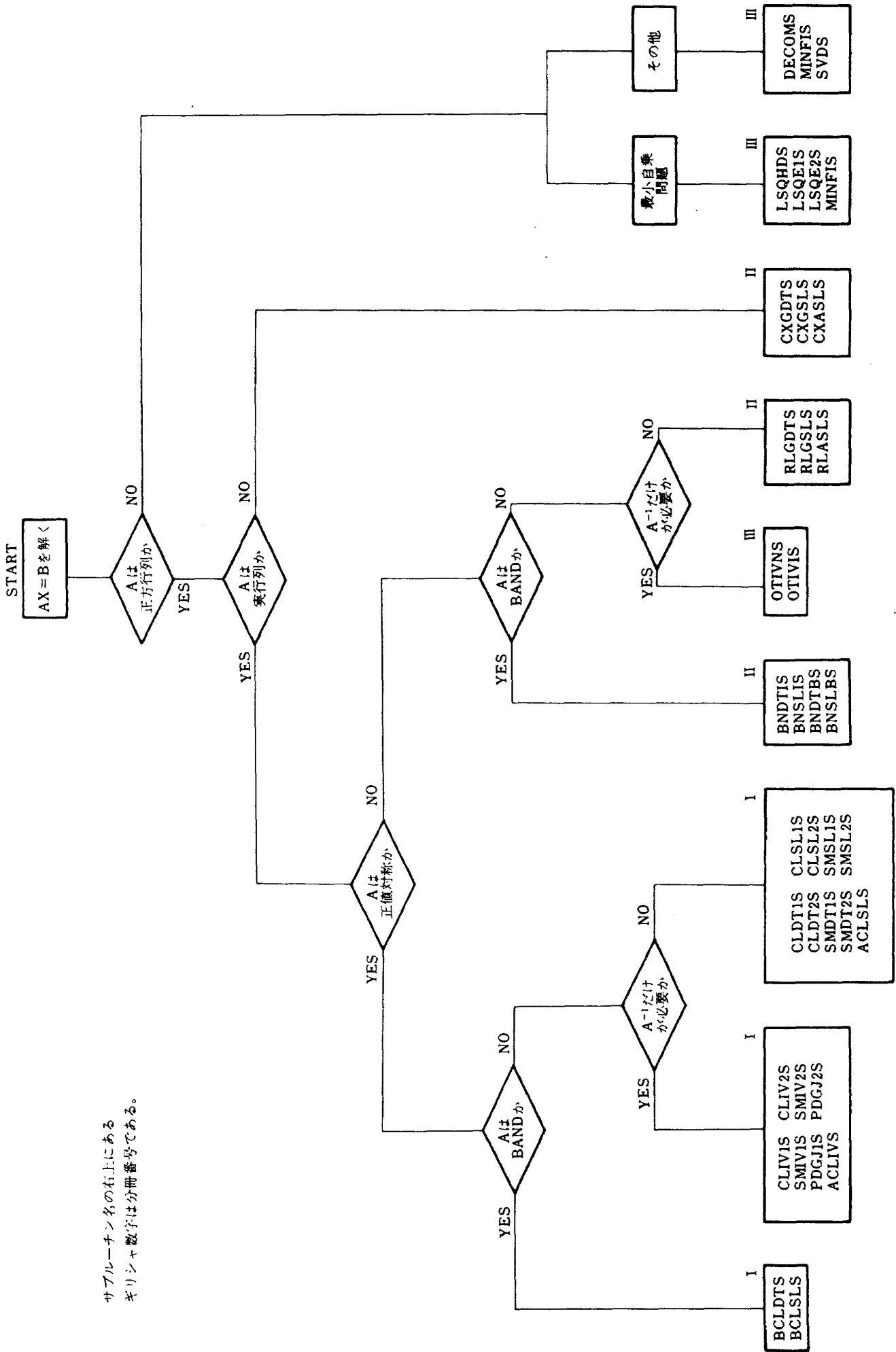
	テスト1	テスト2		テスト1
	倍精度	単精度	倍精度	単精度
$-\log s_m$	8.917	2.965	13.99	6.097
$-\log s_f$	8.809	2.746	13.76	6.035
失われた精度	9.4~9.6	4.9~5.1	4.4~4.6	1.7~1.8

テスト1の単精度は『CLDT2S』の中で計算続行が不可能となった。($ILL=1$ cf. 「1次方程式系の解法 I」) 又、テスト1の右辺 b_1 、 b_2 に対する答として同じ値が出力されてきた。テスト2の右辺 b_1 、 b_2 、テスト4の右辺 b_1 、 b_2 、 b_3 に対しても同様である。さらにテスト2の右辺 b_2 に対しては正確に零ベクトルが出力されてきた。これらのことと、『ACLSLS』を使用した場合に全部正しい値が得られたこととは、 $A^T A$ 、 $A^T b$ の計算が桁落ちなく正確に行われた結果によるものと思われる。

本報告によって1次方程式系を取り扱うサブルーチンについては一応終了する。といっても完了する性格のものではなく、新しい手法の開発や異なる問題へのアプローチ等といった研究もあり、持続的研究が要求されるものである。なおこの後、行列の固有値、固有ベクトルを求める問題、さらには特殊函数の値を求めるルーチンの作成へと進む予定である。

参 考 文 献

- 1) 福田正大他；1次方程式系の解法 I，航技研資料 TM-277 (1975.5)



サブルーチン名の右上にある
ギリシヤ数字は分冊番号である。

航空宇宙技術研究所資料 315号

昭和51年10月発行

発行所 航空宇宙技術研究所
東京都調布市深大寺町1880
電話武蔵野三鷹(0422)47-5911(大代表)〒182
印刷所 株式会社 共 進
東京都杉並区久我山4-1-7(羽田ビル)
