

IT部落窝论坛、中国会计视野论坛共同推荐

大数据时代，不懂数据分析，那你OUT了！

职场中，还不会用Excel进行数据分析，那你已经掉队了！

生活中，是不是也觉得数据分析与你没关系？那你错了！

实际上，我们的生活和工作都已经离不开数据分析了！

既然如此，那就让我们和书中的主人公一起学习数据分析吧！

36.4小时教学视频 + 120个操作实例 + 8个完整案例

拒绝枯燥乏味的说教，代之以轻松有趣的讲解

以“职场人物情景对话”的风格贯穿全书

用企业中的数据分析实战案例引导读者学习

门槛很低，几乎不介绍复杂的统计原理和公式！

数据说服力

—菜鸟学Excel数据分析（职场进阶版）

马军 等编著



清华大学出版社



数据说话 说服力

—菜鸟学Excel数据分析（职场进阶版）

马军 等编著



清华大学出版社

北京

内 容 简 介

本书以一位刚工作的大学生学习数据分析的过程为背景，向读者介绍数据分析的基本概念和实际操作。在具体分析操作时，没有使用复杂的专业分析软件，而是用大家在办公中都会用到的 Excel 来进行操作。本书将数据分析和 Excel 基础操作融入职场办公的场景中，主要用企业的销售数据来演示各类数据分析方法的使用，阅读起来轻松幽默，不再像大部分数据分析类图书那样枯燥乏味。另外，本书配套光盘中提供了专门录制的多媒体教学视频，便于读者高效而直观地学习。

全书共分 15 章，从数据分析的基本概念开始，不仅介绍了 Excel 在数据分析中的基本操作，还结合案例重点介绍了数据收集、数据加工、描述分析、趋势分析、对比分析、相关分析、结构分析、在 Excel 中展示分析数据以及编写分析报告等内容。为了让读者在数据分析时能进行更高层次的操作，书中还介绍了 Excel 的模拟分析功能及 VBA 在数据分析中的应用，最后还简单介绍了数据挖掘的相关概念。

本书适合各类企事业单位从事市场营销、金融、人力资源管理、财务管理等需要进行数据分析的职场人士阅读，也适合企业经营、管理人员和广大数据分析爱好者阅读。

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目（CIP）数据

数据说服力——菜鸟学 Excel 数据分析（职场进阶版）/马军等编著. —北京：清华大学出版社，2013
ISBN 978-7-302-33803-1

I. ①数… II. ①马… III. ①表处理软件 IV. ①TP391.13

中国版本图书馆 CIP 数据核字（2013）第 212285 号

责任编辑：夏兆彦

封面设计：欧振旭

责任校对：徐俊伟

责任印制：杨 艳

出版发行：清华大学出版社

网 址：<http://www.tup.com.cn>, <http://www.wqbook.com>

地 址：北京清华大学学研大厦 A 座 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈：010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者：清华大学印刷厂

装 订 者：三河市李旗庄少明印装厂

经 销：全国新华书店

开 本：185mm×260mm 印 张：22.75 字 数：570 千字

附光盘 1 张

版 次：2014 年 1 月第 1 版 印 次：2014 年 1 月第 1 次印刷

印 数：1~4000

定 价：55.00 元

产品编号：054567-01

前　　言

数据分析该怎么做？

数据分析有什么作用？

该用什么方法进行数据分析？

用什么工具进行数据分析？

这可能是很多人看到“数据分析”马上就会想起的一连串问题。很多人对“数据分析”的认识还停留在难以理解的统计知识、复杂的计算公式和难以操作的分析软件阶段。

其实，在计算机技术快速发展的今天，数据分析已经不再需要你去推算那些复杂的公式了。而且，对于工作中大部分的数据分析任务都只是进行一些简单的、容易理解的数据处理和分析操作而已。

为了使很多还不了解数据分析的人员能掌握基本的数据分析知识，本书以一位刚入职场的大学生学习数据分析的过程，向读者介绍数据分析的基本概念和实际操作。在具体分析操作时，没有使用复杂的专业分析软件 SPSS 或 SAS，而是选择大家在办公中都会用到的 Excel 来进行相关的数据收集、整理、分析和生成图表等操作。这样，就可让所有办公人员都能用本书介绍的方法进行数据分析操作了。

关于本书

本书中模拟了职场中老人带新人学习数据分析的情景。在正式阅读本书之前，有必要对书中涉及的人物角色及本书主要涉及的内容做个简单介绍。

李双双：女，应届大学毕业生，刚应聘到 WM 电器销售公司做总经理陈总的秘书。

沈栋：男，WM 电器销售公司信息部的数据分析师。

主角李双双大学毕业后应聘到 WM 电器销售公司，当总经理秘书，为总经理的日常工作做好辅助工作。由于大部分文字工作由行政部负责，总经理秘书需要做的主要工作是收集分公司经营数据，并在总经理需要时提供相关数据。

对于刚毕业的李双双，没有实际工作经验，更别提数据分析，因此，李双双找到分公司的同事——信息部的沈栋，向其学习数据分析的基础知识。

全书就在李双双向沈栋请教学习的过程中展开，沈栋一共分为 15 次课程，循序渐进地向李双双传授数据分析方面的知识。

在第 1 次的学习中，沈栋向李双双讲解了一些数据分析的入门知识。包括如何从海量数据中发现问题，数据分析的前世今生，怎么进行数据分析，常用分析模型，常用数据分析方法，以及不能不知的统计概念等内容。通过第一天的学习，李双双对数据分析有了一

个大概的认识，不再觉得很神秘。

第 2 次学习时，沈栋从数据分析的武器装备角度向李双双介绍了常见的分析软件 SPSS、SAS 和马克威分析系统。不过，重点还是介绍 Excel 2013 的使用方法，由于李双双有一点 Excel 操作基础，因此对于沈栋介绍的内容学得很快。这次课学得很轻松。

第 3 次学习时，沈栋就开始介绍数据分析流程中的第一环：数据收集。主要介绍了数据收集方法、调查问卷设计，以及如何将这些数据输入 Excel 表格中。

第 4 次学习时，沈栋介绍了怎样在 Excel 中对收集的数据进行加工。包括对数据的审核，对重复值、缺失值、离群值的数据的处理，数据排序、分组，以及通过 Excel 的数据透视表处理数据的方法和技巧。

第 5 次学习时，沈栋以案例的形式介绍数据分析的实际操作，对公司的工资数据进行分析，介绍了描述分析的概念，并制作出相应的分析图表。

第 6 次学习时，李双双通过对销售数据的分析，学习并掌握了趋势分析的相关内容，并对公司的销售数据进行了同比、环比分析，计算移动平均数据，制作各类分析图表。

第 7 次学习时，李双双通过收集、整理数据，分析了公司在连锁零售行业中的地位，学会了对比分析的相关概念。

第 8 次学习时，李双双学习了相关分析的概念，并使用相关分析软件分析了广告投入是否影响销量。

第 9 次学习时，沈栋介绍了什么是结构分析法，然后使用结构分析法进行了 GDP 结构指标的变动趋势分析，还对公司销售额在行业中所占的份额、公司各品类商品对利润的贡献进行了分析。

第 10 次学习时，沈栋介绍了在 Excel 中展示分析数据的相关技巧，包括 Excel 2013 中新增的一些突出显示数据的技巧，图表选择技巧，以及图表美化技巧等。

第 11 次学习时，沈枕介绍了如何编写分析报告，并让人看懂你的分析。主要从数据分析报告的作用和结构方面进行介绍，最后还给李双双演示了一个简单的分析报告。

第 12 次学习时，沈栋着重介绍了 Excel 的其他常用分析功能，包括单变量求解、模拟运算表、方案管理以及规划求解等。通过 Excel 的这些工具，可为数据分析提供更多的手段。

第 13 次学习时，沈栋介绍了 Excel 的高级功能——VBA 程序设计，包括 VBA 的基本概念、程序结构控制以及 Excel 对象的访问等内容。

第 14 次学习时，沈栋从用 VBA 扩展数据分析的角度演示了 VBA 程序的强大功能，主要演示了在数据输入和整理这两个阶段用 VBA 进行功能的扩展。

第 15 次学习时，沈栋向李双双介绍了数据挖掘的基础知识，了解了什么是数据挖掘，数据挖掘的功能和流程以及数据挖掘的应用，并简单介绍了几款常见的数据挖掘软件。

本书特色

1. 化繁为简

本书不介绍数据分析中复杂的统计原理和数学公式，只介绍数据分析必须了解的概念，以解决实际问题为第一要务。

2. 案例实用

本书中的案例都是以一个销售公司的经营数据为例进行介绍，分析的案例也是企业经常要用到的，读者可将这些方法直接应用到工作中。

3. 学习轻松

本书以轻松的写作风格进行编写，书中主人公之间口语化的对白，让你阅读不再枯燥乏味，感觉就像看小说一样就将相关知识学到了。

4. 视频教学

为了帮助读者高效、直观地学习本书内容，本书重点内容都专门录制了配套的多媒体教学视频，这些视频和书中涉及的其他资料一起收录于本书配书光盘中。

适用读者

阅读本书的读者不需要具有数据分析和 Excel 操作的经验，全书零起点，可适用以下各类人员阅读：

- 市场营销人员；
- 财务管理人员；
- 人力资源管理人员；
- 金融从业人员；
- 数据统计和分析人员；
- 即将踏入职场的大中专院校学生；
- 需要提升自身职场竞争力的职场新人；
- 常阅读经营分析、市场调研报告的管理人员。

本书作者

本书由洛阳理工学院的马军老师主持编写。其他参与编写的人员有陈晓建、陈振东、程凯、池建、崔久、崔莎、邓凤霞、邓伟杰、董建中、耿璐、韩红轲、胡超、黄格力、黄缙华、姜晓丽、李学军、刘娣、刘刚、刘宁、刘艳梅、刘志刚、司其军、滕川、王连心、沃怀凯、闫玉宝。

您在阅读本书的过程中若有疑问，请发 E-mail 和我们联系。E-mail 地址：bookservice2008@163.com。

编者

目 录

| | |
|-----------------------------------|-----------|
| 第1章 让数据说话 | 1 |
| 1.1 大海捞针，如何从海量数据中发现问题 | 1 |
| 1.1.1 找不着信息的尴尬 | 2 |
| 1.1.2 如何处理海量数据 | 3 |
| 1.2 数据分析的前世今生 | 3 |
| 1.2.1 数据的记录 | 4 |
| 1.2.2 数据分析的目的 | 4 |
| 1.2.3 数据分析的现状 | 5 |
| 1.2.4 数据分析的前景 | 7 |
| 1.2.5 数据分析不是浮云 | 8 |
| 1.3 怎么进行数据分析 | 8 |
| 1.3.1 明确数据分析目的 | 9 |
| 1.3.2 数据搜集 | 10 |
| 1.3.3 数据处理 | 10 |
| 1.3.4 数据分析 | 10 |
| 1.3.5 数据分析报告 | 11 |
| 1.4 常用分析模型 | 12 |
| 1.4.1 波特五力分析模型 | 13 |
| 1.4.2 SWOT 分析模型 | 14 |
| 1.4.3 5W2H 分析模型 | 14 |
| 1.4.4 PEST 分析模型 | 15 |
| 1.4.5 SCP 分析模型 | 16 |
| 1.5 常用数据分析方法 | 16 |
| 1.5.1 对比分析法 | 16 |
| 1.5.2 分组分析法 | 19 |
| 1.5.3 结构分析法 | 20 |
| 1.5.4 高级数据分析法 | 21 |
| 1.6 不能不知的统计概念 | 23 |
| 1.6.1 总体与样本 | 23 |
| 1.6.2 平均数、中位数和众数 | 24 |
| 1.6.3 番与倍的概念 | 26 |
| 1.6.4 百分数与百分点的区别 | 26 |
| 第2章 数据分析的武器装备——Excel | 28 |
| 2.1 不可不知的数据分析软件 | 28 |
| 2.1.1 SPSS 简介 | 28 |

| | |
|-----------------------------|-----------|
| 2.1.2 SAS 简介 | 29 |
| 2.1.3 马克威分析系统简介 | 31 |
| 2.2 不能不会的 Excel | 32 |
| 2.2.1 认识 Excel 2013 | 32 |
| 2.2.2 Excel 2013 的功能区 | 34 |
| 2.2.3 功能区的使用 | 35 |
| 2.3 Excel 的基本操作 | 38 |
| 2.3.1 必须了解的概念 | 38 |
| 2.3.2 选择操作区域 | 39 |
| 2.3.3 操作单元格 | 41 |
| 2.3.4 行列操作 | 43 |
| 2.4 美化表格 | 45 |
| 2.4.1 让文字更好看 | 45 |
| 2.4.2 让文字排好队 | 47 |
| 2.4.3 不能只使用白背景 | 48 |
| 2.4.4 边框必须有 | 50 |
| 2.4.5 设置数据格式 | 53 |
| 2.5 Excel 的公式与函数 | 56 |
| 2.5.1 基本概念 | 56 |
| 2.5.2 运算符 | 57 |
| 2.5.3 运算符优先级 | 59 |
| 2.5.4 快速计算方法 | 60 |
| 2.5.5 输入公式 | 61 |
| 2.5.6 函数 | 63 |
| 第 3 章 数据从哪里来 | 66 |
| 3.1 数据收集方法 | 66 |
| 3.1.1 间接来源 | 66 |
| 3.1.2 直接来源 | 67 |
| 3.1.3 普查还是抽查 | 68 |
| 3.2 调查问卷设计 | 71 |
| 3.2.1 调查问卷结构 | 71 |
| 3.2.2 调查问卷设计步骤 | 72 |
| 3.2.3 问题设计的注意事项 | 74 |
| 3.2.4 用 Excel 设计调查问卷 | 75 |
| 3.3 Excel 中的数据 | 83 |
| 3.3.1 数据类型 | 83 |
| 3.3.2 字段与记录 | 86 |
| 3.3.3 Excel 表 | 86 |
| 3.4 在 Excel 中输入数据 | 91 |
| 3.4.1 各类数据的输入 | 91 |
| 3.4.2 快速输入数据的方法：填充 | 94 |
| 3.5 在 Excel 中导入数据 | 98 |
| 3.5.1 导入文本数据 | 98 |
| 3.5.2 导入网站数据 | 103 |

| | |
|-----------------------------------|------------|
| 3.5.3 导入数据库数据..... | 104 |
| 第 4 章 数据加工 | 107 |
| 4.1 你的数据来源可靠么..... | 107 |
| 4.1.1 对数据进行审核..... | 107 |
| 4.1.2 审核直接数据..... | 108 |
| 4.1.3 审核间接数据..... | 108 |
| 4.2 你的数据完整么..... | 109 |
| 4.2.1 处理重复数据..... | 109 |
| 4.2.2 处理缺失数据..... | 116 |
| 4.2.3 处理离群值..... | 121 |
| 4.3 你的数据有序么..... | 122 |
| 4.3.1 排序找出趋势..... | 122 |
| 4.3.2 用筛选剔除数据..... | 129 |
| 4.4 你是哪一组 | 133 |
| 4.4.1 按品质标志分组..... | 134 |
| 4.4.2 按数量标志分组..... | 134 |
| 4.4.3 如何快速统计各分组..... | 135 |
| 4.4.4 分类汇总..... | 142 |
| 4.5 Excel 神器：数据透视表 | 143 |
| 4.5.1 创建数据透视表..... | 144 |
| 4.5.2 透视表的计算和汇总..... | 145 |
| 4.5.3 透视表筛选和排序..... | 146 |
| 第 5 章 你的工资被平均了吗——描述分析..... | 148 |
| 5.1 收集整理工资数据..... | 148 |
| 5.1.1 要收集哪些数据..... | 148 |
| 5.1.2 整理计算数据..... | 149 |
| 5.1.3 快速找出错误数据..... | 150 |
| 5.2 用数据透视表汇总工资数据..... | 154 |
| 5.2.1 查表合并数据..... | 154 |
| 5.2.2 更快的方法有没有..... | 156 |
| 5.2.3 向功能区增加新按钮..... | 157 |
| 5.2.4 数据透视表现身..... | 158 |
| 5.2.5 让数据透视表为我所用..... | 160 |
| 5.3 什么是描述分析..... | 164 |
| 5.3.1 集中趋势..... | 164 |
| 5.3.2 离散趋势..... | 166 |
| 5.3.3 偏度与峰度..... | 167 |
| 5.3.4 对工资数据进行描述分析..... | 168 |
| 5.4 工资数据分析图表..... | 170 |
| 5.4.1 总额和人均变化图..... | 170 |
| 5.4.2 部门工资分布示意图..... | 172 |
| 5.4.3 工资与部门的关系..... | 174 |
| 5.4.4 工资与职称的关系..... | 175 |

| | |
|-------------------------------------|------------|
| 5.4.5 工资与学历的关系..... | 175 |
| 5.5 发现公司薪酬体系问题..... | 175 |
| 第 6 章 了解过去，展望未来——趋势分析..... | 177 |
| 6.1 销售数据采集 | 177 |
| 6.1.1 选择数据的思路..... | 177 |
| 6.1.2 整理数据..... | 178 |
| 6.2 什么是趋势分析..... | 182 |
| 6.2.1 什么是时间序列..... | 182 |
| 6.2.2 什么是线性趋势..... | 183 |
| 6.3 公司销售增长情况分析..... | 184 |
| 6.3.1 再说同比和环比..... | 184 |
| 6.3.2 计算销售数据同比增速..... | 185 |
| 6.3.3 计算销售数据环比增速..... | 186 |
| 6.3.4 计算发展速度与增长速度..... | 188 |
| 6.4 用移动平均过滤波动..... | 191 |
| 6.4.1 怎么计算移动平均..... | 192 |
| 6.4.2 销售数据移动平均数..... | 193 |
| 6.5 制作分析图表 | 196 |
| 6.5.1 月销售金额图表..... | 196 |
| 6.5.2 销售金额折线图..... | 197 |
| 第 7 章 公司在行业中的地位——对比分析..... | 200 |
| 7.1 收集行业数据 | 200 |
| 7.1.1 选择数据的思路..... | 200 |
| 7.1.2 收集内部数据..... | 201 |
| 7.1.3 收集外部数据..... | 201 |
| 7.2 什么是对比分析..... | 204 |
| 7.2.1 横向比较和纵向比较..... | 204 |
| 7.2.2 对比的对象要有可比性..... | 205 |
| 7.2.3 对比的口径要一致..... | 206 |
| 7.2.4 对比的指标类型要统一..... | 207 |
| 7.3 公司能进多少强？ | 207 |
| 7.3.1 按销售总额对比..... | 207 |
| 7.3.2 单店销售额对比..... | 210 |
| 第 8 章 隐藏在数据后面的内容——相关分析 | 214 |
| 8.1 什么的相关分析..... | 214 |
| 8.1.1 相关关系有哪些分类..... | 215 |
| 8.1.2 相关表和相关图..... | 216 |
| 8.1.3 相关系数..... | 217 |
| 8.2 广告投入对销量的影响..... | 220 |
| 8.2.1 收集整理广告投入与销量数据..... | 220 |
| 8.2.2 通过折线图查看广告对销量的影响..... | 222 |
| 8.2.3 通过散点图查看广告对销量的影响..... | 224 |

| | |
|------------------------------------|------------|
| 8.2.4 广告与销量的相关系数..... | 225 |
| 第 9 章 公司能分多少蛋糕——结构分析..... | 227 |
| 9.1 什么是结构分析法..... | 227 |
| 9.1.1 结构分析法的概念..... | 227 |
| 9.1.2 结构指标..... | 228 |
| 9.1.3 结构分析的作用..... | 229 |
| 9.2 GDP 结构指标的变动趋势分析..... | 229 |
| 9.2.1 从国家统计局网站获取数据..... | 229 |
| 9.2.2 通过结构指标分析变动趋势..... | 232 |
| 9.3 份额结构分析 | 235 |
| 9.3.1 收集行业、地区数据..... | 235 |
| 9.3.2 计算公司所占份额..... | 235 |
| 9.4 各品类商品的贡献分析..... | 238 |
| 9.4.1 收集各类商品销售额、成本数据..... | 239 |
| 9.4.2 计算各类商品的贡献..... | 239 |
| 第 10 章 在 Excel 中展示分析数据..... | 241 |
| 10.1 突出表格中的关键数据..... | 241 |
| 10.1.1 通过底纹突出数据 | 241 |
| 10.1.2 用条件格式突出数据 | 243 |
| 10.1.3 用数据条区分数据 | 246 |
| 10.1.4 用色阶区分数据 | 246 |
| 10.1.5 用图标集区分数据 | 247 |
| 10.2 重新认识 Excel 的图表 | 249 |
| 10.2.1 图表的组成 | 249 |
| 10.2.2 如何选择 Excel 图表 | 251 |
| 10.2.3 根据数据关系选择图表 | 252 |
| 10.2.4 修改图表设计 | 254 |
| 10.3 给图表美容 | 256 |
| 10.3.1 选择图表中的对象 | 257 |
| 10.3.2 向图表中插入形状 | 257 |
| 10.3.3 改变形状样式 | 258 |
| 10.3.4 新的设置方式 | 260 |
| 第 11 章 让人看懂你的分析 | 263 |
| 11.1 数据分析报告有啥用 | 263 |
| 11.1.1 为什么要写数据分析报告 | 263 |
| 11.1.2 数据分析报告的分类 | 265 |
| 11.1.3 数据分析报告编写流程 | 266 |
| 11.1.4 数据分析报告的原则 | 267 |
| 11.2 编写数据分析报告 | 268 |
| 11.2.1 起一个好的标题 | 268 |
| 11.2.2 报告前言 | 269 |
| 11.2.3 报告正文 | 270 |

| | |
|-------------------------------------|------------|
| 11.2.4 报告结尾，给出结论、建议或展望 | 271 |
| 11.3 一个简单的分析报告 | 272 |
| 第 12 章 Excel 的其他常用分析功能 | 277 |
| 12.1 单变量求解 | 277 |
| 12.2 模拟运算表 | 279 |
| 12.2.1 单变量模拟运算表 | 279 |
| 12.2.2 双变量模拟运算表 | 281 |
| 12.3 方案管理器 | 283 |
| 12.3.1 创建方案前要准备什么 | 283 |
| 12.3.2 创建方案 | 286 |
| 12.3.3 查看方案的结果 | 288 |
| 12.3.4 生成方案预测报表 | 289 |
| 12.4 用 Excel 进行运筹规划 | 290 |
| 12.4.1 规划求解概述 | 290 |
| 12.4.2 规划求解在哪里 | 291 |
| 12.4.3 规划求解的操作步骤 | 293 |
| 12.4.4 计算最大利润 | 297 |
| 第 13 章 Excel 的高级功能 | 301 |
| 13.1 了解宏和 VBA | 301 |
| 13.1.1 什么是宏 | 301 |
| 13.1.2 怎么调用宏 | 303 |
| 13.1.3 什么是 VBA | 304 |
| 13.2 VBA 编程基础 | 305 |
| 13.2.1 书写规范 | 306 |
| 13.2.2 VBA 基本要素 | 307 |
| 13.2.3 控制程序流程 | 310 |
| 13.2.4 VBA 过程 | 316 |
| 13.3 VBA 中的 Excel 对象 | 318 |
| 13.3.1 什么是对象 | 319 |
| 13.3.2 Application 对象 | 321 |
| 13.3.3 Workbook 对象 | 322 |
| 13.3.4 Worksheet 对象 | 322 |
| 13.3.5 Range 对象 | 324 |
| 第 14 章 用 VBA 扩展数据分析 | 326 |
| 14.1 用 VBA 代码扩展数据输入 | 326 |
| 14.1.1 删不掉的公式 | 326 |
| 14.1.2 选不中的单元格 | 328 |
| 14.1.3 可控的数据验证设置 | 329 |
| 14.1.4 禁止输入重复值 | 332 |
| 14.2 用 VBA 代码扩展数据整理 | 333 |
| 14.2.1 相同数据标识相同颜色 | 334 |
| 14.2.2 删除多个空行 | 338 |

| | |
|-----------------------------|------------|
| 14.2.3 每人的最高日销售额 | 339 |
| 第 15 章 数据挖掘是什么 | 341 |
| 15.1 了解数据挖掘 | 341 |
| 15.1.1 什么是数据挖掘 | 341 |
| 15.1.2 数据分析与数据挖掘 | 342 |
| 15.2 数据挖掘的功能和流程 | 343 |
| 15.2.1 数据挖掘有什么功能 | 343 |
| 15.2.2 数据挖掘流程 | 344 |
| 15.3 数据挖掘的应用 | 345 |
| 15.3.1 金融行业应用数据挖掘 | 345 |
| 15.3.2 电信行业应用数据挖掘 | 346 |
| 15.3.3 网站应用数据挖掘 | 346 |
| 15.4 数据挖掘软件 | 347 |
| 15.4.1 数据挖掘软件的分类及选择 | 347 |
| 15.4.2 常用数据挖掘软件简介 | 348 |

第1章 让数据说话

数据分析，这个词对刚入职的李双双来说还很陌生。

对于刚毕业的李双双，没有实际工作经验，更别提数据分析了，可现在却必须面对这个难题。

李双双今年刚毕业，应聘到WM电器销售公司做总经理秘书。上班的第一天，陈总就对李双双的工作进行了安排：由于分公司的大部分文字工作由行政部负责，因此，总经理秘书需要做的主要工作是收集分公司各门店的经营数据，并在总经理需要时提供相关数据。

数据分析该怎么做？数据分析有什么作用？该用什么方法进行数据分析？用什么工具进行数据分析？……李双双带着这一连串的问号，开始了人生的第一份工作。

正所谓“车到山前必有路，柳暗花明又一村”，上班后，李双双知道公司设有一个信息部，专门从事数据分析工作，这个部门的员工基本上都是科班出身学习数据统计分析的。“对，找他们拜师学艺！”李双双心里又开始了女孩特有的幻想，“如果遇到一位帅哥就更好了！嘻嘻。”

想到就做，这是李双双的性格。第二天上班，她抽空到信息部办公室，一个小伙子从办公桌后站了起来：你好，请问你找谁？

李双双定睛一看，光洁白皙的脸庞，乌黑深邃的眼眸，那浓密的眉、高挺的鼻，无不表示这是一位帅哥。李双双脸红了一下：你好，我叫李双双，刚到公司上班。

帅哥：哦，欢迎加入，我叫沈栋，有什么可以帮忙的？

李双双嘴很甜，马上说道：沈哥，你好！是这样的，我现在做陈总的秘书。陈总要求我做一些数据收集、分析的工作，可我对数据分析还是门外汉，听说信息部都是科班出身的数据分析师，所以今天想来拜师学艺。

沈栋：呵呵，拜师就说不上了，有什么问题咱们可以探讨。

李双双：要说问题，我现在感觉是一头雾水，还找不到从哪里说起，想请你从零开始教我。不过，这不着急，沈老师，今天拜了师，中午该我请老师吃饭。

沈栋：不用客气，你是才进公司的新人，该我请客。

1.1 大海捞针，如何从海量数据中发现问题

搞定拜师的事情之后，李双双决定趁热打铁，下午到信息部办公室看了几次，准备瞅准沈栋空闲时马上就去找他上课。

看到李双双在办公室门口转了几次，沈栋把手上的工作忙完后说：小李，进来吧。我看你在门口转了几次了，学习的热情很高啊。

李双双顽皮地笑了笑：我现在是求知若渴啊。老师，什么是数据分析？

1.1.1 找不着信息的尴尬

沈栋：什么是数据分析？通俗地讲，就是对数据进行分析。

李双双：这么简单？

沈栋：其实，从专业角度来解释数据分析的话，有很长一段话，不过，我觉得不用去记那么多专业的描述，对你来说，关键是了解数据分析，掌握常用方法。

李双双：正好，我也不喜欢记那些苦涩难懂的术语。

沈栋：假设咱们公司每天只销售 10 笔商品，这样，就会产生 10 条销售数据，这会是什么情况？

李双双：只销售 10 笔商品，那咱们就失业了。

沈栋：哈哈，你想得挺远，不过确实是这样的。如果只有 10 笔销售数据，老板自己就可逐一检查处理了，要咱们干啥。但是，像现在这种情况，咱们分公司一年销售额几十亿，商品数据、销售数据、客户数据每年都达几千万、上亿笔，这么多数据，老板自己一个人就看不完了，并且面对这么多数据，也不是一眼就能看出规律的。这就需要专业的数据分析师对这些海量数据进行分析，也就有我的一口饭吃。

时至今日，许多企业在数据管理中所面临的主要挑战之一就是如何从海量数据中获得更多的价值，尤其是从企业自身一点一滴辛苦积累起来的数据中获取价值。

当公司的数据很多时，就会出现问题：不知道整个公司的信息系统中保存了多少数据，保存了什么数据，这些数据保存在哪里，各系统中的数据之间存在哪些冲突……

这就导致了很多企业都会出现找不着信息的尴尬（图 1-1 所示）。



图 1-1 找不着信息的尴尬

在信息高度发达的今天，数据信息呈爆炸式增长，每天全世界产生着巨量的数据，例如百度网站每天就有上亿人次访问，访问数据将达到几亿条。这还只是一个网站产生的数据，而互联网上运行的网站数不胜数。根据一项统计显示，全世界 2005 年创造的数据是 1500 亿 GB，到 2020 年创造的数据将超过 3500 万亿 GB。随着这股信息洪流不断增加，面对这种呈几何级数增长的数据，从中提取并存储有用的数据将变得十分困难，已经出现数据泛滥的状况。

1.1.2 如何处理海量数据

李双双：按你这么说，对于这么庞大的海量数据，该怎么处理呢？

沈栋：海量数据是发展趋势，随着电脑存储容量越来越大，数据量肯定也会不断地呈几何级增长。面对这种情况，如何从海量数据中提取有用信息就显得重要而紧迫了。如图 1-2 所示是海量数据处理的过程，在这个过程中，将收集到的海量数据通过数据处理过程，最后找到金子（需要的数据）。

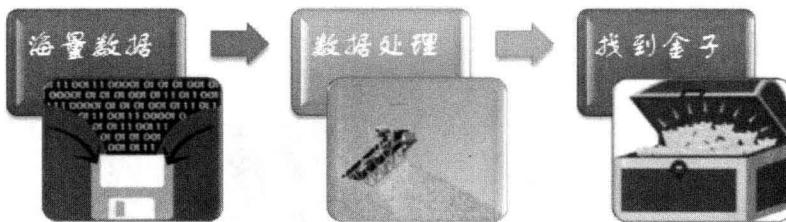


图 1-2 海量数据处理

对海量数据的处理，通常应做到以下几点：

- 处理要准确；
- 精度要高；
- 处理时间要短；
- 得到有价值信息要快。

对于像咱们公司这种有海量数据需要处理的公司，通常都需要投入巨资来建立信息技术系统，并成立相关的技术支持部门和信息分析部门。

对于海量数据的处理并没有通用的方法，需要数据分析师长期工作经验的积累，下面我简单介绍一些通用的原理和规则。

- 选用优秀的数据库工具：好的数据库系统可成倍提高数据处理的速度；
- 优良的程序：良好的程序代码应该包含好的算法，包含好的处理流程，包含好的效率，包含好的异常处理机制等；
- 对海量数据进行分区：对数据进行分区，将数据分散开，减小磁盘 I/O，减小了系统负荷，而且还可以将日志和索引等放于不同的分区中；
- 建立数据索引：对海量的数据处理，对大表建立索引是必需的；
- 建立缓存机制：根据数据量大小建立合适的缓存，提高处理效率。

1.2 数据分析的前世今生

李双双：老师啊，你讲的这些太复杂了，头痛 ing。

沈栋笑了：确实，你现在要理解海量数据的处理是有点难。现在给你说这些只是让你了解数据处理的重要性，不用去详细理解具体的内容，等你把数据分析知识掌握到一定程度后，就能理解了。现在先给你讲一下数据分析的前世今生。

1.2.1 数据的记录

沈栋端起茶杯喝了一口，接着说：要说数据分析，首先得说数据。在日常生活中，人们为了回想过去发生过的事，对进行的各种活动进行记录，就形成了数据。人类记录数据通常具有以下几种方式（如图 1-3 所示）：

- 口头记录；
- 图画记录；
- 文字记录；
- 音频记录；
- 视频记录。



图 1-3 数据的记录

随着人类科技的进步，数据的记录形式也在不断变化，从最初的口口相传发展到现在，已经能用各种介质来保存信息，例如文字、图片、声音和视频等。到现在电子数据存储的软、硬件技术得到很大的发展，已经可以将文本、图片、声音和视频等信息存储在一起，并能方便地检索这些信息。

李双双马上接口说：我知道了，像咱们平常收集的销售数据就是文本类的数据，而广告公司为咱们设计的海报就是图片类的数据，而声音、视频数据就是看电视、电影等。

沈栋竖起了大拇指：聪明，学会总结了。

李双双顽皮地一笑：一般一般，世界第三。

1.2.2 数据分析的目的

沈栋接着说：面对这么多种类、海量的数据，咱们该怎么办呢？

李双双：就需要数据分析？