



中国原子能科学研究院科学技术丛书

数据的分析和处理

赵志祥 编著



中国原子能出版社

中国原子能科学研究院科学技术丛书

数据的分析和处理

赵志祥 编著

中国原子能出版社

图书在版编目(CIP)数据

数据的分析和处理 / 赵志祥编著. —北京: 中国

原子能出版社, 2012. 12

ISBN 978-7-5022-5770-5

I. ①数… II. ①赵… III. ①统计数据—统计分析
(数学)②数据处理 IV. ①0212.1②TP274

中国版本图书馆 CIP 数据核字(2012)第 272329 号

数据的分析和处理

出版发行 中国原子能出版社(北京市海淀区阜成路 43 号 100048)

责任编辑 黄厚坤

责任校对 冯莲凤

责任印制 潘玉玲

印刷 北京盛通印刷股份有限公司

经销 全国新华书店

开本 787mm×1092mm 1/16

印张 24.5

字数 466 千字

版次 2012 年 12 月第 1 版 2012 年 12 月第 1 次印刷

书号 ISBN 978-7-5022-5770-5

印数 1—2500

定价 88.00 元

网址: <http://www.aep.com.cn>

E-mail: atomep123@126.com

发行电话: 010-68452845

版权所有 侵权必究

《中国原子能科学研究院科学技术丛书》

出版委员会

主任 万 钢

副主任 周刘来 柳卫平

委员 (按姓氏笔画为序)

王国保	尹忠红	叶宏生	叶国安	刘振华	刘峻岭
刘森林	李育成	李和香	李树源	张天爵	张东辉
张生栋	陈 凌	陈东风	邵焕会	罗志福	岳维宏
赵志祥	姜兴东	柯国土			

编审委员会

顾问 (按姓氏笔画为序)

王方定	王乃彦	方守贤	阮可强	张焕乔	周永茂
徐 铼	钱绍钧	樊明武	潘自强		

主任 赵志祥

副主任 柳卫平

委员 (按姓氏笔画为序)

于伟翔	王明政	王仲文	王修龙	尹邦跃	刘一兵
刘大鸣	刘森林	汤秀章	杜开文	李志宏	杨红义
杨启法	吴继宗	何 辉	何高魁	肖雪夫	张天爵
张伟国	林 敏	赵永刚	赵守智	罗志福	季松涛
周培德	郑卫芳	姜 山	胡石林	侯 龙	顾忠茂
黄 晨	韩世泉	葛智刚	喻 宏		

丛书办公室

主任 尹忠红

副主任 李来霞

成 员 (按姓氏笔画为序)

马英霞	王宝金	伍险峰	张小庆	张徐璞	骆淑莉
韩翠娥					

《中国原子能科学研究院科学技术丛书》

编辑工作委员会

主任 侯惠群

副主任 杨树录

委员 (按姓氏笔画为序)

丁怀兰 王艳丽 刘 朔 李 宁 杨树录 张关铭
张铨清 赵志军 侯惠群 谭 俊

编辑工作小组

组 长 杨树录

副组长 丁怀兰 赵志军

成 员 (按姓氏笔画为序)

丁怀兰 王艳丽 刘 朔 李 宁 杨树录 张关铭
张铨清 赵志军 谭 俊

总 序

中国原子能科学研究院创建于1950年,是我国核科学技术的发祥地和先导性、基础性、前瞻性的综合性核科学技术研究基地。

在党中央和上级部门的关怀和指导下,中国原子能科学研究院为我国的国防建设、国民经济建设和核科学技术的发展做出了重要贡献,造就了7位“两弹一星”功勋科学家和60多位两院院士,培养了大批科技人才,在核物理、核化学与放射化学、反应堆工程技术、加速器工程技术、同位素技术、核电子学与核探测技术、辐射防护、放射性计量等学科形成了自己的特色和优势,并拥有核科学与技术 and 物理学两个一级学科硕士、博士学位授予权。

为了系统地总结原子能院在核科学技术相关优势学科积累的知识和经验,吸收和借鉴国内外核科学技术最新成果,促进我国核科技事业的发展,我院决定组织出版《中国原子能科学研究院科学技术丛书》,并选定王淦昌、肖伦、丁大钊、王乃彦、阮可强等院士编著的《惯性约束核聚变》、《放射性同位素技术》、《中子物理学——原理、方法与应用》、《新兴的强激光》、《核临界安全》5本专著首批出版,今后还将组织撰写更多的学术专著纳入本丛书系列。

谨以此套丛书献给为我国核科技事业献身的人们!

《中国原子能科学研究院科学技术丛书》出版委员会

2005年9月1日

前 言

在一切做定量测量和判断的科学研究及其应用的领域,都要认真考虑有关量值的误差及这些误差在量值传递过程中的性质和行为。典型的例子是在核科学技术及其应用的领域,在那里,不仅数据本身,而且数据的误差都受到了越来越多的重视。数据和数据的误差紧密地与核工程的经济性和安全性等重大问题联系起来。为了确保核工程与核产品的安全性,在设计上必须要留出余量,而为了获得更好的经济性,这个余量又不能留得太多。在核工业的许多工程领域,关键数据的误差哪怕减少一点,都会带来巨大的经济效益。在安全性和经济性之间的这种博弈,使得对数据的误差分析和处理过程的科学研究工作具有重大的现实意义。

从20世纪70年代以来,由于核工程设计能力和水平的提高,微观核数据的误差往往居于工程设计计算量误差的主导地位。这一事实不仅推动了核数据工作的发展,而且导致了世界范围的核数据专门队伍的形成,促进了核数据处理和误差分析研究工作的开展,而大型计算机的发展和快速计算能力的提高使得上述进展成为可能。

作者从20世纪70年代末期开始进入了核数据处理和误差分析这一领域,并在这一领域持续工作了三十多年。因而,本书在传统的章节中在下列方面相当多地融入了作者多年的研究成果和观点:

关于误差的分类 数据的误差可以分成偶然误差和关联误差两大类,它们在量值传递过程中的性质有着明显的差异。对于数据评价和使用者来说,还可能遇到一种不包括在报道误差中的疏失误差,它多半也是一种关联误差,只是由于难以在一个测量系统上发现而被疏失了。

宏观测量系统和微观测量系统 在实验数据的数学处理中,三种不同性质的误差在一定前提下可以认为都具有偶然误差的随机性质。微观测量系统和宏观测量系统的引进有助于描述和理解三类误差性质。在同一微观系统上对同一物理量做重复测量,测量的关联误差不因测量

次数的增加而减少,在同一微观系统上对某一物理曲线做测量时,测量曲线的关联误差不会因测量点的加密而减少,在同一微观系统上甚至不能发现可能存在的关联误差。而每一微观系统中的关联误差在宏观系统中具有偶然误差的性质。基于这一点,尽可能地设计独立于其他微观系统的实验可以使得关联误差随机化。

协方差矩阵 20世纪70年代中期以前,关联误差在进行数据处理和误差分析时常常被忽略。这不仅反映在有关的专著中,也反映在一些实际使用的数据处理程序中。这之后,情况发生了变化。在工程设计中,不考虑数据的关联误差便不能正确地估计计算量的误差。这就要求在数据处理的全部过程中考虑全部的误差信息—协方差矩阵,同时要求实验工作者做完整的误差报告,即报告包括各种关联误差的协方差矩阵。

协方差矩阵的引入对过去常用的数据处理和误差分析方法的适用性提出了挑战。本书建立了适用于用协方差矩阵表示的关联误差存在时的数据处理方法和公式体系。

当处理带协方差矩阵的数据的合并与拟合时,将会遇到一个高阶矩阵的求逆问题,计算上费时很多,而且容易出现病态矩阵问题,本书给出了一些在适当近似下解决高阶矩阵求逆的方法和技巧。

数据空间和导出量空间 在数据空间和导出量空间应用最小二乘法得到的结果是等价的,但是迭代的方式不同。按照人们的传统习惯在导出量空间处理数据非常容易犯错误,因此要特别小心。作者强烈建议尽可能在数据空间处理数据。

导出量处理中的 PPP 问题 数据处理的对象常常是导出量,其为直接测量量的函数。当考虑关联误差时,使用现成的误差传递公式和最小二乘拟合公式处理导出量的合并及拟合时,常常发生极其荒谬的现象。这一由美国 ORNL 实验室的 R. Peelle 所遇到并提请国际核数据界注意的问题,即所谓 PPP(Peelle's Pertinent Puzzle)问题,困扰了国际核数据界十余年。PPP 总是在导出量空间发生,是不正确地理解和使用误差传递公式和最小二乘法造成的。作者对这一问题的解决,使相关数据的处理实际成为可能,受到了国际核数据界的重视。

贝叶斯方法 对贝叶斯方法的实际应用,历史上一直存在争论,相

当长的一个时期贝叶斯方法受到“冷落”，一些统计学家甚至主张完全放弃它。近三十年来，贝叶斯方法有了很大的进展，由于群论方法的引进，采用贝叶斯假设即无差别原理确定验前分布所引起的悖论均被解决。在包括核数据等在内的许多领域，贝叶斯方法得到了越来越多的应用。本书介绍了现代贝叶斯统计推断的基本概念和应用实例。

模型理论计算值的误差 在本书中讨论这一问题是因为在某种意义上它可以等价为一个实验数据的数学处理问题。和曲线拟合一样，需要解决的是如何由实验数据的协方差矩阵估计参数的协方差矩阵，不同的是，还必须考虑到因理论模型的缺陷带来的不确定性对理论预言值的影响。本书基于作者的工作经验对这一问题进行了初步的探讨。

以上可以认为是体现本书特色的一些方面。

为了使读者易于理解本书所述的概念和结论，作者给出了一些数据测量和评价的例子。这些实例大多是实际问题的简化。本书的大部分章节曾作为教材由作者在核工业研究生院讲授过。我们希望本书成为数据测量、评价和使用者掌握数据处理方法的手册，它应当也是其他领域科技工作者的一本有价值的参考书。

由于水平有限，本书难免存在缺点和错误，恳望读者不吝指正。应当指出的是，本书是按照实验测量和数据评价的观点和要求写成的。有些问题的表述在数学家看来也许是不严格的。比如当写出一个积分时，并不去说明它是否收敛，当对一个函数求导时，也并不去证明这个函数的导数一定存在等。作者希望本书能给读者一个清晰的思路和实验测量、数据评价工作者一个能够理解的结论。

作 者

2011年10月

目 录

第 1 章 随机事件和随机变量	1
1.1 随机事件及其概率	1
1.1.1 随机事件	1
1.1.2 随机事件概率的定义	1
1.1.3 事件之间的关系	3
1.1.4 随机事件的几个概率公式	3
1.2 随机变量及其概率分布	6
1.2.1 随机变量	6
1.2.2 离散分布及概率分布列	6
1.2.3 连续分布及其概率密度函数	7
1.2.4 多元随机变量及联合概率密度函数	8
1.2.5 随机变量的变换:随机变量函数	9
1.2.6 概率分布的数字表征	13
1.2.7 随机变量的特征函数	18
第 2 章 几种常见的概率分布	20
2.1 二项分布	20
2.2 泊松分布	21
2.2.1 泊松分布的形式	21
2.2.2 二项分布的泊松近似	22
2.2.3 泊松变量的再现性	24
2.3 均匀分布	25
2.3.1 等可能性和均匀分布	25
2.3.2 均匀分布的应用	27
2.4 正态分布	30
2.4.1 正态分布的形式	30
2.4.2 正态分布成立的条件:中心极限定理	31
2.4.3 正态变量的线性函数:正态变量的再现性	35
2.5 多元正态分布	35
2.5.1 多元正态分布的形式及一些性质	35

2.5.2	关于二元正态分布的一些讨论	37
2.5.3	N 维正态向量的线性变换	39
2.6	其他常用的导出分布	43
2.6.1	χ^2 分布	43
2.6.2	t 分布	45
2.6.3	F 分布	47
2.6.4	指数分布	47
2.6.5	瑞利分布	49
2.6.6	截尾分布	49
第 3 章	实验测量的误差分析	51
3.1	实验测量的误差	51
3.1.1	微观系统和宏观系统	51
3.1.2	误差的定义	51
3.1.3	误差的分类	52
3.1.4	测量的不确定度及其合成	56
3.2	误差的传递	57
3.2.1	正态性假定	57
3.2.2	线性函数的误差传递	58
3.2.3	线性变换的误差传递	58
3.2.4	非线性函数的误差传递	59
3.2.5	非线性变换的误差传递	62
3.3	实验数据的协方差矩阵	64
3.3.1	协方差矩阵: 实验数据的完整误差报道	64
3.3.2	实验测量值的协方差矩阵的构造	64
3.3.3	构造协方差矩阵要注意的几个问题	66
3.4	含二阶修正的误差传递公式	68
第 4 章	总体参数的估计	72
4.1	估计量及其性质	72
4.1.1	总体和样本	72
4.1.2	估计量及其分布	72
4.1.3	估计量的一些性质	73
4.1.4	估计结果的报道	75
4.2	总体参数的估计值: 点估计	75
4.2.1	矩法估计值	75

4.2.2	最小方差估计值	88
4.2.3	最大似然估计值	93
4.2.4	对不独立观测样本的进一步讨论	101
4.3	区间估计	107
4.3.1	置信区间和置信区域	107
4.3.2	正态总体平均值的区间估计	109
4.3.3	任意分布总体参数的区间估计	114
4.3.4	大样本时的区间估计	117
4.3.5	未知总体分布平均值的区间估计	117
第 5 章	假设检验	119
5.1	假设检验的基本概念	119
5.1.1	统计假设	119
5.1.2	实际判断原理:小概率原理	120
5.1.3	假设检验的基本步骤	120
5.2	参数性假设检验	123
5.2.1	正态总体参数的检验	123
5.2.2	非正态总体参数的检验	130
5.3	非参数性假设检验	132
5.3.1	χ^2 检验	132
5.3.2	符号检验; N 检验	136
5.3.3	游程数检验; R 检验	138
5.3.4	总体分布的拟合性检验	141
5.4	最佳检验	143
5.4.1	临界值和拒绝域的不唯一性	143
5.4.2	假设检验的两类错误	145
5.4.3	似然比检验	147
5.4.4	最佳检验中样本容量和测量误差的影响	152
5.5	异常数据的舍弃	154
5.5.1	实验数据的异常	154
5.5.2	实验数据系综的物理分析	155
5.5.3	实验数据系综的统计学分析	156
第 6 章	最小二乘法	159
6.1	实验观测曲线的光滑和拟合	159
6.2	公式类型的确定	160

6.2.1	四类公式	160
6.2.2	一般多项式	160
6.2.3	正交多项式	161
6.2.4	勒让德多项式	166
6.3	最小二乘原理	167
6.4	曲线的光滑	168
6.4.1	滑动平均法	168
6.4.2	样条函数拟合	173
6.5	曲线拟合:无约束最小二乘法	177
6.5.1	线性函数	177
6.5.2	非线性函数	182
6.6	曲线拟合:约束条件下的最小二乘法	185
6.6.1	不等式约束	185
6.6.2	线性约束	187
6.6.3	非线性约束	193
6.7	最小二乘法的递推公式	195
6.8	考虑自变量不确定性的最小二乘拟合	199
6.8.1	参数变换法	199
6.8.2	拉格朗日乘子法	203
6.8.3	近似方法	206
6.9	多元变量的拟合	207
6.10	导出量的最小二乘拟合	210
6.10.1	最小二乘拟合中的 PPP 现象	210
6.10.2	数据空间和导出量空间	212
6.10.3	数据空间的解	214
6.10.4	导出量空间的解	216
6.10.5	两个空间的解的比较	219
6.10.6	计算实例	221
6.11	最小二乘拟合的质量	226
6.11.1	拟合优度检验(χ^2 检验)	228
6.11.2	拟合多项式的最佳阶数	229
6.12	最小二乘估计值协方差矩阵的概率意义	233
6.12.1	直线拟合	234
6.12.2	线性函数的置信界限	241

6.13	曲线拟合中高阶矩阵的求逆问题	244
6.14	系统学研究中经验公式选择的一些方法	249
6.14.1	相关分析法	249
6.14.2	图示分析法	250
6.14.3	量纲分析法	251
第7章	统计推断的贝叶斯方法	256
7.1	经典统计学和贝叶斯统计学	256
7.2	随机参数及其分布	257
7.3	验前分布的确定	265
7.3.1	经验贝叶斯方法	266
7.3.2	用理论给出验前分布	266
7.3.3	最大信息熵方法	267
7.3.4	群论方法	269
7.4	随机参数的估计	273
7.4.1	估计的损失函数和风险函数	273
7.4.2	极大验后估计值	274
7.4.3	最小方差估计值	276
7.4.4	区间估计	278
7.5	贝叶斯方法在曲线拟合中的应用	280
7.6	贝叶斯方法在数据评价中的应用	282
7.7	贝叶斯方法在数据调整中的应用	282
7.8	随机参数的假设检验	284
7.9	贝叶斯方法的其他应用	286
第8章	模型理论计算值的不确定性	290
8.1	模型理论及参量化	290
8.2	特适参数及其协方差矩阵	291
8.3	普适参数及其协方差矩阵	298
8.3.1	普适参数不确定性的来源	298
8.3.2	标度因子法	301
8.3.3	残差估计法	301
8.3.4	矩法估计	302
8.3.5	计算实例	303

第 9 章 数据处理中的数值计算方法	307
9.1 超越方程求根	307
9.1.1 牛顿法	307
9.1.2 简单迭代法	309
9.1.3 初值的选取	310
9.2 函数插值	311
9.2.1 线性插值	311
9.2.2 拉格朗日插值	312
9.3 数值微分	313
9.3.1 插值微分公式	313
9.3.2 结点处的导数值	313
9.4 数值积分	314
9.4.1 插值求积公式	314
9.4.2 等距结点	315
9.4.3 不等距结点:高斯型求积公式	317
9.5 概率统计在计算方法上的一些应用	321
9.5.1 蒙特卡罗方法和随机数	321
9.5.2 随机数的检验	324
9.5.3 任意分布随机数的产生	326
9.5.4 定积分的概率计算方法	331
9.5.5 随机游动问题的模拟	336
第 10 章 实验设计初步	338
10.1 实验设计的基本概念	338
10.1.1 实验的种类	338
10.1.2 实验设计的要素	339
10.1.3 实验设计的原则	340
10.2 单因素优化实验设计	341
10.2.1 均分法	341
10.2.2 对分法	342
10.2.3 黄金分割法	342
10.3 多因素优化实验设计	343
10.3.1 因素轮换法	343
10.3.2 随机化安排实验	344
10.3.3 拉丁方设计	344

10.4	正交设计	346
10.4.1	正交性和正交表	346
10.4.2	正交表的方差分析	348
10.4.3	有交互作用的正交设计	351
10.5	均匀设计	355
附表 1	标准正态分布表	361
附表 2	χ^2 分布表	363
附表 3	t 分布表	365
附表 4	F 分布表	367
附表 5	符号检验表	373
参考文献		375

第 1 章 随机事件和随机变量

1.1 随机事件及其概率

1.1.1 随机事件

随机事件是指统计意义下的偶然事件,它在一次试验中可能发生,也可能不发生。我们无法肯定地回答一个随机事件在一次试验中是否能够发生。例如,我们无法肯定地回答“某一放射性原子核在时间 t 内是否衰变”或“某一测量值是否高于真值”这样一些问题。在这两个例子中,“某一放射性原子核在时间 t 内衰变”或“某一测量值高于真值”都是随机事件。随机事件并不是毫无规律可循,所谓“随机”是指一次试验而言。如果多次重复试验,随机事件就会显示出很好的统计规律性。有时把随机事件简称为事件。

1.1.2 随机事件概率的定义

定量地描述随机事件 A 发生的可能性大小的量称为随机事件 A 的概率,记为 $Pr(A)$ 。

在概率论发展的历史上,曾经有几种定义概率和计算概率的方法。

古典概率

设某一随机试验中,可能发生的全体基本事件总数是有限的,且具有相等的可能性,则对于任意事件 A ,若基本事件总数为 n ,事件 A 包含的基本事件数目为 k ,则事件 A 发生的概率由下式计算

$$Pr(A) = \frac{k}{n} \quad (1.1.1)$$

上式称为**古典概率**。

例 1.1 短袜问题

2 只黄色短袜和 8 只白色短袜混放在抽屉中,随机地抽出三只,得到一双黄色短袜的概率有多大?

$$\text{基本事件总数: } n = C_{10}^3 = \frac{10 \times 9 \times 8}{1 \times 2 \times 3} = 120$$

抽出 2 只黄色短袜加任何一只白色短袜的基本事件数目: