

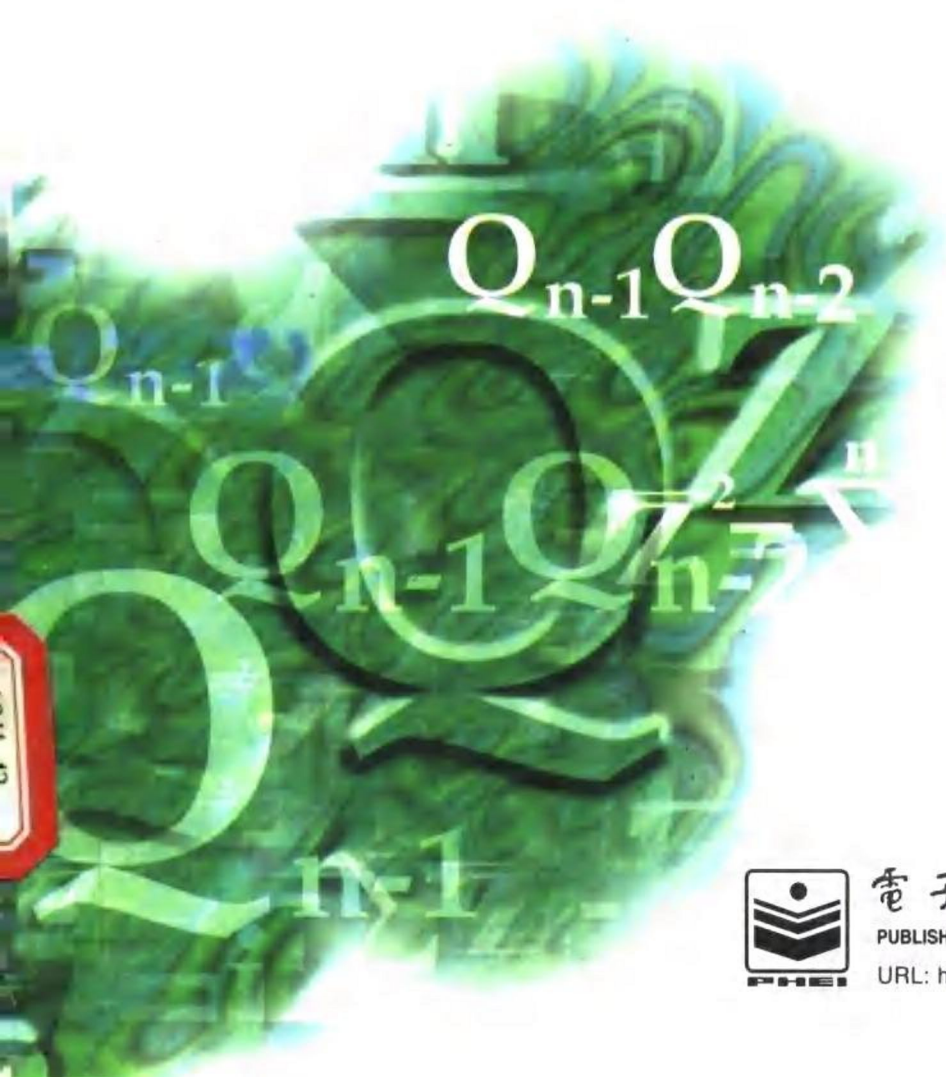


高等专科学校教材

中国计算机学会大专教育学会推荐出版

计算方法

吴筑筑 谭信民 邓秀勤 编



电子工业出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

URL: <http://www.phei.co.cn>

高等专科学校教材

741/207/25

计 算 方 法

吴筑筑 谭信民 邓秀勤 编

电子工业出版社

Publishing House of Electronics Industry

内 容 简 介

本书是根据计算机专业(专科)教学大纲并参照《计算机学科教学计划 1993》而编写的。着重介绍电子计算机上常用的计算方法。内容包括误差、一元非线性方程的解法、线性代数计算方法、插值法、曲线拟合、数值积分、常微分方程数值解法等方面的基础知识。

全书共分七章,内容力求精练,叙述由浅入深,文字通俗易懂,注重实用。各章常用的算法给出计算步骤和框图,并配有较多的例题和习题,部分主要算法给出了用C语言编写的参考程序,便于自学和上机应用。该书既适合作为专科院校开设60学时计算方法课程(包括上机实验10学时)的教材,也适合科技人员自学或参考。

丛 书 名: 高等专科学校教材

书 名: 计算方法

编 者: 吴筑筑 谭信民 邓秀勤

责任编辑: 张凤鹏

特约编辑: 侯维垣

印 刷 者: 北京牛山世兴印刷厂

装 订 者: 三河市路通装订厂

出版发行: 电子工业出版社出版、发行 URL: <http://www.phei.co.cn>

北京市海淀区万寿路173信箱 邮编 100036 发行部电话: 68214070

经 销: 各地新华书店经销

开 本: 787×1092 1/16 印张: 8 字数: 204.8千字

版 次: 1998年1月第一版 1998年1月第一次印刷

书 号: ISBN 7-5053-4265-7
G·340

定 价: 11.00元

凡购买电子工业出版社的图书,如有缺页、倒页、脱页者,本社发行部负责调换

版权所有·翻印必究

出版说明

根据国务院关于高等学校教材工作的有关规定,在电子工业部教材办的组织与指导下,按照教材建设适应“三个面向”的需要和贯彻国家教委关于“以全面提高教材质量水平为中心、保证重点教材,保持教材相对稳定,适当扩大教材品种,逐步完善教材配套”的精神,大专计算机专业教材编审委员会与中国计算机学会教育专业委员会大专教育学会密切合作,于1986~1995年先后完成了两轮大专计算机专业教材的编审与出版工作,共出版教材48种,从而较好地解决了全国高等学校大专层次计算机专业教材需求问题。

为及时使教材内容更适应计算机科学与技术飞速发展的需要以及在管理上适应国家实施“双休日”后的教学安排;在速度上适应市场经济发展形势的需要,在电子工业部教材办的指导下,大专计算机专业教材编委会、中国计算机学会大专教育学会与电子工业出版社密切合作,从1994年7月起经过两年的努力制定了1996~2000年大专计算机专业教材编审出版规划。

本书就是规划中配套教材之一。

这批书稿都是通过教学实践,从师生反映较好的讲义中经学校选报,编委会评选择优推荐或认真遴选主编人,进行约编的。广大编审者,编委和出版社编辑为确保教材质量和如期出版,作出了不懈的努力。

限于水平和经验,编审与出版工作中的缺点和不足在所难免,望使用学校和广大师生提出批评建议。

中国计算机学会教育委员会大专教育学会
电子工业出版社

附：先后参加全国大专计算机教材编审工作和参加全国大专计算机教育学会学术活动的学校名单：

上海科技高等专科学校
上海第二工业大学
上海科技大学
上海机械高等专科学校
上海化工高等专科学校
复旦大学
南京大学
上海交通大学
南京航空航天大学
扬州大学工学院
济南交通专科学校
山东大学
苏州市职工大学
国营 734 厂职工大学
南京动力高等专科学校
南京机械高等专科学校
南京金陵职业大学
南京建筑工程学院
长春大学
哈尔滨工业大学
南京理工大学
上海冶金高等专科学校
杭州电子工业学院
上海电视大学
吉林电气化专科学校
连云港化学矿业专科学校
电子工业部第 47 研究所职工大学
福建漳州大学
扬州工业专科学校
连云港职工大学
沈阳黄金学院
鞍钢职工工学院
天津商学院
国营 738 厂职工大学

北京广播电视大学
天津职业技术师范学院
天津市计算机研究所职工大学
山西大众机械厂职工大学
河北邯郸大学
沈阳机电专科学校
北京燕山职工大学
国营 761 厂职工大学
山西太原市太原大学
大连师范专科学校
江苏无锡江南大学
上海轻工专科学校
上海仪表职工大学
常州电子职工大学
国营 774 厂职工大学
西安电子科技大学
电子科技大学
河南新乡机械专科学校
河南洛阳大学
郑州粮食学院
江汉大学
武钢职工大学
湖北襄樊大学
郑州纺织机电专科学校
河北张家口大学
河南新乡纺织职工大学
河南新乡市平原大学
河南安阳大学
河南洛阳建材专科学校
开封大学
湖北宜昌职业大学
中南工业大学
国防科技大学
湖南大学

湖南计算机高等专科学校
中国保险管理干部学院
湖南税务高等专科学校
湖南二轻职工大学
湖南科技大学
湖南怀化师范专科学校
湘穗电脑学院
湖南纺织专科学校
湖南邵阳工业专科学校
湖南湘潭机电专科学校
湖南株洲大学
湖南岳阳大学
湖南商业专科学校
长沙大学
长沙基础大学

湖南零陵师范专科学校
湖北鄂州职业大学
湖北十堰大学
贵阳建筑大学
广东佛山大学
广东韶关大学
西北工业大学
北京理工大学
华中工学院汉口分院
烟台大学计算机系
安徽省安庆石油化工总厂职工大学
湖北沙市卫生职工医学院
化工部石家庄管理干部学院
西安市西北电业职工大学
湖南邵阳师范专科学校

前 言

本书是全国大专计算机应用专业 1996~2000 年出版系列教材之一,由全国大专计算机教材编审委员会负责征稿、审定、推荐出版的。

电子计算机的应用日益广泛,为工程技术和科学实验进行科学计算提供了强有力的工具。面对种类繁多的数值计算问题,如何选择合适的计算方法,并正确地在计算机上实施,以及如何估计结果的可靠程度,这都是科技工作者需要掌握的基本知识。

本教材系根据编者多年的教学经验并参考了多种相关教材和资料编写而成。由于高等专科学校教学时数较少,本教材充分考虑到专科学生的基础知识水平和教学基本要求,根据少而精、注重实际应用的原则,力求以较少的篇幅系统地介绍常见的、基本的数值计算方法及其基本原理。内容的叙述由浅入深,文字通俗简练,可读性强。对各种算法着重强调其构造思想及用法,同时也注意介绍算法的基本原理和方法误差,并尽量简化有关结论的理论推导。各章有较多的例题、习题和上机实验题目,可供选用。常用算法除给出了计算框图,也给出了一定数量的上机实验参考程序。考虑到 C 语言的应用已日益普及,故所有程序均用 Turbo C 语言编写并已上机运行通过,以便于读者学习和应用。

本书共分七章。第一章介绍误差的有关知识、计算机上进行数值运算的特点以及算法的稳定性等概念;第二章介绍一元非线性方程的几种基本解法,主要是迭代方法;第三章介绍线性代数方程组的直接解法,以顺序高斯消去法为基础,介绍了常用的选主元消去法、三角分解法、高斯-约当消去法以及求解三对角方程组的追赶法,并简要介绍了行列式的求法以及病态方程组的概念;第四章介绍线性代数计算中的矩阵迭代方法,其中着重介绍求解线性代数方程组的雅可比迭代法、高斯-赛德尔迭代法以及收敛条件,并介绍了计算矩阵特征值问题两种常用方法;第五章主要介绍代数多项式插值、样条插值和曲线拟合法;第六章为数值积分,介绍了牛顿-柯特斯求积公式的构造原理、常用的复合求积公式以及龙贝格求积公式,并简要介绍了变步长求积法;第七章主要介绍求解一阶常微分方程初值问题的几种数值方法,如欧拉法、预测-校正法、以及龙格-库塔法和阿达姆斯法,最后简要介绍求解二阶线性常微分方程边值问题的一种数值方法。全书参考学时为 60,其中包括 10 学时的上机实验。对于学时数不足 60 的专业可略去书中打“*”的选学内容,或根据需要自行选择。

本书由吴筑筑主编。第一、二章由邓秀勤编写,第三、四、五章由吴筑筑编写,第六、七章由谭信民编写。

本书由王文章教授主审。本教材在编写过程中,得到了李邦荣副教授、刘亚哲副教授、骆耀祖高级工程师、严庭栋讲师、谭建生讲师以及韶关大学计算机系实验室、韶关大学计算中心的热情支持和关心,在此表示诚挚的谢意。

由于我们的水平有限,书中错误和不妥之处在所难免,敬请读者批评指正。

编 者

1996 年 12 月

目 录

第一章 误差	(1)
第一节 浮点数及其运算特点.....	(1)
第二节 科学计算中误差的来源.....	(3)
第三节 误差的有关概念.....	(4)
一、绝对误差和绝对误差限.....	(4)
二、相对误差和相对误差限.....	(4)
三、有效数字.....	(5)
第四节 数值运算中误差的传播.....	(6)
一、利用微分估计误差.....	(6)
二、加减运算.....	(6)
三、乘除运算.....	(7)
第五节 算法的数值稳定性.....	(7)
一、算法的数值稳定性概念.....	(7)
二、设计算法的若干原则.....	(8)
习题一.....	(11)
第二章 一元非线性方程的解法	(12)
第一节 引言.....	(12)
第二节 二分法.....	(13)
第三节 迭代法的一般知识.....	(15)
一、迭代法的基本思想及几何意义.....	(15)
二、迭代法的收敛条件及误差估计式.....	(17)
* 三、迭代法的收敛阶概念.....	(20)
第四节 牛顿迭代法.....	(20)
第五节 弦截法(割线法).....	(22)
第六节 埃特金迭代法.....	(23)
第七节 上机实验参考程序.....	(24)
习题二.....	(26)
第三章 线性代数方程组的直接解法	(28)
第一节 顺序高斯消去法.....	(28)
一、顺序高斯消去法举例.....	(29)
二、一般情况的计算过程.....	(29)
第二节 选主元高斯消去法.....	(32)
一、列主元高斯消去法.....	(32)
二、全主元高斯消去法.....	(33)
第三节 高斯-约当消去法.....	(36)

第四节	解实三对角线性方程组的追赶法	(38)
第五节	三角分解法	(40)
一、	高斯消去法和矩阵的三角分解	(40)
二、	解方程组的三角分解法	(42)
三、	乔累斯基分解法	(44)
第六节	上机实验参考程序	(47)
习题三		(49)
第四章	线性方程组和矩阵特征值的迭代解法	(52)
第一节	线性代数方程组的迭代解法	(52)
一、	简单迭代法的一般形式	(52)
二、	雅可比迭代法	(53)
三、	高斯-赛德尔迭代法	(54)
第二节	迭代法的收敛性	(56)
一、	向量和矩阵的范数	(56)
二、	迭代法收敛的充分条件	(57)
第三节	矩阵特征值问题的计算方法	(59)
一、	雅可比方法	(59)
* 二、	QR 方法简介	(61)
第四节	上机实验参考程序	(63)
习题四		(64)
第五章	插值法和曲线拟合	(66)
第一节	插值法的基本理论	(66)
一、	插值问题及代数多项式插值	(66)
二、	插值多项式的误差	(67)
第二节	拉格朗日插值多项式	(68)
一、	线性插值和二次插值	(68)
二、	n 次拉格朗日插值	(70)
第三节	牛顿均差插值多项式	(71)
一、	均差及均差表	(71)
二、	牛顿均差插值多项式	(72)
第四节	差分及等距基点的牛顿插值公式	(74)
一、	差分及其性质	(74)
二、	牛顿前差和后差插值多项式	(74)
第五节	三次样条插值	(76)
一、	三次样条插值函数的定义	(76)
二、	三次样条插值函数的求法	(77)
第六节	曲线拟合的最小二乘法	(79)
一、	曲线拟合的最小二乘法	(79)
二、	超定方程组的最小二乘解	(80)
三、	代数多项式拟合	(81)
第七节	上机实验参考程序	(83)
习题五		(85)

第六章 数值积分	(87)
第一节 牛顿-柯特斯求积公式	(87)
一、牛顿-柯特斯求积公式	(87)
二、求积公式的代数精度	(89)
三、梯形公式和抛物线公式的误差估计	(90)
第二节 复合求积公式及其误差	(92)
一、复合梯形公式及其误差	(92)
二、复合抛物线公式及其误差	(92)
三、变步长的梯形公式和抛物线公式	(93)
第三节 龙贝格 (Romberg) 求积法	(95)
第四节 上机实验参考程序	(97)
习题六	(99)
第七章 常微分方程数值解法	(101)
第一节 引言	(101)
一、研究常微分方程数值解的必要性	(101)
二、建立数值方法的一些途径	(101)
第二节 欧拉法和改进的欧拉法	(103)
一、欧拉法及其截断误差	(103)
二、改进的欧拉法及预测-校正公式	(104)
第三节 龙格-库塔法	(106)
一、二阶的龙格-库塔公式	(106)
二、四阶的龙格-库塔公式	(107)
第四节 线性多步法	(108)
一、四阶阿达姆斯 (Adams) 外插公式	(108)
二、四阶阿达姆斯 (Adams) 内插公式	(109)
三、初始出发值的计算	(110)
四、阿达姆斯预测-校正公式	(110)
* 第五节 二阶线性常微分方程边值问题的数值解法	(111)
第六节 上机实验参考程序	(114)
习题七	(115)
参考文献	(116)

第一章 误差

科学实验方法、科学理论方法和科学计算方法是现代社会的三类科学方法。本书的目的是介绍一些常用的、基本的科学计算方法,也就是为科学和技术领域中常见的各种数学问题提供求解的数值计算方法。这些计算方法所给出的答案一般是所求真解的某些近似值,我们必须了解并估计近似值与真解的准确值之间的差异,即误差。由于科学计算的主要工具是数字电子计算机,因此,我们还应当了解在电子计算机上如何实施数值计算,它与严格的数学计算有什么区别。为此,本章简要介绍误差的基本理论以及算法的数值稳定性概念,作为学习以后各章的准备。

第一节 浮点数及其运算特点

在科学计算中常常把数,例如

$$0.003\ 120\ 7, 0.091\ 650, 293.704\ 8$$

等,分别表示成

$$0.312\ 07 \times 10^{-2}, 0.916\ 50 \times 10^{-1}, 0.293\ 704\ 8 \times 10^3$$

这样一来,一个数的数量级就一目了然了。在这种表示方法中,小数点的位置决定于后边那个10的指数。这种允许小数点位置浮动的表示方法,称为**数的浮点表示法**。用浮点表示法表示的数称为**浮点数**。一个浮点数由两部分组成。如上述各数中0.312 07, 0.916 50和0.293 704 8,称为浮点数的**尾数**,后边的 10^{-2} , 10^{-1} 和 10^3 ,称为**定位部**,是用来确定小数点的位置的。10称为**基底**, -2 , -1 和 3 称为**阶码**。基底一般是事先规定的。因此,一个无符号的浮点数由尾数和阶码两部分确定。

数的浮点表示方法是现代数字电子计算机通用的表示法,也是我们研究数值方法的基础。计算机中能够表示的浮点数个数是有限的,我们把计算机中浮点数的全体组成的集合记作 F ,称为**浮点数系**,则 F 中的浮点数具有以下形式:

$$\begin{aligned} x &= \pm 0.d_1d_2\cdots d_t \cdot \beta^p \\ &= \pm (d_1 \cdot \beta^{-1} + d_2 \cdot \beta^{-2} + \cdots + d_t \cdot \beta^{-t}) \cdot \beta^p \end{aligned} \quad (1-1)$$

其中 β 为浮点数的基底。若取十进制, $\beta=10$;若取二进制, $\beta=2$;若取十六进制, $\beta=16$,等等。 $0.d_1d_2\cdots d_t$ 是 β 进制小数,称为浮点数的尾数。 d_1, d_2, \cdots, d_t 为整数,满足

$$0 \leq d_i \leq \beta - 1, \quad i = 1, 2, \cdots, t$$

并规定当 $x \neq 0$ 时 $d_1 \neq 0$ 。自然数 t 为计算机的字长, p 为浮点表示的阶码,它有固定的下限 L 和上限 U ,即

$$L \leq p \leq U \quad (1-2)$$

t, L 和 U 随计算机而异。

这里当 $x \neq 0$ 时,规定 $d_1 \neq 0$ 是为了保证浮点数表示式的唯一性。这样的浮点数称为**规格**

化浮点数。

浮点数系 F 是一个离散的有限集合,在利用计算机进行计算时,初始数据和中间结果都可能不在 F 中,于是便发生用 F 中的数来近似地表示相应数据的问题。设实数 x 不属于 F ,计算机用 F 中最接近 x 的一个浮点数作为 x 的近似值,记这个浮点数为 $fl(x)$,它一般用“舍入”法来确定。

例如,设 $t=4, \beta=10$, 则

$$fl(0.20456 \times 10^{12}) = 0.2046 \times 10^{12}$$

$$fl(15.732) = 0.1573 \times 10^2$$

一般说来,若将非零实数 x 写成

$$x = \pm 0. a_1 a_2 \cdots a_t a_{t+1} \cdots \times 10^b, 0 \leq a_i \leq 9, a_1 > 0 \quad (1-3)$$

$$fl(x) = \begin{cases} \pm 0. a_1 a_2 \cdots a_t \times 10^b, & \text{若 } 0 \leq a_{t+1} \leq 4 \\ (\pm 0. a_1 a_2 \cdots a_t + 10^{-t}) \times 10^b, & \text{若 } a_{t+1} \geq 5 \end{cases} \quad (1-4)$$

其中 $fl(x)$ 与 x 符号相同。

也就是说,对于十进制数,一般还是采用“四舍五入”的方法。但对于二进制数,一般采用“零舍一入”的方法,其余类似。

数 0 在计算机中用尾数为 0 的浮点数表示。事实上,当一个浮点数的尾数为 0 时,不论阶码为何值,计算机都把该浮点数看成零值,又称“机器零”。

当实数 x 大于 F 中的最大数,或小于 F 中的最小非零数时, F 中找不到一个浮点数等于 $fl(x)$,这时计算机就不能继续进行下去,这种现象就是“溢出”(上溢或下溢)。

下面简要介绍一下计算机中浮点数的运算特点。设 x, y 都是规格化的浮点数,即 $x, y \in F$ 。它们的算术运算的精确结果不一定是 F 中的浮点数,计算机自动把运算结果用 F 中的规格化浮点数表示出来,我们称这个过程为“规格化”。此外,当两个数量级不同的数相加减时需要“对阶”,将阶码统一为较大者,然后才能将尾数相加减。例如,设

$$t=4, \beta=10, x=0.3127 \times 10^{-6}, y=0.4153 \times 10^{-4}$$

$$\begin{aligned} \text{则 } x+y &\approx 0.0031 \times 10^{-4} + 0.4153 \times 10^{-4} && \text{(对阶)} \\ &= 0.4184 \times 10^{-4} && \text{(规格化)} \end{aligned}$$

而 $x+y$ 的精确结果是 0.418427×10^{-4} ,它不在 F 中。

又如在四位十进制计算机上计算

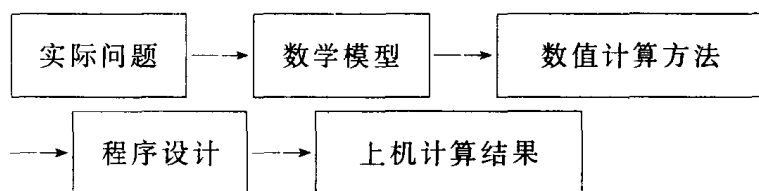
$$\begin{aligned} &0.8961 \times 10^3 + 0.4688 \times 10^{-5} \\ &\approx 0.8961 \times 10^3 + 0.0000 \times 10^3 && \text{(对阶)} \\ &= 0.8961 \times 10^3 && \text{(规格化)} \end{aligned}$$

其结果大数“吃掉”了小数。

在计算机中进行浮点数的运算时,通常实数加法的结合律、乘法对加法的分配律也不成立。由于浮点运算具有种种与实数运算不同的特点,每一步运算都可能产生误差,因此在设计算法和实际计算中必须注意对误差的估计。

第二节 科学计算中误差的来源

用计算机解决科学计算问题通常经历以下过程：



因此误差的来源主要有以下四类。

(一) 模型误差

在将实际问题转化为数学模型的过程中,为了使数学模型尽量简单,以便于分析或计算,往往要忽略一些次要的因素,进行合理的简化。这样,实际问题与数学模型之间就产生了误差,这种误差称为**模型误差**。由于这类误差难于作定量分析,所以在计算方法中,总是假定所研究的数学模型是合理的,对模型误差不作深入的讨论。

(二) 观测误差

在数学模型中,一般都含有从观测(或实验)得到的数据,如温度、时间、速度、距离、电流、电压等等。但由于仪器本身的精度有限或某些偶然的客观因素,会引入一定的误差,这类误差叫做**观测误差**。通常根据测量工具或仪器本身的精度,可以知道这类误差的上限值,所以无须在计算方法中作过多的研究。

(三) 截断误差(方法误差)

当数学模型得不到精确解时,要用数值计算方法求它的近似解,由此产生的误差称为**截断误差**或**方法误差**。譬如在数值计算中,常用收敛的无穷级数的前几项来代替无穷级数进行计算,即抛弃了无穷级数的后段,这样就产生了截断误差。例如

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots, \quad -\infty < x < +\infty$$

当 $|x|$ 很小时,常用 x 代替 $\sin x$,其截断误差大约为 $\frac{1}{6}x^3$ 。截断误差的大小,直接影响数值计算的精度,所以它是数值计算中必须十分重视的一类误差。

(四) 舍入误差

如上节所述,由于计算机字长有限,原始数据的输入及浮点运算过程中都可能产生误差。而事实上,无论用电子计算机、计算器计算还是笔算,都只能用有限位小数来代替无穷小数或用位数较少的小数来代替位数较多的有限小数,这样产生的误差叫做**舍入误差**。在数值计算中,往往要进行成千上万次四则运算,因而就会有成千上万个舍入误差产生,这些误差一经叠加或传递,对精度可能有较大的影响。所以,作数值计算时,对舍入误差应予以足够的重视。

显然,上述四类误差都会影响计算结果的准确性,但模型误差和观测误差往往需要会同各有关学科的科学工作者共同研究,因此在计算方法课程中,主要研究截断误差和舍入误差(包括初始数据的误差)对计算结果的影响。

第三节 误差的有关概念

一、绝对误差和绝对误差限

定义 1 假设某一量的准确值为 x , 近似值为 x^* , 则 x 与 x^* 之差叫做近似值 x^* 的**绝对误差**(简称**误差**), 记为 $\epsilon(x)$, 即

$$\epsilon(x) = x - x^* \quad (1-5)$$

$|\epsilon(x)|$ 的大小标志着 x^* 的精确度。一般地, 在同一量的不同近似值中, $|\epsilon(x)|$ 越小, x^* 的精确度越高。当 $|\epsilon(x)|$ 较小时, 由微分和增量的关系知 x^* 的绝对误差 $\epsilon(x) \approx dx$, 故我们可以利用微分估计误差。

由于准确值 x 一般不能得到, 于是误差 $\epsilon(x)$ 的准确值也无法求得, 但在实际测量计算时, 可根据具体情况估计出它的大小范围。也就是指定一个适当小的正数 ξ , 使

$$|\epsilon(x)| = |x - x^*| \leq \xi \quad (1-6)$$

我们称 ξ 为近似值 x^* 的**绝对误差限**。有时也用

$$x = x^* \pm \xi \quad (1-7)$$

表示近似值的精度或准确值的所在范围。在实际问题中, 绝对误差一般是有量纲的。例如, 测得某一物体的长度为 5m, 其误差限为 0.01m, 通常将准确长度 s 记为

$$s = 5 \pm 0.01$$

即准确值在 5m 左右, 但不超过 0.01m 的误差限。

二、相对误差和相对误差限

绝对误差的大小并不能确定近似程度的好坏。例如, 有两个温度计, 其一测量 1000°C 时的绝对误差限为 5°C , 而另一测量 100°C 时的绝对误差限为 1°C , 虽然后者绝对误差限的数值较小, 但第一种温度计更为精确。可见, 决定一个量的近似值的精确度除了要看绝对误差的大小外, 还要考虑到该量本身的大小。据此, 我们引进相对误差的概念。

定义 2 我们把绝对误差与准确值之比

$$\epsilon_r(x) = \frac{\epsilon(x)}{x} = \frac{x - x^*}{x}, \quad x \neq 0 \quad (1-8)$$

称为 x^* 的**相对误差**。

由于准确值 x 往往是不知道的, 因此在实际问题中, 当 $|\epsilon_r(x)|$ 较小时, 常取

$$\epsilon_r(x) = \frac{\epsilon(x)}{x^*}$$

一般地, 在同一量或不同量的几个近似值中, $|\epsilon_r(x)|$ 小者精确度高。相对误差是一个无量纲量。

在实际计算中, 由于 $\epsilon(x)$ 与 x 都不能准确地求得, 因此相对误差 $\epsilon_r(x)$ 也不可能准确地得到, 我们只能估计它的大小范围。即指定一个适当小的正数 η , 使

$$|\epsilon_r(x)| = \frac{|\epsilon(x)|}{|x|} \leq \eta \quad (1-9)$$

称 η 为近似值 x^* 的**相对误差限**。当 $|\epsilon_r(x)|$ 较小时, 可以用下式来计算 η :

$$\eta = \frac{\xi}{|x^*|} \quad (1-10)$$

显然,上例第一种温度计的相对误差限为 $5/1\,000$;而第二种的相对误差限 $1/100$,它是前者的两倍。

由式(1-4)确定的 $fl(x)$ 作为 x 的浮点数近似值,相对误差满足不等式

$$\left| \frac{fl(x) - x}{x} \right| \leq \frac{5 \times 10^{-(t+1)} \times 10^b}{0.1 \times 10^b} = 5 \times 10^{-t}$$

故 $fl(x)$ 的相对误差限为 5×10^{-t} ,称之为计算机的精度。

三、有效数字

我们知道,当 x 有很多位数字时,常常按照“四舍五入”原则取前几位数字作为 x 的近似值 x^* 。

例 1 设 $x = \pi = 3.141\,592\,6\dots$

取 $x_1^* = 3$ 作为 π 的近似值,则 $|\epsilon_1(x)| = 0.141\,5\dots \leq \frac{1}{2} \times 10^0$;

取 $x_2^* = 3.14$,则 $|\epsilon_2(x)| = 0.001\,59\dots \leq \frac{1}{2} \times 10^{-2}$;

取 $x_3^* = 3.141\,6$,则 $|\epsilon_3(x)| = 0.000\,007\,34\dots \leq \frac{1}{2} \times 10^{-4}$ 。

它们的误差都不超过末位数字的半个单位。

定义 3 若近似值 x^* 的绝对误差限是某一位上的半个单位,该位到 x^* 的第一位非零数字一共有 n 位,则称近似值 x^* 有 n 位有效数字,或说 x^* 精确到该位。

准确数本身有无穷多位有效数字,即从第一位非零数字以后的所有数字都是有效数字。如例 1 中的 x_1^*, x_2^*, x_3^* , 分别有 1, 3, 5 位有效数字。

实际上,用四舍五入法取准确值 x 的前 n 位(不包括第一位非零数字前面的零)作为它的近似值 x^* 时, x^* 有 n 位有效数字。

例 2 设 $x = 4.269\,72$, 则按四舍五入法,取 2 位, $x_1^* = 4.3$, 有效数字为 2 位;取 3 位, $x_2^* = 4.27$, 有效数字为 3 位;取 4 位, $x_3^* = 4.270$, 有效数字为 4 位。

值得注意的是,近似值后面的零不能随便省去。如例 2 中 4.27 和 4.270,前者精确到 0.01, 其有 3 位有效数字;而后者精确到 0.001, 其有 4 位有效数字。可见,它们的近似程度完全不同。

定义 3 换一种说法就是:设 x 的近似值 x^* 表示成

$$x^* = \pm 0. a_1 a_2 \dots a_n \dots \times 10^p \quad (1-11)$$

若其绝对误差限

$$|\epsilon(x)| = |x - x^*| \leq \frac{1}{2} \times 10^{p-n} \quad (1-12)$$

则称近似数 x^* 具有 n 位有效数字。这里 p 为整数, a_1, a_2, \dots, a_n 是 0 到 9 中的一个数字且 $a_1 \neq 0$ 。

例如,若 $x^* = 0.231\,56 \times 10^{-2}$ 是 x 的具有五位有效数字的近似值,则绝对误差限是

$$|x - x^*| \leq \frac{1}{2} \times 10^{-2-5} = \frac{1}{2} \times 10^{-7}$$

定义 3 或(1-12)式建立了绝对误差限和有效数字之间的关系。由于 n 越大, 10^{p-n} 的值越小,所以有效数字位越多,则绝对误差限越小。

下面给出有效数字与相对误差的关系。

定理 1 若近似数 x^* 具有 n 位有效数字, 则其相对误差为

$$|\epsilon_r(x)| \leq \frac{1}{2\alpha_1} \times 10^{-(n-1)} \quad (1-13)$$

其中 $\alpha_1 \neq 0$ 是 x^* 的第一位有效数字。

定理 1 说明有效数字位越多, 相对误差限越小。

定理 2 形式如(1-11)的近似数 x^* , 若其相对误差满足

$$|\epsilon_r(x)| \leq \frac{1}{2(\alpha_1+1)} \times 10^{-(n-1)} \quad (1-14)$$

则 x^* 至少有 n 位有效数字。

由此可知, 有效数字位数可刻画近似数的精确度, 相对误差限与有效数字的位数有关。

第四节 数值运算中误差的传播

一、利用微分估计误差

数值运算的误差估计情况较复杂, 通常可利用微分估计误差。设数学问题的解 y 与变量 x_1, x_2 有关, $y=f(x_1, x_2)$ 。若 x_1, x_2 的近似值为 x_1^*, x_2^* , 相应解为 y^* , 则当数据误差较小时解的绝对误差

$$\begin{aligned} \epsilon(y) &= y - y^* = f(x_1, x_2) - f(x_1^*, x_2^*) \\ &\approx dy = \frac{\partial f(x_1, x_2)}{\partial x_1} \epsilon(x_1) + \frac{\partial f(x_1, x_2)}{\partial x_2} \epsilon(x_2) \end{aligned} \quad (1-15)$$

$$\text{解的相对误差} \quad \epsilon_r(y) \approx \frac{dy}{y} = \sum_{i=1}^2 \frac{\partial f(x_1, x_2)}{\partial x_i} \cdot \frac{x_i}{f(x_1, x_2)} \cdot \epsilon_r(x_i) \quad (1-16)$$

设 $y=f(x)$ 为一元函数, 则计算函数值的误差为

$$\begin{aligned} \epsilon(y) &\approx dy = f'(x) dx \approx f'(x) \epsilon(x) \\ \epsilon_r(y) &\approx \frac{dy}{y} \end{aligned}$$

利用微分可得到两数和、差、积、商的误差估计。

二、加减运算

设 $y=f(x_1, x_2)=x_1+x_2$, 利用式(1-15)可得

$$\epsilon(x_1+x_2) \approx \epsilon(x_1) + \epsilon(x_2) \quad (1-17)$$

$$|\epsilon(x_1+x_2)| = |\epsilon(x_1) + \epsilon(x_2)| \leq |\epsilon(x_1)| + |\epsilon(x_2)| \quad (1-18)$$

因此, 任何两个数之和的绝对误差等于两个数的绝对误差之和, 任何两个数之和的绝对误差限为这两个数的绝对误差限之和, 且可推广到有限多个数相加的情形。所以作大量加减运算后的绝对误差限是绝不可以忽视的。

由于加法和减法是互为逆运算关系, 所以减法可以化为加法情形讨论。

利用式(1-16)可得两个数之和的相对误差为

$$\epsilon_r(x_1+x_2) = \frac{x_1}{x_1+x_2} \epsilon_r(x_1) + \frac{x_2}{x_1+x_2} \epsilon_r(x_2) \quad (1-19)$$

当 x_1 和 $-x_2$ 相当接近时, $x_1 - x_2 \approx 0$, $\left| \frac{x_1}{x_1 + x_2} \right|$ 和 $\left| \frac{x_2}{x_1 + x_2} \right|$ 都将很大, 所以相近两数之差的相对误差将很大, 即原始数据的误差会对计算结果产生很大的影响。

例 1 用四位有效数字计算 $y = \sqrt{1\,001} - \sqrt{1\,000}$ 的值。

解 如果直接计算, 则

$$y = \sqrt{1\,001} - \sqrt{1\,000} = 31.64 - 31.62 = 0.02$$

由于 y 的准确值是 $0.015\,807\,437\,4\dots$, 可见直接计算所得的近似值仅有一位有效数字, 其相对误差大于 26% 。有时若作适当变形后计算, 可以避免相近两数相减的计算, 如

$$y = \sqrt{1\,001} - \sqrt{1\,000} = \frac{1}{\sqrt{1\,001} + \sqrt{1\,000}} = \frac{1}{63.26} \approx 0.015\,81$$

所得结果与准确值比较可知具有四位有效数字, 其相对误差不超过 0.02% 。所以在数值计算中, 必须避免相近两数相减, 以免损失有效数字的位数。

三、乘除运算

利用微分可得两数积的绝对误差为

$$\varepsilon(x_1 x_2) \approx d(x_1 x_2) \approx x_2 \varepsilon(x_1) + x_1 \varepsilon(x_2) \quad (1-20)$$

相对误差为
$$\varepsilon_r(x_1 x_2) = \frac{\varepsilon(x_1 x_2)}{x_1 x_2} \approx \varepsilon_r(x_1) + \varepsilon_r(x_2) \quad (1-21)$$

两数商的绝对误差为

$$\varepsilon\left(\frac{x_1}{x_2}\right) \approx d\left(\frac{x_1}{x_2}\right) \approx \frac{x_2 \varepsilon(x_1) - x_1 \varepsilon(x_2)}{x_2^2}, \quad (x_2 \neq 0) \quad (1-22)$$

相对误差为
$$\varepsilon_r\left(\frac{x_1}{x_2}\right) = \varepsilon\left(\frac{x_1}{x_2}\right) \cdot \frac{x_2}{x_1} \approx \varepsilon_r(x_1) - \varepsilon_r(x_2), \quad (x_2 \neq 0) \quad (1-23)$$

从而得出: 两数乘积的相对误差, 可看作是各乘数的相对误差之和; 两数商的相对误差, 可看作是被除数与除数的相对误差之差。

上述误差积累规律, 对多个近似数的运算也是成立的。通常, 任意多次连乘连除所得结果的相对误差限, 可看作是各乘数和除数的相对误差限之和。

第五节 算法的数值稳定性

一、算法的数值稳定性概念

所谓**算法**, 是指对一些数据按某种规定的顺序进行的运算序列。在实际计算中, 对于同一问题我们选用不同的算法, 所得结果的精度往往大不相同。这是因为初始数据的误差或计算中的舍入误差在计算过程中的传播, 因算法不同而异, 于是就产生了算法的数值稳定性问题。一个算法, 如果计算结果受误差的影响小, 就称这个算法具有较好的**数值稳定性**。否则, 就称这个算法的数值稳定性不好。例如在第四节的例 1 中, 第一种算法就是不稳定的, 而第二种算法是数值稳定的。下面再举一个例子。

例 1 一元二次方程

$$x^2 + 2px + q = 0 \quad (1-24)$$

的两个根分别为 $x_1 = -p + \sqrt{p^2 - q}$ 和 $x_2 = -p - \sqrt{p^2 - q}$, 当 $p = -0.5 \times 10^5$, $q = 1$ 时, 方程的