

§4 计算机算术

§4.2 取整误差分析

§4.2.1 单数取整

Def 4.24 舍入 $f_l: \mathbb{R} \rightarrow F \cup \{\infty, -\infty, NaN\}$, 默认舍入到最近整数; 距离相等情况下选择舍入到最近偶数.

Def 4.25 若 $|x| > OFL(F)$, $f_l(x) = NaN$, 称上溢; 若 $|x| \in (0, UFL(F))$, $f_l(x) = 0$, 称下溢.

Def 4.26 F 的单位舍入 $\epsilon_u = \frac{1}{2}\epsilon_M = \frac{1}{2}\beta^{1-p}$

Thm 4.27 对 $x \in R(F)$, 有 $f_l(x) = x(1+\delta)$, $|\delta| < \epsilon_u$

PF: $x_L, x_R \in F$, $x_L \leq x \leq x_R$ 为 x 相邻两个浮点数, $|f_l(x) - x| \leq \frac{1}{2}|x_R - x_L| \leq \epsilon_u \min(|x_L|, |x_R|) < \epsilon_u |x|$ (由 Lem 4.23)

Thm 4.28 对 $x \in R(F)$, 有 $f_l(x) = \frac{x}{1+\delta}$, $|\delta| \leq \epsilon_u$

PF: $|f_l(x) - x| \leq \frac{1}{2}|x_R - x_L| \leq \epsilon_u \min(|x_L|, |x_R|) \leq \epsilon_u |f_l(x)|$

Ex 4.29 $\frac{2}{3} = (0.1010\dots)_2 = (1.010101\dots)_2 \times 2^{-1}$

$$\begin{aligned} x_L &= (1.010\dots 10)_2 \times 2^{-1} & x - x_L &= \frac{2}{3} \times 2^{-24}, \\ x_R &= (1.010\dots 11)_2 \times 2^{-1} & x_R - x &= 2^{-24} \end{aligned} \Rightarrow f_l(x) = x_R$$

§4.2.2 二进制浮点数运算

Def 4.30 (加减) 对 $a, b \in F$, $a = M_a \times \beta^{e_a}$, $b = M_b \times \beta^{e_b}$, $M_a = \pm m_a$, $M_b = \pm m_b$, 不妨设 $|a| > |b|$, 在精度至少为 $2p$ 的寄存器中计算 $c = f_l(a+b) \in F$.

① 比较指数: 若 $e_a - e_b > p+1$, 相差太大, 返回 $C=a$; 否则 $e_c \leftarrow e_a$, $M_b \leftarrow M_b / \beta^{e_a - e_b}$ (对 M_b 右移)

② 加法: $M_c \leftarrow M_a + M_b$, 舍入后精度为 $2p$

③ 正规化: 若 $M_c = 0$, 返回 0; 移位使 $M_c \in [1, \beta)$: $|M_c| \in (\beta, \beta^2)$, $M_c \leftarrow \frac{M_c}{\beta}$; $|M_c| \in (0, \beta - \epsilon_u(p))$, 一直 $M_c \leftarrow M_c p$ 直到同时改 e_c . 若 $|M_c| \in [\beta - \epsilon_u(p), \beta]$, $|M_c| \leftarrow 1.0$. ($|M_c| \in [1, \beta)$)

④ 检查范围 (ec): $NaN/0$

⑤ 舍入到精度 p ⑥ $C \leftarrow M_c \times \beta^{e_c}$

Lem 4.34 对 $a, b \in F$, $a+b \in R(F)$ 有 $f_l(a+b) = (a+b)(1+\delta)$, $|\delta| < \epsilon_u$

PF: 例数第 2 步 ④ 的舍入误差占主要, 由 Thm 4.27 即证.

Def 4.35 (乘) 对 $a, b \in F$, $a = M_a \times \beta^{e_a}$, $b = M_b \times \beta^{e_b}$, $M_a = \pm m_a$, $M_b = \pm m_b$, 在精度至少为 $2p$ 的寄存器中计算 $f_l(ab)$

① $e_c \leftarrow e_a + e_b$ ② $M_c \leftarrow M_a M_b$ 在寄存器中舍入.

③ 正规化: $|M_c| \in (\beta, \beta^2)$, $M_c \leftarrow \frac{M_c}{\beta}$, $e_c \leftarrow e_c + 1$ ④ $|M_c| \in [\beta - \epsilon_u(p), \beta]$, $|M_c| \leftarrow 1.0$, $e_c \leftarrow e_c + 1$

⑤ 检查范围 $NaN/0$ ⑥ 舍入到精度 p ⑦ $C \leftarrow M_c \times \beta^{e_c}$

Lem 4.37 对 $a, b \in F$, $ab \in R(F)$ 有 $f_l(ab) = (ab)(1+\delta)$, $|\delta| < \epsilon_u$

PF: 同样只有 ⑤ 误差占主要地位.

Def 4.38 对 $a, b \in F$, $a = M_a \times \beta^{e_a}$, $b = M_b \times \beta^{e_b}$, $M_a = \pm m_a$, $M_b = \pm m_b$, 在精度至少为 $2p+1$ 寄存器中计算 $fl(\frac{a}{b}) \in F$.

- ① 若 $m_b = 0$, 返回 NaN; ② 否则 $e_c \leftarrow e_a - e_b$; ③ $M_c \leftarrow \frac{M_a}{M_b}$, 在寄存器中舍入
④ Normalization: 若 $|M_c| < 1$, $M_c \leftarrow M_c \cdot \beta$, $e_c \leftarrow e_c - 1$ ⑤ 检查范围 $|e_c|$: NaN/0
⑥ 舍入到精度 p ; ⑦ 返回 $c \leftarrow M_c \times \beta^{e_c}$

Lem 4.39 对 $a, b \in F$, $\frac{a}{b} \in R(F)$, 有 $fl(\frac{a}{b}) = \frac{a}{b}(1+\delta)$, $|\delta| < \epsilon_u$

PF: 当 $|M_a| = |M_b|$ 时, 没有误差, 考虑 $|M_a| > |M_b|$, $|M_a|, |M_b| \in [1, \beta)$, $|\frac{M_a}{M_b}| \geq \frac{\beta - \epsilon_m}{\beta - 2\epsilon_m} > 1 + \beta^{-1}\epsilon_m$

此时 ④ 不用进行, 设有 $p+k$ 精度, 则舍入单位 $\frac{1}{2}\beta^{p+k} = \frac{1}{2}\beta^{p+1}\beta^k \beta^{p-1-k} = \epsilon_u \epsilon_m \beta^{p-1-k}$

令 $k = p+1$, 则舍入单位为 $\epsilon_u \epsilon_m \beta^{-2}$, 记第 ③ 步结果为 M_{c1} , ④ 结果为 M_{c2} , 从而:

$$M_{c2} = M_{c1} + \delta_2 \quad (|\delta_2| < \epsilon_u) = \frac{M_a}{M_b} + \delta_1 + \delta_2 \quad (|\delta_1| < \epsilon_u \epsilon_m \beta^{-2}) = \frac{M_a}{M_b}(1+\delta)$$

$$|\delta| = \left| \frac{\delta_1 + \delta_2}{M_a/M_b} \right| < \frac{\epsilon_u(1 + \epsilon_m \beta^{-2})}{1 + \epsilon_m \beta^{-1}} < \epsilon_u$$

考虑 $|M_a| < |M_b|$, 类似有 $|\frac{M_a}{M_b}| \leq \frac{\beta - 2\epsilon_m}{\beta - \epsilon_m} = 1 - \frac{\epsilon_m}{\beta - \epsilon_m} < 1 - \beta^{-1}\epsilon_m$, 此时 ④ 一定会进行

$$M_{c1} = M_a/M_b + \delta_1, \quad |\delta_1| < \epsilon_u \epsilon_m \beta^{-2}, \quad M_{c2} = \beta M_{c1} + \delta_2 = \beta \frac{M_a}{M_b} (1 + \frac{\beta \delta_1 + \delta_2}{\beta M_a/M_b}), \quad |\delta_2| < \epsilon_u$$

$$\beta |\frac{M_a}{M_b}| \geq \frac{\beta}{\beta - \epsilon_m} = 1 + \frac{\epsilon_m}{\beta - \epsilon_m} > 1 + \beta^{-1}\epsilon_m, \quad \text{代入得到 } |\delta| = \left| \frac{\beta \delta_1 + \delta_2}{\beta M_a/M_b} \right| < \frac{\epsilon_u \epsilon_m \beta^{-1} + \epsilon_u}{1 + \beta^{-1}\epsilon_m} = \epsilon_u$$

Thm 4.40 记 F 为精度 p 的正规 FPN 系统, 则 $\forall \odot = +, -, \times, /$, $\forall a, b \in F$ $a \odot b \in R(F)$, $fl(a \odot b) = (a \odot b)(1+\delta)$, $|\delta| < \epsilon_u$
当且仅当在 $2p+1$ 精度的寄存器中运算.

§4.2.3 舍入误差的传递

Thm 4.41 若 $i=0, 1, \dots, n$, $a_i \in F$, $a_i > 0$, 则 $fl(\sum_{i=0}^n a_i) = (1+\delta_n) \sum_{i=0}^n a_i$, 其中 $|\delta_n| < (1+\epsilon_u)^n - 1 \approx n\epsilon_u$

PF: 令 $s_k = \sum_{i=0}^k a_i$, $s_0^* = a_0$, $s_{k+1}^* = fl(s_k^* + a_{k+1})$, $\delta_{k+1} = \frac{s_{k+1}^* - s_{k+1}}{s_{k+1}}$, $\epsilon_k = \frac{s_{k+1}^* - (s_k^* + a_{k+1})}{s_k^* + a_{k+1}}$

由归纳法, 则有 $\delta_{k+1} = \frac{s_{k+1}^* - s_{k+1}}{s_{k+1}} = \frac{(s_k^* + a_{k+1})(1+\epsilon_k) - s_{k+1}}{s_{k+1}} = \frac{(s_k(1+\delta_k) + a_{k+1})(1+\epsilon_k) - s_k - a_{k+1}}{s_{k+1}} = \frac{(\epsilon_k + \delta_k + \epsilon_k \delta_k)s_k + \epsilon_k a_{k+1}}{s_{k+1}}$
 $= \frac{\epsilon_k s_{k+1} + \delta_k(1+\epsilon_k)s_k}{s_{k+1}} = \epsilon_k + \delta_k(1+\epsilon_k) \frac{s_k}{s_{k+1}}, \quad \forall s_k < s_{k+1}, |\epsilon_k| < \epsilon_u \therefore |\delta_{k+1}| < |\epsilon_k| + |\delta_k|(1+\epsilon_u) < \epsilon_u + |\delta_k|(1+\epsilon_u)$
 $\therefore \forall k \in \mathbb{N}, |\delta_{k+1}| < \epsilon_u \sum_{i=0}^k (1+\epsilon_u)^i = \epsilon_u \frac{(1+\epsilon_u)^{k+1} - 1}{1+\epsilon_u - 1} = (1+\epsilon_u)^{k+1} - 1$

Ex 4.42 对 $a_i > 0$, 按从小到大的顺序进行加法, 可以最小化舍入误差.

Ex 4.43 类似 $fl(\sum_{i=1}^n a_i)$, $fl(a_1 b_1 + a_2 b_2 + a_3 b_3) = (1+\delta) Ans$ 且 $|\delta| < (1+\epsilon_u)^b - 1$

Thm 4.44 对给定 $M \in \mathbb{R}^+$, 及正整数 $n \leq \lfloor \frac{\ln 2}{\ln M} \rfloor$, 设 $|\delta_i| \leq M$, $i=1, 2, \dots, n$. 则 $1 - nM \leq \prod_{i=1}^n (1+\delta_i) \leq 1 + nM + (nM)^2$

或者等价地对 $I_n = [-\frac{1}{1+nM}, 1]$ 有, $\exists \theta \in I_n$, s.t. $\prod_{i=1}^n (1+\delta_i) = 1 + \theta(nM + nM^2)$

PF: 由 $|\delta_i| \leq M$, $i=1, 2, \dots, n$ 可得 $(1-M)^n \leq \prod_{i=1}^n (1+\delta_i) \leq (1+M)^n \leq e^{nM} < 1+M \leq e^M$

将 $f(\mu) = (1-\mu)^n$ 在 $\mu=0$ 展开 $(1-\mu)^n = 1 - n\mu + \frac{1}{2}n(n-1)\mu^2 \geq 1 - n\mu$

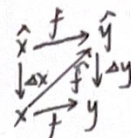
而 $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots \leq 1 + x + \frac{x^2}{2} e^x$ 令 $x = nM \leq \ln 2$ 则有 $e^{nM} \leq 1 + nM + (nM)^2$

由于 $\prod_{i=1}^n (1+\delta_i)$ 在连续函数 $f(T) = 1 + T(nM + nM^2)$, $T \in I_n$ 的值域中, 由中值定理即证.

§4.3 精确性和稳定性.

§4.3.1 Avoiding catastrophic cancellation

Def 4.45 设 \hat{x} 为 $x \in \mathbb{R}$ 的近似, 绝对误差为 $E_{abs}(\hat{x}) = |x - \hat{x}|$, 相对误差为 $E_{rel}(\hat{x}) = \frac{|\hat{x} - x|}{|x|}$



Def 4.46 通过 $\hat{y} = \hat{f}(x)$ 近似 $y = f(x)$, 向前误差为 \hat{y} 近似 y 的相对误差; 向后误差是用 $f(\hat{x}) = \hat{f}(x)$ 中 \hat{x} 近似 x 的最小相对误差.

Def 4.47 计算 $y = f(x)$ 的算法 $\hat{y} = \hat{f}(x)$ 是精确的如果存在 $c > 0$ 且很小, $\forall x \in \text{dom}(f)$, $E_{rel}(\hat{f}(x)) = \left| \frac{\hat{f}(x) - f(x)}{f(x)} \right| \leq c \epsilon_u$

Ex 4.48 (Catastrophic cancellation), $x, y \in \mathbb{R}(F)$, 由 Thm 4.27/40. $f(f(x) \odot f(y)) = (x(1+\delta_1) \odot y(1+\delta_2))(1+\delta_3)$, $|\delta_i| \leq \epsilon_u$

由 Thm 4.40/44 $f(f(x) \times f(y)) = xy(1+\delta_1)(1+\delta_2)(1+\delta_3) = xy[1 + \theta(3\epsilon_u + 3\epsilon_u^2 + \epsilon_u^3)]$, $\theta \in [-1, 1]$

类似的, $f\left(\frac{f(x)}{f(y)}\right) = \frac{x(1+\delta_1)}{y(1+\delta_2)}(1+\delta_3) = \frac{x}{y}(1+\delta_1)(1-\delta_2+\delta_2^2-\dots)(1+\delta_3) \approx \frac{x}{y}(1+\delta_1)(1-\delta_2)(1+\delta_3)$

然而加/减可能不精确: $f(f(x) + f(y)) = (x(1+\delta_1) + y(1+\delta_2))(1+\delta_3) = (x+y)(1+\delta_3 + \frac{x\delta_1+y\delta_2}{x+y} + \delta_3 \frac{x\delta_1+y\delta_2}{x+y})$

当 $x+y \rightarrow 0$ 时, 相对误差可以任意大

Thm 4.49 设 $x, y \in F$, $x > y > 0$ 且 $\beta^l \leq 1 - \frac{y}{x} \leq \beta^s$, 在计算 $x-y$ 时最大有效数字的位数最多丢失 l , 至少丢失 s .

PF: 令 $x = m_x \times \beta^n$, $y = m_y \times \beta^m$, $m_x, m_y \in [1, \beta)$, 首先对 y 移位 $y = (m_y \times \beta^{m-n}) \times \beta^n \Rightarrow x-y = (m_x - m_y \times \beta^{m-n}) \times \beta^n$

$\Rightarrow m_x - y > m_x - m_y \times \beta^{m-n} = m_x \left(1 - \frac{m_y \times \beta^n}{m_x \times \beta^n}\right) = m_x \left(1 - \frac{y}{x}\right) \Rightarrow \beta^l \leq m_x - y < \beta^{l-s}$, 正规化时在左移至少 l 位, 至少 s 位.

△当 x, y 很接近时, 使用减法会出现上述问题. 应尽量避免.

§4.3.2 Backward stability and numerical stability

Def 4.52, 计算 $y = f(x)$ 的算法 $\hat{y} = \hat{f}(x)$ 是向后稳定的, 若 $\exists c > 0$ s.t. $\forall x \in \text{dom}(f)$, $\exists \hat{x} \in \text{dom}(f)$, $\hat{f}(x) = f(\hat{x}) \Rightarrow E_{rel}(\hat{x}) \leq c \epsilon_u$

Def 4.53 $\hat{y} = \hat{f}(x_1, x_2)$ 是向后稳定的, 若 $\exists c_1, c_2 > 0$ s.t. $\forall (x_1, x_2) \in \text{dom}(f)$, $\exists (\hat{x}_1, \hat{x}_2) \in \text{dom}(f)$ s.t.

$\hat{f}(x_1, x_2) = f(\hat{x}_1, \hat{x}_2) \Rightarrow E_{rel}(\hat{x}_1) \leq c_1 \epsilon_u$, $E_{rel}(\hat{x}_2) \leq c_2 \epsilon_u$.

lem 4.54 对 $f(x_1, x_2) = x_1 - x_2$, $x_1, x_2 \in \mathbb{R}(F)$, 算法 $\hat{f}(x_1, x_2) = f(f(x_1) - f(x_2))$ 是向后稳定的.

PF: $\hat{f}(x_1, x_2) = (x_1(1+\delta_1) - x_2(1+\delta_2))(1+\delta_3) = x_1(1+\delta_1+\delta_3+\delta_1\delta_3) - x_2(1+\delta_2+\delta_3+\delta_2\delta_3) := \hat{x}_1 - \hat{x}_2$

Ex 4.55 $\hat{f}(x) = f(1.0 + f(x))$ 不是向后稳定, 当 $x \in [0, \epsilon_u)$ 时, $E_{rel}(\hat{x}) = 1$ ($\hat{x} = 0$)

Def 4.56 计算 $y = f(x)$ 的算法 $\hat{y} = \hat{f}(x)$ 是稳定的(或数值稳定) $\Leftrightarrow \exists c, c_f > 0$, s.t. $\forall x \in \text{dom}(f)$, $\exists \hat{x} \in \text{dom}(f)$

lem 4.57 如果一个算法向后稳定, 则其数值稳定.

s.t. $\left| \frac{\hat{f}(x) - f(x)}{f(x)} \right| \leq c_f \epsilon_u$, $E_{rel}(\hat{x}) \leq c \epsilon_u$

Ex 4.58 $\hat{f}(x) = f(1.0 + f(x))$ 是稳定的.

若 $|x| < \epsilon_u$, 取 $\hat{x} = x$, $\left| \frac{\hat{f}(x) - f(x)}{f(x)} \right| = \left| \frac{x}{1+x} \right| < 2\epsilon_u$; 若 $|x| \geq \epsilon_u$, $\hat{f}(x) = 1 + \delta_3 + x(1+\delta_1+\delta_2+\delta_1\delta_2)$, $|\delta_1|, |\delta_2| < \epsilon_u$

令 $\hat{x} = x(1+\delta_1+\delta_2+\delta_1\delta_2)$, $E_{rel}(\hat{x}) = |\delta_1+\delta_2+\delta_1\delta_2| < 3\epsilon_u \Rightarrow \left| \frac{\hat{f}(x) - f(x)}{f(x)} \right| = \left| \frac{\delta_3}{1+\delta_2+x(1+\delta_1+\delta_2+\delta_1\delta_2)} \right| < |\epsilon_u|$

§4.3.3 Condition numbers: scalar functions

Def 4.59 函数 $y = f(x)$ 的相对条件数是输入的微小变化对输出的度量 $C_f(x) = \left| \frac{x f'(x)}{f(x)} \right|$

Def 4.60 条件数小的问题称为良态, 否则称为病态.

Ex 4.61, 由 Def 4.59, $\frac{Erel(\hat{y})}{Erel(\hat{x})} = \left| \frac{x(f(x)-f(\hat{x}))}{f(x)(x-\hat{x})} \right| \approx C_f = \left| \frac{x f'(x)}{f(x)} \right| \Rightarrow Erel(\hat{y}) \approx C_f Erel(\hat{x})$

可见病态问题前向误差大.

lem 4.63 考虑在单根 r 附近求解 $f(x)=0$ 即 $f(r)=0, f'(r) \neq 0$, 设有微小扰动 $F=f+eg, f, g \in C^2, g(r) \neq 0$ 且满足

$|eg(r)| \ll |f'(r)|$, 则 F 有根 $r+h$, 其中 $h \approx -\epsilon \frac{g(r)}{f'(r)}$.

PF: $0 = F(r+h) = f(r) + hf'(r) + \epsilon[g(r) + hg'(r)] + o(h^2)$ 从而有 $h \approx -\epsilon \frac{g(r)}{f'(r) + \epsilon g'(r)} \approx -\epsilon \frac{g(r)}{f'(r)}$

Ex 4.64 (Wilkinson) 定义 $f(x) = \prod_{k=1}^p (x-k)$, $g(x) = x^p$. 根 $x=p$ 受扰动后 $h \approx -\epsilon \frac{p^p}{(p-1)!}$, 当 p 很大时, h 几乎不可能求解高阶的根

§4.3.4 Condition numbers: vector functions.

Def 4.65 向量函数 $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ 的条件数为 $Cond_f(x) = \frac{\|x\| \|\nabla f\|}{\|f(x)\|}$, 其中 $\|\cdot\|$ 表示欧几里德范数 (如 $1, 2, \infty$)

Def 4.66 非奇异方阵 A 的条件数为 $cond A = \|A\| \|A^{-1}\|$. (由 $Au=b$ 求解 $u=A^{-1}b$ 得到)

lem 4.68 在 2-范数下, 非奇异方阵 A 的条件数为 $cond_2 A = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_{\max}}{\sigma_{\min}}$, 其中 $\sigma_{\max}, \sigma_{\min}$ 为最大/小奇异值.

若 A 是 normal 的, $cond_2 A = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}$, 为最大/小特征值.

Thm 4.70 A 为可逆阵, \hat{A} 为 A 的一个扰动 $\|A^{-1}\| \|\hat{A} - A\| < 1$, 其中 $Ax=b, \hat{A}\hat{x}=\hat{b}$, 则 \hat{x} 关于 x 的相对误差满足

$$\frac{\|\hat{x} - x\|}{\|x\|} \leq \frac{cond A}{1 - (cond A) \frac{\|\hat{A} - A\|}{\|A\|}} \left(\frac{\|\hat{b} - b\|}{\|b\|} + \frac{\|\hat{A} - A\|}{\|A\|} \right)$$

PF: $\hat{A} = A[I + A^{-1}(\hat{A} - A)], \hat{A}^{-1} = [I + A^{-1}(\hat{A} - A)]^{-1} A^{-1} \therefore \|\hat{A}^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\hat{A} - A\|} \Rightarrow \hat{x} - x = \hat{A}^{-1}[\hat{b} - b - (\hat{A} - A)x]$

$$\Rightarrow \frac{\|\hat{x} - x\|}{\|x\|} \leq \frac{cond A}{1 - \|\hat{A} - A\| \|A^{-1}\|} \left(\frac{\|\hat{b} - b\|}{\|A\| \|x\|} + \frac{\|\hat{A} - A\|}{\|A\|} \right), \text{ 由 } \|A\| \|x\| \geq \|b\| \Rightarrow \frac{\|\hat{x} - x\|}{\|x\|} \leq \frac{cond A}{1 - (cond A) \frac{\|\hat{A} - A\|}{\|A\|}} \left(\frac{\|\hat{b} - b\|}{\|b\|} + \frac{\|\hat{A} - A\|}{\|A\|} \right)$$

(由 Thm E.145 $\|I - A^{-1}(\hat{A} - A)\| \leq \frac{1}{1 - \|A^{-1}\| \|\hat{A} - A\|}$ (其中 $\|A^{-1}\| \|\hat{A} - A\| < 1$))

Def 4.71 向量函数 $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ 的条件数为 $Cond_f(x) = \frac{\|x\| \|\nabla f\|}{\|f(x)\|}$, 分量条件数为 $cond_f(x) = \|A(x)\|$

其中 $A(x) = [A_{ij}(x)], A_{ij}(x) = \left| \frac{x_j \frac{\partial f_i}{\partial x_j}}{f_i(x)} \right|$, 也就是 A_{ij} 是 f_i 关于 x_j 的条件数.

Ex 4.72 对于向量函数 $f(x) = [x_1 + x_2, x_1 - x_2]$

其 Jacobi 矩阵为 $\nabla f = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$, $C_1 = \begin{bmatrix} \frac{x_2}{x_1+x_2} & \frac{x_1}{x_1+x_2} \\ \frac{x_2}{x_1-x_2} & \frac{x_1}{x_1-x_2} \end{bmatrix}$ 为 $Cond_f(x)$ 中的 $A(x)$. 可见当 $x_1 \pm x_2 \approx 0$ 时会导致

病态. $C_1 = \frac{\|x\| \|\nabla f\|}{\|f\|} = \frac{|x_1| + |x_2|}{|x_1 x_2|} \frac{2 \max(x_1^2, x_2^2)}{|x_1 + x_2| + |x_1 - x_2|}$, 然而当 $x_1 \pm x_2 \approx 0$ 时 $C_1 \approx 2$, 看不出病态.

注: $\forall A \in \mathbb{R}^{n \times n}, \|A\|_1 = \max_j \sum_i |A_{ij}|$

Ex 4.73 t_1, \dots, t_n 在 $[-1, 1]$ 中等距分布, Vandermonde 阵的条件数 (基于 ∞ 范数) 为 $cond_\infty V_n \sim \frac{1}{\pi} e^{-\pi/4} e^{n/4(\pi+1)n/2}$

Def 4.74 Hilbert 矩阵 $H_n \in \mathbb{R}^{n \times n}$ 为 $h_{ij} = \frac{1}{i+j-1}$

Ex 4.79 希尔伯特矩阵基于 2-范数的条件数为 $cond_2 H_n \sim \frac{(1/2)^{n+4}}{2^{1/4} \sqrt{\pi n}}$

§4.3.5 Condition numbers: algorithms

Def 4.76 向量函数 $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ 通过 $f_A: F^m \rightarrow F^n$ 逼近, 设 $x \in F^m, \exists x_A \in \mathbb{R}^m$, s.t. $f_A(x) = f(x_A)$

则 f_A 的条件数定义为 $\text{cond}_A(x) = \frac{1}{\epsilon_u} \inf_{\{x_A\}} \frac{\|x_A - x\|}{\|x\|}$

Ex 4.77 考虑算法 A 计算 $y = \ln x$, 设 $x > 0$, 计算得到 $y_A = (1+\delta)\ln x, |\delta| \leq \epsilon_u$. 求 A 的条件数.

$$y_A = \ln x_A, x_A = x^{1+\delta}, E_{\text{rel}}(x_A) = \left| \frac{x_A - x}{x} \right| = |x^\delta - 1| = |e^{\delta \ln x} - 1| \approx |\delta \ln x| \leq \epsilon_u |\ln x|$$

因此除了 $x \rightarrow 0^+$, A 是良态的.

Thm 4.78 设光滑函数 $f: \mathbb{R} \rightarrow \mathbb{R}$ 由算法 $A: F \rightarrow F, f_A(x) = f(x)(1+\delta(x)), |\delta(x)| \leq \varphi(x) \epsilon_u$ 逼近, 若 $\text{cond}_f(x)$

有界非 0, 则有 $\forall x \in F, \text{cond}_A(x) \leq \frac{\varphi(x)}{\text{cond}_f(x)}$.

PF: 设有 x_A 满足 $f(x_A) = f_A(x)$. 将 x_A 写作 $x(1+\epsilon_A)$, 则 $f(x)(1+\delta) = f(x_A) = f(x(1+\epsilon_A)) = f(x) + x\epsilon_A f'(x) + O(\epsilon_A^2)$

$$\Rightarrow \epsilon_A = \frac{f(x)}{x f'(x)} \delta. \therefore \left| \frac{x_A - x}{x} \right| = |\epsilon_A| = \left| \frac{f(x)}{x f'(x)} \right| |\delta(x)| \Rightarrow \frac{1}{\epsilon_u} \left| \frac{x_A - x}{x} \right| = \left| \frac{\delta(x)}{\epsilon_u \text{cond}_f(x)} \right| \leq \frac{\varphi(x)}{\text{cond}_f(x)}$$

Ex 4.79 假设 $\sin x, \cos x$ 在计算机中以舍入误差计算. (用截断 Taylor 展开)

$$f_A = f \left[\frac{f(1-f(\cos x))}{f(1(\sin x))} \right] \text{ 用计算 } f(x) = \frac{1-\cos x}{\sin x}, x \in (0, \frac{\pi}{2})$$

$$\text{sol: } \text{cond}_f(x) = \frac{x}{\sin x}, f_A(x) = \frac{(1-\cos x(1+\delta_1))(1+\delta_2)}{\sin x(1+\delta_3)} (1+\delta_4), |\delta_i| \leq \epsilon_u, \text{ 忽略 } O(\epsilon_i^2), \text{ 则}$$

$$f_A(x) = \frac{1-\cos x}{\sin x} \left\{ 1+\delta_2+\delta_4-\delta_3-\delta_1 \frac{\cos x}{1-\cos x} \right\}, \text{ 则有 } \varphi(x) = 3 + \frac{\cos x}{1-\cos x} \Rightarrow \text{cond}_A(x) = \frac{\sin x}{x} \left(3 + \frac{\cos x}{1-\cos x} \right)$$

因此, $x \rightarrow 0$ 时 $\text{cond}_A(x)$ 无界, $x \rightarrow \frac{\pi}{2}$, $\text{cond}_A(x) \rightarrow \frac{\pi}{6}$

§4.3.6 Overall error of a computer solution

$$x \xrightarrow{\text{rounding}} x^* \xrightarrow{\text{Algorithm } f_A} f_A(x^*)$$

Thm 4.81 $f: \mathbb{R}^m \rightarrow \mathbb{R}^n, y = f(x)$. 计算机输入输出表示为 $x^* \approx x, y_A^* = f_A(x^*)$

$$\text{则 } E_{\text{rel}}(y_A^*) \leq E_{\text{rel}}(x_A^*) \text{cond}_f(x) + \epsilon_u \text{cond}_f(x^*) \text{cond}_A(x^*)$$

$$\text{PF: } \frac{\|y_A^* - y\|}{\|y\|} = \frac{\|f_A(x^*) - f(x)\|}{\|f(x)\|} \leq \frac{\|f(x^*) - f(x)\|}{\|f(x)\|} + \frac{\|f_A(x^*) - f(x^*)\|}{\|f(x)\|}$$

$$\text{由 } E_{\text{rel}}(\hat{y}) \leq C_f E_{\text{rel}}(\hat{x}), \frac{\|f(x^*) - f(x)\|}{\|f(x)\|} \leq \text{cond}_f(x) \frac{\|x^* - x\|}{\|x\|} = E_{\text{rel}}(x^*) \text{cond}_f(x)$$

$$\frac{\|f_A(x^*) - f(x^*)\|}{\|f(x)\|} = \frac{\|f(x_A^*) - f(x^*)\|}{\|f(x)\|} \approx \frac{\|f(x_A^*) - f(x^*)\|}{\|f(x^*)\|} \leq \text{cond}_f(x^*) \frac{\|x_A^* - x^*\|}{\|x^*\|} = \epsilon_u \text{cond}_A(x^*) \text{cond}_f(x^*)$$

$$E_{\text{rel}}(y_A^*) \leq E_{\text{rel}}(x_A^*) \text{cond}_f(x) + \epsilon_u \text{cond}_f(x^*) \text{cond}_A(x^*)$$