

KPMG

# FROM DATA TO INSIGHTS

毕马威亚太区及中国主席陶匡淳先生

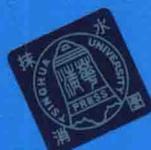
倾情作序推荐

## 洞见数据价值 大数据挖掘要案纪实

毕马威中国大数据团队◎著

企业与个人大数据简明读本

清华大学出版社



FROM DATA

# TO INSIGHTS

洞见数据价值 大数据挖掘要素纪实

大数据要素

## 洞见数据价值 大数据挖掘要素纪实

陈旭东 著

清华大学出版社

清华大学出版社



FROM DATA  
TO INSIGHTS

# 洞见数据价值

大数据挖掘要案纪实

毕马威中国大数据团队◎著

清华大学出版社  
北京

## 内容简介

本书是国际知名咨询公司毕马威的大数据团队的集大成之作，内容丰富，观点新颖，贴近大众生活、工作与学习实际场景，内容不拘泥于技术阐释，通过寓教于乐的方式，以丰富翔实的案例来解析大数据挖掘，尽量通过常见的场景来阐述数据的价值与意义。

第1章介绍大数据在银行业、征信业、审计、传统制造业、互联网行业、舆情监控、影视业、环保产业以及体育产业等多个领域的应用方案和前景。第2章重点介绍大数据分析在商业工作和营销推广中的作用。第3章介绍大数据挖掘过程中涉及的数据的前期准备工作，重点介绍数据准备工作的要点和诀窍。第4章结合业务、生活、娱乐，寓教于乐，介绍大数据的实际应用方式。附录部分介绍一位数据工作者的成长之路，向感兴趣的读者介绍从事数据工作应该具备的素质和掌握的技能。

本书可作为企业管理人员、营销主管、分析人员、IT人员等理解大数据、应用大数据为企业创造价值的指引，同时，本书也可供统计学、应用数学及计算机专业学者和研究人员参考学习。

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

洞见数据价值：大数据挖掘要案纪实 / 毕马威中国大数据团队著. —北京：清华大学出版社，2018

ISBN 978-7-302-49180-4

I. ①洞… II. ①毕… III. ①数据采集—研究 IV. ①TP274

中国版本图书馆CIP数据核字(2017)第330716号

责任编辑：秦健  
封面设计：李召霞  
责任校对：胡伟民  
责任印制：沈露

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](mailto:c-service@tup.tsinghua.edu.cn)

质量反馈：010-62772015，[zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 刷 者：北京亿浓世纪彩色印刷有限公司

装 订 者：三河市金元印装有限公司

经 销：全国新华书店

开 本：148mm×210mm 印 张：9.75 字 数：245千字

版 次：2018年2月第1版 印 次：2018年2月第1次印刷

印 数：1~4000

定 价：69.00元

---

## 谢邦昌教授

台北医学大学管理学院及大数据研究中心院长/主任

多年来我一直关注及致力于大数据在各领域的应用及其产生的巨大作用，而毕马威中国大数据团队的著作正是一本绝佳的案例集，它告诉我们数据挖掘如何为行业的最前沿带来价值，面向高维数据，从需求出发，使用丰富的方法，解决多样化的问题，可以看出作者深厚的理论功底和实践经验。一定可以让读者获益匪浅。

---

## 车品觉

红杉资本中国专家合伙人

大数据的概念从提出到落地经历了七年时间，作为分析、算法、人工智能的“新石油”，给国家和企业都带来了前所未有的机遇。我一直认为，每一种新技术的实践，最难的环节永远在于固有思想的改变，因此案例变得格外重要，好的案例起着启蒙作用。毕马威中国大数据团队的这本书正好为国内企业填补了包括互联网、金融、媒体、制造业等行业中“存管用”的实例。2018

年作为大数据下半场的开端，我相信大数据将会不断被各种创新终端所重新定义及重构。有趣的是，这种新能源还是在壮大中的婴儿而已。

---

## 张波教授

中国人民大学统计学学院副院长

在人大统计学学院执教多年，我希望高校始终保持着对商界、业界实际应用需求的敏感，因为只有真实的数据和真实的场景中，理论才能以最快的速度驱动、革新。如今，数据科学的快速发展特别是大数据的热潮正在改变着统计学发展的方向，由我的学生带领的毕马威中国大数据团队的这本《洞见数据价值：大数据挖掘要案纪实》，既有学术的理论与严谨，也有业界的经验与热情。我愿意把这本书推荐给有志于探索数据价值的高校师生和希望用数据解决实际问题的各界精英。

---

## 刘政博士

SAS 中国研发中心总经理

这些取自“KPMG 大数据挖掘”微信公众号的文章是由经验丰富的一线咨询团队的成员奉献的、理论联系实际范文。内容概括面广，针对性强，具有前瞻性，也有很强的实操性。很多文章都是实际经验的积累和创新，对于目前行业中缺少实际经验的数据分析人员来说，具有很好的指导和借鉴作用。

---

## 刘贤荣博士

中国建设银行数据管理部副总经理

数据已经成为移动互联时代企业的关键生产要素，无论是传统的制造行业还是新兴的生产性服务业，企业的设计、生产、营销、管理各个方面正与大数据高度融合，尤其是数据密度极高的银行业，大数据能力已经成为竞争的核心能力。古语云，“绝知此事要躬行”，毕马威中国大数据团队依托丰富的项目实战经验，将数据挖掘的专业知识用轻快而不失专业的语言和生动的案例呈现，在故事般愉快的阅读中深刻理解大数据的实际案例操作，读来是一种享受。

---

## 施奕明

前海征信副总经理，SAS 中文论坛创始人

第一次看到这个公众号，觉得它也就是繁星中的一颗，人来人往，沧海一粟，并没有太在意。三个月后，我惊讶地发现在并没有推手和炒作的情况下，身边的人都在转发和推荐这个公众号。于是我仔细阅读了每篇文章，确实能够感受到这群小伙伴们的情、专业和坚持。正是在这种正能量的激励下，我也模仿和学习他们，创建了我们自己类似的公众号，于是在大数据挖掘的前进道路上不再孤单。从某种意义上，这些可爱的小伙伴是我的老师。强烈推荐大家订阅他们的文章，在大数据挖掘内容服务领域他们绝对是一流的先锋，我就是他们的忠实粉丝。衷心希望我们能一

起品读他们的这份不沾染任何商业气息的热情、专业与坚持，在当今这个浮躁的社会中这份珍贵显得如此与众不同。

---

## 张磊博士

前 SAS 软件（中国）有限公司首席咨询顾问

从数据到价值，是数据挖掘最令人着迷之处，而其间却要历经种种的曲折反复。毕马威中国大数据团队将理论和实际紧密联系起来，带你领略数据挖掘旅途中的一路风景……

---

## 常国珍博士

中银消费金融大数据与 AI 实验室主任

2013 年维克托的《大数据时代》风靡全国，由此带来的大数据喧嚣甚至超过了楼市。曾几何时，没有数据科学家光环的大数据工作者慎言笃行，透过数据洞察商业先机，默默为企业创造价值。本人有幸在大数据元年前与本书部分作者相识、交流，更有幸成为本书的早期读者。本书以“纪实”的形式再现了鲜活的大数据案例，既有知识，又有典故，让作为读者的我收获颇丰。

---

## 蒋顺利

北京仁科互动网络技术有限公司市场副總裁

我可能是“KPMG 大数据挖掘”微信公众号最早的粉丝之一。从 2015 年 10 月开始，几乎每篇文章都会仔细阅读。它里

面的文章内容丰富，涵盖的范围广，包括银行业、制造业、互联网、影视等不同的行业，以及营销、风险、反欺诈、分析等不同的应用场景。虽包罗万象，但万变不离其宗，都紧密围绕着“大数据挖掘”这个热门的主题。虽是技术型微信公众号，但不故作高深，文字生动有趣，可读性强。无论是回顾历史还是预测未来，与大数据自然结合，娓娓道来，给人启发。欣闻微信公众号上的文章要结集成书，可喜可贺，也是众多读者的福音。无论是对大数据感兴趣的初学者，还是资深的数据科学家，这都是一本非常难得的大数据读物。

越是基本的理念，往往越能成为时代的标志、价值的度量，数据正是如此。

2020年，全球数据总量预计将超过44千万亿兆字节，数据之大，正如我们日常呼吸的空气一般不可缺少。这意味着世界上的一切都在产生数据，一切都在使用数据；万事万物都可以数据化，而数据也成为价值的新载体。大数据时代——也就是说，这样一个以数据衡量价值的时代，我们已经置身其中。

今日，“大数据”已不仅仅是一个新兴概念。国家、机构、企业都拥有数据，也急需使用数据，而物联网等理念的兴起，云计算等技术的应用，让我们能够帮助他们探索、挖掘、利用数据中的价值，能够存储、流通、关联、交换、使用大数据，开发每一个环节中的能量。大数据是资源，是一座亟待开掘的金矿。

毫无疑问，大数据拥有巨大的商业潜力和创造力，而这样的潜力和创造力，当然也属于毕马威中国大数据团队。

本书正是毕马威中国大数据团队的心血之作，汇聚了其微信公众号中的精品文章。该大数据团队建立几年以来，我一直非常欣赏团队的业务能力和业绩表现，同样，文如其人，他们在这本书中展现出的实力、动力和潜力也令人赞叹。

学术成果一向是新兴行业的驱动力，而业务能力是学术与实际工作的衔接点：大数据团队是一个集合了这两方面高端人才的

团队，本书中的作品，专注学术、扎根业务，也体现了不俗的行业发展眼光。何况，写这样的文章，一篇或许容易，但能在繁忙的日常工作之余坚持一年之久则殊为不易；如果不是对数据领域有深切的热爱，又怎么会有这样的动力呢？这本著作见证了团队的成长，更是团队潜力的明证。我相信，本书的读者们也一定会有同样的感觉。

大数据时代，谁掌握了数据，谁就掌握了价值，也掌握了世界的节奏。面对数据之潮，毕马威中国一直力求一马当先，而大数据团队正是公司内的行业先驱。“大鹏一日同风起，抟摇直上九万里”，我不由得心生期待，想要看看这样一个出色的团队还能给我和读者们带来怎样的惊喜。

**陶匡淳**

毕马威亚太区及中国主席

## 01

## 第 1 章

## 大数据在各行各业 // 001

## 1.1 什么是大数据? // 002

## 1.1.1 非常流行的大数据概念 // 002

## 1.1.2 不那么流行的大数据概念 // 006

## 1.1.3 也许会带给你灵感的大数据概念 // 007

## 1.2 大数据在银行业 // 029

## 1.2.1 业界展望: 大数据, 银行业未来的核心动力 // 029

## 1.2.2 创新方向: 大数据助力银行网点实现转型 // 036

## 1.3 大数据在征信业 // 041

1.3.1 业界展望: FICO 评分与芝麻信用, 传统征信向大数据  
征信的转变 // 0421.3.2 创新方向一: 从拒绝推断看个人征信业的  
大有可为 // 046

## 1.3.3 创新方向二: 论大中型客户数字化授信的可行性 // 054

## 1.4 大数据在审计业 // 057

1.4.1 业界展望：大数据分析如何支撑审计工作 // 057

1.4.2 创新方向：大数据能否代替传统审计？ // 061

## 1.5 大数据在传统制造业 // 065

业界展望：数字化企业进阶指南 // 066

## 1.6 大数据在互联网行业 // 069

创新方向：从滴滴收购优步看垄断企业的马太效应 // 069

## 1.7 大数据在舆情行业 // 076

创新方向：数据分析帮你掌握话语权 // 077

## 1.8 大数据在汽车行业 // 086

业界展望：征服汽车后市场，大数据与你同行 // 086

## 1.9 大数据在影视业 // 089

创新方向：星期几上映的电影最具有票房号召力？ // 090

## 1.10 大数据在环保产业 // 098

创新方向：北京治霾，能为你做些什么？ // 098

## 1.11 大数据在体育产业 // 104

创新方向：欧洲杯，跟着西班牙队学数据挖掘！ // 105

## 02

## 第 2 章

## 大数据在商业领域的应用 // 111

- 2.1 推荐算法在传统销售渠道中的应用模式 // 112
- 2.2 巧用运筹优化，提升整合营销管理水平 // 116
- 2.3 关联规则的应用 // 121
  - 2.3.1 小谈关联规则 // 121
  - 2.3.2 购物篮分析：绝不只是“啤酒与尿布” // 124
  - 2.3.3 创新方向：靠关联规则重获新生的东北小馆 // 128
- 2.4 智能荐食模型：大数据告诉你今天吃什么 // 133
  - 智能荐食模型 // 134
- 2.5 顾客时空模型：其实天下没有免费的 WiFi // 138
  - 2.5.1 无处不在的免费 WiFi // 138
  - 2.5.2 顾客时空模型 // 139
  - 2.5.3 进一步挖掘 // 141
- 2.6 社会网络分析法，助力信贷反欺诈 // 142
- 2.7 数据可视化利器：SAS Visual Analytics // 145
  - 2.7.1 为什么需要数据可视化？ // 145
  - 2.7.2 数据可视化的几个常见例子 // 146
- 2.8 文本挖掘，帮你识别网购评论是真是假 // 151

2.9 路径优化：如何改良快递送货路线？ // 156

## 03

### 第 3 章

## 数据前期准备 // 171

3.1 从抗日武装的发展谈到数据治理 // 172

3.1.1 数据质量问题 // 173

3.1.2 数据应用问题 // 174

3.1.3 实施策略和路径问题 // 176

3.2 如何生成你需要的基础数据？ // 182

3.3 如何利用数据仓库优化数据分析？ // 189

3.4 二分类变量的数据缺失插补 // 193

3.5 数据离散化，如何避免丢失信息？ // 201

3.6 如何避免数据离散化影响自变量的重要性？ // 204

3.7 二分类模型中如何应对分类自变量取值过多？ // 206

## 04

### 第 4 章

## 技术案例 // 211

4.1 建模变量太多怎么办？ // 212

4.2 信用评级模型怎么评估？ // 215

4.3 观察窗口怎么选？ // 219

- 4.4 K 折交叉验证怎么做? // 223
- 4.5 如何衡量变量之间的相关性? // 230
- 4.6 决策树算法真的越复杂越好吗? // 235
- 4.7 如何精选分类模型指标? // 245
- 4.8 当数据分析遇上超级奶爸 // 250
- 4.9 深度挖掘, 你的工资拖后腿了吗? // 253
- 4.10 用分位数回归看你的工资水平 // 258

## APP

### 附录 A

## 一位数据工作者的成长之路 // 265

- A.1 数据分析师入门攻略 // 266
- A.2 如何做一名“称职”的数据专家? // 269
- A.3 一个数据仓库转型者眼中的数据挖掘 // 271
- A.4 预测科学: 三点经验谈实际应用 // 276
- A.5 数据模型多了, 应该怎么管? // 277
- A.6 手握数据挖掘模型, 你一定要知道怎么用 // 281
- A.7 浅谈以史为鉴与数据分析 // 286

## 后记 // 297

01

第1章

---

# 大数据在各行各业

---

FROM

DATA

TO

INSIGHTS

---