

不只介绍R语言

更深入数据挖掘的本质：探寻数据背后的逻辑，挖掘人们的欲望、需求及态度

Broadview®
www.broadview.com.cn



探寻数据背后的逻辑

R语言数据挖掘之道

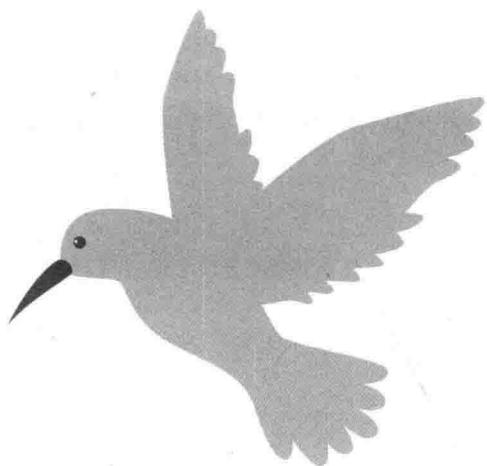
宋云生 张坚洪 黎新年◎著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>



探寻数据背后的逻辑

R语言数据挖掘之道



宋云生 张坚洪 黎新年〇著

电子工业出版社
Publishing House of Electronics Industry
北京•BEIJING

内 容 简 介

数据分析、数据挖掘的本质是探寻数据背后的逻辑，挖掘人们的欲望、需求、态度等。本书不仅仅教会读者如何掌握数据挖掘相关技能，更教会读者如何从数据挖掘结果中分析出更深层次的逻辑。

本书主要介绍使用 R 语言进行数据挖掘的过程。具体内容包括 R 软件的安装及 R 语言基础知识、数据探索、数据可视化、回归预测分析、时间序列分析、算法选择流程及十大算法介绍、数据抓取、社交网络关系分析、情感分析、话题模型、推荐系统，以及数据挖掘在生物信息学中的应用。另外，本书还介绍了 R 脚本优化相关内容，使读者的数据挖掘技能更上一层楼。

本书适合从事数据挖掘、数据分析、市场研究的工作者及学生群体，以及对数据挖掘和数据分析感兴趣的初级读者。

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

探寻数据背后的逻辑：R 语言数据挖掘之道 / 宋云生, 张坚洪, 黎新年著. —北京：电子工业出版社, 2018.8
ISBN 978-7-121-33861-8

I. ①探… II. ①宋… ②张… ③黎… III. ①数据采集—研究 IV. ①TP274

中国版本图书馆 CIP 数据核字(2018)第 049541 号

策划编辑：王 静

责任编辑：石 倩

印 刷：三河市双峰印刷装订有限公司

装 订：三河市双峰印刷装订有限公司

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：787×1092 1/16 印张：27 字数：742.4 千字

版 次：2018 年 8 月第 1 版

印 次：2018 年 8 月第 1 次印刷

定 价：89.00 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, 88258888。

质量投诉请发邮件至 zlts@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。

序言

提纲挈领式心诀： 一名数据挖掘工程师的成长之路

我的学习之路

不知不觉毕业两年多了，有一些大音如霜工作室的读者总想了解一下我是怎么学习数据挖掘、数据分析的，下面就综合大家常见的问题分享一下自己的经历、经验。

我不是学数学的，也不是学计算机的，研究生的专业是植物学，而且方向是植物分类，可以说很难和数据挖掘、市场研究等领域扯上关系。唯一能扯上关系的也就是我的舍友是做生物信息学研究的。

说一句丢人的话，在读本科时上的 SPSS 课我都不知道在讲什么。那时没考过计算机等级考试，原因是我每分钟打汉字的速度都不过关，讲这么多，只是为了告诉读者，我的基础并不扎实。

需要说明的是，我的英语还不错，在大一和大二分别通过了英语四、六级考试（而其他科目则学得比较一般，因为我每学期只有一两个主要学习目标），在大四我读了很多英文文献。因为在读研究生期间需要查阅大量文献，我需要给这些文章建立一个数据库，于是年少无知的我就选择了 Access。选择 Access 的原因并不是我比较熟悉它，而是我的老师用它，我至今也不会太多的操作。这应该算是我开始接触数据分析了。

使用高级语言时，记不住函数不要紧，但是你要有很强的搜索能力。

之所以讲这一段经历，不是为了说明我起步晚，而是为了说明建立 Access 文献数据库锻炼了我的英文搜索能力。我一碰到问题，就在 Google 里搜索，很快就能找到答案。于是 Google 几乎成了我的眼睛，真正做到了用 Google 搜索、发邮件、社交、阅读和写作。在公司里曾经传说，如果是连我都搜不到的内容，那么别人更不可能搜到。有些年轻人就怕英文，我并不是崇洋媚外，客观地想一想，现在的很多知识都是从欧美起源的，如果你连这门世界语言都不掌握，那么你获得的资料永

远都是二手资料。另外，无论你是找函数还是找包、模块，抑或是为问题寻找答案，使用 Google 进行英文搜索会为你节省很多时间。掌握这门语言并不需要你听、说、读、写样样精通，而是将其作为一种工具，应用起来比较方便就可以了。

要善用英文搜索，原因很简单，你所用的编程语言或软件大多是外国人构建的，并且在国外已经普及，相关的问答社区早已完善，你碰到的问题可能早就有人解决了。

在搜索文献的过程中，我喜欢上了《经济学人》的 *Graphic Details* 栏目，发现其绘制的图表非常漂亮、专业，于是我就开始学习 Excel，尽自己所能将 Excel 图表做得更漂亮、更专业，这些经历为我日后做数据可视化打下了坚实的基础：我知道了商务色彩搭配及图表要简洁、易读等原则，我知道怎么使自己的图表特色鲜明。后来看了大前研一先生的著作，了解了专业精神，我曾经写下这样一句话，以勉励自己：

所谓专业，即每一个细节都经得起推敲。

有一天，舍友看到我用 Excel 作图，嘲笑我孤陋寡闻，推荐我学习 R 语言，然后我就开始搜寻一些 R 语言入门读物进行阅读，慢慢地知道了关于这门语言的粗浅知识。

这个时候已经到研二下半学期了，我需要为自己未来的工作做打算了：是步入园林行业还是就此转行？必须做一个决断。我发现我真的对植物分类不感兴趣，而我做家教的学生的妈妈是星空传媒的一个经理，平时待我很好，她说毕业可以介绍我去做市场研究。我了解了一下市场研究，发现其中涉及一些数据分析的内容（现在看起来很简单），于是，我从此决定踏上数据分析这条“不归路”。

为了快速上手、熟悉统计学知识，我并没有马上深入地学习 R 语言，而是像以往一样懒懒散散地学习（后悔当时没有实战学习）。我通过搜索发现，市场研究的岗位大多将熟练使用 SPSS 作为硬性要求，偶尔也会要求熟悉 R 语言，但 SPSS 对我来说更容易上手，于是就开始学习 SPSS。SPSS 帮助我巩固了统计学知识，当学习完简单的统计学知识后，我发现 SPSS 不够灵活，很多功能不够用，做出的图表很难看（这对于我来说是无法忍受的），因此，网络上有一些人鄙视 SPSS，但很推崇 R 语言。于是我决定要深入地学习 R 语言。我先将 SPSS 的功能在 R 中做了一遍，有了一些自己的理解后，我开始在自己的论文里做一些数据分析的内容。

现在想来，如果我直接在实战中学习可能会节省更多的时间。

实战更能锻炼技能水平。

研二快结束了，开始找工作了。我找工作的目的很明确，如果工作不是做数据分析、数据研究的，那么我宁愿放弃这个工作的机会。非数据研究的岗位我也不去面试，这样又省下了大量的时间学习。

在工作中学习

2013 年毕业后，我去了一家医药市场研究公司，当时的工作并不太忙，我有大量的时间学习。但这时也暴露了我的弱点，公司的数据并不是很规整（raw data）的，往往需要标准化等，而且数据规模也不再是之前练习时那么小，在面对这些脏数据、大一点的数据时，我的数据清洗水平显得捉襟见肘。周围的人都是 Excel 高手，如果跟着他们学，估计也能成为高手，但是我一定要在 R 中做数据清洗整理，反正公司的工作不是很多，我就一点点地学习和积累，这样我的数据处理能力就逐渐扎实起来了。其间我用两天读完了《异类》这本书，感触很深，阅读经历已经写成一篇文章在我们的公众号里分享了。

任何一个工具在刚开始学习时都会觉得它很糟糕，其实这并不是工具的问题，而是自己的知识体系跟不上节奏，或者是它的很多方法与自己原有的认知相反，这时不要急于否定它，而是要深入地学习它。知识体系是一个积累过程，为自己准备一万个小时计划吧。

我们公司当时在做 BI（商业智能），于是我接触了市面上常见的 BI 工具，包括 Tableau、QV 等，我熟悉它们的优、劣势，也熟悉它们的数据可视化效果。因为需要将 R 语言的页面融入 BI 中，所以我熟悉了 shiny 包，做了一些页面，但我渐渐看到 R 语言在做这些通用语言的工作时所暴露的缺点，于是开始接触 Python。

后来，我们的合作公司的总经理听说我比较熟悉 R 语言，就向我请教，我们一起讨论了 R 语言和数据挖掘。得知他们在做文本挖掘，于是在我闲暇时间开始学习中文文本挖掘的内容。没有成型的数据和书，我就看帖子，去一个个地实现，然后积累经验，这时我对 R 的操作算得上非常熟练了，从实现到速度优化（并行计算等）等也已经非常熟练，积累的代码也非常多了。

后来，那个经理找我做医院处方数据挖掘工作，之后，他请我去负责法院文本数据挖掘，我没去，但成了他们的外援，仍然没收过钱，他们搭建的一台服务器也帮助我了解了不少 Linux 的知识。

刚开始，锻炼自己的机会远远比钱重要，反正自己闲着也是闲着，但是这种情况只适用于刚开始。

后来，我们公司推出了微信公众号平台，我开始给公司的公众号写文章。其间我为公司的公众号写了多篇关于综合排名的文章，阅读量最高达到 4 万多人次，当时公众号的粉丝才 2 000 人左右。后来我又制作了评价医院市场趋势的综合指标体系，现在公司也一直在沿用这套指标，这些工作中的小点子都是我在公交车上想出来的。

除要把工作当成一种谋生手段外，还必须有极大的兴趣，要么不做，要么做好。

另外，我在公交车上读完了 *Data Mining with R learning by case studies*、*Machine Learning for Hackers*、*R Graphics Cookbook* 等书籍，之所以提这三本书，是因为我不止读过一遍，这三本书很有特色，第一本帮助我学习了各种算法，第二本帮助我接触了实际应用中的知识，第三本帮我熟练了 ggplot 的函数及图表元素结构。我开始学会利用零散的时间，坚持积累，也开始学习高度自律。

古之成大事者，不唯有超世之才，亦必有坚韧不拔之志。

——苏轼

其实，我一直幻想着有一个自己想写什么就写什么的平台，于是，我和小伙伴们开通了微信公众号，直到现在，我们更注重文章的可读性、趣味性，而不仅仅是技术，但是每一篇文章都可以作为一个小项目让希望学习数据分析的读者能锻炼一下自己的技能。

经常有读者问学数据分析就一定要学编程吗？以及为什么要看英文资料？针对这两个问题，我写下了这样一段对话，希望你能在对话中找到答案。

为什么学习数据分析？

赚钱！

什么样的人容易赚钱？

技能比别人高的！

英语是不是一般人的难关？是不是大家都想学习傻瓜式操作软件？

是！

那么如果大家都这么想，你应该怎么做？

很明显，你要做其他人不愿意做的事情，才能赚到别人不能赚的钱！

作为数据分析师，一定要将自己和技术区分开，分析数据、挖掘数据本质上是探寻数据背后的人心，挖掘人们的欲望、需求、态度等，所以数据分析师还要尽量拓宽自己的视野和知识结构，尽自己所能博览群书。

我的经历大体如此，中间会有各种迷茫、各种苍白无力，但是如果你缺少什么，就去搜集资料，做出判断，努力去争取，这一点总不会错，千万不要一味地否定你不了解的东西，这也是我对待未知领域的态度。

作为一名技术人员，要让自己的知识时刻在进步！这是一种宿命。

作 者

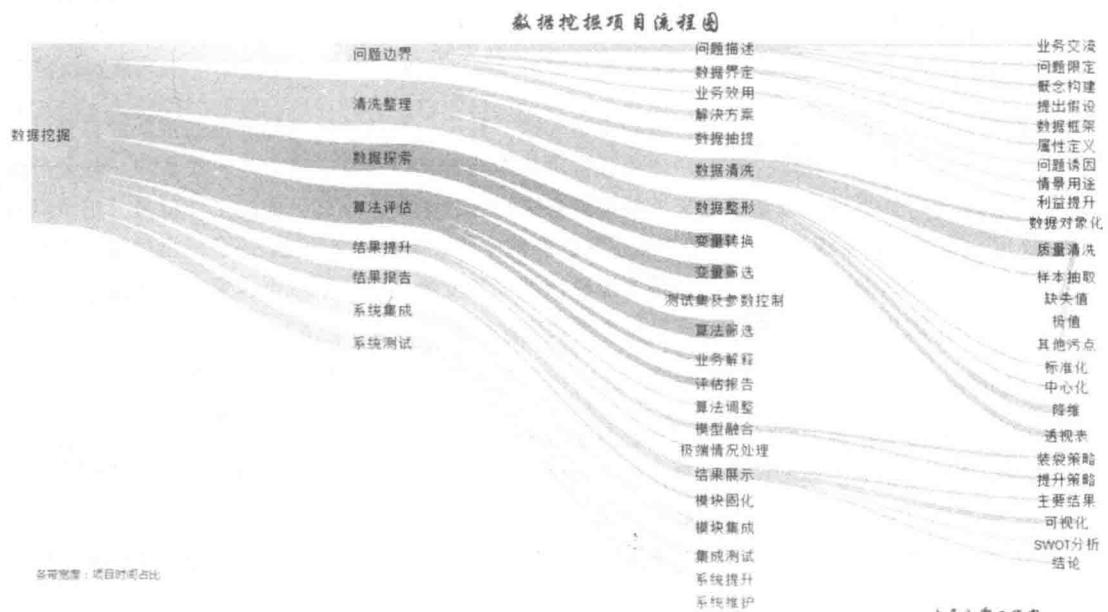
前言

什么是规范化的数据挖掘流程

人总是被自己日常从事的工作所蒙蔽双眼，看不到事态发展的整体面貌，为了手里的工作而工作，这就是所谓的迷失吧。一个数据挖掘项目不仅仅是数据挖掘工程师手中的一部分工作（虽然它是工作的核心），作为一个力求向上的人，要跳出来看看项目的全貌。只有对项目全程有了足够的了解，才能更加有效地使用数据挖掘、机器学习、数据分析的工具。

数据挖掘项目一般可以分为问题边界、清洗整理、数据探索、算法评估、结果提升、结果报告、系统集成和系统测试 8 个主要的模块。一般而言，可以尽量将这些模块合并，但无论怎么合并，它们在项目中都是不可或缺的（当然，有些项目并没有其中某些模块的需求，比如市场研究项目，它们可能就不需要系统集成）。从右图中可以看到数据清洗、数据探索、算法评估占据了项目的大部分时间，这也说明它们是项目的核心内容，缺了这 3 项，就不能再称为是数据挖掘项目了。

数据挖掘项目流程如下图所示。



问题边界

问题边界一般是项目的开头部分，可以分为 4 点。首先要和业务部门细致沟通，从业务背景中提炼出对业务问题的描述，限定项目要解决的问题，便于组织力量集中对这些问题设计解决方案。然后根据解决构想将业务问题转换为数据语言，限定将要使用的数据界限，搞清楚要牵涉哪些数据。之后为了吸引业务部门必须整理出业务效用，告诉业务部门如果解决这些问题能够得到哪些改善，完成业务部门哪些具体目标。最后要将以上问题整理成一个可行的解决方案。很多人忽略了这个阶段，其后果就是业务部门觉得挖掘出来的结果不是他们想要的，或者节外生枝补充各种相关的或不相关的业务问题，最终扭曲了项目本身，使工作反复无常。所以，在项目实施之前，非常仔细地沟通并制定一个完善的问题边界非常重要。

清洗整理

清洗整理是数据挖掘工程师非常熟悉的工作，但是，很少有人认识到这是项目中花费时间最多的部分，很多人会以为算法评估部分才会花费最多的时间。其实不然，如果数据清洗进行得不顺利，则将直接影响后面的工作和模型的效果。首先要设计畅通、高效的数据抽取程序，将数据从各种数据平台抽取出来供数据挖掘工具使用，然后进行数据清洗，将数据转化为数据挖掘工具便于处理的对象类型（在 R 里指 list、data.frame、array 等），再进行质量清洗，包括处理缺失值、异常值、其他污点（在文本挖掘中多如牛毛）等。之后要对数据整形，包括一些统计变化，例如中心化、标准化、降维等，更重要的是数据形状的变化。

还有一项就是数据抽样，面对大数据，在数据处理阶段就要进行抽样，不能因为要清洗一个点就清洗全量的数据，那样会花费大量的时间。不如抽取小样本进行测试，等进行完数据清洗程序后，再进行全量数据的整体清理，这样反而更加省时省事，这里的样本量需要尽量保证抽到足够多的问题数据，同时要让程序运行起来非常轻松、高效。

数据探索

数据探索要完成两个目标：变量转换和变量筛选。其中变量转换既包括变量的重新计算变形，也包括概念变量的构建，比如，在客户流失预警项目中要定义什么样的客户是流失客户，就会产生出一个新的变量。如果这个变量的定义不能用业务进行合理解释，那么下面的工作就是“瞎子点灯白费油”了。有些变量不仅不会对模型产生正向的影响，而且除了影响速度，还会降低模型的效果。显然进行变量筛选就非常重要了。谷歌预测流感模型筛选变量足以证明数据探索多么重要，而且在大数据环境下，数据探索已不再是一件轻而易举就能完成的事情了。

算法评估

算法评估是数据挖掘项目的灵魂。算法评估首先要求我们充分了解算法或模型的参数意义，然后需要预留测试数据集。模型评估不是仅仅比较模型结果的准确性是否存在差异（别忘了统计学教导我们比较差异时要判断差异的显著水平），所以，模型比较是对不同模型准确性均值的比较。算法筛选完成后，工作就告一段落了，这时要和业务部门一起对结果进行业务解释，不能进行业务解释

的数据挖掘结果就是为了数据挖掘而数据挖掘，这显然就是迷失在了项目中，遗忘了项目要解决的问题边界。最后要对结果进行完整的评估报告。评估报告是必需的，因为除了将它给领导看，更重要的是它能帮助你总结发现这个过程中可以改善的节点。

结果提升

首先要判断是否需要调整算法或模型，包括更换算法或调整参数。如果模型调整没有必要，那么就要考虑使用模型融合提高模型效率。模型融合的方法包括装袋（Bagging）和提升（Boosting）等，有些方式可能用业务解释起来比较困难，这也是数据挖掘工程师要考虑的问题之一，显然，有些问题选择可解释模型比较好。在项目中对一些极端情况最好另做处理。

结果报告

“丑媳妇也要见公婆，”分析结果报告最终要给业务部门的同事学习，教他们如何使用数据挖掘的结果进行业务分析和部署，其中主要成果要突出，吸引他们的眼球，一定要联系业务具体的困境或具体的业务情景，即所谓的对症下药。规律和结果必须通过易读的方式传达给受众，充满技巧的数据可视化是不二之选，将美妙的可视化图表嵌入具体的应用情景中进行宣讲，往往能达到事半功倍的效果，因此，在此处无论多么努力都不为过。SWOT 分析是业务部门最喜欢的分析方式，我们当然不能放过，以对方熟悉的方式表达自己的诉求，是交流的法宝。

模块固化这一步工作的快慢取决于之前的工作，如果之前已经考虑到后面要进行模块固化，那么就会将代码写得比较规范、注释良好，这种情况下就很容易将数据清洗、数据整形、变量转换、模型构建、结果输出等模块的内容固定下来，成为一个数据有进口及有出口的脚本文件。

系统集成

将固化下来的模块按照一定的秩序集成在一起，就成为一个分析的脚本体系。在这个体系中，有输入就有产出，中间不需要人工干预，是一个有序的自动化脚本体系。这一步考验数据挖掘工程师对每一步任务的理解。良好的模块集成可以提升整个系统的速度，减少后期维护的时间和次数。

模块集成后要与其他系统集成在一起，首先要和数据平台（数据库、Hive、Hadoop）对接，为分析模块提供数据来源和存储分析结果，同时要和前台展示对接，将结果可视化，让结果真正接触受众，即所谓的为决策者提供支持。

系统测试

这么一个“五脏俱全”的系统需要维护在所难免，总有一些极端情况会导致数据分析模块宕机，所以，代码一定要写得尽量规范，注释要尽量清晰，否则在维护时会有一种再造系统的感觉。关于规范请参看 *Google's R Style Guide*。

目录

第1章 万事不只开头难.....	1
1.1 工欲善其事，必先利其器：安装.....	1
1.1.1 安装 R 和 RStudio	1
1.1.2 安装数据包	3
1.1.3 数据包加载、卸载、升级，查看帮助文档.....	5
1.1.4 什么样的 R 包值得相信	7
1.2 了解 R 的对象	8
1.2.1 如何进行常见的算术运算	8
1.2.2 R 语言的三大数据类型	10
1.2.3 向量及其运算	12
1.2.4 因子变量鲜有人知的秘密	15
1.2.5 矩阵相关运算及神奇的特征值	17
1.2.6 数据框及其筛选、替换、添加、排序、去重	18
1.2.7 与数组（array）相比，表单（list）的用处更加广泛.....	22
1.2.8 如何进行数据结构之间的转化	23
1.3 R 语言的重器：函数	26
1.3.1 自编函数	26
1.3.2 有用的 R 字符串函数	29
1.4 控制流在 R 语言里只是一种辅助工具	31
1.4.1 判断	32
1.4.2 循环	33
1.5 数据的读入与输出	35
1.5.1 常见数据格式的输入 / 输出（CSV、TXT、RDATA、XLSX）	35
1.5.2 数据库连接：Oracle、MySQL 及 Hive	37
1.5.3 乱码就像马赛克一样让人讨厌	39

第 2 章 数据探索，招招都是利器	41
2.1 不要在工作后才认识“脏数据”	41
2.1.1 以老板信服的方式处理缺失数据	42
2.1.2 异常值预警	48
2.1.3 字符处理正则表达式不再是天书	49
2.2 数据透视、数据整形、关联融合与批量处理	50
2.2.1 还忘不掉 Excel 的数据透视表吗	50
2.2.2 你能给数据做整形手术吗：long 型和 wide 型	52
2.2.3 关联合并表	54
2.2.4 数据批处理：R 语言里最重要的一个函数家族：*ply	55
2.3 一招完成数据探索报告	58
2.4 拯救你的很多时候是基础理论	61
2.4.1 参数检验及非参检验	62
2.4.2 学了很多算法却忘了方差分析	68
2.4.3 多因素方差分析及协方差作用	70
2.4.4 很多熟悉的数据处理方法已经成笑话，工具箱该换了	73
第 3 章 从商务气质的数据可视化说起	84
3.1 说说数据可视化的专业素养	84
3.1.1 数据可视化历史上有多少背影等你仰望	84
3.1.2 商务图表应该具有哪些素质	87
3.1.3 那些你不知道的图表误导性伎俩	94
3.1.4 如何快速解构著名杂志的图表	98
3.2 ggplot2 包：一个价值 8 万美元的态度	103
3.2.1 一张图学会 ggplot2 包的绘图原理	105
3.2.2 基础绘图科学：ggplot2 包的主题函数继承关系图（关系网络图）	127
3.2.3 基础图表一网打尽	132
3.2.4 古老的地图焕发新颜	151
3.3 将静态图转为 D3 交互图表：plotly	156
3.4 从基础到进阶的变形图表	157
3.4.1 马赛克图（分类变量描述性分析）	157
3.4.2 Sankey 图和 chordDiagram 图	158
第 4 章 分位数回归模拟股票指数风险通道	163
4.1 用线性回归预测医院的药品销售额	163

4.2 多项式回归及常见回归方程的书写	168
4.3 Lasso 回归和回归评价的常见指标.....	170
4.4 分位数回归拟合上证指数风险通道	175
第 5 章 时间序列分析	181
5.1 时间序列分析：分析带有时间属性的数列	181
5.2 不是所有序列都叫时间序列	181
5.3 时间序列三件宝：趋势、周期、随机波动	183
5.3.1 趋势	183
5.3.2 周期	184
5.3.3 随机波动	186
5.4 预测分析	186
5.4.1 指数平滑法	186
5.4.2 ARIMA 模型预测	188
第 6 章 选择什么算法也有一套流程	192
6.1 重新审视一下这几个模型	192
6.1.1 Logistic 回归	192
6.1.2 我要的不是一棵树，而是整座森林：随机森林	195
6.1.3 神奇的神经网络	196
6.2 银行信用卡评估模型之变量筛选	197
6.2.1 变量构建	197
6.2.2 Logistic 回归变量筛选	198
6.2.3 随机森林变量筛选	203
6.2.4 人工神经网络建模	204
6.3 必须面对的模型评估	204
第 7 章 深入浅出十大算法	208
7.1 C5.0 算法	208
7.1.1 一个重要的概念：信息熵	208
7.1.2 非列变量选择的实例	209
7.1.3 C5.0 算法的 R 实现	210
7.2 K-means 算法	212
7.2.1 K-means 算法的 R 实现	212
7.2.2 怎么确定聚类数	213

7.3 支持向量机 (SVM) 算法	213
7.3.1 通俗理解 SVM	214
7.3.2 SVM 的 R 实现	216
7.4 Apriori 算法	216
7.4.1 举例说明 Apriori	217
7.4.2 Apriori 算法的 R 实现	219
7.5 EM 算法	220
7.5.1 举例说明 EM 算法	221
7.5.2 EM 算法的 R 实现	222
7.6 PageRank 算法	223
7.7 AdaBoost 算法	224
7.8 KNN 算法与 K-means 算法有什么不同	226
7.9 Naive Bayes (朴素贝叶斯) 算法	227
7.10 CART 算法	228
第 8 章 数据抓取	231
8.1 数据挖掘工程师不可抱怨“巧妇难为无米之炊”	231
8.2 抓取股市龙虎榜数据，碰碰运气	232
8.2.1 了解 XML 和 Html 树状结构，才能庖丁解牛	233
8.2.2 了解 RCurl 包和网页解析函数	234
8.2.3 抓取股票龙虎榜	235
8.2.4 资金流入分析	237
8.3 抓取某家医药信息网站全站药品销售数据	240
8.3.1 所有医药公司名称一网打尽	240
8.3.2 为什么抓取数据时可以使用 For 循环	242
8.3.3 不要把代码写复杂	244
8.3.4 用 Sankey 数据流描绘医药市场份额流动	248
第 9 章 不可不说的社交网络关系	254
9.1 社交网络图	254
9.1.1 社交网络图告诉你和谁交朋友	254
9.1.2 这几个基本概念你需要抓牢	256
9.1.3 还有比本章任务更有趣的数据挖掘吗	259
9.2 你还要装备几个评价指标	260
9.2.1 社交网络大小	260

9.2.2 社交网络关系的完备性	261
9.2.3 节点实力评价	262
9.3 全球某货物贸易中的亲密关系	263
9.3.1 全球某货物贸易数据整合清洗	263
9.3.2 分组和社交网络中心	267
9.3.3 全球某货物交易圈：寻找各自的小伙伴	270
9.4 中国电影演艺圈到底有没有“圈”	276
9.4.1 数据清洗与整形	276
9.4.2 看看演艺圈长什么样	279
9.4.3 谁才是演艺圈的“关系户”	281
9.4.4 用 Apriori 算法查查演艺圈合作的“朋友”关系	283
9.4.5 给范冰冰推荐合作伙伴	284
第 10 章 情感分析：一种准确率高达 90%的新方法？	287
10.1 情感分析及其应用：这是老生常谈	287
10.1.1 情感分析的用途	287
10.1.2 情感分析的方法论	288
10.1.3 有关情感分析的一些知识和方向	289
10.2 文本分析的基本武器：R	290
10.2.1 RJava 包配置	290
10.2.2 Rwordseg 包安装	291
10.2.3 jieba 分词包安装	291
10.3 基于词典的情感分析的效果好过瞎猜吗	292
10.3.1 数据整理及词典构建	292
10.3.2 分词整理	297
10.3.3 情感指数计算	299
10.3.4 方法评价：优、缺点分析	300
10.4 监督式情感分析：挑选训练数据集是所有人心中的痛	301
10.4.1 TFIDF 指标	301
10.4.2 构建语料库	302
10.4.3 随机森林模型	304
10.4.4 算法评估：随机森林应该建多少棵树	308
10.5 一种准确率高达 90%的新方法	316
10.5.1 拿来主义的启示	316
10.5.2 情感词典和规则构建	317

10.5.3 朴素贝叶斯情感分析器	329
10.5.4 支持向量机 (SVM)、决策树等情感分析器	330
10.5.5 如何选择支持 SVM 的核函数	339
10.5.6 情感分类器方法评价	343
10.6 谈谈情感分析的下一步思考	344
 第 11 章 话题模型：很多牛人过不去的坎儿	346
11.1 话题模型与文案文本集	346
11.1.1 任务仍然是以处理 dirty data 开始	347
11.1.2 数据清洗	348
11.2 话题模型中几个重要的数据处理步骤	350
11.2.1 中文分词	350
11.2.2 数据整型	352
11.2.3 怎样设定“阈值”	353
11.3 上帝有多少个色子：话题数量估计	356
11.3.1 通俗地说一遍话题模型	356
11.3.2 主题数估计与交叉检验	357
11.3.3 如何使用复杂度、对数似然值确定主题数	362
11.4 LDA 话题模型竟然能输出这么多关系	368
11.4.1 输出主题——词汇及其概率矩阵	368
11.4.2 输出主题——文档归属及其概率矩阵	369
11.5 话题之间也有社交（衍生）关系吗	370
11.6 话题模型的几个强大衍生品	372
11.6.1 话题模型提取特征词	372
11.6.2 三种方法确定聚类的类数和文本层次聚类	373
11.6.3 漂亮的文本聚类树和批量绘制大类词云图	375
 第 12 章 排名就是简单的推荐系统吗？	378
12.1 全球宜居城市综合实力排行	378
12.1.1 综合实力排行：专家法 VS 数据驱动法	379
12.1.2 怎么比较两个排名结果	382
12.2 协同过滤推荐系统	383
12.2.1 基于商品的协同过滤系统（ItemCF）	386
12.2.2 基于用户的系统过滤系统（UserCF）	388
12.2.3 推荐系统效果评比	390

第 13 章 生物信息学中的数据挖掘案例	392
13.1 生物信息学与 R 语言	392
13.2 生物信息学中常用的软件包	392
13.2.1 软件包简介	392
13.2.2 数据表示方式——对象类 (class)	393
13.2.3 生物信息学 R 包简介：Bioconductor 和 CRAN	393
13.2.4 ape 包	394
13.2.5 读懂你的对象	404
13.2.6 修改工具包中的函数以适应新情况	407
第 14 章 产品化：关于内存、速度和自动化	411
14.1 不同终端调用、自动化执行 R 脚本及参数传递	411
14.2 与速度、内存、并行相关的程序优化	414

轻松注册成为博文视点社区用户 (www.broadview.com.cn)，扫码直达本书页面。

- **下载资源：**本书如提供示例代码及资源文件，均可在 [下载资源](#) 处下载。
- **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）。
- **交流互动：**在页面下方 [读者评论](#) 处留下您的疑问或观点，与我们和其他读者一同学习交流。

页面入口：<http://www.broadview.com.cn/33861>

