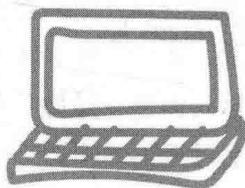


数据科学家 养成手册

高扬◎编著



电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

作为认知科学的延伸，数据科学一方面应该越来越引起广大数据工作者的重视，另一方面也要撩开自己的神秘面纱，以最为亲民的姿态和每位大数据工作者成为亲密无间的战友，为用科学的思维方式进行工作做好理论准备。本书从众多先贤及科学家的轶事讲起，以逐步归纳和递进的脉络总结出科学及数据科学应关注的要点，然后在生产的各个环节中对这些要点逐一进行讨论与落实，从更高、更广的视角回看科学及数据科学在各个生产环节的缩影。本书并不以高深的数学理论研究作为目的，也不以某一种计算机语言编程作为主线脉络，而是在一个个看似孤立的故事与工程中不断拾遗，并试着从中悟出一些道理。

本书适合大数据从业人员和对大数据相关知识感兴趣的人，初级和中级程序员、架构师及希望通过数据的感知改进工作的人，产品经理、运营经理、数据分析师、数据库开发工程师等对数据分析工作敏感的人，以及所有对数据科学感兴趣并希望逐步深入了解数据科学知识体系的人阅读。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目（CIP）数据

数据科学家养成手册 / 高扬编著. —北京：电子工业出版社，2017.5

ISBN 978-7-121-31304-2

I. ①数… II. ①高… III. ①数据管理—手册 IV. ①TP274-62

中国版本图书馆 CIP 数据核字（2017）第 071293 号

策划编辑：潘 昕

责任编辑：潘 昕

印 刷：三河市良远印务有限公司

装 订：三河市良远印务有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×980 1/16 印张：23 彩插：2 字数：500 千字

版 次：2017 年 5 月第 1 版

印 次：2017 年 5 月第 1 次印刷

定 价：79.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：（010）88254888。

质量投诉请发邮件至 zllts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

本书咨询联系方式：（010）51260888-819，faq@phei.com.cn。



图1-3

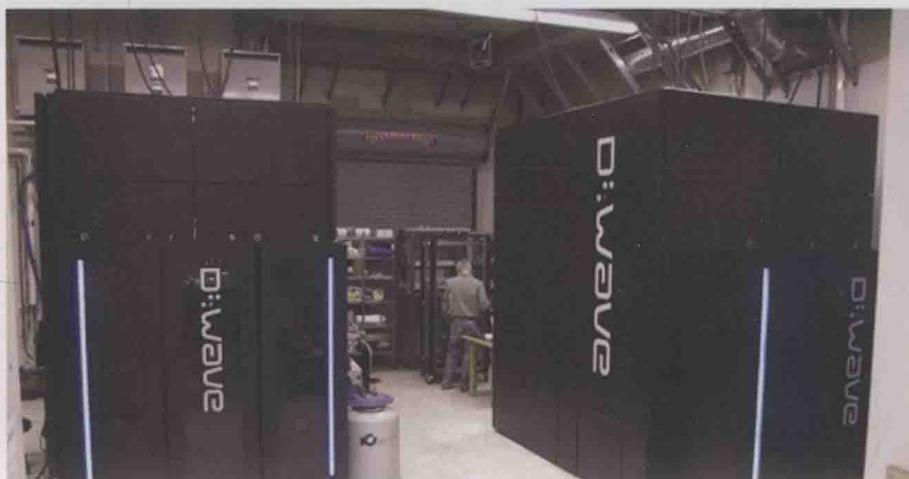


图4-12

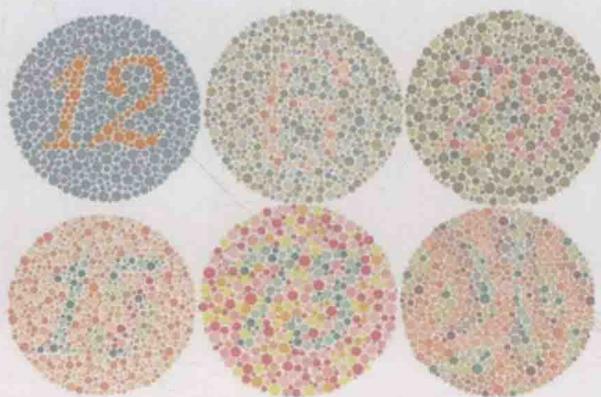


图7-7

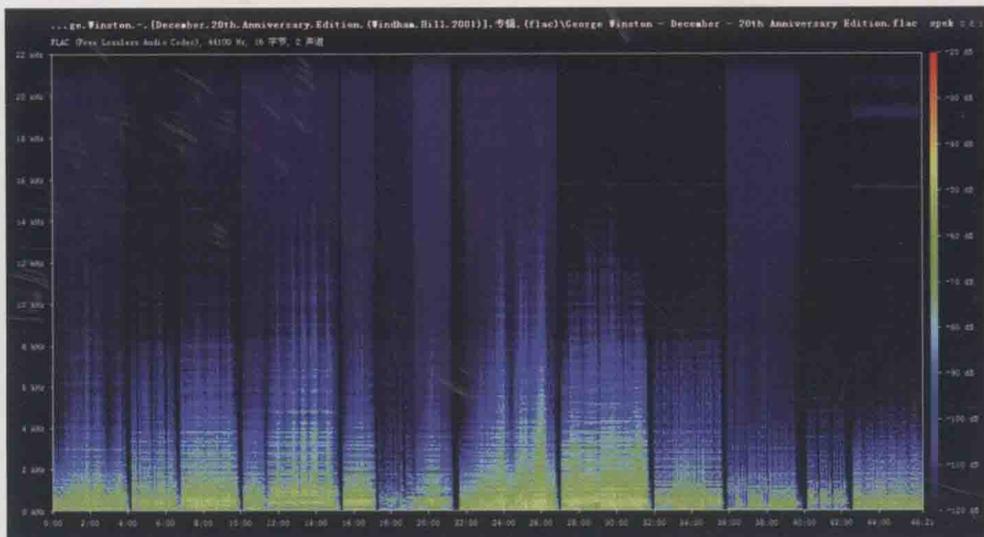


图9-10

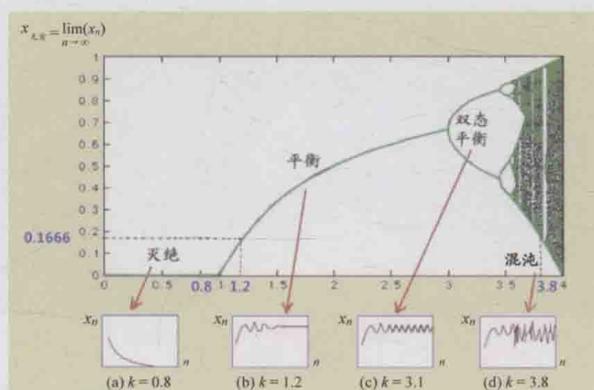


图10-4

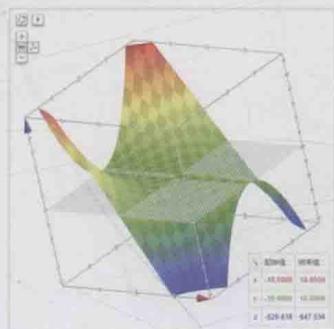


图10-5

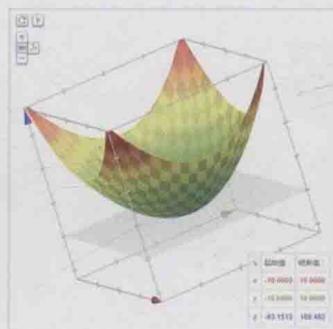


图11-18

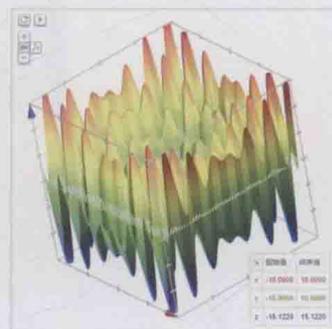


图11-20

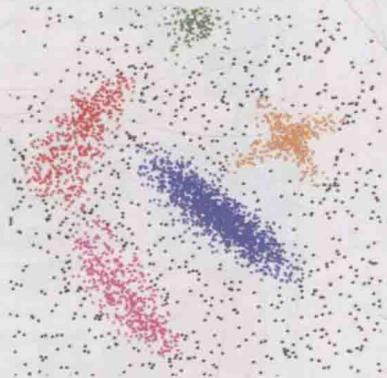
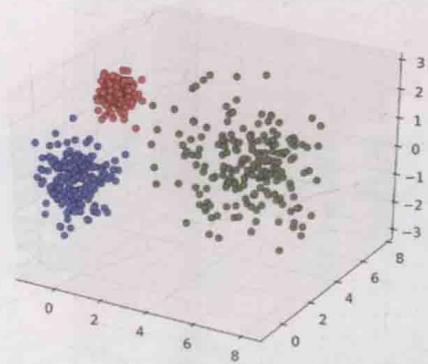


图11-22

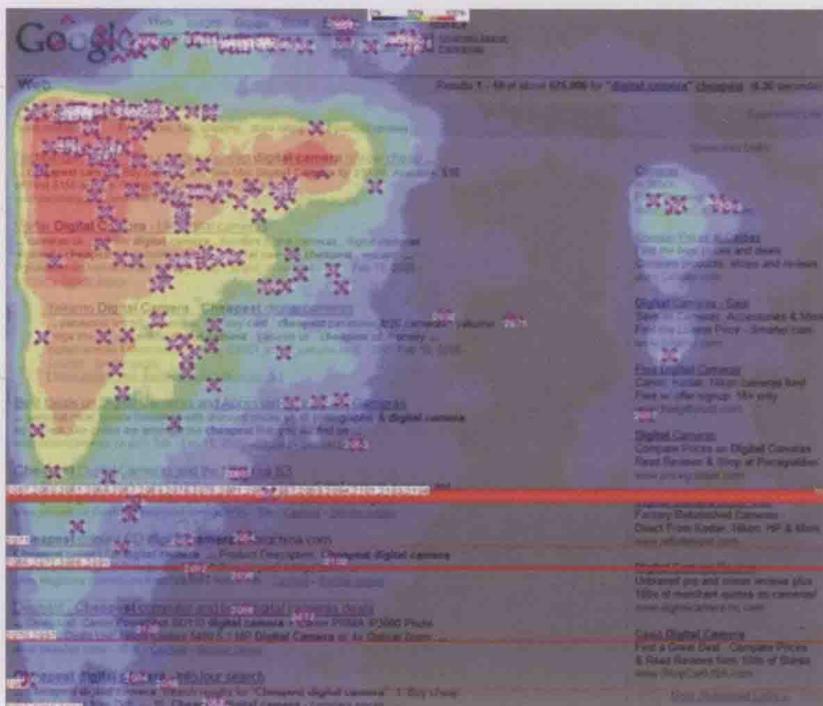


图16-6

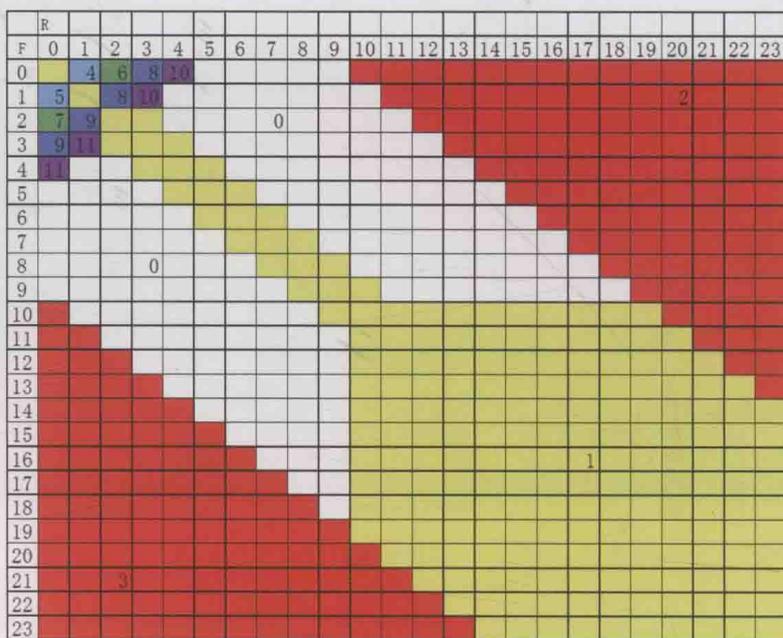


图18-21

轻松注册成为博文视点社区用户 (www.broadview.com.cn)

互动交流 在页面下方【读者评论】处留下您的疑问或观点，与作者和其他读者一同学习交流。

提交勘误 您对书中内容的修改意见可在【提交勘误】处提交，若被采纳，将获赠博文视点社区积分（在您购买电子书时，积分可用来抵扣相应金额）。

页面入口 <http://www.broadview.com.cn/31304>



序

十几二十年前，读书是学习新技术的不二法门。当时如果要学习一门技术，都需要买上几本“砖头书”，一边阅读，一边动手，一页一页“啃”下来。很多在今天叱咤风云的高手，当年都是用这种方式打下基础的。

最近几年，技术学习的方式发生了深刻的变化，大量的在线视频课程、交互式学习环境、开箱即用的工具箱，使技术学习的效率大幅度提升。特别是在动手能力方面，培训效率有了质的飞跃。最近一年，受人工智能领域突破性进展的鼓舞，机器学习和数据科学成为技术圈中的显学，而在线学习成为主流的学习方式。在这种情况下，大批学习者仅仅看过一些视频教程，按要求在 Jupyter Notebook 中做过一些练习，就基本具备动手解决问题的能力，可以上岗了。

这当然很好。但是，倘若你想在某一个领域取得真知，读书仍然是不可或缺的手段。中国信息安全领域的领军人物冯登国院士曾经说，以他的经验，想要真的搞懂某一个领域，非得深入“啃”至少一本书不可。读书的效率相对于听课、看视频要低得多，而多维的知识体系通过单维的文字表达出来，也给理解带来了挑战。然而，唯其有这种挑战，才需要读者进入深度思考状态，使读书成为一个推敲、琢磨、设问和破解的过程。不经过这个过程，我们所学到的知识一般来说只能是浮于表面的，很难达到“知其然知其所以然”的高度。正因如此，我们已经开始发现，仅通过在线视频和动手练习的学习者，对于相关领域的理论掌握经常是肤浅的。可以说，到目前为止，读书作为一种学习手段，依然是其他方式无法取代的。

机器学习和数据科学领域有几本非常重要的著作，每一个有野心的学习者都应该选择至少一本深入研究。Christopher Bishop 于 2006 年出版的 *Pattern Recognition and Machine Learning*，Kevin Murphy 2012 年的巨著 *Machine Learning: A Probabilistic Perspective*，斯坦福大学两位机器学习泰斗 Trevor Hastie 和 Robert Tibshirani 及其学生合著的 *An Introduction to Statistical Learning*，当然还有 Ian Goodfellow 和 Yoshua Bengio 最近出版的 *Deep Learning*——称这几本书为这个领域的“四书五经”，应该没有争议。

但是，这几本书有一个共同的问题——都是按照教材的体例编写的，所以都是尽全力系统化地介绍知识，对这个领域丰富多彩的应用、历史、人文和故事却很少展开论述。而要成为一名数据科学家，仅有知识和动手能力是不够的，还需要有相应的素养，这包括特有的思维方式、

价值观，对相关历史背景和掌故的了解，以及对数据科学社区的认知和互动——这恰恰是本书的最大价值。

作者把数据科学放在一个更广阔的背景之中，从数学、统计学、方法论甚至认知论的层面出发，讨论数据科学的内涵和外延，内容丰富，旁征博引，语言生动，灵活有趣，帮助读者站在一个更丰富的势场中认识数据科学，理解数据科学的基本思想。尤为令人欣喜的是，作者将信息论、混沌理论纳入讨论之中，表明作者敏锐地注意到数据科学与系统科学和认识论的深层联系，这是难能可贵的。从这个角度来解说数据科学的书，应该说在中国是第一本，即使在全球范围内也是独具特色的。为此，我们愿意向读者推荐本书，并相信读者一定能从中获得非常有价值的启发。

CSDN、AI100 创始人 蒋 涛

AI100 合伙人 孟 岩

前 言

为什么要写这本书

随着计算机科学和数据科学的发展，越来越多的人开始把目光投向其中最为耀眼的互联网、物联网、大数据、人工智能等高新技术领域，并且有相当多的高级技术人才已经在这些领域获得了令人瞩目的成就。

在追逐信息技术发展浪潮的过程中，数据科学成为人们在信息技术海洋中遨游所沉淀下来的理论与科学基础。我们都渴望通过对数据科学的理解来对生产工作进行指导和改善，这种工作的意义与其他各种在信息技术产业一线工作所创造价值的意义一样非比寻常。它给我们更广的辩证思考的空间，更高的观察事物的眼界，更多的自新的维度与动力。它是那么神秘且有趣。

今天，数据科学已经渗透到我们每个人的工作和生活之中。在你早上起来赶公车或者地铁的时候，你其实正在享受由数据科学辅助进行的精确调度服务；在你阅读工作报表的时候，你其实正在享受由数据科学辅助进行的大数据统计服务；在你吃午餐的时候，你其实正在享受由数据科学辅助进行的外卖快餐数据分发或食堂菜品改良服务；当你晚上回到家，在网上尽情购物的时候，你其实正在享受由数据科学辅助进行的高效电子商务和智能推荐服务。驾车出行有智能导航，就医问药有分诊机器人……也许你的家人或者朋友现在就在自己的工作岗位上，作为一名普通的销售人员、产品经理、人力资源师、售后服务人员、商务代表等，通过数据决策系统、数据库甚至电子表格来观察数据，作出判断，开展工作。数据科学给我们带来的红利已经紧紧把我们包围。

这本书绝无说教的想法，而是希望以书为媒，用谈天说地的方式，以激发每个人的思考为主要手段，归纳总结数据科学的实质及成就一位数据科学家所需要的基本素养。

遗憾的是，越是基础性、本源性的学科，与变成现实利益的距离就越远，让人觉得似乎不够实惠，不够亲近。至少读完这本书没办法帮你直接在第二天变出米饭、房子和汽车。不过我认为，楼房再高再漂亮，也需要人们看不到的深厚地基来支持；花儿再芬芳再娇艳，也需要在土壤之下吮吸养分的丰富根系来供能。这些看不到的东西，往往起着我们无法想象的巨大作用，

而这才是我希望与你一同讨论并思考的。

我们热爱生活，我们热爱所做的工作，我们希望在不断的攀登中看到更深更远的世界并去伪存真。那就让我们在点点滴滴的知识片段中一起开始慢慢思索、细细揣摩这一养成过程吧。

本书特色

本书从众多先贤及科学家的轶事开始讲起，以逐步归纳和递进的脉络总结出科学及数据科学应关注的要点，然后在生产的各个环节中对这些要点逐一进行讨论与落实，将这本书变成一本具有一定思维升华价值的参考书，从更高、更广的视角回看科学及数据科学在各个生产环节的缩影。

本书并不以高深的数学理论研究作为目的，也不以某一种计算机语言编程作为主线脉络，而是在一个个看似孤立的故事与工程中不断拾遗，并试着从中悟出一些道理。

简洁与深刻并重是本书的另一大特点。作为认知科学的延伸，数据科学一方面应该越来越引起广大数据工作者的重视，另一方面也要撩开自己的神秘面纱，以最为亲民的姿态和每位大数据工作者成为亲密无间的战友，为用科学的思维方式进行工作做好理论准备。

读者对象

- (1) 大数据从业人员和对大数据相关知识感兴趣的人。
- (2) 初级和中级程序员、架构师，以及希望通过对数据的感知改进工作的人。
- (3) 产品经理、运营经理、数据分析师、数据库开发工程师等对数据分析工作敏感的人。
- (4) 希望在思维方式领域进行拓展的高校毕业生和希望接触并了解数据科学的社会人员。
- (5) 所有对数据科学感兴趣并希望逐步深入了解数据科学知识体系的人。

如何阅读本书

本书分为3篇，分别是认知篇、分化篇和实践篇。

认知篇（第1章~第7章）

归纳了什么是科学，数据科学的范围、定义与实践价值，以及辩证思维、哲学和实验的关系等问题。这些是认知观点的基石。

分化篇（第8章~第11章）

重点介绍了数据科学中与现代社会各行业联系最为紧密的统计学、信息论、算法学，另外把混沌论作为一个知识点进行了补充。这些是认知观点在不同细分学科中所形成的一些具体解决问题的思维方式 and 科学观点。

实践篇（第 12 章 ~ 第 19 章）

沿着数据生命周期进行演进。任何行业的数据生命周期都是按照采集、存储、统计与建模、算法、可视化与分析、决策支持的沿革来进行的，本篇对各个环节的注意事项和思维方式都做了详细的讨论，并在第 18 章介绍了两个具体的数据分析案例。

在本书的最后，补充了过去与同行们讨论过的，并在会议演讲及日常分享的过程中总结出来的一些精彩问答。

如果你希望读完这本书后能够在数学方面有很大的提升，在工程代码能力方面有巨大的进步，这本书恐怕帮不上什么大忙。但我相信，在读完这本书后，你会在一些以前并不熟知的领域有所了解和感悟，并逐步完善理解和分析问题的视角。如果你不是数据研究人员，也可以把这本书当成一个休闲读本。这本书里既没有太多的公式，也没有太过高深的理论，有的只是我在和你攀谈的过程中与你一起发现的新视角。

特别致谢

感谢绘麟社相辉先生和李晓林女士对本书的插画助力。

参加本书编写工作的有高扬、卫峥、左妍、尹会生、杨艺、陈钢、肖力。

勘误和支持

由于作者的水平有限，编写时间仓促，书中难免会出现一些错误或者不准确的地方，恳请读者批评指正。如果您有更多的宝贵意见，欢迎扫描本页的二维码，关注“奇点大数据”微信公众号与我们进行互动讨论。本书后续的代码上传及勘误等相关更新内容都会在这个微信公众号发布。关注大数据尖端技术发展，关注奇点大数据。



同时，您也可以通过邮箱 77232517@qq.com 与我联系，期待能够得到您的真挚反馈，在技术之路上互勉共进。

高扬

2017年1月于珠海

目 录

认知篇

第 1 章 什么是科学家	2	3.2 数学的奥妙	29
1.1 从太阳东升西落开始	2	3.2.1 《几何原本》	29
1.1.1 农历	2	3.2.2 《九章算术》	30
1.1.2 公历	5	3.2.3 高等数学	34
1.1.3 小结	7	3.3 本章小结	37
1.2 阿基米德爱洗澡?	7	第 4 章 数据科学的使命	38
1.3 托勒密的秘密	10	4.1 走近数据科学	38
1.4 牛顿为什么那么牛	11	4.1.1 介质	38
1.4.1 苹果和三大定律	11	4.1.2 从信息到数据	41
1.4.2 极限和微积分	12	4.1.3 数据科学的本质	43
1.5 高斯——高，实在是高	15	4.2 万能的数据科学	44
1.6 离经叛道的爱因斯坦	17	4.2.1 测量	44
1.7 本章小结	20	4.2.2 统计计算	47
第 2 章 什么是科学	23	4.2.3 指标	52
2.1 科学之科	23	4.3 使命必达	53
2.2 边界的迷茫	23	4.3.1 高效生产	53
2.3 科学之殇	26	4.3.2 破除迷信	56
2.4 本章小结	27	4.3.3 目标一致与不一致	57
第 3 章 数据与数学	28	4.4 本章小结	58
3.1 什么是数据	28	第 5 章 矛盾的世界	59
		5.1 古希腊——学者高产的国度	59

5.2 矛盾无处不在	61
5.3 世界究竟是否可知	63
5.4 薛定谔的“喵星人”	64
5.5 本章小结	67

第6章 实验和哲学

6.1 朴素的认知方法	68
6.1.1 眼见为实	69
6.1.2 归纳与总结	70
6.2 哲学靠谱吗	71

分化篇

第8章 统计学

8.1 数理统计鼻祖——阿道夫·凯特勒	86
8.2 统计就是统共合计	88
8.3 数据来源	90
8.4 抽样	91
8.5 对照实验	91
8.6 误差	94
8.6.1 抽样误差	94
8.6.2 非抽样误差	96
8.7 概括性度量	97
8.7.1 集中趋势度量	98
8.7.2 离散程度度量	100
8.7.3 小结	100
8.8 概率与分布	100
8.8.1 数学期望	102
8.8.2 正态分布	103
8.8.3 其他分布	106
8.9 统计学与大数据	107

6.3 数学的尽头是哲学	72
6.4 本章小结	73

第7章 辩证思维

7.1 要不要辩证有多大区别	74
7.2 谁对谁错	76
7.3 做到客观不容易	77
7.4 观念的存弭	79
7.5 本章小结	82

第9章 信息论

9.1 模拟信号	109
9.2 信息量与信息熵	110
9.3 香农公式	111
9.4 数字信号	112
9.5 编码与压缩	113
9.5.1 无损压缩	114
9.5.2 有损压缩	117
9.6 本章小结	126

第10章 混沌论

10.1 洛伦兹在想什么	128
10.2 罗伯特·梅的养鱼计划	129
10.3 有限的大脑，无限的维	130
10.4 “谋杀上帝”的拉普拉斯	132
10.5 庞加莱“不是省油的灯”	134
10.6 未知居然还能做预测	137
10.7 本章小结	137

第 11 章 算法学	139	11.8.2 监督学习	164
11.1 离散的世界	139	11.8.3 强化学习	175
11.2 成本的度量	142	11.9 神经网络——深度学习	177
11.3 穷举法——暴力破解	143	11.9.1 神经元	177
11.4 分治法——化繁为简	152	11.9.2 BP 神经网络	180
11.5 回溯法——能省则省	154	11.9.3 损失函数	180
11.6 贪心法——局部最优	155	11.9.4 非线性分类	182
11.7 迭代法——步步逼近	156	11.9.5 激励函数	186
11.7.1 牛顿法	157	11.9.6 卷积神经网络	188
11.7.2 梯度下降法	158	11.9.7 循环神经网络	190
11.7.3 遗传算法	159	11.9.8 小结	193
11.8 机器学习——自动归纳	161	11.10 本章小结	194
11.8.1 非监督学习	162		

实践篇

第 12 章 数据采集	196	13.1.2 读少写多	211
12.1 数据的源头	196	13.1.3 读写都多	212
12.2 日志的收集	197	13.2 进快还是出快	213
12.2.1 实时上传	197	13.2.1 最快写入	213
12.2.2 延时上传	201	13.2.2 读出最快	215
12.2.3 加密问题	202	13.3 文件还是数据库	215
12.2.4 压缩问题	203	13.4 要不要支持事务	216
12.2.5 连接方式	204	13.5 表分区和索引	218
12.2.6 消息格式	206	13.5.1 表分区	219
12.2.7 维度分解	207	13.5.2 索引	219
12.3 这只是不靠谱的开始	208	13.6 稳定最重要	222
12.4 本章小结	209	13.7 安全性和副本	223
第 13 章 数据存储	210	13.7.1 RAID	223
13.1 读写不对等	210	13.7.2 软冗余	225
13.1.1 读多写少	211	13.8 本章小结	226

第 14 章 数据统计	227	16.6.3 巧测圆周率.....	253
14.1 此“统计”恐非彼“统计”.....	227	16.7 仁者见仁,智者见智.....	255
14.2 要精确还是要简洁.....	231	16.8 永恒的困惑.....	256
14.3 统计是万能的吗.....	232	16.9 本章小结.....	257
14.4 注意性能.....	234	第 17 章 数据决策	258
14.5 本章小结.....	234	17.1 决策就是“拍脑袋”.....	258
第 15 章 数据建模	235	17.2 哪里有物质,哪里就有数据.....	259
15.1 模型是宝贵的财富.....	236	17.2.1 目的的统一.....	259
15.2 量化是关键.....	236	17.2.2 数据胜于雄辩.....	260
15.3 该算法出马了.....	237	17.3 这是风险博弈.....	260
15.3.1 统计学模型.....	238	17.3.1 性价比优先.....	261
15.3.2 线性关系.....	238	17.3.2 小迭代至上.....	262
15.3.3 复杂的非线性关系.....	238	17.3.3 不要“输不起”.....	262
15.4 算法的哲学.....	240	17.3.4 留得青山在.....	263
15.5 本章小结.....	241	17.4 本章小结.....	264
第 16 章 数据可视化与分析	242	第 18 章 案例分析	266
16.1 看得见,摸得着.....	242	18.1 K线图里的秘密.....	266
16.2 颜色很重要.....	242	18.1.1 什么是市场.....	267
16.3 别说布局没有用.....	244	18.1.2 谁在控制价格.....	267
16.3.1 由上而下,由简而繁.....	244	18.1.3 货币价格的形成.....	270
16.3.2 总-分,分-总,总-分-总.....	245	18.1.4 零和博弈.....	271
16.3.3 毗邻吸引.....	246	18.1.5 涨跌都盈利.....	272
16.4 有图就别要表格.....	248	18.1.6 价格的预测.....	273
16.5 分析的内涵.....	249	18.1.7 形态.....	274
16.5.1 相关性分析.....	250	18.1.8 K线图周期.....	276
16.5.2 预测分析.....	251	18.1.9 造市商与点差.....	277
16.5.3 其他分析.....	252	18.1.10 科学分析.....	277
16.6 有趣的统计应用.....	252	18.1.11 小结.....	310
16.6.1 不规则图形的面积.....	252	18.2 数学能救命.....	310
16.6.2 套出你的实话.....	253	18.2.1 阴云下的大西洋.....	311