

“十三五”

国家重点图书

大数据科学丛书

# 数据质量征途

[美] Yang W. Lee

Leo L. Pipino

James D. Funk

Richard Y. Wang

著

黄伟 王嘉寅 苏秦 冯耕中等 编译

高等教育出版社

“十三五”国家重点图书  
大数据科学丛书

Journey to Data Quality

# 数据质量征途

[美]Yang W. Lee, Leo L. Pipino,  
James D. Funk, Richard Y. Wang 著

黄伟 王嘉寅 苏秦 冯耕中等 编译

高等教育出版社·北京

Translation from English Language edition:

*Journey to Data Quality*

by Yang W. Lee, Leo L. Pipino, James D. Funk, and Richard Y. Wang

Copyright © The MIT Press 2006

All Rights Reserved

### 图书在版编目(CIP)数据

数据质量征途 / (美)杨·W.李(Yang W.Lee)等著;  
黄伟等编译. -- 北京: 高等教育出版社, 2017.8

ISBN 978-7-04-047797-9

I. ①数… II. ①杨… ②黄… III. ①数据处理  
IV. ①TP274

中国版本图书馆 CIP 数据核字(2017)第 126583 号

Shuju Zhiliang Zhengtu

策划编辑 冯 英

责任编辑 冯 英

封面设计 张志奇

版式设计 童 丹

责任校对 吕红颖

责任印制 耿 轩

出版发行 高等教育出版社

社 址 北京市西城区德外大街 4 号

邮政编码 100120

印 刷 北京市白帆印务有限公司

开 本 787mm×1092mm 1/16

印 张 13

字 数 330 千字

购书热线 010-58581118

咨询电话 400-810-0598

网 址 <http://www.hep.edu.cn>

<http://www.hep.com.cn>

网上订购 <http://www.hepmall.com.cn>

<http://www.hepmall.com>

<http://www.hepmall.cn>

版 次 2017 年 8 月第 1 版

印 次 2017 年 8 月第 1 次印刷

定 价 59.00 元

本书如有缺页、倒页、脱页等质量问题, 请到所购图书销售部门联系调换

版权所有 侵权必究

物 料 号 47797-00

# 序一

大数据是数字时代的新型战略资源,也是服务创新、驱动发展的重要抓手。大数据是数据科学的一个应用,也是数据科学重要的发展方向。近年来,大数据的热潮与数据科学的发展互为促进,正改变着人们的生产、生活方式。而对于学者和业界来说,是否能够抓住机遇、深入研究、形成解决方案,就显得非常必要和紧迫。

由于大数据具有分散存储、整合使用,分析处理的时间、空间复杂度高,以及数据整体及其关系协同呈现高价值的三大特征,数据质量往往难以保障。但是数据质量对于使用、用好大数据起到决定性的作用。数据质量低不仅会降低决策质量,更可能带来难以估量的灾难性后果,所以保障和提高大数据质量的工作迫在眉睫。

*Journey to Data Quality* 一书堪称数据质量领域的经典之作。该书从数据质量的概念入手,结合案例和分析工具,深入浅出地总结了美国学术界和产业界十余年的成果和经验,具有很强的指导性和实用性。由黄伟教授统筹,联合了国内外几位学者编译该书,并将最近几年的研究成果融入其中,是一件很有意义的工作。对于国内致力于数据质量的学者和业界来说,本书可以提供基础性的介绍和指导,为解决大数据环境下的数据质量问题指出方向。

徐宗本教授  
中国科学院院士  
2015年5月

## 序二

大数据时代的到来,不仅不断改变着人们认知世界的方式,而且正以其独特的影响力,推动着各行各业发展进步的新潮流。伴随着各种大数据技术的不断涌现,无论是在生命科学、医学,还是在金融、管理等诸多领域,大数据都显现出越来越重要的作用。大数据产业和管理将信息技术广泛而深入地应用于不同的行业,使后者拥有快速获取、高效存储、精准分析、正确判断各类数据和信息的能力,从而实现组织的科学决策。

积极抓住大数据带来的技术创新和产业变革的契机,成为国家间竞争和取得领先优势的关键。美国的大数据研发计划法案提出加强大数据相关领域的研究,将大数据及其产业作为国家推动的一个战略方向。党的十八大明确提出了创新驱动发展战略,大数据产业的发展以其内生的创新优势也必将成为重中之重。如何利用大数据提高国家的竞争力,如何发挥大数据的作用,成为政府、各行业都十分关注的问题——这其中,数据质量是一切的保障。

由黄伟教授牵头,联合西安交通大学和美国圣路易斯华盛顿大学的几位学者共同编译了 *Journey to Data Quality* 一书。该书不仅在数据质量研究领域享有盛誉,而且得到产业界的高度赞许和推崇,这是非常难能可贵的。书中总结了数据质量的研究成果和产业界多年的实践经验,并且通过对经典案例的深入剖析实现了理论与实践应用的统一。本书对于研究者、实践者,特别是管理者来说,将会带来重要的启迪。

汪应洛教授

中国工程院院士

西安交通大学管理学院名誉院长

2015年5月

# 前言

本书汇集了作者和本领域众多学者的研究成果,也包括他们在政府和相关行业实践工作中累积的经验。本书尝试将分别来自于学术期刊和学术实践性会议中的诸多观点、概念予以总结和升华,向读者展现这些观点和概念是如何被许多组织采纳并用作数据质量管理和实践的原则、政策和技术工具的。进一步地,作者将通过具体的现实例子和来自行业的案例对本书中的理论观点和方法加以讨论。

本书的读者群主要是企业的管理层、从事数据质量工作的人员、数据质量领域的研究者和学生。对于业界人员,本书有助于深入理解他们所从事工作的理论基础,为将来更好地解决问题并付诸实践做好准备。研究人员能够通过本书了解数据质量理论是怎样被应用到实践中的,进而有助于更加专注于未来的研究领域。而对于学生来说,本书能够提供对于这个领域的宏观认知,为今后在这一领域的学习和研究奠定坚实基础。管理层人员则可能会对本书的前几章和第11章(数据质量政策)更感兴趣并得以裨益。

# 致谢

很多人都为本书的出版做出了贡献。我们感谢匿名评审,是你们看到本书第一稿中的知识和实际潜力,并用心帮助我们形成连贯的书。我们感谢同事们审阅本书,并提出宝贵的意见。在之前的研究和咨询工作中,许多研究和从业人员与我们一起合作,从而奠定了这本书的基础。他们是 Stuart Madnick、Don Ballou、Harry Pazer、Giri Tayi、Diane Strong、Beverly Kahn、Elizabeth Pierce、Stanley Dobbs 和 Shobha Chengular-Smith。Raïssa Katz-Haas 和 Bruce Davidson 为本书中的数据质量实践案例提供了很多素材。

此外,我们感谢产业界的实践者,开放他们的实践环节和机构,并允许我们将此作为我们的实验场。

本书中的一些成果借鉴了其他作者之前的出版物。感谢下列期刊允许我们使用相关内容: *Communications of the ACM*、*Sloan Management Review*、*Prentice Hall*、*Journal of Management Information Systems*、*Journal of Database Management*、*Information and Management*、*International Journal of Healthcare Technology and Management* 和 *IEEE Computer*。

麻省理工学院(Massachusetts Institute of Technology, MIT)出版社的 Doug Sery 在整个出版过程中不遗余力地为本书提供独特的见解和毫无保留的支持。

如果没有卓有成效的工作环境,这本书也不可能完成,感谢麻省理工学院信息质量项目、剑桥(Cambridge)研究小组和麻省理工学院全面数据质量管理(TDQM)项目,以及睿智和富于奉献精神同事们: Tony Nguyen、Karen Tran、Rith Peou、Andrey Sutanto、John Maglio 和 Jeff Richardson,感谢东北大学(Northeastern University)工商管理学院的大力支持,感谢马萨诸塞大学洛厄尔分校(University of Massachusetts Lowell)管理学院的支持。

最后,我们感谢自己的家人,谢谢他们在编写这本书的漫长过程中付出的爱、支持和理解。Albina Bertolotti、Marcella Croci、Karen Funk、Karalyn Smith、Marc Smith、Jacob Smith、Cameron Smith、Kirsten Ertl、Phil Ertl、Austin Ertl、Logan Ertl、Lorenzina Gustafson、Laura Gustafson、Juliana Gustafson 和 Fori Wang,给我们的生活带来这么多快乐和幸福。特别地,我们要感谢我们的父母,谢谢他们培养我们对学习的热爱。

# 目录

<b>第 1 章 引言</b> .....	1
1.1 信息可以被共享吗 .....	2
1.2 新系统不是解决办法 .....	2
1.3 开启数据质量之旅 .....	4
1.4 成功开始的故事 .....	4
1.5 CEO 领导的旅程 .....	6
1.6 数据质量之旅面临的挑战 .....	6
1.7 数据质量为什么重要 .....	7
1.8 本书概览 .....	8
<b>第 2 章 成本—效益分析</b> .....	11
2.1 挑战性 .....	11
2.2 成本—收益的权衡 .....	13
2.3 一个案例 .....	15
2.4 高级成本—效益分析技术 .....	17
2.5 本章小结 .....	20
<b>第 3 章 数据质量评估（一）</b> .....	21
3.1 评估技术和相关方法 .....	21
3.2 实际中的评价方法 .....	22
3.3 差距分析技术 .....	28
3.4 数据完整性评价 .....	30
3.5 本章小结 .....	31
附录 数据质量评价调查(IQA)问卷 .....	31
<b>第 4 章 数据质量评估（二）</b> .....	43
4.1 科德完整性约束 .....	43
4.2 数据质量指标 .....	44



4.3	自动化的测量方法 .....	48
4.4	嵌入过程的数据整体性方法 .....	51
4.5	本章小结 .....	53
<b>第5章</b>	<b>保证数据质量的抽样方法 .....</b>	<b>55</b>
5.1	基本概念 .....	55
5.2	选择抽样过程 .....	57
5.3	确定样本量 .....	58
5.4	交易数据库的抽样 .....	59
5.5	环境扩展:分布式数据库和数据仓库 .....	62
5.6	本章小结 .....	62
<b>第6章</b>	<b>数据质量问题及其模式剖析 .....</b>	<b>65</b>
6.1	数据质量问题的十大根源 .....	65
6.2	数据质量问题的表现 .....	73
6.3	数据质量问题的转换 .....	83
6.4	本章小结 .....	85
<b>第7章</b>	<b>识别数据质量问题的根本原因——一个医疗保健组织案例 .....</b>	<b>87</b>
7.1	案例:好感觉健康系统公司 .....	87
7.2	识别问题 .....	88
7.3	组建跨部门的团队 .....	91
7.4	采用一种框架:建立并测试假设 .....	92
7.5	关键信息 .....	92
7.6	找出数据质量问题的诱因 .....	93
7.7	本章小结 .....	98
<b>第8章</b>	<b>数据的产品化管理 .....</b>	<b>99</b>
8.1	数据产品 .....	99
8.2	四个案例 .....	100
8.3	四个原则 .....	101
8.4	把数据当成副产品来管理是无效的 .....	103
8.5	本章小结 .....	106

<b>第 9 章</b>	<b>开发数据产品地图</b>	107
9.1	数据产品地图的概念、定义和符号	107
9.2	绘制数据产品地图的步骤	111
9.3	建立数据产品地图的一个案例	112
9.4	本章小结	117
附录	基于 IPMAP 的图形化编辑软件	117
<b>第 10 章</b>	<b>数据质量实践——一家大型教学医院的案例</b>	123
10.1	LTH 健康系统案例研究	123
10.2	提交数据质量改进项目	127
10.3	数据产品地图	129
10.4	改进方案:当前的处理过程和未来计划	135
10.5	本章小结	136
<b>第 11 章</b>	<b>数据质量政策</b>	139
11.1	十大政策指引	140
11.2	本章小结	147
附录 1	数据质量岗位介绍	148
附录 2	来自全球制造公司的数据架构政策示例	152
附录 3	数据质量实践与产品评估工具	153
<b>第 12 章</b>	<b>旅途结束了吗</b>	159
12.1	要点回顾	159
12.2	面临的挑战和威胁	160
12.3	对数据质量特征的规范定义	161
12.4	公司家族化	164
12.5	数据挖掘	166
12.6	数据集成	167
12.7	安全性	168
12.8	有线和无线的世界	168
12.9	后记	169

附录 一种基于期望失验理论的信息质量评估指标体系·····	171
F.1 引言·····	171
F.2 文献回顾·····	172
F.3 信息质量的概念·····	175
F.4 信息质量的指标体系·····	178
F.5 讨论·····	180
参考文献·····	183
编译者后记·····	191

# 第1章 引言

Jane Fine 是一家全球制造公司的信息系统主管。她正坐在自己的位子上思考着下一步应该怎么办：对于销售副总裁提出的“对某个大客户的销售总量是多少？”的问题，她很难做出回答。因为她知道公司的数据库中存在重复的客户标识和重复的产品代码，这样不恰当的数据库结构及其带来的不准确内容，令她难以回答销售总量的问题。在千里之外的一家大型教学医院，负责医疗事务的常务副总裁 Jim Brace 也正面临着不同的窘境。他刚刚参加了一个由医院负责人召集的会议，州政府的监管机构对医院的数据报告中部分数据的真实性提出了质疑，这导致委员会拒绝批准医院的报告，而医院将因此失去州政府的财政补贴，这一情况必须尽快得到处理。与此同时，政府人力资源部门的高级信息专员 Dean Gary 在审核人事文件时，发现一些文件中的数据与该部门最新的人事报告中的一些汇总数字存在出入，数据存在不一致的现象。尽管这个问题目前还没有被人事工作报告的使用者发现，也暂时没有从其他方面体现出来，但是不可能永远不被发现。

在以上案例中，前两个是信息系统的上层管理者发现的数据质量问题。如果问题继续存在，组织的高层管理者很可能基于低质量的数据做出糟糕的决策。例如，在大型教学医院的案例中，医院负责人对数据质量问题的关注是由外部监管机构的质疑引起的。

在大多数组织中，如何让管理层相信数据质量存在问题，并制定一个正式的数据质量计划，是一项挑战。比如，在全球制造公司中，高层通常要在一个可接受的时间内获取他们要求的数据。对他们而言，这些数据不存在问题。他们也看不到下面的人为了满足他们的数据需求，清理不一致数据所做的额外工作。尽管在修正数据质量问题时花费的时间并不多，但是这些时间本可以用于更加高效地完成其他任务。

在以上三个案例中，管理者们不得不面对他们组织中的数据质量问题。可以猜测大多数组织应该也面临着相似的问题。

## 1.1 信息可以被共享吗

一些组织在执行跨业务的流程,或者尝试跨系统、跨组织交互时,常常难以充分地利用信息。当这些组织相信他们拥有完成业务功能的数据但却不能顺利开展业务时,组织内部就容易产生挫败情绪。例如,某公司希望做一些趋势分析,以便与客户和其他合作者构建更紧密的关系,但是该公司的信息技术部门却经常不能提供客户所要求的整合性信息,或者无法按客户要求的时间提供其所需的信息,这导致公司错过了利用这些收集和储存的信息的最佳时机。更糟糕的是,竞争对手却能迅速反应,战略性地应用类似的信息。

很多组织长期以来都面临这些问题。数据质量问题还可能表现在其他方面,具体包括:

- 许多跨国公司难以管理其全球的数据,虽然这些数据可以用于解决公司当前以及未来全球性的、区域性的业务问题。
- 外部检查使组织内部的数据质量问题浮出水面,正如前面提到的监管机构对医院的医保报销和患者投诉进行审查的事例。
- 信息系统项目也能够揭示出存在的数据问题,特别是一些涉及跨业务的、多数据来源的数据质量问题。
- 组织成员在工作中发现了数据质量问题,却只使用某些变通方法临时满足数据需求,而不是使用或创建永久性的、持续性的解决方案。

## 1.2 新系统不是解决办法

每一个组织都希望自己拥有高质量的数据,但是常常不知道如何实现这个目标。一类常见的做法是开发一个新系统来取代旧系统,然而常常会在实施之后立即后悔。这是因为公司实施此类方案时,总是重建一套全新的系统,却很少在第一时间考虑原系统存在困难的真正原因——数据质量问题。比如信息系统部门往往热衷于使用最新的技术,开发更流行或更常见的软、硬件解决方案,我们将这种方法称为系统驱动型解决方案。此时,公司采取的方案的真实目标退化为开发新系统,而非修正数据质量问题以提供高质量的数据。显然,这种舍本逐末的新系统非但不能解决原有问题,而且很有可能加剧数据质量问题。即使某个解决方案偶尔会有成效,通常真正造成问题的原因却更容易被掩盖或进一步隐藏。

许多公司误以为使用了最新的软件,比如企业资源规划(enterprise resource planning,ERP)系统或者紧跟潮流地引入一个数据仓库(data warehouse,DW)就会坐享更高质量的数据。公司希望通过这些系统更好地实现公司范围内的信息共享。然而,信息技术部门在整合多来源数据的过程中越来越清醒地认识到,数据定义、数据格式以及数据的值都可能存在大量的不一致现象,但是时间等各方面的压力迫使他们依然继续使用之前存在的同样糟糕的数据。

许多公司感到失望的是,在数据仓库上付出大量努力却没有得到较好的商业价值。在众多案例中,许多采用了 ERP 系统和数据仓库的公司并没有获得最初承诺的预期商业价值。

依然以前面提到的全球制造公司为例,公司试图整合全球范围的销售信息,尽管公司有全部的原始数据,却依然需要花费几个月的时间才能实际提供某个指定客户商业需求的一套有用的数据。存在的问题包括:同一个客户对应多个标识,以及多个客户被赋予同一个标识。此外,子系统中储存的数据在公司层面没有合理的定义和记录,物理数据库并不是一直可以访问的,公司内部没有对概念和术语的定义实行标准化,与标准不同的内容没有被记录或不能被共享,等等。

在此阶段,公司已经在这个项目上花费了大量的预算,管理层还愿意增加预算吗?或者管理层会终止这个项目吗?如果管理层增加预算,但是依旧没有对基础业务和数据质量问题予以足够的关注,上述状况会有实质性的改变吗?公司是否应该致力于创建另一个业务流程,但是与新业务流程相关的费用会不会导致商业价值没有增加?如果管理层终止项目,那么公司共享信息的努力也将终止。此时无论增加预算还是终止项目,公司整体的数据质量都将不会有任何提高。

大多数组织总是狭隘地关注系统层面的问题,却一再忽视数据层面的问题。解决数据质量问题可以增加组织内部共享信息的能力;反之,如果忽视数据问题,大多数系统层的解决方案最终都将失败。那么,如果公司发现自身可能存在严重的数据问题应该怎么办呢?

一些组织通过使用基础性的数据清理软件来尝试改进数据质量。在全球制造公司的案例中,通过初始的数据清理建立了一个可用的数据仓库。然而,随着时间的推移,数据仓库内的数据质量再次急剧下降。在更普遍的案例中,企业通常会指派个别人员去解决特定的数据质量问题,或者某个、某些对数据质量问题感兴趣的人主动解决了其中的问题。但是,不论这些问题是怎样被发现的,无论某个人如何成为问题的负责人,最初的调查和解决方案通常都是临时方案。

许多企业已经应用了多种多样的临时方法,却依然得不到尽如人意的结果。此时,数据质量项目可能会被迅速终止。这种情况下,重新开始新的数据质量方案将变得非常困难。本书中,我们将提供更系统、更全面的基础性解决思路 and 方案。

## 1.3 开启数据质量之旅

让我们重温上述案例的场景,看看他们采取了什么后续行动。全球制造公司的 Jane Fine 做了某些调研,在自己的数据库管理经验和数据质量领域知识的基础上,她开始了解业界和学术界的前沿发展和行业状况并参加了多个研讨会。她广泛地搜索外部的资源,试图获取解决问题所需的知识,通过部署技术和流程来改善公司的这一问题。当然,她仍然面临着很多挑战。

在医院负责人的支持下,Jim Brace 编写了一个独立的内部软件来尝试解决问题。在此基础上,他获得了一些反馈建议并得以实施。但是数据质量问题依然存在,所以 Brace 通过查阅全面数据质量管理方面的文献,采用测量的方法实现了数据质量的改善。该方法取得了一定的成功,得益于外部的技术和知识,Brace 采取更主动、更全面的方法设计出一个可持续并切实可行的数据质量方案。

人力资源部门的 Dean Gary 利用经典数据库理论中数据整合的概念和数据挖掘技术解决了他的问题。他使用一种技术手段力求识别出不同类型的数据错误,这可以帮助他解决当前的问题。在数据分析前清理所接收的数据,使之能够提供有效的报告。然而,他无法识别并消除数据不一致现象发生的根源,所以他在每次收到新的数据时,都不得不重复进行数据清理,这促使他开始寻求其他改进数据质量的方法。

三位管理者都有意或无意地踏上了数据质量之旅。有许多不同的路可供选择。如果选择了某条合适的路径,伴随着旅程,数据质量将会不断提升,即使在旅程中会不断遇到新的数据质量问题。这种发现问题、提高质量、解决问题的过程形成一个周而复始的循环。在经历了几个循环之后,低质量的数据对组织的影响将会快速降低。然而,重复过程仍将会继续。问题的解决方案会产生新的问题,这些问题会激发新的需求,而新的需求又会产生新的问题。

## 1.4 成功开始的故事

许多组织已经走上了数据质量之旅,但多数仅仅是为了尽早地完结这一旅

程,而从未意识到在数据质量上持续努力的益处。过早地完结数据质量项目会导致组织反复遭受数据质量问题的困扰。导致数据质量项目过早完结的主要原因是缺乏对这一旅程预期的理解,进而提早泄气并失去坚持的动力。组织应该认识到,数据状况不是一夜之间形成的,也就不能期盼一蹴而就、一劳永逸的解决。

回顾上述案例我们发现,每个管理者从不同的点出发、经由不同的路径,最后都走向相同的目标——一种长期的、可持续的数据质量改善,并采用自己的方法使之适应具体组织的特点。

基于对环境的调研,全球制造公司的 Jane Fine 引入了一个提升数据质量意识的项目。此外,通过测量组织对数据质量的主观评价,进一步测量数据库中的数据的完整度,分析的结果使她开始关注数据不一致的问题。在已经获得某些成功的基础上,她现在能够回答一些问题——比如某一产品在全球范围内对某个大客户的销售总量是多少。这一成功得到了高层管理者的认可,高层管理者任命她领导一个全公司范围内的数据质量项目, Fine 需要考虑她下一步的行动。

在组织人力资源部门的 Dean Gray 能够清理接收的数据并用清理后的数据来准备他的报告。然而,他不得不重复这样的工作。这促使他探寻根除数据错误的解决方案,他发现这些问题的根源涉及超出他控制范围的外部数据资源系统。他现在面临的挑战是选择一个可行的解决方法。

三人之中,医院负责人 Jim Brace 在数据质量之旅中走得最远。由于使用了传统的全面数据质量管理方法,Brace 能够向管理层展示像管理产品一样管理信息的概念,这种管理类似于制造业中实体产品的管理过程。高层管理者赞同这一观点,并委任他领导一个数据质量管理工作小组。此外,高层管理者确定了以奖励为基础的数据质量目标。作为回应,Brace 提出了数据质量政策,得到了董事会的一致认可。得益于明确的数据质量政策,他和工作组进一步绘制了数据产品地图(information product map, IPMAP),进一步帮助他们处理监管机构发现的数据质量问题。

正如以上三个案例所述,每个管理者都面临着不同的挑战,也在积极地应对这些挑战。然而,要想保持稳定的数据质量水平并促使数据质量的长期提高,仅靠指定个别管理者对组织的数据质量进行改进是远远不够的,公司的高层必须参与并指导这个过程。令人十分鼓舞的是,企业的高层管理者已经越来越深入地参与到数据质量项目中来,公司的首席执行官(CEO)参与到数据质量项目的重要性怎么强调都不为过。



## 1.5 CEO 领导的旅程

CEO 必须有踏入数据质量之旅的愿景。高层管理者很容易轻视数据质量项目,把它们归为低优先级的项目,而更看重资源稀缺的竞争项目。CEO 通常难以意识到低质量数据引起的问题的严重性。具有讽刺意味的是,在 CEO 身后经常有一批人致力于解决数据质量问题,而这无疑是一项昂贵的开支。CEO 有很多理由可以意识到支持数据质量项目的必要性,以及使用低质量的数据很可能会带来的危机。例如在 Jim Brace 的案例中,当监管机构质疑数据的有效性并因此拒绝了医院的管理报告时,相关决定一定会交到医院 CEO 的手中。由于医院 CEO 和其他高层管理者以医院优质的服务和良好的声誉为荣,他们将亲自批准具有战略目标的数据质量项目。当医院 CEO 无法获得所需数据,或者数据冲突浮现,或者使用错误数据制定决策时,他们可能会更直接地遇到问题。

除非遭遇危机或者灾难性的事件,需要反复沟通,对案例进行严谨的分析和推理,得出可以说服 CEO 开展数据质量项目的结论。要取得支持,这种分析必须包括有说服力的价值分析或成本-效益分析。然而,一旦拥有 CEO 的全力支持,将对成功的数据质量之旅产生极大的帮助。

## 1.6 数据质量之旅面临的挑战

数据质量旅程面临的挑战将是众多且棘手的。我们已经认识到第一个,也许也是最艰巨的挑战——获得 CEO 对数据质量项目的认可和支持;其次,我们也需要采用强有力的经济角度的论证、成本-效益分析来支持数据质量项目的实施;在整个组织内,宣传推广数据质量意识则是第三个挑战。

获得高层管理者的支持是提高组织内的数据质量意识的第一步,因为提升数据质量意识需要从传统数据质量的视角中扩展和分离出新的视角。

这需要清晰地理解数据质量对组织意味着什么,以及数据质量的重要性。数据质量领域的研讨会常常提到数据与信息的区别、信息与知识的区别,数据、信息和知识是三个不同的概念。当把它们纳入一个层次结构,则知识包含信息,信息包含数据。但是武断地区分数据和信息会使我们偏离首要的工作,而且会妨碍我们对复杂的数据产品系统的理解。

管理者普遍使用的区分数据和信息的传统方法是:数据由原始事实或资料构成,而信息是经过加工的数据。然而,一个人的数据可能是另一个人的信息。