

多元统计分析

——R与Python的实现

Multivariate Statistical Analysis
with R and Python

吴喜之 编著

k-Means Clustering of 2 Colours



k-Means Clustering of 4 Colours



k-Means Clustering of 16 Colours



k-Means Clustering of 32 Colours





基于R应用的统计学丛书

多元统计分析

——R与Python的实现

*Multivariate Statistical Analysis
with R and Python*

吴喜之 编著

中国人民大学出版社
· 北京 ·

图书在版编目 (CIP) 数据

多元统计分析: R 与 Python 的实现 / 吴喜之编著. — 北京: 中国人民大学出版社, 2019.1
(基于 R 应用的统计学丛书)
ISBN 978-7-300-26655-8

I. ①多 … II. ①吴 … III. 多元分析 - 统计分析 IV. ①O212.4

中国版本图书馆 CIP 数据核字 (2019) 第 016562 号

基于 R 应用的统计学丛书

多元统计分析——R 与 Python 的实现

吴喜之 编著

Duoyuan Tongji Fenxi——R yu Python de Shixian

出版发行	中国人民大学出版社		
社 址	北京中关村大街 31 号	邮政编码	100080
电 话	010-62511242(总编室)		010-62511398(质管部)
	010-82501766(邮购部)		010-62514148(门市部)
	010-62515195(发行公司)		010-62515275(盗版举报)
网 址	http://www.crup.com.cn http://www.ttrnet.com (人大教研网)		
经 销	新华书店		
印 刷	北京鑫丰华彩印有限公司		
规 格	185mm×260mm 16 开本	版 次	2019 年 1 月第 1 版
印 张	20.75 插页 1	印 次	2019 年 1 月第 1 次印刷
字 数	484 000	定 价	39.80 元

版权所有	侵权必究	印装差错	负责调换
------	------	------	------

前言

目前,多数国内多元分析教材有下面几个特点:(1)着重于数学理论和推导.这往往使得读者忘记了分析数据的本来目的,把数据科学当成数学来学.(2)关于多元正态分布的性质占据很大篇幅.这实际上过分强调了对数据的狭隘数学假定,而忽略了真实数据的复杂性和多样性.(3)对真实的数据分析强调得不够.一些典型的教材数据往往和现实世界差得很远,对读者具有误导性.(4)在所涉及的统计软件应用方面,大多只限于商业软件.笔者认为,使用盗版国外软件不仅犯法,而且最终会形成对这些盗版软件的依赖性而不能自拔.使用“傻瓜式”商业软件对于学习编程没有任何好处,反而会阻碍人们获得编程能力,而是否会编程将如是否文盲一样影响人们的未来发展.

本书力图用简单通俗的语言阐述有关的基本概念,使得非数学背景的读者亦可以较容易地理解.考虑到许多读者的主要目标是应用,本书通过案例来帮助理解各种方法的概念及计算.本书也包括了传统多元分析课本所具有的数学公式,目的在于方便读者查询.

本书基于实际数据分析,在正文中通过 R 软件解释方法及计算输出,同时在每一章都给出用 Python 3 计算这一章数据例子的代码.

本书包括传统多元分析教材通常包括的内容:主成分分析、因子分析、聚类分析、判别分析、典型相关分析以及非传统的对应分析方法.一些多元分析教材还包括完整的多重回归内容,但这远远超出一个学期的课程内容.本书对多元回归予以回顾,介绍了部分机器学习回归方法,判别分析则安排在“分类”一章,并且介绍了部分机器学习分类方法.关于回归和分类的更多内容及细节,请读者参看笔者编写的由中国人民大学出版社出版的《应用回归及分类——基于 R》一书.

应用统计方法的学习,应该是拉动式学习,也就是遵从问题驱动或者数据驱动模式:从数据出发,需要什么就学什么.如此,所有学习都是目标明确的,学的都是有用的知识和能力.有的教师喜欢以“打基础”的名义给学生一大堆书去读,这样的教学方法事倍功半.

和其他自然科学一样,统计需要从数学课程中获得严密的逻辑思维能力,还需要计算机编程能力.这不仅需要学会针对一两个统计软件的编程,还需要获得能迅速掌握任何开源软件的**泛型编程能力**.有了来自数学的思维能力、想象力和逻辑,以及计算机编程能力,再加上学习自然科学所必需的批判性思维及对所面对实际问题的理解,你就拥有了研究统计最重要的必要条件.

本书的排版是笔者用 L^AT_EX 软件实现的,一切错误由笔者负责.

吴喜之

目 录

前 言

第 1 章 引 言

1.1 数据科学	1
1.1.1 统计是数据科学吗?	1
1.1.2 计算机学科在数据科学中的地位	2
1.1.3 问题驱动应成为数据科学的基本思维方式	2
1.2 多元分析的对象	2
1.3 需要的工具	3
1.4 各章的安排	3
1.5 软件和编程	4
1.6 如何教学	4

第 2 章 矩阵代数回顾

2.1 矩 阵	6
2.1.1 基本定义	6
2.1.2 基本矩阵运算	7
2.1.3 行列式	8
2.1.4 矩阵的逆	9
2.1.5 矩阵的广义逆	9
2.1.6 Kronecker 积	10
2.1.7 幂等矩阵	10
2.1.8 向量空间	11
2.1.9 正交性	11
2.1.10 矩阵的秩	11
2.1.11 矩阵的迹	12
2.1.12 特征值	12
2.1.13 广义特征值	13
2.1.14 分块矩阵	13
2.2 矩阵的分解	15
2.2.1 矩阵的特征值分解	15
2.2.2 奇异值分解及广义奇异值分解	16
2.2.3 QR 分解	18
2.2.4 Cholesky 分解	19

2.3 二次型	20
2.3.1 定义	20
2.3.2 二次型和矩阵的定性	20
2.3.3 椭球	21
2.4 矩阵的导数	22
2.4.1 向量关于数量的偏导数	22
2.4.2 数量关于向量的偏导数	22
2.4.3 向量关于向量的偏导数	22
2.4.4 矩阵关于数量的偏导数	22
2.4.5 数量关于矩阵的偏导数	23
2.4.6 有关内积、二次型的导数	23
2.4.7 函数的偏导数	23
2.5 本章矩阵简单运算的 R 和 Python 代码	24
2.6 习 题	27

第 3 章 回 归

3.1 经典回归模型基本要素	29
3.1.1 描述数据	30
3.1.2 线性回归模型	31
3.1.3 最小二乘回归	32
3.1.4 最小二乘回归参数的估计	34
3.2 交叉验证	36
3.3 经典线性回归的基本假定及根据假定所得到的结论	38
3.3.1 经典线性回归的数学假定	38
3.3.2 最小二乘回归参数估计量的一些性质及有关检验	39
3.3.3 例 3.1 波士顿住房数据的最小二乘回归	42
3.3.4 最小二乘回归结果的解释	43
3.3.5 有关拟合的 R^2 及 AIC	46
3.4 自变量有分类变量(定性变量)的情况*	48
3.5 机器学习回归简介及案例	52
3.5.1 例子和机器学习方法的优势	52
3.5.2 决策树回归	54
3.5.3 bagging 回归	57

3.5.4	随机森林回归	58	4.6	机器学习模型简介及案例	106
3.6	经典线性回归与各种机器学习回 归的比较	62	4.6.1	决策树分类模型	107
3.6.1	对例 3.1 波士顿住房数据的九种 方法的比较	62	4.6.2	随机森林分类模型	109
3.6.2	对例 3.1 波士顿住房数据使用 caret 包比较各种方法 *	64	4.6.3	adaboost 分类模型	114
3.7	多元多重回归回顾 *	69	4.6.4	adaboost 拟合例 4.1 数字笔迹 识别全部数据	115
3.7.1	模型	69	4.6.5	adaboost 的变量重要性	115
3.7.2	最小二乘估计及对估计量的检验	70	4.6.6	adaboost 的交叉验证	115
3.7.3	多元回归检验例子	71	4.6.7	各种分类模型的比较	116
3.8	本章计算的 Python 代码	76	4.7	本章计算的 Python 代码	118
3.8.1	例 3.1 波士顿住房数据的线性回 归	76	4.7.1	产生 Z 折交叉验证集的代码	118
3.8.2	例 3.1 波士顿住房数据的决策树 回归、bagging 回归和随机森林回归	79	4.7.2	线性判别分析及交叉验证	120
3.8.3	把分类变量转换成哑元变量	84	4.7.3	二次判别分析	121
3.9	习题	86	4.7.4	Logistic 回归	122
			4.7.5	对例 4.2 献血数据做 3 个经典 分类方法的交叉验证	124
			4.7.6	随机森林分类	125
			4.7.7	对例 4.1 数字笔迹识别数据做 5 个机器学习方法分类的交叉验证比较	125
			4.8	习题	128
第 4 章 分类			第 5 章 主成分分析		
4.1	目的和基本概念	89	5.1	基本思想和例子	129
4.1.1	目的和例子	89	5.1.1	数据中变量之间的关系与降维的 可能性	129
4.1.2	分类和聚类的区别	90	5.1.2	从例子中产生的问题	130
4.2	经典判别分析	90	5.2	问题和计算	133
4.2.1	线性判别分析和二次判别分析	90	5.2.1	求协方差矩阵及样本相关阵的特 征值和特征向量	134
4.2.2	Fisher 判别分析	92	5.2.2	求各个成分的贡献率和累积贡献 率	135
4.3	例 4.1 数字笔迹识别的判别分析 实践	93	5.2.3	根据贡献率选取少数重要成分	138
4.3.1	例 4.1 数字笔迹识别的线性判别 分析的计算	93	5.2.4	主成分和原始变量的相关系数计 算与载荷图	138
4.3.2	例 4.1 数字笔迹识别的二次判别 分析的计算	94	5.2.5	主成分得分	139
4.4	判别分析的一些概念和某些数学 细节的补充 *	95	5.3	主成分分析在图像处理中的应用	142
4.4.1	数据总体分布已知时的判别	95	5.3.1	图像压缩案例 (例 5.2 江景图 片)	142
4.4.2	线性判别分析的某些细节	97	5.3.2	主成分分析运用于图像识别	145
4.4.3	Fisher 线性判别分析的某些细 节	98	5.4	案例: 例 5.3 数据的主成分分析	145
4.5	二分类因变量的 Logistic 回归	99	5.5	主成分分析的一些数学知识 *	148
4.5.1	Logistic 回归模型	99	5.5.1	目标	148
4.5.2	Logistic 回归模型对实际数据 的拟合	101	5.5.2	过程	148
4.5.3	Logistic 回归对例 4.2 献血数 据的分类	102	5.5.3	累积贡献率	150
4.5.4	ROC 等描述性曲线	103	5.5.4	总体主成分的一些性质	151
4.5.5	三种分类方法对例 4.2 献血数据 的交叉验证比较	106	5.5.5	样本主成分的一些性质	151
			5.5.6	主成分的解释	152
			5.5.7	标准化主成分分析	152

5.6	本章计算的 Python 代码	153	7.4	基于模型的聚类 *	205
5.6.1	例 5.1 教师数据	153	7.4.1	直观描述	205
5.6.2	图像压缩例子	158	7.4.2	E-M 算法	206
5.7	习 题	159	7.4.3	基于模型聚类的计算实现	208
第 6 章 因子分析			7.5	聚类数目的选择	211
6.1	基本内容	161	7.5.1	Gap 方法 *	211
6.1.1	概述	161	7.5.2	轮廓法 *	212
6.1.2	模型	163	7.5.3	计算聚类数目的不同软件	215
6.1.3	性质	164	7.6	实例的计算	220
6.1.4	旋转	165	7.6.1	种子数据	220
6.1.5	因子得分	165	7.6.2	胎心宫缩监护数据	225
6.1.6	计算及分析步骤	166	7.6.3	例 5.1 教师数据的聚类	228
6.2	例子和计算	166	7.7	图像色彩的聚类	230
6.2.1	例 6.1WHO 数据分析	166	7.7.1	对图像色彩聚类的例子	230
6.2.2	例 6.2 豌豆数据的因子分析	172	7.7.2	颜色的十六进制表示形式 *	233
6.2.3	例 5.1 教师数据的因子分析	174	7.8	本章计算的 Python 代码	234
6.3	因子分析计算基于的原理 *	176	7.8.1	例 7.2 美国各州数据的聚类	235
6.3.1	主成分方法	176	7.8.2	例 7.6 种子数据的聚类	238
6.3.2	主因子方法	177	7.8.3	图像色彩的聚类	241
6.3.3	最大似然法	177	7.9	习 题	243
6.3.4	因子个数的选择	178	第 8 章 典型相关分析		
6.3.5	旋转	178	8.1	基本内容	245
6.3.6	因子得分	179	8.1.1	问题的描述和例子	245
6.3.7	建议和评论	180	8.1.2	典型相关及特征值问题	247
6.4	本章计算的 Python 代码	180	8.1.3	载荷	252
6.5	习 题	182	8.1.4	典型变量的选择	253
第 7 章 聚类分析			8.2	实例计算: 例 5.1 教师数据	257
7.1	目的及基本概念	184	8.3	典型相关分析的不同角度推导 *	262
7.1.1	目的	184	8.3.1	完全利用特征值问题的性质	262
7.1.2	聚类和分类方法没有必然联系	184	8.3.2	利用 Cauchy-Schwarz 不等式	263
7.1.3	距离	186	8.3.3	利用奇异值分解	264
7.1.4	两个数据例子	191	8.3.4	利用广义特征值问题及 Cholesky 分解	265
7.1.5	集群倾向的度量	193	8.3.5	各种公式的总结	265
7.2	分层聚类	195	8.4	本章计算的 Python 代码	266
7.2.1	对观测值分类: 例 7.2 美国各州数据	195	8.5	习 题	267
7.2.2	对观测值聚类和对变量聚类: 例		第 9 章 对应分析		
7.3	花卉数据	198	9.1	基本内容 and 应用	268
7.3	K 均值聚类和基于密度的聚类	200	9.2	二元对应分析	268
7.3.1	K 均值聚类的基本思想	200	9.2.1	二元对应分析的数学原理	271
7.3.2	K 均值聚类中类别数目的确定	201	9.2.2	对应分析行列相关的显著性 *	272
7.3.3	基于密度聚类的思想 *	203	9.2.3	两个解释图形	273

9.3 多元及联合对应分析	274	10.2 多维尺度变换分类	285
9.3.1 例 9.1 收入数据多元及联合对应 分析实践	274	10.3 距离函数	285
9.3.2 多元对应分析方法的数学 * ..	277	10.4 mMDS 的目的	286
9.3.3 基于指标矩阵和 Burt 矩阵的方 法	277	10.5 mMDS 方法的原理	287
9.3.4 从 Burt 矩阵到联合对应分析	280	10.6 nMDS 的基本思维	288
9.4 本章计算的 Python 代码	280	10.7 本章计算的 Python 代码	290
9.4.1 二元对应分析	280	10.8 习题	292
9.4.2 多元对应分析	282	附录 1 R 简介——通过运行来领悟	293
9.5 习题	282	附录 2 Python 简介——通过运行来 领悟	305
第 10 章 多维尺度变换		参考文献	319
10.1 目的	284		

第 1 章 引 言

1.1 数据科学

1.1.1 统计是数据科学吗？

按照《不列颠百科全书》所述，统计是收集、分析、展示及解释数据的科学 (Statistics, the science of collecting, analyzing, presenting, and interpreting data).¹ 也就是说，统计是数据科学 (the science of data)。早在 20 年前，华裔统计学家吴建福 (C.F. Jeff Wu) 就多次说过“统计应该称为数据科学”。

但是，无论口头如何说，实际上统计界许多人是在搞数学，并冠之以“理论统计”或“数理统计”的名字。这是有历史原因的。计算机出现之前，人们只能处理少量数据，为了计算方便和弥补信息量的不足，以数学家为主体的学者创造了很多基于主观假定的数学模型。统计因此被数学所主导，统计方面的研究文章大部分是数学研究。《中国大百科全书》明确定义，数理统计是“数学的一个分支学科”。²

事实上，在前计算机时代，数据分析一直依赖于这些数学家所创造的数学模型。因此，模型驱动而非问题驱动一直是统计界的主流思维方式，模型驱动思维就是数学式的演绎思维。一直到 20 世纪 90 年代，计算机还仅仅是各个领域数据分析者的工具，是手和笔的延伸。如果在那个时候，统计学家认为自己是数据科学家，而把计算机视为自己不可缺少的工具包括在数据科学之内，统计学的进程就完全不一样了。可是，以数学模型为主导的统计界没有这么做。当今，统计界无法和真正的数学界比数学之美，也无法和计算机群体比处理数据问题的效率和能力，统计失去了很多热门领域和优秀人才。

机器学习或者基于机器学习方法的人工智能目前发展迅速，但有多少人会把这些发展和本应是“数据科学”的传统数理统计相联系呢？尽管现在一部分“大数据中心”(或“研究院”“学院”等)的领导者都是专注于数学模型的学者，但那些计算机领域的人会臣服吗？

从深度、难度以及美学的角度来看，统计所用的数学和纯粹数学根本无法相比，与数学家在一起可能会令某些统计学人感到不舒服。其实，如果这些统计学人展示的不是数学而是分析数据的能力，就不会有这种感觉了，问题是这些统计学人究竟有没有数据处理的能力。类似地，部分统计学人和搞计算机的同事在一起时也会有不快的感觉。如果他们不会编程计算却又专注于一些没有实践背景的数学模型，自然会感到不爽。

对于整个社会来说，传统统计的边缘化并不见得是个坏事情。对于统计学科本身，危机感能够促进变化，把统计界从封闭的小圈子里解放出来。

¹不列颠百科全书。 <https://www.britannica.com/science/statistics>。

²中国大百科全书总编辑委员会。中国大百科全书：数学卷。北京：中国大百科全书出版社，1988：593。

1.1.2 计算机学科在数据科学中的地位

由于统计界的主流实际上不愿意采取问题驱动或数据驱动的思维，而计算机学科借助数据源、数据清理、数据处理及算法改进等方面的优势，已经成为数据科学的主要实体，而且占有越来越重要的地位。

实际上，主要的机器学习方法的创造者大都不是沉溺于数学模型的“统计学家”，也不是计算机界的学者。他们多是有数学背景或其他领域背景的科学家。他们依赖的并不是某些专业课程知识，而是自己对实际问题的理解、对随机性的理解、抽象思维能力和想象力，而这些思维能力很多是从严格的纯粹数学训练中获得的。

在当今世界，人们离不开计算机，计算机是数据科学最重要的基石之一，也是工作量最大的部分，没有计算机，所有的机器学习、人工智能、大数据分析就都没有意义。但是必须清醒地认识到，单纯的“软件 + 硬件”能力不会自动产生新的算法，工程师式的工作代替不了科学家的创造。

1.1.3 问题驱动应成为数据科学的基本思维方式

笔者觉得目前的统计学科、数学或计算机学科没有一个能够单独代表数据科学。数据科学本身包括了所有可能的知识积累和人们的才智。知识不是由学科划分的，当前科技进步的进程不断突破人们有限的想象力，试图当老大或者试图把学科标签固化的思维曾经阻碍了一些学科的发展，一些故步自封的学科将来也会消亡，但社会总是在不断前进。

很多人认为，学完了统计、数学、计算机及各领域的一些课程就可以解决数据科学的问题。实际上，课程仅仅是历史的一些记录，人们可以从中得到一些经验、启发及知识，而不可能得到能力。有人觉得，根据分析某些“教科书式的数据”的经验，按照现成模型可能会找到所谓“标准答案”，但是，在实际问题中不可能存在“标准答案”，针对一个问题可能有很多不同的答案，而每种答案都可能具有某种合理性。试图寻求一个“标准答案”本身就限制了人们的想象力、能力及对实际问题复杂性的判断。更为大量的实际问题是无法限定于一种或若干种方法（或组合）的，认为只要选择合适的模型和参数就应该得到“正确”或“精准”答案的想法是把世界简单化，真实世界比围棋机器人所面对的围棋世界要复杂得多。

因此，面对复杂的真实问题，我们对问题本身的意义及有关数据的理解、知识积累、想象力和处理复杂问题的能力就是关键。能力是从不断解决实际问题中获得的，是从无数失败中悟出来的，而不是从书本上学来的。

1.2 多元分析的对象

本书介绍的方法主要适用于规范的横截面数据，观测值及变量可以按照行列排成方阵。这里讨论的变量都是相关的，而且每个观测值都是同时度量的。根据机器学习术语，回归和分类属于**有指导（或有监督）学习**（supervised learning），也就是有因变量作为目标变量，主要目的是建立可以预测的模型，有指导学习的回归和分类是学习的核心，也是人工智能的基本组成部分；其他方法都属于**无指导（或无监督）学习**（unsupervised learning），其中没有目标变量，目的是揭示相关变量所提供的潜在重叠信息所代表数据的背景结构，因此，简化、降维、汇总是其主要特点。这部分多元统计分析基本上是探索性的数据分析，它可能提出一些假设，但检验往往并不是其主要目的。

从名称上看,只要多于一个变量的数据分析都应该属于多元分析.这实际上包括了除收集数据之外的几乎整个数据科学的内容.但是,传统的多元分析要狭窄得多,除了往往放在回归课程中讲授的回归分析之外,多元分析传统上是由无指导学习中的一组方法所构成的.这包括主成分分析、因子分析、聚类分析、典型相关分析和判别分析等,这些方法传统上多假定多元正态分布,但由于大多是描述性的,在方法及结果的解释方面并不完全依赖分布假定.再拓宽一些,还可能包括本书没有包括在内的多元方差分析及结构方程模型的协方差方法等,它们严重依赖于对数据的各种假定.

除了可以应用于所有领域并成为人工智能基础的有指导学习之外,属于无指导学习的经典多元分析技术的很大一部分最早应用于行为科学和生物科学,而现在已经扩展到各个领域,包括物理、化学、生物、医学、心理学、社会学、经济学、教育学、法律、商学、文学、传媒学、护理学、语言学、地质学等等.

读者可以通过本书理解各种多元分析的方法、目的、适用范围、局限性及数据分析的具体计算机实现,以及各种方法的关联及区别.在动手操作方面,期望读者能够根据具体数据选择合适的方法并且解释多元分析方法的计算机输出.希望读者能够通过编程来熟悉本书的概念,同时通过分析数据来熟悉编程.

1.3 需要的工具

学习本书,需要掌握一定的数学知识,主要是线性代数或矩阵理论等数学,这些知识在本书的第2章专门进行了回顾,可以自给自足;微积分用得不多(只有个别地方应用了导数的概念).读者有一些初等统计知识就够了.本书强调对模型和方法的理解,而不是数学证明和推理的细节,大多数被展示的结果未给出证明,因为直观的解释更加有说服力.本书所使用的R或者Python等编程语言完全可以在使用中学习.学习软件最快的方法就是使用,有问题可从网络上寻求答案,不必专门在课堂上学习.

1.4 各章的安排

本书从第2章开始的各章内容安排如下:第2章“矩阵代数回顾”,提供本书所需要的一些代数知识,供需要时查阅,这些内容一般不用在课堂上专门介绍(由于后续各章多少涉及该内容,所以没有打星号).第3、4章的内容为有指导学习:第3章“回归”,回顾经典回归分析的内容,同时简单介绍少数机器学习的回归方法,这一章可以自学,也可以在课堂上做介绍,该章的多元多重回归一节仅供参考之用;第4章“分类”,介绍经典判别分析,同时介绍少数机器学习分类方法.从第5章开始介绍无指导学习的概念、方法、目的及案例分析,包括第5章“主成分分析”、第6章“因子分析”、第7章“聚类分析”、第8章“典型相关分析”、第9章“对应分析”、第10章“多维尺度变换”.本书的最后设置了关于R和Python语言的附录.

此外,本书在正文中使用R语言来实现数据分析的计算并解释其输出的结果,从第2章开始的每一章都提供了例题计算的Python代码,以满足学习或使用Python语言读者的需要.

除了第1章,每章后面都提供了习题,帮助读者掌握各章的内容或计算.

1.5 软件和编程

在现代社会，一个人是否会编程密切关系到其生存状况，笔者认为，编程不限于一个软件的语言，而是要培养明白编程要点的泛型编程能力。

编程必须在应用中学，上专门的软件理论课学习编程远远不如针对具体的来完成一个数据分析任务效果好。在编程中要不怕犯错误，不断纠正错误是熟悉软件的最快途径。

R 软件的程序大多是统计学家编写的，对搞统计分析的人很友好。R 网站上不断更新并持续增加的 1 万多个程序包³展示了各种最新的数据分析方法。绝大多数新方法的计算机实现均首先出现在 R 网站上。这些程序包的程序都是透明的，可供审视和改写，这是任何商业软件无法比拟的，透明是对抗错误的最好方式。R 软件的优点是容易编写，帮助系统很优秀，而且对各种统计方法有比较细致的描述，输出结果丰富；缺点是软件本身处理较大数据的速度较慢，但目前许多快速编程软件都有针对 R 的接口，而且还有一些针对 R 包装的软件，使得速度慢不会成为其发展及受欢迎的障碍。

不像 R 软件主要用于数据分析，Python 是通用软件，可以完成任何计算机能够做的任务。由于 Python 比 C++ 更易编写，有限的速度差距又已经被一些包装软件所克服，Python 很快成为 C++ 最强有力的替代。Python 的编程非常灵活，虽然其帮助系统不如 R 软件，但由于它是开源软件，网上帮助很好用。Python 不是专门的统计软件，其统计方法资源远不如 R 软件多，而且对于一些模型，它的输出结果也不如 R 软件那么丰富。

人们通常想知道应该学习什么编程软件，笔者认为，学习任何编程软件都有益处，如果学会第一个编程语言花了一个星期，学会第二个可能只需三四天的时间。这里所说的“学会”是指能够完成你手中的任务，而不是指精通。软件是工具，我们能够使用它完成任务就行了，没有必要完全精通。学习编程的关键是遇到问题时知道去哪里寻求帮助，网络上的信息很全面，显然任何教师都无法独自掌握如此海量信息，在网上寻求帮助是编程者的惯例。软件在不断发展，我们必须跟上时代的潮流。

鉴于 R 和 Python 的区别，以及各种模型的选项和程序包或模块的编写过程及意图的不同，本书中使用两种软件处理同一个数据的输出结果可能不同。我们对数据和程序包或模块采用了一些比较随意的选择方式，希望这不会对可能有“洁癖”或“强迫症”心理的人造成不适。

1.6 如何教学

人都是在不断学习中成长的，而且绝大多数能力是自学的。回想一下，你有哪一能力是课堂上当场学会的（不用做作业）？教师应该引导学生产生兴趣及疑问、展开独立思考并理解世界，让学生关注世上的最新事物，还是应该告诉学生世上有什么？学生难道不会自己查阅资料和自学吗？离开学校就是知识获取过程的终结吗？教学是否等于播放陈旧的录音？不看讲义而能不错一个符号地把书上的内容抄写在黑板上是好老师的标准吗？教师自己有必要不断学习吗？学生挑战教师好不好？没有学生挑战说明什么？教师应如何面对学生（在许多方面）比自己强这样的现实？不带疑问和怀疑去上课的学生会有收获吗？课后学生没有疑问的课堂教学正常吗？

³在 2009 年底，R 网站上的程序包不到 1000 个。

思考了上面一些问题之后,关于本书的教学方式也就十分清楚了.最主要的是弄明白每个模型的直观含义及拟合数据的过程和原理.至于有关的数学公式,只需要介绍其目的及意义,而不必做过度的数学推导,这些数学公式大都是有百年历史的既定结论.实际上,解释这些模型的原理并不比推导数学公式更容易.本书中标有*号的部分可作为参考材料,教师可以根据需要来确定讲哪些内容.但是,部分没有标*号的内容也可以用于自学或作为参考材料. **教科书仅仅是教学的参考和工具,教科书的使用应该服从于获取知识和能力的需要,教科书不应是限制师生思维的教条和框框.**

我相信,任何读者只要愿意努力,就有可能成为一个名副其实的数据科学家.

第 2 章 矩阵代数回顾

本章的内容是非常经典的数学知识, 在几乎所有的高等代数书中都可以找到. 在本章安排的这些内容, 主要是方便读者查阅, 它们只是工具, 不是必须要在课堂上讲授. 由于后续各章多少会涉及本章内容, 所以本章内容没有标星号 (*).

2.1 矩阵

2.1.1 基本定义

矩阵是数值或数值变量的一个矩形数字阵, m 行 n 列 (亦称 $m \times n$ 维或 $m \times n$ 阶) 的矩阵 \mathbf{A} 的形式为:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

\mathbf{A} 的表示方法或记号有很多, 比如显示维数的 $\mathbf{A}_{(m \times n)} = \{a_{ij}\}$ 或者显示元素的 $\mathbf{A} = \{a_{ij}\}$. 此外, 根据不同的需要及方便起见, 元素下标的次序并没有统一规定, 比如, 元素 a_{ij} 的两个下标不一定必须是 i 代表行及 j 代表列, 也可能相反.

前述矩阵 \mathbf{A} 的**转置**记为 \mathbf{A}^\top (这是本书采用的记号) 或 \mathbf{A}' , 定义为:

$$\mathbf{A}^\top = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix}.$$

一个 $m \times 1$ 的矩阵称为**列向量**, 而一个 $1 \times n$ 的矩阵称为**行向量**. 比如 $1 \times m$ 矩阵 $\mathbf{a} = [a_1, a_2, \dots, a_m]$ 是一个行向量, 也记为 $\mathbf{a} = (a_1, a_2, \dots, a_m)$, 其元素为 a_1, a_2, \dots, a_m , 而其转置为列向量, 即

$$\mathbf{a}^\top = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix}.$$

也可以用列向量或行向量表示矩阵, 比如, 前述矩阵 \mathbf{A} 的 n 个列向量记为 $\mathbf{a}_j = (a_{1j}, a_{2j}, \dots, a_{mj})^\top$ ($j = 1, 2, \dots, n$), 那么矩阵 \mathbf{A} 可表示为 $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$; 如果

前述矩阵 A 的 m 个行向量记为 $\mathbf{a}_i = (a_{i1}, a_{i2}, \dots, a_{in})^\top$ ($i = 1, 2, \dots, m$), 那么, 矩阵 A 可表示为 $A^\top = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]$.

对于一个方阵 $A = \{a_{ij}\}$, 如果 $a_{ij} = a_{ji}$, $\forall i, j$ ¹, 则称该矩阵为**对称矩阵**.

如果一个方阵 $U = \{u_{ij}\}$ 的元素满足 $u_{ij} = 0$, $\forall j < i$, 则称该矩阵为**上三角矩阵**; 如果一个方阵 $L = \{l_{ij}\}$ 的元素满足 $l_{ij} = 0$, $\forall j > i$, 则称该矩阵为**下三角矩阵**.

如果一个方阵 $D = \{d_{ij}\}$ 的元素满足 $d_{ij} = 0$, $\forall j \neq i$, 则称该矩阵为**对角矩阵**, 这时可记 d_{ii} 为 d_i , 以及 $D = \text{diag}(d_1, d_2, \dots, d_m)$. 对角线元素全部为 1 的对角矩阵 $I = \text{diag}(1, 1, \dots, 1)$ 称为**单位矩阵**; 如果其维度为 n , 往往记为 I_n 或 $I_{(n \times n)}$. 元素全部是 1 的向量记为 $\mathbf{1} = (1, 1, \dots, 1)^\top$.

2.1.2 基本矩阵运算

两个同样维数的矩阵 A 和 B 之和定义为一个新矩阵 $C = A + B$, 其元素等于 A 和 B 相应元素之和, 即 $C = \{c_{ij}\} = \{a_{ij} + b_{ij}\}$. 同样定义两个矩阵之差. 一个数量 c 和一个矩阵 $A = \{a_{ij}\}$ 相乘意味着逐项相乘, 即 $B = cA = \{ca_{ij}\}$. $-A$ 定义为 A 的每个元素乘 -1 . 显然有下面的性质:

$$A + B = B + A, \quad A + (B + C) = (A + B) + C,$$

$$A - B = A + (-B) = -(B - A), \quad A - A = \mathbf{0},$$

$$(A + B)^\top = A^\top + B^\top, \quad (b + c)A = bA + cA,$$

$$A + \mathbf{0} = A, \quad cA = Ac, \quad -(-A) = A,$$

$$c(A + B) = cA + cB, \quad \mathbf{0}A = \mathbf{0}, \quad \mathbf{1}A = A.$$

如果一个矩阵的列数和另一矩阵的行数相同 (称为可相乘的), 比如 $A_{(m \times n)}$ 和 $B_{(n \times p)}$, 则它们的积定义为一个 $m \times p$ 矩阵 $C = AB = \{c_{ij}\}$, 这里

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}.$$

两个相同维列向量 $\mathbf{a} = (a_1, a_2, \dots, a_n)^\top$ 和 $\mathbf{b} = (b_1, b_2, \dots, b_n)^\top$ 的内积 (又称点积或点乘) $\mathbf{a}^\top \mathbf{b}$, 也记为 $\mathbf{a} \cdot \mathbf{b}$ 或者 $\langle \mathbf{a}, \mathbf{b} \rangle$, 定义为数量

$$\mathbf{a}^\top \mathbf{b} = \sum_{i=1}^n a_i b_i.$$

¹符号“ \forall ”是“对于所有”(for all)的简写.

关于矩阵乘法有下列性质:

$$\begin{aligned}(A+B)C &= AC+BC, \quad A(B+C) = AB+AC, \\ A_{(m \times n)}I_n &= I_m A_{(m \times n)}, \quad A_{(m \times n)}\mathbf{0}_{(n \times p)} = \mathbf{0}_{(m \times p)}, \\ \mathbf{0}_{(q \times m)}A_{(m \times n)} &= \mathbf{0}_{(q \times n)}, \quad AB^T = (BA^T)^T.\end{aligned}$$

分块矩阵 (见2.1.14) 的运算和矩阵运算类似, 只不过把块 (子矩阵) 当成矩阵元素, 比如, 对分块矩阵

$$A = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix} = \left[\begin{array}{ccc|cc} a_{11} & a_{12} & a_{13} & b_{11} & b_{12} \\ a_{21} & a_{22} & a_{23} & b_{21} & b_{22} \\ \hline c_{11} & c_{12} & c_{13} & d_{11} & d_{12} \\ c_{21} & c_{22} & c_{23} & d_{21} & d_{22} \\ c_{31} & c_{32} & c_{33} & d_{31} & d_{32} \end{array} \right],$$

有

$$A^2 = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix} \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix} = \begin{bmatrix} M_{11}^2 + M_{12}M_{21} & M_{11}M_{12} + M_{12}M_{22} \\ M_{21}M_{11} + M_{22}M_{21} & M_{21}M_{12} + M_{22}^2 \end{bmatrix}.$$

2.1.3 行列式

一个 2×2 矩阵

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

的行列式定义为

$$|A| = \det(A) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{21}a_{12},$$

而一个 3×3 矩阵

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

的行列式定义为:

$$|A| = \det(A) = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}.$$

一般的 $n \times n$ 矩阵 $\mathbf{A} = \{a_{ij}\}$ 的行列式定义为:

$$\begin{aligned} |\mathbf{A}| &= \det(\mathbf{A}) = a_{k1}\alpha_{k1} + a_{k2}\alpha_{k2} + \cdots + a_{kn}\alpha_{kn} = \sum_{j=1}^n a_{kj}\alpha_{kj}, \quad \forall k \\ &= \sum_{i=1}^n a_{ik}\alpha_{ik}, \quad \forall k, \end{aligned}$$

这里

$$\alpha_{ij} = (-1)^{i+j}\beta_{ij},$$

而 β_{ij} 为 \mathbf{A} 中没有第 i 行及第 j 列元素的 $(n-1) \times (n-1)$ 维子矩阵的行列式. 这里的 α_{ij} 称为 a_{ij} 的余子式.

对于 n 阶方阵 \mathbf{A} 和 \mathbf{B} , 以及 $c \in \mathbb{R}$, 有下列性质:

$$|c\mathbf{A}| = c^n|\mathbf{A}|, \quad |\mathbf{AB}| = |\mathbf{BA}| = |\mathbf{A}||\mathbf{B}|.$$

2.1.4 矩阵的逆

一个 n 阶方阵 \mathbf{A} 的逆也是 n 阶方阵, 记为 \mathbf{A}^{-1} , 满足 $\mathbf{A}^{-1}\mathbf{A} = \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}_n$. 当然, 矩阵的逆不一定存在, 如果 \mathbf{A} 的逆存在, 则称 \mathbf{A} 为非奇异的, 否则称它为奇异的. 方阵 \mathbf{A} 存在逆的充分必要条件是其行列式 $|\mathbf{A}| \neq 0$.

关于矩阵的逆, 有下列性质:

$$\begin{aligned} \mathbf{I}^{-1} &= \mathbf{I}, \quad (\mathbf{A}^{-1})^{-1} = \mathbf{A}, \quad (\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top \\ (\mathbf{AB})^{-1} &= \mathbf{B}^{-1}\mathbf{A}^{-1}, \quad (c\mathbf{A})^{-1} = c^{-1}\mathbf{A}^{-1} \\ \text{diag}^{-1}(d_1, d_2, \dots, d_n) &= \text{diag}(1/d_1, 1/d_2, \dots, 1/d_n), \end{aligned}$$

这里, \mathbf{A} 和 \mathbf{B} 为非奇异方阵, c 为数量.

2.1.5 矩阵的广义逆

如果 \mathbf{A} 为 $n \times m$ 矩阵, $\mathbf{y} \in \mathbb{R}^n$ 及 $\mathbf{x} \in \mathbb{R}^m$, 考虑方程组

$$\mathbf{Ax} = \mathbf{y}, \tag{2.1.1}$$

则在 $m = n$ 时, 如果 \mathbf{A} 的逆存在, 则方程组 (2.1.1) 有解 $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$; 如果 \mathbf{A} 不可逆或者 $m \neq n$, 则无解或者解不唯一. 这时可以定义 \mathbf{A} 的广义逆.

矩阵 \mathbf{A} 的广义逆 \mathbf{A}^- 满足下列条件:

$$\mathbf{AA}^-\mathbf{A} = \mathbf{A}. \tag{2.1.2}$$

容易验证, 如果方阵 \mathbf{A} 的逆 \mathbf{A}^{-1} 存在, 则 $\mathbf{A}^{-1} = \mathbf{A}^-$. 这种情况之外的广义逆并不唯一.