

# 数值计算方法

唐旭清 主编

王林君 方伟 副主编

蔡日增 主审

NUMERICAL  
CALCULATION



科学出版社

# 数值计算方法

主编 唐旭清

副主编 王林君 方伟

主审 蔡日增



科学出版社

北京

## 内 容 简 介

本书参考国内外相关文献，结合教育部关于“数值计算方法”课程的基本要求，从基本概念、基本理论和方法系统介绍数值分析与计算的相关内容和观点。本书既注重理论的严谨性，又注重方法的实用性，重点阐明数值分析和各种算法构造的基本思想与原理。其主要内容包括：绪论、线性方程组的直接解法、解线性方程组的迭代法、矩阵的特征值和特征向量计算、插值法、曲线拟合、数值微分与数值积分、非线性方程和方程组的数值解法、常微分方程数值解法、瞬时扩散方程的差分解法简介和 Matlab 软件介绍等。全书重点突出，各篇章相互衔接，每章均附有应用实例与习题。

本书内容精炼，由浅入深，循序渐进，易于教学。适用于理工科大学硕士研究生及高年级本科生的“数值计算方法”课程的教学，也可供从事工程应用与计算的技术人员参考。

---

### 图书在版编目(CIP)数据

---

数值计算方法/唐旭清主编. —北京：科学出版社, 2015.6

ISBN 978-7-03-044616-9

I. ①数… II. ①唐… III. ①数值计算—计算方法 IV. ①O241

中国版本图书馆 CIP 数据核字(2015) 第 126074 号

---

责任编辑：周丹 曾佳佳 / 责任校对：郑金红

责任印制：赵博 / 封面设计：许瑞

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>

北京市文林印务有限公司 印刷

科学出版社发行 各地新华书店经销

\*

2015 年 6 月第 一 版 开本：787×1092 1/16

2015 年 6 月第一次印刷 印张：18 1/4

字数：433 000

定价：49.00 元

(如有印装质量问题，我社负责调换)

## 前　　言

随着计算机的广泛使用和科学技术的迅速发展,科学计算已经成为继理论分析和科学实验之后的第三种重要的科学研究方法。数值计算方法是介绍各类数学问题的近似求解的最基本、最常用的方法,它既具有数学各专业课程的抽象性和严谨性,又具有解决实际问题的实用性和实验性的技术特征,是理工科相关专业本科生和硕士生的一门重要专业基础课程。

本教材是在江南大学蔡日增教授编写的数值分析讲义基础上,参照教育部关于“数值计算方法”课程的基本要求为理工科各专业研究生及高年级本科生编写的。其基本内容包括数值代数、数值分析和微分方程数值解法等。同时,利用 Matlab 应用软件的数值计算和绘图的基本功能,进行各类计算方法的程序构造与实现。本书力求全面、系统地介绍求各类数学问题近似解的基本方法,重点阐明算法构造的基本思想与原理,图文并茂,突出教育部“重概念、重方法、重应用、重能力的培养”的精神,以提高学生的实践能力,加深对数值计算方法的理解。

本教材内容可供理工科相关专业研究生及理科部分专业高年级本科生教学选用或工程技术人员参考。根据编者及有关教师的教学实践,授完本书全部理论内容需 72 个学时。若略去部分理论推证和相对独立的章,如第 8 章、第 10 章等,亦可安排在 32~48 个学时讲授。

本书由唐旭清主编,王林君、方伟为副主编,蔡日增主审。参加编写工作的还有过榴晓、殷萍和程浩。其中王林君编写第 2 章、第 3 章和第 8 章,方伟编写第 4 章、第 5 章和第 6 章,殷萍编写第 9 章和第 10 章,程浩编写第 7 章,过榴晓编写附录和各章节中 Matlab 应用部分的内容。唐旭清进行了第 1 章编写及全书的统稿。本书在编写过程中,得到了江南大学理学院和江苏大学理学院部分教师和研究生的帮助,在此表示感谢。

本书得到了科技部国际合作项目(2011DFR70500)和江苏省优秀研究生课程教改项目(1145210232141730)的资助。特别是科学出版社的大力支持,在此深表感谢。

在本书编写过程中,作者虽然力求突出重点,内容系统而精炼,兼顾科学性和实用性,但因时间和水平有限,书中缺点和错误在所难免,敬请读者批评指正。

作　　者

2015 年 4 月

# 目 录

|                           |    |
|---------------------------|----|
| <b>第 1 章 绪论</b>           | 1  |
| 1.1 数值计算方法的任务与基本方法        | 1  |
| 1.2 误差及有关概念               | 2  |
| 1.2.1 误差的来源及分类            | 2  |
| 1.2.2 误差的描述               | 3  |
| 1.3 数值计算中的误差传播            | 6  |
| 1.3.1 基本运算中的误差估计          | 6  |
| 1.3.2 算法的数值稳定性            | 7  |
| 1.4 设计算法应注意的问题            | 9  |
| 1.4.1 避免两个相近的数相减          | 9  |
| 1.4.2 绝对值太小的数不宜作除数        | 10 |
| 1.4.3 避免大数“吃”小数的现象        | 10 |
| 1.4.4 简化计算步骤, 提高计算效率      | 11 |
| 本章小结                      | 11 |
| 习题                        | 12 |
| <b>第 2 章 线性方程组的直接解法</b>   | 13 |
| 2.1 引言                    | 13 |
| 2.2 Gauss 消去法及计算量         | 14 |
| 2.2.1 Gauss 消去法           | 14 |
| 2.2.2 Gauss 消去法的计算量       | 17 |
| 2.3 Gauss 主元素消去法          | 17 |
| 2.3.1 列主元素法               | 18 |
| 2.3.2 全主元素法               | 19 |
| 2.4 矩阵三角分解及其在解方程组中的应用     | 20 |
| 2.4.1 Gauss 消去过程的矩阵表示     | 20 |
| 2.4.2 矩阵的三角分解             | 22 |
| 2.4.3 线性方程组的直接三角分解法       | 25 |
| 2.4.4 解三对角方程组的追赶法         | 26 |
| 2.5 平方根法与改进的平方根法          | 29 |
| 2.5.1 平方根法 (Cholesky 分解法) | 30 |
| 2.5.2 改进的平方根法             | 31 |
| 2.6 矩阵、向量和连续函数的范数         | 33 |
| 2.6.1 范数的一般概念             | 33 |

---

|                                  |           |
|----------------------------------|-----------|
| 2.6.2 连续函数范数 .....               | 36        |
| 2.6.3 向量范数 .....                 | 36        |
| 2.6.4 矩阵范数 .....                 | 38        |
| 2.7 线性方程组的误差分析 .....             | 43        |
| 2.7.1 线性方程组的性态与条件数 .....         | 43        |
| 2.7.2 线性方程组解的误差估计 .....          | 46        |
| 2.8 应用实例 .....                   | 47        |
| 本章小结 .....                       | 51        |
| 习题 .....                         | 51        |
| <b>第 3 章 解线性方程组的迭代法 .....</b>    | <b>54</b> |
| 3.1 迭代法的基本概念 .....               | 54        |
| 3.1.1 迭代法的一般形式 .....             | 54        |
| 3.1.2 向量序列与矩阵序列的收敛性 .....        | 55        |
| 3.2 几种常用的单步定常线性迭代法 .....         | 56        |
| 3.2.1 Jacobi 迭代法 .....           | 56        |
| 3.2.2 Gauss-Seidel 迭代法 .....     | 59        |
| 3.2.3 超松弛 (SOR) 迭代法 .....        | 60        |
| 3.3 迭代法的收敛条件及误差分析 .....          | 62        |
| 3.3.1 迭代法的一般收敛条件 .....           | 62        |
| 3.3.2 几类特殊类型的迭代法收敛性判别 .....      | 64        |
| 3.3.3 简单迭代法的误差估计 .....           | 69        |
| 3.4 最速下降法与共轭梯度法 .....            | 70        |
| 3.4.1 最速下降法 .....                | 70        |
| 3.4.2 共轭梯度法 .....                | 71        |
| 3.5 应用实例 .....                   | 73        |
| 本章小结 .....                       | 75        |
| 习题 .....                         | 75        |
| <b>第 4 章 矩阵的特征值和特征向量计算 .....</b> | <b>78</b> |
| 4.1 幂法和反幂法 .....                 | 78        |
| 4.1.1 幂法 .....                   | 78        |
| 4.1.2 幂法的收敛加速 .....              | 82        |
| 4.1.3 反幂法 .....                  | 86        |
| 4.2 Jacobi 方法 .....              | 87        |
| 4.3 QR 方法 .....                  | 93        |
| 4.3.1 基本 QR 方法 .....             | 93        |
| 4.3.2 Householder 变换 .....       | 95        |
| 4.3.3 化一般矩阵为拟三角阵 .....           | 96        |
| 4.3.4 拟上三角矩阵的 QR 分解 .....        | 97        |

|                                     |            |
|-------------------------------------|------------|
| 4.3.5 带原点移位的 QR 方法——QR 加速收敛方法 ..... | 100        |
| 4.4 广义特征值问题的计算方法 .....              | 100        |
| 4.5 应用实例 .....                      | 101        |
| 本章小结 .....                          | 103        |
| 习题 .....                            | 103        |
| <b>第 5 章 插值法 .....</b>              | <b>105</b> |
| 5.1 多项式插值问题的一般描述 .....              | 106        |
| 5.1.1 多项式插值问题 .....                 | 106        |
| 5.1.2 插值多项式的误差估计 .....              | 106        |
| 5.2 几种常用插值多项式求法 .....               | 107        |
| 5.2.1 Lagrange 插值公式 .....           | 107        |
| 5.2.2 Newton 插值公式 .....             | 109        |
| 5.2.3 Hermite 插值 .....              | 117        |
| 5.3 分段低次插值 .....                    | 120        |
| 5.3.1 分段线性插值 .....                  | 121        |
| 5.3.2 分段三次 Hermite 插值 .....         | 123        |
| 5.3.3 三次样条 .....                    | 126        |
| 5.4 应用实例 .....                      | 131        |
| 本章小结 .....                          | 136        |
| 习题 .....                            | 137        |
| <b>第 6 章 曲线拟合 .....</b>             | <b>139</b> |
| 6.1 数据拟合的最小二乘法 .....                | 139        |
| 6.1.1 多项式拟合 .....                   | 140        |
| 6.1.2 可化为多项式拟合类型 .....              | 141        |
| 6.1.3 线性最小二乘法的一般形式 .....            | 143        |
| 6.2 正交多项式 .....                     | 147        |
| 6.2.1 正交多项式基本概念与性质 .....            | 147        |
| 6.2.2 正交多项式一般方法 .....               | 148        |
| 6.3 函数的最佳平方逼近 .....                 | 151        |
| 6.4 应用实例 .....                      | 156        |
| 本章小结 .....                          | 159        |
| 习题 .....                            | 159        |
| <b>第 7 章 数值微分与数值积分 .....</b>        | <b>161</b> |
| 7.1 Newton-Cotes 求积公式 .....         | 161        |
| 7.1.1 数值积分的基本思想 .....               | 161        |
| 7.1.2 Newton-Cotes 求积公式 .....       | 162        |
| 7.1.3 求积公式的误差估计 .....               | 165        |
| 7.2 复合求积公式 .....                    | 168        |

|                                      |            |
|--------------------------------------|------------|
| 7.2.1 复合梯形公式 .....                   | 168        |
| 7.2.2 复合 Simpson 公式 .....            | 169        |
| 7.2.3 复合 Cotes 公式 .....              | 170        |
| 7.2.4 复合求积公式的逐次分半算法 .....            | 173        |
| 7.3 Romberg 求积公式 .....               | 175        |
| 7.3.1 Richardson 外推法 .....           | 175        |
| 7.3.2 Romberg 求积公式 .....             | 176        |
| 7.4 Gauss 型求积公式 .....                | 178        |
| 7.4.1 Gauss 型求积公式的一般提法 .....         | 178        |
| 7.4.2 Gauss 点与正交多项式的关系 .....         | 180        |
| 7.4.3 Gauss 型求积公式的稳定性和收敛性 .....      | 182        |
| 7.4.4 常用 Gauss 型求积公式 .....           | 183        |
| 7.4.5 Gauss 型求积公式余项 .....            | 187        |
| 7.5 数值微分 .....                       | 188        |
| 7.5.1 插值型求导公式 .....                  | 188        |
| 7.5.2 外推法 .....                      | 190        |
| 7.5.3 用三次样条函数求数值导数 .....             | 191        |
| 7.6 应用实例 .....                       | 191        |
| 本章小结 .....                           | 193        |
| 习题 .....                             | 194        |
| <b>第 8 章 非线性方程和方程组的数值解法 .....</b>    | <b>197</b> |
| 8.1 引言 .....                         | 197        |
| 8.1.1 问题的背景 .....                    | 197        |
| 8.1.2 一元方程的搜索法 .....                 | 197        |
| 8.1.3 二分法 .....                      | 198        |
| 8.2 一元方程的基本迭代法 .....                 | 200        |
| 8.2.1 基本迭代法及其收敛性 .....               | 200        |
| 8.2.2 局部收敛性和收敛阶 .....                | 203        |
| 8.2.3 收敛性的改善 —— Steffensen 迭代法 ..... | 206        |
| 8.3 一元方程 Newton 迭代法 .....            | 207        |
| 8.3.1 Newton 迭代法及其收敛性 .....          | 207        |
| 8.3.2 重根时的 Newton 迭代改善 .....         | 210        |
| 8.3.3 离散 Newton 法 .....              | 211        |
| 8.4 非线性方程组的解法 .....                  | 212        |
| 8.4.1 不动点的迭代法 .....                  | 212        |
| 8.4.2 Newton 迭代法 .....               | 216        |
| 8.4.3 最速下降法 .....                    | 219        |
| 8.5 应用实例 .....                       | 220        |

|                                   |            |
|-----------------------------------|------------|
| 本章小结 .....                        | 221        |
| 习题 .....                          | 221        |
| <b>第 9 章 常微分方程数值解法 .....</b>      | <b>223</b> |
| 9.1 Euler 方法与改进的 Euler 方法 .....   | 224        |
| 9.1.1 Euler 方法 .....              | 224        |
| 9.1.2 Euler 方法的误差估计 .....         | 226        |
| 9.1.3 改进的 Euler 方法 .....          | 227        |
| 9.2 Runge-Kutta 法 .....           | 229        |
| 9.3 单步法的稳定性 .....                 | 233        |
| 9.3.1 相容性与收敛性 .....               | 233        |
| 9.3.2 稳定性 .....                   | 235        |
| 9.4 线性多步法 .....                   | 237        |
| 9.4.1 线性多步公式的导出 .....             | 237        |
| 9.4.2 常用的线性多步公式 .....             | 238        |
| 9.4.3 预测-校正系统 .....               | 241        |
| 9.5 一阶微分方程组与高阶方程的数值解法 .....       | 245        |
| 9.5.1 一阶微分方程组的数值解法 .....          | 245        |
| 9.5.2 高阶微分方程的数值解法 .....           | 246        |
| 9.5.3 差分方程解常微分方程边界问题 .....        | 247        |
| 9.6 应用实例 .....                    | 248        |
| 本章小结 .....                        | 251        |
| 习题 .....                          | 252        |
| <b>第 10 章 瞬时扩散方程的差分解法简介 .....</b> | <b>255</b> |
| 10.1 引言 .....                     | 255        |
| 10.2 差分格式建立 .....                 | 256        |
| 10.2.1 显式格式 .....                 | 256        |
| 10.2.2 隐式格式 .....                 | 256        |
| 10.2.3 Crank-Nicolson 格式 .....    | 257        |
| 10.3 局部截断误差与收敛性 .....             | 259        |
| 10.3.1 局部截断误差 .....               | 259        |
| 10.3.2 差分格式的收敛性 .....             | 260        |
| 10.4 应用实例 .....                   | 263        |
| 习题 .....                          | 266        |
| <b>参考文献 .....</b>                 | <b>267</b> |
| <b>附录 Matlab 软件简介 .....</b>       | <b>268</b> |

# 第1章 绪 论

## 1.1 数值计算方法的任务与基本方法

随着计算机科学与技术的不断发展及计算机应用的普及, 继实验方法和理论方法之后, 科学计算已成为科学实践的第三种重要手段。它主要在物理、力学、化学、生命科学、天文学、环境科学、经济科学及社会科学等领域中得到了广泛的应用, 成为不可缺少的重要工具。因此适用于计算机的数值计算方法已成为相关专业硕士研究生及本科生的必修课程。

利用计算机进行数值计算, 其实质上就是对具有一定数位的数值进行加、减、乘、除等算术运算以及一些逻辑运算。而研究怎样把各种数学问题的求解运算归结为对有限数位的四则计算, 以求得各种数学问题的数值解或近似数值解, 正是数值计算方法的根本课题。由四则运算及运算顺序的规定所构成的完整的解题步骤, 称之为算法。数值计算方法的根本任务就是研究算法, 即包括算法构成与算法分析。事实上, 数值计算方法是与数学问题(或数学模型)密不可分的, 而我们所研究的各种与数值相关的数学问题最终归结为各类数学模型, 表 1.1.1 给出的是数值计算方法按所描述的客观现象、数学模型、数学工具及所属的数值逻辑范畴。

表 1.1.1 描述客观现象、数学模型、数学工具和数值逻辑对应表

| 客观现象   | 数学模型      | 数学工具                 | 数值逻辑           |
|--------|-----------|----------------------|----------------|
| 确定性现象  | 白箱(或机械模型) | 经典数学                 | 经典逻辑(或 0-1 逻辑) |
| 随机性现象  | 黑箱(或随机模型) | 概率论与数理统计等            |                |
| 非确定性现象 |           |                      |                |
| 非随机性现象 | 灰箱        | Fuzzy 数学、灰色控制系统、可拓学等 | 非经典逻辑          |

关于描述非确定性随机现象的数学模型中的数值计算方法将由另外的课程去介绍, 如“正交实验与数据处理”等。本课程所涉及的内容是介绍描述确定性现象的数学模型中的数值计算方法, 因此本课程的基本任务就是要研究描述确定性现象的各种数学问题的算法, 包括算法的构成及算法分析等。

构造算法的原则就是要以计算机所能执行的运算为依据, 尽可能节省机器内存和运算工作量。在构造算法时, 常常采用近似替代, 而在数值计算方法中, 函数的近似替代称为函数逼近。在函数逼近中, 被逼近函数一般比较复杂, 或只知在若干点的值, 难以计算和分析。逼近的函数往往比较简单, 如多项式、有理函数、分段多项式等。利用函数的近似替代, 可以计算函数的积分、导数、极值及零点等。利用积分和导数的近似公式, 可以把微分方程或积分方程化为代数方程组。微分方程或积分方程的解本来是连续变量, 数值计算方法常常只算它在某点处的值, 这些值当然是连续解的离散的结果。把求连续变量问题转化为离散变量问题, 称为离散化。离散化后得到的代数方程组, 往往用来获取数值间的递推关系, 利用递推关系

编写计算机程序去求解. 这种使用递推公式求一系列的近似解, 并使它们越来越接近真实解的算法, 称为迭代法或逐次逼近法. 上述这一整套办法都是数值方法求解各种数学问题的基本方法.

所谓算法分析, 就是分析算法的理论依据、应用范围、收敛性、稳定性、误差估计及计算的空间和时间复杂度等.

本书将在《数学分析》、《空间解析几何》和《高等代数》(或《高等数学》和《线性代数》)的基础上, 不仅介绍求各类数学问题近似解的最基本、常用的数值方法, 而且注重阐明构造算法的基本思想和原理. 既注重介绍算法的构造和使用, 也注重算法的分析与研究.

## 1.2 误差及有关概念

### 1.2.1 误差的来源及分类

一个物理量和实际计算的值往往不同, 它们之差就称为这个物理量的误差. 产生误差的原因是多方面的, 因此一个物理量的误差具有多种来源.

首先, 数学模型是通过对实际问题进行抽象与简化(即忽略了一些次要因素)而得到, 其一般过程如图 1.2.1 所示. 因而即使数学问题能准确求解, 它与实际问题的解之间也有误差. 这种数学模型与实际问题之间出现的误差称为模型误差. 同时, 数学模型往往包含若干个由观测得到的参量, 如温度、时间、电压等. 这些观测得到的数据也有误差, 这种由观测产生的误差称为观测误差.



图 1.2.1 实际问题的近似求解过程图

其次, 根据实际问题建立起来的数学模型, 在许多情况下很难得到精确性, 需要选取适当的数值计算方法将其简化为较易求解的数值计算问题. 这种由数值计算方法产生的误差称为截断误差.

例如, 用  $e^x$  的幂级数展开式  $e^x = \sum_{i=0}^{\infty} \frac{x^i}{i!}$  来计算  $e^x$  的值时, 由于算法的有限性, 故只能截取其部分和  $p_n(x) = \sum_{i=0}^n \frac{x^i}{i!}$  来近似替代, 由此产生的误差为截断误差. 根据微积分可知 (Taylor 公式), 该截断误差为

$$R_n(x) = e^x - p_n(x) = \frac{e^\xi x^{n+1}}{(n+1)!}, \quad |\xi| < |x|$$

最后, 利用计算机为计算工具时, 由于计算机的字长有限, 只能用有限位进行取值和运算, 因此, 原始数据在计算机上表示时会产生误差, 而计算过程中又可能产生新的误差. 这种由计算机字长而产生的误差称为舍入误差.

例如, 有一台计算机只能表示 6 位十进制数, 圆周率  $\pi$  在计算机上表示为 3.14159, 从而产生误差  $R_1 = \pi - 3.14159 = 0.0000026\ldots$ . 又如 3.14159 与 9.21000 在计算机上进行加法运

算, 得到

$$s = 3.14159 + 9.21000 \approx 1.23516 \times 10$$

所产生的误差是:  $R_2 = 12.3516 - 12.35159 = 0.00001$ .

一般地, 将模型误差和观测误差统称为系统误差, 而将截断误差和舍入误差统称为方法误差. 本课程所涉及的误差分析仅指方法误差(即截断误差和舍入误差).

### 1.2.2 误差的描述

#### A. 绝对误差与绝对误差限

设  $x^*$  是准确值(或精确值)  $x$  的一个近似值, 则称  $e^* = x^* - x$  为近似值  $x^*$  的绝对误差, 简称误差.

一般而言,  $e^*$  的准确值很难求出或不能准确知道. 但可以根据测量工具或计算的具体情况估计出它的取值范围, 即存在某个正数  $\varepsilon^* > 0$ , 使得

$$|e^*| = |x^* - x| \leq \varepsilon^* \quad (1.2.1)$$

这个  $\varepsilon^*$  就称为近似值  $x^*$  的绝对误差限.

显然一个近似值的绝对误差限是不唯一的, 且若已知它的一个绝对误差限  $\varepsilon^*$ , 就可以知道准确值  $x$  的取值范围:

$$x^* - \varepsilon^* \leq x \leq x^* + \varepsilon^* \quad \text{或} \quad x = x^* \pm \varepsilon^* \quad (1.2.2)$$

对于同一个准确值  $x$  而言,  $e^*$  或  $\varepsilon^*$  越小, 近似值  $x^*$  就越精确, 但是对于不同的数  $x$  和  $y$  而言, 误差  $e^*$  或误差限  $\varepsilon^*$  的大小不能完全反映出近似值  $x^*$  和  $y^*$  中哪一个近似程度更好. 例如, 有两个测量值  $x = 15 \pm 2$  和  $y = 1000 \pm 5$ , 其中  $x$  和  $y$  的近似值分别为  $x^*=15$  与  $y^*=1000$ , 误差限分别是 2 和 5. 单从误差限来看: 前者小, 而后者大. 但是, 不能下这样的结论: 前者的测量精度高于后者, 这是为什么? 因为绝对误差或绝对误差限仅考虑了误差本身的大小, 没有考虑准确值的大小. 为了更好地反映近似值的精确程度, 引入相对误差的概念.

#### B. 相对误差与相对误差限

设  $x^*$  是准确值  $x$  的一个近似值,  $e^*$  是它的绝对误差, 则称:

$$e_r^* = \frac{e^*}{x} = \frac{x^* - x}{x} \quad (1.2.3)$$

为近似值  $x^*$  的相对误差.

在实际计算时, 由于真值  $x$  往往是不知道的, 因此相对误差  $e_r^*$  不能准确知道. 特别地, 当  $\left| \frac{e^*}{x^*} \right|$  较小时, 通常取

$$e_r^* = \frac{e^*}{x^*} = \frac{x^* - x}{x^*} \quad (1.2.4)$$

作为  $x^*$  的相对误差. 这是因为当  $\left| \frac{e^*}{x^*} \right|$  较小时,

$$\frac{e^*}{x} - \frac{e_r^*}{x^*} = \frac{e^*(x^* - x)}{xx^*} = \frac{(e^*)^2}{x^*(x^* - e^*)} = \frac{\left(\frac{e^*}{x^*}\right)^2}{1 - \frac{e^*}{x^*}} \approx (e_r^*)^2 \quad (1.2.5)$$

这是  $e_r^*$  的高阶无穷小量.

与绝对误差一样, 相对误差也只能估计其上限. 如果存在正数  $\varepsilon_r^*$ , 使得

$$|e_r^*| = \left| \frac{e^*}{x^*} \right| \leq \varepsilon_r^* \quad (1.2.6)$$

则称  $\varepsilon_r^*$  为近似值  $x^*$  的相对误差限.

如在上面提到的近似值  $x^* = 15$  和  $y^* = 1000$  的相对误差限分别为

$$\varepsilon_r^*(x^*) = \frac{2}{15} = 13.3\%, \quad \varepsilon_r^*(y^*) = \frac{5}{1000} = 0.5\%$$

由此可见,  $y^*$  近似  $y$  的程度比  $x^*$  近似  $x$  的程度要高.

### C. 精确位数与有效数字

有效数字是近似值的一种表示法, 它既能表示近似值的大小, 又能表示其精确程度. 一个数值的有效数字与精确位数密切相关. 当精确值  $x$  有无限多位数时, 常常按照“四舍五入”的原则取前  $n$  位数  $x^*$  作为  $x$  的近似值.

例如,  $x = 1.41421356237\cdots$ , 若取前五位数得  $x^* = 1.4142$ , 相应地误差为  $0.00001356237\cdots$ , 误差限为  $0.00005 = \frac{1}{2} \times 10^{-4}$ , 此时称  $x^*$  准确到小数点后的第 4 位, 并称由此算起的前五位数字 1.4142 为  $x^*$ , 下面给出有效数字的准确定义.

**定义 1.2.1** 如果近似值  $x^*$  的误差绝对值不超过某一位数字所在数位的半个单位, 且该数字到  $x^*$  的第一位非零数字共有  $n$  位, 则称用  $x^*$  近似  $x$  时具有  $n$  位有效数字, 简称  $x^*$  有  $n$  位有效数字.

例如, 因圆周率  $\pi$  分别满足:

$$|\pi - 3.1416| \leq \frac{1}{2} \times 10^{-4}, \quad |\pi - 3.14159| \leq \frac{1}{2} \times 10^{-5}$$

则  $\pi$  的近似值 3.1416 具有 5 位有效数字, 而 3.14159 具有 6 位有效数字.

在计算机中参加运算的数值通常进行标准化表示, 即将数值表示成如下形式:

$$x = \pm 0.a_1 a_2 \cdots a_n \cdots \times 10^m \quad (1.2.7)$$

其中  $m$  为整数,  $a_1, a_2, a_3, \dots$  为  $0, 1, \dots, 9$  中的数字, 且  $a_1 \neq 0$ . 此种表示法也称为科学计数法.

**例 1.2.1**  $x_1 = 0.0025 = 0.25 \times 10^{-2}$ ,  $x_2 = -387.8001 = -0.3878001 \times 10^3$ .

在式 (1.2.7) 中, 如果

$$x^* = \pm 0.a_1 a_2 \cdots a_n \times 10^m \quad (1.2.8)$$

是对  $x$  的第  $n+1$  位数字进行四舍五入后得到的近似值, 则  $x^*$  具有  $n$  位有效数字, 且其误差的绝对值不超过  $\frac{1}{2} \times 10^{m-n}$ , 即

$$|x^* - x| \leq \frac{1}{2} \times 10^{m-n} \quad (1.2.9)$$

其中  $\frac{1}{2} \times 10^{m-n}$  是  $x^*$  的一个特殊误差限. 它之所以为特殊误差限, 是因为它可以被用来判定近似值的有效数字的位数. 因此, 有效数字还可以作如下定义.

**定义 1.2.2** 如果  $x$  的近似值  $x^*$  满足不等式 (1.2.9), 则  $x^*$  具有  $n$  位有效数字. ■

如果例 1.2.1 中  $x_1$  和  $x_2$  都是四舍五入得到的有效数字, 那么  $x_1$  和  $x_2$  分别具有 2 位和 7 位有效数字. 根据式 (1.2.9), 它们的误差绝对值分别满足:

$$|x_1^* - x_1| \leq \frac{1}{2} \times 10^{-2-2} = \frac{1}{2} \times 10^{-4}$$

$$|x_2^* - x_2| \leq \frac{1}{2} \times 10^{3-7} = \frac{1}{2} \times 10^{-4}$$

从定义 1.2.1 可知: 有效数字的位数与小数点的位置无关. 因此, 精确小数点后  $n$  位, 不能反映它有效位数的多少, 只有经四舍五入得到的数字或式 (1.2.7) 的规格化形式后, 小数点后的有效位数才能反映出有效位数的多少.

一般地, 一个数位经四舍五入后得到的近似值, 每位有效数字都是唯一确定, 但必须注意有效数字定义中出现“等号”所带来的问题, 即此时会出现有效数字不唯一的特殊情况.

**例 1.2.2** 对准确值  $x = 3.95$  最后一位进行四舍五入后得到  $x_1^* = 4.0$ , 但若将最后一位 5 舍掉就得到  $x_2^* = 3.9$ , 它们的误差绝对值都不超过最后一位所在位的半个单位, 即

$$|x_1^* - x| \leq |4.0 - 3.95| \leq \frac{1}{2} \times 10^{-1}$$

$$|x_2^* - x| \leq |3.9 - 3.95| \leq \frac{1}{2} \times 10^{-1}$$

由定义 1.2.2 知:  $x_1^*, x_2^*$  都具有 2 位有效数字. ■

例 1.2.2 说明了近似数中的有效数字不一定都是通过四舍五入得到的, 从而, 通过对某个数进行四舍五入取近似值可得到它的有效数字, 但是并不是所有的有效数字都应通过四舍五入而得到.

**例 1.2.3** 设  $x = 1000$ , 它的两个近似值分别为:  $x_1^* = 999.9$  和  $x_2^* = 1000.1$ , 其误差绝对值均为  $|x_1^* - x| = |x_2^* - x| = 0.1$ , 但  $x_1^* = 999.9 = 0.9999 \times 10^3$ , 从而  $m = 3$ ; 而  $|x_1^* - x| = 0.1 \leq \frac{1}{2} \times 10^0 = \frac{1}{2} \times 10^{3-n}$ , 可得  $n = 3$ , 由定义 1.2.2 知:  $x_1^*$  具有 3 位有效数字. 同理可知,  $x_2^*$  具有 4 位有效数字. ■

例 1.2.3 说明某个数的近似值, 如果不是通过四舍五入得到, 那么它的数字并不都是有效数字, 同时它的数字位数并不等于该数的有效数字的位数. 下面定理给出了相对误差限与有效数字的关系.

**定理 1.2.1** 设  $x^*$  是  $x$  的近似值, 它的表达式为式 (1.2.8), 则  $x^*$  的有效数字与  $x^*$  的相对误差之间有如下关系:

(1) 若  $x^*$  具有  $n$  位有效数字时,  $x^*$  的相对误差

$$|e_r^*| \leq \frac{1}{2a_1} \times 10^{-n+1} \quad (1.2.10)$$

(2) 若  $x^*$  的相对误差限

$$\varepsilon_r^* \leq \frac{1}{2(a_1 + 1)} \times 10^{-n+1} \quad (1.2.11)$$

则  $x^*$  至少具有  $n$  位有效数字.

证明

(1) 由式 (1.2.9) 可得  $|e^*| = |x - x^*| \leq \frac{1}{2} \times 10^{m-n}$ , 从而有

$$|e_r^*| = \left| \frac{e^*}{x^*} \right| \leq \frac{\frac{1}{2} \times 10^{m-n}}{0.a_1 \cdots a_n \times 10^m} \leq \frac{1}{2a_1} \times 10^{-n+1}$$

即  $\frac{1}{2a_1} \times 10^{-n+1}$  是  $x^*$  的相对误差限.

(2) 若  $\varepsilon_r^* \leq \frac{1}{2(a_1 + 1)} \times 10^{-n+1}$ , 则由  $|e_r^*| = \left| \frac{e^*}{x^*} \right|$  可得

$$|e^*| = |x^* e_r^*| \leq 0.a_1 \cdots a_n \times 10^m \varepsilon_r^* \leq (a_1 + 1) \times 10^{m-1} \times \frac{1}{2(a_1 + 1)} \times 10^{-n+1} = \frac{1}{2} \times 10^{m-n}$$

由定义 1.2.2,  $x^*$  至少有  $n$  位有效数字. ■

这个定理表明, 近似值的有效数字越多 (即  $n$  越大), 相对误差 (限) 就越小; 反之, 相对误差 (限) 越小, 则式 (1.2.11) 右端项中的  $n$  就有可能越大, 有效数字位数就有可能越多.

在今后的数值问题中, 如果没有特别申明, 都可认为所有的原始数据均是有效数. 计算值具有有效数字的位数多少, 是数值方法的根本, 也是评定算法好坏的主要标准之一.

### 1.3 数值计算中的误差传播

#### 1.3.1 基本运算中的误差估计

本节中所讨论的基本运算是指四则运算与一些常用函数的计算.

由微积分知识可知: 当自变量的改变量 (误差) 很小时, 函数的微分作为函数的改变量的主要线性部分可近似表示函数的改变量. 因此, 利用微分运算公式可导出误差运算公式. 设数值计算中求得解  $y$  与参量  $x_1, \dots, x_n$  的函数关系为

$$y = f(x_1, x_2, \dots, x_n) \quad (1.3.1)$$

记  $(x_1, x_2, \dots, x_n)$  的近似值为  $(x_1^*, x_2^*, \dots, x_n^*)$ , 相应的解为

$$y^* = f(x_1^*, x_2^*, \dots, x_n^*) \quad (1.3.2)$$

假设  $f$  在点  $(x_1^*, x_2^*, \dots, x_n^*)$  处可微, 则当数据误差较小时, 解的绝对误差为

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

即有

$$e^*(y^*) \approx \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i} e^*(x_i^*) \quad (1.3.3)$$

其相对误差为

$$\begin{aligned} e_r^*(y^*) &= \frac{e^*(y^*)}{y^*} \approx d \ln f(x_1, x_2, \dots, x_n) |_{(x_1^*, x_2^*, \dots, x_n^*)} \\ &= \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i} \cdot \frac{e^*(x_i^*)}{y^*} \\ &= \sum_{i=1}^n \frac{\partial f(x_1^*, x_2^*, \dots, x_n^*)}{\partial x_i} \cdot \frac{x_i^*}{f(x_1^*, x_2^*, \dots, x_n^*)} \cdot e_r^*(x_i^*) \end{aligned} \quad (1.3.4)$$

特别地, 由式 (1.3.3) 和式 (1.3.4) 可得和、差、积和商的误差公式如下:

$$\begin{cases} e^*(x_1 \pm x_2) = e^*(x_1) \pm e^*(x_2) \\ e_r^*(x_1 \pm x_2) = \frac{x_1}{x_1 \pm x_2} e_r^*(x_1) \pm \frac{x_2}{x_1 \pm x_2} e_r^*(x_2) \end{cases} \quad (1.3.5)$$

$$\begin{cases} e^*(x_1 x_2) \approx x_2 e^*(x_1) + x_1 e^*(x_2) \\ e_r^*(x_1 x_2) \approx e_r^*(x_1) + e_r^*(x_2) \end{cases} \quad (1.3.6)$$

$$\begin{cases} e^*\left(\frac{x_1}{x_2}\right) \approx \frac{1}{x_2} e^*(x_1) - \frac{x_1}{x_2^2} e^*(x_2) \\ e_r^*\left(\frac{x_1}{x_2}\right) \approx e_r^*(x_1) - e_r^*(x_2) \end{cases} \quad (1.3.7)$$

式 (1.3.5) ~ 式 (1.3.7) 表明: 两个数值之和或差的误差分别为各自误差的和或差, 两个数值之积或商的相对误差分别为各自相对误差的和或差. 进一步有

$$|e^*(x_1 \pm x_2)| = |e^*(x_1) \pm e^*(x_2)| \leq |e^*(x_1)| + |e^*(x_2)| \quad (1.3.8)$$

$$|e_r^*(x_1 x_2)| \approx |e_r^*(x_1) + e_r^*(x_2)| \leq |e_r^*(x_1)| + |e_r^*(x_2)| \quad (1.3.9)$$

$$\left| e_r^*\left(\frac{x_1}{x_2}\right) \right| \approx |e_r^*(x_1) - e_r^*(x_2)| \leq |e_r^*(x_1)| + |e_r^*(x_2)| \quad (1.3.10)$$

因此, 两个数值之和或差的绝对误差限不超过各自的绝对误差限之和, 两个数值之积或商的相对误差限不超过各自的相对误差限之和.

**例 1.3.1** 设  $y = x^n$ , 求  $y$  的相对误差与  $x$  的相对误差之间的关系.

**解** 由式 (1.3.4) 知:  $e_r^*(y) \approx d \ln x^n = n d \ln x = n e_r^*(x)$ , 即  $x^n$  的相对误差是  $x$  的相对误差的  $n$  倍.

### 1.3.2 算法的数值稳定性

计算一个数学问题, 需要预先设计好由已知数据到计算问题结果的运算顺序, 即数学问题的算法. 对于给定的数学问题之算法, 需要判定此算法的好与不好, 即算法的数值稳定性. 为说明此问题, 先来介绍一个具体的数值计算问题.

例 1.3.2 计算积分值:  $I_n = \int_0^1 \frac{x^n}{x+5} dx (x = 0, 1, 2, \dots)$ .

解 由关系式

$$I_n + 5I_{n-1} = \int_0^1 \frac{x^n}{x+5} dx + \int_0^1 \frac{5x^{n-1}}{x+5} dx = \int_0^1 x^{n-1} dx = \frac{1}{n}, \quad n \geq 1$$

$$\text{且当 } n \geq 1 \text{ 时, } \frac{1}{6(n+1)} = \frac{1}{6} \int_0^1 x^n dx < I_n < \frac{1}{5} \int_0^1 x^n dx = \frac{1}{5(n+1)}.$$

于是可设计如下的两种算法:

$$\text{算法一} \quad \begin{cases} I_0 = \int_0^1 \frac{1}{x+5} dx = \ln 1.2 \approx 0.182321550 = I_0^* \\ I_n = \frac{1}{n} - 5I_{n-1} \end{cases} \quad (1.3.11)$$

取 8 位有效数字, 其中  $n = 1, 2, \dots$ , 于是可得  $I_1, I_2, \dots, I_n, \dots$

$$\text{算法二} \quad \begin{cases} I_n \approx \frac{1}{2} \left[ \frac{1}{5(n+1)} + \frac{1}{6(n+1)} \right] = \frac{11}{60(n+1)} \\ I_{k-1} = \frac{1}{5} \left( \frac{1}{k} - I_k \right) \end{cases} \quad (1.3.12)$$

其中  $k = n, n-1, \dots, 1$ , 于是可得  $I_{n-1}, I_{n-2}, \dots, I_1, I_0$ .

如果取  $n = 14$ , 即在算法二中  $I_{14} \approx \frac{11}{60 \times 15} \approx 0.012222222 = I_{14}^*$ , 按算法一和算法二计算所获得结果见表 1.3.1.

表 1.3.1 例 1.3.2 中算法一与算法二计算结果比较表

| $n$ | $I_n^* \text{ (按算法一)}$ | $I_n^* \text{ (按算法二)}$ |
|-----|------------------------|------------------------|
| 0   | 0.182321550            | 0.182321550            |
| 1   | 0.088392250            | 0.088392216            |
| 2   | 0.058038750            | 0.058038920            |
| 3   | 0.043139580            | 0.043138734            |
| 4   | 0.034302100            | 0.034306330            |
| 5   | 0.028489500            | 0.028468352            |
| 6   | 0.024219170            | 0.024324908            |
| 7   | 0.021761290            | 0.021232602            |
| 8   | 0.016993550            | 0.018836988            |
| 9   | 0.030143360            | 0.016926172            |
| 10  | -0.050716800           | 0.015369139            |
| 11  | 0.344493090            | 0.014063394            |
| 12  | -1.63914220            | 0.013016361            |
| 13  | 8.27263410             | 0.011841270            |
| 14  | -41.4346000            | 0.012222222            |

由表中结果可见: 按算法一可得到  $I_{10}^* = -0.050716800 < 0$ , 这显然是错误的. 这是因为当  $n \geq 0$  时,  $I_n > \frac{1}{6(n+1)} > 0$ ; 而按算法二计算, 尽管  $I_{14}^* = 0.012222222$  的精度不高, 其误