



**DUT**

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 数值分析

大连理工大学国际信息与软件学院本科课程



**DUT**

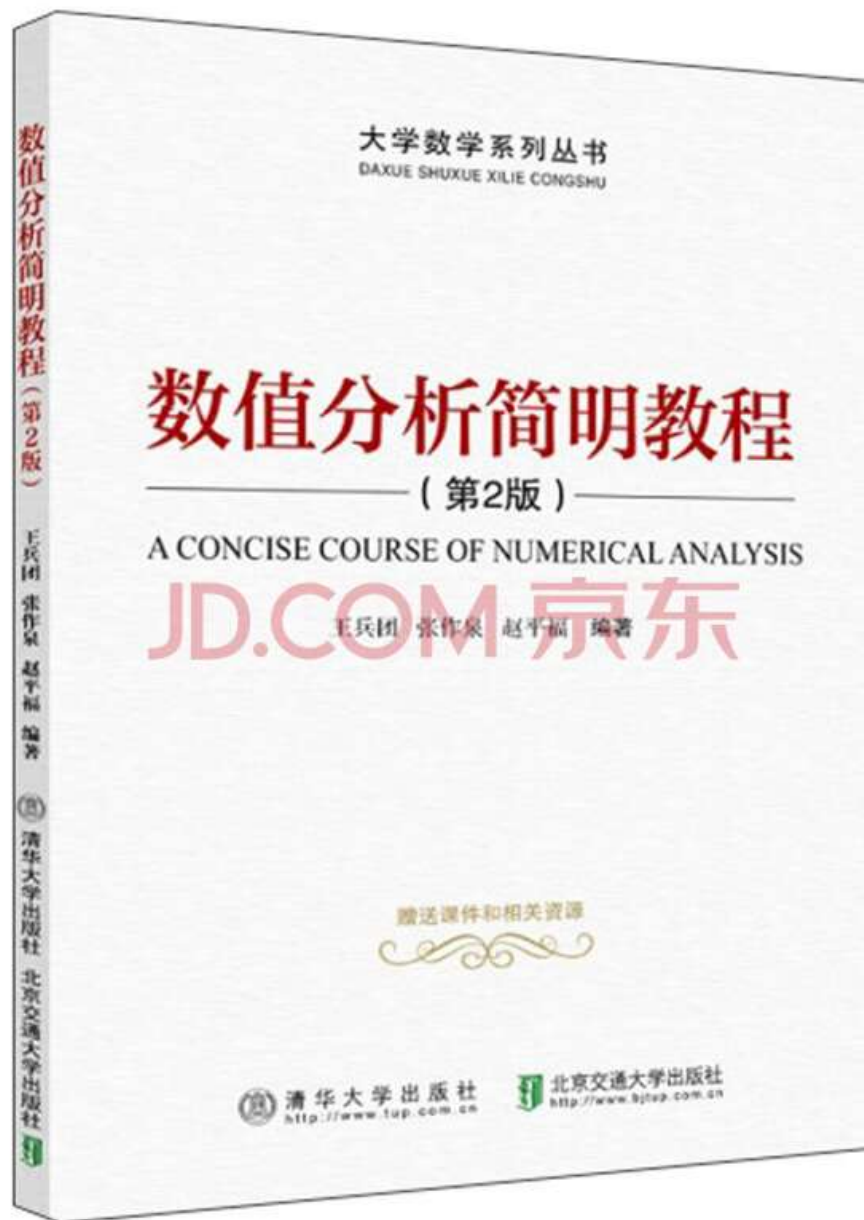
大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

## 授课教师基本信息

- 姓 名：郑晓朋
- 办公地点：信息楼319A室
- EMAIL: zhengxp2017@qq.com

# 主 讲 教 材





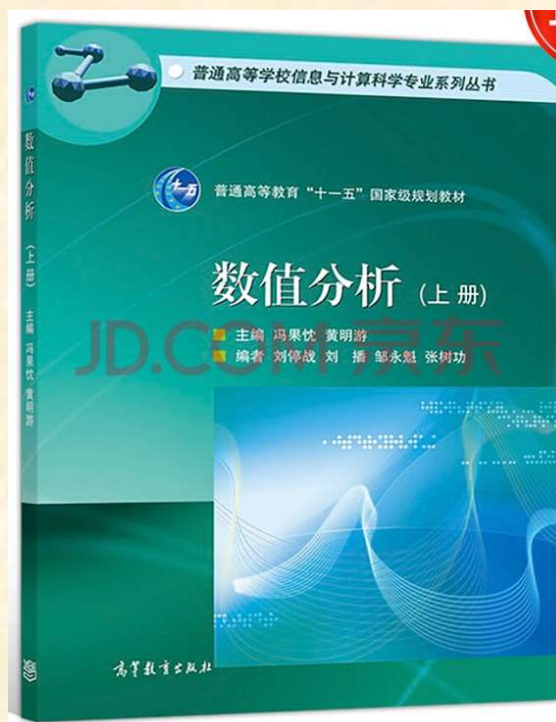
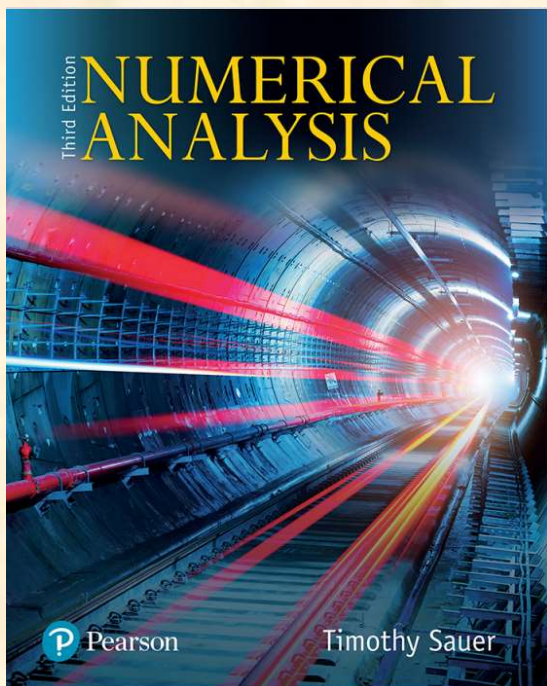
DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY



## 参考书目 (Reference)







DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

## 考核要求

课程的总成绩 { 平时作业 → 占20%;  
数值实验 → 占10%;  
期末考试 → 占70%;



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 绪论

## 什么是数值分析?

# 数值分析

*Numerical Analysis*

科学计算的桥梁！

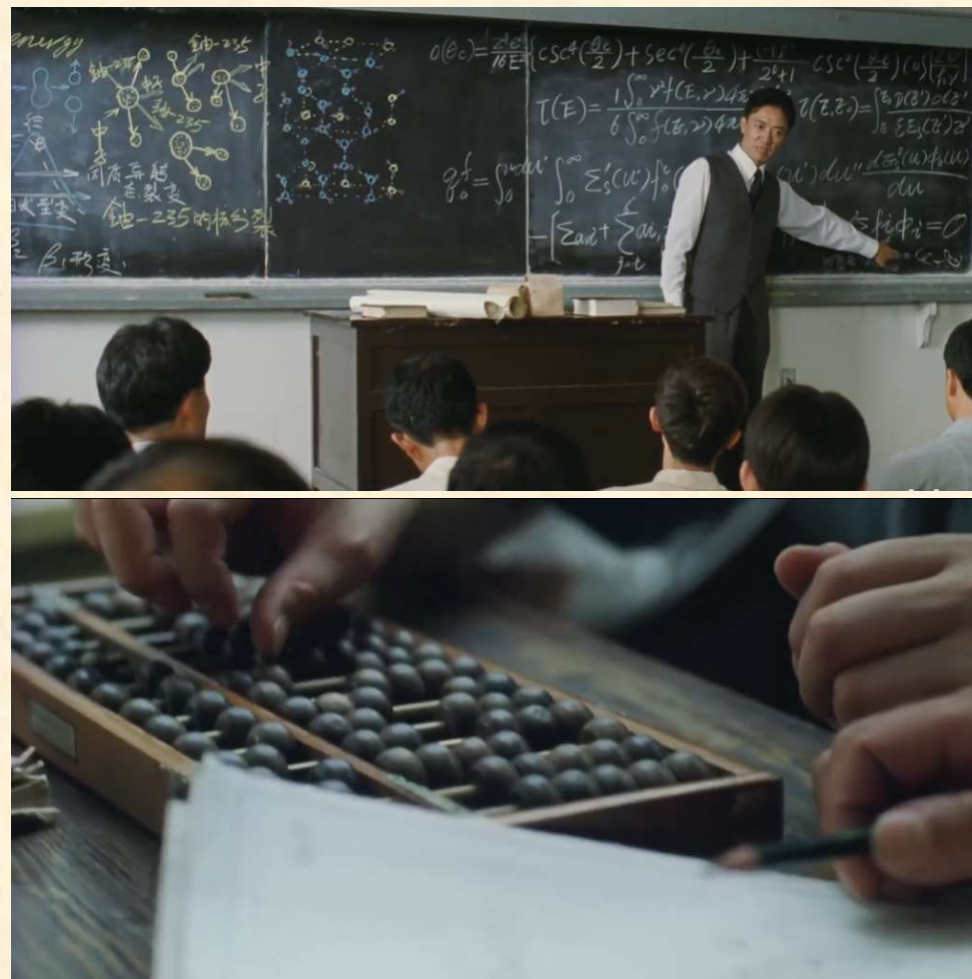
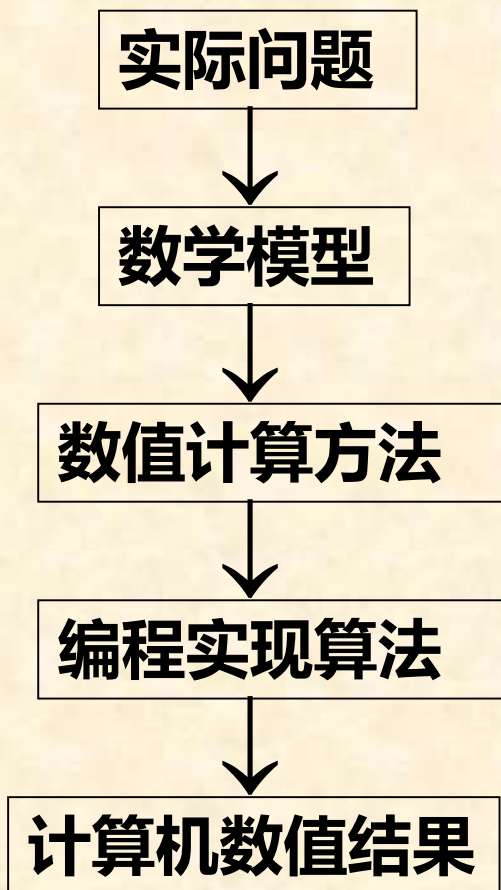
**科学计算**是指利用计算机再现、预测和发现客观世界运动规律和演化特征的全过程，具体是指应用计算机处理科学研究和工程技术中所遇到的数学计算。

**数值分析**主要研究使用计算机求解各种数学问题的方法、理论分析及其程序实现，是科学计算的重要理论支撑。它既有纯粹数学的高度抽象性和严密科学性，又有着具体应用的广泛性和实际实验的技术性，是一门与计算机使用密切结合的实用性很强的数学课程。





# 计算机科学计算的流程图

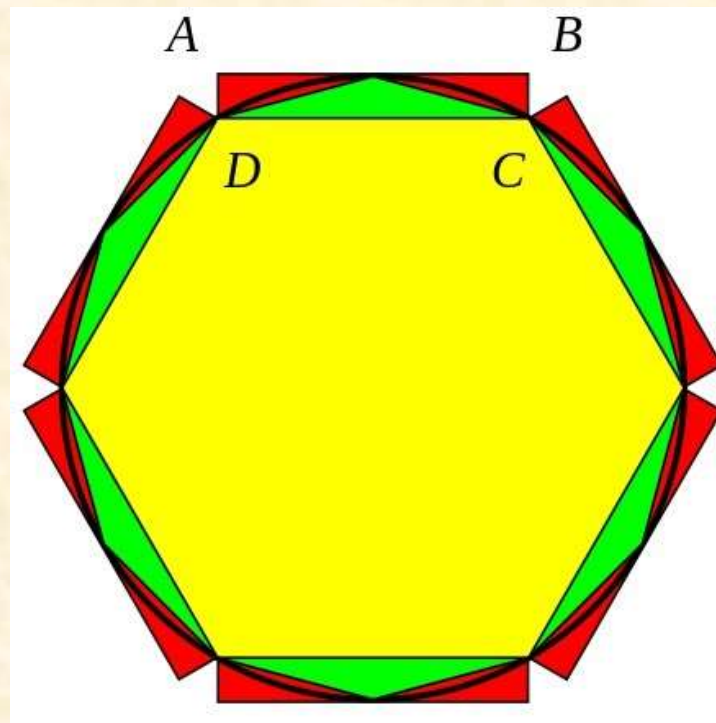
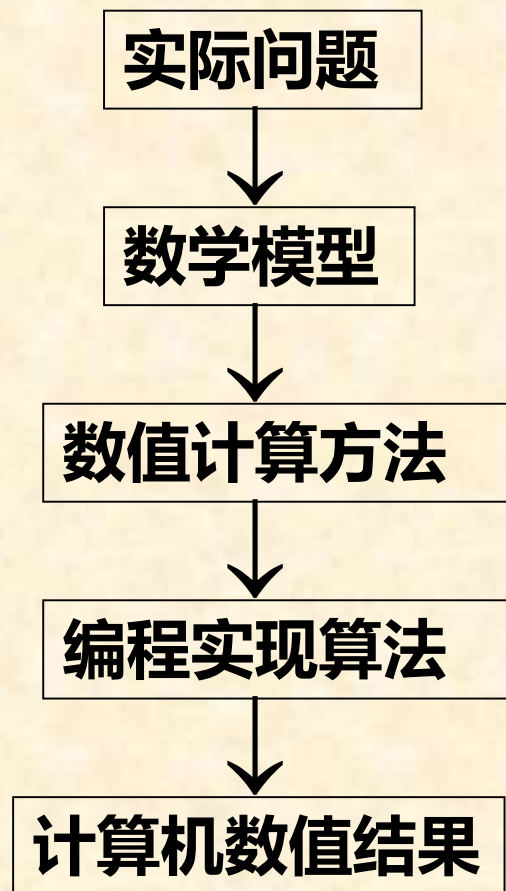


【李幼斌、李雪健】横空出世【1999】



1. 《周髀算经》(约公元前 200 年)中的 "径一而周三" 的记载 ( $\pi \approx 3$ ).
2. 汉朝张衡的  $\pi \approx \sqrt{10}$ .

计算机科学计算的流程图



数值  
3.14159265...

3. 公元 263 年刘徽的 "割圆术" 算法 (一种基于正多边形的迭代算法) 近似得到  $\frac{3927}{1250} \approx 3.1416$ .

4. 约公元 480 年, 祖冲之利用同样的方法, 确定了圆周率  $3.1415926 < \pi < 3.1415927$



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY



考察，线性方程组的解法

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

早在18世纪**Cramer**已给出了求解法则：

Cramer's Ruler

# Cramer's Ruler

$$x_i = \frac{D_i}{D} \quad i = 1, 2, \dots, n, \quad (D \neq 0)$$

$$D = \det(A) = \begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{vmatrix}$$

$$D_i = \det(A_i) = \begin{vmatrix} a_{11} & \cdots & b_1 & \cdots & a_{1n} \\ a_{21} & \cdots & b_2 & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & \cdots & b_n & \cdots & a_{nn} \end{vmatrix}$$

第  
 $i$   
列

这一结果理论上是非常漂亮的，它把线性方程组的求解问题归结果为计算 $n+1$ 个 $n$ 阶行列式问题。

对于行列式的计算，理论上有着著名的**Laplace**展开定理。  
这样理论上我们就有了一种非常漂亮的求解线性方程组的方法。

$$D = \det(A) = a_{i1}A_{i1} + a_{i2}A_{i2} + \cdots + a_{in}A_{in}$$

其中 $A_{ij}$ 表示元素 $a_{ij}$ 的代数余子式。

然而我们做一简单的计算就会发现，由于这一方法的运算量大得惊人，以至于完全不能用于实际计算。

设计算 $k$ 阶行列式所需要的乘法运算的次数为 $m_k$ ，则容易推出

$$m_k = k + k m_{k-1}$$

于是，我们有

$$\begin{aligned} m_n &= n + n m_{n-1} = n + n \left[ (n-1) + (n-1) m_{n-2} \right] \\ &= n + n(n-1) + n(n-1)(n-2) + \cdots + n(n-1) \cdots 3 \cdot 2 \\ &> n! \end{aligned}$$



这样，利用**Cramer**法则和**Laplace**展开定理来求解一个 $n$ 阶线性方程组，所需的乘法运算次数就大于

$$(n+1) n! = (n+1) !$$

以求解**25**阶线性方程组为例，如果用**Cramer**法则求解，在算法中运用行列展开计算，则总的乘法运算次数将达：

$$26! = 4.0329 \times 10^{26} \text{ (次)}$$

若使用每秒千亿次浮点乘法运算的串行计算机计算，一年可进行的运算应为：

$$365(\text{天}) \times 24(\text{小时}) \times 3600(\text{秒}) \times 10^{11} \approx 3.1536 \times 10^{18} \text{ (次)}$$

共需要耗费时间为：

$$(4.0329 \times 10^{26}) \div (3.1536 \times 10^{18}) \approx 1.2788 \times 10^9 \approx 13(\text{亿年})$$

它远远超出目前所了解的人类文明历史！



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

**Cramer** 算法是“实际计算不了”的。

而著名的 **Gauss**消元法，它的计算过程已作根本改进，成为有效算法，使得可在不到一秒钟之内即可完成上述计算任务。

随着科学技术的发展，出现的数学问题也越来越多样化，有些问题用消去法求解达不到精度，甚至算不出结果，从而促使人们对消去法进行改进，又出现了主元消去法，大大提高了消去法的计算精度。

这就是研究数值方法的必要性。

将数列

$$I_n = \int_0^1 \frac{x^n}{x+5} dx$$

写成递推公式形式，并计算数列：  $I_1, I_2, \dots$

解：  $\because I_n = \int_0^1 \frac{x^n + 5x^{n-1} - 5x^{n-1}}{x+5} dx$

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

$$\therefore I_n = \frac{1}{n} - 5I_{n-1} \quad n = 1, 2, \dots \quad (1.1)$$

由  $I_0 = \int_0^1 \frac{1}{x+5} dx = \ln \frac{6}{5}$  和式(1.1)可以依次算出  $I_1, I_2, \dots$



但科学计算中要具体数值！若把  $I_0 = \ln \frac{6}{5}$

**取为准确到小数点后8位的近似值作并在字长为8的计算机上编程计算，可出现**

$$I_{12} = -0.32902110 \times 10^2$$

这显然是错误的！  $\left( I_n = \int_0^1 \frac{x^n}{x+5} dx \right)$





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

- 基础数学  $\leftrightarrow$  计算数学
- 理想的  $\leftrightarrow$  落地使用的
- 实数  $\leftrightarrow$  计算过程中使用的数

数值分析是一门帮助我们将数学用起来，用好的一门学科



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

## 本课程关注数值计算的可行性、精度、效率和稳健性

- 构造计算机可行的**有效算法**
- 给出可靠的理论分析，即对任意逼近达到精度要求，保证数值算法的**收敛性**和**数值稳定性**，并可进行**误差分析**。
- 有好的计算复杂性，既要**时间复杂性**好，是指节省时间，又要**空间复杂性**好，是指节省存储量，这也是建立算法要研究的问题，它关系到算法能否在计算机上实现。
- **数值实验**，即任何一个算法除了从理论上要满足上述三点外，还要通过数值试验证明是行之有效的。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 机器数系及其运算

# 1.计算机的数系

## •数学中的实数（10进制）

$$x = \pm 10^c \times 0.a_1a_2a_3 \cdots$$

其中  $a_i \in \{0, 1, 2, \cdots, 9\}$ ,  $c$  为整数。  $x$  称为十进制浮点数。

---

## • $\beta$ 进制的浮点数

$$x = \pm \beta^c \times 0.a_1a_2a_3 \cdots a_t \cdots,$$

$$a_i \in \{0, 1, 2, \cdots, \beta - 1\}$$



## •计算机中的实数

$$x = \pm \beta^c \times 0.a_1a_2a_3 \cdots a_t \quad \text{其中 } a_i \in \{0, 1, \dots, \beta-1\}$$

**比较**  $x = \pm \beta^c \times 0.a_1a_2a_3 \cdots a_t \cdots,$

t: **字长**, 正整数;

$\beta$ : **进制**, 一般取为2,8,10和16;

$c$ : **阶码**, 整数, 满足 $L \leq c \leq U$ , L和U为固定整数;

## •机器数系

$$F(\beta, t, L, U) = \left\{ \pm \beta^c \times 0.a_1a_2a_3 \cdots a_t \mid a_k \in \{0, 1, \dots, \beta-1\}, L \leq c \leq U \right\}$$

# 机器数系的特点

- 1、是有限的离散集；
  - 2、有绝对值最大非零数( $M$ )和最小非零数( $m$ )；
  - 3、数绝对值大于 $M$ ，计算机产生上溢错误，绝对值小于 $m$ ，则计算机产生下溢错误；
  - 4、上溢时，计算机中断程序处理；下溢时，用零表示该数并继续执行程序；
- 无论是上溢，还是下溢，都称为溢出错误。
- 5、计算机把尾数为0且阶数最小的数表示数零。

## 2.计算机对数的接收与处理

### •计算机对数 $x$ 的接收

- 1) 若  $x \in F(\beta, t, L, U)$  , 原样接收  $x$  ;
- 2) 若  $x \notin F(\beta, t, L, U)$ , 但  $m \leq |x| \leq M$  , 则用  $F(\beta, t, L, U)$  中最属于靠近  $x$  的数  $fl(x)$  表示并记录  $x$  。

### •计算机对数的运算处理

- 1) 加减法: 先对阶, 后运算, 再舍入;
- 2) 乘除法: 先运算, 再舍入。

例，某计算机数系 $F(10, 4, -99, 99)$ 的两个数

$$x_1 = 0.2337 \times 10^{-1}, x_2 = 0.3364 \times 10^2$$

则运算过程如下：

$$\underline{fl(x_1 + x_2)} = \underline{fl(0.2337 \times 10^{-1} + 0.3364 \times 10^2)}$$

对阶

$$= \underline{fl(0.0002337 \times 10^2 + 0.3364 \times 10^2)}$$

运算

$$= \underline{fl(0.3366337 \times 10^2)}$$

舍入

$$= 0.3366 \times 10^2$$





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 误差

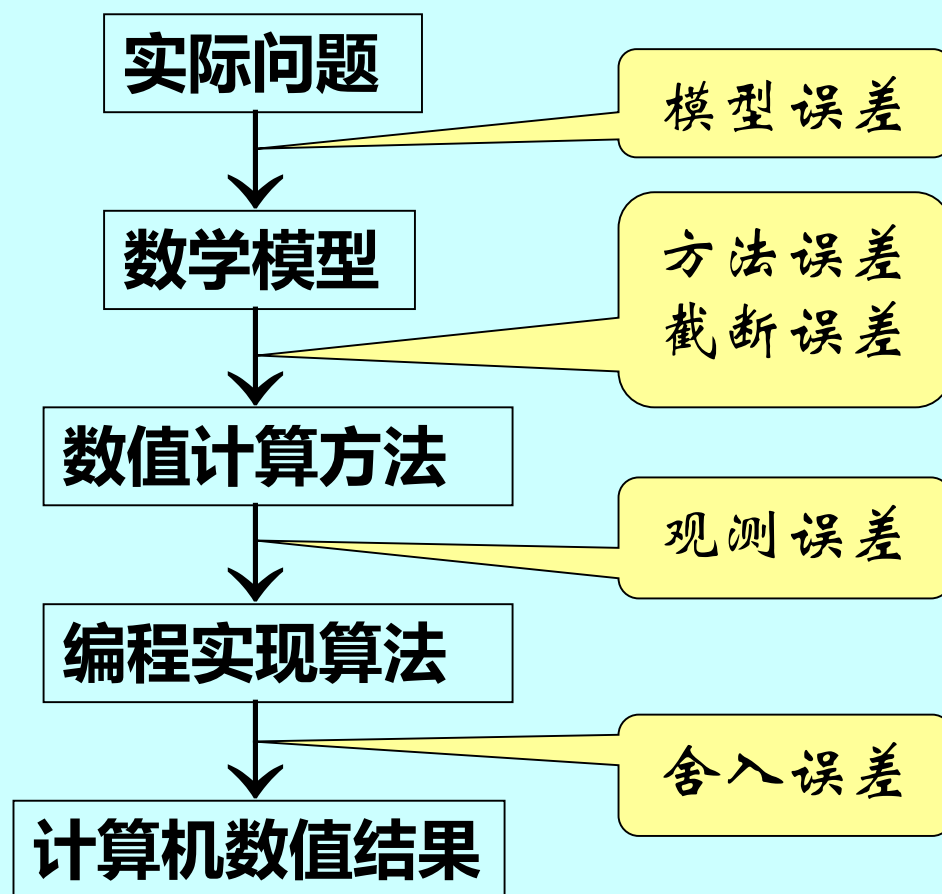
准确值与近似值的差异就是误差！



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY



## 误差来源

- 连续问题离散化
- 以有限代替无限
- 数值近似



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

1. **模型误差** 由实际问题抽象出数学模型，要简化许多条件，这就不可避免地要产生误差。实际问题的解与数学模型的解之间的误差

2. **截断误差** 从数学问题转化为数值问题的算法时所产生的误差，如用有限代替无限的过程所产生的误差

**截断误差**通常是指用一个基本表达式替换一个相当复杂的算术表达式时所引起的误差。这一术语从用截断**Taylor**级数替换一个复杂的算术表达式的技术中衍生而来。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

例如，给定  $x$  求  $e^{x^2}$  的值的运算，我们可用无穷级数：

$$e^{x^2} = 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \cdots + \frac{x^{2n}}{n!} + \frac{x^{2(n+1)}}{(n+1)!} + \cdots$$

我们可用它的前  $n+1$  项和

$$s(x) = 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \cdots + \frac{x^{2n}}{n!}$$

近似代替函数  $e^{x^2}$ ，则数值方法的误差是

$$R_n(x) = e^{x^2} - s(x) = \frac{e^{(\theta x)^2}}{(n+1)!} x^{2(n+1)}, \quad 0 < \theta < 1$$

截断误差



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

3. **观测误差** 初始数据大多数是由观测而得到的。由于观测手段的限制，得到的数据必然有误差

4. **舍入误差** 以计算机为工具进行数值运算时，由于计算机的字长有限，原始数据在计算机上的表示往往会有误差，在计算过程中也可能产生误差

例如，用1.4142近似代替 $\sqrt{2}$ ，产生的误差

$$\begin{aligned} E &= \sqrt{2} - 1.4142 = 1.4142135\cdots - 1.4142 \\ &= 0.0000135\cdots \end{aligned}$$

就是**舍入误差**。





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

模型和观测两种误差不在本课程的讨论范围

这里主要讨论算法的截断误差与舍入误差，而截断误差将结合具体算法讨论

分析初始数据的误差通常也归结为舍入误差

研究计算结果的误差是否满足精度要求就是：

**误差估计问题**

## 2、误差的定义（数学定义）

**定义1** 设 $x$ 是准确值， $x^*$ 是 $x$ 的一个近似值，称差 $x^* - x$ 为 $x^*$ 的**绝对误差**，简称**误差**，记为 $e^*$ 或 $e(x^*)$ ，即：

$$e(x^*) = x^* - x$$

---

**定义2** 称满足

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

的正数 $\varepsilon^*$ 为近似值 $x^*$ 的误差限。

$$x^* - \varepsilon^* \leq x \leq x^* + \varepsilon^* \quad \text{简记 } x = x^* \pm \varepsilon^*$$



例如，用毫米刻度的米尺测量一长度 $x$ ，读出和该长度接近的刻度 $a$ ， $a$ 是 $x$ 的近似值，它的误差界是 $0.5mm$ ，于是有

$$|x - a| \leq 0.5mm$$

如若读出的长度为 $765mm$ ，则有，

$$|x - 765| \leq 0.5$$

绝对误差界

虽然从这个不等式不能知道准确的 $x$ 是多少，但可知

$$764.5 \leq x \leq 765.5,$$

结果说明 $x$ 在区间 $[764.5, 765.5]$ 内。

对于一般情形  $|x - a| \leq e_a$  既可以表示为

$$a - e_a \leq x \leq a + e_a,$$

也可以表示为

$$x = a \pm e_a。$$

但要注意<sup>Ⓢ</sup>的是，误差的大小并不能完全表示近似值的好坏。

**定义3** 设  $x$  是准确值,  $x^*$  是  $x$  的近似值, 称

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

为近似值  $x^*$  的相对误差, 记为

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

**重要结论!**

**相对误差绝对值越小, 近似程度越高。**

例

有两个量  $x=3.000$ ,  $a=3.100$ , 则其绝对误差:

$$x - a = -0.1$$

绝对误差

其相对误差为:

$$\frac{x - a}{x} = \frac{-0.1}{3.00} = -0.333 \times 10^{-1},$$

相对误差

又有两个量  $x=300.0$ ,  $a=310.0$ , 则其绝对误差:

$$x - a = -0.1 \times 10^2,$$

绝对误差

其相对误差为:

$$\frac{x - a}{x} = \frac{-0.1 \times 10^2}{0.3 \times 10^4} = -0.333 \times 10^{-1}$$

相对误差



上例说明绝对误差有较大变化，相对误差相同。作为精确值的度量，绝对误差可能会引起误会，而相对误差由于考虑到准确值本身的大小而更有意义。



## 定义4 满足

$$\left| e_r^* \right| = \left| \frac{x^* - x}{x} \right| \leq \varepsilon_r^*$$

的正数  $\varepsilon_r^*$  称为  $x^*$  的相对误差限。

$$\text{实用中, } \varepsilon_r^* = \frac{\varepsilon^*}{|x^*|}$$

## 例

已知  $e = 2.71828182\cdots$  其近似值  $a = 2.718$  , 求  $a$  的绝对误差界和相对误差界。

解:  $e - a = 0.00028182\cdots$  , 因此其绝对误差界为:

$$|e - a| \leq 0.0003$$

相对误差界为:

$$\frac{|e - a|}{|a|} = \frac{0.0003}{2.718} \approx 0.0001110375 \leq 0.0002。$$

此例计算中不难发现, 绝对误差界和相对误差界并不是唯一的。 我们要注意它们的作用。

## 误差界的取法

当准确值 $x$ 位数比较多时，人们常常按四舍五入的原则得到 $x$ 的前几位近似值 $a$ ，例如

$$x = \pi = 3.14159265 \dots$$

取3位：  $a_1 = 3.14$ ,  $\pi - a_1 = 0.00159265 \dots$

取5位：  $a_2 = 3.1416$ ,  $\pi - a_2 = -0.00000735 \dots$

那么，它们的误差界的取法应为：

$$|\pi - 3.14| \leq \frac{1}{2} \times 10^{-2}, \quad |\pi - 3.1416| \leq \frac{1}{2} \times 10^{-4}.$$

# 有效数字



科学计算中常用有效数字来估计和处理误差，有效数字易算且与误差有密切关系。

有效数字是误差的定量化。



# 有效数字的定义

**定义5：**若近似数  $x^*$  的误差限是其某一位上数字的半个单位，就说近似数  $x^*$  准确到该位；由该位自右向左数到  $x^*$  的第一个非零数字若有  $n$  位，就称近似数  $x^*$  有  $n$  位有效数字。

---

$$x^* = 216.503$$

$$= 2 \times 10^2 + 1 \times 10^1 + 6 \times 10^0 + 5 \times 10^{-1} + 0 \times 10^{-2} + 3 \times 10^{-3}$$

2	1	6	□	5	9	3
$10^2$	$10^1$	$10^0$		$10^{-1}$	$10^{-2}$	$10^{-3}$

$$\left| e(x^*) \right| \leq \varepsilon = \frac{1}{2} \times 10^{-3} \Rightarrow n = 6 \quad \left| e(x^*) \right| \leq \varepsilon = \frac{1}{2} \times 10^1 \Rightarrow n = 2$$



# 有效数字的数学描述

设  $x^* = \pm 0.a_1a_2 \cdots a_k \times 10^m$

$$a_1 \neq 0, a_l \in \{0, 1, 2, \dots, 9\}, k \geq n, m \in \mathbb{Z}$$

如果成立

$$\left| e(x^*) \right| = \left| x^* - x \right| \leq 0.5 \times 10^{m-n}$$

则称近似数 $x^*$ 有 $n$ 位有效数字，此时可写

$$x^* = \pm 0.a_1a_2 \cdots a_n \times 10^m$$

**有效数字越多，绝对误差和相对误差就越小，因此近似数就越准确！**

### 例3 求圆周率 $\pi$ 的两个近似值

$$x_1 = 3.14, \quad x_2 = 3.141$$

**的有效数字。**

解： $\because \pi = 3.1415926\cdots, x_1 = 0.314 \times 10^1, m = 1$

$$|\pi - x_1| = 0.015926\cdots$$

$$= 10^{-2} \times 0.15926\cdots < 0.5 \times 10^{-2}$$

**有 $m-n = -2$ ，得 $n = 3$ ， $x_1$ 有3位有效数字；**

$$|\pi - x_2| = 0.0005926\cdots \quad (x_2 = 0.3141 \times 10^1, m = 1)$$

$$= 10^{-3} \times 0.5926\cdots < 0.5 \times 10^{-2}$$

**有 $m-n = -2$ ，得 $n = 3$ ， $x_2$ 有3位有效数字；**

可以证明：

果十进制  $x$  经过四舍五入得到近似数  $x^*$ ，  
则  $x^*$  的有效数字位为将  $x^*$  写为规格化浮点数  
后的尾数的位数。

例 0.00345 四舍五入得  $0.0035 = 0.35 \times 10^{-2}$  可知近似数有2位有效数字。

$\pi$  的近似数 3.14（四舍五入数）有3位有效数字；  
3.142（四舍五入数）有4位有效数字。  
而近似数 3.141（不是四舍五入数）有3位有效数字，虽然它有4位数。

# 一般来说，绝对误差与小数位数有关， 相对误差与有效数字位数有关

有两个量  $x=3.000$ ,  $a=3.100$ , 则其绝对误差:

$$x - a = -0.1$$

绝对误差

其相对误差为:

$$\frac{x - a}{x} = \frac{-0.1}{3.00} = -0.333 \times 10^{-1},$$

相对误差

又有两个量  $x=300.0$ ,  $a=310.0$ , 则其绝对误差:

$$x - a = -0.1 \times 10^2,$$

绝对误差

其相对误差为:

$$\frac{x - a}{x} = \frac{-0.1 \times 10^2}{0.3 \times 10^4} = -0.333 \times 10^{-1}$$

相对误差



**例4 已知近似数 $x^*$ 有5位有效数字，试求其相对误差限。**

**解： 因为  $x^*$ 有5位有效数字，可以设**

$$x^* = \pm 0.a_1 a_2 \cdots a_5 \times 10^m, a_1 \geq 1$$

**由题意，有 $n=5$ 和  $|x^* - x| \leq 0.5 \times 10^{m-5}$**

$$\begin{aligned} \therefore \frac{|x^* - x|}{|x^*|} &\leq \frac{0.5 \times 10^{m-5}}{0.a_1 a_2 \cdots a_5 \times 10^m} \leq \frac{5 \times 10^{-5}}{a_1} \\ &\leq \frac{1}{2a_1} \times 10^{-4} < \frac{1}{2} \times 10^{-4} \Rightarrow \underline{\underline{\varepsilon_r^* = \frac{1}{2} \times 10^{-4}}} \end{aligned}$$

# 有效数字与相对误差的关系

**定理3 设10进制近似数**

$$x^* = \pm 0.a_1 a_2 \cdots a_k \times 10^m, a_1 \neq 0, k \geq n$$

**1) 若 $x^*$ 有 $n$ 位有效数字, 则有**

$$\left| e_r(x^*) \right| = \frac{|x^* - x|}{|x^*|} \leq \frac{1}{2a_1} \times 10^{1-n}$$

**2) 若 $x^*$ 的相对误差**

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

**则 $x^*$ 有 $n$ 位有效数字。**



证

下面用 $a$ 表示 $x^*$ , 由有效数字定义可得到

$$a_1 \times 10^{k-1} \leq |a| \leq (a_1 + 1) \times 10^{k-1} \quad (1-6)$$

所以如果 $a$ 有 $n$ 位有效数字, 那么

$$\frac{|x-a|}{|a|} = |x-a| \times \frac{1}{|a|} \leq \frac{1}{2} \times 10^{k-n} \times \frac{1}{a_1 \times 10^{k-1}} = \frac{1}{2a_1} \times 10^{1-n},$$

结论 1) 成立。 针对结论2), 结合 (1-6) 有

$$\frac{|x-a|}{|a|} \leq \frac{1}{2(a_1+1)} \times 10^{1-n} \quad |x-a| \leq \frac{(a_1+1) \times 10^{k-1}}{2 \times (a_1+1)} \times 10^{1-n} = \frac{1}{2} \times 10^{k-n},$$

由定义知,  $a$ 具有 $n$ 位有效数字。

**例5 为保证某算式的计算精度，要求其中的  $\sqrt[3]{23}$  的近似值  $x^*$  的相对误差小于0.1%，请确定  $x^*$  至少要取几位有效数字才能达到要求。**

解：  $\because 2 < \sqrt[3]{23} < 3$

$$\Rightarrow \sqrt[3]{23} = 2.a_2a_3\cdots = 0.2a_2a_3\cdots \times 10^1 \Rightarrow a_1 = 2$$

**假设  $x^*$  至少要取  $n$  位有效数字才能保证相对误差小于0.1%，由定理3，选择**

$$\frac{1}{2a_1} \times 10^{1-n} = \frac{1}{2 \times 2} \times 10^{1-n} < 0.1\% \Rightarrow \underline{n \geq 4}$$

$$\left| e_r(x^*) \right| = \frac{|x^* - x|}{|x^*|} \leq \frac{1}{2a_1} \times 10^{1-n}$$

### 3、数值计算的误差

**定理1** 假设 $x^*$ 和 $y^*$ 分别是准确值 $x$ 和 $y$ 的一个近似值，则有四则运算的绝对误差估计：

$$(1) e(x^* \pm y^*) = e(x^*) \pm e(y^*)$$

$$(2) e(x^* \cdot y^*) \approx y^* e(x^*) + x^* e(y^*)$$

$$(3) e\left(\frac{x^*}{y^*}\right) \approx \frac{y^* e(x^*) - x^* e(y^*)}{(y^*)^2}$$



## 证明 (2)

$$e(x^* y^*) = x^* y^* - xy$$

$$= x^* y^* - \underline{xy^*} + \underline{xy^*} - xy$$

$$= y^* (x^* - x) + x(y^* - y)$$

$$= y^* e(x^*) + x e(y^*)$$

$$\approx y^* e(x^*) + x^* e(y^*) \quad (\because x \approx x^*)$$

- 微分与误差的对应

$$e^*(x) = x^* - x = dx$$

$$e_r(x^*) = \frac{x^* - x}{x} = \frac{dx}{x} = d \ln x$$

- 绝对误差和相对误差与微分的关系

$$(1) \quad dx = e(x^*)$$

$$(2) \quad d \ln x = e_r(x^*)$$

## 例2 求函数 $y = x^n$ 与自变量 $x$ 的相对误差关系。

解：取对数有

$$\ln y = n \ln x$$

取微分有

$$d \ln y = n d \ln x$$

$$\therefore e_r \left( (x^*)^n \right) = n e_r (x^*)$$



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 数值方法的稳定性



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

## 数值方法的稳定性

用某一种数值方法求一个问题的数值解，如果在方法的计算过程中舍入误差在一定条件下能够得到控制（或者说舍入误差的增长不影响产生可靠的结果），则称该方法是数值稳定的；否则，出现与数值稳定相反的情况，则称之为数值不稳定的。



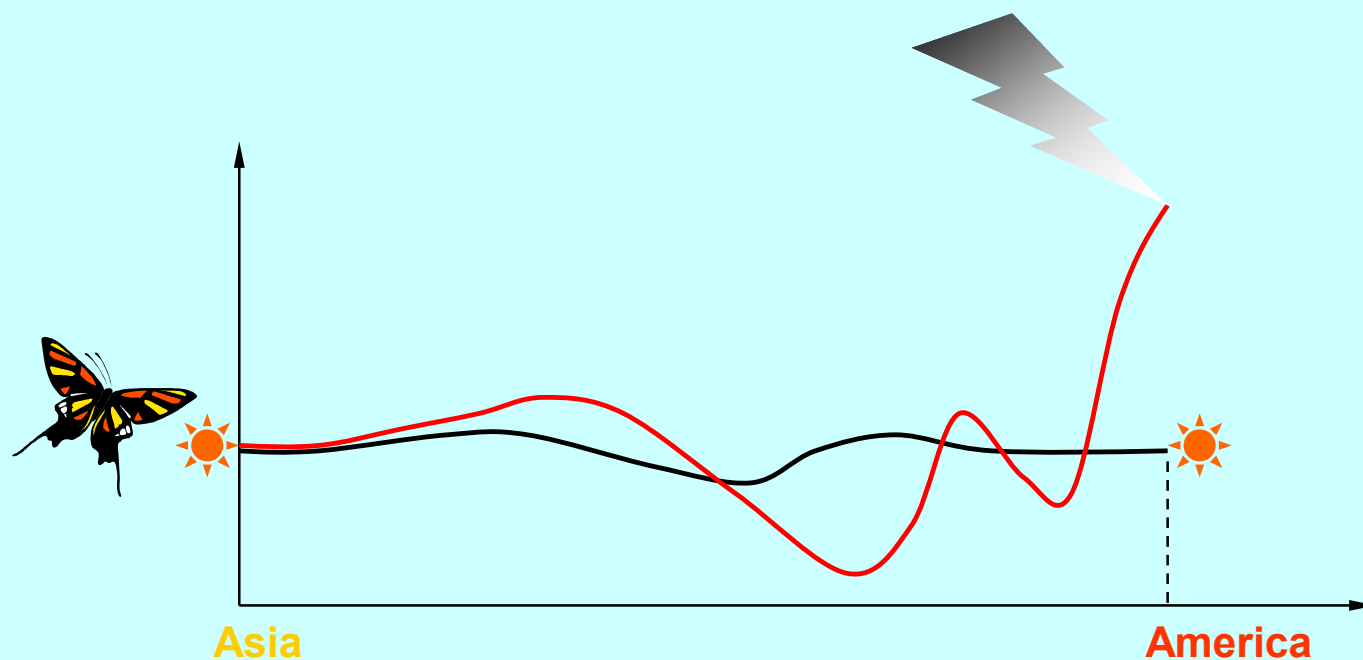


DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

蝴蝶效应——亚洲蝴蝶拍拍翅膀，将使风和日丽的美洲几个月后出现狂风暴雨？！





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

### 什么是蝴蝶效应？

美国麻省理工学院气象学家洛伦兹 (Lorenz) 为了预报天气，他用计算机求解仿真地球大气的13个方程式。为了更细致地考察结果，他把一个中间解取出，提高精度再送回。而当他喝了杯咖啡以后回来再看时竟大吃一惊：本来很小的差异，结果却偏离了十万八千里！计算机没有毛病，于是，洛伦兹 (Lorenz) 认定，他发现了新的现象：“对初始值的极端不稳定性”，即：“混沌”，又称“蝴蝶效应”





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

初始数据的微小变化，导致计算结果的**剧烈变化**的问题称为**病态问题**。

例如

$$\begin{cases} x_1 + \frac{1}{2}x_2 + \frac{1}{3}x_3 = \frac{11}{6} \\ \frac{1}{2}x_1 + \frac{1}{3}x_2 + \frac{1}{4}x_3 = \frac{13}{12} \\ \frac{1}{3}x_1 + \frac{1}{4}x_2 + \frac{1}{5}x_3 = \frac{47}{60} \end{cases}$$

$$x_1 = x_2 = x_3 = 1$$

$\Downarrow$

$$x_1 = -6.2, x_2 = 38, x_3 = -34$$

扰动  $\rightarrow$

$$\begin{cases} x_1 + 0.50x_2 + 0.33x_3 = 1.8 \\ 0.50x_1 + 0.33x_2 + 0.25x_3 = 1.1 \\ 0.33x_1 + 0.25x_2 + 0.20x_3 = 0.78 \end{cases}$$







DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

**初始数据的微小变化只引起计算结果的微小变化的计算问题称为良态问题。**

$$\begin{cases} 2x_1 - x_2 = 6 \\ x_1 + 2x_2 = -2 \end{cases} \xrightarrow{\text{扰动}} \begin{cases} 2x_1 - x_2 = 6 \\ x_1 + 2x_2 = -2.005 \end{cases}$$

$$x_1 = 2, x_2 = -2 \Rightarrow x_1 = 1.999, x_2 = -2.002$$

**病态问题的计算或求解应使用专门的方法或将其转化为非病态问题来解决。**

**数值分析主要研究良态问题数值解法。**



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

例6

计算积分  $I_n = \int_0^1 \frac{x^n}{x+5} dx, n = 0, 1, 2, \dots, 7$

解： 由于

$$\begin{aligned} I_n + 5I_{n-1} &= \int_0^1 \frac{x^n}{x+5} dx + 5 \int_0^1 \frac{x^{n-1}}{x+5} dx = \int_0^1 \frac{x^n + 5x^{n-1}}{x+5} dx \\ &= \int_0^1 x^{n-1} dx = \frac{1}{n} \end{aligned}$$

则递归算法如下：

1.  $I_n = \frac{1}{n} - 5I_{n-1}$  , 由  $I_0 = \ln \frac{6}{5}$  计算出  $I_1, \dots, I_7$
2.  $I_{n-1} = \frac{1}{5} \left( \frac{1}{n} - I_n \right)$  , 由  $I_7 = 0.0210$  计算出  $I_7, \dots, I_0$





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

$n$	$I_n$	方法1	方法2
0	0.1820	0.1820	0.1820
1	0.0880	0.0900	0.0880
2	0.0580	0.0500	0.0580
3	0.0431	0.0830	0.0431
4	0.0343	-0.0165 ?	0.0343
5	0.0284	1.0250 ??	0.0284
6	0.0240	-4.9580 ?!	0.0240
7	0.0210	24.933 !!	0.0210

What  
happened  
?!





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

设  $I_0$  的近似值为  $\bar{I}_0$ ，然后按方法1计算  $I_1, \dots, I_7$  的近似值  $\bar{I}_1, \dots, \bar{I}_7$ ，如果最初计算时误差为： $E_0 = I_0 - \bar{I}_0$  递推过程的舍入误差不记，并记  $E_n = I_n - \bar{I}_n$ ，则有

$$E_7 = I_7 - \bar{I}_7 = (-5)E_6 = (-5) \cdot (-5)E_5 = \dots = (-5)^7 E_0$$

由此可见，用该方法计算  $I_1, \dots, I_7$  时，当计算  $I_0$  时产生的舍入误差为  $E_0$ ，那么计算  $I_7$  时产生的舍入误差放大了  $5^7 = 78,125$  倍，因此，该方法是数值不稳定的。

按方法2计算时，记初始误差为  $E_7 = I_7 - \bar{I}_7$ ，则有

$$E_0 = I_0 - \bar{I}_0 = \left(-\frac{1}{5}\right)E_1 = \left(-\frac{1}{5}\right) \cdot \left(-\frac{1}{5}\right)E_2 = \dots = \left(-\frac{1}{5}\right)^7 E_7$$

**由此可知，使用公式2计算时不会放大舍入误差。**

因此，该方法是数值稳定的。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

# 避免误差危害的基本原则



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

为了用数值方法求得数值问题满意的近似解，在数值运算中应注意下面两个基本原则。

## (I) 避免有效数字的损失

在四则运算中为避免有效数值的损失，应注意以下事项：

- (1) 在做加法运算时，应防止“大数吃小数”；
- (2) 避免两个相近数相减；
- (3) 避免小数做除数或大数做乘数。





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

**例5** 在五位十进制的计算机上计算  $x = 63015 + \sum_{i=1}^{1000} \delta_i$ ,  $\delta_i = 0.4$

解 计算机作加减法时, 先将所相加数阶码对齐, 根据字长舍入, 再加减。如果用63015依次加各个 $\delta_i$ , **那么上式用规范化和阶码对齐后的数表示为:**

$$x = 0.63015 \times 10^5 + \underbrace{0.000004 \times 10^5 + \cdots + 0.000004 \times 10^5}_{1000 \uparrow}$$

因其中  $0.000004 \times 10^5$  的舍入结果为0, 所以上式的计算结果是

**$0.63015 \times 10^5$** 。这种现象被称为“大数吃小数”。如果改变运算

次序, 先把1000个  $\delta_i$  相加, 再和63015相加, 即

$$\begin{aligned} x &= \underbrace{0.4 + 0.4 + \cdots + 0.4}_{1000} + 0.63015 \times 10^5 = 0.4 \times 10^3 + 0.63015 \times 10^5 \\ &= 0.004 \times 10^5 + 0.63015 \times 10^5 = \mathbf{0.63415 \times 10^5} \end{aligned}$$

后一种方法的结果是正确的, 前一种方法的舍入误差影响太大。





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

例7

一元二次方程  $ax^2 + 2bx + c = 0$  ( $a \cdot b \cdot c \neq 0$ )

有两个根，其求根公式为

$$x_1 = \frac{-b + \sqrt{b^2 - ac}}{a}, \quad x_2 = \frac{-b - \sqrt{b^2 - ac}}{a}$$

如果  $b^2 \gg |ac|$ ，则  $\sqrt{b^2 - ac} \approx |b|$ ，用上述公式计算时

如果  $b > 0$ ，则有  $x_1 \approx \frac{-b + |b|}{a} = 0$

如果  $b < 0$ ，则有  $x_2 \approx \frac{-b - |b|}{a} = 0$

总之，两者其中之一必将会损失有效数字。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

解一般二次方程  $ax^2+2bx+c=0$ , ( $a, b, c$ 均不为零), 应取

$$x_1 = \frac{-b - \operatorname{sgn}(b)\sqrt{b^2 - ac}}{a}, \quad x_2 = \frac{c}{ax_1};$$

其中  $\operatorname{sgn}(b) = \begin{cases} 1, & \text{当 } b > 0 \text{ 时} \\ -1, & \text{当 } b < 0 \text{ 时} \end{cases}$  是  $b$  的符号函数。

如果  $b^2 \gg |ac|$ , 则  $\sqrt{b^2 - ac} \approx |b|$ , 用上述公式计算时

如果  $b > 0$ , 则有  $x_1 \approx \frac{-b - |b|\sqrt{b^2 - ac}}{a}$   $x_2 = \frac{c}{ax_1};$

如果  $b < 0$ , 则有  $x_1 \approx \frac{-b + |b|\sqrt{b^2 - ac}}{a}$



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

例

求方程  $x^2 - 16x + 1 = 0$  的根

由习惯的公式:  $x_1 = 8 + \sqrt{63}$ ,  $x_2 = 8 - \sqrt{63}$ 。

若取三位有效数字计算, 有  $\sqrt{63} \approx 7.94$ 。

$x_1 = 8 + \sqrt{63} \approx 8.00 + 7.94 = 15.9$ , 有三位有效数字。而

$x_2 = 8 - \sqrt{63} \approx 8.00 - 7.94 = 0.06$ , 只有一位有效数字。

其原因为在计算  $x_2$  时发生了两个相近数相减, 造成有效数字损失。

而  $x_2$  的精确值是  $0.062746\cdots$ 。如果改用公式:

$$x_2 = \frac{c}{ax_1} = \frac{1}{x_1}$$

计算得  $x_2 \approx 0.062746$ , 具有三位有效数字。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

## (II) 减少运算次数

例如，多项式求值运算，设  $p_n(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$

如果直接逐项求和计算，需要大约  $2n$  次乘法运算 即

$$x \cdot x \rightarrow x^2 \cdot x \rightarrow x^3 \cdot x \rightarrow \cdots \rightarrow x^{n-1} \cdot x \rightarrow x^n \quad \bullet \quad \bullet \quad \bullet \quad n-1 \text{次}$$

$$a_n \cdot x^n, a_{n-1} \cdot x^{n-1}, \cdots, a_1 \cdot x \quad \rightarrow n \text{次}$$

若取  $t_k = x^k$ ,  $u_k = a_0 + a_1 x + \cdots + a_k x^k$ ,

则有递推公式：

$$\begin{cases} t_k = x \cdot t_{k-1} \\ u_k = u_{k-1} + a_k \cdot t_k \end{cases} \quad k = 1, 2, \cdots, n, \quad \begin{cases} t_0 = 1 \\ u_0 = a_0 \end{cases}$$

$p_n(x) = u_n$  就是所求的值。总的计算量需进行  $2n$  次乘法。



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

若将公式变成如下递推公式，即令

$$\begin{aligned} p_n(x) &= (a_n x + a_{n-1}) x^{n-1} + \cdots + a_1 x + a_0 \\ &= ((a_n x + a_{n-1}) x + a_{n-2}) x^{n-2} + a_{n-3} x^{n-3} \cdots + a_1 x + a_0 \\ &= \cdots = \\ &= (\cdots (a_n x + a_{n-1}) x + a_{n-2}) x + \cdots + a_2) x + a_1) x + a_0 \end{aligned}$$

若令  $s_k = (\cdots (a_n x + a_{n-1}) x + a_{n-2}) x + \cdots + a_{k+1}) x + a_k$

则有递推公式：

$$\begin{cases} s_n = a_n \\ s_k = x \cdot s_{k+1} + a_k \end{cases} \quad k = n-1, n-2, \cdots, 2, 1, 0$$

$p_n(x) = s_0$  就是所求的值。总的计算量为  $n$  次乘法。





DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

$$P_5(x) = 5x^5 + 0x^4 + 3x^3 + 3x^2 + 1x - 1$$



$$P_5(x) = (((((5x + 0)x + 1)x - 3)x + 1)x - 1$$

$$5 \quad 0 \quad 1 \quad -3 \quad 1 \quad -1$$

令  $x=2$

---


$$10 \quad 20 \quad 42 \quad 39 \quad 138 \quad 157 = P_5(2)$$

以上计算过程称之为 秦九韶算法



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

**例10** 利用  $\ln(1+x) = \sum_{n=1}^{\infty} (-1)^{n+1} \frac{x^n}{n}$  计算  $\ln 2$ , 若要精确到  $10^{-5}$  要计算十万项的和, 计算量很大, 另一方面舍入误差的积累也十分严重。

如果改用级数

$$\ln \frac{1+x}{1-x} = 2 \left( x + \frac{x^3}{3} + \frac{x^5}{5} + \cdots + \frac{x^{2n+1}}{2n+1} + \cdots \right)$$

$$\text{取 } x = \frac{1}{3}, \quad \ln 2 = \ln \frac{1+\frac{1}{3}}{1-\frac{1}{3}} = 2 \left( \frac{1}{3} + \frac{\left(\frac{1}{3}\right)^3}{3} + \frac{\left(\frac{1}{3}\right)^5}{5} + \cdots + \frac{\left(\frac{1}{3}\right)^{2n+1}}{2n+1} + \cdots \right)$$

只须计算前9项的和, 截断误差便小于  $10^{-10}$



DUT

大连理工大学

DALIAN UNIVERSITY OF TECHNOLOGY

一个好的、有效的数值方法的评价标准

- (1) 运算次数少
- (2) 运算过程具有规律性（如递归性），便于编程
- (3) 需存储的中间结果少
- (4) 数值稳定性好（能控制误差的传播和积累）

其核心为：“快”、“准”





**DUT**

**大连理工大学**

**DALIAN UNIVERSITY OF TECHNOLOGY**

## 数值分析学习方法

- 注意掌握各种方法的基本原理
- 注意各种方法的构造手法
- 重视各种方法的误差分析
- 做一定量的习题
- 注意与实际问题的联系
- 了解各种方法的算法与程序实现