

第一章 数值计算中的误差

科学的迅猛发展,使计算机已成为科学研究与工程设计的主要工具之一。在实际工作中,由于受计算机字长的限制,我们不得不将无限小数截取有限位,转化为计算机能表示出的近似值,并进行这些近似值间的运算,这就是数值计算。将待求解问题的数学模型(即能描述并等价于实际问题的数学问题)简化为一系列算术运算、逻辑运算等,以便在计算机上求出问题解的近似值(常称为数值解)的方法,叫做数值计算方法,简称算法。我们的问题是,如何规定运算的次序、运算方式等,提出相应算法,并对算法的收敛性、稳定性,特别是误差加以分析?

本章的目的是帮助我们认识误差,了解误差对运算结果的影响,掌握误差分析的基本方法。

§ 1.1 误差的种类及来源

产生误差的原因多种多样,但常见的有以下几种。

1.1.1 模型误差

解决工程技术实际问题常要首先建立对实际问题的数学描述——数学模型。建立模型的时候,为便于分析和计算,需要对问题进行合理的简化,忽略掉一些次要因素,使建立的模型既能解决实际问题,又简单易行。但无论如何,合理简化的结果必会使建立的模型与实际问题间产生一定的差异,使模型仅能近似地描述实际问题,由此而产生的误差称为模型误差。

模型误差过大时,则要考虑修改模型。

1.1.2 观测误差

为研究实际问题,常要首先搜集有关的数据。实验观察、测量等都是获取数据的重要手段。受观测手段(实验仪器与测量仪器等)的限制,得到的数据只能是一些近似值,带有一定误差,这种误差称为测量误差,或观测误差。

通常,观测误差可根据测量工具、实验仪器本身的精度来确定。但有时也受到某种干扰而带有随机性,也很难估计。

1.1.3 舍入误差

数据计算过程中,要对无限小数(如 $\pi = 3.14159265\cdots$, $\frac{1}{3} = 0.3333\cdots$ 等)进行舍入。常见的舍入规则有:“只舍不入”、“只入不舍”与“四舍五入”。至于要用哪种舍入规则常由实际问题来定。由舍入而引起的误差称为舍入误差。

本书中仅考虑由“四舍五入”而引起的舍入误差。

1.1.4 截断误差

计算机只能完成有限次算术与逻辑运算。对于一些超越运算(如微分、积分、无穷级数求和等),则要化为一系列有限次的算术运算或逻辑运算来处理,处理的结果也会产生一定的误差。例如,利用函数 $\sin x$ 的幂级数展开式:

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

计算 $\sin 0.5$ 的近似值时,取

$$\sin x \approx x - \frac{x^3}{3!} + \frac{x^5}{5!}$$

则会产生误差

$$-\frac{x^7}{7!} + \frac{x^9}{9!} - \dots \triangleq e_4(x)$$

(在此为 $e_4(0.5)$,是由甩掉无穷级数自第四项开始后的所有项产生的)。计算 $\sin 0.5$ 时产生的误差满足

$$|e_4(0.5)| \leq \frac{0.5^7}{7!} \approx 1.55 \times 10^{-4}$$

这种由仅保留了超越运算过程中的有限次运算部分而引起的误差,称为截断误差。被舍去的部分称为余项。

一般来说,截断误差常可根据算法所舍去的余项来估计。

上述四种误差,尽管单独看来常是微不足道的,但科学计算与工程设计常要进行成千上万次的计算;计算中的多次舍入、截断等,均会造成误差的积累,甚至使结果严重偏离真值(这在后面例子中会看到),因而要对误差予以足够的重视。

通常,模型误差和观测误差往往是计算工作者不能独立解决的。本书主要讨论截断误差与舍入误差。

§ 1.2 误差的表示

要认识误差,分析误差,则应给出能体现误差的量,本节的绝对误差与相对误差描述了误差在量上的本质特征。

1.2.1 绝对误差

设 a 是精确值 x (即真值)的一个近似值,常称

$$e(x) = x - a$$

为近似值 a 的绝对误差,简称误差。又若 ϵ 是 $e(x)$ 的一个上界,即

$$|e(x)| \leq \epsilon$$

则称 ϵ 为近似值 a 的绝对误差限或绝对误差界,简称误差限。

例如,用刻有毫米的米尺测量一个桌子,测得桌子长为 985mm,设 x 是桌子的真正长度,则 985 实际上为 x 的一个近似值,由于测量误差不会超过半个毫米,即

$$|985 - x| \leq 0.5(\text{mm})$$

从而

$$985 - 0.5 \leq x \leq 985 + 0.5$$

这说明 x 的真值在区间 [984.5, 985.5] 内, 通常也记为

$$x = 985 \pm 0.5$$

上例中, 由于真值 x 无法测得, 绝对误差 $985 - x$ 实际上是根本无法知道的, 但我们却很容易得到了其绝对误差限。实用中, 往往能得到误差限也就够了。

1.2.2 相对误差

用绝对误差及其误差限来刻划近似值的精确程度是有局限性的。例如, 假设你要买煤, 由于秤的原因, 无论你买多少煤, 都会少给你 5 千克, 那么你会选择买 50 千克还是选择买一吨煤? 当然你会毫不犹豫地选择买一吨煤。为什么呢? 因为你认为买 50 千克煤“差的太多了”。这里两种情况下的绝对误差都是 5 千克, “差的太多”不能用绝对误差来刻划, 它是相对于被衡量值本身的大小而言的。

我们称绝对误差与被衡量值真值之比, 即

$$\frac{e(x)}{x} = \frac{x - a}{x} \triangleq e_r(x) \quad (1.1)$$

为近似值 a 的相对误差。(注意, 今后我们常使用符号 \triangleq 表示“记为”)若有一个正数 δ 使
 $|e_r(x)| \leq \delta$

则称 δ 为近似值 a 的相对误差限。注意, 相对误差与相对误差限是无量纲的。

前述例子中, 选择买一吨煤的相对误差为 $5/1000 = 1/200$, 而买 50 千克的相对误差为 $5/50 = 1/10$, 为前种相对误差的 20 倍!

由(1.1)式可知, 相对误差与绝对误差间有如下关系:

$$e(x) = e_r(x)x$$

实际上, (1.1)式蕴含着, 单位真值中所含的绝对误差就是相对误差。因此, 在进行误差分析时, 相对误差常比绝对误差更能准确描述近似值的精确程度。然而, x 的真值通常总是无法知道的, 即便可以估计出绝对误差与绝对误差限, 也无法由(1.1)式给出相对误差与相对误差限。实际应用中, 常以

$$e'_r(x) = \frac{e(x)}{a} \quad (1.2)$$

作为相对误差。容易证明, $e'_r(x)$ 与 $e_r(x)$ 仅差一个 $e_r(x)$ 的高阶无穷小量, 因而用 $e'_r(x)$ 近似代替 $e_r(x)$ 不会引起更明显的误差。

另外, 相对误差也常用百分数表示, 如本节开始的例子中, 两个相对误差可以分别表示为 0.5% 与 10%。

最后, 需要指明的是, 分析误差常要同时考虑相对误差与绝对误差。

1.2.3 有效数字

设 x 是一个实数, 用四舍五入规则舍入后, 近似值的误差限为其末位的半个单位。如, $\pi = 3.1415926\cdots$, 四舍五入取四位小数, 近似值为 3.1416, 此时

$$|\pi - 3.1416| < \frac{1}{2} \times 10^{-4}$$

若 x 的近似值 a 的误差限是其某一位的半个单位, 则从这位起直到左边第一个非零数字为止的所有数字都称为 a 的有效数字。

具体地说, 若 x 的近似值

$$a = (\pm d_1 \times 10^{m-1} + d_2 \times 10^{m-2} + \cdots + d_{l-1} \times 10^{m-l+1} + \cdots + d_s \times 10^{m-s})$$

其中 $s \geq l$, 简记为 ($d_1 \neq 0$)

$$a = \pm 0.d_1d_2\cdots d_l\cdots d_s \times 10^m \quad (1.3)$$

若误差满足

$$|e(x)| = |x - a| \leq \frac{1}{2} \times 10^{m-l}$$

则称近似值 a 具有 l 位有效数字, 且 d_1, d_2, \dots, d_l 即为 a 的 l 位有效数字。

例如, 3.1416 是 π 具有 5 位有效数字的近似值, 3.1415 是 π 具有 4 位有效数字的近似值。

用计算机进行数值计算时, 受计算机字长的限制, 参加数值计算的数有一定的位数, 计算结果也只能保留有限位, 所保留下来的各位数字不一定都是有效数字; 同时, 即使保留下的是有效数字, 也往往只是计算舍去一些有效数字后的结果。

1.2.4 有效数字与误差的关系

设 a 是 x 的近似值且具有式(1.3)的形式, 则

$$|a| \geq d_1 \times 10^{m-1}$$

故, 由(1.2)式, 忽略 $e_r(x)$ 与 $e'(x)$ 的差别

$$|e_r(x)| = |\frac{e(x)}{a}| \leq \frac{0.5 \times 10^{m-l}}{d_1 \times 10^{m-1}} = \frac{1}{2d_1} \times 10^{-l+1}$$

因此, 当 $a = \pm 0.d_1d_2\cdots d_l\cdots d_s \times 10^{m-1}$ 作为 x 的近似值具有 l 位有效数字时, 相对误差满足

$$|e_r(x)| \leq \frac{1}{2d_1} \times 10^{-l+1}$$

反之, 若数 $a = 0.d_1d_2\cdots d_l\cdots d_s \times 10^{m-1}$ (或 $a = -0.d_1d_2\cdots d_l\cdots d_s \times 10^{m-1}$) 是 x 的一个近似值, 相对误差限为:

$$\frac{1}{2(d_1 + 1)} \times 10^{-l+1}$$

则易证明, a 至少具有 l 位有效数字。这就是有效数字与误差间的关系。

利用此关系, 不仅可以估计近似值的相对误差, 而且还可以求得为使精度(相对误差)在一定范围内时应取的有效位数。

例如, 用 3.1416 作为 π 的近似值, 其相对误差必满足:

$$|e_r(\pi)| \leq \frac{1}{2 \times 3} \times 10^{-3+1} \approx 1.667 \times 10^{-5}$$

又如, 为使 π 的近似值 a 的相对误差不超过 0.2% , 注意到 $\pi = 3.14159\dots$, 应使

$$|e_r(\pi)| = \frac{1}{2 \times (3 + 1)} \times 10^{-m+1} \leq 0.2\%$$

可求出 $m \geq 3$, 从而, 取三位有效数字即可, 即应取 $\pi = 3.14$ 。

§ 1.3 调差分析

1.3.1 算术运算中误差的传播

由误差概念与微分定义比较可以发现, 真值 x 处的绝对误差可以看做此处的微分, 即

$$e(x) = dx$$

则

$$e_r(x) = \frac{1}{x} e(x) = \frac{1}{x} d(x) = d(\ln x)$$

这说明, x 的近似值 a 的相对误差可以看做 $\ln x$ 的微分。可见, 微分是描述相对误差和绝对误差的一个有效工具。

利用微分与绝对误差、相对误差的上述关系, 容易得到

$$\begin{aligned} e(x \pm y) &= e(x) \pm e(y) \\ e(xy) &= xe(y) + ye(x) \\ e\left(\frac{x}{y}\right) &= \frac{1}{y} e(x) - \frac{x}{y^2} e(y) = \frac{x}{y} [e_r(x) - e_r(y)] \\ e_r(x+y) &= \frac{e(x)+e(y)}{x+y} = \frac{xe_r(x)+ye_r(y)}{x+y} \\ e_r(x-y) &= \frac{xe_r(x)-ye_r(y)}{x-y} \\ e_r(xy) &= e_r(x) + e_r(y) \\ e_r\left(\frac{x}{y}\right) &= e_r(x) - e_r(y) \end{aligned}$$

由这些公式可知

1. 大小相近的两个同号数相减时, 结果的相对误差可能会很大, 而使计算结果的有效数字严重丢失, 精度降低。
2. 两个数相乘, 其中一个数很大时, 乘积的绝对误差可能会增大。
3. 当除数的绝对值很小时, 商的绝对误差可能增大。

因此, 应尽量避免大小相近的两个同号数相减, 避免乘数的绝对值很大以及除数接近于零。如何避免呢? 下面对此进行讨论。

1.3.2 减小误差的方法

误差是不可避免的, 这是由我们所使用的计算工具——计算机所决定的。但根据经验, 通过误差分析, 避免出现较大误差, 保证精度在允许范围内是可以做到的。关于误差分析方法, 将在下一节中讨论; 在此我们介绍几种常用的小误差的方法。

1. 对于前边 1.3.1 中三种误差可能增大的情况, 常用将公式变形的方法来解决。例如, 当 x 很小时, $\cos x$ 接近等于 1。若我们要计算 $1 - \cos x$ 的值, 则可变形为

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

或用 $\cos x$ 的泰勒(Taylor)展开式

$$1 - \cos x = \frac{x^2}{2!} - \frac{x^4}{4!} + \dots$$

来计算。

2. 调整运算次序,常可以防止大数“吃”小数现象。为了对位,计算机中常要进行小数形式 $0.00153, 0.153, 15.3$ 等与形式 $1.53 \times 10^{-3}, 1.53 \times 10^{-1}, 1.53 \times 10^1$ 等的互相转换。这种允许小数点移位的表示方式称为数的浮点表示。若用五位浮点十进制数计算:

$$54321 + 0.4 + 0.3 + 0.3$$

计算机上先对位后运算,结果为

$$\begin{array}{r} 54321 \\ + 0.3 + 0.4 + 0.3 \\ \hline \end{array}$$

小数被大数吃掉了,为避免这种现象,改用计算顺序由小到大

$$0.3 + 0.3 + 0.4 + 54321$$

则可得到精确结果 54322。

3. 减少运算次数,可以减少误差积累。因此,尽量寻找运算次数较少的算法,是算法设计中常要考虑的,例如,计算多项式

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

在某点的值,常采用

$$a_n + x(a_{n-1} + x(\dots + x(a_1 + x)))\dots$$

由里向外计算的方式进行。注意,这样做不仅减少了误差积累,同时也减少了运算量。

1.3.3 误差分析方法

由误差与微分间关系,利用函数的全微分公式,可以进行误差分析。

例如,当计算函数 $p(x, y) = 5x^3 + 4xy + 7y^2$ 在 $x=1.414 \pm 0.5 \times 10^{-4}, y=2 \pm 0$ 处的函数值时,误差为:

$$|e(1.414, 2)| = |dp(x, y)|_{(1.414, 2)} = |(\frac{\partial p}{\partial x})_{(1.414, 2)}| |dx| + |(\frac{\partial p}{\partial y})_{(1.414, 2)}| |dy|$$

注意到 $dx=e(x)=0.5 \times 10^{-4}, dy=e(y)=0$ 并将各导数值求出,代入上式

$$|e(1.414, 2)| \leq [15 \times 1.414^2 + 4 \times 2] \times 0.5 \times 10^{-4} = 1.899547 \times 10^{-3}$$

取小数点后 4 位小数计算结果为:

$$P(1.414, 2) = 53.4478$$

不难发现,上述误差估计并未考虑到计算过程中产生的舍入误差,这就必然使估计的误差偏小。当计算过程简单时,计算过程中产生的误差也许很小;但当计算过程复杂时,其计算过程中的舍入误差就会对结果影响很大,因此必须加以考虑。实际中还采用以下方法:

1. 向前误差分析

从初始数据的误差出发,随着计算过程把逐次发生的计算误差累积起来,获得关于计算结果的误差或误差界,这种方法称为向前误差分析法。

例如,计算 $x = \sum_{i=1}^n x_i$ 时,若每个 x_i 的误差 e_i 满足

$$|e_i| \leq 0.5 \times 10^{-7} \quad (i = 1, 2, \dots, n)$$

其中 r 为某一正整数, 则和的误差

$$|e(x)| = |e_1 + e_2 + \dots + e_r| \leq \frac{n}{2} \times 10^{-r}$$

这种方法实际上很难实施, 且对于复杂的计算来说, 分析过程也会很令人乏味。

2. 向后误差分析

为估计计算误差, 假定初始值存在某一误差, 并使这一误差的影响等效于计算过程中所产生计算误差的影响, 并估计出误差的方法, 称为向后误差分析法。

例如, 解方程

$$ax = b \quad (1.4)$$

由于求解过程中发生了误差, 最终使求得的解 \bar{x} 产生了一个误差(也常称为扰动) $x^* - \bar{x}$ (x^* 代表方程的精确解), 将其等效地看成仅由数据存在误差 δb 而引起的, 方程求解过程本身不再产生误差, 则 \bar{x} 必为方程

$$ax = b + \delta b$$

的精确解, 即

$$a\bar{x} = b + \delta b \quad (1.5)$$

将 x^* 代入(1.4)式减(1.5)式可得

$$x^* - \bar{x} = -\delta b/a$$

这种方法常用于解线性方程组的误差分析, 也很有效。

3. 舍入误差的统计分析方法

向前误差分析方法, 尽管计算起来有时过于乏味, 但最终总能给出误差界。然而, 问题却常会是误差界过于保守而失去估计的意义。究其原因则是忽视了带有不同符号的误差, 在运算过程中相互抵消。

如何才能恰当估计出误差总积累的值? 对于成千上万次的计算, 计算的舍入误差会遵从统计规律。将计算的舍入误差看成随机变量, 设其概率密度函数为 $p(x)$, 则可以用数理统计的方法估计误差。以下通过例子说明误差估计的统计分析方法。

考虑 n 个数的相加 $x = \sum_{i=1}^n x_i$ 。用小数点后 r 位的浮点数表示每个 x_i , 则 x_i 的舍入误差 e_i 相互独立, 且

$$|e_i| \leq 0.5 \times 10^{-r}$$

e_i 在区间 $[-0.5 \times 10^{-r}, 0.5 \times 10^{-r}]$ 内等可能地取每一个值。因此, e_i 的分布密度函数 $p_i(x) = 10^r$ ($i = 1, 2, \dots, n$), 均值为零, 方差 σ_i^2 为:

$$\int_{-\infty}^{\infty} x^2 p_i(x) dx = \int_{-e_i}^{e_i} x^2 p_i(x) dx = \frac{e_i^2}{3} \quad (i = 1, 2, \dots, n)$$

其中 $e = 0.5 \times 10^{-r}$ 。 n 个数相加的总体误差

$$e(x) = \sum_{i=1}^n e_i$$

方差为

$$\sigma^2 = \sum_{i=1}^n \sigma_i^2 = \frac{n}{3} \epsilon^2$$

当 n 充分大时,由概率论的中心极限定理,相加的总误差 $\epsilon(x)$ 服从均值为零,方差为 $\frac{n}{3} \epsilon^2$ 的正态分布,即 $\epsilon(x)$ 的分布密度函数为

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

其中 $\sigma^2 = \frac{n}{3} \epsilon^2$ 。令 $P(|\epsilon(x)| \leq \epsilon) = \frac{1}{2}$, 即使 $\epsilon(x)$ 取值于 $(-\epsilon, \epsilon)$ 的概率为 0.5, 则有

$$\epsilon = 0.6745\sigma = 0.3894\epsilon\sqrt{n}$$

这就是 $\epsilon(x)$ 的一个大致误差限, 它与 \sqrt{n} 成正比, 与向前误差分析中估计的上限值 $\frac{n}{2} \times 10^{-6}$ 比较可知, n 较大时, 统计方法估计出的误差限小得多。

用向前误差分析方法求得的误差限常与运算次数 n 成正比。实际应用的经验告诉我们, 用 \sqrt{n} 代替向前误差分析的误差限中 n , 可以得到一个较可靠的误差限。

§ 1.4 算法的稳定性与收敛性

计算机上解数学问题常用逼近方法。当误差上限太大, 不可接受时, 往往采用修正目前估计值的方法加以改善, 直到误差满足要求。这样就产生了一列逼近值 $a_1, a_2, \dots, a_n, \dots$ 。若当 n 趋于无穷大时 a_n 的极限为要逼近的真值 a^* , 则称上述修正算法为收敛的。

1.4.1 收敛性

收敛的算法, 使我们有了求出满足精度要求的解的希望。理论上讲, 求出一个较好的近似解是不成问题的。但事实却并非如此。

例如, 给定级数

$$1 - \sum_{i=1}^{\infty} \frac{2}{16i^2 - 1} \triangleq a_n$$

易知,

$$\lim_{n \rightarrow \infty} a_n = \frac{\pi}{4}$$

取前 n 项逼近时的余项为 $\frac{1}{4n+3}$ 。现我们用它来计算 π 的精确到 10^{-6} 的近似值, 则令

$$\frac{4}{4n+3} \leq 10^{-6}$$

得到 n 至少为 10^6 , 也就是要计算约 10^6 项的和。这么多次运算后, 舍入误差的总积累也许已超过了 10^{-6} 。此算法用来计算 π 的值时, 收敛太慢了, 实际上根本不能用。

因此, 实用中, 不仅要用收敛的算法, 而且要用收敛更快的算法。

1.4.2 稳定性

一个算法, 若其对舍入误差不敏感时, 称为稳定的。换句话说, 算法的稳定性是指其计算结果受舍入误差影响的大小。受影响小或不受影响的算法为稳定的, 否则为不稳定的。

例如,给定递推公式

$$f_{n+1}(x) = (n+1)f_n(x) - x^{n+1} \quad (1.6)$$

这里 $f_0(x) = e^x - 1$ 。用五位浮点十进制数计算 $x=1$ 时 $f_n(x)$ 的值,可以得到

n	$f_n(1)$
0	1.7183
1	0.71830
2	0.43660
3	0.30980
4	0.23920
5	0.19600
6	0.17600
7	0.23200 (真值 $f_7(1)=0.14042$)

其中下边画线的为已失去的有效位。 $n=7$ 时,已无有效位数字。

同样为上述递推公式,将其改写为

$$f_n(x) = \frac{1}{n+1}(f_{n+1}(x) + x^{n+1})$$

反推,取 $f_{15}(1)=0$,则有

n	$f_n(x)$
15	0.00000
14	0.06667
13	0.07619
12	0.08278
11	0.09023
10	0.09911
9	0.10991
8	0.12332
7	0.14042

这里,有效位数很快增加到小数点后 5 位。

同一递推公式,用不同的递推次序,即取不同的算法,相差却很大,什么原因呢? 设递推公式(1.6)的精确值为 $f_n^*(x)$,初始值取为 $f_0(x) = e^x$ 的近似值 $\bar{f}_0(x)$ 时,迭代值为 $f_n(x)$ (假设不计计算中的舍入误差),注意到 $f_n^*(x)$ 实际上为

$$f_n^*(x) = n!(e^x - 1 - x - \frac{x^2}{2!} - \cdots - \frac{x^n}{n!})$$

用(1.6)式的迭代值 $f_n(x)$ 为

$$f_n(x) = n!(\bar{f}_0(x) - 1 - x - \frac{x^2}{2!} - \cdots - \frac{x^n}{n!})$$

从而有

$$f_n^*(x) - f_n(x) = n!(e^x - \bar{f}_0(x)) = n(f_{n-1}^*(x) - f_{n-1}(x))$$

可见,初始值的舍入误差为 $e^x - \bar{f}_0(x)$,即使计算过程中不产生误差,计算结果的误差也会急

剧增大。准确地说，每次迭代产生结果的误差为上次迭代产生误差的 n 倍，此算法不稳定。但若采用递推算法(1.7)式反推时有

$$f_n^*(x) - f_n(x) = \frac{1}{n+1} (f_{n+1}^*(x) - f_{n+1}(x))$$

可见每次迭代产生误差为逐次减小的，此时，算法则是稳定的。

因此，稳定性也是我们考察算法好坏的一个重要方面。

第二章 非线性方程的解

在科学的研究和工程设计中，常常会遇到高次方程，例如

$$x^7 - 8x + 9 = 0$$

高次代数方程通常没有求根的一般公式，很难求出它的精确解。又如

$$e^{-x} - x + 1 = 0$$

是一个超越方程，我们同样无法推导出根的精确表达式。

一般地说，我们可以将这类方程写成

$$f(x) = 0$$

如上所述，在许多场合不能求出其根的精确表达式。但是在实际工作中，常常只要求计算根的近似值。本章将要讨论这类方程根的各种近似计算方法，以及它们如何满足所给定的精度要求。方程 $f(x)=0$ 的根，通常也称为函数 $f(x)$ 的根。

§ 2.1 二分法

2.1.1 概念与算法

假设 $f(x)$ 在区间 $[a, b]$ 上连续，其中 $a < b$ ，且 $f(a)$ 与 $f(b)$ 异号，则根据连续函数的零点定理，方程 $f(x)=0$ 在 (a, b) 内至少有一个实根。

二分法求根的基本思想是将含根的区间平分为两个小区间，然后判断这两个小区间哪一个是含根的区间，再进一步将含根的小区间平分为两个更小的区间，再判断其中哪一个是含根的区间，进一步对它平分；如此继续下去。这样，含根的区间越分越小，若取含根区间的中点作为根的近似值，则随着含根区间的缩小，根的近似值就越来越接近根的精确值。

具体地说，设最初的含根区间为 $[a_0, b_0]$ ，即 $f(a_0)$ 与 $f(b_0)$ 异号，今取区间的中点

$$m = (a_0 + b_0)/2$$

并计算 $f(m)$ ，找出新的含根区间 $[a_1, b_1]$ ：

若 $f(a_0)f(m) < 0$ ，则 $[a_0, m]$ 为含根区间，即

$$a_1 = a_0, b_1 = m$$

若 $f(m)f(b_0) < 0$ ，则 $[m, b_0]$ 为含根区间，即

$$a_1 = m, b_1 = b_0$$

将此含根区间 $[a_1, b_1]$ 再平分，并找出下一个含根区间 $[a_2, b_2]$ ，继续平分后判断根所在的更小区间 $[a_3, b_3]$ ，如此周而复始，直到求出满足给定精度要求的近似根为止。如图 2-1 所示。

下面讨论误差分析。

由于最初的含根区间 $[a_0, b_0]$ 的宽度为 $b_0 - a_0$ ，经过 n 次平分后的含根区间 $[a_n, b_n]$ 的宽度为

$$b_n - a_n = (b_0 - a_0)/2^n$$

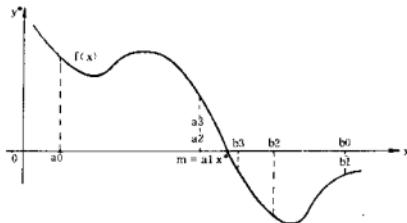


图 2-1 二分法

我们取第 n 次平分后的含根区间 $[a_n, b_n]$ 的中点作为根的近似值

$$x = (a_n + b_n) / 2$$

则其误差至多为 $(b_n - a_n) / 2$ 。也就是说,假设根的精确值为 x^* ,则有

$$|x^* - x| \leq (b_n - a_n) / 2$$

或

$$|x^* - x| \leq (b_0 - a_0) / 2^{n+1}$$

如果预先给定根的容许误差为 ϵ ,利用上式可以确定二分法中所需的步数,即

$$\epsilon \geq (b_0 - a_0) / 2^{n+1}$$

两边取对数,化简得到

$$n \geq \frac{\log(b_0 - a_0) - \log(2\epsilon)}{\log 2} \quad (2.1)$$

这个不等式可以用来确定二分法中所需要的步数。

在计算过程中,当(2.1)式满足时,即 $(b_n - a_n) / 2 \leq \epsilon$ 时,迭代将停止。或者预先给定一个充分靠近零的正数 δ 作为函数值的容许误差,当 $|f(a_n)|$ 或 $|f(b_n)| < \delta$ 时,取 a_n 或 b_n 作为根的近似值,迭代也停止。

二分法的算法

(首先判断初始区间是否是含根区间)

若 $f(a_0)f(b_0) > 0$, 终止

对 $n=0, 1, 2, \dots$ 做循环

若 $|f(a_n)|$ 或 $|f(b_n)| \leq \delta$, 则 $x \leftarrow a_n$ 或者 $x \leftarrow b_n$, 终止

否则 $m \leftarrow (a_n + b_n) / 2$

若 $|a_n - b_n| \leq \epsilon$, 则 $x \leftarrow m$ 终止

否则若 $f(a_n)f(m) < 0$, 则 $a_{n+1} \leftarrow a_n, b_{n+1} \leftarrow m$

否则 $a_{n+1} \leftarrow m, b_{n+1} \leftarrow b_n$

2.1.2 计算过程举例

例 2.1 用二分法求下列方程

$$x^3 - x - 1 = 0$$

在区间 $[1, 1.5]$ 内的一个实根,要求绝对误差不超过 0.001。

〔解〕 这里 $a_0=1, b_0=1.5$,

$$f(a_0) = -1, f(b_0) = 0.875$$

函数 $f(x)=x^3-x-1$ 在闭区间 $[1, 1.5]$ 上连续且 $f(a_0)$ 与 $f(b_0)$ 异号, 故方程在该区间内确实存在一个实根。于是根据给定误差, 利用公式(2.1)可计算出二分法中所需的步数

$$n \geq \frac{\log(1.5 - 1) - \log(2 \times 0.001)}{\log 2}$$

计算得到 $n \geq 7.97$, 即最多 8 步便能获得所要求的精度。

二分法的计算结果如下:

n	a	b	m	$f(m)$
0	1.0000	1.5000	1.2500	-0.2969
1	1.2500	1.5000	1.3750	0.2246
2	1.2500	1.3750	1.3125	-0.0515
3	1.3125	1.3750	1.3438	0.0828
4	1.3125	1.3438	1.3282	0.0149
5	1.3125	1.3282	1.3204	-0.0183
6	1.3204	1.3282	1.3243	-0.0018
7	1.3243	1.3282	1.3263	0.0068
8	1.3243	1.3263	1.3253	0.0025

这个方程根的精确值为 $1.324718\cdots$, 显然 $n=8$ 所得到的近似值能够满足给定的精度要求。如果在本例中要求绝对误差不超过 0.005, 则根据公式(2.1)计算出 $n \geq 5.6$, 即最多 6 步便能获得所要求的精度。计算结果表明, $n=5$ 时所得到的近似值已经满足给定的精度要求, 这表明公式(2.1)给出的结果是指最多所需的步数。

2.1.3 程序

- C 本算法用二分法求解方程 $F(x)=0$ 的根。
- C 算法输入为: FUN(X) — 函数 $F(x)$ 的表达式, A—起始搜索区间左端点, B—起始搜索区间右端点, DETA 和 EPS—求解精度控制、MAX—迭代控制。
- C 算法结束后, X 返回值为所求方程的根。

```
SUBROUTINE BISECT(A,B,DETA,EPS,MAX,X)
F=FUN(A)
G=FUN(B)
IF (SING(1.0,F).EQ.SING(1.0,G)) THEN
    WRITE(*,100)
100   FORMAT(' 数据参数有误! ')
    RETURN
END IF
IF (ABS(F).LT.ABS(B)) THEN
```

```

X=A
V=F
ELSE
  X=B
  V=G
END IF
DO 20 I=1,MAX
  IF((ABS(V).LE.DETA).OR.(ABS(A-B).LT.EPS)) RETURN
  X=(A+B)/2.0
  V=FUN(X)
  IF (SIGN(1.0,F).EQ.SIGN(1.0,V)) THEN
    A=X
    F=V
  ELSE
    B=X
    G=V
  END IF
20  CONTINUE
  WRITE(*,100)
100 FORMAT('在给定的迭代控制下未能达到所要求的求解精度!')
END

```

{
本算法用二分法求解方程 $F(x)=0$ 的根。
算法输入为:f—指向函数 $F(x)$ 的指针, a一起始搜索区间左端点, b一起始搜索区间右端点, deta 和
eps—求解精度控制, it_max—迭代控制。

算法结束后, 函数返回值为 ERROR_CODE 时表明数据参数有误; 函数返回值为所求方程的根。

```

}
type
  f_type=function(x,double);double;
function bisection(f,f_type;a,b,deta,eps,double;
  label
  EXIT,
  var
  fa,fb,m,fm;double;
  it;integer;
begin
  fa:=f(a);fb:=f(b);it:=0;
  if (abs(fa)<deta)
  then begin
    bisection:=a;goto EXIT
  end;
  if (abs(fb)<deta)

```

```

then begin
    bisection:=b;goto EXIT
end;
if (fa * fb>0)
then begin
    bisection:=ERROR_CODE;goto EXIT
end;
m:=(a+b) * 0.5;fm:=f(m);
while ((abs(a-b)>eps)and(abs(fm)>data)) do
begin
    if (fm * fa<0)
    then begin
        b:=m;fb:=fm;
    end
    else begin
        a:=m;fa:=fm;
    end;
    m:=(a+b) * 0.5;fm:=f(m);inc(it);
end;
bisection:=m;
EXIT:
end;

/*
本算法用二分法求解方程 F(x)=0 的根。
算法输入为:f—指向函数 F(x)的指针，a—起始搜索区间左端点，b—起始搜索区间右端点,data 和
eps—求解精度控制,it_max- 迭代控制。
算法结束后,函数返回值为 ERROR_CODE 时表明数据参数有误;函数返回值为所求方程的根。
*/
double bisection(double (*f)(double),double a,double b,
                  double data,double eps)
{
double fa,fb,m,fn;
fa=(*f)(a);fb=(*f)(b);it=0;
if (fabs(fa)<data) return(a);
if (fabs(fb)<data) return(b);
if (fa * fb>0) return(ERROR_CODE);
while ((fabs(a-b)>eps)&&(fabs(fm=(*f)(m=(a+b) * .5))>data))
{
    if (fm * fa<0) {b=m;fb=fn;}
    else {a=m;fa=fn;}
    it++;
}

```

```

    return(m);
}

```

§ 2.2 试位法

用二分法求方程 $f(x)=0$ 的根时,总是将含根区间 $[a, b]$ 平分,然而,方程的根很可能更靠近该区间的一个端点,在此端点处的函数值的绝对值可能要比另一端点处的函数值的绝对值小。也就是说,如果 $|f(a)| > |f(b)|$,那么方程的根很可能更靠近 b ,因此我们用过两点 $(a, f(a))$ 和 $(b, f(b))$ 的割线与 x 轴的交点来代替含根区间的中点,将含根区间分为两个小区间,从中找出新的含根区间。这种方法称为试位法或线性内插法。它有可能比二分法收敛得还要快,这要依赖含根区间内函数 $f(x)$ 的变化情况如何而定。

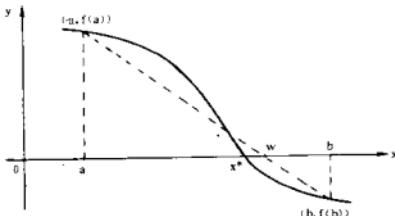


图 2-2 试位法

2.2.1 试位法算法与举例

如图 2-2 所示,过函数 $y=f(x)$ 曲线上两点 $(a, f(a))$ 和 $(b, f(b))$ 的割线为

$$y = f(a) + [(f(b) - f(a))/(b - a)](x - a)$$

令 $y=0$,解出割线与 x 轴的交点横坐标为

$$x = [f(b)a - f(a)b]/[f(b) - f(a)] \quad (2.2)$$

记为 W ,这实质上是区间 $[a, b]$ 的加权分割。 W 分含根区间 $[a, b]$ 为两个小区间,根据 $f(W)$ • $f(a)$ 的正或负判断出新的含根区间。这样重复上述过程,含根区间的长度将越来越小,特别是 $|f(W)|$ 将越来越小。当预先给定含根区间长度的容许误差 ϵ 和函数值的容许误差 δ 其中之一被满足时,计算过程终止,取 w 为根的近似值。

试位法算法

检验 $f(a_0)f(b_0) \leq 0$;否则在 $[a_0, b_0]$ 上无根

若 $|f(a_0)|$ 或 $|f(b_0)| \leq \delta$,则 $x \leftarrow a_0$ 或者 $x \leftarrow b_0$,终止

对 $n=0, 1, 2, \dots$ 做循环

$$w \leftarrow [f(b_n)a_n - f(a_n)b_n]/[f(b_n) - f(a_n)]$$

若 $|f(w)| < \delta$,则 $x \leftarrow w$,终止

否则若 $f(a_n)f(w) < 0$,则 $a_{n+1} \leftarrow a_n, b_{n+1} \leftarrow w$

否则 $a_{n+1} \leftarrow w, b_{n+1} \leftarrow b_n$

若 $|a_{n+1} - b_{n+1}| < \epsilon$, 则 $x \leftarrow w$, 终止

例 2.2 用试位法求下列方程

$$x^3 - x - 1 = 0$$

在区间 $[1, 1.5]$ 内的一个实根, 要求 $\epsilon = \delta = 10^{-4}$.

【解】计算结果如下

n	a_n	b_n	$f'(a_n)$	$f(b_n)$	w	$f(w)$
0	1.0000	1.5000	-1.0000	0.8750	1.2667	-0.2342
1	1.2667	1.5000	-0.2342	0.8750	1.3160	-0.0369
2	1.3160	1.5000	-0.0369	0.8750	1.3234	-0.0056
3	1.3234	1.5000	-0.0056	0.8750	1.3245	-0.0009
4	1.3245	1.5000	-0.0009	0.8750	1.3247	-0.0001

由于函数值的容许误差已经满足, 计算到此终止, 但此时含根区间的长度并没有满足精度要求。往往也会出现相反的情形, 即含根区间的长度已经满足精度要求, 而函数值的容许误差还没满足, 计算也终止。有时计算过程陷入无限循环, 因此在算法中设置一个允许迭代次数的上限往往是必要的。

从上面的例子看到, 虽然精度要求提高了, 试位法比二分法所需的迭代次数还要小。下面我们将进一步讨论改进试位法, 使其收敛速度更快。

2.2.2 外推试位法

从例 2.2 的计算结果中发现, 含根区间的一个端点 b_n 始终未变, 即割线的一个端点 $(b_n, f(b_n))$ 始终没动。如果将此端点改为 $(b_n, f(b_n)/2)$, 则割线的斜率将减小将近一半, 这时新的割线与 x 轴的交点将更靠近方程的根。因此收敛速度将加快。如图 2-3 所示。

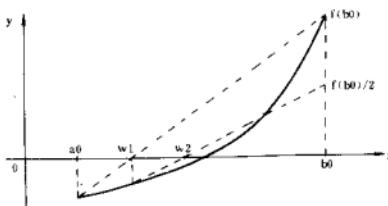


图 2-3 外推试位法

这种改进的试位法称为外推试位法或加速试位法。

外推试位法算法

设置最大迭代次数 max

$F \leftarrow f(a_0)$

$G \leftarrow f(b_0)$