

舍入误差分析

教师: 胡俊

北京大学数学科学学院

October 18, 2019

- 1 基本运算的舍入误差分析
- 2 列主元Gauss消去法的舍入误差分析
- 3 计算解的精度估计和迭代改进

1 基本运算的舍入误差分析

2 列主元Gauss消去法的舍入误差分析

3 计算解的精度估计和迭代改进

计算机中浮点数 f 表示为:

$$f = \pm w \times \beta^J, \quad L \leq J \leq U,$$

这里 β 是机器所用浮点数的基底, J 是阶, w 是尾数. 尾数 w 一般可表示为

$$w = 0.d_1d_2\cdots d_t,$$

其中 t 是尾数位, 称为字长, $0 \leq d_i < \beta$. 若 $d_1 \neq 0$, 则称该浮点数为规格化浮点数.

若用 \mathcal{F} 表示一个系统的浮点数全体构成的集合, 则

$$\mathcal{F} = \{0\} \cup \{f : f = \pm 0.d_1d_2\cdots d_t \times \beta^J, 0 \leq d_i < \beta, \\ d_1 \neq 0, L \leq J \leq U\}$$

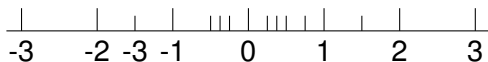
显然, 集合 \mathcal{F} 可用四元数组 (β, t, L, U) 来刻画. 机器不同, 这四个值亦不同, 较典型的值是 $(2, 56, -63, 64)$.

集合 \mathcal{F} 是一个包含 $2(\beta - 1)\beta^{t-1}(U - L + 1) + 1$ 个数的有限集, 这些数对称地分布在区间 $[m, M]$ 和 $[-M, -m]$ 中, 其中

$$m = \beta^{L-1}, \quad M = \beta^U(1 - \beta^{-t}). \quad (1)$$

值得注意的是, 这些数在 $[m, M]$ 和 $[-M, -m]$ 中的分布是不等距的. 例如, 若 $\beta = 2, t = 2, L = -1$ 和 $U = 2$, 则 \mathcal{F} 中17个数的分布如图1所示.

Figure



记实数 x 的浮点数表示为 $fl(x)$. 若 $x = 0$, 则 $fl(x)$ 取为零.

若 $m \leq |x| \leq M$, 则当使用舍入法时, 取 $fl(x)$ 为 \mathcal{F} 中最接近于 x 的数 f ;
当使用截断法时, 取 $fl(x)$ 为满足 $|f| \leq |x|$ 且最接近于 x 的数 f . 例如,
对 $(\beta, t, L, U) = (10, 3, 0, 2)$ 和实数 $x = 5.45627$, 若用舍入法, 则
有 $fl(x) = 0.546 \times 10$; 若用截断法, 则有 $fl(x) = 0.545 \times 10$.

定理 2.3.1: 设 $m \leq |x| \leq M$, 其中 m 和 M 由 (1)式定义, 则

$$fl(x) = x(1 + \delta), \quad |\delta| \leq \mathbf{u} \quad (2)$$

其中 \mathbf{u} 为 机器精度, 即

$$\mathbf{u} = \begin{cases} \frac{1}{2}\beta^{1-t}, & \text{用舍入法} \\ \beta^{1-t}, & \text{用截断法} \end{cases}$$

证明： 现不妨假定 $x > 0$ (因若 $x < 0$, 证明完全类似). 设 α 是满足

$$\beta^{\alpha-1} \leq x < \beta^\alpha \quad (3)$$

的唯一整数. 在 $[\beta^{\alpha-1}, \beta^\alpha)$ 中浮点数的阶为 α , 所以在这个区间中所有 t 位的浮点数以间距 $\beta^{\alpha-t}$ 分布. 对于舍入法, 根据(3)式, 有

$$|fl(x) - x| \leq \frac{1}{2}\beta^{\alpha-t} = \frac{1}{2}\beta^{\alpha-1}\beta^{1-t} \leq \frac{1}{2}x\beta^{1-t},$$

即

$$\frac{|fl(x) - x|}{x} \leq \frac{1}{2}\beta^{1-t}; \quad (4)$$

对于截断法, 有

$$|fl(x) - x| \leq \beta^{\alpha-1}\beta^{1-t} \leq x\beta^{1-t}, \quad (5)$$

即

$$\frac{|fl(x) - x|}{x} \leq \beta^{1-t}, \quad (6)$$

这就证明了(2)式.

设 $a, b \in \mathcal{F}$ 是两个给定的浮点数, 用 \circ 来表示 $+, -, \times, /$ 中任意一种运算. $fl(a \circ b)$ 的意义是先进行运算, 得到精确的实数, 再按舍入规则表示成浮点数, 在运算中, 若出现 $|a \circ b| > M$ 或 $0 < |a \circ b| < m$, 则就是发生了上溢或下溢. 在不发生溢出的情况下, 由定理2.3.1立即得到如下定理:

定理 2.3.2: 设 $a, b \in \mathcal{F}$, 则

$$fl(a \circ b) = (a \circ b)(1 + \delta), |\delta| < u.$$

例2.3.1 设 x, y 是两个由浮点数构成的 n 维向量. 试估计

$$|fl(x^T y) - x^T y|$$

的上界.

解: 令

$$S_k = fl\left(\sum_{i=1}^k x_i y_i\right),$$

则由定理2.3.2, 可得

$$S_1 = x_1 y_1 (1 + \gamma_1), \quad |\gamma_1| \leq \mathbf{u},$$

$$S_k = fl(S_{k-1} + fl(x_k y_k))$$

$$= [S_{k-1} + x_k y_k (1 + \gamma_k)] (1 + \delta_k), \quad |\delta_k|, |\gamma_k| \leq \mathbf{u},$$

于是有

$$\begin{aligned} fl(x^T y) = S_n &= \sum_{i=1}^n x_i y_i (1 + \gamma_i) \prod_{j=i}^n (1 + \delta_j) \\ &= \sum_{i=1}^n (1 + \varepsilon_i) x_i y_i, \end{aligned} \tag{7}$$

其中

$$1 + \varepsilon_i = (1 + \gamma_i) \prod_{j=i}^n (1 + \delta_j),$$

这里定义 $\delta_1 = 0$. 这样, 就有

$$|fl(x^T y) - x^T y| \leq \sum_{i=1}^n |\varepsilon_i| |x_i y_i| \leq 1.01n\mathbf{u} \sum_{i=1}^n |x_i y_i|,$$

其中最后一个不等式用了下面将要证明的定理2.3.3的结论. 注意上式表明, 若 $|x^T y| \ll \sum_{i=1}^n |x_i y_i|$, 则 $fl(x^T y)$ 的相对误差可能会很大.

向量内积运算舍入误差分析

定理 **2.3.3**: 若 $|\delta_i| \leq \mathbf{u}$ 且 $n\mathbf{u} \leq 0.01$, 那么

$$1 - n\mathbf{u} \leq \prod_{i=1}^n (1 + \delta_i) \leq 1 + 1.01n\mathbf{u},$$

或者写成

$$\prod_{i=1}^n (1 + \delta_i) = 1 + \delta, \quad |\delta| \leq 1.01n\mathbf{u}.$$

证明: 因为 $\delta_i \leq \mathbf{u}$, 故有

$$(1 - \mathbf{u})^n \leq \prod_{i=1}^n (1 + \delta_i) \leq (1 + \mathbf{u})^n. \quad (8)$$

先估计 $(1 - \mathbf{u})^n$ 的下界. 利用函数 $(1 - x)^n$ ($0 < x < 1$)的Taylor展开

$$(1 - x)^n = 1 - nx + \frac{n(n-1)}{2}(1 - \theta x)^{n-2}x^2,$$

得到

$$1 - nx \leq (1 - x)^n.$$

由此立即得到

$$1 - 1.01n\mathbf{u} \leq 1 - n\mathbf{u} \leq (1 - \mathbf{u})^n \quad (9)$$

下面估计 $(1 + \mathbf{u})^n$ 的上界. 由 e^x 的幂级数展开

$$\begin{aligned} e^x &= 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots \\ &= 1 + x + \frac{x}{2} \cdot x \cdot \left(1 + \frac{x}{3} + \frac{2x^2}{4!} + \cdots \right) \end{aligned}$$

立即知道, 当 $0 \leq x \leq 0.01$ 时, 有

$$1 + x \leq e^x \leq 1 + x + \frac{0.01}{2}xe^x \leq 1 + 1.01x, \quad (10)$$

这里用到了 $e^{0.01} < 2$ 这一事实. 令 $x = \mathbf{u}$, 由(10)式左端不等式得

$$(1 + \mathbf{u}) \leq e^{\mathbf{u}} \quad (11)$$

令 $x = n\mathbf{u}$, 由(11)式右端不等式得

$$e^{n\mathbf{u}} \leq 1 + 1.01n\mathbf{u}. \quad (12)$$

综合(11)式和(12)式得

矩阵运算舍入误差分析

下面分析一下矩阵基本运算的舍入误差, 这些结果对以后的讨论是有用的. 为此, 引进记号

$$|E| = [|e_{ij}|], \quad E \in \mathbf{R}^{n \times n}$$

并且规定

$$|E| \leq |F| \quad \text{当且仅当} \quad |e_{ij}| \leq |f_{ij}|, i, j = 1, \dots, n.$$

设 A, B 是由 \mathcal{F} 中的元素构成的 $n \times n$ 矩阵, 且 $\alpha \in \mathcal{F}$, 则由定理2.3.2易知

$$\begin{aligned} fl(\alpha A) &= \alpha A + E, \quad |E| \leq \mathbf{u}|\alpha A|, \\ fl(A + B) &= (A + B) + E, \quad |E| \leq \mathbf{u}|A + B| \end{aligned}$$

此外, 应用例2.3.1可得

$$fl(AB) = AB + E, \quad |E| \leq 1.01n\mathbf{u}|A||B|.$$

注意, $|AB|$ 可能比 $|A||B|$ 小得多, 因此矩阵乘积的相对误差未必很小.

基于这个原因, 通常计算内积时是先用双精度(即字长为 $2t$)计算, 最后再把计算结果舍入为单精度数(即字长为 t).

上述的这三个矩阵基本运算的舍入误差界, 是通过估计计算解与精确解之间的误差得到的, 舍入误差的界与精确解有关. 这种误差分析的方法称为**向前误差分析法**.

实际上常用的是**向后误差分析法**. 为了说明这种误差分析方法, 看一个简单的例子. 假定上面所述的矩阵 A, B 是 2×2 的上三角矩阵, 则由定理2.3.1可知

$$fl(AB) = \begin{pmatrix} a_{11}b_{11}(1 + \varepsilon_1) & \widetilde{a}_{12} \\ 0 & a_{22}b_{22}(1 + \varepsilon_5) \end{pmatrix}$$

其中

$$\begin{aligned} \widetilde{a}_{12} &= (a_{11}b_{12}(1 + \varepsilon_2) + a_{12}b_{22}(1 + \varepsilon_3))(1 + \varepsilon_4), \\ |\varepsilon_i| &\leq \mathbf{u}, \quad i = 1, 2, 3, 4, 5. \end{aligned}$$

若令

$$\widetilde{A} = \begin{pmatrix} a_{11} & a_{12}(1 + \varepsilon_3)(1 + \varepsilon_4) \\ 0 & a_{22}(1 + \varepsilon_5) \end{pmatrix},$$
$$\widetilde{B} = \begin{pmatrix} b_{11}(1 + \varepsilon_1) & b_{12}(1 + \varepsilon_2)(1 + \varepsilon_4) \\ 0 & b_{22} \end{pmatrix}$$

则易证

$$fl(AB) = \widetilde{A}\widetilde{B},$$

而且

$$\begin{aligned} \widetilde{A} &= A + E, & |E| &\leq 3\mathbf{u}|A|, \\ \widetilde{B} &= B + F, & |F| &\leq 3\mathbf{u}|B|. \end{aligned}$$

换句话说, 计算得到的乘积 $fl(AB)$ 是有了微小扰动的两个矩阵 \widetilde{A} 和 \widetilde{B} 的精确的乘积.

- 1 基本运算的舍入误差分析
- 2 列主元Gauss消去法的舍入误差分析
- 3 计算解的精度估计和迭代改进

矩阵三角分解的舍入误差分析

引理 2.4.1: 设 $n \times n$ 浮点数矩阵 $A = [a_{ij}]$ 有三角分解, 且 $1.01n\mathbf{u} \leq 0.01$, 则用 Gauss 消去法计算得到的单位下三角矩阵 L 和上三角矩阵 U 满足

$$\widetilde{L}\widetilde{U} = A + E \quad (14)$$

其中

$$|E| \leq 2.05n\mathbf{u}|\widetilde{L}||\widetilde{U}|. \quad (15)$$

证明: 设 $\widetilde{U} = [\widetilde{u}_{ij}]$, $\widetilde{L} = [\widetilde{l}_{ij}]$, 由 Gauss 消去法的具体实现知, $\widetilde{u}_{ij} (i \leq j)$ 是从 a_{ij} 中依次减去 $\widetilde{l}_{ik}\widetilde{u}_{kj} (k = 1, \dots, i-1)$ 而得到的, 即

$$\begin{aligned} a_{ij}^{(0)} &= a_{ij}, \\ a_{ij}^{(k)} &= fl(a_{ij}^{(k-1)} - fl(\widetilde{l}_{ik}\widetilde{u}_{kj})), \quad k = 1, \dots, i-2, \\ \widetilde{u}_{ij} &= a_{ij}^{(i-1)} = fl(a_{ij}^{(i-2)} - fl(\widetilde{l}_{i,i-1}\widetilde{u}_{i-1,j})) \end{aligned}$$

由基本运算舍入误差分析的基本结果可得

$$a_{ij}^{(k)} = \left(a_{ij}^{(k-1)} - (\tilde{l}_{ik} \tilde{u}_{kj}) (1 + \gamma_k) \right) (1 + \varepsilon_k),$$

其中 $|\gamma_k| \leq \mathbf{u}, |\varepsilon_k| \leq |\mathbf{u}| (k = 1, \dots, i-1)$. 由此出发, 利用与例2.3.1同样的推理方法可得

$$\tilde{u}_{ij} = a_{ij} (1 + \delta_i) - \sum_{k=1}^{i-1} (\tilde{l}_{ik} \tilde{u}_{kj}) (1 + \delta_k) \quad (16)$$

其中 $|\delta_k| \leq 1.01n\mathbf{u} (k = 1, \dots, i)$. 从(16)中将 a_{ij} 解出来, 便有

$$\begin{aligned} a_{ij} &= \frac{\tilde{u}_{ij}}{1 + \delta_i} + \sum_{k=1}^{i-1} (\tilde{l}_{ik} \tilde{u}_{kj}) \frac{1 + \delta_k}{1 + \delta_i} \\ &= \sum_{k=1}^i \tilde{l}_{ik} \tilde{u}_{kj} - e_{ij} \end{aligned} \quad (17)$$

其中 $\widetilde{l}_{ii} = 1$, 而

$$e_{ij} = (\widetilde{l}_{ii}\widetilde{u}_{ij}) \frac{\delta_i}{1 + \delta_i} + \sum_{k=1}^{i-1} (\widetilde{l}_{ik}\widetilde{u}_{kj}) \frac{\delta_i - \delta_k}{1 + \delta_i} \quad (18)$$

注意到 $|\delta_k| \leq 1.01n\mathbf{u} < 1.01$, 有

$$|e_{ij}| \leq \frac{2.02n\mathbf{u}}{1 - 0.01} \sum_{k=1}^i |\widetilde{l}_{ik}| |\widetilde{u}_{kj}| \leq 2.05n\mathbf{u} \sum_{k=1}^i |\widetilde{l}_{ik}| |\widetilde{u}_{kj}|. \quad (19)$$

此外, 由Gauss消去法的具体实现知, $\widetilde{l}_{ij} (i > j)$ 的计算过程如下

$$\begin{aligned} a_{ij}^{(0)} &= a_{ij}, \\ a_{ij}^{(k)} &= fl\left(a_{ij}^{(k-1)} - fl(\widetilde{l}_{ik}\widetilde{u}_{kj})\right), \quad k = 1, 2, \dots, j-1 \\ \widetilde{l}_{ij} &= fl\left(a_{ij}^{(j-1)} / \widetilde{u}_{jj}\right). \end{aligned}$$

由基本运算的舍入误差分析的基本结果可得

$$\begin{aligned} a_{ij}^{(k)} &= \left(a_{ij}^{(k-1)} - (\widetilde{l}_{ik} \widetilde{u}_{kj}) (1 + \gamma_k) \right) (1 + \varepsilon_k) \\ \widetilde{l}_{ij} &= a_{ij}^{(j-1)} / [\widetilde{u}_{jj} (1 + \delta).] \end{aligned}$$

其中 $|\delta| \leq \mathbf{u}, |\gamma_k| \leq \mathbf{u}, |\varepsilon_k| \leq \mathbf{u}$ ($k = 1, \dots, j-1$) 由此出发, 完全类似(17)式和(19)式的证明, 可证

$$a_{ij} = \sum_{k=1}^j \widetilde{l}_{ik} \widetilde{u}_{kj} - e_{ij} \quad (20)$$

其中

$$|e_{ij}| \leq 2.05n\mathbf{u} \sum_{k=1}^j |\widetilde{l}_{ik}| |\widetilde{u}_{kj}| \quad (21)$$

综合上面所证, 即知引理的结论成立.

注意到交换矩阵的行或列并不引进舍入误差, 由引理2.4.1立即得到如下结论:

推论2.4.1 设 $n \times n$ 浮点数矩阵 A 是非奇异的, 且 $1.01n\mathbf{u} \leq 0.01$, 则用列主元Gauss消去法计算得到的单位下三角阵 \widetilde{L} , 上三角阵 \widetilde{U} 以及排列方阵 \widetilde{P} 满足

$$\widetilde{L}\widetilde{U} = \widetilde{P}A + E,$$

其中

$$|E| \leq 2.05n\mathbf{u}|\widetilde{L}||\widetilde{U}|.$$

当对 A 完成 LU 分解之后, 求解线性方程组 $Ax = b$ 的问题就归为求解两个三角形方程组

$$\widetilde{L}y = \widetilde{P}b \quad \text{和} \quad \widetilde{U}x = y$$

的问题. 所以现在来估计求解三角形方程组的舍入误差.

求解三角方程舍入误差分析

引理2.4.2: 设 $n \times n$ 浮点数三角阵 S 是非奇异的, 并且假定 $1.01n\mathbf{u} \leq 0.01$, 则用第一章所介绍的解三角方程组的方法求解 $Sx = b$ 所得到的计算解 \tilde{x} 满足

$$(S + H)\tilde{x} = b, \quad (22)$$

其中

$$|H| \leq 1.01n\mathbf{u}|S|, \quad (23)$$

证明: 对 n 用数学归纳法. 不失一般性, 设 $S = L$ 是下三角阵.

当 $n = 1$ 时, 引理显然成立. 假设对所有的 $n - 1$ 阶下三角形方程组已证引理成立. 考虑 n 阶的情形. 假定用前代法解 $Lx = b$ 的计算解为 \tilde{x} , 并将 L, b 和 \tilde{x} 分块如下:

$$b = \begin{bmatrix} b_1 \\ c \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}, \quad \tilde{x} = \begin{bmatrix} \tilde{x}_1 \\ \tilde{y} \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}.$$

$$L = \begin{bmatrix} l_{11} & 0 \\ l_1 & L_1 \end{bmatrix} \begin{matrix} 1 \\ n-1 \\ 1 & n-1 \end{matrix}$$

由定理2.3.2有

$$\widetilde{x}_1 = fl\left(\frac{b_1}{l_{11}}\right) = \frac{b_1}{l_{11}(1 + \delta_1)}, \quad |\delta_1| \leq \mathbf{u} \quad (24)$$

此外, 注意到是用前代法求解 $n - 1$ 阶方程组

$$L_1 y = fl(c - \widetilde{x}_1 l_1)$$

所得到的计算解, 由归纳法假设即有

$$(L_1 + H_1)\widetilde{y} = fl(c - \widetilde{x}_1 l_1),$$

其中

$$|H_1| \leq 1.01(n - 1)\mathbf{u} |L_1| \quad (25)$$

$$fl(c - \widetilde{x}_1 l_1) = fl(c - fl(\widetilde{x}_1 l_1)) = (I + D_\gamma)^{-1} (c - \widetilde{x}_1 l_1 - \widetilde{x}_1 D_\delta l_1),$$

其中

$$D_\gamma = \text{diag}(\gamma_2, \dots, \gamma_n), \quad D_\delta = \text{diag}(\delta_2, \dots, \delta_n) \\ |\gamma_i| \leq \mathbf{u}, \quad |\delta_i| \leq \mathbf{u}, \quad i = 2, \dots, n.$$

于是,

$$\widetilde{x}_1 l_1 + \widetilde{x}_1 D_\delta l_1 + (I + D_\gamma)(L_1 + H_1)\widetilde{y} = c,$$

从而有

$$(L + H)\widetilde{x} = b,$$

其中

$$H = \begin{bmatrix} \delta_1 l_{11} & 0 \\ D_\delta l_1 & H_1 + D_\gamma(L_1 + H_1) \end{bmatrix}.$$

由(24)式和(25)式得

$$\begin{aligned} |H| &\leq \begin{bmatrix} |\delta_1| |l_{11}| & 0 \\ |D_\delta| |l_1| & |H_1| + |D_\gamma| (|L_1| + |H_1|) \end{bmatrix} \\ &\leq \begin{bmatrix} u |l_{11}| & 0 \\ \mathbf{u} |l_1| & |H_1| + \mathbf{u} (|L_1| + |H_1|) \end{bmatrix} \\ &\leq u \begin{bmatrix} |l_{11}| & 0 \\ |l_1| & (1.01(n-1) + 1 + 1.01(n-1)\mathbf{u}) |L_1| \end{bmatrix} \\ &\leq 1.01n\mathbf{u}|L|, \end{aligned}$$

其中最后一个不等式用到了假设条件 $1.01n\mathbf{u} \leq 0.01$.

列主元高斯消去法舍入误差分析

应用引理2.4.2得到三角形方程组

$$\widetilde{L}y = \widetilde{P}b \quad \text{和} \quad \widetilde{U}x = y$$

即知最后得到的计算解 \widetilde{x} 应满足

$$(\widetilde{L} + F)(\widetilde{U} + G)\widetilde{x} = \widetilde{P}b$$

即

$$(\widetilde{L}\widetilde{U} + F\widetilde{U} + \widetilde{L}G + FG)\widetilde{x} = \widetilde{P}b \quad (26)$$

其中

$$|F| \leq 1.01n\mathbf{u}|\widetilde{L}|, \quad |G| \leq 1.01n\mathbf{u}|\widetilde{U}|. \quad (27)$$

再将 $\widetilde{L}\widetilde{U} = \widetilde{P}A + E$ 带入(26)式, 得

$$(A + \delta A)\widetilde{x} = b,$$

这里

$$\delta A = \widetilde{P}^T(E + F\widetilde{U} + \widetilde{L}G + FG).$$

由(27)式和推论2.4.1 得

$$|\delta A| \leq 4.09n\mathbf{u}\widetilde{P}^T|\widetilde{L}||\widetilde{U}|.$$

注意到 \tilde{L} 的元素的绝对值均不超过1, 故有

$$\|\tilde{L}\|_{\infty} \leq n.$$

为了给出 $\|\tilde{U}\|_{\infty}$ 的估计, 定义

$$\rho = \max_{i,j} |\tilde{u}_{ij}| / \max_{i,j} |a_{ij}|,$$

通常称之为列主元 **Gauss** 消去法的增长因子. 于是,

$$\|\tilde{U}\|_{\infty} \leq n \max_{i,j} |\tilde{U}_{ij}| = n\rho \max_{i,j} |a_{ij}| \leq n\rho \|A\|_{\infty}.$$

这样, 就得到了本节的主要定理:

定理2.4.1: 设 $n \times n$ 浮点数矩阵 A 是非奇异的, 且 $1.01n\mathbf{u} \leq 0.01$, 则用列主元 **Gauss** 消去法解线性方程组 $Ax = b$ 所得到的计算解 \tilde{x} 满足

$$(A + \delta A)\tilde{x} = b \tag{28}$$

其中

$$\|\delta A\|_{\infty} / \|A\|_{\infty} \leq 4.09n^3 \rho \mathbf{u}. \tag{29}$$

定理2.4.1表明, 由于消去法求解过程中引进舍入误差而产生的计算解相当于系数矩阵作某些扰动而得到的扰动方程组的精确解. 一般来说, δA 的元素比起 A 的元素的初始误差 (数据的测量误差、数学模型误差等) 来是很小的. 在这个意义上讲, 列主元Gauss消去法是数值稳定的.

最后需指出的是, 理论上可以证明 $\rho \leq 2^{n-1}$, 而且上界可以达到. 但在实际运算时, 常遇到的问题其 ρ 很小, 一般来讲不会超过 n . 此外, 定理2.4.1中所给出的上界 $4.09n^3\rho u$ 一般要比真正的 $\|\delta A\|_\infty/\|A\|_\infty$ 大很多. 在实际计算时, 常遇到的问题几乎都有 $\|\delta A\|_\infty/\|A\|_\infty \approx u$.

- 1 基本运算的舍入误差分析
- 2 列主元Gauss消去法的舍入误差分析
- 3 计算解的精度估计和迭代改进

设用某种计算方法求解线性代数方程组 $Ax = b$ 得到的计算解为 \hat{x} .
令

$$r = b - A\hat{x},$$

则有

$$r = Ax - A\hat{x} = A(x - \hat{x}).$$

于是

$$\|x - \hat{x}\| = \|A^{-1}r\| \leq \|A^{-1}\| \|r\|.$$

再注意到

$$\|b\| \leq \|A\| \|x\|,$$

即有

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq \|A^{-1}\| \|A\| \frac{\|r\|}{\|b\|}.$$

特别地, 在上式中取 ∞ 范数便有

$$\frac{\|x - \hat{x}\|_{\infty}}{\|x\|_{\infty}} \leq \kappa_{\infty}(A) \frac{\|r\|_{\infty}}{\|b\|_{\infty}}. \quad (30)$$

这样可在实际计算时通过计算

$$\kappa_{\infty}(A) \frac{\|r\|_{\infty}}{\|b\|_{\infty}}$$

来给出计算解的精度估计, 而上式中除了 $\kappa_{\infty}(A)$ 外, 其余的量都是容易计算的. 注意到

$$\kappa_{\infty}(A) = \|A^{-1}\|_{\infty} \|A\|_{\infty},$$

而 $\|A\|_{\infty}$ 又是易于计算的, 便知用(30)式来估计计算解精度的关键在于如何估计 $\|A^{-1}\|_{\infty}$. 现在已有不少的实用方法可以给出这一估计. 这里介绍LAPACK所采用的一种优化方法, 该方法就是著名的“盲人爬山法”的一个具体应用.

设 $B \in \mathbf{R}^{n \times n}$, 下面估计 $\|B\|_1$. 定义

$$f(x) = \|Bx\|_1 = \sum_{i=1}^n \left| \sum_{j=1}^n b_{ij} x_j \right|, \\ \mathcal{D} = \{x \in \mathbf{R}^n : \|x\|_1 \leq 1\}$$

则易证 f 是凸函数, \mathcal{D} 是凸集, 而且求 $\|B\|_1$ 的问题就等价于求凸函数 f 在凸集 \mathcal{D} 上的最大值问题.

设 f 在 x 点的梯度 $\nabla f(x)$ 存在, 则由凸函数的性质可知

$$f(y) \geq f(x) + \nabla f(x)^T (y - x), \quad y \in \mathbf{R}^n.$$

现假定 $x_0 = (x_j^{(0)}) \in \mathbf{R}^n$ 满足 $\|x_0\|_1 = 1$, 使得

$$\sum_{j=1}^n b_{ij} x_j^{(0)} \neq 0, \quad i = 1, \dots, n.$$

令

$$\xi_i = \operatorname{sgn} \left(\sum_{j=1}^n b_{ij} x_j^{(0)} \right),$$

则在 x_0 附近有

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n \xi_i b_{ij} x_j.$$

因此, 有

$$\begin{aligned} \nabla f(x_0) &= \left(\frac{\partial f(x_0)}{\partial x_1}, \dots, \frac{\partial f(x_0)}{\partial x_n} \right)^T \\ &= \left(\sum_{i=1}^n \xi_i b_{i1}, \dots, \sum_{i=1}^n \xi_i b_{ij}, \dots, \sum_{i=1}^n \xi_i b_{in} \right)^T. \\ &= B^T v \end{aligned}$$

其中 $v = (\xi_1, \dots, \xi_n)^T$.

再令

$$w = Bx_0, \quad z = B^T v.$$

于是, 有下面结论成立:

定理 2.5.1: 假定 B, x_0, v, w 和 z 如上所述, 则有

- 1 若 $\|z\|_\infty = z^T x_0$ (**Holder不等式成立**),
则 $\|w\|_1 = \|Bx_0\|_1$ 是 $f(x) = \|Bx\|_1$ 在 \mathcal{D} 中的局部极大值;
- 2 若 $\|z\|_\infty > z^T x_0$, 则 $\|Be_j\|_1 > \|Bx_0\|_1$, 其中 j 满足

$$|z_j| = \|z\|_\infty.$$

证明: (1) 由于在 x_0 附近 $f(x)$ 是 x 的线性函数, 因此有

$$f(x) = f(x_0) + \nabla f(x_0)^T (x - x_0).$$

这样, 只需证在 x_0 附近有

$$\nabla f(x_0)^T (x - x_0) = z^T (x - x_0) \leq 0$$

即可. 事实上, 对于 $\|x\|_1 \leq 1$, 有

$$\begin{aligned} z^T (x - x_0) &= z^T x - z^T x_0 \\ &\leq \|z\|_\infty \|x\|_1 - z^T x_0 \\ &\leq \|z\|_\infty - z^T x_0 \leq 0. \end{aligned}$$

(2) 取 $\tilde{x} = e_j \operatorname{sgn}(z_j)$, 则有

$$\begin{aligned}\|Be_j\|_1 &= \|B\tilde{x}\|_1 = f(\tilde{x}) \\ &\geq f(x_0) + \nabla f(x_0)^T (\tilde{x} - x_0) \\ &= \|Bx_0\|_1 + z^T \tilde{x} - z^T x_0 \\ &= \|Bx_0\|_1 + |z_j| - z^T x_0 \\ &= \|Bx_0\|_1 + \|z\|_\infty - z^T x_0 \\ &> \|Bx_0\|_1\end{aligned}$$

即(2)成立.

基于这一定理可设计算法如下:

算法2.5.1: 估计矩阵的1 范数的优化法

设 $k = 1$

while $k = 1$

$w = Bx$; $v = \mathbf{sign}(w)$; $z = B^T v$

if $\|z\|_\infty = z^T x$

$v = \|w\|_1$

$k = 0$

else

$x = e_j$, 其中 $|z_j| = \|z\|_\infty$

$k = 1$

end

end

该算法初始的 x 可选任意满足 $\|x\|_1 = 1$ 的向量, 例如可取

$$x_i = 1/n, \quad i = 1, \dots, n.$$

假如已经计算好矩阵 A 的列主元三角分解: $PA = LU$, 则利用上述算法仅用 $O(n^2)$ 的运算量就可给出 $\|A^{-1}\|_\infty$ 的一个估计值.

由于 $\|A^{-1}\|_\infty = \|A^{-T}\|_1$, 因此只需应用算法2.5.1于矩阵 $B = A^{-T}$ 上即可, 此时计算 $w = Bx$ 和 $z = B^T v$ 就相当于解方程组 $A^T w = x$ 和 $Az = v$, 利用 A 的三角分解这两个方程组是很容易求解的.

综合上述讨论, 可以按如下的步骤来估计一个计算解 \hat{x} 的精度:

- 1 应用算法2.5.1于 $B = A^{-T}$ 上得到 $\|A^{-1}\|_{\infty}$ 的一个估计值 \tilde{v} ;
- 2 分别计算 $\|r\|_{\infty}$, $\|b\|_{\infty}$ 和 $\|A\|_{\infty}$ 得到它们的计算值 $\tilde{\gamma}$, $\tilde{\beta}$, 和 $\tilde{\mu}$;
- 3 计算 $\tilde{\rho} = \frac{\tilde{v}\tilde{\mu}\tilde{\gamma}}{\tilde{\beta}}$, 则数 $\tilde{\rho}$ 就可以作为计算解 \hat{x} 的相对误差的一个估计.

对于绝大多数问题这一方法常常可以给出计算解相对误差的相当好的估计, 但是也有一些特殊的问题利用该方法得到的 $\tilde{\rho}$ 远远小于计算解的实际相对误差. 这主要是由两个方面的原因引起的: 一是由于舍入误差的影响使得计算得到的 $\tilde{\gamma}$ 远远小于 $\|r\|_{\infty}$ 的真值; 二是当 A 十分病态时计算得到的三角分解已经相当不准确, 以至于应用它去估计出的 \tilde{v} 要比 $\|A^{-1}\|_{\infty}$ 的真值小得多. 虽然现在已经有了一些部分解决这些问题的方法, 但怎样才能使给出的估计更好仍然是一个值得进一步深入研究的问题.

若计算的解 \hat{x} 的精度太低, 可将 \hat{x} 作为初值, 应用Newton迭代法于函数 $f(x) = Ax - b$ 上, 来改进其精度. 具体计算过程可按如下步骤进行:

- (1) 计算 $r = b - A\hat{x}$ (用双精度和原始矩阵 A).
- (2) 求解 $Az = r$ (利用 A 的三角分解).
- (3) 计算 $x = \hat{x} + z$.
- (4) 若 $\frac{\|x - \hat{x}\|_{\infty}}{\|x\|_{\infty}} \leq \varepsilon$, 则结束; 否则, 令 $\hat{x} = x$, 转步(1).

实际计算的经验表明, 当 A 的病态并不是十分严重时, 利用这一方法最终可使其解的计算精度达到机器精度. 可是, 当 A 十分病态时, 这样做对解的精度并不会有太大的改进.