

计算方法

黄明游 梁振珊 编

吉林大学出版社

计 算 方 法

黄明游 梁振珊 编

吉林大学出版社

计算方法
黄明游 梁振珊 编

责任编辑、责任校对：赵洪波

封面设计：孙 泓

吉林大学出版社出版

吉林省新华书店发行

(长春市东中华路 29 号)

长春大学印刷厂印刷

开本：850×1168 毫米 1/32

1994 年 6 月第 1 版

印张：7.75

1994 年 6 月第 1 次印刷

字数：189 千字

印数：1—1100 册

ISBN 7-5601-1581-0/O · 182

定价：4.50 元

前　　言

本书是适应为非计算数学的有关专业开设〈计算方法〉课程的需要而编写的。计算方法的内容极为广泛，目前计算数学专业一般是分三门课程（数值逼近、数值代数、微分方程数值解法）讲授其内容，需要花费近200个学时。然而，对于其他有关专业，上述作法是既不可能，同时也并非是可取的。编写这本教材的目的，是便于在一个学期（约用68个学时）内讲授科学与工程计算中遇到的一些基本问题的计算方法，介绍计算方法的原理和选讲一些典型的、便于在计算机上应用的数值计算方法。

这本书的内容大纲曾在理科力学专业教材会议（1985，兰州）上与北京大学、复旦大学、兰州大学和中山大学的同事们进行过讨论，他们提供了宝贵的修改意见。根据这个大纲，经本书两名编者的合作，于1986年编写出一本‘实用计算方法’讲义。该讲义曾作为吉林大学数学系力学专业、应用数学专业和计算机科学系讲授《计算方法》课的教材使用过多遍。本书就是在这个讲义和历年来该课教学实践的基础上，经过重新修改、补充之后写成的。其中第一至第四章由梁振珊同志编写，第五、六两章由黄明游同志编写，大纲的修订和全书内容的校核是黄明游同志完成的。

吉林大学计算数学教研室常玉堂、于加艾、马富明等同志曾先后利用上述教材讲授过计算方法课，他们提供过许多宝贵的修改意见，在此向这些同事表示亲切的谢意。特别地，这里要向参加本教材大纲讨论并给予指导的外校和本校的老师们表示衷心的感谢。

由于水平和经验的原因，本书在取材、衔接、表述等方面一定还存在许多不足乃至某些错误，恳请批评、指正。

编 者
1994 年 4 月

目 录

绪言	(1)
数的浮点表示和浮点四则运算	(1)
研究计算方法的必要性	(3)
计算方法的主要内容	(6)
第一章 线性代数方程组的解法	(8)
§ 1 消元法和矩阵分解	(8)
§ 1.1 Gauss 消元法的消元过程	(8)
§ 1.2 回代过程	(10)
§ 1.3 计算量和存贮	(11)
§ 1.4 矩阵的 LU 分解	(12)
§ 1.5 Doolittle 分解	(14)
§ 1.6 对称正定矩阵的 LDL^T 分解	(16)
§ 2 消元法在计算机上的实现	(18)
§ 2.1 选取主元的必要性	(18)
§ 2.2 选主元的方法	(19)
§ 2.3 迭代改善	(20)
§ 2.4 行列式和逆矩阵的计算	(20)
§ 3 条件数	(21)
§ 3.1 病态方程组举例	(21)
§ 3.2 向量模	(22)
§ 3.3 矩阵模	(23)
§ 3.4 矩阵的谱半径	(26)
§ 3.5 矩阵的条件数	(26)
§ 3.6 矩阵的平衡	(28)
§ 4 迭代法	(29)

§ 4.1 一般迭代法	(29)
§ 4.2 几种常用的迭代法	(31)
§ 4.3 对角占优矩阵和不可约矩阵	(35)
§ 4.4 迭代法收敛的充分条件.....	(38)
§ 5 共轭斜量法	(39)
§ 5.1 等价的极值问题	(40)
§ 5.2 最速下降法	(40)
§ 5.3 共轭斜量法	(41)
§ 5.4 共轭斜量法的收敛性	(44)
提要和注记	(46)
习题	(47)
第二章 矩阵特征值和特征向量的计算	(51)
§ 1 乘幂法	(51)
§ 1.1 乘幂法	(51)
§ 1.2 反幂法	(54)
§ 1.3 原点位移	(54)
§ 1.4 加速技巧	(55)
§ 1.5 压缩	(57)
§ 2 Jacobi 方法	(60)
§ 2.1 旋转变换	(60)
§ 2.2 Jacobi 方法	(62)
§ 2.3 Jacobi 方法在计算机上的实现	(64)
§ 3 Givens-Householder 方法	(65)
§ 3.1 三对角化过程	(65)
§ 3.2 求对称三对角矩阵特征值的对分法	(71)
§ 3.3 特征向量的计算	(75)
§ 4 QR 算法	(76)
§ 4.1 QR 算法的基本思想.....	(76)
§ 4.2 Hessenberg 型矩阵的 QR 算法	(77)
§ 4.3 具有原点位移的 QR 算法	(80)
§ 4.4 双步 QR 算法	(81)

§ 4.5 特征向量的计算	(83)
提要和注记	(84)
习题	(85)
第三章 非线性方程(组)的解法和最优化问题的	(88)
计算方法	(89)
§ 1 非线性方程的解法	(89)
§ 1.1 对分法	(89)
§ 1.2 迭代法	(90)
§ 1.3 Newton 法	(91)
§ 2 解非线性方程组的 Newton 法	(95)
§ 3 一维搜索	(98)
§ 3.1 搜索区间	(98)
§ 3.2 黄金分割法	(99)
§ 3.3 抛物线插值法	(101)
§ 4 梯度法	(104)
§ 4.1 最速下降法	(104)
§ 4.2 共轭斜量法	(105)
§ 4.3 无约束最优化问题的 Newton 法	(108)
§ 4.4 变尺度法	(109)
§ 5 直接搜索法	(112)
§ 5.1 步长加速法(Hooke-Jeeves 方法)	(112)
§ 5.2 单纯形法	(115)
提要和注记	(119)
习题	(120)
第四章 函数插值和数值积分	(123)
§ 1 Lagrange 插值公式和 Hermite 插值公式	(123)
§ 1.1 Lagrange 插值公式	(123)
§ 1.2 Hermite 插值公式	(125)
§ 2 Newton 插值公式	(127)
§ 2.1 差商和 Newton 插值公式	(127)

§ 2.2 差分和等距节点的 Newton 插值公式	(129)
§ 3 样条函数插值	(132)
§ 4 数值积分	(137)
§ 4.1 Newton-Cotes 求积公式	(138)
§ 4.2 梯形求积公式和抛物线求积公式	(139)
§ 4.3 分半加速算法(Romberg 算法)	(142)
§ 4.4 Gauss 型求积公式	(144)
§ 5 快速富氏变换	(150)
§ 5.1 离散富氏变换	(150)
§ 5.2 离散富氏变换的快速算法	(152)
提要和注记	(157)
习题	(158)
第五章 常微分方程初值问题的数值解法	(161)
§ 1 单步方法	(163)
§ 1.1 Euler 方法及其改进	(164)
§ 1.2 Runge-Kutta 方法	(166)
§ 1.3 Richardson 外推法	(170)
§ 2 收敛性和稳定性	(173)
§ 2.1 协调性	(174)
§ 2.2 收敛性	(175)
§ 2.3 稳定性	(179)
§ 2.4 整体误差	(181)
§ 3 多步方法	(182)
§ 3.1 Adams 外插方法	(182)
§ 3.2 Adams 内插方法	(184)
§ 3.3 一般的线性多步方法	(186)
§ 3.4 多步方法的应用技巧	(188)
习题	(191)
第六章 偏微分方程的数值解法	(193)
§ 1 边值问题的差分法	(193)
§ 1.1 边值问题的差分逼近	(194)

§ 1.2 差分解的存在、唯一性和收敛性	(197)
§ 1.3 差分方程的求解-逐次超松弛法	(202)
§ 2 初值问题的差分法	(205)
§ 2.1 抛物型方程	(205)
§ 2.2 差分格式的稳定性	(210)
§ 2.3 双曲型方程	(215)
§ 3 有限元方法	(221)
§ 3.1 变分原理	(221)
§ 3.2 Ritz-Galerkin 方法	(226)
§ 3.3 有限元逼近	(228)
习题	(234)

绪 言

用电子计算机进行科学计算，解决生产实际问题，其过程大致如下：



我们所说的计算方法，是指在数字电子计算机上使用的计算方法。计算机所执行的算法由算术运算和逻辑运算组成。把数学模型归结为计算机能执行的有效算法，便是计算方法课程讨论的课题。为此，初步了解一下数的计算机表示和运算是必要的。

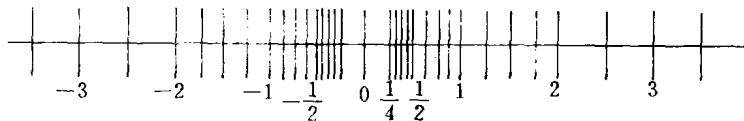
数的浮点表示和浮点四则运算

目前，数字电子计算机大多使用二进制浮点计算系统。在这种二进制浮点计算机上，每个实数 x 具有如下形式：

$$x = \pm 0.d_1d_2\cdots d_k \times 2^p. \quad (1)$$

这里， $d_1=1$ ， d_i 为 0 或 1， $i=2, \dots, k$ 。整数 p 满足 $-m \leq p \leq M$ 。零的浮点表示可能是不同的，通常用 $0=\pm 0.00\cdots 0 \times 2^{-m}$ 表示。正整数 k, m, M 随计算机而定。因此，一部给定的二进制浮点计算机，只能表示所有形如 (1) 的有限数集 $S=S(k, m, M)$ 。这是实数轴上的不等距的有限点集。为了说明问题，我们考虑一个特殊的计算机，它的 $k=3, m=1, M=2$ ，那么这个计算机所能表示的浮点数的集合 $S(3, 1, 2)$ 只有如下图的 33 个

点：



一个实数 x ，在计算机上只能按一定规则用 S 中最接近 x 的数来近似表示。这就带来了数的表示误差。在我们这个小计算机上， $|x| > 3.5$ 就溢出了（上溢出）， $|x| < 0.25$ 就是机器零（下溢出）。

计算机只能对浮点数集 S 中的数做加、减、乘、除四则运算。由于用二进制数作四则运算我们不习惯，所以这里不妨用 10 进制的 4 位浮点数的四则运算来说明，因为它们的情况是一样的。设 S_0 是所有形如

$$\pm 0.d_1 d_2 d_3 d_4 \times 10^p$$

的 4 位 10 进制浮点数的集合，其中 $1 \leq d_i \leq 9$, $0 \leq d_i \leq 9$, $i = 2, 3, 4$, 整数 p 满足 $-9 \leq p \leq 10$. 下面举例说明 S_0 上的算术运算，“ \rightarrow ” 表示在运算器上的运算过程和结果。

- 1) $0.1984 \times 10^4 + 0.2008 \times 10^2$
 $\rightarrow 0.1984 \times 10^4 + 0.0020 \times 10^4$
 $\rightarrow 0.2004 \times 10^4$
- 2) $0.1984 \times 10^4 + 0.9876 \times 10^{-1}$
 $\rightarrow 0.01984 \times 10^4 + 0.0000 \times 10^4$
 $\rightarrow 0.1984 \times 10^4$
- 3) $0.1984 \times 10^{-6} - 0.1976 \times 10^{-6}$
 $\rightarrow 0.0008 \times 10^{-6}$
 $\rightarrow 0.8000 \times 10^{-9}$
- 4) $(0.5678 \times 10^3) \times (0.6789 \times 10^{-5})$
 $\rightarrow (0.5678 \times 0.6789) \times 10^{-2}$

$$\begin{aligned}
 &\rightarrow 0.38547942 \times 10^{-2} \\
 &\rightarrow 0.3855 \times 10^{-2} \\
 5) \quad &(0.5678 \times 10^4) \div (0.4567 \times 10^{-5}) \\
 &\rightarrow 0.5678 \times 10^9 \div 0.4567 \\
 &\rightarrow 0.05678000 \times 10^{10} \div 0.4567 \\
 &\rightarrow 0.12432669 \times 10^{10} \\
 &\rightarrow 0.1243 \times 10^{10}
 \end{aligned}$$

由这些例子可以看出，计算机上的算术运算不可避免地带来舍入误差。尤其应当注意，尽量避免象 2)、3)、5) 那样的情况发生：绝对值大小悬殊的两个数相加、减；相近的数相减；相对被除数来说绝对值很小的数作除数。这些是数值计算中的基本原则。在 5) 中当除数为 0.4567×10^{-6} 时，商 0.1243×10^{11} 已不属于浮点数集 S_0 ，这就是所谓“溢出”。

研究计算方法的必要性

了解了数的计算机表示和计算机的算术运算，就不难理解为什么要讨论计算方法。这里，我们通过几个简单的例题，来进一步说明这个问题。

例 1 在使用上述浮点数集 S_0 的 4 位 10 进制浮点计算机上，解一元二次方程

$$x^2 - 18x + 1 = 0.$$

按求根公式，两个根为

$$x_1 = 9 + \sqrt{80}, \quad x_2 = 9 - \sqrt{80}.$$

在我们的计算机上， $\sqrt{80} = 0.8944 \times 10^1$ 。于是得

$$x_1 = 0.1794 \times 10^2, \quad x_2 = 0.5600 \times 10^{-1}.$$

但是换一种办法计算 x_2 ，由 $x_1 x_2 = 1$ 得

$$x_2 = \frac{1}{9 + \sqrt{80}} = 0.5574 \times 10^{-1}.$$

它的误差不超过 $0.0001 \times 10^{-1} \times 0.5 = 0.000005$ (末位有效数字的一半); 而前面按求根公式算得的 x_2 , 其误差不超过 $0.01 \times 10^{-1} \times 0.5 = 0.0005$. 显然, 在数学上公式 $x_2 = 9 - \sqrt{80}$ 和 $x_2 = 1 / (9 + \sqrt{80})$ 是等价的, 但在计算机上它们是不同的. 从计算方法的观点看来, 后者远优于前者, 前者的缺欠在于两个相近数相减.

例 2 考虑对于给定 x 计算多项式

$$P_n(x) = a_0x^n + a_1x^{n-1} + \cdots + a_{n-1}x + a_n \quad (2)$$

的值. 显然, 它等价于

$$P_n(x) = (\cdots((a_0x + a_1)x + a_2)x + \cdots + a_{n-1})x + a_n \quad (3)$$

但从计算方法的观点看来, 两者却大不相同: 公式 (2) 需先算出 x^2, x^3, \dots, x^n , 共需 $n-1$ 次乘法; 还要保存下来, 需增加 $n-1$ 个存贮单元; 然后再相乘 (n 个乘法)、相加 (n 个加法), 共需 $2n-1$ 次乘法, n 次加法. 公式 (3) 只需 n 个乘法, n 个加法, 不用增加存贮单元. 公式 (2) 和 (3) 的差别远不止这些. 实际上, 在我们上述的计算机上, $x=10$ 时按 (2) 就不能计算, 因为 10^{10} 已经溢出.

例 3 考虑解线代数方程组

$$Ax = b. \quad (4)$$

设系数矩阵 A 是 $n \times n$ 方阵, 其行列式 $D = \det(A) \neq 0$. 按 Gramer 法则, (4) 的解为

$$x_i = D_i/D, \quad i = 1, 2, \dots, n.$$

共需计算 $n+1$ 个 n 阶行列式. 用 Laplace 展开计算 n 阶行列式, 需要

$$n! + \frac{n!}{2!} + \frac{n!}{3!} + \cdots + \frac{n!}{(n-1)!}$$

$$= n! \left(1 + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{(n-1)!} \right)$$

个乘法, 姑且算作 $n!$ 个乘法. 这样, 不计加法, 用 Gramer 法

则解 n 阶方程组, 用 Laplace 展开计算行列式, 共需 $(n+1)n!$
 $= (n+1)!$ 以上的乘法. 对于一个 20 阶的方程组, 就需要 $21! \doteq 5.11 \times 10^{19}$ 以上的乘法. 设用每秒可作百万次乘法的计算机, 一年可作 $365 \times 24 \times 3600 \times 10^6 \doteq 3.15 \times 10^{13}$ 次乘法. 所以, 在每秒作百万次乘法的计算机上, 用 Gramer 法则解 20 阶的线代数方程组, 需要的时间在 $(5.11 \times 10^{19}) \div (3.15 \times 10^{13}) = 1.62 \times 10^6 \doteq 162$ 万年以上!

例 4 考虑计算积分

$$I_n = \int_0^1 x^n e^{-x} dx. \quad (5)$$

它满足递推关系

$$I_n = 1 - n I_{n-1}. \quad (6)$$

我们首先算出 I_0 的近似值 $\bar{I}_0 = 0.6321$. 利用递推关系 (6) 依次算得:

\bar{I}_0	\bar{I}_1	\bar{I}_2	\bar{I}_3	\bar{I}_4	\bar{I}_5	\bar{I}_6	\bar{I}_7
0.6321	0.3680	0.2640	0.2080	0.1680	0.1600	0.0400	0.7200

显然 I_7 的递推计算结果 \bar{I}_7 是错误的. 因为

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

所以

$$I_7 < \frac{1}{8} = 0.1250.$$

实际上, 序列 $\{I_n\}$ 是一个下降序列. 为什么会产生这样面目全非的错误结果呢? 这是由于在计算 \bar{I}_0 时的原始误差 (表示误差和舍入误差) 以及计算 \bar{I}_1 的误差, 在按(6)递推时, 逐次乘以因子 $2, 3, \dots, 7$, 使误差急剧地增长和积累, 从而产生错误的没有意义的结果. 然而, 如果将(6)改写成

$$I_{n-1} = (1 - I_n)/n, \quad (7)$$

并先计算出 I_7 的近似值 $\bar{I}_7=0.1124$, 从 \bar{I}_7 出发按 (7) 递推算得:

\bar{I}_7	\bar{I}_6	\bar{I}_5	\bar{I}_4	\bar{I}_3	\bar{I}_2	\bar{I}_1	\bar{I}_0
0.1124	0.1269	0.1455	0.1708	0.2073	0.2643	0.3680	0.6320
0.1124	0.1268	0.1456	0.1709	0.2073	0.2642	0.3679	0.6321

表中末行是 I_n 的精确值的舍入结果. 可见按 (7) 递推计算出的 \bar{I}_n ($n=6, 5, \dots, 0$) 是令人满意的. 有趣的是, 即使从 $\bar{I}_7=0$ 开始按 (7) 递推, 也可得到 $\bar{I}_0=0.6320$. 这是由于原始误差每步依次乘以 $\frac{1}{7}, \frac{1}{6}, \dots, 1$, 从而被大大地缩小了.

通过上述例题可以看出, 即使有了数学模型, 即使数学上已经有了完善的结果, 仍然存在能不能在计算机上实现和如何实现的问题. 所以我们必须研究计算方法, 提供计算机上的有效算法.

计算方法的主要内容

对于给定的数学问题, 常常可以提出各式各样的算法. 怎样评价这些算法的优劣呢? 这要具体问题具体分析. 一般说来, 从计算方法的观点考虑, 自然提出以下要求:

第一, 计算结果可靠. 即精确度高、误差小. 通常原始数据是通过仪器观测得到的, 含有观测误差; 数据送入计算机, 产生计算机的表示误差; 计算机进行算术运算, 产生舍入误差; 计算机只能作有限的算术四则运算, 微商、积分、微分方程等只能通过算术运算给出近似公式. 这就产生所谓截断误差或余项.

观测误差和表示误差取决于仪器的精度和计算机，这里无须讨论。浮点运算的舍入误差分析，比较复杂，超出我们课程的范围。截断误差或余项，在我们讨论的算法中，一般都要给出。

从上面例 4 可以看到，计算结果的可靠性不仅和误差有关，而且和算法有关。像递推公式（6）那样的算法，将原始误差不断扩大（即使按公式精确计算也是如此），这样的算法说是不稳定的。像递推公式（7）那样的算法，使原始的误差得到控制，不再增长，这样的算法称作稳定的。为了保证计算结果可靠，我们使用的算法必须是稳定的。

第二，算法简单，运算量少。

第三，存贮量小。计算机的存贮设备有内存和外存，调用外存比使用内存慢得多，使用内存中的数据才能实现快速计算。而内存总是有限的（外存几乎可以看成是无限的），所以我们提出的算法，占用内存单元应当尽可能少。

在电子计算机飞速发展的今天，计算机的应用已经深入到科研、生产和生活的各个领域，计算机上使用的计算方法已浩如烟海。本书只限于介绍科学计算方面的最基本的计算方法。主要内容有：线代数方程组的数值解法和矩阵特征值、特征向量的计算，非线性方程和非线性方程组的解法，最优化问题的计算方法，函数插值和数值积分，常微分方程初值问题的数值解法和偏微分方程的数值解法。