

白话大数据 与机器学习

高扬 卫峥 尹会生◎著
万娟◎插画设计



资深大数据专家多年实战经验总结，拒绝晦涩，开启大数据与机器学习妙趣之旅
以降低学习曲线和阅读难度为宗旨，重点讲解了统计与概率、数据挖掘算法、实际应用案例、
数据价值与变现，以及高级拓展技能，清晰勾勒出大数据技术路线与产业蓝图



机械工业出版社
China Machine Press

白话大数据 与机器学习

高扬 卫峥 尹会生◎著

万娟◎插画设计



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

白话大数据与机器学习 / 高扬等著. —北京: 机械工业出版社, 2016.6 (2016.9 重印)

ISBN 978-7-111-53847-9

I. 白… II. 高… III. ①数据处理 ②机器学习 IV. ① TP274 ② TP181

中国版本图书馆 CIP 数据核字 (2016) 第 115280 号

白话大数据与机器学习

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 高婧雅

责任校对: 殷虹

印刷: 三河市宏图印务有限公司

版次: 2016 年 9 月第 1 版第 2 次印刷

开本: 186mm × 240mm 1/16

印张: 21.75 (含 0.25 印张彩插)

书号: ISBN 978-7-111-53847-9

定价: 69.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88379426 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzit@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

为什么要写这本书

不知从何时开始我们已周身没入大数据时代的潮流，不知不觉被卷入了大数据时代。

无论是每天上网看网页、聊QQ、聊微信，或者登录银行、网购、买票，或者出行、投宿，甚至是出入任何公众场合、驾车、用水用电……我们无时无刻不在生产着各种数据。而同时我们也在消费着其他人生产的数据，我们使用的众多家电产品，每一个设计细节都融入了设计者对用户体验数据的调查与分析；我们使用的每一部手机、每一台电脑，每一个部件的产出都融入着多得无法想象的指标数据控制下的生产与监控；我们访问的每一个网页、每一个软件，每一次享受到的贴心的产品改动和服务的升级，无不浸透着无数的数据汇集与精细的分析和反馈。这是一场慢慢到来的、贯穿所有产业的革命，这是一次润物细无声的各行业精耕细作的开端。

不管我们是不是愿意，不管我们有没有意识到，我们现在已经身处大数据时代的奇点，而未来要迎接的是大数据奇点爆炸给我们带来的冲击力。我们需要力量来驾驭浪里的航船，我们需要乘风破浪前进的动力。

在这一次远航中，我们不必担心自己的能力水平无法感知数据这种磅礴之力的气魄，不必担心晦涩难懂的公式定理会让我们感到阻力。

请相信我，这是一本通俗易懂的大数据图书，这是一本轻松愉悦的数据挖掘和机器学习的读本，这是一本没有门槛的机器学习实战手册。让我们一起扬帆远航吧！

本书特色

从行为脉络来看，本书基本上是从数据统计、数据指标理解、数据模型、聚类/分类与机器学习、数据应用、大数据框架补充知识，以及扩展讨论这样的角度来层层深入完成的。

这种方式会给读者比较好的带入感，让大家——尤其是不擅长数学的读者降低对大数据与机器学习算法的恐惧感。如果读者朋友对排列组合、统计分布这些基础知识比较了解，完全可以考虑跳过这些部分直接去读后面更感兴趣的内容。

为了调节阅读气氛，我们还尝试加入了一些漫画插图。为了让读者朋友能够更快地进行实践，我们几乎在每一个算法讲解后都配有 Python 或者 SQL 语言的实现部分。相信这些能够帮助大家更快、更轻松地完成阅读。

读者对象

- (1) 对大数据感兴趣但是完全不了解的技术人员。
- (2) 对机器学习和数据挖掘比较感兴趣的技术人员。
- (3) 大数据初级从业人员。

如何阅读本书

本书一共分为 18 章。

第 1 章~第 5 章为入门所需基础知识及对数据指标运营的阐述。

第 6 章~第 10 章是对数据挖掘基础知识与算法的介绍。

第 11 章~第 18 章为生产应用与高级扩展。

其中，第 1 章~第 15 章正文内容，以及第 17 章、第 18 章的正文内容由高扬编写。

全书所有的 Python 代码由卫峥编写与补充整理。

第 16 章、附录全部由尹会生编写。

全书所有的漫画插画由万娟创作完成。

勘误和支持

由于水平有限，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。如果你有更多的宝贵意见，欢迎扫描下方二维码，关注“奇点大数据”微信公众号和我们进行互动讨论。关注大数据尖端技术发展，关注“奇点大数据”。

同时，你也可以通过邮箱 77232517@qq.com 联系到我，期待能够得到你的真挚反馈，在技术之路上互勉共进。



致谢

特别感谢：万娟女士为本书做的漫画插画内容。

万娟女士现任深圳星盘科技有限公司 UI 设计师，是我在多年工作中遇到过的最敬业的 UI 设计师之一，在 2013 年一起合作的过程中给我留下了非常深刻的印象。

她多次参加全国和国际艺术比赛，曾获得全国青少年绘画大赛铜奖，中国-新加坡国际青少年绘画比赛优秀奖，以及全国大学生工业设计大赛三等奖。从小酷爱绘画，理想是开一个属于自己的画室。

她给我留下的最深刻的印象用两个词可以描述：一个词是“敬业”，不管是在过去共事期间的合作，还是在为本书创作插画的过程中，为了保证进度带病坚持创作，都让我非常感动；另一个词是“唯美”，不仅人长得美，作品设计风格也透出现代与时尚的气息。

此外还要对所有支持和关心本书成书的各界朋友表示由衷的感谢：

衷心感谢北京邮电大学软件学院杨谈老师对本书的审校工作。

衷心感谢腾讯公司数据分析师彭瑶女士对本书的审校工作。

衷心感谢重庆工商大学黄辉老师、杨艺老师对本书的大力支持。

衷心感谢机械工业出版社华章公司对本书的支持与帮助。

衷心感谢“奇点大数据”微信群友对本书的关注与支持。



高 扬

目 录 *Contents*

前 言

第1章 大数据产业	1
1.1 大数据产业现状	1
1.2 对大数据产业的理解	2
1.3 大数据人才	3
1.3.1 供需失衡	3
1.3.2 人才方向	3
1.3.3 环节和工具	5
1.3.4 门槛障碍	6
1.4 小结	8
第2章 步入数据之门	9
2.1 什么是数据	9
2.2 什么是信息	10
2.3 什么是算法	12
2.4 统计、概率和数据挖掘	13
2.5 什么是商业智能	13
2.6 小结	14
第3章 排列组合与古典概型	15
3.1 排列组合的概念	16

3.1.1	公平的决断——扔硬币	16
3.1.2	非古典概型	17
3.2	排列组合的应用示例	18
3.2.1	双色球彩票	18
3.2.2	购车摇号	20
3.2.3	德州扑克	21
3.3	小结	25
第4章	统计与分布	27
4.1	加和值、平均值和标准差	27
4.1.1	加和值	28
4.1.2	平均值	29
4.1.3	标准差	30
4.2	加权均值	32
4.2.1	混合物定价	32
4.2.2	决策权衡	34
4.3	众数、中位数	35
4.3.1	众数	36
4.3.2	中位数	37
4.4	欧氏距离	37
4.5	曼哈顿距离	39
4.6	同比和环比	41
4.7	抽样	43
4.8	高斯分布	45
4.9	泊松分布	49
4.10	伯努利分布	52
4.11	小结	54
第5章	指标	55
5.1	什么是指标	55
5.2	指标化运营	58

5.2.1	指标的选择	58
5.2.2	指标体系的构建	62
5.3	小结	63
第6章	信息论	64
6.1	信息的定义	64
6.2	信息量	65
6.2.1	信息量的计算	65
6.2.2	信息量的理解	66
6.3	香农公式	68
6.4	熵	70
6.4.1	热力熵	70
6.4.2	信息熵	72
6.5	小结	75
第7章	多维向量空间	76
7.1	向量和维度	76
7.1.1	信息冗余	77
7.1.2	维度	79
7.2	矩阵和矩阵计算	80
7.3	数据立方体	83
7.4	上卷和下钻	85
7.5	小结	86
第8章	回归	87
8.1	线性回归	87
8.2	拟合	88
8.3	残差分析	94
8.4	过拟合	99
8.5	欠拟合	100
8.6	曲线拟合转化为线性拟合	101
8.7	小结	104

第9章 聚类	105
9.1 K-Means 算法	106
9.2 有趣模式	109
9.3 孤立点	110
9.4 层次聚类	110
9.5 密度聚类	113
9.6 聚类评估	116
9.6.1 聚类趋势	117
9.6.2 簇数确定	119
9.6.3 测定聚类质量	121
9.7 小结	124
第10章 分类	125
10.1 朴素贝叶斯	126
10.1.1 天气的预测	128
10.1.2 疾病的预测	130
10.1.3 小结	132
10.2 决策树归纳	133
10.2.1 样本收集	135
10.2.2 信息增益	136
10.2.3 连续型变量	137
10.3 随机森林	140
10.4 隐马尔可夫模型	141
10.4.1 维特比算法	144
10.4.2 前向算法	151
10.5 支持向量机 SVM	154
10.5.1 年龄和好坏	154
10.5.2 “下刀”不容易	157
10.5.3 距离有多远	158
10.5.4 N 维度空间中的距离	159
10.5.5 超平面怎么画	160

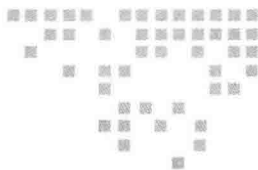
10.5.6	分不开怎么办	160
10.5.7	示例	163
10.5.8	小结	164
10.6	遗传算法	164
10.6.1	进化过程	164
10.6.2	算法过程	165
10.6.3	背包问题	165
10.6.4	极大值问题	173
10.7	小结	181
第11章	关联分析	183
11.1	频繁模式和 Apriori 算法	184
11.1.1	频繁模式	184
11.1.2	支持度和置信度	185
11.1.3	经典的 Apriori 算法	187
11.1.4	求出所有频繁模式	190
11.2	关联分析与相关性分析	192
11.3	稀有模式和负模式	193
11.4	小结	194
第12章	用户画像	195
12.1	标签	195
12.2	画像的方法	196
12.2.1	结构化标签	196
12.2.2	非结构化标签	198
12.3	利用用户画像	203
12.3.1	割裂型用户画像	203
12.3.2	紧密型用户画像	204
12.3.3	到底“像不像”	204
12.4	小结	205

第13章 推荐算法	206
13.1 推荐思路	206
13.1.1 贝叶斯分类	206
13.1.2 利用搜索记录	207
13.2 User-based CF	209
13.3 Item-based CF	211
13.4 优化问题	215
13.5 小结	217
第14章 文本挖掘	218
14.1 文本挖掘的领域	218
14.2 文本分类	219
14.2.1 Rocchio 算法	220
14.2.2 朴素贝叶斯算法	223
14.2.3 K-近邻算法	225
14.2.4 支持向量机 SVM 算法	226
14.3 小结	227
第15章 人工神经网络	228
15.1 人的神经网络	228
15.1.1 神经网络结构	229
15.1.2 结构模拟	230
15.1.3 训练与工作	231
15.2 FANN 库简介	233
15.3 常见的神经网络	235
15.4 BP 神经网络	235
15.4.1 结构和原理	236
15.4.2 训练过程	237
15.4.3 过程解释	240
15.4.4 示例	240
15.5 玻尔兹曼机	244

15.5.1	退火模型	244
15.5.2	玻尔兹曼机	245
15.6	卷积神经网络	247
15.6.1	卷积	248
15.6.2	图像识别	249
15.7	深度学习	255
15.8	小结	256
第16章	大数据框架简介	257
16.1	著名的大数据框架	257
16.2	Hadoop 框架	258
16.2.1	MapReduce 原理	259
16.2.2	安装 Hadoop	261
16.2.3	经典的 WordCount	264
16.3	Spark 框架	269
16.3.1	安装 Spark	270
16.3.2	使用 Scala 计算 WordCount	271
16.4	分布式列存储框架	272
16.5	PrestoDB——神奇的 CLI	273
16.5.1	Presto 为什么那么快	273
16.5.2	安装 Presto	274
16.6	小结	277
第17章	系统架构和调优	278
17.1	速度——资源的配置	278
17.1.1	思路一：逻辑层面的优化	279
17.1.2	思路二：容器层面的优化	279
17.1.3	思路三：存储结构层面的优化	280
17.1.4	思路四：环节层面的优化	280
17.1.5	资源不足	281
17.2	稳定——资源的可用	282

17.2.1	借助云服务	282
17.2.2	锁分散	282
17.2.3	排队	283
17.2.4	谨防“雪崩”	283
17.3	小结	285
第18章	数据解读与数据的价值	286
18.1	运营指标	286
18.1.1	互联网类型公司常用指标	287
18.1.2	注意事项	288
18.2	AB 测试	289
18.2.1	网页测试	290
18.2.2	方案测试	290
18.2.3	灰度发布	292
18.2.4	注意事项	293
18.3	数据可视化	295
18.3.1	图表	295
18.3.2	表格	299
18.4	多维度——大数据的灵魂	299
18.4.1	多大算大	299
18.4.2	大数据网络	300
18.4.3	去中心化才能活跃	301
18.4.4	数据会过剩吗	302
18.5	数据变现的场景	303
18.5.1	数据价值的衡量的讨论	303
18.5.2	场景 1：征信数据	307
18.5.3	场景 2：宏观数据	308
18.5.4	场景 3：画像数据	309
18.6	小结	310
附录A	VMware Workstation的安装	311

附录B	CentOS虚拟机的安装方法	314
附录C	Python语言简介	318
附录D	Scikit-learn库简介	323
附录E	FANN for Python安装	324
附录F	群众眼中的大数据	325
	写作花絮	327
	参考文献	329



大数据产业

1.1 大数据产业现状

大数据是近几年来都一直非常火热的一个名词，似乎是伴随着“互联网”的逐渐发展所出现的一个新名词。我们在天天听着“互联网+”的同时也在听说“大数据+”。

大数据其实是一个比较抽象和笼统的概念，应该说这个词是为了涵盖性地表达一系列生产和业务行为的一个统称。但是也正是由于这种抽象和过于简略的称谓方式，让每个人都容易对这个词产生见仁见智的不同视角的印象或者看法。

大数据是一个以数据为核心的产业，是一个围绕大数据生命周期不断循环往复的生产过程，同时也是由多种行业分工和协同配合而产生的一个复合性极高的行业。

在我看来，大数据产业生产流程从数据的生命周期的传导和演变上可以分为这样几个部分：**数据收集、数据存储、数据建模、数据分析、数据变现。**

其中每个环节都是非常重要的数据生命环节，每个环节的生产加工行为都是有其价值的，并且每个环节做到极致都可以成就一个伟大的公司。整个完整的产业生态圈就是大数据，它的缩影也渗透在任何一家以数据作为运营基础的公司中。

根据麦肯锡 2011 年发布的一份研究报告，到 2018 年世界范围内将会出现高达 14 万~19 万的“大数据”岗位空缺。而艾瑞咨询集团在“2014 年会”上曾指出，全球数据量每 18 个月翻一番，到 2015 年，中国专用数据分析人员预计缺口 1 400 万。

可以看到，在仅仅三四年的时间间隔上，两家咨询公司做出的预测都很大胆，但是两个估算数字相差也确实非常悬殊。究竟哪个数字更贴近“事实”并不好判断，因为大家对“大数据”的概念边界理解可能有很大的偏差，估算出现偏差是必然的，但是有一点可以肯定，大数据人才缺口一定是未来几年非常显著的问题。

2015年12月21日，全球第一家大数据交易所——贵阳大数据交易所经过半年多的发展，交易金额已突破6000万元人民币，会员数量超过300家，接入贵阳大数据交易所的数据源公司超过100家，数据总量超过10PB，已发生实际交易的会员超过70家。预计在未來3~5年，交易所日交易额将突破100亿元^①。

截至目前，中国境内除了贵阳大数据交易所以外，还有长江大数据交易所、武汉东湖大数据交易中心、崇州大数据交易所等十余家大数据交易所挂牌营业。

2016年1月，阿里云的“数加”大数据平台和金山云的KMR平台等国内大品牌云产品供应商的重磅产品先后登场，几乎所有有远见的云产品巨擘资本都在向大数据产业链集中。但是这块蛋糕似乎有点太大了，只能边烘焙边分割，谁也没办法一下子全吃掉。

1.2 对大数据产业的理解

“大数据”这个人造词汇其实很容易产生不少误解，尤其是这个“大”字，很容易让人感觉，数据量必须大，而且特别大，越大越能形成产业，也越有价值。其实这真的是“大数据”给人带来的误导。大数据产业的存在其实和其他产业并无二致，本身是为了给其他产业提供服务。

做个假设，假如现在给石油产业冠以“大石油”产业的名字，那么会影响石油行业本身对其他行业的服务样态吗？应该不会。

在“大石油”产业里，同样有人从事着这样的工作内容：**石油勘探、石油开采、石油运输、石油提炼、石油产品销售**等多个细分领域和环节。

最后提供给社会的是由大量人工和智慧凝结在石油产品上的服务，而这些服务极大地方便并满足了社会各领域对于工业能源、建筑材料、食品包装、服装面料、模型器具、日杂用品等多种制造与使用的需求。试想如果没有石油，也就没有廉价汽车与航空动力，尤其是没有乙烯等重要化工原料的来源，是否存在塑料这样一种廉价的工业制造材料都很难说，那么各个产业则需要用其他造价更为高昂的材料对其进行取代，更不用提家用的天然气和液化石油气了，人们只能再去寻找其他能源：要么不洁净——如柴火和煤炭，要么价格昂贵——如氢气。人们之所以选用石油作为整个产业链的根源，并把它发展成一个完整的产业也是由于这样的原因，大概这个逻辑是比较容易理解的。

类比一下“大数据”产业，数据收集、数据传输、数据存储、数据建模、数据分析、数据交易贯穿了大数据产业的完整产业链。在这个产业链里同样蕴含着和“大石油”一样的东西，这个东西是什么？

数据通过各种软件进行收集，通过网络进行传输，通过云数据中心进行存储，通过数据科学家或者行业专家进行建模和加工，最后数据分析得到的是一种知识，是一种人们通

^① 来自《新华网》的报道。