



普通高等教育“十一五”国家级规划教材

计算方法引论

(第4版)

*Introduction
to Numerical
Calculation
Method*

徐萃薇 孙绳武 编著

高等教育出版社

普通高等教育“十一五”国家级规划教材

计算方法引论

Jisuan Fangfa Yinlun

(第4版)

徐萃薇 孙绳武 编著

高等教育出版社·北京

内容提要

本书为普通高等教育“十一五”国家级规划教材。本书服务于多层次、多专业、多学科的教学需要,在选材上考虑普适性,涉及现代数字电子计算机上适用的各类数学问题的数值解法及必要的基础理论;在材料组织安排上给讲授者根据教学要求和学生情况适当裁剪的自由,一些内容还可作为阅读材料。

本次改正了之前各版中发现的各种错误和不当之处,并对全书整理、修改,增加了一些内容,重写了某些章节。第三章增加了 Chebyshev 多项式对函数逼近的应用等内容;第五章增加了自适应数值积分技术一节;微分方程数值解的内容做了较大调整,改写了第十二、十三章;第十四章增加了节点编序方法,使方程组的写法更加完整。

本书算法描述不拘一格,或用自然语言,或用某种形式语言(以描述某些细节),便于理解,也便于编程,可作为工科非计算数学专业本科生学习“计算方法”课程的教材,也可作为科技人员进修、自学的参考用书。

图书在版编目(CIP)数据

计算方法引论 / 徐萃薇,孙绳武编著. --4 版. --
北京:高等教育出版社,2015.3
ISBN 978-7-04-041889-7

I. ①计… II. ①徐… ②孙… III. ①计算方法-高等学校-教材 IV. ①O241

中国版本图书馆 CIP 数据核字(2015)第 014347 号

策划编辑 时 阳 责任编辑 时 阳 封面设计 张 志 版式设计 童 丹
插图绘制 郝 林 责任校对 杨凤玲 责任印制 毛斯璐

出版发行	高等教育出版社	网 址	http://www.hep.edu.cn
社 址	北京市西城区德外大街4号		http://www.hep.com.cn
邮政编码	100120	网上订购	http://www.landaco.com
印 刷	三河市骏杰印刷有限公司		http://www.landaco.com.cn
开 本	850mm×1168mm 1/16	版 次	1982年11月第1版
印 张	24		2015年3月第4版
字 数	540千字	印 次	2015年3月第1次印刷
购书热线	010-58581118	定 价	33.90元
咨询电话	400-810-0598		

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换

版权所有 侵权必究

物 料 号 41889-00

第4版序言

本版依旧服务于多层次、多专业、多学科的教学。各部分内容的讲述力求由浅入深,逐步展开,把问题讲清说透,便于自学,也便于教师根据学生情况调节进度。

这次修订审阅了全书,改正了此前发现的各种错误和不当之处,并对全书整理、修改,增加了一些内容,重写了某些章节。

第一至五章,数值分析部分是基础内容,大都是必修的。二、三版都有修改,此次第三章增加了 Chebyshev 多项式对函数逼近的应用等内容,第五章增加了自适应数值积分技术一节。这部分叙述较为细致,在学生能接受时进度可以适当快些。

第六至十章,数值代数部分有一些内容是基础的,也有一些近代的内容,也有不少理论结果。正如第三版序言所指出的,代数计算在应用中所占份额较大,是比较活跃的领域,第三版有较大改动。应该指出,这部分介绍的算法大都比较成熟,是一些软件所采用的,在运算量、存储量方面多有推敲,学习时要仔细揣摩。数值代数软件是最早成熟的数学软件。C. Moler 是在研制数值代数软件 LINPACK 后开发 MATLAB 的,起初只是为了更好地应用 LINPACK。此次修订,这部分内容改动较少,只是适当增加了习题,冀能为大规模计算积累经验。

第十一至十五章,微分方程数值解的内容此次做了比较大的调整,改写了第十二、十三章。我们把介绍初(边)值问题的第十一至十三章做了统一考虑,试图把初值问题的差分法和有关概念形成系统。抛物型方程一章前移于常微分方程后,增加了半离散化以应用常微分方程的方法,接着一般地叙述差分法稳定性和收敛性的概念。限于本书的要求,常微分方程相应叙述没有严格展开,双曲型方程也只限于简单的应用。可在此基础上作进一步的阅读。顺便指出,初边值问题可以把边条件纳入许可函数中,从而不加区别地作为初值问题。第十四章加上了节点编序方法,方程组的写法也就完整了。另外,在椭圆型差分方程组的解法方面,第十二章讲了抛物型方程的 ADI 方法,可用于解椭圆型差分方程组,本章从第六章迭代改善角度引入了多网格法。

徐萃薇教授编写的《计算方法引论》第一版成书已近三十年,如果加上在北京大学编写的讲义就还要多四五年。遗憾的是,徐萃薇教授辞世已十五年,未见第二版以后的成书。好在第4版在各方的支持与努力下也终于出版了,似可聊慰她在天之灵吧。

最后,诚挚地感谢三十年来众多高校的老师选用本书并提出宝贵意见。感谢关爱、支持本书,对本书历次版本作出贡献的诸多同仁,恕不细举。

孙绳武

2014年2月9日

第三版序言

《计算方法引论》第三版(孙绳武主编)申报普通高等教育“十一五”国家级规划教材获准,现在终于付梓了。出版社方面力促其事,需求催生了第三版。众多读者和老师的关爱与支持也是我们努力工作的动力,感谢他们。

本版一仍旧旨,服务于多层次、多专业、多学科的教学。各部分内容的讲述力求由浅入深逐步展开,把问题讲清说透。便于自学,也便于教师根据学生情况调节进度。

在选材上,考虑到普适性,教材涉及了现代数字电子计算机上适用的、行之有效的各类数学问题的数值解法以及必要的基础理论。我们试图在广度和深度间把握适度,使学生面对科学与工程计算问题时知道如何求解,怎么解更好,能够做出一定的理论分析,也为其以后发展打好基础。

在材料组织安排上,我们也给讲授者根据教学要求和学生情况适当剪裁的自由,一些内容还可作为阅读材料。

作为基础课教材固然应当相对稳定,但是,随着科技的进步,应用计算方法的需求在增加,作为计算工具的计算机在发展,计算方法本身也在发展。一方面新方法涌现,一方面原有方法也随着计算机系统的变化而有不同评价。因此,计算方法教材内容的更新、重新组织也是必要的。

新版全书经过整理、润色,多处内容有所修改,乃至重写。考虑到代数计算在应用中所占份额较大,是比较活跃的领域,六至十章改动较大。

新增共轭斜量法、预善共轭斜量法、拟 Newton 法等。

改进了例题设置,增加数量,加强例题间联系。以图给读者示范,告诉他们怎么算,积累处理计算问题的经验,有助于设计算法。对于比较同一问题不同解法也能有感性的了解。

新增习题参考答案。

每章增写评述,总结该章内容,述说参考文献。

参考文献收集了国内外内容结构与本书相近的、有影响的、包括新近面世的一些书籍,并按大学生教材和研究生教材或专著分列,可供读者加深理解和进一步提高,有些对研究工作亦不无裨益。

重编索引,删除冗余,节省篇幅。

本书算法描述不拘一格,或用自然语言,或用某种形式语言(以描述某些细节),便于理解,也便于编程。一个长的过程分成几个部分,每一部分完成一定的计算任务(也许它还要再分细)。所给算法虽有某些软件设计思想,终非软件算法。建议结合例题和习题,利用成熟软件如 MATLAB 以交互方式、程序方式进行练习。通过这种训练形成自顶而下的、模块化的程序设计习惯。

学时仍可依第一版序言所建议的安排。作者经验,五六十学时的课程可只讲授一至十一章。各章内容还可适当选择,例如正交多项式等古典题材,四、七、九章中复杂的算法只作简介,理论题材少讲。对于72学时的课程,在学生接受程度较好的情况下,可适当加快前面的进度,或略讲六至十一章的一些理论,选讲十二至十五章部分内容。

时光荏苒,感谢二十年来众多院校的老师选用本书并提出宝贵意见。感谢“十一五”规划教材评审专家们对本书的信任。感谢高等教育出版社等有关方面的支持和努力。策划编辑倪文慧和责任编辑李华英精心编排、仔细校勘为本书增色不少。

清华大学李庆扬教授审阅全书。感谢他对本书历次版本所作出的杰出贡献。

最后,仍望各方同仁继续支持,不吝赐教,以使本书能不断改善、与时俱进。

孙绳武

2006年7月1日

第二版序言

徐萃薇教授编写的《计算方法引论》自出版以来已被很多院校用做“计算方法”课教材,颇受欢迎。本书出版前高等教育出版社计算机编辑室请我审阅,使我有机会拜读这本教材。我认为这本教材在不大的篇幅中将计算数学学科中数值分析、数值代数和微分方程数值解三大部分的基本内容深入浅出地做了介绍,实属难能可贵。教材特点是取材恰当,“少而精”,既重视算法和基本概念的描述,又有理论分析。针对非计算数学专业学生特点,很好地掌握了理论深度的分寸,做到理论联系实际,又不过分追求理论证明的完整性,使学生既学到数值计算方法,又能进行必要的理论分析,为解决科学与工程计算问题打下了良好基础,也为进一步提高奠定基础。

此外,本书适用面较广,可用于非计算数学专业48~72学时的不同类型“计算方法”课的教材,对学时较少的专业,只要适当削减书中部分章节仍可使用。当然也适合于在职科技人员进修、自学和参考。

本书修订版除保持了原版特点外,增加和改写了一些章节,使内容更充实,加强了算法描述与实现,有相当广度和一定深度,文字叙述更流畅,能适应新世纪“计算方法”课的教学要求,是一本好教材。

徐萃薇教授是我北大数学系师姐,高我两届,1958年我们同在东北人民大学(现吉林大学)进修,听苏联专家梅索夫斯基讲授“计算方法”,是国内最早从事计算数学教学的老师。以后几十年,由于都从事计算数学教学与研究,常有往来,探讨教学和学术问题,到她退休才很少见面,没想到她重病不起,未完成本书的修改就走了。好在孙绳武教授遵照她的遗愿,出色完成了修改任务,使本书修改版得以面世。我写这篇序言,也作为对徐萃薇教授的悼念。

李庆扬

2001年3月

第一版序言

本书是作者 1976 年—1980 年在北京大学为非计算数学专业的学生讲授计算方法课程时所写的教材。

由于电子计算机的迅速发展,国内外关于计算方法的教程和专著日益增多。对于非计算数学专业的学生来说,虽然他们需要学习计算方法,但并不需要在计算数学的理论问题上花费过多的时间。这样,写一本适合于物理系、计算机系、地球物理系、力学系等大学生和研究生的计算方法教材就成为十分必要的了。作者本人讲授这门课程四次,每讲完一次都对教材进行了修改。另外,北京的几个兄弟院校(北京工业大学、北京工业学院、北京化工学院、北京师范大学和北京师范大学等)从 1980 年起也采用此书作为教材,反映也比较好,并曾根据使用情况对此教材的内容提出了许多宝贵意见,本书就是在采纳了大家提出的修改意见后定稿的。

本书内容共分三部分。第一部分是数值分析,第二部分是数值代数,第三部分是常微分方程数值解法和偏微分方程数值解法。全书讲授学时为 72~80 学时,如果学时少于 72 学时,可少讲或不讲偏微分方程数值解法。如果学时多于 80 学时,可根据本书参考书目,增加所需的内容。

学习本书所必需的数学基础是微积分和线性代数,以及常微分方程和偏微分方程的基本概念。这是一般理工科大学的学生都具备的。对于自学过这些课程的青年,如果他们想进一步自学计算方法,以解决一些实际应用问题,本书对他们也是适宜的。本书每章都附有一定数量的习题,通过这些习题可以加深对各章内容的理解,掌握必要的解题技巧。

作者特别感谢北京大学计算数学教研室主任胡祖炽教授,他仔细地审阅了原稿,提出了大量的宝贵意见和建议;还要感谢北京师范大学沈嘉骥同志和北京大学王莲芬同志,他们根据使用这本教材的情况,给作者提出了不少有价值的修改意见。

徐萃薇

1982 年冬识于

北京大学燕东园

目录

第一章 误差	1	习题	65
1.1 误差的来源	1		
1.2 浮点数, 误差、误差限和有效数字	2	第四章 快速 Fourier 变换	67
1.3 相对误差和相对误差限	5	4.1 三角函数插值或有限离散 Fourier 变换 (DFT)	67
1.4 误差的传播	7	4.2 快速 Fourier 变换(FFT)	69
1.5 在近似计算中需要注意的一些现象	8	评述	76
评述	12	习题	77
习题	13	第五章 数值积分	78
第二章 插值法与数值微分	14	5.1 Newton - Cotes 公式	78
2.1 线性插值	14	5.2 梯形求积公式和抛物线求积公式的误差估计	81
2.2 二次插值	17	5.3 复化公式及其误差估计	85
2.3 n 次插值	22	5.4 逐次分半法	88
2.4 分段线性插值	28	5.5 加速收敛技巧与 Romberg 求积	91
2.5 Hermite 插值	33	5.6 Gauss 型求积公式	96
2.6 分段三次 Hermite 插值	35	5.7 自适应数值积分技术	103
2.7 样条插值函数	38	评述	106
2.8 数值微分	41	习题	107
评述	44	第六章 解线性代数方程组的直接法	110
习题	44	6.1 Gauss 消去法	110
第三章 数据拟合法	47	6.2 主元素消去法	119
3.1 问题的提出及最小二乘原理	47	6.3 LU 分解	123
3.2 多变量的数据拟合	52	6.4 对称正定矩阵的平方根法和 LDL ^T 分解	128
3.3 非线性曲线的的数据拟合	54	6.5 误差分析	131
3.4 正交多项式拟合	58		
评述	64		

评述	139	10.2 迭代法	213
习题	140	10.3 迭代收敛的加速	216
第七章 线性方程组最小二乘问题	142	10.4 Newton 法	220
7.1 矩阵的广义逆	142	10.5 弦位法	222
7.2 用广义逆矩阵讨论方程组的解	144	10.6 抛物线法	223
7.3 几个正交变换	146	10.7 解非线性方程组的 Newton 法 和拟 Newton 法	225
7.4 算法: A 列满秩	152	10.8 最速下降法	232
7.5 算法: 奇异值分解	159	评述	236
评述	161	习题	236
习题	162		
第八章 解线性方程组的迭代法	164	第十一章 常微分方程初值问题的 数值解法	239
8.1 几种常用的迭代格式	164	11.1 几种简单的数值解法	239
8.2 迭代法的收敛性及误差估计	170	11.2 R-K 方法	244
8.3 判别收敛的几个常用条件	174	11.3 线性多步法	248
8.4 收敛速率	176	11.4 预估-校正公式	252
8.5 共轭斜量法	178	11.5 常微分方程组和高阶微分 方程的数值解法	254
评述	186	11.6 自动选取步长的需要和事后 估计	256
习题	187	11.7 Stiff 方程	259
第九章 矩阵特征值和特征向量的 计算	190	评述	262
9.1 幂法	190	习题	262
9.2 幂法的加速与降阶	195	第十二章 抛物型方程的差分解法	265
9.3 反幂法	196	12.1 微分方程的差分近似	265
9.4 平行迭代法	197	12.2 边界条件的差分近似	268
9.5 QR 算法	200	12.3 几种常用的差分格式	270
9.6 Jacobi 方法	204	12.4 差分格式的稳定性 and 收敛性	273
评述	208	12.5 二维和三维热传导方程	279
习题	209	评述	284
第十章 非线性方程及非线性方程组 解法	211	附录	284
10.1 求实根的对分区间法	211	习题	286

第十三章 双曲型方程的差分解法·····	288	习题·····	311
13.1 差分格式的建立·····	289		
13.2 差分格式的收敛性·····	291	第十五章 有限元方法·····	313
13.3 差分格式的稳定性·····	293	15.1 通过一个例子看有限元方	
13.4 利用特征线构造差分格式·····	297	法的计算过程·····	313
评述·····	298	15.2 一般二阶常微分方程边值	
附录·····	299	问题的有限元解法·····	323
习题·····	301	15.3 平面有限元·····	329
		评述·····	338
		习题·····	339
第十四章 椭圆型方程的差分解法·····	302	部分习题参考答案·····	340
14.1 差分方程的建立·····	302		
14.2 差分方程组解的存在唯一性		参考文献·····	358
问题·····	305	索引·····	360
14.3 差分方法的收敛性与误差			
估计·····	307		
评述·····	311		

第一章 误差

1.1 误差的来源

用数学工具来解决实际问题时,在哪些地方会产生误差呢?首先,用数学模型来描述具体的物理现象要作许多简化,因此,数学模型本身就包含着误差,这种误差叫做“模型误差”.在数学模型中,通常总要包含一些观测数据,这种观测结果不是绝对准确的,因此还有“观测误差”.考虑到数据还可能由先前的计算等途径得到,也可概称之为“数据误差”.现举例说明.

例1 设一根铝棒在温度为 t 时的实际长度为 L_t ,在 $t=0$ 时的实际长度为 L_0 ,用 l_t 来表示铝棒在温度为 t 时的长度近似值,并建立一个数学模型

$$l_t = L_0(1 + \alpha t)$$

其中, α 是由实验观测到的常数

$$\alpha = (0.000\ 023\ 8 \pm 0.000\ 000\ 1) 1/^\circ\text{C}$$

则称 $L_t - l_t$ 为“模型误差”. $0.000\ 000\ 1$ 是 α 的“观测误差”.

例2 通常用

$$s(t) = \frac{1}{2}gt^2, \quad g \approx 9.81 \text{ m/s}^2$$

来描述自由落体下落时距离和时间的关系.设自由落体在时间 t 的实际下落距离为 \tilde{s}_t ,则把 $\tilde{s}_t - s(t)$ 叫做“模型误差”.

在解决实际问题时,数学模型往往很复杂,因而不易获得分析解,这就需要建立一套行之有效的近似方法或数值方法.模型的准确解与用数值方法求得的解之差称为“方法误差”,或者叫做“截断误差”.

例3 一个无穷级数

$$\sum_{k=0}^{\infty} \frac{1}{k!} f^{(k)}(x_0)$$

在实际计算时,只能取前面有限项(如 n 项)

$$\sum_{k=0}^{n-1} \frac{1}{k!} f^{(k)}(x_0)$$

来代替,这就抛弃了无穷级数的后半段,因而出现了误差,这种误差就是一种“截断误差”.对这个问题来说,截断误差是

$$\sum_{k=0}^{\infty} \frac{1}{k!} f^{(k)}(x_0) - \sum_{k=0}^{n-1} \frac{1}{k!} f^{(k)}(x_0) = \sum_{k=n}^{\infty} \frac{1}{k!} f^{(k)}(x_0)$$

最后,还有一类误差是因为在计算时总是只能取有限位数字进行运算而引起的,这种误差称为“舍入误差”。

例 4 $\pi = 3.141\ 592\ 6\dots$, $\sqrt{2} = 1.414\ 213\ 56\dots$, $\frac{1}{3} = 0.333\ 3\dots$, 等等。在计算机上运算时只能取有限位小数,如取小数点后 4 位数字,则

$$\rho_1 = 3.141\ 6 - \pi = +0.000\ 007\ 3\dots$$

$$\rho_2 = 1.414\ 2 - \sqrt{2} = -0.000\ 013\dots$$

$$\rho_3 = 0.333\ 3 - \frac{1}{3} = -0.000\ 033\dots$$

就是“舍入误差”。

概括起来,误差一般有模型误差、数据误差、截断误差和舍入误差。在计算方法中,主要讨论的是截断误差和舍入误差。

1.2 浮点数, 误差、误差限和有效数字

1. 浮点数

任何一个浮点数均可表示为

$$x = \pm w \times \beta^J = \pm 0.\alpha_1\alpha_2\dots\alpha_t \times \beta^J, L \leq J \leq U$$

其中, β 叫做这个数的基,如常见的十进制数, $\beta = 10$; 计算机上使用的二进制数, $\beta = 2$ 。 J 是阶,是一个整数,取正数、负数或零。 w 是尾数,由 t 位小数构成, $0 \leq \alpha_i \leq \beta - 1$ ($i = 1, 2, \dots, t$)。若 $\alpha_1 \neq 0$, 则称该浮点数为规格化浮点数。

若用 \mathcal{F} 来表示一个系统的浮点数的集合, 则

$$\mathcal{F} = \{x: x = \pm 0.\alpha_1\dots\alpha_t \times \beta^J, 0 \leq \alpha_i \leq \beta - 1, \alpha_1 \neq 0, L \leq J \leq U\} \cup \{0\}$$

显然,集合 \mathcal{F} 可以用四个元素的数组 (β, t, L, U) 来刻画。对于不同的机器,这四个值不一定相同,常见的有 $(16, 14, -64, 63)$ 。它表示一个十六进制数集,每个数有 14 位小数,阶码为 $-64 \sim +63$ 。

对于一个特定的机器来说,尾数的位数 t 是固定的,也称其机器精度有 t 个 β 进位数字。浮点数中阶的上界 (U) 和下界 (L) 分别称为上溢限和下溢限。若在计算过程中产生的数的指数 (J) 超出上溢限或下溢限,便无法在该机器上表示,这种现象称为上溢或下溢。

集合 \mathcal{F} 是包含 $2(\beta - 1)\beta^{t-1}(U - L + 1) + 1$ 个数的有限集,这些数分布在区间 $[m, M]$, $[-M, -m]$ 和 $\{0\}$ 中,其中

$$m = \beta^{L-1}, M = \beta^U(1 - \beta^{-t})$$

这些数在 $[m, M]$ 和 $[-M, -m]$ 中的分布是不等距的。例如, $\beta = 2, t = 2, L = -1, U = 2$, 则 \mathcal{F} 中 17 个点的分布如图 1-1 所示。

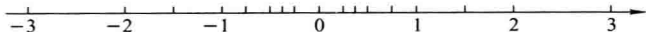


图 1-1

既然只是一个有限集,它就不可能将 $[m, M]$ 和 $[-M, -m]$ 中的任意实数表示出来,这就决定了计算机中的浮点运算是无法精确进行的.

2. 误差、误差限和有效数字

如何来定义误差?若用 x^* 来表示 x 准确值的一个近似值,则此近似值 x^* 和 x 准确值的差称为误差,用 e^* 来表示,即

$$e^* = x^* - x$$

这样定义后,就有 $x = x^* - e^*$,即近似值去掉(减去)它的误差就是准确值.因此,把误差的相反数 $-e^*$ 叫做近似值 x^* 的“修正值”.或者说,近似值加上它的修正值就是准确值.

误差可正可负.当误差为正时,近似值偏大,叫做“强近似”;当误差为负时,近似值偏小,叫做“弱近似”.

由于在一般情况下准确值 x 是不知道的,所以误差 e^* 的准确值也不可能求出.但根据具体测量或计算的情况,可以事先估计出误差的绝对值不能超过某个正数 ε^* , ε^* 叫做误差绝对值的“上界”,或称“误差限”.

定义 如果

$$|e^*| = |x^* - x| \leq \varepsilon^*$$

ε^* 就叫做近似值 x^* 的“误差限”.误差限一定是一个非负数.

因为在任何情况下都有

$$|x^* - x| \leq \varepsilon^*$$

即

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

这就表明 x 在 $[x^* - \varepsilon^*, x^* + \varepsilon^*]$ 这个区间内,用

$$x = x^* \pm \varepsilon^*$$

来表示近似值 x^* 的精确度,或准确值所在的范围.

例 5 用一把有毫米刻度的米尺测量桌子的长度,读出的长度 $x^* = 1\,235$ mm,是桌子实际长度 x 的一个近似值,由米尺的精度可知,这个近似值的误差不会超过 0.5 mm,则有

$$|x^* - x| = |1\,235 - x| \leq \frac{1}{2}$$

即

$$1\,234.5 \leq x \leq 1\,235.5$$

这表明 x 在 $[1\,234.5, 1\,235.5]$ 这个区间内,写成

$$x = (1\,235 \pm 0.5) \text{ mm}$$

例 6 光速 c 的近似值目前公认的是

$$c^* = 2.997\,925 \times 10^{10} \text{ cm/s}$$

通常记为

$$c = (2.997\,925 \pm 0.000\,001) \times 10^{10} \text{ cm/s}$$

取 x 的近似值 x^* , 通常用四舍五入的方法取前面几位.

例 7 $x = \pi = 3.141\,592\,65\dots$, 按四舍五入的原则

取 1 位: $x_1^* = 3$, $e_1^* \approx -0.14$.

取 3 位: $x_3^* = 3.14$, $e_3^* \approx -0.001\,6$.

取 5 位: $x_5^* = 3.141\,6$, $e_5^* \approx +0.000\,007$.

取 6 位: $x_6^* = 3.141\,59$, $e_6^* \approx -0.000\,003$.

定义 如果 x 的近似值 x^* 的误差限是某一位上的半个单位, 该位到 x^* 的第一位非零数字共有 n 位, 就说 x^* 有“ n 位有效数字”, 或者说 x^* 准确到该位.

用四舍五入法取准确值的前 n 位作为近似值 x^* , 则 x^* 有 n 位有效数字. 上面例子中的 $x_3^* = 3.14$ 以 3 位有效数字来表示 π , 它的误差限为

$$|x_3^* - \pi| \leq \frac{1}{2} \times 10^{-2}$$

$x_5^* = 3.141\,6$ 以 5 位有效数字来表示 π , 它的误差限为

$$|x_5^* - \pi| \leq \frac{1}{2} \times 10^{-4}$$

定义给出了有效数字位和误差限之间的关系.

若用 x^* 表示 x 的近似值, 并将 x^* 表示成规格化浮点数

$$x^* = \pm 0.\alpha_1\alpha_2\dots\alpha_n \dots \times 10^p \quad (1.1)$$

若

$$|x^* - x| \leq \frac{1}{2} \times 10^{p-n}$$

则近似值 x^* 具有 n 位有效数字, 这里 p 是一个整数, $\alpha_1, \alpha_2, \dots, \alpha_n, \dots$ 都是 $0 \sim 9$ 中的一个数字, 而且假定 $\alpha_1 \neq 0$.

从这里也可以看出误差限和有效数字位之间的关系, 并可以通过有效数字位来刻画误差限.

例 8 若 $x^* = 3\,587.64$ 是 x 的具有 6 位有效数字的近似值, 那么它的误差限是

$$|x^* - x| \leq \frac{1}{2} \times 10^{4-6} = \frac{1}{2} \times 10^{-2}$$

若 $x^* = 0.002\,315\,6$ 是 x 的具有 5 位有效数字的近似值, 则误差限是

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

当然, 也可由定义直接读出两数的误差限: 0.005 和 $0.000\,000\,05$.

1.3 相对误差和相对误差限

上节引入的误差和误差限的概念不能说明近似的好坏程度,它是有量纲单位的.例如,工人甲平均每生产一百个零件有一个次品,而工人乙则平均每生产五百个零件有一个次品.他们的次品都是一个,但显然乙的技术水平要比甲高.这就启发人们除了要看次品的多少外,还必须注意产品的合格率,甲的次品率是百分之一,而乙的次品率是五百分之一.把近似数的误差与准确值的比值定义为“相对误差”,记为 e_r^* .

定义 记

$$e_r^* = \frac{e^*}{x} = \frac{x^* - x}{x}$$

为近似数 x^* 的相对误差.在实际计算中,由于准确值 x 总是不知道的,所以也把

$$e_r^* = \frac{e^*}{x^*} = \frac{x^* - x}{x^*}$$

记为近似值 x^* 的相对误差,条件是 e_r^* 比较小.

与前面引入的误差一样,相对误差可正可负.相对误差绝对值的上界叫做“相对误差限”,记为

$$\varepsilon_r^* = \frac{\varepsilon^*}{|x|} \quad \text{或} \quad \varepsilon_r^* = \frac{\varepsilon^*}{|x^*|}$$

其中, ε^* 是 x^* 的误差限.

为了区别相对误差与 1.2 节中讲的误差,也把 1.2 节中讲的误差叫做“绝对误差”.

例 9 $c = (2.997\,925 \pm 0.000\,001) \times 10^{10}$ cm/s,这时 $c^* = 2.997\,925 \times 10^{10}$ cm/s 的相对误差限是

$$\varepsilon_r^* = \frac{0.000\,001}{2.997\,925} \approx 0.000\,000\,3$$

c^* 是 c 的很好的近似值,如果取

$$c^{**} = 3 \times 10^{10} \text{ cm/s}$$

作为光速的近似值,则有

$$\varepsilon_r^{**} \approx \frac{0.002\,1}{3} = 0.000\,7$$

相对误差限不到千分之一, $c^{**} = 3.00 \times 10^{10}$ cm/s 是用四舍五入法取 c 的前 3 位数的近似值,它有 3 位有效数字.

同样,可以给出相对误差限和有效数字位的关系.

设形如(1.1)式的近似数 x^* 具有 n 位有效数字,则其相对误差限可取为

$$\frac{1}{2\alpha_1} \times 10^{-(n-1)}$$

要说明这一点并不困难. 由定义知 x^* 是具有 n 位有效数字的近似值, 因此

$$|x^* - x| \leq \frac{1}{2} \times 10^{p-n}$$

$$\frac{|x^* - x|}{|x^*|} \leq \frac{\frac{1}{2} \times 10^{p-n}}{|x^*|}$$

而 $|x^*| \geq \alpha_1 \times 10^{p-1}$, 所以

$$\frac{|x^* - x|}{|x^*|} \leq \frac{\frac{1}{2} \times 10^{p-n}}{\alpha_1 \times 10^{p-1}} = \frac{1}{2\alpha_1} \times 10^{-(n-1)}$$

结论得证. 要注意的是, 用相对误差限来得到有效数字位时, 其关系略有不同.

形如(1.1)式的近似数 x^* , 相对误差限满足关系式

$$\varepsilon_r^* \leq \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)}$$

则 x^* 至少具有 n 位有效数字.

和前面关系不同的是, 分母为 $2(\alpha_1 + 1)$. 这是因为若要 x^* 具有 n 位有效数字, 则只要能证明

$$|x^* - x| \leq \frac{1}{2} \times 10^{p-n}$$

即可, 而

$$|x^* - x| \leq |x^*| \varepsilon_r^* \leq |x^*| \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)}$$

因为 $|x^*| \leq (\alpha_1 + 1) \times 10^{p-1}$, 所以

$$\begin{aligned} |x^* - x| &\leq (\alpha_1 + 1) \times 10^{p-1} \times \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)} \\ &= \frac{1}{2} \times 10^{p-n} \end{aligned}$$

例 10 用 $x^* = 2.72$ 来表示 e 具有 3 位有效数字的近似值, 则相对误差限是

$$\frac{1}{2 \times 2} \times 10^{-(3-1)} = \frac{1}{4} \times 10^{-2}$$

数 x 进入计算机时要转化成规格化浮点数并舍入成 t 位尾数, 记为 $fl(x)$, 它有 t 位有效数字, 并且不难得到关系式

$$fl(x) = x(1 + \varepsilon), |\varepsilon| \leq \frac{1}{2} \beta^{-(t-1)} \quad (1.2)$$