

大数据时代，R语言让你从数据中发掘价值

R语言与大数据编程实战

李倩星 ◎ 编著



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>

高级大数据人才培养丛书

R语言与大数据编程实战

李倩星 编著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

本书是一本 R 语言入门读物，它旨在帮助读者迅速构建起与数据分析相关的知识体系，并学习如何使用 R 软件实现数据分析方法。无论有无深厚的编程基础或数学基础，本书都能帮助读者成长为一名合格的数据分析师。

本书全面介绍了来自统计分析、机器学习、人工智能等领域的多种数据分析算法，在讲解与之相关的 R 代码时，还讨论了这些算法的原理、优缺点与适用背景。本书按照由易到难的原则组织章节主题，读者将获得最好的阅读体验。通过阅读本书，读者将对 R 语言在数据分析领域的应用有一个全面的认识。这种认识不被特定行业所局限，任何行业的读者都能利用本书介绍的数据分析方法解决本行业的数据分析问题。

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

R 语言与大数据编程实战 / 李倩星编著 . —北京：电子工业出版社，2017.9
(高级大数据人才培养丛书)

ISBN 978-7-121-32634-9

I . ① R … II . ①李… III . ①程序语言—程序设计 IV . ① TP312

中国版本图书馆 CIP 数据核字（2017）第 215708 号

策划编辑：李 冰

责任编辑：李 冰

特约编辑：彭 瑛 赵海军等

印 刷：三河市华成印务有限公司

装 订：三河市华成印务有限公司

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：787×1092 1/16 印张：20 字数：512千字

版 次：2017年9月第1版

印 次：2017年9月第1次印刷

定 价：59.00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：libing@phei.com.cn。

前　　言

R 语言是如今最热门的编程语言之一，它由统计学家开发，在解决数据分析问题时具有先天优势。R 语言是一门新兴的语言，掌握它，就是掌握了一门高效的数据分析软件。随着大数据概念的普及，R 语言能够实现的功能越来越丰富，越来越多的数据分析从业人员产生了学习 R 语言的需求。本书迎合时代潮流，讲解了大数据时代下 R 语言渗透最广泛的几个领域，全面介绍了如何使用 R 语言完成数据挖掘工作。对 R 语言编程人员来说，本书是一本不可或缺的工具书。

本书特色

1. 通俗易懂、实用性强，适合各层次读者学习

本书对读者的数学基础或编程基础不做任何要求。在讲解知识点时，本书采用了平实的语言，对每个疑难点都加以详细解释。此外，本书以实用为主旨，秉承“看得懂、学得会、用得上”的编写原则，精心选取了流行于行业前沿的 18 个主题，不仅通俗易懂，还确保读者所学的知识具有实际应用价值。通过阅读本书，读者都能迅速掌握 R 语言的编程技巧及相关的数据分析知识，并在实际工作中立刻应用它们。

2. 条理清晰、结构巧妙，全面盘点数据分析常用算法

数据分析是一个涉及多领域的交叉学科，R 软件的触角同样也能伸展到多个领域。本书选取了统计分析、机器学习、人工智能等多个学科的流行算法作为主题，讲解了如何使用 R 语言实现它们。这些算法有些偏重数学思维，有些偏重编程技巧，本书主要遵循由易到难的顺序排列主题，并尽量把起源于同一学科的算法放在一起。读者可以按照顺序阅读本书，也可以优先选择感兴趣的部分。此外，本书还穿插介绍了与 R 软件相关的一些其他编程主题，这些主题共同形成知识网络，帮助读者迅速成长为能够独当一面的数据科学家。

3. 知识点丰富、可拓展性强，满足读者的多重需求

本书涉及多个学科，全面介绍了 R 软件能够实现的多种算法，满足了读者的三大需求：首先，使用通俗易懂的语言介绍 R 软件，帮助读者实现零基础入门；其次，囊括多种数据分析算法，带领读者全面认识 R 软件的强大之处，帮助读者成长为合格的数据科学家；最后，本书具备较强的可拓展性，从事任何行业的读者都能够从本书中获取适合其行业的知识。本书还给出了 R 语言进阶的线索，无论想向哪一方面进阶，本书都能为读者打造最坚实的基础。

本书内容及体系结构

本书总共 18 章，分别为 R 的基本介绍、原始数据的探索与预处理、R 的数据可视化、R 中参数的估计和检验、R 中的方差分析、R 中的相关分析和回归分析、更高级的数据可视化、R 中的聚类分析和判别分析、R 中的主成分分析和因子分析、R 中的广义线性回归模型、R 中的时间序列模型、R 中的最优化问题、使用 R 绘制地理信息图形、使用 R 构建支持向量机、实现更高效的流程控制和高级循环、R 代码的调试与优化、构建电影评分预测模型、贝叶斯垃圾邮件过滤器模型。这 18 章进一步又分为 5 个部分。

第一部分为本书的第 1~6 章。其中前 3 章展示了 R 软件的一些入门功能，如数据预处理和数据可视化等，后 3 章则介绍了三种基础的统计分析方法，即参数的估计和检验、方差分析、相关分析、回归分析。这 6 章围绕初级的统计方法展开，是数据分析师必备的基本知识。

第二部分为本书的第 7~11 章，这 5 章介绍了更高级的统计方法。其中，第 7 章为第 3 章的延伸，介绍了数据可视化的高级方法，第 8~11 章则介绍了 6 种高级统计分析方法，这部分的内容与第一部分互为补充。

第三部分为本书的第 12~14 章，这部分内容围绕机器学习展开。第 12 章的主题为最优化，是机器学习的基本理论。第 13 章介绍了如何使用矢量化的思想绘制地图。第 14 章则介绍了支持向量机，它是最典型的机器学习算法之一。这部分讲解了更高深的 R 语言编程技巧，讨论了一些 R 软件能够解决的高难度问题。

第 15、16 章可视为本书的第四部分。这两章围绕如何优化 R 代码展开，系统地讨论了如何写出错误较少的、运行速度较快的代码。这部分内容帮助读者建立良好的编程习惯，以及与其他 R 用户更好地协同工作。

第 17、18 章则为本书的最后一部分，这两章分别讨论了一个完整的数据挖掘项目。其中电影评分预测的案例着重于表现数据挖掘的完整流程，包括繁复的数据预处理与反复的模型比较等工作；垃圾邮件过滤的案例则引出 R 软件能够处理的另一个主题——文本分析。

上述划分方法仅为一个参考，本书的 18 章既互相联系又彼此独立，读者可按照上述划分方法阅读本书，也可优先阅读某些章节，如将第 3、7、13 章等与数据可视化相关的三个章节放在一起阅读。

本书读者对象

- 想要了解 R 语言的数据分析从业人员。
- 统计学、金融学、计算机技术与科学等专业的学生。
- 想要提高 R 语言编程能力的数据分析师。
- 希望系统学习统计分析方法的从业人员。
- 其他对 R 语言有兴趣的各类人员。

目 录

第 1 章 R 的基本介绍	1
1.1 强大的 R	1
1.2 R 语言在大数据中的应用	2
1.2.1 R 语言用户行为分析	2
1.2.2 R 语言处理金融大数据	3
1.2.3 R 语言天气数据可视化	4
1.2.4 R 语言医疗大数据分析	4
1.3 R 的安装与启动	5
1.3.1 安装并启动 R	6
1.3.2 安装并启动一个 IDE	7
1.4 R 的向量、矩阵和数组	8
1.4.1 向量的操作方法和固有属性	8
1.4.2 矩阵的操作和运算	10
1.4.3 数组中的维度函数	13
1.5 R 的列表和数据框	14
1.5.1 列表的特性和编辑方法	14
1.5.2 数据框的创建和基本操作	17
1.6 R 数据文件的载入和载出	19
1.6.1 结构化纯文本文件的读取和输出	19
1.6.2 其他文件的读取和输出	22
1.7 向 R 中安装包	23
第 2 章 原始数据的探索与预处理	26
2.1 度量数据集的集中程度	26
2.2 度量数据集的分散程度	27
2.2.1 极值、方差和标准差	27
2.2.2 标准误和偏度系数、峰度系数	29

2.3 创建一个数值摘要表	30
2.4 异常值的观测与说明	32
2.4.1 利用箱线图观测异常值并处理	32
2.4.2 异常值检测的其他情况和说明	34
2.5 缺失值的填补与处理	35
2.5.1 删除缺失值或对其进行简单填补	36
2.5.2 按照相关性对空缺值进行填补	38
第 3 章 R 的数据可视化	40
3.1 plot() 函数和常用的图形参数	40
3.1.1 设置 plot() 函数中的参数	40
3.1.2 修改散点图的坐标并加入标注	43
3.2 经典的基础图形及用途	45
3.2.1 线图	45
3.2.2 直方图	49
3.2.3 箱线图和茎叶图	52
3.3 将图形组合起来	55
3.4 更高的高水平作图函数	57
3.5 更多的常用作图命令	59
第 4 章 R 中参数的估计和检验	62
4.1 使用 R 进行点估计和区间估计	62
4.1.1 简单的点估计和区间估计	62
4.1.2 估计单侧置信区间	65
4.2 与正态总体有关的参数检验	68
4.3 列联表与独立性检验	71
4.4 几种检验数据分布的函数	72
4.5 对非正态总体的区间估计和检验	75
4.5.1 非正态总体的区间估计	75
4.5.2 非参数检验中的符号检验	76
4.5.3 非参数检验中的秩检验	78

第 5 章 R 中的方差分析	80
5.1 方差分析模型的建立.....	80
5.2 单因素方差分析.....	81
5.2.1 单因素方差分析的数学思想与模型	81
5.2.2 检验样本是否满足方差分析的假设条件	82
5.2.3 构建单因素方差分析模型	84
5.3 多因素方差分析.....	87
5.3.1 多因素方差分析的数学思想与模型	87
5.3.2 不考虑交互作用的双因素方差分析	88
5.3.3 考虑交互作用的双因素方差分析	89
5.4 秩检验和协方差分析.....	91
5.4.1 对控制变量应用秩检验方法	91
5.4.2 协方差分析的假设与应用	92
第 6 章 R 中的相关分析和回归分析	94
6.1 多种相关系数的度量和分析.....	94
6.1.1 简单相关系数的计算和检验	94
6.1.2 散布矩阵图和偏相关系数	96
6.1.3 典型相关分析	98
6.2 线性回归分析及其常规参数.....	99
6.2.1 对数据进行预处理	100
6.2.2 构建第一个回归模型	101
6.2.3 修正方程并检验残差	102
6.3 使用逐步回归筛选自变量.....	104
6.3.1 逐步回归的思想与分类	104
6.3.2 构建逐步回归模型	105
6.4 哑变量和逻辑回归.....	107
6.4.1 哑变量和逻辑回归的思想	107
6.4.2 向线性回归模型中纳入哑变量	108

第 7 章 更高级的数据可视化	110
7.1 基础图形的拓展与延伸	110
7.1.1 绘制分类散点图并添加图标	110
7.1.2 绘制含多种类别的密度分布图	112
7.1.3 复合条形图和堆栈条形图	114
7.2 有关多元分布函数的特殊图形	117
7.2.1 星图和脸谱图	117
7.2.2 轮廓图	120
7.2.3 调和曲线图	122
7.3 建立最简单的 3D 图形	123
7.4 如何让图形更美观	125
7.5 更多的绘图包和系统	128
第 8 章 R 中的聚类分析和判别分析	129
8.1 几种聚类分析的异同	129
8.2 使用 R 实现 KNN 聚类	130
8.2.1 KNN 算法的思想和模型	130
8.2.2 使用 R 实现 KNN 聚类	131
8.3 使用 R 实现系统聚类	133
8.3.1 系统聚类的思想和模型	133
8.3.2 使用 R 实现系统聚类	134
8.4 使用 R 实现快速聚类	136
8.4.1 快速聚类的思想和模型	136
8.4.2 使用 R 实现快速聚类	137
8.5 几种判别分析模型综述	140
8.5.1 距离判别模型	140
8.5.2 Fisher 判别模型	142
第 9 章 R 中的主成分分析和因子分析	145
9.1 主成分分析的实现与应用	145
9.1.1 主成分分析的模型假设和数据处理	145
9.1.2 构造一个主成分分析模型	147

9.1.3 计算主成分的综合得分	149
9.2 因子分析的初次构建与完善	150
9.2.1 构造一个简单的因子分析模型	150
9.2.2 计算因子得分并分析	152
9.3 对因子分析模型进行修正	153
9.3.1 修改因子分析模型中的因子个数	153
9.3.2 基于主成分法和主轴因子法进行因子分析	155
9.4 在降维分析的基础上进行回归分析和聚类分析	157
9.4.1 在降维分析的基础上进行回归分析	157
9.4.2 在降维分析的基础上进行聚类分析	160
第 10 章 R 中的广义线性回归模型	162
10.1 一般的广义线性回归模型	162
10.1.1 使用二次函数拟合线性回归模型	162
10.1.2 拟合更多的广义线性模型	164
10.1.3 比较线性模型的优劣	166
10.2 Logistic 线性回归模型	168
10.2.1 Logistic 模型的原理与构建方法	168
10.2.2 Logistic 模型的显著性检验和优势比	170
10.2.3 修正被警告的 Logistic 模型	171
10.3 泊松回归分析模型	173
10.3.1 拟合第一个泊松回归模型	174
10.3.2 泊松回归模型的过散布检验	176
10.4 广义线性模型的交叉验证	178
第 11 章 R 中的时间序列模型	180
11.1 将数据转换为时间序列格式	180
11.1.1 使用 ts() 函数转换数据格式并绘制时间序列曲线	180
11.1.2 使用 zoo() 函数转换数据格式并绘制时间序列曲线	182
11.2 分解时间序列并检验时间序列的自相关性	185
11.2.1 使用经典方法分解时间序列	185
11.2.2 使用 STL 方法分解时间序列	186

11.3 探究时间序列的自相关性	188
11.3.1 使用月图和季度图探究自相关性	188
11.3.2 使用散点图探究自相关性	189
11.4 构建时间序列并预测	191
11.4.1 均值预测、单纯预测和漂移	192
11.4.2 不考虑长期趋势和季节波动的简单指数平滑	195
11.4.3 在指数平滑中加入长期趋势和季节波动	196
11.4.4 自回归移动平均模型	198
第 12 章 R 中的最优化问题	201
12.1 最优化问题简述	201
12.2 黄金分割法	202
12.2.1 黄金分割法和局部最优解	202
12.2.2 使用 R 实现黄金分割法	203
12.3 牛顿最优化方法	205
12.3.1 牛顿法的算法原理	206
12.3.2 在一维情形下实现牛顿迭代法	207
12.3.3 在多维情形下实现牛顿迭代法	209
12.4 最快上升法	210
12.4.1 利用梯度求解上升最快的相邻点	210
12.4.2 构建最快上升法函数并检验	212
12.5 R 中的最优化函数	213
第 13 章 使用 R 绘制地理信息图形	216
13.1 绘制世界、国家、省市地图	216
13.1.1 使用 map() 函数绘制地图	216
13.1.2 另一种绘制地图的方法	217
13.1.3 分省市绘制地图	218
13.2 向地图中添加颜色	220
13.2.1 向地图中添加颜色前的准备工作	220
13.2.2 在地图上添加颜色	221
13.3 向地图中添加标签和线条	222
13.3.1 向地图中添加标签前的准备工作	222

13.3.2 在地图上添加标签	224
13.3.3 在地图上添加线条	225
13.4 使用其他格式的文件优化地图.....	226
第 14 章 使用 R 构建支持向量机	230
14.1 构建一个简单的支持向量机.....	230
14.1.1 支持向量机的算法原理	230
14.1.2 构建一个简单的支持向量机	232
14.1.3 使用其他核函数构建支持向量机	235
14.2 优化支持向量机的参数.....	237
14.2.1 优化参数 degree.....	238
14.2.2 优化参数 cost.....	241
14.2.3 优化参数 gamma	243
14.3 比较支持向量机与 Logistic 回归的优劣.....	246
14.4 比较支持向量机和 KNN 聚类算法的优劣.....	249
第 15 章 实现更高效的流程控制和高级循环	251
15.1 R 中的流程控制	251
15.1.1 if 语句的多种实现方法	251
15.1.2 ifelse 语句与花括号的结合	252
15.1.3 适合多分支情况的 switch 语句	254
15.2 R 中的 for 循环、while 循环和 repeat 循环.....	256
15.2.1 R 中的 for 循环和 while 循环	256
15.2.2 R 中的 repeat 循环	258
15.3 apply 家族中的循环函数	260
15.3.1 R 中的 apply() 函数	260
15.3.2 R 中的 lapply() 函数和 sapply() 函数	263
15.3.3 R 中的 tapply() 函数	265
15.3.4 R 中的 mapply() 函数	268
15.4 更多的高级循环函数	270
15.4.1 R 中的 replicate() 函数和 sweep() 函数.....	270
15.4.2 R 中的 aggregate() 函数	273

第 16 章 R 代码的调试与优化.....	276
16.1 R 代码的常见信息与警告	276
16.1.1 R 代码的正常信息与警告.....	276
16.1.2 R 代码中的警告处理方法.....	278
16.2 R 代码中的错误与错误处理方法	279
16.2.1 使用 try() 函数处理错误信息	279
16.2.2 将 try() 函数与循环相结合	281
16.3 调试 R 代码	282
16.3.1 查看调用栈或暂停代码	282
16.3.2 修改 error 选项	284
16.4 向量化编程方法.....	285
16.4.1 向量化编程思想	285
16.4.2 比较循环和向量的运行速度	286
第 17 章 构建电影评分预测模型	289
17.1 获取数据并探索	289
17.2 利用 recommenderlab 包处理数据	291
17.3 建立模型并评估	293
17.3.1 模型的选择与建立	293
17.3.2 模型之间的比较和评估	295
第 18 章 贝叶斯垃圾邮件过滤器模型	297
18.1 贝叶斯模型中的条件概率	297
18.2 复杂的数据预处理过程	298
18.2.1 利用 for 循环读入多封邮件正文	298
18.2.2 利用 tm 包进一步转换数据格式	300
18.2.3 将 TDM 转换成真正有用的数据框	301
18.3 利用 occurence 值构造分类器	303
18.3.1 完成理论准备并处理测试邮件和普通邮件	303
18.3.2 创建一个函数用于比较概率	305

第1章 R 的基本介绍

作为一门新兴的编程语言，R 是如今值得学习的语言。由统计学家开发出的 R 语言具有许多奇特性质，本章将较为全面地介绍 R 的特性和用途，并讲解 R 的安装方法、变量类型、从其他数据源读取数据、程序包等基本知识。本章帮助读者对 R 形成整体印象，同时本章内容也是后续章节的基石。

1.1 强大的 R

R 语言脱胎于 S 语言，是一门专门用于处理数据探索、统计分析等任务的编程语言。它由统计学家开发完成，在数据分析方面具有天然的优势，运行 R 程序的 R 软件是如今最流行的统计软件之一。

与其他统计软件相比，R 软件最特别的地方在于它是开源的。这同时意味着：第一，R 是免费的；第二，R 的用户能够自由地参与到 R 的开发中。R 社区将它的忠实用户聚合在一起，这些用户主要由统计学家、计算机学家、数据分析师等组成，不同领域的用户在 R 社区中交流碰撞，协助 R 核心团队丰富和完善 R 的功能。

R 的用户之间具有非常紧密的联系，他们最大的贡献是创建了形形色色的程序包，这些程序包分别封装了一些具有特定作用的函数。如今，R 软件已经内置了非常丰富的各类函数库，能够满足绝大多数统计人员的各类需求，它的制图功能也远超其他统计软件。

R 的另一个特点在于它支持混合型的编程范式。R 是一种解释型的语言，当用户在 R 软件中编写好一条代码后，R 会立即执行它。这种做法的好处在于用户可以即时地看到程序的返回结果，在作图时尤其方便。R 是一种面向对象的语言，同时它也支持函数式编程，即用户可以在 R 中调用现成的或自己编写的函数，这一点与 C 语言较为相似，但 R 要比 C 语言更加灵活。

尽管 R 的优点很突出，但它也同样具有局限之处。首先，R 语言的编程原理较为传统，在处理数据时，R 需要将数据全部载入内存，这一点极大地影响了 R 的运行效率，尽管如今的计算机内存做得越来越大，但在有些大规模数据集的处理工作中，R 还是会显得不够得力。其次，R 软件的保密性不如 SAS 等统计软件好，这限制了 R 在大型商业项目中的应用。最后，由于 R 软件是由统计学家开发的，因此其语法设计并不特别严谨，有时它会出现一些奇怪的错误。

随着大数据时代的到来，R 语言正被越来越多的人关注，不仅是统计分析和数据挖掘，一些研究机器学习和模式识别的专家同样关注到 R 的发展。根据 TIOBE 提供的编程语言排行榜，R 语言的流行程度在近几年内已经飙升至前十名，其火爆程度只有 Python 才能与其比肩，而同为统计软件的 SAS 和 MATLAB 则一直徘徊于二三十名的位置。

R 的优点使它广泛地流行于统计人员和中小型商业公司中。Google、百度等互联网巨头则将 R 语言看作一个沙盘，使用软件验证各种数据模型的可行性，并最终使用其他语言实现。随着 R 的用户越来越多样化，其可扩展能力进一步强化，能够解决的问题也越来越丰富。如今，金融、医药、教育、社会科学等每一个需要数据分析的领域都需要精通 R 的人才。

1.2 R 语言在大数据中的应用

R 语言的起源是统计学家为解决数据分析领域问题而开发的语言，所以和 MATLAB、Python 等可用于数据处理的语言相比，在数据分析处理方面具有一些独特优势，本节将讲述 R 语言在大数据领域中的典型应用。

1.2.1 R 语言用户行为分析

近几年，淘宝、京东等几家电商的价格战打得不亦乐乎，而从电商发布的战果来看，几败具伤的价格战已经无法保证电商的利润，他们开始转向利用大数据分析工具对用户行为进行分析，通过对大数据的充分使用和挖掘在商战中获胜。

何为“用户行为分析”？简单的理解就是对用户在网站上发生的所有行为进行分析，找到里面的规律和用户感兴趣的信息。这些行为如搜索商品信息、浏览新闻、购物评价和打分、美团点评、加入收藏列表、加入购物车、购买、使用特价购物券、换货和退货等；除此之外，还包括在第三方网站上的相关行为分析，如比价、看相关评测、参与讨论、社交媒体上的交流、与好友互动等。电子商务的突出特点就是可以通过后台收集到大量客户在购买前的行为信息，而这些信息对于分析用户行为至关重要。

全球电子商务的创始者之一亚马逊公司结合大数据分析工具，以迅雷不及掩耳之势，彻底颠覆了很多行业的市场规则及竞争关系。亚马逊取胜的根本原因在于对数据的战略性认识和使用，亚马逊通过传统门店无法比拟的互联网手段，获取了极其丰富的用户行为信息，并且进行深度分析与挖掘。

电商通过对用户行为的分析，可以制定更加贴心的服务。例如，当客户浏览了多款手机而没有做购买的行为时，系统会把适合客户的品牌、价位和类型的多款手机信息推送到客户的账户，只要客户登录系统，就可以看到这些推送信息。这样的个性化推荐服务往往会起到非常好的效果，不仅可以提高客户购买的意愿，缩短购买的路径和时间，还可以在比较恰当的时机捕获客户的最佳购买冲动，提升用户体验，是一个一举多得的好方式。

在电商领域中，用户的行为信息量之大令人难以想象，据不完全统计，用户在电商网站购买一个商品前，平均会浏览 3~5 个网站（淘宝、京东、1 号店等）、30~36 个页面，统计起来对于一个一天有近百万访问量的中型电商，一天就会有 1TB 左右的活跃数据。

除了电商，爱奇艺、优酷、土豆网等各类视频网站已成为人们娱乐和学习的重要平台，视频网站成功与否的最重要衡量标准就是用户的满意度。由于 Web 应用能够以很细的粒

度、很高的频度不断记录用户的行为轨迹，这些数据中隐藏着用户的访问习惯、兴趣偏好及情绪变化等信息，同时也隐含着用户群体行为的规律和发展趋势。挖掘深藏在数据背后的知识，能够发现用户习惯的观看流程，访问网站的时间及喜好的视频，各视频间存在的关联关系等。掌握了这些知识，就能科学解决用户跳出等问题，根据用户访问习惯改进网站服务流程，以及针对用户进行个性化服务，变革传统的网站管理和运营模式，主动提升用户的体验以促进视频网站的快速发展。

纵观国内外成功的电商和视频网站等互联网企业，对用户行为信息的分析和使用，无不在这个必争之地做了大量投入。他们对数据战略性的高度认识和使用，非常值得国内的电商学习和借鉴。

R 语言作为能够进行交互式数据分析和探索的强大平台，具备一套完整的数据处理、计算和制图软件系统，在电商用户行为数据分析和挖掘领域应用广泛。基于 R 语言对视频网站的用户在线行为数据进行分析挖掘，通过对采集到的用户行为数据进行引流指标、黏性指标的分析和挖掘，可以得到网页被访问的频繁度、停留时间、用户观看视频的喜好等信息。将分析的结果应用到网站运营与管理中，不仅能够为网站个性化服务、精准推销和开发新型业

务模式提供技术和理论支撑，而且能够把握当前网络关注的热点问题，从而正确引导用户的网络舆论方向。

1.2.2 R 语言处理金融大数据

最近几年，数据分析方法在商业和金融市场上的重要性持续增加，因为我们有丰富的数据环境，经济和金融市场的数据相比以前更加综合和完整，在许多国家成百上千个变量的数据可以更系统、更精确地搜集，计算处理上的便利和统计软件包的使用使得对复杂的高维金融数据的分析成为可能，通过互联网可以很容易地应用开源软件包下载公开的金融数据，如 R 语言和软件开发环境，所有这些软件的特性和功能免费公开，因而被广泛使用。

2008 年的金融危机在某种程度上是由错误的金融模型造成的，既有模型过于简单的原因，又有模型过于复杂的原因，房地产经纪人和买家依赖于一个隐式模型，它表明价格已经在高位，且还会继续上涨。贷款人使用统计模型来对打包的按揭产品进行分析设计，这似乎可以奇迹般地降低风险，然而最后的结果是灾难性的，在 4 年之后仍然可以感受到房地产泡沫的影响。

那么，如何进行有用的并且没有危险的金融分析呢？首先应该对金融数据有一个基本的理解，尤其是时间序列数据，因为不确定性是主要的影响因素，比如可以用概率模型来描述资产收益率的频率分布，利用时间序列模型对数据进行描述平滑和季节调整。

使用 R 语言进行金融分析时，可以利用 R 语言的优势，研究分析基于 Hadoop 存储证券的日内交易数据，通过 RHive 连接 R 语言与 Hive，建立相关性算法模型，在历史数据中回测，构建投资决策组合，并生成可视化结果用于展示。

R 语言 quantmod 包是股市金融建模应用比较多的一个包。该包从多个数据源获取历史数据，绘制金融数据图表，以及在金融数据图表中添加各种技术指标，通过多种金融

模型分析，辅助股票筛选和判断。该包获取数据的来源主要有两个：Yahoo! 和 Google。最常用的是 Yahoo! 中的各种数据。但该包只能获取股票的历史交易记录信息，如最大值、最小值、开盘价、收盘价及成交量。在此基础上利用股票的历史数据，通过 R 语言建立模型，并对数据进行分析，从模型的检验决定未来的交易行为。

1.2.3 R 语言天气数据可视化

R 语言天气数据可视化，就是通过获取天气的历史大数据，使用 R 语言和可视化技术，展示中国每个省份的天气情况，给准备旅游的朋友提供一种出行提示。除了旅游，在传统零售行业，雨天大概会影响相对于晴天 30% ~ 40% 的销售业绩，所以从网上获取天气数据进行分析，并根据天气数据做出预测，提前做好预防措施和提醒业务人员，把损失减少到最低就显得十分重要。要实现天气数据可视化，需要实现的功能和遇到的问题如下。

- 天气数据：数据从哪里找到；如何下载；如何存储。
- 定时任务：天气数据需要每日更新，图片需要每日新生成。
- 地图和天气可视化：要把中国行政区划图和天气数据（包括风力方向可视化）结合在一起绘图，让用户一眼就能看明白。
- Web 展示：通过可视化技术，我们生成的只是一张静态图片，要如何发布到 Web 端进行展示。
- 微博：结合新浪微博，让更多的用户看到并使用这个应用。
- 用户交互：用户可以查看不同日期、不同类型的图片，用户还可以通过微博分享。

从上面的描述中，单独使用一种语言不容易实现这些功能。单独用 PHP 开发，做一个 Web 网站非常容易，连接新浪微博也有现成的 SDK 可以调用，获取数据及存储也不麻烦，但是如何实现地图和天气数据的可视化？这个是 R 语言的强项，所以，可以将 R 语言和 PHP 语言相结合，发挥 R 语言的优势，用 R 语言的 rvest 包就可以方便获取天气数据，并实现天气数据的可视化。

1.2.4 R 语言医疗大数据分析

医疗大数据是相对于一般数据而言的，指的是人们从医疗系统的软件系统中捕捉大容量数据，通过大数据分析获得新的认知，从而创造新的价值来源。医疗大数据几乎包含公民所有个人信息，包括医疗、饮食、住所、旅行登记等。在临床操作方面，有 5 个主要场景的大数据应用。根据麦肯锡公司的估计，如果这些应用被充分采用，针对美国一个国家的医疗健康开支一年就将减少 165 亿美元。主要场景的大数据应用包括 5 个方面。

（1）比较效果研究。

通过全面分析病人特征数据和疗效数据，然后比较多种干预措施的有效性，可以找到针对特定病人的最佳治疗途径。世界各地的很多医疗机构，如德国 IQWIG（德国医疗质量和效率研究所）、加拿大普通药品检查机构等，已经开始了 CER（比较效果研究）项目，并取得了初步成功。