

数值计算方法

数值计算中的误差分析

张晓平

2019 年 8 月 30 日

武汉大学数学与统计学院

Table of contents

1. 数值计算简介
2. 误差
3. 数值计算的误差估计
4. 选用和设计算法应遵循的原则

数值计算简介

数值计算简介

现代科学的三种研究方式

现代科学的三种研究方式

当今世界科学活动的三种主要方式：

- 理论
- 实验
- 科学计算

数值计算简介

解决科学与工程问题的几个步骤

解决科学工程问题的几个步骤

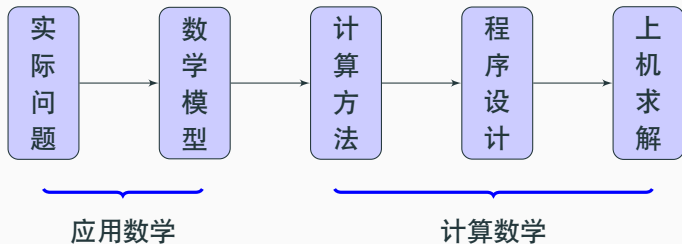


图 1: 解决科学工程问题的步骤

数值计算简介

研究数值方法的必要性

研究数值方法的必要性

例

解线性方程组

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n$$

定理：Crammer 法则

A 非奇异 \implies 此方程组有唯一解，且 $x_i = \frac{|A_i|}{|A|}, \quad i = 1, 2, \dots, n.$

该结论从理论上讲非常漂亮，它把线性方程组的求解问题归结为计算 $n+1$ 个 n 阶行列式的计算问题。

而对于行列式，可以采用 Laplace 展开定理进行计算：

定理：Laplace 展开定理

$$|A| = a_{i1}|A_{i1}| + a_{i2}|A_{i2}| + \cdots + a_{in}|A_{in}|, \quad A_{ij} \text{ 为 } a_{ij} \text{ 的代数余子式}$$

研究数值方法的必要性

实际操作中，该方法的运算量大的惊人，以至于完全不能用于实际计算。事实上，设 k 阶行列式所需乘法运算的次数为 m_k ，则

$$m_k = k + km_{k-1},$$

于是有

$$\begin{aligned} m_n &= n + nm_{n-1} \\ &= n + n[(n-1) + (n-1)m_{n-2}] \\ &= \dots \\ &= n + n(n-1) + n(n-1)(n-2) + \dots + n(n-1)\dots 3 \cdot 2 \\ &> n! \end{aligned}$$

故用 Crammer 法则和 Laplace 展开定理求解一个 n 阶线性方程组，所需乘法运算的次数就大于

$$(n+1)n! = (n+1)!.$$

在一台百亿次的计算机上求解一个 25 阶线性方程组，则至少需要

$$\frac{26!}{10^{10} \times 3600 \times 24 \times 365} \approx \frac{4.0329 \times 10^{28}}{3.1526 \times 10^{17}} \approx 13 \text{ 亿年}$$

而用下章介绍的消去法求解，则需要不到一秒钟。

数值计算简介

数值计算方法的研究对象

数值计算方法的研究对象

数值分析的研究对象为：

- 线性代数
- 曲线拟合
- 数值逼近
- 微积分
- 微分方程
- 积分方程
- ...

数值计算简介

数值计算方法的研究任务

数值计算方法的研究任务

数值分析的研究任务为：

- 算法设计
- 理论分析
 - 算法的收敛性
 - 稳定性
 - 误差分析
- 复杂度分析
 - 时间复杂度
 - 空间复杂度

▪

数值计算简介

数值计算方法的特点

数值计算方法的特点

数值分析的特点为：

- 既有数学的抽象性与严格性，又有广泛的应用性
- 有自身的研究方法和理论系统
- 与计算机紧密结合，实用性很强

误差

误差

误差来源与分类

误差来源与分类

误差按照其来源可分为四类，即

模型误差

数学模型只是复杂客观现象的一种近似，它与实际问题总会存在一定误差。

观测误差

由于测量精度和手段的限制，观测或实验得来的物理量总会与实际量之间存在误差。

截断误差

数学模型的精确解与由数值方法求出的近似解之间的误差。如

$$e^x \approx 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^{10}}{10!}$$

$$R_{10}(x) = \frac{\xi^{11}}{11!}$$

误差来源与分类

舍入误差

由于计算机的字长有限，进行数值计算的过程中，对计算得到的中间结果数据要使用“四舍五入”或其他规则取近似值，因而使计算过程有误差。

典故

1990 年 2 月 25 日，海湾战争期间，在沙特阿拉伯宰赫兰的爱国者导弹防御系统因浮点数舍入错误而失效，该系统的计算机精度仅有 24 位，存在 0.0001% 的计时误差，所以有效时间阈值是 20 个小时。当系统运行 100 个小时以后，已经积累了 0.3422 秒的误差。这个错误导致导弹系统不断地自我循环，而不能正确地瞄准目标。结果未能拦截一枚伊拉克飞毛腿导弹，飞毛腿导弹在军营中爆炸，造成 28 名美国陆军死亡。

误差

误差与有效数字

定义：绝对误差与绝对误差限

设某个量的精确值为 x ，其近似值为 x^* ，则称

$$E(x) = x - x^*$$

为近似值 x^* 的**绝对误差**，简称**误差**。若存在 $\eta > 0$ ，使得

$$|E(x)| = |x - x^*| \leq \eta$$

则称 η 为近似值 x^* 的**绝对误差限**，简称**误差限**或**精度**。

定义：绝对误差与绝对误差限

设某个量的精确值为 x ，其近似值为 x^* ，则称

$$E(x) = x - x^*$$

为近似值 x^* 的**绝对误差**，简称**误差**。若存在 $\eta > 0$ ，使得

$$|E(x)| = |x - x^*| \leq \eta$$

则称 η 为近似值 x^* 的**绝对误差限**，简称**误差限**或**精度**。

η 越小，表示近似值 x^* 的精度越高。

$$x^* - \eta \leq x \leq x^* + \eta \quad \text{or} \quad x = x^* \pm \eta$$

例

用毫米刻度的直尺量一长度为 x 的物体，测得其近似值为 $x^* = 84\text{mm}$ 。

因直尺以 mm 为刻度，其误差不超过 0.5mm，即有

$$|x - 84| \leq 0.5 \text{ mm} \quad \text{or} \quad x = 84 \pm 0.5 \text{ mm}.$$

例

测量 100m 和 10m 的两个长度，若它们的绝对误差均为 1cm，显然前者的测量更为精确。

例

测量 100m 和 10m 的两个长度，若它们的绝对误差均为 1cm，显然前者的测量更为精确。

由此可见，决定一个量的近似值的精确度，除了绝对误差外，还必须考虑该量本身的大小，为此引入相对误差的概念。

定义：相对误差与相对误差限

近似值 x^* 的相对误差是绝对误差与精确值之比，即

$$E_r(x) = \frac{E(x)}{x} = \frac{x - x^*}{x}.$$

实际中精确值 x 一般无法知道，通常取

$$E_r(x) = \frac{E(x)}{x^*} = \frac{x - x^*}{x^*}.$$

若存在 $\delta > 0$ ，使得

$$|E_r(x)| = \left| \frac{x - x^*}{x^*} \right| \leq \delta,$$

则称 δ 为近似值 x^* 的相对误差限。

例

$|x - x^*| \leq 1\text{cm}$ 时, 测量 100m 物体时的相对误差为

$$|E_r(x)| = \frac{1}{10000} = 0.01\%,$$

测量 10m 物体时的相对误差为

$$|E_r(x)| = \frac{1}{1000} = 0.1\%.$$

定义：有效数字

若近似值 x^* 的绝对误差限是某一位的半个单位，就称其精确到这一位，且从该位直到 x^* 的第一位非零数字共有 n 位，则称近似值 x^* 有 n 位有效数字。

给定数	具有 5 位有效数字的近似值
358.467	358.47
0.00427511	0.0042751
8.000034	8.0000
8.000034×10^3	8000.0

任何一个实数 x 经四舍五入后得到的近似值 x^* 都可写成

$$x^* = \pm(\alpha_1 \times 10^{-1} + \alpha_2 \times 10^{-2} + \cdots + \alpha_n \times 10^{-n}) \times 10^m.$$

当其绝对误差限满足

$$|x - x^*| \leq \frac{1}{2} \times 10^{m-n}$$

时, 则称近似值 x^* 具有 n 位有效数字, 其中 m 为整数, α_1 为 1 到 9 中的一个数字, $\alpha_1, \cdots, \alpha_n$ 是 0 到 9 中的数字。

例

验证 $e = 2.718281828459046 \cdots$ 的近似值 2.71828 具有 6 位有效数字

解

$$2.71828 = (2 \times 10^{-1} + 7 \times 10^{-2} + 1 \times 10^{-3} + 8 \times 10^{-4} + 2 \times 10^{-5} + 8 \times 10^{-6}) \times 10$$

$$\implies m = 1, n = 6.$$

$$\therefore |e - 2.71828| = 0.000001828 \cdots < \frac{1}{2} \times 10^{-5},$$

\therefore 它具有 6 位有效数字

定理

设某数 x 的近似值为 x^* , 则

$$x^* \text{ 有 } n \text{ 位有效数字} \implies |x - x^*| \leq \frac{1}{2} \times 10^{m-n}.$$

注

当 m 一定时, 有效数字位数 n 越多, 其绝对误差限越小。

定理

设某数 x 的近似值为 x^* 有如下形式

$$x^* = \pm(\alpha_1 \times 10^{-1} + \alpha_2 \times 10^{-2} + \cdots + \alpha_n \times 10^{-n}) \times 10^m,$$

则

- 若 x^* 有 n 位有效数字, 则

$$|E_r(x)| \leq \frac{1}{2\alpha_1} \times 10^{-(n-1)}.$$

- 若

$$|E_r(x)| \leq \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)}$$

则 x^* 至少有 n 位有效数字。

证明

由 $x^* = \pm(\alpha_1 \times 10^{-1} + \alpha_2 \times 10^{-2} + \cdots + \alpha_n \times 10^{-n}) \times 10^m$ 可知

$$\alpha_1 \times 10^{m-1} \leq |x^*| \leq (\alpha_1 + 1) \times 10^{m-1}$$

所以,

$$|E_r^*(x)| = \frac{|x - x^*|}{|x^*|} \leq \frac{\frac{1}{2} \times 10^{m-n}}{\alpha_1 \times 10^{m-1}} = \frac{1}{2\alpha_1} \times 10^{-(n-1)}$$

反之, 由

$$\begin{aligned} |x - x^*| &= |x^*| \cdot |E_r^*(x)| \leq (\alpha_1 + 1) \times 10^{m-1} \cdot \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)} \\ &= \frac{1}{2} \times 10^{m-n} \end{aligned}$$

知, x^* 至少具有 n 位有效数字。

例

为使 $\sqrt{20}$ 的近似值的相对误差小于 1%，问至少应取几位有效数字？

例

为使 $\sqrt{20}$ 的近似值的相对误差小于 1%，问至少应取几位有效数字？

解

$\sqrt{20}$ 的近似值的首位非零数字是 $\alpha_1 = 4$ ，故

$$|E_r^*(x)| = \frac{1}{2 \times 4} \times 10^{-(n-1)} < 1\%$$

可解得 $n > 2$ ，故取 $n = 3$ 即可满足要求。

数值计算的误差估计

数值计算的误差估计

设有可微函数 $y = f(x_1, x_2, \dots, x_n)$, 其中 $x_1^*, x_2^*, \dots, x_n^*$ 分别是 x_1, x_2, \dots, x_n 的近似值, 记 y 的近似值为 $y^* = f(x_1^*, x_2^*, \dots, x_n^*)$.

$$\begin{aligned} E(y^*) &= y - y^* = f(x_1, x_2, \dots, x_n) - f(x_1^*, x_2^*, \dots, x_n^*) \\ &\approx \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*} \cdot (x_i - x_i^*) = \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*} \cdot E(x_i^*), \end{aligned}$$

$$\begin{aligned} E_r(y^*) &= \frac{E(y^*)}{y^*} \approx \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*} \cdot \frac{x_i^*}{y^*} \cdot \frac{E(x_i^*)}{x_i^*} \\ &= \sum_{i=1}^n \frac{x_i^*}{y^*} \cdot \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*} \cdot E_r(x_i^*). \end{aligned}$$

定义

称

$$\frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*}, \quad \frac{x_i^*}{y^*} \cdot \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i^*}$$

为各个 x_i^* 对 y^* 的绝对误差和相对误差的增长因子。

数值计算的误差估计

例

试估计 $y = f(x_1, x_2, \dots, x_n) = x_1 + x_2 + \dots + x_n$ 的绝对误差与相对误差。

数值计算的误差估计

例

试估计 $y = f(x_1, x_2, \dots, x_n) = x_1 + x_2 + \dots + x_n$ 的绝对误差与相对误差。

解

因 $\frac{\partial f}{\partial x_i} = 1$, 故

- 和的绝对误差不会超过各项的绝对误差之和, 因

$$E(y^*) = \sum_{i=1}^n E(x_i^*) \implies |E(y^*)| \leq \sum_{i=1}^n |E(x_i^*)|,$$

即

数值计算的误差估计

例

试估计 $y = f(x_1, x_2, \dots, x_n) = x_1 + x_2 + \dots + x_n$ 的绝对误差与相对误差。

解

因 $\frac{\partial f}{\partial x_i} = 1$, 故

- 和的绝对误差不会超过各项的绝对误差之和, 因

$$E(y^*) = \sum_{i=1}^n E(x_i^*) \implies |E(y^*)| \leq \sum_{i=1}^n |E(x_i^*)|,$$

即

- 和的相对误差不会超过各项中最大的相对误差, 因

$$|E_r(y^*)| \leq \sum_{i=1}^n \left| \frac{x_i^*}{y^*} \right| |E_r(x_i^*)| \leq \max_{1 \leq i \leq n} |E_r(x_i^*)| \sum_{i=1}^n \frac{|x_i^*|}{|y^*|} = \max_{1 \leq i \leq n} |E_r(x_i^*)|.$$

数值计算的误差估计

例

测得某桌面的长 a 的近似值为 $a^* = 120\text{cm}$, 宽 b 的近似值为 $b^* = 60\text{cm}$, 若已知 $|a - a^*| \leq 0.2\text{cm}$, $|b - b^*| \leq 0.2\text{cm}$, 试求近似面积 $S^* = a^* b^*$ 的绝对误差限和相对误差限。

解

因 $S = ab$, $\frac{\partial S}{\partial a} = b$, $\frac{\partial S}{\partial b} = a$, 故

$$E(S^*) \approx \frac{\partial S^*}{\partial a^*} E(a^*) + \frac{\partial S^*}{\partial b^*} E(b^*) = b^* E(a^*) + a^* E(b^*)$$

从而有

$$|E(S^*)| \leq |60 \times 0.2| + |120 \times 0.1| = 24 \text{ cm}^2.$$

$$|E_r(S^*)| = \left| \frac{E(S^*)}{S^*} \right| = \frac{E(S^*)}{a^* b^*} \leq \frac{24}{7200} \approx 0.33\%.$$

选用和设计算法应遵循的原则

选用和设计算法应遵循的原则

一、选用数值稳定的计算公式，控制舍入误差的传播

若算法不稳定，则数值计算的结果就会严重背离数学模型的真实结果。因此在选择数值计算公式来进行近似计算时，应特别注意选用那些在计算过程中不会导致误差迅速增长的计算公式。

选用和设计算法应遵循的原则

例

计算积分

$$I_n = e^{-1} \int_0^1 x^n e^x dx, \quad n = 0, 1, 2, \dots$$

选用和设计算法应遵循的原则

算法 1:

$$\begin{cases} I_n = 1 - nI_{n-1}, \\ I_0 = 1 - e^{-1} \approx 0.6321. \end{cases}$$

matlab code

```
t0 = 0.6321;  
for i = 1:9  
    fprintf('%10.5f\n', t0);  
    t1 = 1 - i * t0;  
    t0 = t1;  
end
```

running result

```
0.63210  
0.36790  
0.26420  
0.20740  
0.17040  
0.14800  
0.11200  
0.21600  
-0.72800
```

误差分析

$$0 < I_n < e^{-1} \max_{0 \leq x \leq 1} (e^x) \int_0^1 x^n dx = \frac{1}{n+1}$$

知

$$I_7 < \frac{1}{8} = 0.1250, \quad I_8 < \frac{1}{9} \approx 0.1111,$$

I_0 本身有不超过 0.5×10^{-4} 的舍入误差，此误差在运算中传播、积累很快，传播到 I_7 和 I_8 时，该误差已放大了 7 与 8 倍，从而使得 I_7 和 I_8 的结果面目全非。

选用和设计算法应遵循的原则

算法 2:

$$I_{n-1} = \frac{1}{n}(1 - I_n)$$

matlab code

```
t0 = 0.1124;  
for i = 7:-1:0  
    fprintf('%10.5f\n', t0);  
    t1 = (1 - t0) / i;  
    t0 = t1;  
end
```

running result

```
0.11240  
0.12680  
0.14553  
0.17089  
0.20728  
0.26424  
0.36788  
0.63212
```

误差分析

由

$$I_n > e^{-1} \min_{0 \leq x \leq 1} (e^x) \int_0^1 x^n dx = \frac{e^{-1}}{n+1}$$

知

$$\frac{e^{-1}}{n+1} < I_n < \frac{1}{n+1}.$$

$$I_7 \approx 0.1124.$$

选用和设计算法应遵循的原则

定义：数值稳定

在数值计算中，误差不会增长的计算格式称为是数值稳定的，否则就是不稳定的。

二、尽量简化计算步骤以便减少运算次数

节省计算量，提高计算速度，简化逻辑结构，减少误差积累。

例

计算多项式

$$P_n(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$$

选用和设计算法应遵循的原则

- 逐项计算

共需 $1 + 2 + \cdots + (n-1) + n = \frac{1}{2}n(n+1)$ 次乘法和 n 次加法

- 秦九韶算法

$$\begin{cases} u_0 = a_n, \\ u_k = u_{k-1}x + a_{n-k}, \quad k = 1, 2, \cdots, n. \end{cases}$$

共需 n 次乘法和 n 次加法。

选用和设计算法应遵循的原则



秦九韶 (1208 年 - 1261 年)，南宋官员、数学家，与李冶、杨辉、朱世杰并称宋元数学四大家。

数学成就

- 数学九章
- 大衍求一术：比高斯的同余理论早 554 年
- 任意次方程的数值解法：比英国人霍纳早提出 572 年
- 三斜求积术：海伦公式 (公元 50 年左右)
- 秦九韶公式
- ...

三、避免两个相近的数相减

数值计算中，两个相近的数相减会造成有效数字的严重丢失。常用的处理办法有：

- 因式分解
- 分子分母有理化
- 三角函数恒等式
- Taylor 展开式
- ...

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$\sqrt{x+1} - \sqrt{x} \quad (x = 1000)$$

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$\sqrt{x+1} - \sqrt{x} \quad (x = 1000)$$

- 直接计算

$$\sqrt{1001} - \sqrt{1000} \approx 31.64 - 31.62 = 0.02$$

只有一个有效数字，损失了三位有效数字

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$\sqrt{x+1} - \sqrt{x} \quad (x = 1000)$$

- 直接计算

$$\sqrt{1001} - \sqrt{1000} \approx 31.64 - 31.62 = 0.02$$

只有一个有效数字，损失了三位有效数字

- 分子有理化

$$\sqrt{x+1} - \sqrt{x} = \frac{1}{\sqrt{x+1} + \sqrt{x}} \approx 0.01581$$

没有损失有效数字

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$A = 10^7(1 - \cos 2^\circ) \quad (\cos 2^\circ = 0.9994)$$

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$A = 10^7(1 - \cos 2^\circ) \quad (\cos 2^\circ = 0.9994)$$

- 直接计算

$$A \approx 10^7(1 - 0.9994) = 6 \times 10^3$$

只有一个有效数字

选用和设计算法应遵循的原则

例

计算（取 4 位有效数字）

$$A = 10^7(1 - \cos 2^\circ) \quad (\cos 2^\circ = 0.9994)$$

- 直接计算

$$A \approx 10^7(1 - 0.9994) = 6 \times 10^3$$

只有一个有效数字

- 三角恒等式

$$1 - \cos x = 2 \sin^2 \frac{x}{2}$$

$$A = 2 \times (\sin 1^\circ)^2 \times 10^7 \approx 2 \times 0.01745^2 \times 10^7 \approx 6.09 \times 10^3$$

有三位有效数字

四、绝对值太小的数不宜做除数

数值计算中，用绝对值很小的数作除数，会使商的数量级增加，甚至在计算机中造成“溢出”停机，而且当很小的除数稍有一点误差，会对计算结果影响很大。

例

$$\begin{array}{rcl} \frac{3.1416}{0.001} & = & 3141.6 \\ \frac{3.1416}{0.001 + 0.0001} & = & 2856 \end{array}$$

五、合理安排运算次序，防止“大数吃小数”

例

计算 a, b, c 的和，其中 $a = 10^{12}$, $b = 10$, $c \approx -a$.

- $(a + b) + c$

结果接近于 0

- $(a + c) + b$

结果接近于 10