

计算方法

吉林大学计算机学院、软件学院
机器学习研究室
计算方法课程组

第1章 绪论

- 1.1 计算方法概述

- 科学计算与计算方法
- 数学模型与计算方法
- 计算方法的特点及学习方法

- 1.2 误差

- 计算机的浮点表示及算术运算
 - 误差来源
 - 误差的基本概念
 - 误差分析
-

应用举例1

例：一个古老的数学问题

问：今有

上禾三秉，中禾二秉，下禾一秉，实三十九斗；

上禾二秉，中禾三秉，下禾一秉，实三十四斗；

上禾一秉，中禾二秉，下禾三秉，实二十六斗。

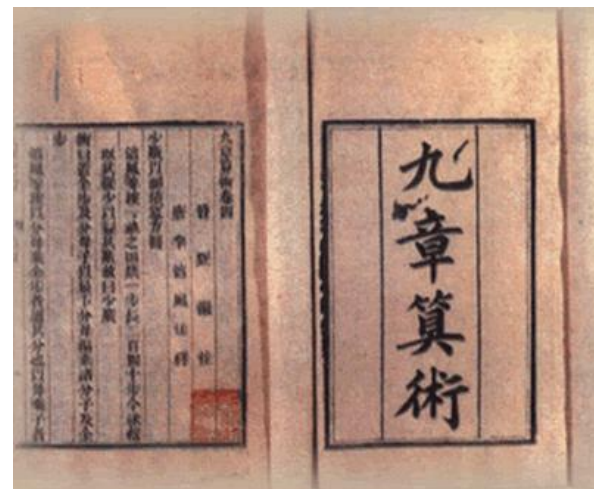
问上、中、下禾实一秉各几何？

$$3x + 2y + z = 39$$

$$2x + 3y + z = 34$$

$$x + 2y + 3z = 26$$

——《九章算术》



应用举例1

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

$$Ax = b$$

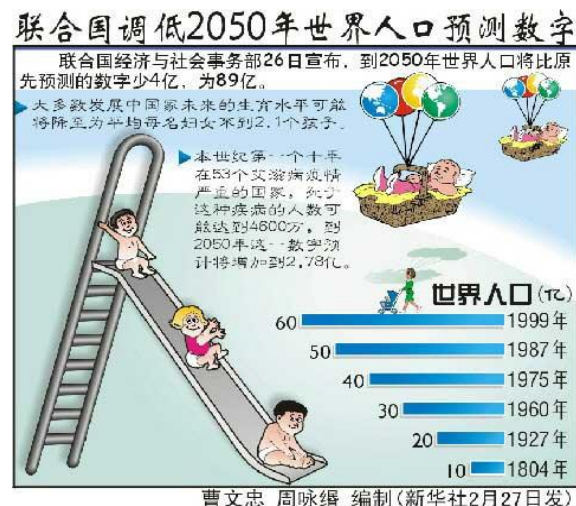
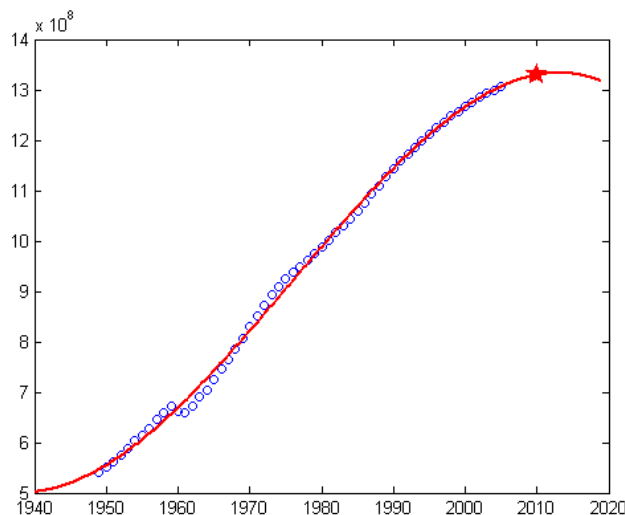
线性方程组数值求解
——第二章

应用举例2

例：人口预测

表格中是我国1950年到2005年的人口数（见中国统计年鉴），试预测未来的人口数

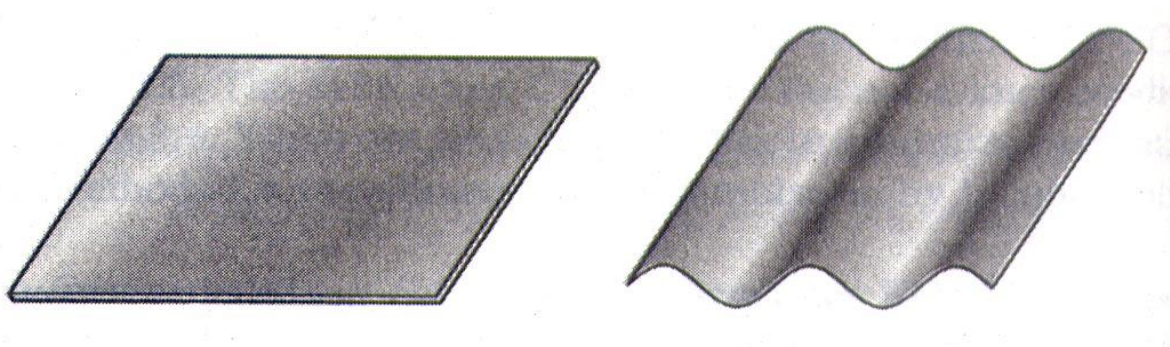
年份	人口(万)
1950	55196
1955	61465
1960	66207
1965	72538
1970	82992
1975	92420
1980	98705
1985	105851
1990	11433
1995	121121
2000	126743
2005	130756



插值与曲线拟合——第五章

应用举例3

例：铝制波纹瓦的长度问题



建筑上用的一种铝制波纹瓦是由机器将一块平整的铝板压制而成。假若要求波纹瓦长 4 英尺，每个波纹的高度(从中心线)为 1 英寸，且每个波纹以近似 2π 英寸为一个周期。求制做一块波纹瓦所需铝板的长度 L 。

应用举例3

这个问题就是要求由函数

$$f(x)=\sin x$$

给定的曲线从 $x=0$ 到 $x=48$ 英寸间的弧长 L ，即：

$$L = \int_0^{48} \sqrt{1 + (f'(x))^2} \, dx = \int_0^{48} \sqrt{1 + (\cos x)^2} \, dx$$

上述积分为第二类椭圆积分，无法用普通方法来计算

数值积分与数值微分
——第六章

应用举例4

Google 搜索引擎

PageRank: 对搜索结果按重要性排序

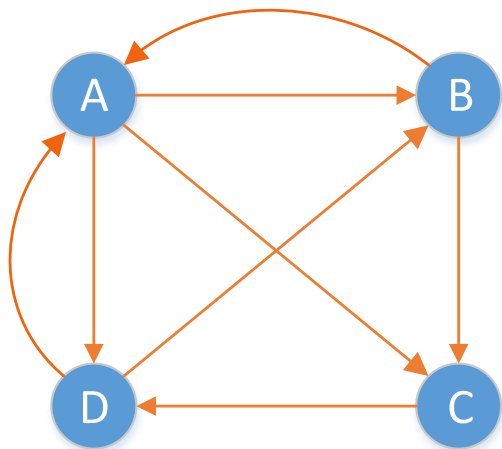
基本原理“从优质的网页链接过来的网页
必定还是优质网页”

超链接 $A \rightarrow B \equiv A$ 对 B 投一票
若 A 的质量高，则该投票分数高



1998年，美国斯坦福大学的Larry Page和Sergey Brin创立了Google公司

应用举例4



$$M = \begin{bmatrix} 0 & 1/2 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \\ 1/3 & 0 & 1 & 0 \end{bmatrix}$$

- 由于所有网页的 \mathbf{Pr} 值可由所有链向它的页面的重要性加和得到
- $\mathbf{M} \cdot \mathbf{Pr} = \mathbf{Pr}$
- $(\mathbf{M} - \mathbf{I}) \cdot \mathbf{Pr} = \mathbf{0}$

矩阵特征值计算——第四章

线性方程组求解问题
——第二章

应用举例5

蝴蝶效应

洛伦兹吸引子(*Lorenz attractor*)是由MIT大学的气象学家Edward Lorenz在1963年给出的, 他给出第一个混沌现象——蝴蝶效应。

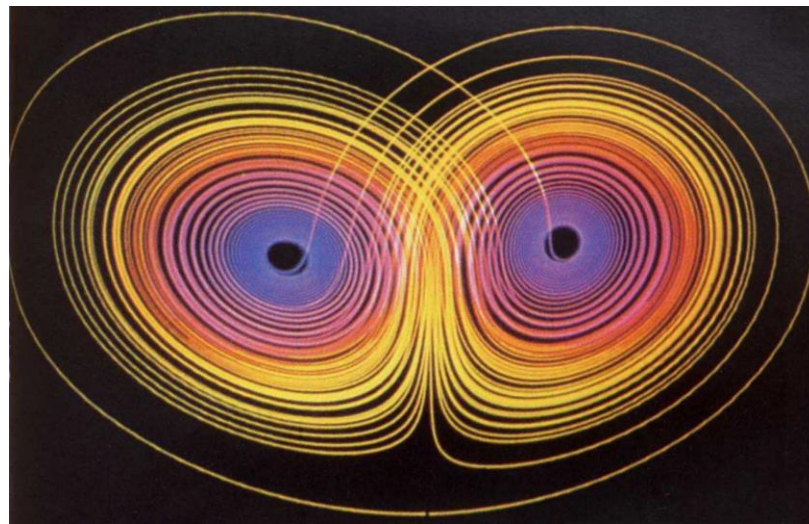


图1 蝴蝶效应示意图

应用举例5

洛伦兹方程是大气流体动力学模型的一个简化的常微分方程组：

$$\begin{cases} \frac{dx}{dt} = -\sigma x + \sigma y \\ \frac{dy}{dt} = rx - y - xz \\ \frac{dz}{dt} = -bz + xy \end{cases}$$

常（偏）微分方程数值解——第七章

1.1.1 科学计算与计算方法

1.1.1 科学计算与计算方法

- 计算工具



1.1.1 科学计算与计算方法

- 计算工具



图一：结绳计数/记事



图二：古代算筹计数的摆法

横式	—	=	≡	≡	≡	⊥	⊥	⊥	≡
纵式						⊥	⊥	⊥	≡
	1	2	3	4	5	6	7	8	9

1.1.1 科学计算与计算方法

• 计算工具



图一：结绳计数/记事



图二：古代算筹计数的摆法

横式	—	=	≡	≡	≡	⊥	⊥	⊥	≡
纵式						⊥	⊥	⊥	≡
	1	2	3	4	5	6	7	8	9



1.1.1 科学计算与计算方法

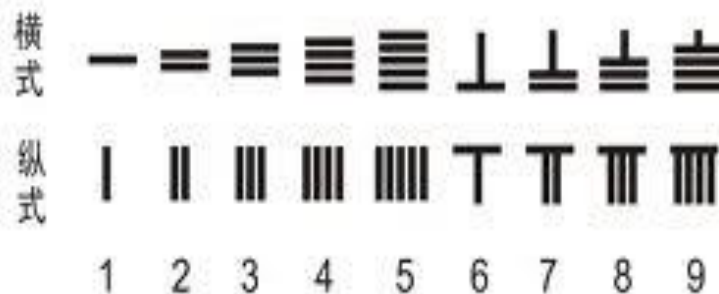
• 计算工具



图一：结绳计数/记事



图二：古代算筹计数的摆法



1.1.1 科学计算与计算方法

计算

- 古老

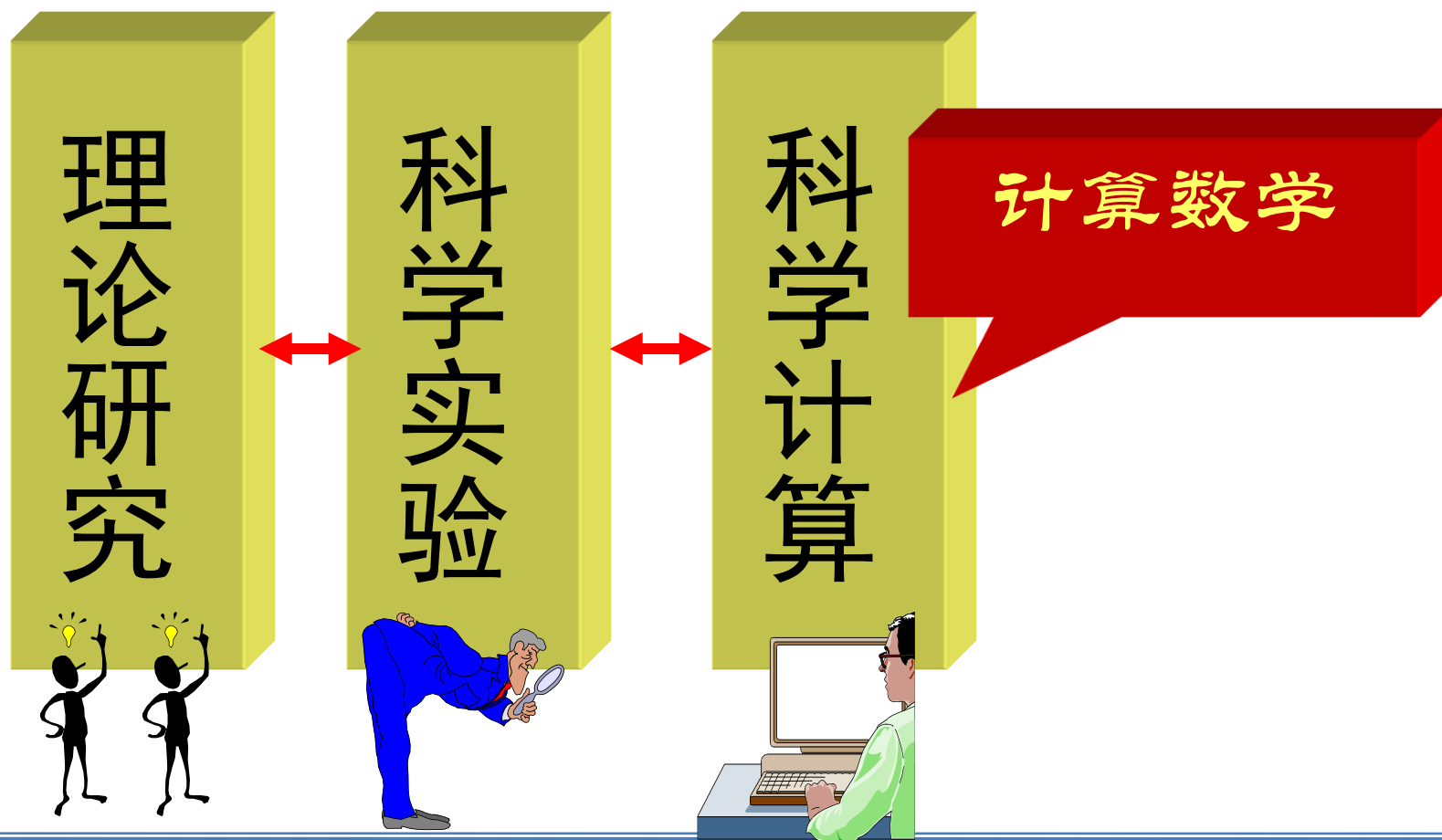
- 年轻

- 与时俱进

计算物理、计算化学、计算力学、
数量经济、计算生物学.....

1.1.1 科学计算与计算方法

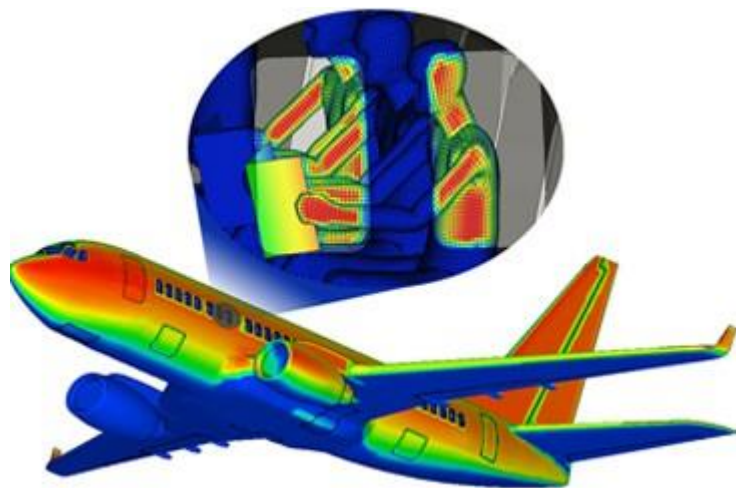
- 科学研究



1.1.1 科学计算与计算方法

- 科学研究

有限元建模



风洞



1.1.1 科学计算与计算方法

- 科学计算

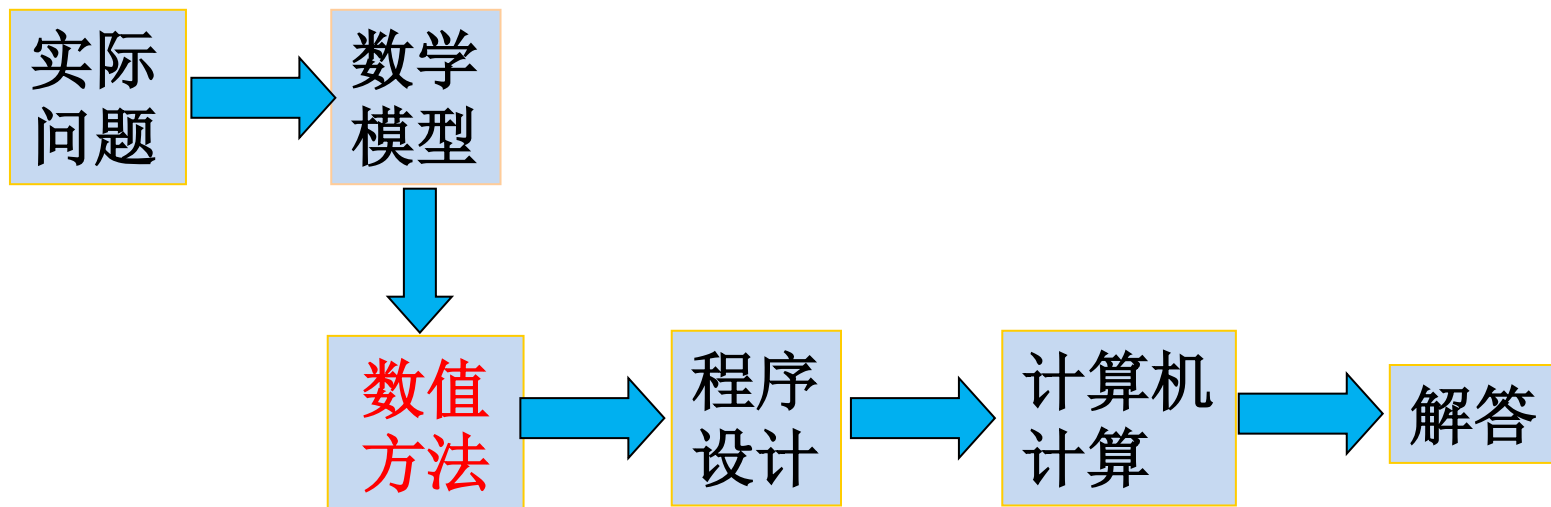
计算机计算能力的提高，衍生了计算方法这门课程



1.1.1 科学计算与计算方法

- 科学计算是人类从事科学活动和解决科学技术问题不可缺少的手段。
 - 计算机科学技术的发展，为科学计算及数据处理提供了高速和高精度的计算工具。
 - 计算机运算: 只能进行加，减，乘，除等算术运算和一些逻辑运算。
-

1.1.1 科学计算与计算方法



科学计算的过程

1.1.1 科学计算与计算方法

- 三维地图构建



1.1.1 科学计算与计算方法

- 三维地图构建

- (1) 实际问题：空中航测（空中连续拍照）方法，构建某地三维地形图。

- (2) 数学模型：建立一个大型超定线性方程组。

- (3) 数值方法：采用最小二乘方法求解该方程组的最小二乘解。

- (4) 程序设计：任何语言。

- (5) 计算机计算

- (6) 问题的解，对解进行解释。

1.1.1 科学计算与计算方法

- 数值计算方法定义

数学的一个分支，它以数字计算机求解数学问题的方法与理论为研究对象，

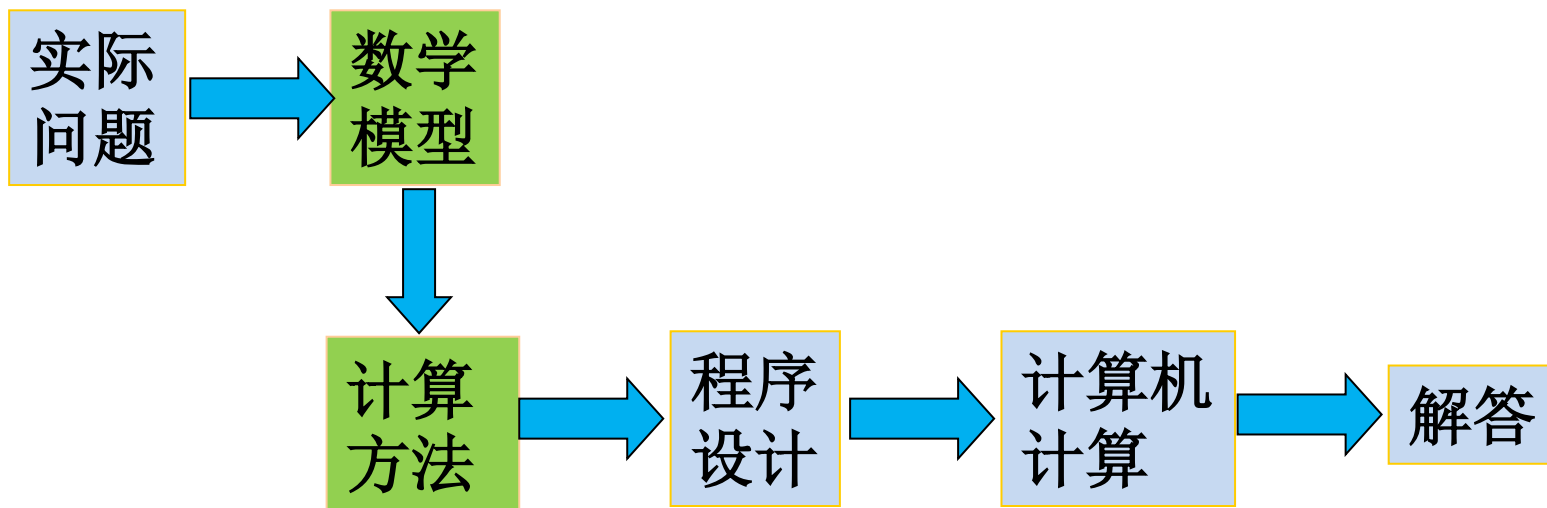
其内容包括：

1.1.1 科学计算与计算方法

- 线性方程组的解法
 - 矩阵特征值与特征向量的计算
 - 函数插值
 - 数值微分与积分
 - 非线性方程(组)的解法与最优化问题的计算方法
 - 常微与偏微分方程的数值解法
 -
 - 有关计算方法可靠性的理论研究, 如方法的收敛性和稳定性分析与误差估计等.
-

1.1.2 数学模型与计算方法

1.1.2 数学模型与计算方法



科学计算的过程



知道了数学模型，直接在计算机上编程实现不就行了吗？为什么要学习计算方法呢？

1.1.2 数学模型与计算方法

以定积分求解为例

- 对于积分

$$I = \int_a^b f(x)dx$$

- 由微积分知识可知：只要找到被积函数 $f(x)$ 的原函数 $F(x)$ ，便有下列牛顿—莱布尼兹公式

$$\int_a^b f(x)dx = F(b) - F(a)$$

1.1.2 数学模型与计算方法

- 原因一：原函数不能用初等函数表示成有限形式

$$\int_a^b \sin x^2 dx, \quad \int_a^b \frac{\sin x}{x} dx$$

- 原因二：原函数过于复杂

$$f(x) = x^2 \sqrt{2x^2 + 3}$$

$$F(x) = \frac{x^3 \sqrt{2x^2 + 3}}{4} + \frac{3x \sqrt{2x^2 + 3}}{16} - \frac{9}{16\sqrt{2}} \ln(\sqrt{2}x + \sqrt{2x^2 + 3})$$

1.1.2 数学模型与计算方法

- 原因三： $f(x)$ 以离散数据点形式给出

x_i	x_0	x_1	\dots	x_n
$y_i = f(x_i)$	y_0	y_1	\dots	y_n

1.1.2 数学模型与计算方法

有了数学模型，并不一定能够直接用计算机求解，因此我们需要学习计算方法！

- 学习计算方法这门课程，能够让我们学会如何构造或选择那些“好”的数值计算方法。

1.1.2 数学模型与计算方法

秦九韶算法

- 考虑对任意给定的 x , 计算代数多项式

$$P_n(x) = a_0x^n + a_1x^{n-1} + \cdots + a_{n-1}x + a_n = \sum_{k=0}^n a_k x^{n-k}$$

的值的问题。

1.1.2 数学模型与计算方法

秦九韶算法

- 考虑对任意给定的 x , 计算代数多项式

$$P_n(x) = a_0x^n + a_1x^{n-1} + \cdots + a_{n-1}x + a_n = \sum_{k=0}^n a_k x^{n-k}$$

的值的问题。显然, 上式等价于

$$P_n(x) = (((a_0x + a_1)x + a_2)x + \cdots + a_{n-1})x + a_n$$

1.1.2 数学模型与计算方法

- 考虑线性代数方程组 $Ax=b$ 的求解计算问题。其中系数矩阵 A 为一 n 阶方阵, $D=\det(A)\neq 0$

- Cramer法则

$$x_i = D_i / D \quad i=1,2,\dots,n$$

-

1.1.2 数学模型与计算方法

设 S_n 为 n 阶行列式所做的乘法次数, 则

$$D = \begin{vmatrix} & & \\ & & \\ & & \end{vmatrix}$$

$$= a_{11} \cdot S_{n-1} + a_{12} \cdot S_{n-1} + \dots + a_{1n} S_{n-1}$$

$$= n \cdot S_{n-1} + n$$

$$= n(S_{n-1} + 1)$$

$$\begin{cases} S_n = n(S_{n-1} + 1) \\ \vdots \\ S_2 = 2 \\ S_1 = 0 \end{cases}$$

1.1.2 数学模型与计算方法

$$S_n = n(S_{n-1} + 1)$$

$$= n[(n-1)(S_{n-2} + 1) + 1]$$

$$= n(n-1)(S_{n-2} + 1) + n$$

$$= n(n-1)(n-2)(S_{n-3} + 1) + n(n-1) + n$$

$$= n(n-1)(n-2)S_{n-3} + n(n-1)(n-2) + n(n-1) + n$$

⋮

$$= n(n-1)(n-2) \cdots 3 \times S_2 + \cdots n(n-1)(n-2) \cdots 3 + \cdots$$

$$= n(n-1)(n-2) \cdots 3 \times 2 S_1 + n(n-1)(n-2) \cdots 2 + \cdots + n$$

$$= n! + \frac{n!}{2!} + \cdots + \frac{n!}{(n-1)!}$$

$$= n! \left(1 + \frac{1}{2!} + \cdots + \frac{1}{(n-1)!} \right)$$

1.1.2 数学模型与计算方法

即需要做大约 $n!$ 次乘法, 若不计加法. 则 Cramer 公式求解线性方程组要做

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

次以上的乘法. 若 $n=20$, 则 $(20)! \approx 5.11 \times 10^{19}$ 以上的乘法

~~/0000 0000 0000 0000~~

天 时 秒

10"

百亿次

1.1.3 计算方法的特点及学习方法

1.1.3 计算方法的特点及学习方法

- 课程特点

计算方法是一门与计算机应用紧密结合、实用性很强的数学课程，它所涉及的数学问题面很广、内容非常丰富，亦有其自身的体系。它既有数学的高度概括，又非常讲究实用性并具有高度的技巧性。

1.1.3 计算方法的特点及学习方法

- 课程特点

第一，面向计算机，研究计算机上用的计算方法.

（算法最终只可包含四则运算和逻辑运算）

第二，要有可靠的理论分析.

（算法收敛性、稳定性及误差分析）

第三，要注重算法的效率.

（计算时间、存储空间）

第四，要重视数值实验.

1.1.3 计算方法的特点及学习方法

- 本课程的学习方法

- ① 掌握构造方法的原理、思想，理解算法，会分析算法精度
- ② 注重算法的效率和适用范围，针对不同情况学会选择和设计优秀算法
- ③ 要重视实践，通过算例和动手计算，学会怎样使用数值方法在计算机上解决各类数学计算问题

1.2 误差

误差

- 计算机的浮点表示及算术运算
 - 误差来源
 - 误差的基本概念
 - 误差估计
-

1. 2. 1 计算机的浮点表示及算术运算

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

— 实数 x 在计算机中被表示为

$$x = \pm 0.d_1d_2\cdots d_k \times 2^p$$

$$d_1=1, d_i \text{ 为 } 0 \text{ 或 } 1, \quad i=2,3,\dots,k, \quad -m \leq p \leq M.$$

— 零的浮点数通常表示为

$$0 = \pm 0.00\cdots 0 \times 2^{-m}$$

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

— 实数 x 在计算机中被表示为

$$x = \pm 0.\underbrace{d_1 d_2 \cdots d_k}_{\text{尾数}} \times 2^{\underbrace{p}_{\text{阶码}}}$$

$$d_1=1, d_i \text{ 为 } 0 \text{ 或 } 1, \quad i=2, 3, \dots, k, \quad -m \leq p \leq M.$$

— 零的浮点数通常表示为

$$0 = \pm 0.00 \cdots 0 \times 2^{-m}$$

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

$$x = \pm 0.d_1d_2\cdots d_k \times 2^p$$

— 给定的二进制浮点计算机，只能表示所有形如上式的有限数集 $S=S(k,m,M)$ ，这是实数轴上的不等距有限点集。

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

$$S=S(3,1,2) \quad k=3, \quad m=1, \quad M=2, \quad -1 \leq p \leq 2$$

那么计算机能表示的浮点数集合是如下的33个点。

$$0 = 0.00 \cdots 0 \times 2^{-1}$$

$$x = \pm 0.1d_2d_3 \times 2^{-1}$$

$$x = \pm 0.1d_2d_3 \times 2^1$$

$$x = \pm 0.1d_2d_3 \times 2^0$$

$$x = \pm 0.1d_2d_3 \times 2^2$$

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

$$x = \pm 0.1d_2d_3 \times 2^{-1}$$

$$0.100 \times 2^{-1} = 1/4$$

$$0.101 \times 2^{-1} = 5/16$$

$$0.110 \times 2^{-1} = 3/8$$

$$0.111 \times 2^{-1} = 7/16$$

$$x = \pm 0.1d_2d_3 \times 2^0$$

$$0.100 \times 2^0 = 1/2$$

$$0.101 \times 2^0 = 5/8$$

$$0.110 \times 2^0 = 3/4$$

$$0.111 \times 2^0 = 7/8$$

1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

$$x = \pm 0.1d_2d_3 \times 2^{-1}$$

$$0.100 \times 2^{-1} = 1/4$$

$$0.101 \times 2^{-1} = 5/16$$

$$0.110 \times 2^{-1} = 3/8$$

$$0.111 \times 2^{-1} = 7/16$$

$$x = \pm 0.1d_2d_3 \times 2^0$$

$$0.100 \times 2^0 = 1/2$$

$$0.101 \times 2^0 = 5/8$$

$$0.110 \times 2^0 = 3/4$$

$$0.111 \times 2^0 = 7/8$$



1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

$$x = \pm 0.1d_2d_3 \times 2^{-1}$$

$$0.100 \times 2^{-1} = 1/4$$

$$0.101 \times 2^{-1} = 5/16$$

$$0.110 \times 2^{-1} = 3/8$$

$$0.111 \times 2^{-1} = 7/16$$

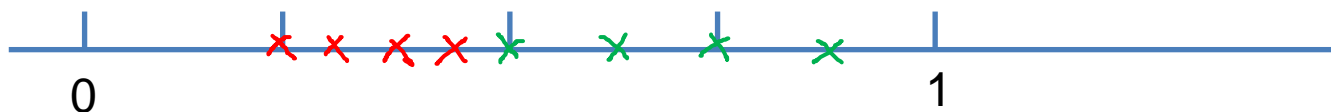
$$x = \pm 0.1d_2d_3 \times 2^0$$

$$0.100 \times 2^0 = 1/2$$

$$0.101 \times 2^0 = 5/8$$

$$0.110 \times 2^0 = 3/4$$

$$0.111 \times 2^0 = 7/8$$



1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

例如单精度实数用32位的二进制表示，其中符号位占1位，尾数占23位，阶数占8位，可以写成如下形式

$$x = \pm 0.d_1d_2\cdots d_{23} \times 2^p \quad |p| \leq 2^7-1 \quad (1.2.2)$$

注意上面的8位阶数中须有1位表示阶数的符号，所以阶数值占7位。

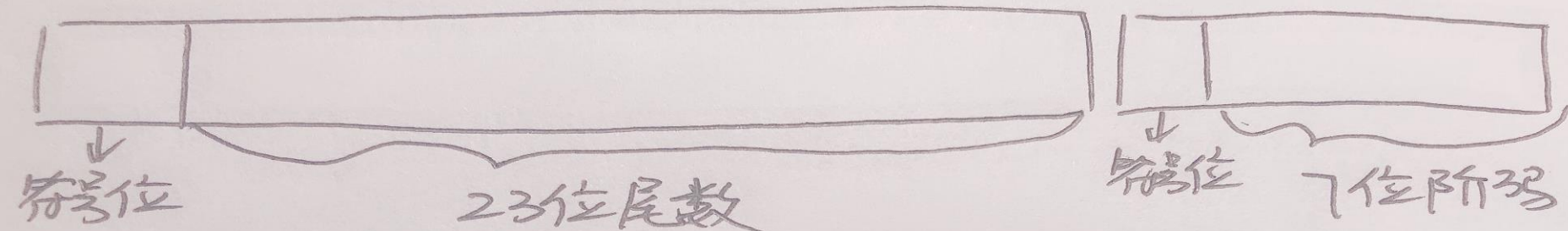
1.2.1 计算机的浮点表示及算术运算

- 数的浮点表示

例如单精度实数用32位的二进制表示，其中符号位占1位，尾数占23位，阶数占8位，可以写成如下形式

$$x = \pm 0.d_1d_2\cdots d_{23} \times 2^p \quad |p| \leq 2^7-1 \quad (1.2.2)$$

注意上面的8位阶数中须有1位表示阶数的符号，所以阶数值占7位。



1.2.1 计算机的浮点表示及算术运算

- 浮点数的四则运算

设是 S_{10} 由所有形如 $\pm 0.d_1d_2d_3d_4 \times 10^p$ 的4位十进制浮点数的集合，其中 $1 \leq d_1 \leq 9$ ， $0 \leq d_i \leq 9$ ， $i = 2, 3, 4$ ，整数 p 满足 $-9 \leq p \leq 10$ 。下面举例说明 S_{10} 上的算术运算。

1.2.1 计算机的浮点表示及算术运算

- 浮点数的计算特点

- 加减法先对阶(将阶码统一为较大者)，后计算，再舍入

- 乘法先运算再舍入

- 不在计算机数系中的数做四舍五入处理

例1.1

- (1) $0.2015 \times 10^4 + 0.1911 \times 10^2$
→ $0.2015 \times 10^4 + 0.0019 \times 10^4$ 对阶
→ 0.2034×10^4 计算
- (2) $0.2015 \times 10^4 + 0.1911 \times 10^{-1}$
→ $0.2015 \times 10^4 + 0.0000 \times 10^4$ 对阶
→ 0.2015×10^4 计算
- (3) $0.2015 \times 10^4 - 0.2008 \times 10^4$
→ 0.0007×10^4 计算
→ 0.7000×10^1 规范化

例1.1

(4) $(0.2015 \times 10^4) \times (0.1911 \times 10^{-5})$

$\rightarrow (0.2015 \times 0.1911) \times 10^{-1}$ 对阶

$\rightarrow (0.3851 \times 10^{-1}) \times 10^{-1}$ 计算

$\rightarrow 0.3851 \times 10^{-2}$ 规范化

(5) $(0.2015 \times 10^4) \div (0.1911 \times 10^{-5})$

$\rightarrow (0.2015 \div 0.1911) \times 10^9$ 对阶

$\rightarrow (0.1054 \times 10^1) \times 10^9$ 计算

$\rightarrow 0.1054 \times 10^{10}$ 规范化

1.2.1 计算机的浮点表示及算术运算

- 浮点数的计算特点
- 计算过程中应该注意
 - 绝对值相差悬殊的两个数做加减，会造成“大数吃小数”的现象；（例2）
 - 非常接近的数相减，会损失掉有效数字；（例3）
 - 相对被除数来说，绝对值很小的数做除数，会产生绝对值很大的数，甚至溢出；（例5）
 - 在运算过程中注意合理安排运算顺序，以便提高运算的精度或保护重要的参数。

例1.2

在前面所述的4位十进制浮点计算机（数集 S_{10} ）上求解如下一元二次方程

$$x^2 - 24x + 1 = 0 \quad (1.2.3)$$

按求根公式，此方程的两个根是

$$x_1 = 12 + \sqrt{143} \quad x_2 = 12 - \sqrt{143}$$

在4位十进制浮点计算机上， $\sqrt{143} = 0.1196 \times 10^2$ ，于是按照上面求根公式有

$$x_1 = 0.2396 \times 10^2, \quad x_2 = 0.4000 \times 10^{-1}$$

下面我们换一种方法进行计算 x_2 ，即

$$x_2 = 12 - \sqrt{143} = \frac{1}{12 + \sqrt{143}} \quad (1.2.4)$$

则 $x_2=0.4174\times 10^{-1}$ 。

事实上， x_2 的精确解应为 $0.0417393\dots$ 。

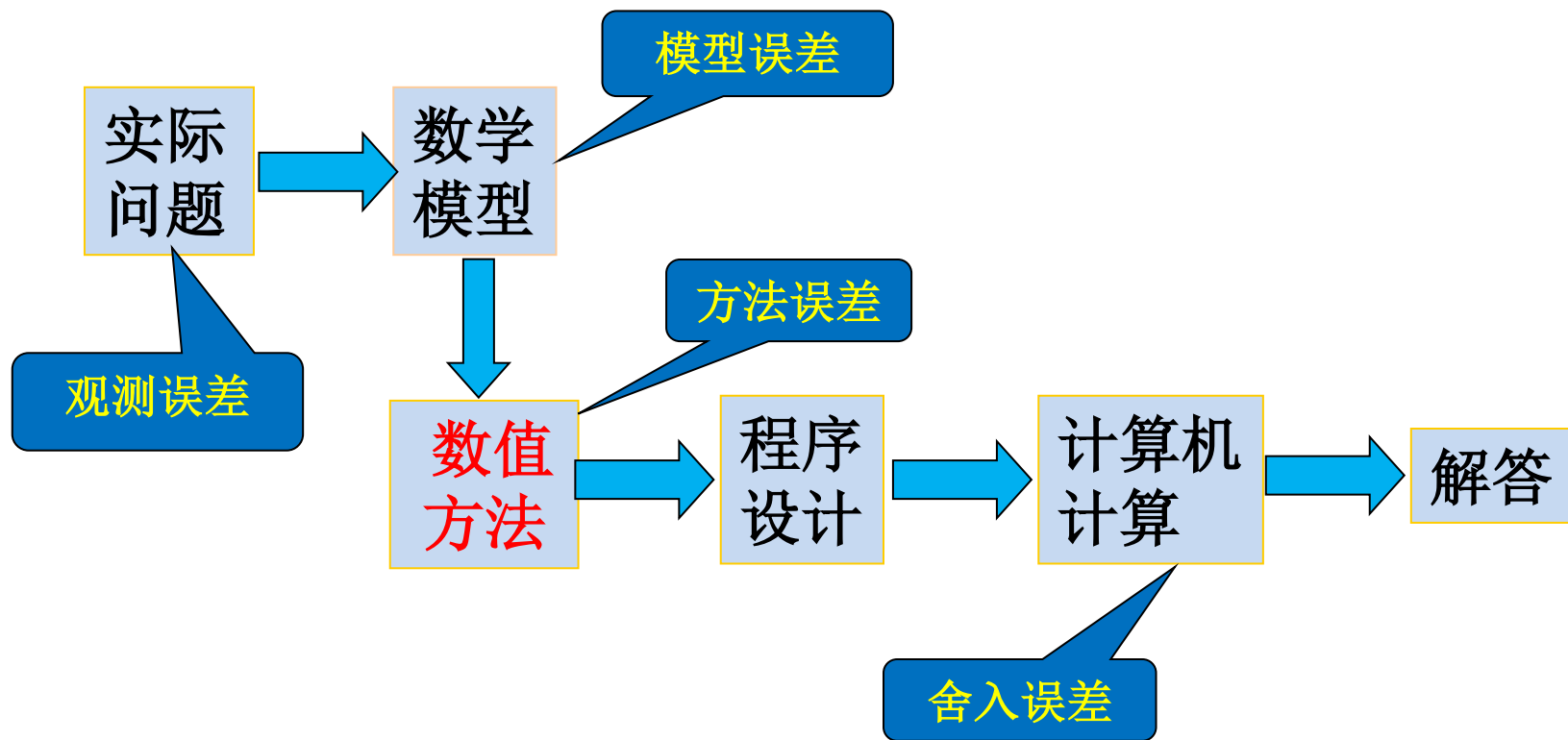
1.2.1 计算机的浮点表示及算术运算

- 问题



1.2.2 误差来源

1.2.2 误差来源



1.2.2 误差来源

计算地球的表面积

$$A=4\pi R^2$$

- 观测误差
 - 模型误差
 - 方法误差
 - 舍入误差
-

例1.3

为了计算函数值 e^x , $|x| < 1$, 我们用Taylor多项式

$$P_n(x) = 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} \quad (1.12)$$

近似代替 e^x , 此时的方法误差 (又称截断误差) 为

$$R_n(x) = e^x - P_n(x) = \frac{x^{n+1}e^\xi}{(n+1)!}, \quad |\xi| < 1 \quad (1.13)$$

1. 2. 3 误差的基本概念

- 绝对误差与绝对误差限

定义1.1 设 x 为准确值， x^* 是 x 的一个近似值，称 $e^*=x^*-x$ 为近似值 x^* 的绝对误差，或简称误差。

定义1.2 设 $\varepsilon^* > 0$ ，并满足

$$|e^*| = |x^* - x| \leq \varepsilon^* \quad (1.2.7)$$

则称 ε^* 为近似值 x^* 的绝对误差限，或简称误差限。

- 相对误差

定义1.3 设 x 为精确值， x^* 为近似值，则称比值

$$\frac{e^*}{x} = \frac{x^* - x}{x} \quad (1.2.8)$$

为近似值 x^* 的**相对误差**，记作 e_r^* （实际应用时，常用 x^* 代替上式分母中的 x ）。

- 相对误差限

定义1.4 设 ε^* 是近似值 x^* 的误差限，则称

$$\varepsilon_r^* = \frac{\varepsilon^*}{|x^*|} \quad (1.2.9)$$

为近似值 x^* 的**相对误差限**，此时，有

$$\frac{|x^* - x|}{|x^*|} \leq \frac{\varepsilon^*}{|x^*|} = \varepsilon_r^* \quad (1.2.10)$$

- 有效数字

定义1.5 如果

$$|e^*| = |x^* - x| \leq 0.5 \times 10^{-n} \quad (1.2.11)$$

则说 x^* 近似表示 x 准确到 10^{-n} 位（小数点后第 n 位），并从此位起直到最左边的非零数字之间的一切数字都称为有效数字，并把有效数字的位数称为有效位数。

例1.4

取 e ($e = 2.718281828459$) 的近似值 $x^* = 2.72$, 则

$$|2.72 - e| \leq 0.001718\cdots \leq 0.5 \times 10^{-2}$$

即 x^* 近似表示 e 准确到 10^{-2} 位, 因此具有3位有效数字。

若取 e 的近似值 $x^* = 2.71828$, 则

$$|2.71828 - e| \leq 0.0000018\cdots \leq 0.5 \times 10^{-5}$$

即 x^* 近似表示 e 准确到 10^{-5} 位, 因此具有6位有效数字。

- 有效数字

定义1.6 将 x 的近似值 x^* 表示为十进制浮点数的标准形式

$$x^* = \pm 0.\alpha_1\alpha_2\dots\alpha_k \times 10^m \quad (\alpha_i=0,1,\dots,9, \alpha_1 \neq 0) \quad (1.2.12)$$

如果

$$|e^*| = |x^* - x| \leq 0.5 \times 10^{m-n} \quad (1.2.13)$$

则说近似值 x^* 具有 n 位有效数字。这里 n 是正整数， m 是整数。

例

- 若 $x^*=3578.64$ 是 x 的具有6位有效数字的近似值，试求 x^* 的误差限。

例

- 若 $x^*=3578.64$ 是 x 的具有6位有效数字的近似值，试求 x^* 的误差限。

$$x^* = 3578.64 = 0.357864 \times 10^4$$

$$n = 6$$

$$|e^*| \leq 0.5 \times 10^{m-n} = 0.5 \times 10^{-2}$$

- 有效数字与相对误差的关系

定理1.1 若近似值 x^* 具有 n 位有效数字，则其相对误差满足

$$|e_r^*| \leq \frac{1}{2\alpha_1} \times 10^{-(n-1)} \quad (1.2.14)$$

反之，如果 x^* 的相对误差 e_r^* 满足

$$|e_r^*| \leq \frac{1}{2(\alpha_1 + 1)} \times 10^{-(n-1)} \quad (1.2.15)$$

则 x^* 至少具有 n 位有效数字。

1.2.4 误差分析

1. 2. 4 误差分析

1.2.4 误差分析

- 将带有误差的数据进行计算时，误差在运算过程中会进行传播，必然导致计算结果出现误差。一般来说，精确值 x 与近似值 x^* 之间都比较接近，其误差可以看作是一个小的增量，即可以把误差看作微分，即

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

$$e_r^* = \frac{e^*}{x} = \frac{dx}{x} = d \ln x$$

这表明： x 的微分表示 x 的误差， $\ln x$ 的微分表示 x 的相对误差

1.2.4 误差分析

根据上式，可以得到算术运算的误差，以 x, y 两数为例

$$e^*(x \pm y) = d(x \pm y) = dx \pm dy = e^*(x) \pm e^*(y)$$

$$e^*(xy) = d(xy) = y dx + x dy = y e^*(x) + x e^*(y)$$

$$e^*\left(\frac{x}{y}\right) = d\left(\frac{x}{y}\right) = \frac{y dx - x dy}{y^2} = \frac{y e^*(x) - x e^*(y)}{y^2} \quad (y \neq 0)$$

$$e_r^*(xy) = d \ln(xy) = d \ln(x) + d \ln(y) = e_r^*(x) + e_r^*(y)$$

$$e_r^*\left(\frac{x}{y}\right) = d \ln\left(\frac{x}{y}\right) = d \ln(x) - d \ln(y) = e_r^*(x) - e_r^*(y) \quad (y \neq 0)$$

1.2.4 误差分析

而更一般的情况是，当自变量有误差时计算函数值时也会产生误差。其误差可以用函数的Taylor展开式进行估计。

$$f(x) - f(x^*) = f'(x^*)(x - x^*) + \frac{f''(\xi)}{2}(x - x^*)^2$$

$$e^*(f(x^*)) = |f(x^*) - f(x)| \approx |f'(x^*)| \cdot e^*(x^*)$$

例1.5

已知

$$(\sqrt{2} - 1)^6 = 99 - 70\sqrt{2} = \frac{1}{99 + 70\sqrt{2}}$$

由于 $\sqrt{2}$ 的精确值未知，取 $\sqrt{2} \approx 1.414$ 进行计算，试问上述3个表达式哪个计算精度最高？

$$f_1(x) = (x-1)^6$$

$$f_2(x) = 99 - 70x$$

$$f_3(x) = \frac{1}{99 + 70x}$$

例1.5

解：记 $|e^*(\sqrt{2})| = |\sqrt{2} - 1.414| \leq \frac{1}{2} \times 10^{-3}$

令 $f_1(x) = (x-1)^6$ ，则 $f_1'(x^*) = 6(x-1)^5$ ，于是

$$\begin{aligned} \left| (\sqrt{2}-1)^6 - (1.414-1)^6 \right| &= \left| f_1(\sqrt{2}) - f_1(1.414) \right| \approx f_1'(1.414) |\sqrt{2} - 1.414| \\ &= \left| 6(1.414-1)^5 \right| \cdot |e^*(\sqrt{2})| \leq 0.073 \cdot |e^*(\sqrt{2})| \end{aligned}$$

同理，令 $f_2(x) = 99 - 70x$ ，则 $f_2'(x) = -70$ ，于是

$$\begin{aligned} \left| (99 - 70\sqrt{2}) - (99 - 70 \times 1.414) \right| &= \left| f_2(\sqrt{2}) - f_2(1.414) \right| \approx f_2'(1.414) |\sqrt{2} - 1.414| \\ &= |-70| \cdot |e^*(\sqrt{2})| \leq 70 \cdot |e^*(\sqrt{2})| \end{aligned}$$

令 $f_3(x) = \frac{1}{99 + 70x}$ ，则 $f_3'(x) = \frac{-70}{(99 + 70x)^2}$ ，于是

$$\begin{aligned} \left| \frac{1}{99 + 70\sqrt{2}} - \frac{1}{99 + 70 \times 1.414} \right| &= \left| f_3(\sqrt{2}) - f_3(1.414) \right| \approx f_3'(1.414) |\sqrt{2} - 1.414| \\ &= \left| \frac{-70}{(99 + 70 \times 1.414)^2} \right| \cdot |e^*(\sqrt{2})| \leq 0.002 \cdot |e^*(\sqrt{2})| \end{aligned}$$

例1.5

例 1.5 已知 $(\sqrt{2}-1)^6 = 99-70\sqrt{2} = \frac{1}{99+70\sqrt{2}}$ ，由于 $\sqrt{2}$ 的精确值未知，取 $\sqrt{2} \approx 1.414$ 计算

该连等式的值，试问用连等式中的哪个表达式计算的精度最高？

解：记 $|e^*(\sqrt{2})| = |\sqrt{2} - 1.414| \leq \frac{1}{2} \times 10^{-3}$

令 $f_1(x) = (x-1)^6$ ，则 $f_1'(x) = 6(x-1)^5$ ，于是

$$\begin{aligned} \left| (\sqrt{2}-1)^6 - (1.414-1)^6 \right| &= \left| f_1(\sqrt{2}) - f_1(1.414) \right| \approx f_1'(1.414) |\sqrt{2} - 1.414| \\ &= \left| 6(1.414-1)^5 \right| \cdot |e^*(\sqrt{2})| \leq 0.073 \cdot |e^*(\sqrt{2})| \end{aligned}$$

同理，令 $f_2(x) = 99-70x$ ，则 $f_2'(x) = -70$ ，于是

$$\begin{aligned} \left| (99-70\sqrt{2}) - (99-70 \times 1.414) \right| &= \left| f_2(\sqrt{2}) - f_2(1.414) \right| \approx f_2'(1.414) |\sqrt{2} - 1.414| \\ &= |-70| \cdot |e^*(\sqrt{2})| \leq 70 \cdot |e^*(\sqrt{2})| \end{aligned}$$

令 $f_3(x) = \frac{1}{99+70x}$ ，则 $f_3'(x) = \frac{-70}{(99+70x)^2}$ ，于是

$$\begin{aligned} \left| \frac{1}{99+70\sqrt{2}} - \frac{1}{99+70 \times 1.414} \right| &= \left| f_3(\sqrt{2}) - f_3(1.414) \right| \approx f_3'(1.414) |\sqrt{2} - 1.414| \\ &= \left| \frac{-70}{(99+70 \times 1.414)^2} \right| \cdot |e^*(\sqrt{2})| \leq 0.002 \cdot |e^*(\sqrt{2})| \end{aligned}$$

由此结果可以看出，第 3 个表达式的精度最高，第 2 个表达式的精度最差。

例1.6

考虑积分

$$I_n = \int_0^1 x^n e^{x-1} dx \quad (1.25)$$

的近似计算.

此积分满足递推关系式

$$I_n = 1 - nI_{n-1}, \quad (1.26)$$

假定我们首先计算出 I_0 的近似值 \bar{I}_0 ，保留3位有效数字，利用递推关系式(1.26)依次算出

\bar{I}_0	\bar{I}_1	\bar{I}_2	\bar{I}_3	\bar{I}_4	\bar{I}_5	\bar{I}_6	\bar{I}_7
0.632	0.368	0.264	0.208	0.168	0.160	0.040	0.720

例1.6

$$\bar{I}_1 = 1 - 1 \times \bar{I}_0 = 1 - 0.632 = 0.368$$

$$\bar{I}_2 = 1 - 2 \times \bar{I}_1 = 1 - 0.736 = 0.264$$

$$\bar{I}_6 = 0.040$$

$$\bar{I}_7 = 1 - 7 \times \bar{I}_6 = 1 - 0.280 = 0.720$$

例1.6

$$0 < x < 1 \Rightarrow 0 < x^{n+1} e^{x-1} < x^n e^{x-1}$$

- 根据积分公式

$$I_{n+1} = \int_0^1 x^{n+1} e^{x-1} dx < \int_0^1 x^n e^{x-1} dx = I_n$$

- 那么

$$\bar{I}_7 < \bar{I}_6$$

例1.6

反之，若将(1.26)式改写成

$$I_{n-1} = (1 - I_n) / n, \quad (1.27)$$

先计算出 I_7 的近似值，再从 \bar{I}_7 开始按(1.27)式递推，可依次算出

\bar{I}_7	\bar{I}_6	\bar{I}_5	\bar{I}_4	\bar{I}_3	\bar{I}_2	\bar{I}_1	\bar{I}_0
0.112	0.127	0.146	0.171	0.207	0.264	0.368	0.632
0.112	0.127	0.146	0.171	0.207	0.264	0.368	0.632

思考题

- 请对式(1.2.4) 的误差进行分析。

习题

- 习题1 下列各数都是经过四舍五入得到的近似数，试指出它们的有效数字的位数，并给出相对误差限
- $x_1^* = 3.1416$; $x_2^* = 0.028$; $x_3^* = 1.250$; $x_4^* = 3 \times 10^4$.
- 习题2 利用四位数学用表求 $1 - \cos 2^\circ$ 的近似值，比较下面几种方法的结果，分析各自误差

习题

- ① 直接按 $1-\cos 2^\circ$ 求解;
- ② 按 $1-\cos 2^\circ = \sin^2 2^\circ / (1+\cos 2^\circ)$ 求解;
- ③ 按 $1-\cos 2^\circ = 2\sin^2 1^\circ$ 求解。

其中，在数学用表中 $\cos 2^\circ = 0.9994$, $\sin 2^\circ = 0.0349$, $\sin 1^\circ = 0.0175$.

- 习题3证明定理1.1

习题

- 习题4 将3.141, 3.14, 3.15, $22/7$ 分别作为 π 近似值, 试确定它们各有几位有效数字, 并确定其相对误差限。
- 习题5 设 $x = 10 \pm 0.05$, 试求函数 $f(x) = \sqrt[n]{x}$ 的相对误差限。

习题

习题6 函数sinx有幂级数展开

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

利用幂级数计算sinx的Octave程序为

```
function s=power sin(x)
% Power series for sin(x)
s=0;
t=x;
n=1;
while s+t!= s;
s=s+t;
t=-x^2/(n+1)/(n+2)*t;
n=n+2;
end while
endfunction
```

- 1) 解释上述程序的终止准则;
- 2) 如果将幂级数在某确定位 N 之后截断, 试分析误差。

实验

- 实验1

问题提出： 计算 $1-0.2-0.2-0.2-0.2-0.2$

实验要求：

得到了什么结果？ 分析原因。

实验2

- 实验要求:

- ① 保留4位有效数字，给出例0.6中(0.26)式的误差估计式。
- ② 其中 I_0 可由(0.25)推出，即 $I_0=1-e^{-1}$ 。首先用Octave实现(0.26)式的递推过程，之后再实现4位十进制数的程序实现。
- ③ Octave的计算精度为15位，因此前者可以看作是后者舍入之前的精确值，比较两者之间的误差是否与你给出的误差估计式相符。

实验3

- 问题提出
- 函数 $f(x)$ 的导数定义为

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

于是可以给出计算 $f(x)$ 在 x_0 点导数的两种计算方法:

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0)}{h} \quad (0.28)$$

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0 - h)}{2h} \quad (0.29)$$

- 实验要求:

(1)选择有代表性的函数 $f(x)$, 分别用上面两种方法求导数;

(2)比较对同样的步长 h , 两种方法的精度如何, 并解释原因

(3) 对同一种方法, 比较步长变化时精度的变化, 比如从 10^{-2} 到 10^{-14} , 之后思考原因。