

大数据人才培养规划教材

以解决实际问题为**学习目标**

以实战案例贯穿为**学习手段**



R语言

商务数据分析实战

R Language Business Data Analysis

韩宝国 张良均 ● 主编

林柳琳 何展鸿 施兴 ● 副主编



中国工信出版集团



人民邮电出版社
POSTS & TELECOM PRESS

大数据人才培养规划教材



R语言 商务数据分析实战

R Language Business Data Analysis

韩宝国 张良均 ● 主编
林柳琳 何展鸿 施兴 ● 副主编

人民邮电出版社
北京

图书在版编目 (C I P) 数据

R语言商务数据分析实战 / 韩宝国, 张良均主编. --
北京: 人民邮电出版社, 2018. 4
大数据人才培养规划教材
ISBN 978-7-115-47448-3

I. ①R… II. ①韩… ②张… III. ①程序语言—程序设计—应用—商业统计—统计数据—统计分析—教材
IV. ①F712.3-39

中国版本图书馆CIP数据核字(2018)第016534号

内 容 提 要

本书以任务为导向, 较为全面地介绍了商务领域中 R 语言数据分析的应用。全书共 9 章, 介绍商务领域不同方向项目的数据分析方法, 具体内容包括 R 语言数据分析概述、商品零售购物篮分析、航空公司客户价值分析、财政收入预测分析、金融服务机构资金流量预测、P2P 信用贷款风险控制、电子商务网站智能推荐服务、电商产品评论数据情感分析、餐饮企业综合分析。除第 1 章外, 本书各章都包含了实训与课后习题, 通过练习和操作实践, 帮助读者巩固所学的内容。

本书可以作为高校大数据专业或商科类专业教材, 也可作为大数据技术爱好者的自学用书。

-
- ◆ 主 编 韩宝国 张良均
副 主 编 林柳琳 何展鸿 施 兴
责任编辑 左仲海
责任印制 马振武
 - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号
邮编 100164 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京市艺辉印刷有限公司印刷
 - ◆ 开本: 787×1092 1/16
印张: 15 2018 年 4 月第 1 版
字数: 332 千字 2018 年 4 月北京第 1 次印刷
-

定价: 45.00 元

读者服务热线: (010)81055256 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

大数据专业系列图书

编写委员会

编委会主任：余明辉 聂哲

编委会成员（按姓氏笔画排序）：

王玉宝	王宏刚	王海	石坤泉	冯健文
刘名军	刘晓玲	刘晓勇	许昊	麦国炫
李红	李怡婷	杨坦	杨征	杨惠
肖永火	肖刚	肖芳	吴勇	邱伟绵
何小苑	何贤斌	何燕	汪作文	张玉虹
张红	张良均	张健	张凌	张敏
张澧生	陈胜	陈浩	林志章	林昆
林碧娴	欧阳国军	易琳琳	周龙	周东平
郑素铃	官金兰	赵文启	胡大威	胡坚
胡洋	钟阳晶	施兴	姜鹏辉	敖新宇
莫芳	莫济成	徐圣兵	高杨	郭信佑
黄华	黄红梅	梁同乐	焦正升	雷俊丽
詹增荣	樊哲			



序

PREFACE

随着大数据时代的到来，移动互联网络和智能手机迅速普及，多种形态的移动互联应用蓬勃发展，电子商务、云计算、互联网金融、物联网等不断渗透并重塑传统产业，大数据当之无愧地成了新的产业革命核心。

未来 5~10 年，我国大数据产业将会是一个飞速发展时期，社会对大数据相关专业人才有着巨大的需求。目前，国内各大高校都在争相设立或准备设立大数据相关专业，以适应地方产业发展对战略性新兴产业的人才需求。

人才培养离不开教材，大数据专业是 2016 年才获批的新专业，目前还没有成套的系列教材，已有教材也存在企业案例缺失等亟须解决的问题。由广州泰迪智能科技有限公司和人民邮电出版社策划，校企联合编写的这套图书，犹如大旱中的甘露，可以有效解决高校大数据相关专业教材紧缺的困境。

实践教学是在一定的理论指导下，通过引导学习者的实践活动，从而传承实践知识、形成技能、发展实践能力、提高综合素质的教学活动。目前，高校教学体系的设置有诸多限制因素，过多地偏向理论教学，课程设置与企业实际应用切合度不高，学生无法把理论转化为实践应用技能。课程内容设置方面看似繁多又各自为“政”，课程冗余、缺漏，体系不健全。本套图书的第一大特点就是注重学生的实践能力培养，根据高校实践教学中的痛点，首次提出“鱼骨教学法”的概念。以企业真实需求为导向，学生学习技能紧紧围绕企业实际应用需求，将学生需掌握的理论知识，通过企业案例的形式进行衔接，达到知行合一、以用促学的目的。

大数据专业应该以大数据技术应用为核心，紧紧围绕大数据应用闭环的流程进行教学，才能够使学生从宏观上理解大数据技术在行业中的具体应用场景及应用方法。高校现有的大数据课程集中在如何进行数据处理、建模分析、调整参数，使得模型的结果更加准确。但是，完整的大数据应用却是一个容易被忽视的部分。本套图书的第二大特点就是围绕大数据应用的整个流程，从数据采集、数据迁移、数据存储、数据

分析与挖掘，最终到数据可视化，覆盖完整的大数据应用流程，涵盖企业大数据应用中的各个环节，符合企业大数据应用真实场景。

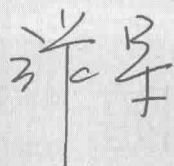
我很高兴看到这套书的出版，也希望这套书能给更多的高校师生带来教学上的便利，帮助读者尽快掌握本领，成为有用之才！

教育部长江学者特聘教授

国家杰出青年基金获得者

IEEE Fellow

华南理工大学计算机与工程学院院长



2017年12月



前言

FOREWORD

本书特色

本书定位于 R 语言数据分析应用教材，深入浅出地介绍 R 语言应用的相关知识，包括 R 语言数据分析的基础知识和商务领域中多个使用 R 语言进行数据分析的实战项目。本书涉及的知识点简要精到，实践操作性强。使用本书能有效指导读者对 R 语言数据分析的学习理解及开发应用。

本书采用了以任务为导向的教学模式，按照解决实际任务的工作流程，逐步展开学习相关的理论知识点，推导生成可行的解决方案，最后落实在任务实现环节。除第 1 章外，本书各章紧扣任务需求展开，不堆积知识点，着重于解决思路的启发与解决方案的实施。通过从任务需求到实现这一完整工作流程的体验，读者可对 R 语言数据分析技术真正理解与掌握。

本书适用对象

- 开设数据分析课程的高校教师和学生

目前国内不少高校将数据分析引入教学中，在计算机、数学、自动化、电子信息、金融等专业开设了与数据分析相关的课程，但目前这一课程教学用的相关教材没有统一，有的使用传统的 SPSS、SAS 等统计工具，并没有使用 R 语言作为数据分析工具。本书提供了 R 语言相关技术的介绍、原理、实践、企业应用等，能有效指导高校教师和学生使用 R 语言解决企业实际问题，为以后的工作打下良好基础。

- 数据分析开发人员

这类人员可以在理解数据分析、应用需求和设计方案的基础上，结合书中提供的 R 语言使用方法快速实现数据分析和应用编程。

- 进行数据分析应用研究的科研人员

许多科研院所为了更好地对科研工作进行管理，纷纷开发了适应自身特点的科研业务管理系统，并在使用过程中积累了大量的科研数据。R 语言可以提供一个优异的环境对这些数据进行分析和应用。

- 关注高级数据分析的人员

R 语言作为一个专业的数据分析软件，能为数据分析人员提供可靠依据。

代码下载及问题反馈

为方便读者实践与练习，书中提供全部项目的数据文件及源代码，读者可登录人民邮电出版社教育社区（www.ryjiaoyu.com）或“泰迪杯”全国数据挖掘挑战赛网站（www.tipdm.org/tj/1307.jhtml）下载。为满足广大教师授课需要，我们还特意提供了 PPT 课件，读者可以从“泰迪杯”数据挖掘挑战赛网站（<http://www.tipdm.org/tj/840.jhtml>）下载申请表，填写后发送至指定邮箱；其他图书资源，读者可通过泰迪大数据挖掘微信公众号（TipDataMining）或者热线电话（40068-40020）进行在线咨询获取。



我们已经尽最大努力避免在文本和代码中出现错误，但是由于水平有限，编写时间仓促，书中难免存在一些不足和疏漏之处。如果您有宝贵意见，欢迎发送邮件至邮箱 13560356095@qq.com，期待得到您真挚的反馈。同时，本书的内容更新将及时在“泰迪杯”全国数据挖掘挑战赛网站上发布，读者可以登录网站或关注泰迪大数据挖掘微信公众号查阅相关信息。

编者

2017年11月

目 录 CONTENTS

第 1 章 R 语言数据分析概述	1	实训 使用 Apriori 算法对西饼屋 订单进行关联分析	25
任务 1.1 认识数据分析	1	课后习题	25
1.1.1 掌握数据分析的概念	1	第 3 章 航空公司客户价值分析	28
1.1.2 熟悉数据分析的流程	2	任务 3.1 了解航空公司现状与 客户价值分析	28
1.1.3 了解数据分析应用场景	4	3.1.1 了解航空公司现状	28
任务 1.2 熟悉 R 语言数据分析 工具	5	3.1.2 了解客户价值分析	30
1.2.1 了解数据分析常用工具	6	3.1.3 熟悉航空客户价值分析的 步骤与流程	30
1.2.2 了解 R 语言数据分析的优势	6	任务 3.2 预处理航空客户数据	31
1.2.3 了解 R 语言数据分析常用的 Packages	7	3.2.1 处理数据缺失值与异常值	31
小结	10	3.2.2 构建航空客户价值分析的 关键特征	31
课后习题	10	3.2.3 标准化 LRFMC 的 5 个特征	35
第 2 章 商品零售购物篮分析	12	3.2.4 任务实现	36
任务 2.1 了解购物篮分析	12	任务 3.3 使用 K-Means 算法进行 客户分群	37
2.1.1 分析商品零售企业现状	12	3.3.1 了解 K-Means 聚类算法	37
2.1.2 了解某商品零售企业基本 数据情况	13	3.3.2 分析聚类结果	38
2.1.3 熟悉购物篮分析的步骤与流程	13	3.3.3 模型应用	41
任务 2.2 分析商品销售状况	14	3.3.4 任务实现	42
2.2.1 分析热销商品	14	小结	43
2.2.2 分析商品结构	15	实训	43
2.2.3 任务实现	17	实训 1 处理信用卡数据异常值	43
任务 2.3 使用 Apriori 关联规则 构建购物篮分析模型	18	实训 2 构造信用卡客户风险 评价关键特征	45
2.3.1 了解 Apriori 算法的基本 原理与使用方法	18	实训 3 构建 K-Means 聚类模型	45
2.3.2 分析结果	23	课后习题	46
2.3.3 任务实现	24		
小结	24		

R 语言商务数据分析实战

第 4 章 财政收入预测分析	48
任务 4.1 了解财政收入预测的背景与方法	48
4.1.1 分析财政收入预测背景	48
4.1.2 了解财政收入预测的方法	50
4.1.3 熟悉财政收入预测的步骤与流程	51
任务 4.2 分析财政收入数据特征的相关性	51
4.2.1 了解相关性分析	51
4.2.2 分析计算结果	52
4.2.3 任务实现	53
任务 4.3 使用 Lasso 回归方法选取财政收入预测的关键特征	53
4.3.1 了解 Lasso 回归方法	53
4.3.2 分析 Lasso 回归结果	54
4.3.3 任务实现	54
任务 4.4 使用灰色预测和 SVR 构建财政收入预测模型	55
4.4.1 了解灰色预测算法	55
4.4.2 了解 SVR 算法	56
4.4.3 分析预测结果	58
4.4.4 任务实现	60
小结	61
实训	61
实训 1 求取企业所得税各特征间的相关系数	61
实训 2 选取企业所得税预测关键特征	62
实训 3 构建企业所得税预测模型	62
课后习题	62
第 5 章 金融服务机构资金流量预测	64
任务 5.1 了解金融服务机构现状与资金流量预测	64
5.1.1 分析金融服务机构现状	64
5.1.2 认识资金流量预测	65
5.1.3 熟悉金融服务机构资金流量预测的步骤与流程	66
任务 5.2 检验数据的平稳性	67
5.2.1 检验平稳性	67
5.2.2 处理非平稳序列	69
5.2.3 任务实现	71
任务 5.3 检验数据的纯随机性	72
5.3.1 了解纯随机性检验	73
5.3.2 检验纯随机性	73
5.3.3 任务实现	74
任务 5.4 建立 ARIMA 模型	74
5.4.1 了解 ARIMA 模型	74
5.4.2 识别模型阶数	75
5.4.3 建立 ARIMA 模型	76
5.4.4 任务实现	81
小结	83
实训	83
实训 1 检验资金赎回数据的平稳性与纯随机性	83
实训 2 识别资金赎回数据集的阶数	83
实训 3 构建 ARIMA 模型	83
课后习题	84
第 6 章 P2P 信用贷款风险控制	85
任务 6.1 认识 P2P 信贷行业的风险控制	85
6.1.1 分析 P2P 信贷行业的现状	86
6.1.2 了解某 P2P 平台数据情况	86
6.1.3 熟悉用户逾期预测的步骤与流程	87
任务 6.2 探索 P2P 信贷用户逾期的相关因素	88
6.2.1 分析用户信息完善程度与逾期率的关系	88
6.2.2 分析用户信息修改情况与逾期率的关系	89
6.2.3 分析用户所在区域经济发展情况与逾期率的关系	90
6.2.4 分析借款月份与逾期率的关系	91
6.2.5 任务实现	92

任务 6.3 预处理 P2P 信贷	
用户数据	95
6.3.1 使用第三方平台信息构建新特征	95
6.3.2 对登录信息表与更新信息表	
进行长宽表转换	95
6.3.3 清洗 P2P 信贷数据	97
6.3.4 任务实现	98
任务 6.4 构建用户逾期还款概率	
预测模型	107
6.4.1 了解 GBM 算法	107
6.4.2 评价 GBM 模型	108
6.4.3 分析结果	109
6.4.4 任务实现	109
小结	111
实训	111
实训 1 探索某银行贷款数据规律	111
实训 2 预处理某银行贷款数据	111
实训 3 使用 GBM 算法构建信贷	
审批模型	111
课后习题	112
第 7 章 电子商务网站智能推荐服务	113
任务 7.1 了解某网站现状与	
智能推荐系统	113
7.1.1 分析某网站现状	113
7.1.2 了解智能推荐服务	115
7.1.3 熟悉网站智能推荐的	
步骤与流程	116
任务 7.2 使用 R 连接数据库	
并提取数据	117
7.2.1 访问数据库	117
7.2.2 任务实现	118
任务 7.3 统计网页整体流量状况	118
7.3.1 分析网页类型	119
7.3.2 分析网页点击次数	122
7.3.3 分析网页排名	123
7.3.4 任务实现	124
任务 7.4 预处理网页浏览数据	130
7.4.1 删除不符合规则的网页	130
7.4.2 还原翻页网址	131
7.4.3 划分正确的网页类别	131
7.4.4 选择用户和用户访问网页记录	132
7.4.5 任务实现	133
任务 7.5 构建智能推荐模型	136
7.5.1 了解协同过滤算法	136
7.5.2 评价智能推荐模型	139
7.5.3 分析模型结果	142
7.5.4 任务实现	142
小结	144
实训 实现 MovieLens 电影数据的	
智能推荐	144
实训 1 清洗 MovieLens 原始数据	144
实训 2 构建 MovieLens 智能	
推荐模型	144
实训 3 评估推荐系统模型	145
课后习题	145
第 8 章 电商产品评论数据情感分析	147
任务 8.1 了解电商企业现状与	
文本情感分析流程	147
8.1.1 分析电商企业现状	147
8.1.2 了解电商产品评论数据	148
8.1.3 实现电商评论数据情感分析的	
步骤与流程	149
任务 8.2 获取电商产品评论数据	149
8.2.1 了解 R 语言获取网络数据的	
方法	149
8.2.2 了解数据获取的方法	151
8.2.3 任务实现	153
任务 8.3 对电商产品评论数据	
进行预处理	156
8.3.1 去除评论数据中的重复数据	156
8.3.2 清洗评论数据	156
8.3.3 对评论数据进行分词	157
8.3.4 去除停用词	158
8.3.5 提取有意义的评论	159
8.3.6 绘制词云查看分词效果	160
8.3.7 任务实现	162

R 语言商务数据分析实战

任务 8.4 评论数据情感倾向分析	163
8.4.1 匹配情感词	164
8.4.2 修正情感倾向	164
8.4.3 检验情感分析效果	164
8.4.4 任务实现	165
任务 8.5 使用 LDA 模型进行主题分析	169
8.5.1 了解 LDA 主题模型	169
8.5.2 寻找最优主题数	171
8.5.3 进行 LDA 主题分析	171
8.5.4 评价主题分析结果	172
8.5.5 任务实现	173
小结	176
实训	176
实训 1 清洗酒店评论原始数据	176
实训 2 对酒店评论数据进行预处理	176
实训 3 使用 LDA 模型建模并分析酒店评论	177
课后习题	177
第 9 章 餐饮企业综合分析	179
任务 9.1 了解餐饮企业分析需求	179
9.1.1 分析餐饮企业现状与需求	180
9.1.2 了解餐饮企业数据基本状况	181
9.1.3 熟悉餐饮企业数据分析的步骤与流程	183
任务 9.2 统计餐饮菜品数据	184
9.2.1 统计每日用餐人数与销售额	184
9.2.2 统计菜品热销度	190
9.2.3 统计菜品的毛利率	191
9.2.4 任务实现	192
任务 9.3 使用 ARIMA 算法预测销售额	194
9.3.1 检验平稳性和纯随机性	194
9.3.2 构建 ARIMA 模型	196
9.3.3 任务实现	198
任务 9.4 使用协同过滤算法实现菜品的智能推荐	201
9.4.1 选取特征	202
9.4.2 使用基于物品的智能推荐算法进行推荐	202
9.4.3 了解基于用户的智能推荐算法	203
9.4.4 分析协同过滤结果	203
9.4.5 任务实现	204
任务 9.5 使用 Apriori 算法实现菜品的关联分析	207
9.5.1 构建 Apriori 模型	207
9.5.2 分析关联规则结果	209
9.5.3 任务实现	210
任务 9.6 使用 K-Means 算法进行客户价值分析	214
9.6.1 构建关键特征	214
9.6.2 构建 K-Means 模型	214
9.6.3 分析 K-Means 模型结果	215
9.6.4 任务实现	217
任务 9.7 用决策树算法实现餐饮客户流失预测	219
9.7.1 了解客户流失	219
9.7.2 了解决策树算法	220
9.7.3 构建客户流失特征	221
9.7.4 分析决策树模型结果	223
9.7.5 任务实现	223
小结	226
实训	226
实训 1 使用 ARIMA 模型预测网站访问量	226
实训 2 使用决策树算法实现运营商客户流失预测	227
实训 3 使用协同过滤算法实现网站的智能推荐	227
实训 4 使用 Apriori 算法实现网站的关联分析	227
实训 5 使用 K-Means 算法实现运营商客户价值分析	228
课后习题	228



第 1 章 R 语言数据分析概述

当今社会，网络和信息技术开始渗透进人类日常生活的方方面面，产生的数据量也呈现指数型增长的态势。如何管理和使用这些数据成为一个全新的领域——数据科学领域。R 语言在最近十年中受到了大量数据科学家的青睐，为数众多的数据科学领域学者和从业者使用 R 语言完成数据科学相关的工作，其中最突出的就是数据分析师。



学习目标

- (1) 掌握数据分析的概念。
- (2) 熟悉数据分析的流程。
- (3) 了解数据分析的应用场景。
- (4) 了解 R 语言在数据分析中的优势。
- (5) 了解 R 语言数据分析常用的 Packages。

任务 1.1 认识数据分析



任务描述

数据分析作为大数据技术的重要组成部分，近年来随着大数据技术逐渐发展和成熟。数据分析技能，被认为是数据科学领域中数据从业人员都需要具备的技能之一。与此同时，数据分析师也成为时下最热门职业之一。数据分析技能的掌握是一个循序渐进的过程。明确数据分析概念、分析流程、分析方法等相关知识是迈出数据分析的第一步。



任务分析

- (1) 掌握广义的数据分析和狭义的数据分析的概念。
- (2) 掌握典型的数据分析流程。
- (3) 了解七大类常见的数据分析应用场景。

1.1.1 掌握数据分析的概念

数据分析是指用适当的分析方法对收集来的大量数据进行分析，提取有用信息和形成

结论并对数据加以详细研究和概括总结的过程。随着计算机技术的全面发展，企业生产、收集、存储和处理数据的能力大大提高，数据量与日俱增。而在现实生活中，需要把这些繁多、嘈杂的数据运用统计分析进行萃取和提炼，以此研究出数据的发展规律，然后帮助企业管理层做出决策。

广义的数据分析包括狭义数据分析和数据挖掘。狭义的数据分析是指根据分析目的，采用对比分析、分组分析、交叉分析和回归分析等分析方法，对收集的数据进行处理与分析，提取有价值的信息，发挥数据的作用，得到一个特征统计量结果的过程。数据挖掘则是从大量的、不完全的、有噪声的、模糊的、随机的实际应用数据中，通过应用聚类、分类、回归和关联规则等技术，挖掘潜在价值的过程。

图 1-1 展示了广义数据分析的概念。综合上述内容，广义数据分析是指依据一定的目标，通过统计分析、聚类和分类等方法发现大量数据中的目标隐含信息的过程。

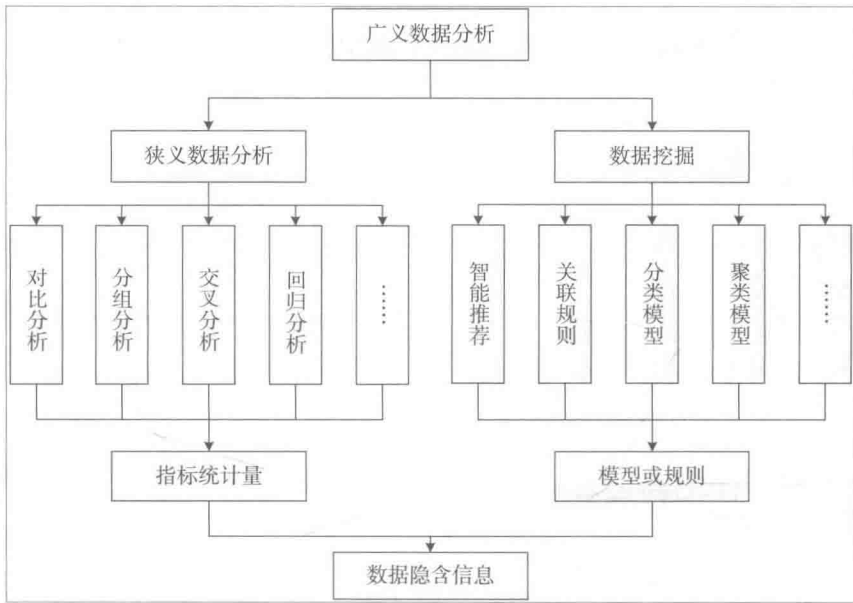


图 1-1 数据分析的概念

1.1.2 熟悉数据分析的流程

数据分析已经逐渐演化为一种解决问题的过程，甚至是一种方法论。虽然每个公司都会根据自身需求和目标创建最适合的数据分析流程，但数据分析的核心步骤是一致的。图 1-2 所示是一个典型的数据分析流程。

1. 需求分析

需求分析一词来源于产品设计，主要是指从用户提出的需求出发，挖掘用户内心的真实意图，并转化为产品需求的过程。产品设计的第一步就是需求分析，也是最关键的一步，因为需求分析决定了产品方向。错误的需求分析，会导致在产品实现过程中走入错误方向，对企业造成损失。

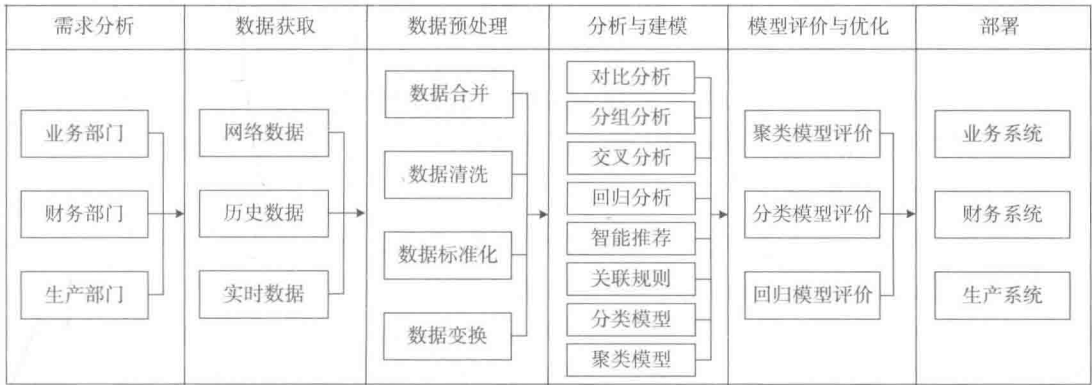


图 1-2 数据分析流程

数据分析中的需求分析是数据分析环节的第一步和最重要的步骤之一，决定了后续分析的方向和方法。数据分析中的需求分析的主要内容是，根据业务、生产、财务等部门的需要，结合现有的数据情况，提出数据分析需求的整体分析方向和分析内容，最终与需求方达成一致意见。

2. 数据获取

数据是数据分析工作的基础，是指根据需求分析的结果提取和收集数据。数据获取主要有两种方式：网络数据与本地数据。网络数据是指存储在在互联网中的各类视频、图片、语音和文字等信息。本地数据则是指存储在本地数据库中的生产、营销和财务等系统的数据。本地数据按照数据时间又可以划分为两部分：历史数据与实时数据。历史数据是指系统在运行过程中遗存下来的数据，其数据量随系统运行时间增加而增长。实时数据是指最近一个单位时间周期（月、周、日、小时等）产生的数据。

在数据分析过程中，具体使用哪种数据获取方式，依据需求分析的结果而定。

3. 数据预处理

数据预处理是指对数据进行数据合并、数据清洗、数据标准化和数据变换，使得整体数据变得干净整齐，从而可以直接用于分析建模这一过程的总称。其中数据合并可以将多张互相关联的表格合并为一张表。数据清洗可以去掉数据中的重复、缺失、异常、不一致的数据。数据标准化可以去除特征间的量纲差异。数据变换则可以通过离散化和哑变量处理等技术使数据满足后期分析与建模的数据要求。在数据分析的过程中，数据预处理的各个过程互相交叉，并没有明确的先后顺序。

4. 分析与建模

分析与建模是指通过对比分析、分组分析、交叉分析、回归分析等分析方法，以及聚类模型、分类模型、关联规则、智能推荐等模型与算法，发现数据中有价值的信息，并得出结论的过程。

分析与建模的方法按照目标不同可以分为几大类。如果分析目标是描述客户行为模式的，可采用描述型数据分析方法，同时还可以考虑关联规则、序列规则、聚类等模型。如果分析目标是量化未来一段时间内某个事件发生概率的，则可以使用两大预测分析模型，

R 语言商务数据分析实战

即分类预测模型和回归预测模型。分类预测模型用于预测离散的目标特征，目标特征通常都是二元数据，例如欺诈与否，流失与否，信用好坏等。回归预测模型用于预测连续的目标特征，常见的目标特征有股票价格预测和违约损失率预测等。分类预测模型和回归预测模型的目标都是训练一个模型，使目标特征预测值与实际值之间的误差达到最小。

5. 模型评价与优化

模型评价是指对于已经建立的一个或多个模型，根据其模型的类别，使用不同的指标评价其性能优劣的过程。常用的聚类模型评价指标有 ARI 评价法（兰德系数）、AMI 评价法（互信息）、V-measure 评分、FMI 评价法和轮廓系数等。常用分类模型的评价指标有准确率（Accuracy）、精确率（Precision）、召回率（Recall）、F1 值（F1 Value）、ROC 和 AUC 等。常用的回归模型评价指标有平均绝对误差、均方误差、中值绝对误差和可解释方差值等。

模型优化则是指模型性能在经过模型评价后已经达到了要求，但在实际生产环境应用过程中，发现模型的性能并不理想，继而对模型进行重构与优化的过程。在多数情况下，模型优化和分析与建模的过程基本一致。

6. 部署

部署是指将通过了正式应用数据分析的结果与结论应用至实际生产系统的过程。根据需求的不同，部署阶段可以是一份包含了对现状具体整改的措施的数据分析报告，也可以是将模型部署在整个生产系统的解决方案。在多数项目中，数据分析师提供的是一份数据分析报告或者是一套解决方案，实际执行与部署的是需求方。

1.1.3 了解数据分析应用场景

企业使用数据分析解决不同的问题，实际应用的数据分析场景主要分为以下七大类。

1. 客户分析（Customer Analytics）

客户分析主要是对客户的基本数据信息进行商业行为分析。首先，界定目标客户，根据客户的需求、目标客户的性质、所处行业的特征以及客户的经济状况等基本信息，使用统计分析法和预测验证法，分析目标客户，提高销售效率。其次，了解客户的采购过程，根据采购类型、采购性质进行分类分析，制定不同的营销策略。最后，还可以根据已有的客户特征进行客户特征分析、客户忠诚度分析、客户注意力分析、客户营销分析和客户收益分析。通过有效的客户分析，能够掌握客户的具体行为特征，将客户细分，使得运营策略达到最优，提升企业整体效益等。

2. 营销分析（Sales and Marketing Analytics）

营销分析囊括了产品分析、价格分析、渠道分析、广告与促销分析这 4 类分析。产品分析主要是竞争产品分析，通过对竞争产品的分析制定自身产品策略。价格分析又可以分为成本分析和售价分析，成本分析的目的是降低不必要的成本，售价分析的目的是制定符合市场的价格。渠道分析的目的是对产品的销售渠道进行分析，确定最优的渠道配比。广告与促销分析则能够结合客户分析，实现销量的提升和利润的增加。

3. 社交媒体分析（Social Media Analytics）

社交媒体分析以不同社交媒体渠道生成的内容为基础，实现不同社交媒体的用户分析、

访问分析和互动分析等。用户分析主要根据用户注册信息、登录平台的时间点和平时发表的内容等用户数据，分析用户个人画像和行为特征；访问分析则是通过用户平时访问的内容，分析用户的兴趣爱好，进而分析潜在的商业价值；互动分析是根据互相关注对象的行为，预测该对象未来的某些行为特征。同时，社交媒体分析还能为情感和舆情监督提供丰富的资料。

4. 网络安全 (Cybersecurity)

大规模网络安全事件（例如 2017 年 5 月席卷全球的 WannaCry 病毒）的发生，让企业意识到网络攻击发生时预先快速识别的重要性。传统的网络安全主要依靠静态防御，处理病毒的主要流程是发现威胁、分析威胁和处理威胁。这种情况下，往往在威胁发生以后才能做出反应。新型的病毒防御系统可使用数据分析技术建立潜在的攻击识别分析模型，监测大量网络活动数据和相应的访问行为，识别可能进行入侵的可疑模式，做到未雨绸缪。

5. 设备管理 (Plant and Facility Management)

设备管理同样是企业关注的重点，设备维修一般采用标准修理法、定期修理法和检查后修理法等。其中标准修理法可能会造成设备过剩修理，修理费用高；检查后修理法解决修理费用成本问题，但是修理前的准备工作繁多，设备停歇时间过长。目前企业通过物联网技术收集和分析设备上的数据流，包括连续用电、零部件温度、环境湿度和污染物颗粒等多种潜在特征，建立设备管理模型，从而预测设备故障，合理安排预防性的维护，以确保设备正常作业，降低因设备故障带来的安全风险。

6. 交通物流分析 (Transport and Logistics Analytics)

物流是物品从供应地向接收地的实体流动。它将运输、储存、装卸搬运、包装、流通加工、配送、信息处理等功能有机结合起来以实现用户要求的过程。人们可以通过业务系统和 GPS 定位系统获得数据，根据客户使用的数据构建交通状况预测分析模型，有效预测实时路况、物流状况、车流量、客流量和货物吞吐量，进而提前补货，制定库存管理策略。

7. 欺诈行为检测 (Fraud Detection)

身份信息泄露及盗用事件的数量逐年增长，随之而来的是欺诈行为和交易的增多。公安机关、各大金融机构和电信部门可利用用户基本信息、用户交易信息和用户通话短信信息等数据，识别可能发生的潜在欺诈交易，做到提前预防，未雨绸缪。以大型金融机构为例，通过分类模型分析方法对非法集资和洗钱的逻辑路径分析，找到其行为特征。聚类模型分析方法可以分析相似价格的运动模式。例如对股票进行聚类，可能发现关联交易及内幕交易的可疑信息。关联规则方法分析可以监控多个用户的关联交易行为，为发现跨账号协同的金融诈骗行为提供依据。

任务 1.2 熟悉 R 语言数据分析工具

任务描述

R 语言自 1993 年首次宣布面向公众开放以来，逐渐成为事实上的计算统计学