

# 1 数值误差的避免

## 1.1 求平均的误差

$N$  数平均的误差来源于求和、除以  $N$  两个过程；在  $N$  较大时，除以大数所引入的误差相对较小，此时求和的误差占主要成分。

两数相加时，引入的相对误差为机器精度  $\frac{\epsilon_M}{2}$ ；记  $x_0 = \max |x_i|$ ，考虑最坏的情况，即可能的误差最大值，这一情形在每个  $x_i \rightarrow x_0$  时取到。不妨设  $x_i$  均为正数，此时求和的上限为：

$$f \circ f \circ \cdots \circ f(x_0) \equiv f^{N-1} \circ (x_0), \quad f(x) = (x + x_0) \left(1 + \frac{\epsilon_M}{2}\right)$$

(1.1)

这里  $f$  是每次数值求和操作的函数表示。

作用于 $x_0$	$\mathcal{O}(1)$ 项	$\mathcal{O}(\frac{\epsilon_M}{2})$ 系数
$f^0 = \mathbb{1}$	$x_0$	0
$f^1$	$2x_0$	$2x_0$
$f^2$	$3x_0$	$5x_0$
$\vdots$	$\vdots$	$\vdots$
$f^k$	$x_0 + kx_0$	$c_k$

$f$  的作用规律：先加  $x_0$ ，再乘以  $(1 + \frac{\epsilon_M}{2})$

考察  $\frac{\epsilon_M}{2}$  的系数，设  $f^k$  作用后的  $\frac{\epsilon_M}{2}$  系数为  $c_k$ ，则不难发现：

$$c_k = c_{k-1} + x_0 + kx_0$$

(1.2)

其中  $kx_0$  源于前一步  $\mathcal{O}(1)$  项的系数。已知  $c_0 = 0$ ，展开此递推关系，即得：

$$c_{N-1} = \frac{(N+2)(N-1)}{2} x_0,$$

(1.3)

均值的误差限：

$$\frac{1}{N} \cdot c_{N-1} \frac{\epsilon_M}{2} = \frac{(N+2)(N-1)}{2N} \frac{\epsilon_M}{2} \max |x_i| \sim \frac{N}{2} \frac{\epsilon_M}{2} \max |x_i|$$

## 1.2 方差计算的稳定性

两种方差计算公式如下：

$$S^2 = \frac{1}{N-1} \left\{ \sum_i x_i^2 - N\bar{x}^2 \right\}$$

(1.4a)

$$= \frac{1}{N-1} \sum_i (x_i - \bar{x})^2$$

(1.4b)

沿用前文给出的估计办法，可以给出两式的误差限；有：

$$\begin{aligned} e_{(a)} &\sim S^2 \frac{\epsilon_M}{2} + \left\{ \frac{N}{2} \frac{\epsilon_M}{2} \max |x_i|^2 + \frac{N}{2} \frac{\epsilon_M}{2} \max |x_i| \times 2 \right\} \\ &= \left( S^2 + \frac{N}{2} \max |x_i|^2 + N \max |x_i| \right) \frac{\epsilon_M}{2}, \end{aligned}$$

(1.5)

$$e_{(b)} \sim \frac{N}{2} \frac{\epsilon_M}{2} \max |x_i - \bar{x}|^2$$

可见，多数情况下第一式 (1.4a) 将带来较大误差；特别是在  $x_i$  很大但方差却很小的情况下，此时将产生大数相消，从而大量损失有效数字。相比之下，第二式 (1.4b) 较为稳定和准确。

## 1.3 递归计算的稳定性

考察  $I_n = \int_0^1 \mathrm{d}x \frac{x^n}{x+5}$ ，首先有  $I_0 = \ln(x+5)|_0^1 = \ln \frac{6}{5}$ ，而：

$$I_k + 5I_{k-1} = \int_0^1 \mathrm{d}x \frac{x^k + 5x^{k-1}}{x+5} = \int_0^1 \mathrm{d}x x^{k-1} = \frac{x^k}{k} \Big|_0^1 = \frac{1}{k}, \quad k = 1, 2, \dots$$

(1.6)

从而可以递归地给出  $I_k$  的值。

关注这一过程的误差传递，设计算值  $\hat{I}_{k-1} = I_{k-1} + \epsilon_{k-1}$ ，则相应地：

$$\begin{aligned} \hat{I}_k &\sim \left( \frac{1}{k} - 5\hat{I}_{k-1} \right) \left( 1 + \frac{\epsilon_M}{2} \right) \\ &\sim I_k - 5\epsilon_{k-1} + I_k \frac{\epsilon_M}{2} \end{aligned}$$

(1.7)

即有：

$$\epsilon_k \sim -\left( 5/I_k \right) \epsilon_{k-1} + \frac{\epsilon_M}{2}$$

系数  $\kappa = |5/I_k|$  是关键；若  $\kappa < 1$ ，则误差将得到控制，不会进一步放大。

然而，不幸的是，本问题中的  $I_n < I_0 < 1$ ，即始终有  $\kappa > 1$ ，初始误差  $\epsilon$  将随递归过程不断（指数）放大，可见这一算法是不稳定的。