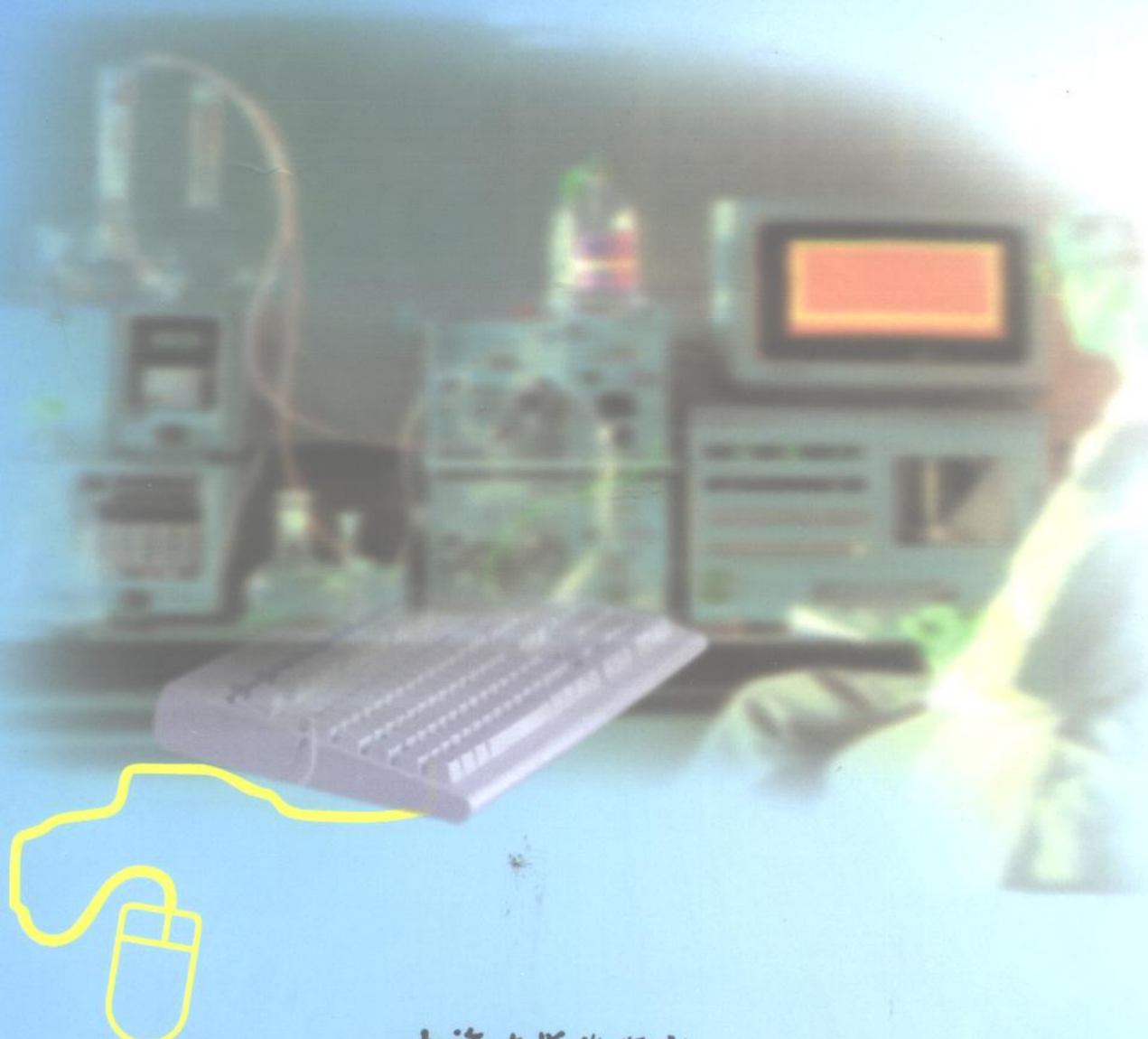


# 计算机在化学化工中的应用

吴若峰 乐之伟 陆文聪 编著



上海大学出版社

计算机在化学化工中的应用

上海大学出版社

# 计算机在化学化工中的应用

吴若峰 乐之伟 陆文聪 编著

上海大学出版社

·上海·

## 内 容 提 要

本书除介绍化学化工领域中常用的曲线拟合、回归分析、高次方程求根、线性方程组求解、数值积分、微分方程求解等经典算法外,特别对一些近年来在化学化工领域中得到广泛应用的计算技术作了详细介绍,如单纯形(simplex)优化技术、模式识别技术、神经网络方法、蒙特卡罗(Monte Carlo)方法、定量结构—活性相关分析(QSAR)等,并以相当篇幅介绍这些方法在解决化学化工实际问题时的应用。此外,对国外一些优秀的科技软件如Origin、ACD/ChemSketch等也略作介绍。附录中编写了常用算法的QBASIC程序。

本书可作为高等院校化学化工及相近专业本科生的教材或教学参考书,也可供化学工作者及其他科技工作者使用。

2P2S/29

### 图书在版编目(CIP)数据

计算机在化学化工中的应用 / 吴若峰编著. — 上海: 上海大学出版社, 2000. 8

ISBN 7-81058-008-6

I. 计... II. 吴... III. ①计算机应用—化学②计算机应用—化学工业 IV. TQ015.9

中国版本图书馆 CIP 数据核字(2000)第 40971 号

上海大学出版社出版发行

(上海市延长路 149 号 邮政编码: 200072)

上海大学印刷厂印刷 各地新华书店经销

开本: 787×1092 1/16 印张 9.25 字数 213 千字

2000 年 8 月第 1 版 2000 年 8 月第 1 次印刷

印数: 1~1100

定价: 17.00 元

# 前 言

当前化学研究已达到分子设计的水平,化工生产和管理也多采用计算机控制。计算机引入化学化工领域后,在帮助深入研究化学基本理论和促进化工生产两个方面都显示了强大作用。化学化工专业的学生必须具备用计算机解决化学化工问题的能力,才能面对 21 世纪的挑战。因此,只学习一般的计算机文化和技术知识显然不够,为此,近年来国内化学化工类专业都已经或准备开设计算机在化学化工中应用的课程,作为化学化工类各专业的必修专业基础课。

上海大学化学化工类专业开设“计算机在化学化工中的应用”课程已有 7 年。该课程是上海市教委重点课程建设项目,也是上海大学“面向 21 世纪高等化学教育课程体系和教学内容改革计划”的教学改革课程。本书是为配合该课程的建设而编写的。

本书除介绍化学化工领域中常用的曲线拟合、回归分析、高次方程求根、线性方程组求解、数值积分、微分方程求解等经典算法外,特别根据近年来计算机技术及其在化学化工领域中应用的发展趋势,对一些新方法和新应用作了详细介绍,如单纯形(simplex)优化技术、模式识别技术、人工神经网络方法、蒙特卡罗(Monte Carlo)方法、定量结构—活性相关分析(QSAR)等,并以相当篇幅介绍这些方法在解决化学化工实际问题时的应用。此外,对国外一些优秀的科技软件如 Origin、ACD/ChemSketch 等也略作介绍。附录中编写了常用算法的 QBASIC 程序。

本书各章的执笔人是:吴若峰(第一、四、五、六、八、九、十章),乐之伟(第二、三章),陆文聪(第七章)。研究生王征编写了附录中的大部分内容。编写工作得到上海大学教务处的关心和支持。

编 者

2000 年 6 月

# 目 录

<b>第一章 最小二乘法和曲线拟合</b> .....	1
1-1 线性最小二乘法.....	2
1-2 用线性最小二乘法确定方程的系数.....	3
1-3 非线性方程变为线性方程.....	5
1-4 非线性最小二乘法.....	7
1-5 线性回归分析.....	7
习题.....	9
上机作业 1 气体反应动力学方程的确定 .....	10
上机作业 2 用 Origin 软件进行曲线拟合 .....	11
<b>第二章 高次方程求解</b> .....	15
2-1 扫描法 .....	15
2-2 对分法和优选法 .....	16
2-3 迭代法 .....	18
2-4 牛顿法 .....	19
习题 .....	19
上机作业 用牛顿迭代法计算化学反应的平衡浓度 .....	20
<b>第三章 线性方程组求解</b> .....	21
3-1 克莱姆法则 .....	21
3-2 消元法 .....	22
3-3 按列选主元素消去法 .....	23
3-4 按列选主元素消去法的计算流程 .....	25
习题 .....	26
上机作业 通过求解线性方程组确定混合物的组成 .....	26
<b>第四章 数值积分</b> .....	28
4-1 矩形和梯形法求积 .....	28
4-2 辛普森法求积 .....	29
4-3 高斯-勒让德法求积 .....	31
4-4 Monte Carlo 方法求积 .....	33
习题 .....	33
上机作业 等压加热过程中的热效应计算 .....	34
<b>第五章 常微分方程的数值解</b> .....	35
5-1 欧拉法 .....	35
5-2 改进的欧拉法——预测—校正法 .....	36
5-3 龙格—库塔法 .....	38

习题 .....	39
上机作业 化学反应动力学计算 .....	40
<b>第六章 单纯形最优化方法 .....</b>	<b>41</b>
6-1 改变单因子法 .....	41
6-2 单纯形法 .....	42
6-2-1 基本单纯形法 .....	42
6-2-2 改良单纯形法 .....	44
6-2-3 初始单纯形的定位 .....	45
6-2-4 单纯形迭代的结束 .....	46
6-2-5 多维单纯形法 .....	46
6-2-6 改良单纯形法的计算程序 .....	46
6-3 单纯形法最优化实例 .....	49
习题 .....	52
<b>第七章 模式识别与工业优化 .....</b>	<b>53</b>
7-1 模式识别与工业数据处理 .....	53
7-2 模式识别方法的原理和基本概念 .....	53
7-3 数据标准化处理 .....	53
7-4 特征参数的抽提 .....	54
7-5 常用模式识别方法简介 .....	54
7-5-1 KNN 法及其衍生法 .....	54
7-5-2 主成分分析方法 .....	55
7-5-3 多重判别矢量法 .....	57
7-5-4 Fisher 判别分析法 .....	58
7-5-5 非线性映照法 .....	59
7-6 人工神经网络简介 .....	60
7-7 模式识别优化的应用 .....	61
<b>第八章 Monte Carlo 方法 .....</b>	<b>63</b>
8-1 Monte Carlo 方法的基本思想 .....	64
8-2 Monte Carlo 方法的特点 .....	65
8-3 Monte Carlo 方法的基本工具——随机数 .....	65
8-4 随机变量的抽样 .....	66
8-4-1 离散型分布随机变量的抽样 .....	67
8-4-2 连续型分布随机变量的抽样 .....	68
8-5 Monte Carlo 方法在高分子研究中的应用 .....	69
8-5-1 共聚反应的模拟 .....	69
8-5-2 高分子邻基反应的模拟 .....	75
8-5-3 高分子降解的模拟 .....	75
8-5-4 高分子链构象的模拟 .....	76
习题 .....	77

上机作业 1 两元共聚物链结构的 Monte Carlo 模拟 .....	77
上机作业 2 用 Monte Carlo 方法计算聚乙烯醇的最大缩醛化率 .....	78
上机作业 3 高分子无规降解的 Monte Carlo 模拟 .....	79
上机作业 4 高分子链构象的 Monte Carlo 模拟 .....	80
<b>第九章 定量构效相关分析方法——QSAR .....</b>	<b>81</b>
9-1 QSAR 方法的基本特点 .....	81
9-2 用分子连接性指数法进行 QSAR 研究 .....	82
9-2-1 分子连接性指数法的特点 .....	82
9-2-2 分子连接性指数的计算方法 .....	82
9-2-3 分子连接性指数的修正 .....	85
9-2-4 分子连接性指数计算示例 .....	87
9-3 QSAR 分析示例 .....	88
习题 .....	89
上机作业 建立含氮杂环化合物结构与生物活性之间的关系 .....	90
<b>第十章 化学结构绘制软件 ACD/ChemSketch 的应用 .....</b>	<b>91</b>
10-1 软件的运行环境 .....	91
10-2 软件的界面及工具按钮 .....	91
10-3 化学结构 3D 显示和动画效果 .....	92
<b>附录 .....</b>	<b>94</b>
A 化学结构绘制软件 ACD/ChemSketch 使用说明 .....	94
B QBASIC 语言简介 .....	101
C 常用 QBASIC 程序 .....	111
<b>参考文献 .....</b>	<b>138</b>

# 第一章 最小二乘法和曲线拟合

化学实验中测量的数据往往是一些条件的函数,如温度、压力、浓度、pH 值、离子强度等等。如果变量只有一个,一般可以在平面图(二维坐标系)中直观地表示,如某种物质的溶解度与温度的关系(图 1-1)、蒸气压与温度的关系、平衡常数与温度的关系,高分子化学中聚合度与转化率的关系、共聚物组成与单体投料比的关系、溶液粘度与浓度的关系、结晶速率与温度的关系、聚合物玻璃化和温度与分子量的关系等,均可在二维坐标系上表示。

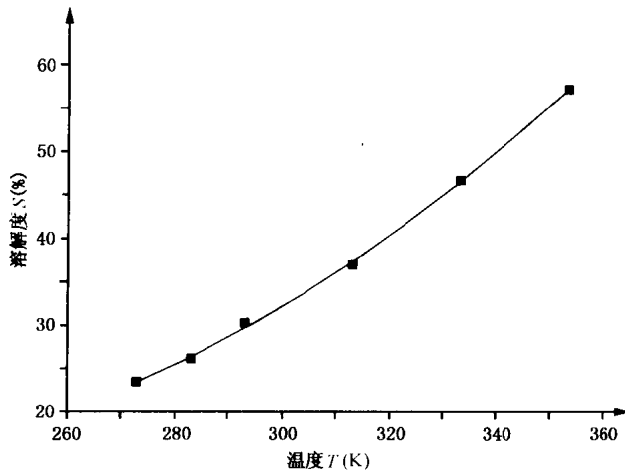


图 1-1 某种物质溶解度  $S$  与温度  $T$  的关系

图 1-1 表示某种物质溶解度与温度的关系,经曲线拟合后得到

$$S = 7.62 \times 10^{-8} T^{3.484} \quad (1-1)$$

化学化工中有很多类似 (1-1) 式的方程,其中大部分是由实验数据通过曲线拟合而得到的。它们一般都具有下列特点: ① 表达式简单; ② 形式相似; ③ 对多变量的函数关系,往往用多项式表示,以便求导、积分及进行计算机处理。

本章主要介绍这类经验方程的确定方法。确定一个经验方程,一般可采用作图法或代数法。这里我们介绍代数法,因为这也是计算机处理的方法基础。如欲确定线性方程  $y = a + bx$  中的系数  $a$  和  $b$ ,用代数法求这样一个含  $a$  和  $b$  两个未知数的方程,需要两个独立的方程,也就是需要两对实验测量值  $x_i, y_i$  ( $i = 1, 2$ ),但在实际上,为减少误差,至少应取五对实验测量值  $x_i, y_i$  ( $i = 1, 2, \dots, 5$ )。这样,用代数法处理时,这五对实验测量值便会构成五个方程,若将它们两两组合,可以产生十个二元一次方程组,解这些方程组所得到的  $a, b$  值可能大相径庭,故这类方程组叫做矛盾方程组。下面是一个简单的例子。

例 1-1 某合成气转化率与压力的关系如表 1-1 所示。



表 1-1 某合成气转化率与压力的关系

压力 $P$	2	4	5	8	9
转化率 $T$	2.01	2.98	3.50	5.02	5.47

先将这些实验测量值在二维坐标系上描点,如图 1-2。

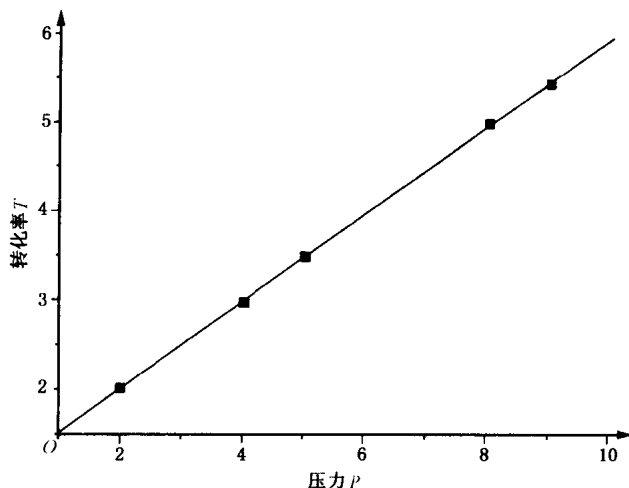


图 1-2 某合成气转化率  $T$  与压力  $P$  的关系

从图 1-2 中可看出,这两种变量之间有一种接近直线的关系,试用直线方程来表示:

$$T = a + bP \quad (1-2)$$

把第一和第二对实验测量值代入上述直线方程,建立下面的方程组:

$$\begin{aligned} a + 2b &= 2.01 \\ a + 4b &= 2.98 \end{aligned} \quad (1-3)$$

解得  $a = 1.0900$ ,  $b = 0.4850$ 。

把第三和第四对实验测量值代入上述直线方程,建立下面的方程组:

$$\begin{aligned} a + 5b &= 3.50 \\ a + 8b &= 5.02 \end{aligned} \quad (1-4)$$

解得  $a = 0.9000$ ,  $b = 0.5200$ 。

由此可见,任选两对实验测量值求解  $a, b$ , 得到的结果并不一致,这时候我们就难以确定这个方程。因此,怎样从给定的实验数据出发,在某个函数类型  $\Phi(x)$  中寻求出一个最好的函数  $\phi(x)$  来拟合这些数据,应当有一个原则。为此建立了最小二乘法。

### 1-1 线性最小二乘法

假如我们已经确定了方程  $y = ax + b$  的系数  $a$  和  $b$ , 则对每一个实验自变量  $x_k$ , 都可以计算出一个  $y_k^*$ :

$$y_k^* = ax_k + b \quad (1-5)$$

式中  $k$  是实验的次数。 $y_k^*$  是  $y_k$  的计算值,而  $y_k$  是由实验变量  $x_k$  所观察到的实验因变量。 $y_k^*$

和  $y_k$  之间的差值称为残差,用  $e_k$  表示。

$$e_k = y_k - y_k^* = y_k - (ax_k + b) \quad (1-6)$$

显然,  $e_k$  的大小,可衡量被确定的系数  $a, b$  的好坏,反过来,好的系数的确定应遵循使  $e_k$  最小这样一个原则。

有这样几种基于  $e_k$  来调节系数的方法可供选择:

1. 调节系数使原来  $e_k$  绝对值最大的一个变得最小,就是使  $T = \max |e_k|$  最小;
2. 调节系数使  $e_k$  的绝对值之和变得最小,就是使  $A = \sum |e_k|$  最小;
3. 调节系数使  $e_k$  的平方和变得最小,就是使  $Q = \sum e_k^2$  最小。

曲线拟合中应用较广泛的是第三种方法。

现在可以看出,曲线拟合的原则是,被确定的曲线要求尽可能从实验点的附近通过,达到在某种平均意义上的误差积累最小。

### 1-2 用线性最小二乘法确定方程的系数

需要说明,下面的方法是针对  $y$  是  $x$  的多元线性函数而建立的,根据上一节,

$$Q = \sum e_k^2 \quad (1-7)$$

$$Q = \sum (y_k - y_k^*)^2 \quad (1-8)$$

而  $y_k^* = f(x_1, x_2, x_3, \dots, x_n; b_1, b_2, b_3, \dots, b_n)$ 。

为简便,用  $b_i$  表示  $b_1, b_2, b_3, \dots, b_n$  中的任一个,用  $x_i$  表示  $x_1, x_2, x_3, \dots, x_n$  中的任一个。要使  $Q$  达到最小,各个  $b_i$  应满足下列方程组:

$$\left. \begin{aligned} \frac{\partial Q}{\partial b_1} &= 0 \\ \frac{\partial Q}{\partial b_2} &= 0 \\ \dots\dots \end{aligned} \right\} \quad (1-9)$$

即残差的平方和对各系数的偏导数全为 0。

如上所述,  $y$  是  $b_i$  的线性函数,即可设

$$y = b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_nx_n \quad (1-10)$$

把  $m$  次实验数据代入(1-10)式得方程组

$$\left. \begin{aligned} y_1 &= b_1x_{11} + b_2x_{21} + b_3x_{31} + \dots + b_nx_{n1} \\ y_2 &= b_1x_{12} + b_2x_{22} + b_3x_{32} + \dots + b_nx_{n2} \\ &\dots\dots \\ y_m &= b_1x_{1m} + b_2x_{2m} + b_3x_{3m} + \dots + b_nx_{nm} \end{aligned} \right\} (m \geq n) \quad (1-11)$$

如果  $b_1$  是常数项,则令  $x_{1k} \equiv 1$ 。这样,

$$Q = \sum_{k=1}^m [y - (b_1x_{1k} + b_2x_{2k} + \dots + b_nx_{nk})]^2$$

$$\frac{\partial Q}{\partial b_i} = 2 \left( b_1 \sum_{k=1}^m x_{1k} x_{ik} + b_2 \sum_{k=1}^m x_{2k} x_{ik} + \cdots + b_n \sum_{k=1}^m x_{nk} x_{ik} - \sum_{k=1}^m x_{ik} y_k \right) = 0$$

令

$$S_{i1} = \sum_{k=1}^m x_{1k} x_{ik}, S_{i2} = \sum_{k=1}^m x_{2k} x_{ik}, \dots, S_{in} = \sum_{k=1}^m x_{nk} x_{ik}, S_{iy} = \sum_{k=1}^m x_{ik} y_k$$

则

$$\frac{\partial Q}{\partial b_i} = (b_1 S_{i1} + b_2 S_{i2} + \cdots + b_n S_{in} - S_{iy}) = 0 \quad (1-12)$$

这样,当  $x_{1k}, x_{2k}, \dots, x_{nk}, y_k$  从实验中得到后,  $S_{i1}, \dots, S_{in}, S_{iy}$  也可求得,于是得到以  $b_i$  为未知数的方程组:

$$\left. \begin{aligned} S_{11}b_1 + S_{12}b_2 + \cdots + S_{1n}b_n &= S_{1y} \\ S_{21}b_1 + S_{22}b_2 + \cdots + S_{2n}b_n &= S_{2y} \\ S_{31}b_1 + S_{32}b_2 + \cdots + S_{3n}b_n &= S_{3y} \\ &\dots\dots\dots \\ S_{m1}b_1 + S_{m2}b_2 + \cdots + S_{mn}b_n &= S_{my} \end{aligned} \right\} \quad (1-13)$$

这一个以  $b_1, b_2, \dots, b_n$  为未知数的  $n$  元联立方程组,称为原方程(1-10)的法方程,解得  $b_1, b_2, \dots, b_n$  后,也就确立了原方程。

例 1-2 某种铝合金的含铝量为  $x\%$ ,其溶解温度为  $y^\circ\text{C}$ ,由实验测得  $x$  与  $y$  的数据如表 1-2 所示,试用最小二乘法建立  $x$  与  $y$  之间的经验公式。

表 1-2 铝合金的含铝量和溶解温度的实验数据

$k$	$x_k$	$y_k$	$x_k^2$	$x_k y_k$
1	36.9	181	1 361. 61	6 678. 9
2	46.7	197	2 180. 89	9 199. 9
3	63.7	235	4 057. 69	14 969. 5
4	77.8	270	6 052. 84	21 006. 0
5	84.0	283	7 056. 00	23 772. 0
6	87.5	292	7 656. 25	25 550. 0
$\Sigma$	396. 6	1 458	28 365. 28	101 176. 3

解: 根据前面的讨论,求解过程如下:

1. 将表 1-2 给出的数据点  $(x_k, y_k)$  描绘在坐标纸上,如图 1-3 所示。
2. 确定拟合曲线的形式。由图 1-3 可以看出,六个点位于一条直线的附近,因此可先用线性函数来拟合这组实验数据。即令

$$y(x) = a + bx \quad (1-14)$$

注意,由于  $a$  为常数项,故  $x_{1k} \equiv 1, x_{2k} = x_k$ 。

3. 建立法方程组。由于问题归结为一次多项式拟合问题,相应的法方程形如:

$$\begin{aligned} S_{11}a + S_{12}b &= S_{1y} \\ S_{21}a + S_{22}b &= S_{2y} \end{aligned}$$

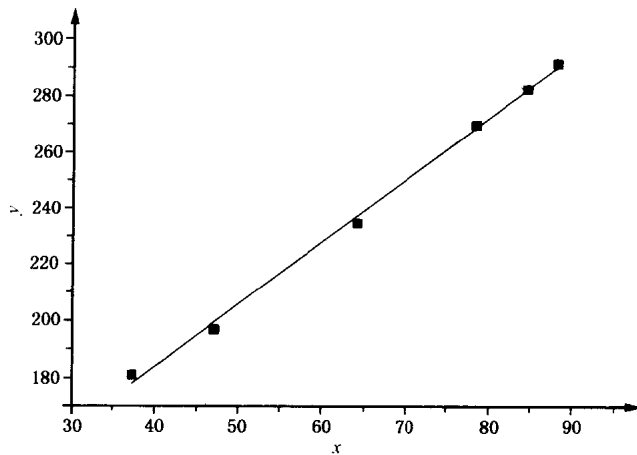


图 1-3 铝合金的熔解温度  $y$  与含铝量  $x$  的关系

其中

$$S_{11} = \sum_{k=1}^6 x_{1k} x_{1k} = 6$$

$$S_{12} = \sum_{k=1}^6 x_{2k} x_{1k} = 396.6$$

$$S_{21} = \sum_{k=1}^6 x_{1k} x_{2k} = 396.6$$

$$S_{22} = \sum_{k=1}^6 x_{2k} x_{2k} = 28\,365.28$$

$$S_{1y} = \sum_{k=1}^6 x_{1k} y_k = 1\,458$$

$$S_{2y} = \sum_{k=1}^6 x_{2k} y_k = 101\,176.3$$

即

$$6a + 396.6b = 1\,458$$

$$396.6a + 28\,365.28b = 101\,176.3 \quad (1-15)$$

4. 解此法方程,得

$$a = 95.3524, b = 2.2337 \quad (1-16)$$

从而得到经验公式

$$y = 95.3524 + 2.2337x \quad (1-17)$$

### 1-3 非线性方程变为线性方程

前面所介绍的是线性函数关系的最小二乘法,有些函数关系虽然不是线性关系,但可以通过代换转变为线性关系,然后再用线性最小二乘法处理。如

1. 双曲线

$$\frac{1}{y} = a + \frac{b}{x}$$

令

$$Y = \frac{1}{y}, X = \frac{1}{x}$$

代换方程为

$$Y = a + bX$$

## 2. 幂函数曲线

$$y = cx^b$$

两边取对数

$$\lg y = \lg c + b \lg x$$

令

$$Y = \lg y, X = \lg x, c' = \lg c$$

代换方程为

$$Y = c' + bX$$

## 3. 指数函数曲线

$$y = ce^{bx}$$

两边取对数

$$\ln y = \ln c + bx$$

令

$$Y = \ln y, c' = \ln c$$

代换方程为

$$Y = c' + bx$$

## 4. S 型曲线

$$y = \frac{1}{a + be^{-x}}$$

变形后得

$$\frac{1}{y} = a + be^{-x}$$

令

$$Y = \frac{1}{y}, X = e^{-x}$$

代换方程为

$$Y = a + bX$$

## 5. $n$ 次多项式

$$y = b_0 + b_1x + b_2x^2 + b_3x^3$$

令

$$X_1 = x, X_2 = x^2, X_3 = x^3$$

代换方程为

$$y = b_0 + b_1X_1 + b_2X_2 + b_3X_3$$

用上述方法将非线性方程转换为线性方程后,即可用 1-2 节所述的线性最小二乘法确

定方程的系数。

#### 1-4 非线性最小二乘法

如非线性方程  $y = f(x_1, x_2, x_3, \dots; b_1, b_2, b_3, \dots)$  不能通过代换法转换成线性方程,因而无法直接求解  $b_i$ ,则需要采用非线性最小二乘法。该方法的实质是用泰勒公式对非线性函数逐次线性化。

首先给每一个待定系数  $b_i$  以初值  $b_i^{(0)}$ ,这初值和真值的差为  $\Delta_i$ ,即

$$b_i = b_i^{(0)} + \Delta_i \quad (1-18)$$

这样,我们把确定  $b_i$  值的问题转化成了求极小  $\Delta_i$  的问题,为此我们对原来的非线性方程在  $b_i^{(0)}$  附近作泰勒展开,略去展开式中的高次项,就使之成为对  $\Delta_i$  的线性方程:

$$f(x_k, b_1, b_2, b_3, \dots, b_n) \approx f_{k0} + \frac{\partial f_{k0}}{\partial b_1} \Delta_1 + \frac{\partial f_{k0}}{\partial b_2} \Delta_2 + \dots + \frac{\partial f_{k0}}{\partial b_n} \Delta_n$$

式中

$$\begin{aligned} \Delta_1 &= b_1 - b_1^{(0)}, \Delta_2 = b_2 - b_2^{(0)}, \dots, \Delta_n = b_n - b_n^{(0)} \\ f_{k0} &= f(x_k, b_1^{(0)}, b_2^{(0)}, b_3^{(0)}, \dots, b_n^{(0)}) \\ \frac{\partial f_{k0}}{\partial b_i} &= \frac{\partial f(x_k, b_1^{(0)}, b_2^{(0)}, b_3^{(0)}, \dots, b_n^{(0)})}{\partial b_i} \end{aligned}$$

这样就得到了  $f(x_k, b_1, b_2, b_3, \dots, b_n)$  的线性近似式,其残差平方和  $Q$  的表达式也不难写出:

$$\begin{aligned} Q &= \sum [y_k - f(x_k, b_1, b_2, b_3, \dots, b_n)]^2 = \\ &= \sum \left[ y_k - \left( f_{k0} + \frac{\partial f_{k0}}{\partial b_1} \Delta_1 + \frac{\partial f_{k0}}{\partial b_2} \Delta_2 + \dots + \frac{\partial f_{k0}}{\partial b_n} \Delta_n \right) \right]^2 \end{aligned} \quad (1-19)$$

显见,  $Q$  已成为  $\Delta_i$  的函数,令  $\frac{\partial Q}{\partial \Delta_i} = 0$ ,通过 1-2 节所述步骤,最后得到原方程的法方程:

$$\left. \begin{aligned} S_{11}\Delta_1 + S_{12}\Delta_2 + \dots + S_{1n}\Delta_n &= S_{1y} \\ S_{21}\Delta_1 + S_{22}\Delta_2 + \dots + S_{2n}\Delta_n &= S_{2y} \\ S_{31}\Delta_1 + S_{32}\Delta_2 + \dots + S_{3n}\Delta_n &= S_{3y} \\ &\dots\dots\dots \\ S_{m1}\Delta_1 + S_{m2}\Delta_2 + \dots + S_{mn}\Delta_n &= S_{my} \end{aligned} \right\} \quad (1-20)$$

将解此法方程所得到的第一套修正值  $\Delta_i^{(0)}$  代入(1-18)式求得  $b_i^{(1)}$ ,再用上述方法求得第二套修正值  $\Delta_i^{(1)}$ ,并求得  $b_i^{(2)}$ ...这样经过  $n$  次迭代后,若  $\Delta_i^{(n)}$  小到一定程度,也就逼近了真值  $b_i$ 。

#### 1-5 线性回归分析

用最小二乘法求得经验方程后,需要给出一个定量的指标,从统计的意义上来说描述变量之间关系的密切程度,这称为回归分析。这里着重介绍一元线性回归分析。首先引出相关系

数  $R$ , 在一元线性回归中,  $R$  用下式计算:

$$R = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \quad (1-21)$$

式中,  $\bar{x}$  和  $\bar{y}$  分别是  $x_i$  和  $y_i$  的平均值。由(1-21)式可知, 不管  $x, y$  为何数值, 分子项的绝对值永远不会大于分母的值, 因此相关系数的取值范围是

$$0 \leq |R| \leq 1$$

$|R|$  值的大小反映了  $x, y$  之间相关的程度。 $|R|$  值为 0, 表示  $x, y$  之间线性无关,  $|R|$  值为 1, 表示  $x, y$  之间严格服从关系式, 是完全线性相关。绝大部分  $|R|$  值在 0 与 1 之间, 参见图 1-4。

在一元线性回归中, 相关系数只表示  $x$  与  $y$  线性关系的密切程度, 当  $R$  很小甚至为 0 时,  $x$  与  $y$  虽无线性关系, 但可能存在其他关系。如图 1-4(b)所示, 尽管  $x$  与  $y$  不呈线性关系, 但却存在明显的抛物线关系。在线性关系分析中, 只有当  $R$  的绝对值大到相当程度时, 方程表示的  $x, y$  的关系才有意义, 我们把这称为相关系数显著。另外, 相关系数  $R$  与样品的抽样数  $n$  有关。

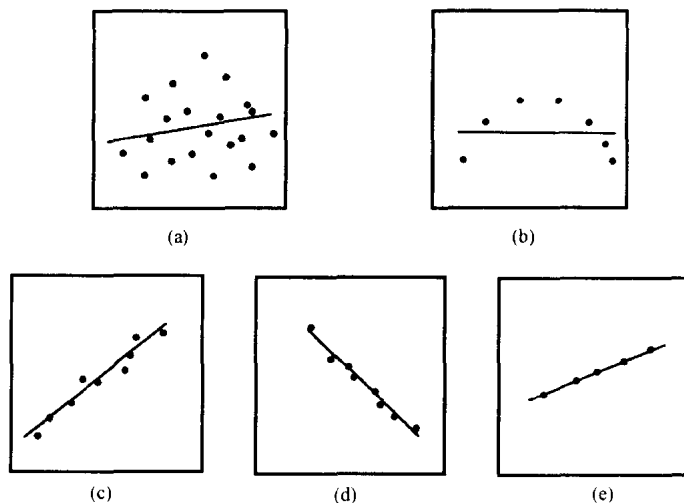


图 1-4 衡量  $x$  和  $y$  相关性好坏的参数——相关系数  $R$

(a)  $|R| \approx 0$ ; (b)  $|R| \approx 0$ ; (c)  $|R| \approx 0.6$ ; (d)  $|R| \approx 0.9$ ; (e)  $|R| \approx 1$

表 1-3 中给出了在不同观察样本数和不同显著性水平  $\alpha$  (0.05 及 0.01) 下的相关系数的最小值, 相关系数  $R$  只有达到某种显著水平的最小值才在该水平上有意义。

表 1-3 相关系数检查表

$n-2$	$\alpha=0.05$	$\alpha=0.01$	$n-2$	$\alpha=0.05$	$\alpha=0.01$	$n-2$	$\alpha=0.05$	$\alpha=0.01$
1	0.997	1.000	7	0.666	0.798	13	0.514	0.641
2	0.950	0.990	8	0.632	0.765	14	0.497	0.623
3	0.878	0.957	9	0.602	0.735	15	0.482	0.606
4	0.811	0.917	10	0.576	0.708	16	0.468	0.590
5	0.754	0.874	11	0.553	0.684	17	0.456	0.575
6	0.707	0.834	12	0.532	0.661	18	0.444	0.561

(续表)

$n-2$	$\alpha=0.05$	$\alpha=0.01$	$n-2$	$\alpha=0.05$	$\alpha=0.01$	$n-2$	$\alpha=0.05$	$\alpha=0.01$
19	0.433	0.549	28	0.361	0.463	80	0.217	0.283
20	0.423	0.537	29	0.355	0.456	90	0.205	0.267
21	0.413	0.526	30	0.349	0.449	100	0.195	0.254
22	0.404	0.515	35	0.325	0.418	125	0.174	0.228
23	0.396	0.505	40	0.304	0.393	150	0.159	0.208
24	0.388	0.496	45	0.288	0.372	200	0.138	0.181
25	0.381	0.487	50	0.273	0.354	300	0.113	0.148
26	0.374	0.478	60	0.250	0.325	400	0.098	0.128
27	0.367	0.470	70	0.232	0.302	1 000	0.062	0.081

例 1-3 分析 14 块某类矿石中 Ni 和  $P_2O_5$  的含量, 结果见表 1-4。

表 1-4 矿石中 Ni 和  $P_2O_5$  含量的分析数据

No.	Ni 含量(%)	$P_2O_5$ 含量(%)	No.	Ni 含量(%)	$P_2O_5$ 含量(%)
1	0.009	4.00	8	0.014	1.70
2	0.013	3.44	9	0.016	2.92
3	0.006	3.60	10	0.014	4.80
4	0.025	1.00	11	0.016	3.28
5	0.022	2.04	12	0.012	4.16
6	0.007	2.74	13	0.020	3.35
7	0.036	0.60	14	0.018	2.20

用最小二乘法对上述数据进行线性回归得到 Ni 和  $P_2O_5$  的回归方程是

$$y = 0.0306 - 0.0043x \quad (1-22)$$

考察对于 Ni 和  $P_2O_5$  含量数据, 共进行了 14 次样品分析,  $n = 14$ ,  $n - 2 = 12$ , 查相关系数检查表,  $R^{0.05} = 0.532$ ,  $R^{0.01} = 0.661$ , 根据 (1-21) 式计算  $R$  值为 0.804, 大于任何一种显著性水平的最小值, 故认为它是高度显著的, 即所确定的 Ni 和  $P_2O_5$  之间的线性回归方程是合理的。

### 习 题

1. 尿中胆色素经处理后, 在 550 nm 处有很强的吸光性, 现测得配制好的不同胆色素浓度的标准溶液的吸光率数据如表 1-5 所示。假定标定曲线可以用  $y = a + bx$  表示, 试计算上述方程参数  $a$ 、 $b$  的值, 并判断线性关系是否密切。

表 1-5 不同胆色素浓度标准溶液的吸光率

胆色素浓度(mg/100ml)	50	75	100	125	150	175	200	225	250
吸 光 率	0.039	0.061	0.087	0.107	0.119	0.163	0.179	0.194	0.213

2. 某催化剂活性  $Y$  与工作持续时间  $t$  的关系为

$$Y = Ae^{(Bt + Ct^2)}$$



将表 1-6 所列的实验数据通过曲线拟合求系数  $A, B, C$ 。

表 1-6 某催化剂活性  $Y$  与工作持续时间  $t$  的关系

$t(\text{h})$	0	27	40	52	70	89	106
$Y(\%)$	100	82.2	76.3	71.8	66.4	63.3	61.3

3. 将下列非线性函数线性化 ( $a, b, c$  是待定系数):

$$y = ae^{\frac{b}{x}}; y = \frac{x}{ax - b}; y = \frac{1}{a + be^{-x}}; y = kx_1^a x_2^b x_3^c$$

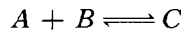
4. 化学中的 Antoine 方程能描述很多物质的蒸气压  $P$  与温度  $T$  之间的关系:

$$\lg P = a - \frac{b}{c + T}$$

如何将此方程线性化以求得参数  $a, b$  和  $c$ ?

### 上机作业 1 气体反应动力学方程的确定

设某气体反应可表示为



其反应动力学方程可用下列非线性方程表示:

$$V = KP_A^{n_1} P_B^{n_2} P_C^{n_3}$$

式中,  $V$  为反应速度,  $K$  为反应速度常数,  $P_A, P_B, P_C$  依次为气体  $A, B$  和  $C$  的分压, 而  $n_1, n_2$  和  $n_3$  为待求解的方程参数。

表 1-7 为实验测定的不同分压下的  $V$  值, 试根据此实验数据计算反应速度常数  $K$  及方程中的参数  $n_1, n_2$  和  $n_3$ 。

表 1-7 不同分压下的反应速度  $V$

$x_1 = \ln P_A$	$x_2 = \ln P_B$	$x_3 = \ln P_C$	$V$	$x_1 = \ln P_A$	$x_2 = \ln P_B$	$x_3 = \ln P_C$	$V$
2.197	2.116	0.993	8.58	1.946	1.361	2.292	2.18
2.104	1.946	1.482	6.05	1.917	1.224	2.282	2.11
2.067	1.825	1.775	4.73	1.872	0.956	2.389	1.88
1.946	1.459	2.104	3.35	2.079	0.788	2.879	1.04

首先将动力学方程的两边取对数:

$$\ln V = \ln K + n_1 \ln P_A + n_2 \ln P_B + n_3 \ln P_C$$

令

$$y = \ln V, x_1 = \ln P_A, x_2 = \ln P_B, x_3 = \ln P_C$$

则上式变为

$$y = \ln K + n_1 x_1 + n_2 x_2 + n_3 x_3$$