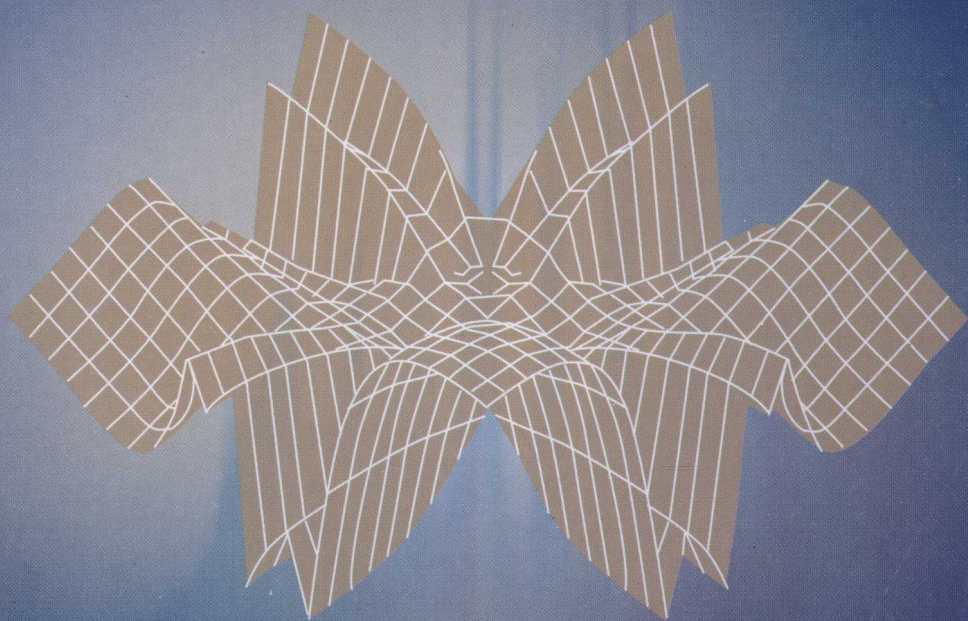


# 数值计算方法与应用

曾喆昭 黄创霞 周富照 编著



科学出版社

013026119

0241  
308

# 数值计算方法与应用

曾喆昭 黄创霞 周富照 编著

国家自然科学基金资助项目 (61040049)  
长沙理工大学学术专著出版基金资助项目



科学出版社

0241/308



北航

C1633007

013058118

## 内 容 简 介

本书详细介绍了科学计算领域中常用的数值计算方法, 主要内容包括插值与逼近、数值积分与数值微分、非线性方程及非线性方程组的数值计算方法、线性方程组的数值计算方法、常微分方程初值问题的数值计算方法等。本书不仅系统介绍了求解各类数学问题的最基本的数值计算方法和相关基础理论, 而且补充和新增了相应的优化计算方法。为了方便教学, 作者给出了相关实例的 MATLAB 源程序, 便于师生上机练习。本书的最大特色是以提出问题—分析问题—解决问题为主线, 先有问题背景后有解决问题的模型、算法和程序设计的教学和教材体系, 体系严密, 系统性强。除第 2 章外每章给出典型例子和一定数量的习题, 并在书后给出了习题解答。

本书可作为高等院校理工科专业本科生和研究生的教材, 也可作为相关科研人员的参考用书。

### 图书在版编目(CIP)数据

数值计算方法与应用/曾喆昭, 黄创霞, 周富照编著. —北京: 科学出版社, 2013

ISBN 978-7-03-036434-0

I. ①数… II. ①曾… ②黄… ③周… III. 数值计算—计算方法 IV. ①O241

中国版本图书馆 CIP 数据核字 (2013) 第 008342 号

责任编辑: 姚莉丽 张中兴 / 责任校对: 朱光兰  
责任印制: 阎磊 / 封面设计: 迷底书装

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

\*

2013 年 1 月第 一 版 开本: 720×1000 B5

2013 年 1 月第一次印刷 印张: 16 1/2

字数: 320 000

定价: 32.00 元

(如有印装质量问题, 我社负责调换)

# 前 言

编者以提出问题(给出问题背景)—分析问题(研究数学模型和算法及其 MATLAB 语言程序设计)—解决问题(将研究算法应用于实际问题并给出计算机仿真结果)为主线,先有问题背景后有解决问题的模型、算法和程序设计的教学和教材体系,体系严密,系统性强.在课程内容上,精简了部分陈旧、繁杂的定理证明内容,补充了解决各种应用问题的先进计算方法,增强了工程应用性,更有利于提高理工科学生分析问题和解决问题的能力.

课程内容紧扣实际应用,既保证理论内容的完整性和严密性,又不拘泥于烦琐和枯燥的理论推导,根据学生的特点,通过对问题背景的分析,向学生传授有关知识,符合学生的认知和学习规律,增强了教学效果.

本书遵循理论与应用相结合、经典算法与现代算法相结合、算法与计算机程序相结合的原则,以提出问题(应用背景)—分析问题(算法研究)—解决问题(程序设计)为主线,提出应用背景,针对问题研究各种可能的数值计算方法和相应的程序,并比较各种算法的计算效率(精度和速度),充分体现学以致用目的,做到易学、易教、易用.理论的价值体现在工程实践的具体应用中,“问题式”的教学方法可以激发学生分析问题、解决问题的能力,培养学生的学习兴趣,提高学习的积极性,便于引导学生找到明确的学习方向.通过本书的学习,学生不仅可以提高理论水平,而且可以提高分析问题和解决问题的能力.

全书共 7 章,内容包括误差分析、MATLAB 语言在数值计算中的应用、插值与逼近、数值积分与数值微分、非线性方程的数值解法、线性方程组的数值解法、常微分方程初值问题的数值解法等常用的数值计算方法及其有关的理论.每章均给出了典型例子和一定数量的习题,并在书后给出了习题解答.大部分例题给出了源程序,便于学生上机实习.

在本书的编写过程中,编者广泛参阅了国内外的相关教材和资料,在此谨向这些教材和资料的作者表示诚挚的感谢.

因编者水平有限,书中难免有疏漏和不妥之处,恳请读者批评指正.

编 者

2012 年 10 月



# 目 录

## 前言

<b>第 1 章 引论</b> .....	1
1.1 数值计算方法的对象、特点和意义 .....	1
1.2 误差分析 .....	2
1.3 数值计算中应注意的问题 .....	6
习题 1 .....	10
<b>第 2 章 MATLAB 在数值计算中的应用</b> .....	11
2.1 MATLAB 语言基础知识 .....	11
2.1.1 MATLAB 文件类型 .....	11
2.1.2 MATLAB 的矩阵、变量与表达式 .....	12
2.2 基本绘图方法 .....	13
2.2.1 直角坐标中的二维曲线 .....	13
2.3 MATLAB 基本运算 .....	16
2.3.1 关系运算 .....	17
2.3.2 逻辑运算 .....	17
2.3.3 特殊运算符 .....	17
2.3.4 矩阵运算 .....	18
2.4 MATLAB 控制语句 .....	18
2.5 自定义函数 .....	20
2.6 数值计算中的常用库函数 .....	20
2.6.1 向量与矩阵常用库函数 .....	20
2.6.2 插值函数 .....	21
2.6.3 多项式计算 .....	23
2.6.4 曲线拟合 .....	24
2.6.5 数值微分与差分 diff .....	25
2.6.6 数值积分函数 quad 和 quad8 .....	26
2.6.7 常微分方程求解函数 ode23 和 ode45 .....	26
2.6.8 非线性方程求解函数 .....	28
<b>第 3 章 插值与逼近</b> .....	30
3.1 问题背景: 人口增长问题 .....	30

3.2 拉格朗日插值 (Lagrange interpolation)·····	31
3.2.1 线性插值·····	31
3.2.2 抛物插值 (也称二次插值)·····	32
3.2.3 $n$ 次插值·····	33
3.2.4 插值余项·····	34
3.3 牛顿插值 (Newton interpolation)·····	37
3.3.1 具有继承性的插值公式·····	38
3.3.2 差商及其性质·····	39
3.3.3 差商形式的插值公式·····	40
3.3.4 差分形式的插值公式·····	41
3.4 埃尔米特插值 (Hermite interpolation)·····	43
3.4.1 二次插值·····	43
3.4.2 三次插值·····	44
3.4.3 $2n+1$ 次插值·····	45
3.4.4 Hermite 插值余项定理·····	47
3.5 三次样条插值·····	50
3.5.1 样条函数的概念·····	50
3.5.2 三次样条插值·····	51
3.5.3 三次样条插值函数的求法·····	52
3.6 曲线拟合的最小二乘法·····	59
3.6.1 直线拟合·····	59
3.6.2 多项式拟合·····	62
3.7 多项式曲线拟合的递归最小二乘法·····	64
习题 3·····	68
<b>第 4 章 数值积分与数值微分</b> ·····	<b>70</b>
4.1 问题背景: PID 调节器·····	70
4.1.1 PID 控制规律 (比例、积分、微分) 的基本形式·····	70
4.1.2 PID 控制规律的物理意义·····	71
4.2 机械求积·····	72
4.2.1 数值积分的基本思想·····	72
4.2.2 求积公式和它的代数精度·····	73
4.2.3 插值型的求积公式·····	74
4.3 牛顿-柯特斯 (Newton-Cotes) 求积公式·····	76
4.3.1 公式的推导·····	76
4.3.2 $n$ 低阶求积公式的代数精度·····	77

4.4	龙贝格 (Romberg) 算法	80
4.4.1	梯形法的递推公式	81
4.4.2	算法步骤	81
4.4.3	MATLAB 源程序	82
4.4.4	龙贝格算法	84
4.5	高斯 (Gauss) 求积算法	88
4.5.1	高精度的求积公式	88
4.5.2	高斯公式的基本特点	90
4.5.3	勒让德多项式	91
4.5.4	高斯求积公式的余项	95
4.5.5	高斯求积公式的稳定性与收敛性	96
4.6	数值积分的神经网络算法	97
4.6.1	余弦基函数神经网络模型	97
4.6.2	数值积分实例	98
4.7	数值微分	101
4.7.1	用插值多项式求数值微分	101
4.7.2	二阶数值微分公式	104
4.7.3	用三次样条函数求数值微分	104
	习题 4	105
<b>第 5 章</b>	<b>非线性方程的数值解法</b>	<b>110</b>
5.1	问题背景: 人口增长问题	110
5.2	二分法 (The Bisection Method)	112
5.2.1	二分法基本思想	112
5.2.2	二分法算法的源程序 (bisection.m)	114
5.2.3	总结	114
5.3	迭代法	115
5.3.1	迭代法的基本思路	116
5.3.2	线性迭代函数的启示	117
5.3.3	压缩映像原理	117
5.3.4	定点迭代法源程序 (fixedp.m)	118
5.3.5	迭代过程的收敛速度	120
5.4	迭代过程的加速收敛方法	121
5.4.1	迭代公式的加工	121
5.4.2	埃特金算法	123
5.4.3	埃特金加速算法的源程序 (aitken.m)	123

5.5	牛顿迭代法	124
5.5.1	牛顿迭代公式的导出	124
5.5.2	牛顿法的收敛性	125
5.5.3	牛顿迭代法源程序 (newtoniter.m)	126
5.5.4	牛顿下山法	127
5.6	弦截法	128
5.6.1	弦截法	128
5.6.2	弦截法的收敛性	129
5.7	求解非线性方程的神经网络算法	131
5.7.1	求解一元非线性方程的神经网络算法	132
5.7.2	神经网络算法收敛性研究	132
5.7.3	神经网络算法步骤	134
5.7.4	算例	134
5.7.5	算法改进	135
5.8	求解非线性方程组的神经网络算法	139
5.8.1	求解非线性方程组的神经网络模型	139
5.8.2	神经网络算法收敛性研究	141
5.8.3	神经网络算法步骤	142
5.8.4	数值试验	143
5.9	求解非线性方程的其他算法	148
5.10	求解非线性方程或代数方程重根的方法	149
5.10.1	算法描述	149
5.10.2	数值实例	150
	习题 5	152
<b>第 6 章</b>	<b>线性方程组的数值解法</b>	<b>154</b>
6.1	问题背景: 电阻网络	154
6.1.1	直接法	155
6.1.2	迭代法	155
6.2	高斯 (Gauss) 消元法	156
6.2.1	高斯消去法的计算过程	156
6.2.2	高斯消去法应注意的问题	157
6.3	三角分解法	158
6.3.1	矩阵 $A = [a_{ij}]_{n \times n}$ 的 Crout 分解	159
6.3.2	矩阵 $A = [a_{ij}]_{n \times n}$ 的 Cholesky 分解 ( $LL^T$ 分解)	163
6.3.3	解三对角线性方程组的三对角算法 (追赶法)	166



---

6.4	向量和矩阵的范数	167
6.4.1	向量的范数	167
6.4.2	向量范数的定义	167
6.4.3	矩阵的范数	169
6.4.4	谱半径、谱范数与方阵的 F-范数	170
6.4.5	方程组的状态与条件数	170
6.4.6	向量、矩阵的范数和条件数的计算	174
6.5	矩阵特征值和特征向量	175
6.5.1	雅可比 (Jacobi) 方法	176
6.5.2	QR 方法	177
6.5.3	计算矩阵特征值和特征向量的库函数	177
6.5.4	计算矩阵行列式值的库函数: $\det(\cdot)$	177
6.6	迭代法	178
6.6.1	雅可比 (Jacobi) 迭代法	178
6.6.2	赛德尔迭代法	180
6.6.3	关于 Jacobi 迭代法与 G-S 迭代法收敛性判据	181
6.6.4	逐次超松弛迭代法 (SOR 法)	182
6.7	共轭斜量 (梯度) 法	184
6.7.1	改善矩阵 $A$ 条件数的方法	185
6.7.2	条件预优共轭梯度算法	186
6.7.3	残差校正方法	188
6.8	基于梯度下降法的神经网络算法	189
6.8.1	基于梯度下降法 (Gradient-descent method) 的神经网络算法 (NN-GDM)	189
6.8.2	应用实例	193
6.9	基于递推最小二乘算法的神经网络计算方法 (NN-RLS)	199
	习题 6	204
<b>第 7 章</b>	<b>常微分方程的初值问题的数值解法</b>	<b>206</b>
7.1	问题背景: RLC 电路网络	206
7.2	欧拉方法	207
7.3	改进的欧拉方法	210
7.3.1	梯形公式	210
7.3.2	改进的欧拉公式	211
7.4	高阶泰勒方法 (Higher-order Taylor Methods)	213
7.5	龙格-库塔方法 (Runge-Kutta Methods)	216

7.5.1	龙格-库塔方法的设计思想	217
7.5.2	二阶龙格-库塔方法	217
7.5.3	三阶龙格-库塔方法	219
7.5.4	四阶龙格-库塔方法	219
7.6	亚当斯方法 (Adams Method)	221
7.6.1	亚当斯格式	221
7.6.2	亚当斯预报-校正系统	223
7.6.3	亚当斯预报-校正系统误差分析	224
7.7	收敛性与稳定性	226
7.7.1	收敛性问题	226
7.7.2	单步法的收敛性	226
7.7.3	单步法的稳定性问题	228
7.8	一阶常微分方程组和高阶微分方程求解	228
7.8.1	一阶方程组	228
7.8.2	高阶常微分方程的初值问题	229
7.9	高阶微分方程边值问题求解	231
7.10	求解常微分方程初值问题的神经网络算法	232
7.10.1	解微分方程初值问题的神经网络算法描述	232
7.10.2	解微分方程初值问题的神经网络算法步骤	239
7.10.3	仿真实例	240
	习题 7	248
	习题答案	250
	参考文献	254

# 第1章 引 论

## 1.1 数值计算方法的对象、特点和意义

### 1. 研究的对象

数值计算方法是近代数学的一个重要分支,是研究各种数学问题的数值方法(近似解法).

数值计算方法也称为计算方法或数值分析,是一门与计算机应用紧密结合的实用性很强的数学课程.

利用计算机解决科学计算的工作称为科学计算,或简称为计算,一般有以下几个过程(如图 1.1):

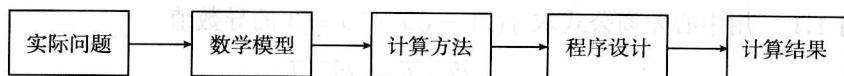


图 1.1

科学计算的应用范围十分广泛,国防尖端的一些科研项目,如核武器的研制、导弹的发射等,始终是科学计算最活跃的领域.今天,科学计算在工农生产的各个部门也正在发挥日益重要的作用.例如,气象资料的汇总、加工并求得天气图像,这方面的计算量大而且时间性很强,要求计算机作出高速或超高速运算,以对天气作出短期及中期的预报.

### 2. 主要特点

数值计算方法是数学的一个重要分支,着重研究求解的数值计算方法及相关的理论,包括方法的收敛性、稳定性及误差分析等.因此,数值计算既有纯数学的高度抽象性与严密科学性的特点,又有应用的广泛性与数值实验的高度技术性的特点.

### 3. 研究算法的意义

尽管计算机的速度越来越高,可以承担大运算量的工作,但是,这并不意味着人们对计算机上的算法可以随意选择.

众所周知,行列式的克拉默(Cramer)法则原则上可以用来求解线性方程组.用这种方法求解一个  $n$  阶的线性方程组,需要计算  $n+1$  个  $n$  阶的行列式的值,为此总共需要做  $n!(n-1)(n+1)$  次乘法.当  $n$  足够大时,其计算量是十分惊人的.比如一个 20 阶不算大的线性方程组,大约要做  $10^{21}$  次乘法,即使用百亿次每秒的计算

机去计算,也得要连续工作三千余年才能完成,这显然是毫无实际意义的.其实,解线性方程组有许多实用的算法(见第 6 章),比如用熟知的高斯消元法,一个 20 阶的线性方程组,即使用计算器也能计算出来.由此可知,研究高效的算法是科学计算成败的关键,在工程实践中具有十分重要的理论意义和实用价值.

本书介绍的各种问题的计算方法,大部分是国内外学者多年来研究成果的结晶,并补充了作者近年来的研究成果.希望读者通过本书的学习,掌握各种问题的科学计算方法,并开动脑筋,提出自己的一些能解决工程实际问题的科学计算方法,变被动学习为主动学习.

## 1.2 误差分析

### 1. 误差分析的重要性

许多数值计算方法给出的解答仅仅是所要求的解析解(精确解)的某种近似,因而研究数值计算方法,必须要注重误差分析、完成来源、误差传播情况以及对计算结果给出合理的误差估计,否则,一个合理的算法也可能得出错误的结果.

**例 1.1** 用中心差商公式求  $f(x) = \sqrt{x}$  在  $x = 3$  的导数值

$$f'(3) = \frac{\sqrt{3+h} - \sqrt{3-h}}{2h}$$

从理论上说,步长  $h$  越小,计算结果就越准确,但上机计算的实际情况将会怎样呢?下面来看看分析的结果.

**解** 众所周知,在计算机上表示的数是受机器字长限制的,假设取五位数字计算,当取  $h = 0.1$  时,得

$$f'(3) = \frac{\sqrt{3.1} - \sqrt{2.9}}{0.2} = \frac{1.7607 - 1.7029}{0.2} = 0.28900$$

与导数的精确值  $f'(3) = 0.28867513 \cdots$  相比,这项计算还是可取的,但是,如果缩小步长,且取  $h = 0.0001$  时,则有

$$f'(3) = \frac{\sqrt{3.0001} - \sqrt{2.9999}}{0.2} = \frac{1.7321 - 1.7320}{0.2} = 0.50000$$

算出的结果反而毫无价值.

由以上分析可知,两个相近的数相减,会造成有效数字的严重抵消,实际计算中要尽可能避免此类情况的发生.

**例 1.2** 求解方程  $x^2 - (10^5 + 1)x + 10^5 = 0$ .

**解** 仍然取五位数字进行计算,并用“ $\cong$ ”标记对阶舍入的计算过程,这里

$$b = -(10^5 + 1), \quad c = 10^5$$

由一元二次方程求根公式可知

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4c}}{2}$$

而  $-b = 10^5 + 1 \cong 10^5$ ,  $\sqrt{b^2 - 4c} = \sqrt{[-(10^5 + 1)]^2 - 4 \times 10^5} \cong 10^5$ , 所以有

$$x_1 = \frac{-b + \sqrt{b^2 - 4c}}{2} \cong \frac{10^5 + 10^5}{2} = 10^5$$

$$x_2 = \frac{-b - \sqrt{b^2 - 4c}}{2} \cong \frac{10^5 - 10^5}{2} = 0$$

与原方程的精确解  $x_1 = 10^5, x_2 = 1$  相比可知, 上面求出的结果出现了严重失真.

由以上分析可知, 加减运算引起了“大数”吃掉小数的后果, 使计算结果出现严重失真, 因此, 在实际计算时不宜采用相差悬殊的两个数做加减运算.

## 2. 误差的来源

误差的来源是多方面的, 但主要有以下几种:

### 1) 模型误差

用计算机解决科学计算问题首先要建立数学模型, 它是对被描述的实际问题进行抽象、简化而得到的, 因此总是近似的, 这就不可避免地产生误差, 通常把这种数学模型的解与实际问题的解之间出现的误差称为**模型误差**.

### 2) 观察误差

在数学模型中, 通常总包含一些观测数据, 如温度、长度、速度、电压等, 这些数据的值一般是由观测 (或实验) 得到的, 由于观测不可能绝对准确, 由此产生的误差称为**观测误差**.

### 3) 截断误差 (也称方法误差)

由实际问题建立起来的数学模型, 在很多情况下要得到准确解是很困难的. 当数学模型不能得到准确解时, 通常要用数值方法求它的近似解. 例如常把无限的计算过程用有限的计算过程来替代, 这种模型的准确解与数值方法的准确解之间的误差称为**截断误差**. 因为截断误差是方法固有的, 所以也常称为**方法误差**. 例如, 指数函数  $e^x$  可展开为幂级数形式

$$e^x = 1 + x + \frac{1}{2!}x^2 + \frac{1}{3!}x^3 + \cdots + \frac{1}{n!}x^n + \cdots$$

但使用计算机求值时, 我们不可能得出右端无穷多项的和, 只能截取有限项进行计算:

$$s_n(x) = 1 + x + \frac{1}{2!}x^2 + \frac{1}{3!}x^3 + \cdots + \frac{1}{n!}x^n$$

这样计算部分和  $s_n(x)$  作为  $e^x$  的值必然会有误差. 根据泰勒 (Taylor) 余项定理, 其截断误差为

$$R_n(x) = e^x - s_n(x) = \frac{x^{n+1}}{(n+1)!}e^\xi, \quad \xi \in (0, x)$$



则  $e^x = s_n(x) + R_n(x)$ . 若取  $e^x \approx s_n(x)$ , 则截断误差为  $R_n(x)$ , 也即用  $e^x$  的泰勒展开式的部分和  $s_n(x)$  来近似函数  $e^x$ , 其余项  $R_n(x)$  就是真值  $e^x$  的截断误差.

#### 4) 舍入误差

由于计算机的字长有限, 原始数据在计算机上表示会产生误差, 每一次运算又可能会产生新的误差, 这种误差称为舍入误差或计算误差.

以上几种误差都会影响计算结果的准确性. 数值计算中除了研究求解数学问题的数值方法外, 还要研究计算结果的误差是否满足精度要求, 这就是误差估计问题. 由于模型误差和观测误差往往并不是由计算工作者独立完成的, 数值计算中用不到描述自然现象, 也用不到观测测量, 因此, 我们的主要任务是研究截断误差和舍入误差对计算结果的影响.

重视误差分析并控制误差扩散是十分重要的, 没有误差分析的数值计算结果是不可信的. 误差就像矛盾一样无处不在, 无时不有. 在工程实际中, 也并不需要精确结果, 只要满足精度要求即可, 过分追求精确结果是不现实的, 就像导弹打飞机, 并不要求导弹非要命中飞机的发动机不可, 只要命中飞机的任何一部位均会击落飞机.

### 3. 绝对误差、相对误差和有效数字

#### 1) 绝对误差和相对误差

**定义 1.1** 设  $x$  为准确值,  $x^*$  为  $x$  的一个近似值, 则称  $e_a = x^* - x$  为近似值  $x^*$  的绝对误差, 并简称为误差.

需要特别指出的是, 通常无法得到准确值  $x$ , 因而不可能得到  $x^*$  的绝对误差  $e_a$  的真值, 只能根据测量的情况, 估计出误差绝对值的一个上限  $\varepsilon_a$ , 即

$$|e_a| = |x^* - x| \leq \varepsilon_a$$

这个正数  $\varepsilon_a$  通常称为近似值  $x^*$  的绝对误差限. 有了绝对误差限, 就可知道真值  $x$  的范围:

$$x^* - \varepsilon_a \leq x \leq x^* + \varepsilon_a$$

绝对误差的大小, 在许多情况下还不能完全刻画一个近似值的准确度. 例如测量 1000 米和 1 米两个长度, 若它们的绝对误差都是 1 厘米, 显然前者的测量比后者准确. 由此可知, 决定一个量的近似值的精确度, 除了考虑绝对误差的大小外, 还要考虑该量本身的大小, 为此引入了相对误差的概念.

**定义 1.2** 设  $x$  为准确值,  $x^*$  为  $x$  的一个近似值, 则称

$$e_r = \frac{e_a}{x} = \frac{x^* - x}{x} \quad (x \neq 0)$$

为近似值  $x^*$  的相对误差.

在实际计算中, 由于真值  $x$  一般是未知的, 但可以证明, 当  $e_r$  较小时,  $e_r$  中的分母  $x$  可用  $x^*$  来代替, 其两者之差是  $e_r$  的高阶无穷小, 因而可以忽略不计.

相对误差可正可负, 它的绝对值上限称为相对误差限, 即如果有正数  $\varepsilon_r$  使成立

$$|e_r| \leq \varepsilon_r$$

则称  $\varepsilon_r$  为  $x^*$  的相对误差限.

在误差分析中, 相对误差比绝对误差更重要.

## 2) 有效数字

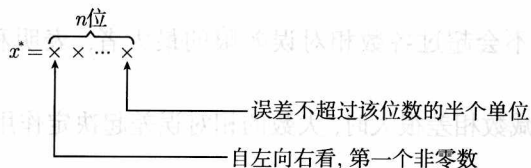
为了可以从近似数的有限位小数表示本身就能知道近似数的精度, 我们引入了有效数字的概念. 大家熟知, 当  $x$  有很多位数字时, 常按照“四舍五入”原则, 取  $x$  的前几位数字作为  $x$  的近似值  $x^*$ . 例如,  $\pi = 3.14159265 \dots$  可以表示为 3.14, 3.1416, 等, 这种表示方法的特点是, 近似数的误差限为其最末一位的半个单位.

若只取到小数后的四位数字, 则得

$$x^* = 3.1416$$

其误差为  $0.000007346 \dots$ , 误差限为  $0.00005 = \frac{1}{2} \times 10^{-4}$ , 此时称  $x^*$  准确到小数后第四位, 并称由此位算起的前五位数字 31416 为  $x^*$  的有效数字.

一般而言, 若近似值  $x^*$  的误差不超过某位数字的半个单位, 而从该位数字到  $x^*$  的最左边的那个非零数字共有  $n$  位, 那么这  $n$  位数字都称为  $x$  的有效数字, 即



并称近似值  $x^*$  具有  $n$  位有效数字.

又如

$$|e - 2.718| < 0.0005 = \frac{1}{2} \times 10^{-3}$$

则  $e$  的近似值 2.718 具有四位有效数字.

具有  $n$  位有效数字的近似值  $x^*$  也可以用指数形式写成

$$x^* = \pm a_1.a_2a_3 \dots a_n \times 10^m$$

其中  $a_1$  是 1 到 9 中的一个数字,  $a_2, \dots, a_n$  是 0 到 9 中的一个数字,  $m$  为整数, 且  $x^*$  的绝对误差为

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

则  $a_1, a_2, \dots, a_n$  是  $x^*$  的  $n$  位有效数字.

### 1.3 数值计算中应注意的问题

1.  $z = x \pm y$

1) 绝对误差

$$|e_a(z)| = |e_a(x \pm y)| = |e_a(x) \pm e_a(y)| \leq |e_a(x)| + |e_a(y)| \leq \varepsilon_x + \varepsilon_y \quad (1.3.1)$$

其中  $|e_a(x)| = |x - x^*| \leq \varepsilon_x$ ,  $|e_a(y)| = |y - y^*| \leq \varepsilon_y$ .

2) 和的相对误差 (假定  $x, y$  同号)

$$\begin{aligned} |e_r(z)| &= \left| \frac{z - z^*}{z} \right| = \left| \frac{(x + y) - (x^* + y^*)}{x + y} \right| \\ &= \left| \frac{x - x^*}{x + y} + \frac{y - y^*}{x + y} \right| = \left| \frac{x - x^*}{x} \times \frac{x}{x + y} + \frac{y - y^*}{y} \times \frac{y}{x + y} \right| \\ &\leq \left| \frac{x - x^*}{x} \right| \left| \frac{x}{x + y} \right| + \left| \frac{y - y^*}{y} \right| \left| \frac{y}{x + y} \right| \\ &\leq \max \left\{ \left| \frac{x - x^*}{x} \right|, \left| \frac{y - y^*}{y} \right| \right\} \times \left\{ \left| \frac{x}{x + y} \right|, \left| \frac{y}{x + y} \right| \right\} \\ &= \max \left\{ \left| \frac{x - x^*}{x} \right|, \left| \frac{y - y^*}{y} \right| \right\} \end{aligned} \quad (1.3.2)$$

即和的相对误差限不会超过各数相对误差限的最大者, 表明和的相对误差增长不快.

3) 当被减数和减数相差很大时, 大数的相对误差起决定作用, 这是因为

$$\begin{aligned} |e_r(z)| &= \left| \frac{z - z^*}{z} \right| = \left| \frac{(x - y) - (x^* - y^*)}{x - y} \right| = \left| \frac{x - x^*}{x - y} - \frac{y - y^*}{x - y} \right| \\ &\leq \left| \frac{x - x^*}{x} \right| \left| \frac{x}{x - y} \right| + \left| \frac{y - y^*}{y} \right| \left| \frac{y}{x - y} \right| \end{aligned} \quad (1.3.3)$$

当  $|x| \gg |y|$  时,  $\left| \frac{y}{x - y} \right| \rightarrow 0$ , 而  $\left| \frac{x}{x - y} \right| \rightarrow 1$ , 所以上式变为

$$|e_r(z)| \approx \left| \frac{x - x^*}{x} \right| = e_r(x)$$

4) 当两个相近数相减时, 由于  $\left| \frac{x}{x - y} \right|, \left| \frac{y}{x - y} \right|$  都将很大, 由 (1.3.3) 式可知, 其差的相对误差也很大. 在这种情况下, 由于此两数中前面的若干位数字相同, 相减后的结果中, 有效数字位大大减少. 例如,  $\cos 2^\circ = 0.9994$  具有四位有效数字, 但

$1 - \cos 2^\circ = 0.0006$  却只有一位有效数字. 为了避免这种情况出现, 常常改变计算公式, 比如

$$\sqrt{x-1} - \sqrt{x} = \frac{1}{\sqrt{x-1} + \sqrt{x}} \quad (\text{当 } x \text{ 很大时})$$

$$\sin(x + \varphi) - \sin x = 2 \cos\left(x + \frac{\varphi}{2}\right) \sin \frac{\varphi}{2} \quad (\text{当 } \varphi \text{ 很小时})$$

$$\ln x_1 - \ln x_2 = \ln \frac{x_1}{x_2} \quad (\text{单 } x_1 \text{ 与 } x_2 \text{ 接近时})$$

$$\arctan(x+1) - \arctan x = \arctan \frac{1}{1+x(x+1)} \quad (\text{当 } x \text{ 充分大时})$$

**例 1.3** 求有效数 3.150950, 15.426463, 568.3758, 7684.388 的和.

**解**  $3.150950 + 15.426463 + 568.3758 + 7684.388 = 8271.341213$  和数的绝对误差限为

$$2 \times (0.5 \times 10^{-6}) + 0.5 \times 10^{-4} + 0.5 \times 10^{-3} \approx 0.5 \times 10^{-3}$$

因此和值 8271.341213 应舍入为 8271.341.

2.  $z = xy$

1) 绝对误差限为

$$\begin{aligned} |e_a(z)| &= |z - z^*| = |\Delta z| = \left| \frac{\partial z}{\partial x} \Delta x + \frac{\partial z}{\partial y} \Delta y \right| = |y \Delta x + x \Delta y| \leq |y| |\Delta x| + |x| |\Delta y| \\ &\leq |y| \varepsilon_x + |x| \varepsilon_y = |y| |e_a(x)| + |x| |e_a(y)| \end{aligned}$$

即

$$|e_a(z)| \leq |y| |e_a(x)| + |x| |e_a(y)| \quad (1.3.4)$$

若取  $\varepsilon_x = \varepsilon_y = \varepsilon$ , 上式则为

$$|e_a(z)| \leq (|x| + |y|) \varepsilon \quad (1.3.5)$$

2) 相对误差限为

$$\begin{aligned} |e_r(z)| &= \left| \frac{e_a(z)}{z} \right| = \left| \frac{x e_a(y) + y e_a(x)}{xy} \right| = \left| \frac{e_a(y)}{y} + \frac{e_a(x)}{x} \right| \\ &= |e_r(x) + e_r(y)| \leq |e_r(x)| + |e_r(y)| \end{aligned}$$

即

$$|e_r(z)| \leq |e_r(x)| + |e_r(y)| \quad (1.3.6)$$

上式表明, 乘积的相对误差限等于各乘数相对误差限之和.

**例 1.4** 求  $z = 12.2 \times 73.56$  的绝对误差限和相对误差限.