



中国科学院规划教材

计算方法

孙文瑜 杜其奎 陈金如 编

$$l_{ik} = \left(a_{ik} - \sum_{j=1}^{k-1} \right)$$

$$u_{kj} = a_{kj} - \sum_{p=1}^{k-1} l_{kp} u_{pj}$$



科学出版社
www.sciencep.com

中国科学院规划教材

计算方法

孙文瑜 杜其奎 陈金如 编

科学出版社

北京

内 容 简 介

本书主要介绍计算方法中的一些基本内容：误差和条件问题、解线性方程组的直接法与迭代法、特征值问题的计算方法、解非线性方程和方程组的迭代法、插值与逼近、数值积分与数值微分以及常微分方程数值解法。本书内容深入浅出，既强调计算方法的基本概念和理论，更注重算法和实践。每章后面都附有一定数量的习题。

本书可作为各类高等院校本科生和研究生“计算方法”或“数值分析”课程的教材，也可作为从事科学工程计算的科技人员的参考书。

图书在版编目(CIP)数据

计算方法/孙文瑜, 杜其奎, 陈金如编. —北京: 科学出版社, 2007
(中国科学院规划教材)

ISBN 978-7-03-018487-0

I. 计… II. ①孙… ②杜… ③陈… III. 计算方法—高等学校—教材 IV. O241
中国版本图书馆 CIP 数据核字(2007) 第 012465 号

责任编辑: 赵 靖 / 责任校对: 邹慧卿

责任印制: 张克忠 / 封面设计: 耕者设计工作室

科学出版社出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

化学工业出版社印刷厂印刷

科学出版社发行 各地新华书店经销

*

2007 年 3 月第 一 版 开本: B5(720×1000)

2007 年 3 月第一次印刷 印张: 16 1/4

印数: 1—4 000 字数: 307 000

定价: 24.00 元

(如有印装质量问题, 我社负责调换(化工))

前　　言

计算方法，又称数值分析、数值计算或科学计算，是研究利用计算机求解科学工程中各种问题的数值方法。当今，科学、工程技术、经济、管理中的许多问题都必须依靠计算机和计算方法来求解。科学计算已经成为当今科学的研究的三种基本手段之一。本书为学生提供了各种常用的现代计算方法。

我国计算数学事业和计算数学教材的发展已经有五十多年的历史。1959年，北京大学、南京大学、吉林大学“计算方法”编写组编写了我国第一本计算数学教材《计算方法》，开启了我国计算数学教学的篇章。1977年高考恢复后，为了适应新的形势，冯康院士、何旭初教授等老一辈计算数学家出版了《数值计算方法》（冯康等，1978）、《计算数学简明教程》（何旭初等，1980）等一批优秀教材。进入21世纪，大量的现代计算机和计算数学方面的教材纷纷面世，展示了我国计算数学教学和研究的一片新气象。

现在，在我国高等学校，除了信息与计算科学专业开设一年的计算方法课程外，其他学科，如理科的数学、物理、化学、天文、生物、地学、计算机科学与技术和各种工科专业，以及经济、金融、管理等学科都开设一学期的计算方法课程。一些学校理科和工科的研究生也要开设一学期的计算方法课。根据这样的需要和特点，我们在自己多年从事计算方法教学的基础上，参考了国内外计算方法教材，编写了这本书，适合一学期“计算方法”课程的教学。

本书覆盖科学计算的主要内容，包括误差和条件问题，解线性方程组的直接法与迭代法，特征值问题的计算方法，解非线性方程和方程组的迭代法，插值与逼近，数值积分与数值微分，常微分方程数值解法等8章内容。通常，偏微分方程数值解法和最优化方法作为本课程的后续课程，故本书没有包括，需要的学校和专业可以接着本课程开设有关的后续课程。

计算方法是一门应用性很强的课程，它来源于科学工程实际，又反过来应用于科学工程实际。因此，无论是数学类的还是理工科或经管类的学生，学好计算方法这门课，将对他（她）本身的专业业务和科学研究有极大的好处。因此，学习计算方法，不是单纯的纸上谈兵，必须要在计算机环境下通过算法分析、算法设计、上机计算等实践环节来完成。因此，希望讲授这门课的老师和学习这门课的学生，都要高度重视编程和上机算题这样的实践环节。本书每章都配有一定量的习题和必要的上机实验题，学生只要掌握了任何一种程序设计语言（如C语言，FORTRAN语言，或MATLAB语言），就不难将本书的算法或伪程序编成相应的程序上机计算，并完成习题和上机实验题。

本书是我们在南京大学、南京师范大学和国外一些大学多年来从事计算方法课程教学的讲稿基础上整理而成的。本书力求通俗、系统、简明、深入浅出。本书虽然强调计算方法的基本概念和基本理论，但对于重要的证明、复杂的定理，只叙述而省略证明，感兴趣的读者可以参考有关文献。我们希望学习本书的本科生和研究生，尤其是理工科的朋友，既要学好每一章的具体内容，了解各种方法的基本思想、推导、算法特性及算法分析，并编程算题，更要高屋建瓴，总体上把握数值分析中大的本质的方面，弄清各种方法之间的相互关系和区别，进而应用这些方法和研究新的方法去解决科学工程中的具体问题。

本书可以作为各类高等院校一学期“计算方法”或“数值分析”课程的教材，每周3学时或4学时(上机实践课时另外安排)，教师可根据具体情况对本书内容进行选择和增删。本书后面附了部分经典的数值分析和计算方法方面的参考文献，供感兴趣的读者选择阅读。

在本书的写作过程中，得到了许多专家和同行的关心、支持和帮助。这里，作者首先要感谢中国科学院数学与系统科学研究院计算数学与科学工程计算研究所石钟慈院士、林群院士、崔俊芝院士、袁亚湘研究员、余德浩研究员，香港城市大学祁力群教授，南京师范大学校长宋永忠教授等。此外，我们的学生对本书初稿提出了不少有益的建议，并协助做了一定的录入工作，谨向他们表示谢意。感谢科学出版社的编辑对本书写作所给予的热情指导和关心，确保本书顺利出版。作者还要感谢国家自然科学基金委员会多年来对作者研究工作的资助。

尽管本书作者多年来一直从事计算方法的教学和研究，但疏漏和不当之处在所难免，恳请专家和读者予以批评指正。

孙文瑜 杜其奎 陈金如

2006年7月15日

于南京师范大学随园

目 录

第 1 章 绪论	1
1.1 误差的基本概念	2
1.1.1 误差的来源	2
1.1.2 绝对误差与相对误差	4
1.1.3 算术运算的相对误差	5
1.1.4 有效数字	6
1.2 算法设计中应注意的问题	7
习题 1	12
第 2 章 解线性方程组的直接法	13
2.1 引言	13
2.2 消去法	14
2.2.1 Gauss 消去法	14
2.2.2 选主元消去法	19
2.3 矩阵的 LU 分解法	27
2.4 平方根法	31
2.5 追赶法	34
2.5.1 带状矩阵	34
2.5.2 追赶法	36
2.6 向量与矩阵的范数	38
2.6.1 向量范数	38
2.6.2 矩阵范数	41
2.7 误差分析	45
习题 2	48
第 2 章上机实验题	52
第 3 章 解线性方程组的迭代法	53
3.1 引言	53
3.2 迭代法的一般格式及收敛性条件	54
3.2.1 迭代法的一般格式	54
3.2.2 迭代法的收敛性条件	55
3.3 Jacobi(雅可比)迭代法	58

3.4 Gauss-Seidel(高斯-赛德尔)迭代法	61
3.5 逐次超松弛迭代法(SOR 方法)	63
3.6 迭代法的收敛性	66
习题 3	71
第 3 章上机实验题	73
第 4 章 特征值问题的计算方法	75
4.1 特征值问题的基本理论	75
4.2 乘幂法与反乘幂法	81
4.3 QR 方法	87
4.3.1 Givens 变换和 Householder 变换	87
4.3.2 化矩阵为上 Hessenberg 矩阵	91
4.3.3 QR 方法	95
4.3.4 对上 Hessenberg 矩阵采用 QR 方法	97
4.3.5 带原点平移的 QR 方法	98
习题 4	100
第 4 章上机实验题	102
第 5 章 解非线性方程和方程组的迭代法	103
5.1 迭代序列收敛的基本概念	103
5.2 不动点迭代	106
5.3 解非线性方程的几个方法	110
5.3.1 二分法	110
5.3.2 牛顿法	113
5.3.3 割线法	118
5.3.4 弦方法	121
5.4 解非线性方程组的牛顿法及其变形	122
5.4.1 解非线性方程组的牛顿法	122
5.4.2 修改牛顿法简介	126
5.5 解非线性方程组的割线法	130
习题 5	135
第 5 章上机实验题	136
第 6 章 插值与逼近	137
6.1 Lagrange 插值	138
6.1.1 插值基函数	139
6.1.2 Lagrange 插值多项式	140
6.1.3 插值余项	141

6.2 Hermite 插值	143
6.3 差分	147
6.3.1 差分及其基本性质	147
6.3.2 高阶差分的表达式	149
6.4 Newton 插值公式	151
6.4.1 逐步插值多项式	151
6.4.2 差商与 Newton 插值公式	152
6.4.3 差商表	153
6.4.4 等距节点插值公式	158
6.4.5* 带重节点差商	160
6.5 分段低次插值	161
6.5.1 分段线性插值	162
6.5.2 分段三次 Hermite 插值	163
6.6* 三次样条插值	165
6.6.1 样条函数的概念	166
6.6.2 三次样条的构造	166
6.6.3 边界条件	168
6.6.4 计算的基本步骤	169
6.7* 正交多项式与最佳平方逼近	169
6.7.1 正交函数系的概念	169
6.7.2 正交多项式	170
6.7.3 用正交多项式作最佳平方逼近	174
习题 6	177
第 6 章 上机实验题	180
第 7 章 数值积分与数值微分	181
7.1 复化矩形公式、复化梯形公式和抛物线公式	182
7.1.1 复化矩形公式、复化梯形公式及其截断误差	182
7.1.2 抛物线公式及其截断误差	184
7.1.3 复化抛物线公式及其截断误差	186
7.2 Newton-Cotes 求积公式	188
7.3 Romberg 求积法	190
7.3.1 Euler-Maclaurin 公式	190
7.3.2 梯形公式的二分技术	191
7.3.3 Richardson 外推法与抛物线公式	192
7.3.4 Romberg 求积法	193

7.4 Gauss 型求积公式	195
7.4.1 Gauss 型求积公式	196
7.4.2 常用的两个 Gauss 型求积公式	199
7.5* 应用样条插值的求积公式	201
7.6 数值微分	202
7.6.1 用插值多项式求数值导数	202
7.6.2 用幂级数展开式求数值导数	205
7.6.3 用外推法求数值导数	206
7.6.4* 用三次样条插值方法求数值导数	208
习题 7	209
第 7 章上机实验题	210
第 8 章 常微分方程数值解法	211
8.1 引言	211
8.2 Euler 方法	214
8.2.1 Euler 格式	214
8.2.2 Euler 格式的误差分析	217
8.2.3 Euler 方法的收敛性与稳定性	219
8.3 预估-校正法	222
8.3.1 改进的 Euler 方法	222
8.3.2 预估-校正法	224
8.4 Runge-Kutta(龙格-库塔)法	230
8.4.1 二阶 Runge-Kutta 法	230
8.4.2 三阶 Runge-Kutta 法	231
8.4.3 四阶 Runge-Kutta 方法	233
8.5 线性多步法	236
8.5.1 线性二步法	237
8.5.2 Adams(亚当斯)外推法	239
8.5.3 Adams 内插法	241
8.6 单步法的收敛性与稳定性	243
8.6.1 单步法的收敛性	244
8.6.2 单步法的绝对稳定性	245
习题 8	247
第 8 章上机实验题	249
参考文献	250

第1章 絮 论

在解决一些实际问题时,通常要将其归结为数值计算问题. 所谓数值计算,就是用计算机等计算工具来求出数学问题的数值解的全过程. 具体地说,是指由一组已知数据(通常称之为输入数据),求出一组数值(称之为输出数据),使得这两组数据之间满足预先指定的某种关系. 目前随着计算机的广泛应用,使得越来越多的实际问题,通过数值计算而得到很好的解决. 数值计算的重要性,使得科学计算与科学理论和科学实验相伴而成为当今世界科学活动的第三种手段.

用数值计算的方法来解决具体的实际问题时,首先必须将具体的问题抽象为数学问题,即建立起能描述并等价代替该实际问题的数学模型(数学建模),然后提出合适的算法,编制出计算机程序,最后上机调试并进行运算,以得到所需的结果. 这里所说的“算法”,是指由基本运算和运算顺序的规定所组成的整个解题方案和步骤.

建立和选择合适的算法是整个数值计算中非常重要的一环. 例如,计算如下多项式

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

的值时,若我们直接计算 $a_i x^i$ ($i = 0, 1, \dots, n$) 后,再逐项相加,需要进行

$$1 + 2 + \cdots + n = \frac{1}{2}n(n+1)$$

次乘法和 n 次加法. 如果使用著名的秦九韶算法,先将多项式写成如下形式计算

$$P_n(x) = (\cdots ((a_0x + a_1)x + \cdots + a_{n-2})x + a_{n-1})x + a_n,$$

只要进行 n 次乘法和 n 次加法即可(总共 $2n$ 次运算). 当 n 较大时,后者的计算量明显地小于前者的计算量. 因此,算法的优劣直接影响到计算的速度和效率. 算法选的不恰当,不仅会影响到计算的速度和效率,同时还会影响计算误差,即计算的精度,甚至会直接影响到计算的成败. 不良的算法,会导致计算的彻底失败.

下面通过一个简单的例子来说明算法选取的重要性. 例如,计算下式的值:

$$a = \left(\frac{\sqrt{2}-1}{\sqrt{2}+1} \right)^3.$$

我们给出如下四种计算公式:

$$(1) a = (\sqrt{2} - 1)^6; \quad (2) a = 99 - 70\sqrt{2};$$

$$(3) a = \left(\frac{1}{\sqrt{2} + 1}\right)^6; \quad (4) a = \frac{1}{99 + 70\sqrt{2}}.$$

如果分别用近似值 $\sqrt{2} \approx 1.4 = \frac{7}{5}$ 和 $\sqrt{2} \approx 1.4166\cdots = \frac{17}{12}$, 按上述四种算法来计算 a 的值, 计算结果如表 1.0.1 所示.

表 1.0.1

算法	计算结果	
	$\sqrt{2} \approx \frac{7}{5}$	$\sqrt{2} \approx \frac{17}{12}$
$(\sqrt{2} - 1)^6$	$\left(\frac{2}{5}\right)^6 = 0.0040930$	$\left(\frac{5}{12}\right)^6 = 0.00523278$
$99 - 70\sqrt{2}$	1	$-\frac{1}{6} = -0.16666667$
$\left(\frac{1}{\sqrt{2} + 1}\right)^6$	$\left(\frac{5}{12}\right)^6 = 0.00523278$	$\left(\frac{12}{29}\right)^6 = 0.00501995$
$\frac{1}{99 + 70\sqrt{2}}$	$\frac{1}{197} = 0.00507614$	$\frac{12}{2378} = 0.00504626$

由表 1.0.1 可见, 按不同的计算公式和近似值所算出的结果五花八门, 各不相同, 有的甚至相差甚远, 还出现了负值. 近似值和算法的选定对计算结果的精度影响很大.

在研究算法的同时, 必须正确掌握误差的基本概念、误差在数值计算中的传播规律、误差分析和算法的稳定性等概念. 本章先对数值计算的误差作一些介绍.

1.1 误差的基本概念

1.1.1 误差的来源

用数值方法求解问题时, 计算过程中不可避免地存在误差, 其来源主要有:

(1) 模型误差. 所谓模型误差是指数学描述和实际问题之间存在的误差. 在用数学模型来描述实际问题时, 往往是通过抓住主要因素, 忽略一些次要因素, 而对问题作某些必要的简化建立起数学模型. 这样建立起来的数学模型实际上只是对实际问题的一种近似, 它与实际问题之间必存在一定的误差.

(2) 观测误差. 数值计算所需的一些原始数据, 一般是由观测或实验获得. 由于受到所用的观测仪器、设备精度的限制, 所得的数据都只能是近似的, 即存在着误

差. 这种误差, 称之为观测误差, 也称之为初值误差.

(3) 截断误差. 在不少数值计算中常常遇到超越计算, 如微分、积分和无穷级数求和等, 它们常需要用极限或无穷过程来实现. 而实际计算只能用有限次运算来完成, 这样就要对某种无穷过程进行“截断”, 即仅保留无穷过程的前有限部分而舍去它后面的无穷小部分, 这就带来了误差, 称之为截断误差. 例如, 求函数 $\sin x$ 和 $\ln(1+x)$ 时, 有表达式

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots, \quad (1.1.1)$$

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots. \quad (1.1.2)$$

当 $|x| < 1$ 时, 取两个函数的近似计算公式, 如取

$$\sin x \approx x - \frac{x^3}{3!} + \frac{x^5}{5!}, \quad (1.1.3)$$

$$\ln(1+x) \approx x - \frac{x^2}{2} + \frac{x^3}{3}. \quad (1.1.4)$$

由于只取前 3 项, 而将第 4 项和以后各项都舍去了, 自然产生了误差. 由数学分析知识, 易知当 $|x| < 1$ 时, 它们的截断误差可分别估计为

$$\left| \sin x - \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} \right) \right| \leq \frac{x^7}{7!}, \quad (1.1.5)$$

$$\left| \ln(1+x) - \left(x - \frac{x^2}{2} + \frac{x^3}{3} \right) \right| \leq \frac{x^4}{4}. \quad (1.1.6)$$

(4) 舍入误差. 在数值计算过程中经常遇到一些无穷小数, 如无理数和有理数中的某些分数化成的无限循环小数

$$\pi = 3.14159265 \dots,$$

$$\sqrt{2} = 1.41421356 \dots,$$

$$\frac{1}{3} = 0.33333333 \dots,$$

$$\frac{1}{3!} = 0.16666666 \dots,$$

.....

由于受计算机字长的限制, 计算机所能表示的数据只能是一定的有限位数, 这时就需要将数据按一定的舍入方式舍入成一定位数的近似有理数来代替. 由此而产生的误差, 称之为舍入误差.

数值计算过程中除了一些可以完全避免的过失误差之外, 还存在难以避免的上述四种误差. 所研究问题的数学模型一旦建立, 进入具体计算时所要考虑和分析的就是截断误差和舍入误差. 在数值分析中所涉及的误差, 通常指的就是舍入误差(含初始数据误差) 和截断误差. 讨论它们在计算过程中的传播和对计算结果的影响; 研究控制它们的影响以保证最终结果具有足够的精度; 既希望解决数值问题的算法简便、有效, 又要使最终结果准确、可靠.

1.1.2 绝对误差与相对误差

假设某个量的准确值(称为真值) x 的近似值为 x^* , 则 x 与 x^* 的差

$$\mathcal{E}(x) := x^* - x \quad (1.1.7)$$

称为近似值 x^* 的绝对误差, 简称为误差.

由于准确值一般是未知的或是无法知道的, 因而 $\mathcal{E}(x)$ 也是未知的. 但我们往往可以估计出绝对误差的上限, 即可以找到一个正数 η , 使

$$|\mathcal{E}(x)| \leq \eta. \quad (1.1.8)$$

满足(1.1.8)式的 η 称为 x^* 的绝对误差限. 有时也用

$$x = x^* \pm \eta \quad (1.1.9)$$

来表示(1.1.8)式. 这时(1.1.9)式右端的两数值 $x^* + \eta$ 和 $x^* - \eta$ 表示了 x 所在范围的上、下限. η 越小, 表示近似值 x^* 的精度越高.

绝对误差不足以刻画近似值的精确程度. 如在测量飞机机翼长度时, 发生1毫米的误差, 和测量飞机机翼的厚度时所发生的1毫米误差, 虽然两者的绝对误差是一样的, 但它们决定近似值的精确程度却是不一样的. 因而, 要评价一个近似值的精确程度, 除了要看它的绝对误差的大小之外, 还必须要考虑该数值本身的小, 这就需要引入相对误差的概念.

绝对误差与准确值的比, 即

$$\mathcal{E}_r(x) := \frac{\mathcal{E}(x)}{x} = \frac{x^* - x}{x} \quad (1.1.10)$$

称为 x^* 的相对误差. 由于准确值 x 往往是不知道的, 我们也常常将相对误差定义为

$$\mathcal{E}_r(x) := \frac{\mathcal{E}(x)}{x^*} = \frac{x^* - x}{x^*}. \quad (1.1.10')$$

(1.1.10)或(1.1.10')式反映了绝对误差和相对误差间的关系. 相对误差可以由绝对误差求得. 反之, 绝对误差也可由相对误差求出

$$\mathcal{E}(x) := x \cdot \mathcal{E}_r(x) \quad \text{或} \quad \mathcal{E}(x) := x^* \cdot \mathcal{E}_r(x). \quad (1.1.11)$$

同样, 相对误差也是无法准确求出的, 因为 (1.1.10) 式中的 $\mathcal{E}(x)$ 和 x 均无法准确求出. 但往往可以估计出相对误差的上限, 即可以找到一个正数 δ , 使

$$|\mathcal{E}_r(x)| \leq \delta. \quad (1.1.12)$$

满足 (1.1.12) 式的 δ 称为 x^* 的相对误差限.

绝对误差是一个有量纲的量, 而相对误差是没有量纲的.

1.1.3 算术运算的相对误差

现在, 我们来讨论数进行加、减、乘、除等运算时, 原始数据的相对误差和计算结果的相对误差之间的关系.

(1) 乘法和除法的情况. 设正数 x 的近似值为 $x + \Delta x$, 绝对误差 Δx 近似地等于 x 的微分, 即 $\Delta x \approx dx$. 相对误差为

$$\mathcal{E}_r(x) = \frac{\Delta x}{x} \approx \frac{dx}{x} = d \ln x, \quad (1.1.13)$$

即 x 的相对误差近似地等于 $\ln x$ 的微分. 由此可得乘和除的相对误差

$$\mathcal{E}_r(x_1 x_2) \approx d \ln(x_1 x_2) = d \ln x_1 + d \ln x_2 \approx \mathcal{E}_r(x_1) + \mathcal{E}_r(x_2), \quad (1.1.14)$$

$$\mathcal{E}_r\left(\frac{x_1}{x_2}\right) \approx d \ln\left(\frac{x_1}{x_2}\right) = d \ln x_1 - d \ln x_2 \approx \mathcal{E}_r(x_1) - \mathcal{E}_r(x_2), \quad (1.1.15)$$

即乘积的相对误差为各乘数的相对误差之和, 商的相对误差是被除数与除数的相对误差的差.

(2) 加法和减法的情况. 加法和减法的运算结果是数的代数和. 设 x_1 和 x_2 的近似值分别为 $x_1 + \Delta x_1$, $x_2 + \Delta x_2$, 则

$$\begin{aligned} \mathcal{E}_r(x_1 + x_2) &= \frac{\Delta(x_1 + x_2)}{x_1 + x_2} = \frac{\Delta x_1 + \Delta x_2}{x_1 + x_2} \\ &= \frac{x_1}{x_1 + x_2} \cdot \frac{\Delta x_1}{x_1} + \frac{x_2}{x_1 + x_2} \cdot \frac{\Delta x_2}{x_2} \\ &= \frac{x_1}{x_1 + x_2} \cdot \mathcal{E}_r(x_1) + \frac{x_2}{x_1 + x_2} \cdot \mathcal{E}_r(x_2). \end{aligned} \quad (1.1.16)$$

若 x_1 与 x_2 同号, 则上式右端 $\mathcal{E}_r(x_1)$ 和 $\mathcal{E}_r(x_2)$ 的系数

$$\frac{x_1}{x_1 + x_2}, \quad \frac{x_2}{x_1 + x_2} \quad (1.1.17)$$

都在 0 和 1 之间, 且它们的和等于 1. 这时, 由 (1.1.16) 式可得

$$\begin{aligned} |\mathcal{E}_r(x_1 + x_2)| &\leq \frac{x_1}{x_1 + x_2} \max \left\{ |\mathcal{E}_r(x_1)|, |\mathcal{E}_r(x_2)| \right\} \\ &\quad + \frac{x_2}{x_1 + x_2} \max \left\{ |\mathcal{E}_r(x_1)|, |\mathcal{E}_r(x_2)| \right\} \\ &= \max \left\{ |\mathcal{E}_r(x_1)|, |\mathcal{E}_r(x_2)| \right\}. \end{aligned} \quad (1.1.18)$$

所以, 当一些数的符号相同时, 它们和的相对误差限小于各数相对误差限中的最大者.

若 x_1 与 x_2 满足条件 $|x_1| \gg |x_2|$ (表示 $|x_1|$ 远大于 $|x_2|$), 则 (1.1.17) 式中前一个数近似地等于 1, 而后一个数的绝对值相当小, 此时有

$$\mathcal{E}_r(x_1 + x_2) \approx \mathcal{E}_r(x_1). \quad (1.1.19)$$

所以, 当两数的绝对值相差很大时, 此二数代数和的相对误差近似地等于绝对值较大者的相对误差.

若 x_1 与 x_2 异号, 则 (1.1.17) 式中两个数的绝对值至少有一个大于 1. 如果这时 x_1 和 $-x_2$ 相当接近, 则 (1.1.17) 式中两个数的绝对值都可能很大. 由 (1.1.16) 式可以看出, 这种情况下, 原始数据的误差会对计算结果产生相当大的影响.

1.1.4 有效数字

当准确值 x 有多位时, 常常按四舍五入的原则得到 x 的前几位的近似值 x^* , 例如

$$\pi = 3.14159265 \dots$$

取 3 位有效数字时, $x_3^* = 3.14$, $|\mathcal{E}_3(x)| \leq 0.002$;

取 5 位有效数字时, $x_5^* = 3.1416$, $|\mathcal{E}_5(x)| \leq 0.00005$.

它们的误差均不超过末位数字的半个单位, 即

$$|\pi - 3.14| \leq \frac{1}{2} \times 10^{-2}, \quad |\pi - 3.1416| \leq \frac{1}{2} \times 10^{-4}.$$

下面我们给出有效数字的概念.

定义 1.1.1 若近似值 x^* 的误差限不超过某一位的半个单位, 该位到 x^* 的第一位非零数字共有 n 位, 则称 x^* 有 n 位有效数字.

如取 $x^* = 3.14$ 作为 π 的近似值, x^* 就有 3 位有效数字; 如取 $x^* = 3.1416$ 作为 π 的近似值, x^* 就有 5 位有效数字.

x^* 有 n 位有效数字可写成如下的标准形式

$$\begin{aligned} x^* &= \pm 10^m \times (a_1 \times 10^{-1} + a_2 \times 10^{-2} + \dots + a_n \times 10^{-n}) \\ &= \pm 10^m \times 0.a_1 a_2 \dots a_n, \end{aligned} \quad (1.1.20)$$

其中, a_1 为 1 到 9 中的某一个数, 而 a_2, a_3, \dots, a_n 为 0 到 9 中的某一个数, m 为正整数. 我们有

$$|\mathcal{E}(x)| \equiv |x^* - x| \leq \frac{1}{2} \times 10^{m-n}. \quad (1.1.21)$$

此式表明, 当 m 相同的情况下, n 越大则 10^{m-n} 就越小, 故有效数位数越多, 绝对误差限越小.

由于

$$a_1 \times 10^{m-1} \leq |x^*| \leq (a_1 + 1) \times 10^{m-1},$$

故当 x^* 有 n 位有效数字时,

$$|\mathcal{E}_r(x)| = \frac{|x^* - x|}{|x^*|} \leq \frac{\frac{1}{2} \times 10^{m-n}}{a_1 \times 10^{m-1}} = \frac{1}{2a_1} \times 10^{1-n},$$

即

$$|\mathcal{E}_r(x)| \leq \frac{1}{2a_1} \times 10^{1-n}. \quad (1.1.22)$$

这也表明, 有效数位数越多, 相对误差限越小.

1.2 算法设计中应注意的问题

数值计算需要设计出好的算法. 衡量算法的标准一般有算法的稳定性, 运算的复杂性, 数值结果的精度等. 当这些要求不能同时满足时, 就应根据需要, 权衡利弊, 综合平衡而作抉择.

在数值计算中每一步都可能产生误差. 而解决一个问题往往要经过成千上万次运算, 我们不可能每步都加以分析, 只能从整体上考虑. 下面我们指出在数值计算中, 为控制误差的传播应注意的几个问题.

1. 简化计算步骤, 减少运算次数

减少运算次数, 不仅可以提高计算速度, 而且能减少误差的积累. 例如, 公式

$$\begin{aligned} ab + ac + ad &= a(b + c + d), \\ ax^3 + bx^2 + cx + d &= ((ax + b)x + c)x + d, \end{aligned}$$

左右两端的值相等, 但每一个公式左右两端的运算次数却是不同的. 前一个公式应用了分配律, 而后一个公式应用了著名的秦九韶算法. 当上述两公式左端的项数很大时, 差别更加显著. 又如, 计算 x^{255} 的值时, 如果逐个相乘则需要 254 次乘法, 但若采用下面的方法只需 14 次乘法即可

$$x^{255} = x \cdot x^2 \cdot x^4 \cdot x^8 \cdot x^{16} \cdot x^{32} \cdot x^{64} \cdot x^{128}.$$

2. 避免相接近的两数相减

两个相近的数相减, 由于它们前几位有效数字相同, 相减之后有效数字就减少了好几位, 从而使得相对误差增大. 这种现象称之为相减抵消. 例如, 考虑

$$x = 0.3721478693, \quad y = 0.3720230572,$$

则

$$x - y = 0.0001248121.$$

如果我们在五位十进制计算机上计算, 它们的近似值为

$$\begin{aligned} x^* &= 0.37215, \quad y^* = 0.37202, \\ x^* - y^* &= 0.00013. \end{aligned}$$

这样, $x^* - y^*$ 只有两位有效数字. 相应的相对误差满足

$$\left| \frac{(x^* - y^*) - (x - y)}{x - y} \right| = \left| \frac{0.00013 - 0.0001248121}{0.0001248121} \right| \approx 4\%.$$

这个相对误差是很大的.

如果遇到两个相接近的数相减时, 通常采用的方法是改变计算公式. 例如, 我们要计算 $1 - \cos x$. 如果 x 很小时, 即 $x \approx 0$, 这个计算导致相减抵消, 损失有效数字, 但我们有恒等式

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

如果我们按上式右边计算, 则可避免这种现象, 从而使误差减小. 类似地, 对于小的 ε , 我们可以把正弦的差化为

$$\sin(x + \varepsilon) - \sin x = 2 \cos\left(x + \frac{\varepsilon}{2}\right) \sin\frac{\varepsilon}{2}.$$

显然, 按上式右边计算, 误差就比较小, 结果就比较精确.

由于

$$\ln x_1 - \ln x_2 = \ln \frac{x_1}{x_2},$$

当 x_1 和 x_2 相接近时, 对上式左端进行计算导致相减抵消, 而用右端的公式来代替左端计算, 有效数字就不会损失. 另外, 当 $|\delta| \ll x$ 时,

$$\sqrt{x + \delta} - \sqrt{x} = \frac{\delta}{\sqrt{x + \delta} + \sqrt{x}},$$