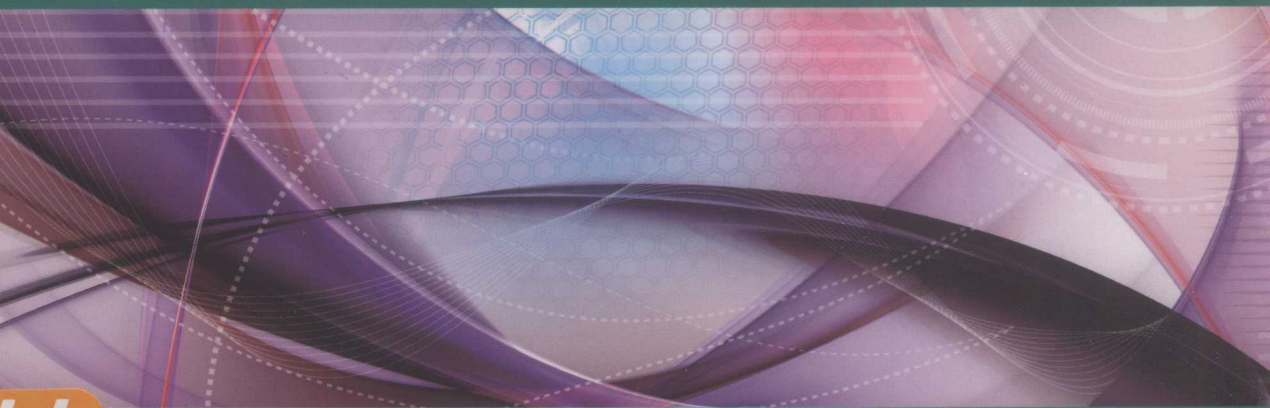




普通高等教育“十二五”规划教材

INTRODUCTION TO DATA SCIENCE

数据科学导论



INTRODUCTION TO DATA SCIENCE

杨 旭 汤海京 丁刚毅 主编

 北京理工大学出版社
BEIJING INSTITUTE OF TECHNOLOGY PRESS

014033116

TP274-43
52

普通高等教育“十二五”规划教材

数据科学导论

杨旭 汤海京 丁刚毅 主编



TP274-43
52

 北京理工大学出版社

BEIJING INSTITUTE OF TECHNOLOGY PRESS



北航

G1721307

内 容 简 介

数据科学，作为一门正在蓬勃发展的新学科，关注的是如何在大数据时代背景下运用各门与数据相关的技术和理论来服务社会。本书系统地讲述了与数据科学相关的各方面知识，着重培养数据工程师所需要的技能与思维。本书从与数据科学相关的概念出发，通过丰富翔实的案例，从各个方面展示数据科学的运用方式，让读者有一个更为直观的认识，也可以从中感受到运用数据科学处理各个领域问题的方法和流程，并且在其中穿插了数据科学研究方式下新的思维模式的讲解。本书还从工程概论的流程角度来讲述数据科学的工程体系架构，展望数据科学的未来发展。本书可作为计算机相关专业的本科生教材，也可供相关专业技术人员阅读参考。

版权专有 侵权必究

图书在版编目 (CIP) 数据

数据科学导论 / 杨旭, 汤海京, 丁刚毅主编. —北京: 北京理工大学出版社, 2014. 3
ISBN 978 - 7 - 5640 - 6384 - 9

I. ①数… II. ①杨… ②汤… ③丁… III. ①数据管理 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2014) 第 011575 号

出版发行 / 北京理工大学出版社有限责任公司

社 址 / 北京市海淀区中关村南大街 5 号

邮 编 / 100081

电 话 / (010) 68914775 (总编室)

82562903 (教材售后服务热线)

68948351 (其他图书服务热线)

网 址 / <http://www.bitpress.com.cn>

经 销 / 全国各地新华书店

印 刷 / 保定市巾画美凯印刷有限公司

开 本 / 787 毫米 × 1092 毫米 1/16

印 张 / 10

字 数 / 210 千字

版 次 / 2014 年 3 月第 1 版 2014 年 3 月第 1 次印刷

定 价 / 29.00 元

责任编辑 / 刘 娟

文案编辑 / 王晓莉

责任校对 / 周瑞红

责任印制 / 李志强

图书出现印装质量问题, 请拨打售后服务热线, 本社负责调换



前 言

我们已经身处于一个数据爆炸的时代，在新的时代背景下，需要运用新的科学研究方式去应对新的挑战。数据科学，作为一门正在蓬勃发展的新学科，所关注的正是在大数据时代背景下，运用各门与数据相关的技术和理论，服务社会，让我们可以更好地利用身边的数据，使生活变得更加美好。

本教材系统性地讲述了与数据科学相关的各方面知识，着重培养数据工程师所需要的技能与思维。本书将从与数据科学相关的概念出发，通过丰富、翔实的案例，从各方面展示数据科学的运用方式，并且在其中穿插数据科学研究方式下新的思维模式的讲解，让读者有一个更为直观的认识，也可以从中感受到运用数据科学处理各个领域问题的方法和流程。本教材还从工程概论的流程角度来讲述数据科学的工程体系架构，并展望数据科学的未来发展。

本教材由北京理工大学软件学院“数据科学与技术”课题组的杨旭老师、汤海京老师，以及北京理工大学软件学院院长丁刚毅老师担任主编。其中第2、3、4章由杨旭老师编写；第1、5、6、7章由汤海京老师编写；第8章由丁刚毅老师编写，并由丁刚毅老师负责全书的统稿。对书中存在的错误及不妥之处，恳请各位读者、同行批评指正。

编 者

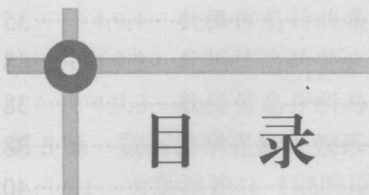


图目录

图 1-1	结绳记事	2
图 1-2	古埃及人用莎草记录数据	2
图 1-3	造纸流程	3
图 1-4	最早的留声机	3
图 1-5	最古老的照相机	4
图 1-6	早期的摄影机	4
图 1-7	世界上第一台电子计算机	4
图 1-8	当前世界上最快的计算机——天河 2 号	5
图 1-9	网络与数据	5
图 1-10	物联网	6
图 1-11	美国著名未来学家阿尔文·托夫勒	7
图 1-12	大数据时代来临	8
图 1-13	大数据的 4V 特征	9
图 1-14	大数据将会带来大变革	10
图 1-15	大数据带来新的洞察力	12
图 1-16	大数据下的企业布局	14
图 1-17	数据科学家	14
图 1-18	数据科学	15
图 2-1	啤酒与尿布	17
图 2-2	美国沃尔玛卖场	21
图 2-3	日本 7-11 卖场	21
图 2-4	Farecast 机票预测	25
图 2-5	Decide 电商比价网站	26
图 2-6	Decide 的比价无远弗届	27
图 2-7	潘多拉音乐盒子	30
图 2-8	潘多拉网络电台	31
图 2-9	潘多拉音乐基因组工程	31
图 2-10	定制化的电台音乐	32
图 2-11	潘多拉电台 App	32

图 2-12	Target 预测怀孕	34
图 2-13	“平静”指数和道琼斯指数对比	37
图 3-1	谷歌预测美国流感趋势	38
图 3-2	美国：流感样疾病（ILI）数据	40
图 3-3	科学家利用大数据监控脑伤病人的恢复情况	49
图 3-4	可穿戴设备	50
图 3-5	Maxim 生命体征测量 T 恤	51
图 3-6	TI 公司的 Health Tech	52
图 3-7	Valencell：可随身穿戴的微型生理监测模块	53
图 3-8	Google Glass	54
图 3-9	苹果的 iWatch	55
图 3-10	BrainLink 意念头箍	55
图 4-1	智慧城市	59
图 4-2	城市建设所面临的社会基础设施问题	62
图 4-3	松岛新城的服务构想	64
图 4-4	韩国政府斥巨资修建的松岛新城	65
图 4-5	美国俄亥俄州的哥伦布市	66
图 4-6	My Columbu 移动应用程序	67
图 4-7	爱沙尼亚的塔林市	71
图 4-8	阿姆斯特丹的 ASC 计划	72
图 4-9	里约热内卢的城市运营中心	74
图 4-10	里约热内卢的智慧城市建设	75
图 4-11	智慧城市建设必需的数据分析技术的发展趋势	79
图 5-1	《纸牌屋》	81
图 5-2	谷歌预测电影票房	86
图 5-3	2012 年票房收入与搜索量的曲线	88
图 5-4	2012 年票房收入和两类搜索量的曲线	88
图 5-5	搜索量与首周票房收入之间的关系	89
图 5-6	提前一周预测票房的效果	89
图 5-7	提前一个月预测票房的效果	90
图 5-8	微软研究院的戴维德·罗斯柴尔德	92
图 5-9	微软研究院的戴维德·罗斯柴尔德博士预测奥斯卡获奖名单	92
图 5-10	奥斯卡投票预测器	94
图 6-1	奥巴马	96
图 6-2	大数据与选举	98
图 6-3	美国“棱镜门”风波	101
图 6-4	谁是“棱镜”计划的帮凶	101
图 6-5	美国“棱镜”计划	106

图 6-6	大数据时代的个人隐私安全	109
图 6-7	Ancestry.com	111
图 6-8	家谱网站帮助寻根问祖	113
图 7-1	数据的分类	124
图 7-2	Google Chart API	136
图 7-3	Raphaël	136
图 7-4	Visual.ly	136
图 7-5	Crossfilter	137
图 7-6	PolyMaps	137
图 7-7	Kartograph	138
图 7-8	Processing 编程环境	138
图 7-9	R 语言编程	139
图 7-10	Weka 编程环境	140
图 7-11	Gephi 做数据可视化	140
图 8-1	麦肯锡对数据科学方面人才需求空缺的预测	142
图 8-2	数据科学从业人员的未来成长性	143
图 8-3	明确数据的优势和不足	145



目 录

第 1 章 引论	1
1.1 序言	1
1.2 数据	1
1.2.1 数据的概念	1
1.2.2 数据的发展史	2
1.2.3 数据、信息与知识	6
1.3 大数据	7
1.3.1 大数据时代的来临	7
1.3.2 大数据的概念	8
1.3.3 大数据的特征	9
1.3.4 大数据对社会所产生的影响	10
1.3.5 迎接大数据时代的挑战	13
1.4 数据科学	15
1.5 本书结构	16
第 2 章 数据科学在商业金融领域的应用	17
2.1 啤酒与尿布	17
2.1.1 案例详析	17
2.1.2 数据挖掘技术	19
2.1.3 购物篮分析法	20
2.1.4 对我们的思维模式启示	21
2.2 比价网站的成功	23
2.2.1 Farecast 案例详析	23
2.2.2 Decide 案例详析	26
2.2.3 对我们的思维模式启示	27
2.3 基于大数据的个性化推荐系统	29
2.3.1 基于亚马逊的个性化推荐系统	29
2.3.2 潘多拉 (Pandora) —— 基于基因的推荐系统	29
2.4 Target 的大数据营销	33
2.4.1 案例详析	33

2.4.2	给我们的思维模式启示	34
2.5	社交网络数据之于对冲基金	35
第3章	数据科学在生物医学领域的应用	38
3.1	流行病预测	38
3.1.1	谷歌的流感预测	38
3.1.2	利用微博来预测流感	40
3.1.3	给我们的思维模式启示——大数据时代的科学伦理问题	41
3.2	大数据与智慧医疗	42
3.2.1	临床操作	43
3.2.2	付款/定价	44
3.2.3	研发	45
3.2.4	新的商业模式	46
3.2.5	公众健康	47
3.2.6	给我们的思维模式启示	47
3.3	疾病监控	48
3.3.1	大数据服务心脏病患者	48
3.3.2	“魔毯”病人的监控	49
3.3.3	大数据监测脑外伤病人恢复	49
3.4	可穿戴技术、大数据与智慧医疗	50
3.4.1	什么是可穿戴技术	50
3.4.2	可穿戴设备简析	51
3.4.3	可穿戴设备与智慧医疗	55
3.4.4	给我们的思维模式启示——可穿戴设备的缺陷	56
第4章	数据科学在智慧城市领域的应用	59
4.1	概述	59
4.1.1	什么是智慧城市	59
4.1.2	产生背景	61
4.1.3	IT企业相继介入智慧城市领域	62
4.1.4	国际实践	63
4.2	韩国的松岛新城	64
4.3	美国的智慧城市建设	65
4.3.1	哥伦布市	65
4.3.2	其他智慧城市建设的举措	68
4.4	英国的智慧城市建设	69
4.5	日本的智慧城市建设	70
4.6	北欧智慧城市——爱沙尼亚	71
4.7	荷兰阿姆斯特丹的智慧城市计划	72
4.8	巴西里约热内卢的智慧城市建设	74

4.9 智慧城市建设中所应用的数据科学技术	78
4.9.1 数据信息的收集: 利用传感网络收集数据信息	78
4.9.2 数据信息的整合: 不同数据信息的整合和统一管理	78
4.9.3 数据信息分析与应用: 大容量、实时性分析技术	79
第5章 数据科学在影视娱乐领域的应用	81
5.1 大数据捧红《纸牌屋》	81
5.1.1 案例详析	81
5.1.2 大数据如何捧红《纸牌屋》	82
5.1.3 给我们的思维模式启示	83
5.2 谷歌预测电影票房	86
5.2.1 案例详析	86
5.2.2 谷歌的预测机理	87
5.2.3 给我们的思维模式启示	90
5.3 利用数据预测奥斯卡奖项	92
第6章 数据科学在其他领域的应用实例	96
6.1 大数据帮助奥巴马赢得大选	96
6.1.1 案例详析	96
6.1.2 给我们的思维模式启示	100
6.2 棱镜门	101
6.2.1 案例详析	101
6.2.2 “棱镜”计划	105
6.2.3 加拿大的“棱镜门”	107
6.2.4 给我们的思维模式启示	108
6.3 大数据帮助寻根问祖	111
6.3.1 案例分析	111
6.3.2 运作机理	113
6.4 大数据与社会治安	115
第7章 数据科学工程概论	116
7.1 科学研究的第四范式——数据密集型研究方法	116
7.1.1 范式和范式的演变	116
7.1.2 科学研究的第四范式	117
7.2 数据密集型科学研究兴起的社会环境	118
7.2.1 数据洪流的到来	118
7.2.2 科学界对海量数据的关注	118
7.2.3 关联数据运动	119
7.2.4 政府数据开放运动	120
7.3 对数据密集型科学研究范式的分析	121
7.3.1 科学数据和科学研究的问题	122

7.3.2	相应的解决方案	122
7.4	数据的收集	123
7.4.1	客观世界 (Matter) 中的数据	123
7.4.2	主观世界 (Mind) 中的数据	124
7.4.3	细谈数据	124
7.5	数据的存储	125
7.6	数据的管理	126
7.6.1	NoSQL 数据库简介	126
7.6.2	NoSQL 数据库的特点	128
7.6.3	开源的 NoSQL 数据库软件	129
7.7	数据的处理	131
7.7.1	Hadoop 的起源	132
7.7.2	优点	132
7.7.3	架构	133
7.7.4	MapReduce 流程	134
7.8	数据的可视化	135
7.8.1	Excel	135
7.8.2	Raphaël	136
7.8.3	Visual.ly	136
7.8.4	Crossfilter	137
7.8.5	PolyMaps	137
7.8.6	Kartograph	137
7.8.7	Processing	138
7.8.8	R	138
7.8.9	Weka	139
7.8.10	Gephi	140
第 8 章	数据科学的未来展望	141
8.1	从业前景广阔	141
8.2	对未来数据科学发展的探讨	144
8.2.1	提防进入数据误区	144
8.2.2	数据不是万能的	144

第1章

引 论

1.1 序 言

“数据科学”从出现到现在已经有三十多年的历史了，其中涉及了很多方面的内容，涵盖数学、统计学、数据工程、模式识别、机器学习、高性能计算、可视化、数据仓库以及数据建模等多个领域的技术和理论。数据科学的最终目的就是从数据中挖掘出有用的信息，让数据增值。

虽然已有三十年历史了，但数据科学仍然是一门新兴的学科，尽管之前运用较多的是在计算智能或者是商业分析方面，但已经慢慢地深入到了人类社会的各个方面。之所以要开设数据科学这门学科，就是为了培养这门学科的专业人员，使他们运用所有可以得到的数据，寻找其背后的故事，从而找到办法让这些数据所蕴含的意义可以轻易地被人们所理解，即便他们不具备数据科学的相关知识。

本章将会首先从数据开始谈起，向读者讲述数据的发展概况和现状，从而感受学习数据科学的重要性。

1.2 数 据

1.2.1 数据的概念

数据科学这门学科研究的核心内容就是数据，那究竟什么是数据呢？一提到数据，我们首先想到的会是数字。但数据并不局限于数字，文本、音频、图像、视频都可以是数据。在本书里，我们对数据给出如下的定义：

数据是指以定性或者定量的方式来描述事物的符号记录，是可定义为意义的实体，它涉及事物的存在形式。

简单说来，数据就是人为创造的一种对事物的表示方式，是通过观察或者实验得来的对现实世界中的地方、事件、对象或概念的描述和反映。

数据可以是连续的值，例如声音，称为模拟数据；也可以是不连续（离散）的值，例如成绩，称为数字数据。

1.2.2 数据的发展史

人类历史上最早的有记录的数据，可以追溯到穴居的原始人时期。当时的人类，会在作为居处的洞穴墙壁上，以石器或者骨器刻画来记录数据。这些被记录的数据，或者是简单的记录日期的刻痕，或者是形象化地记载一些日常发生事件的壁画。

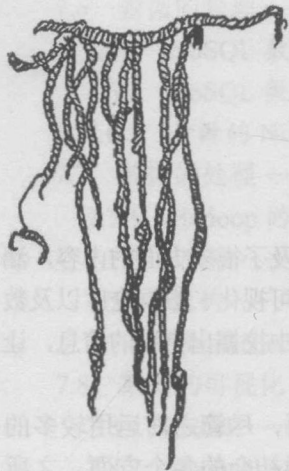


图 1-1 结绳记事

后来，人们创造了结绳记事的方式来记录数据。《周易·系辞下》中有云：“上古结绳而治，后世圣人易之以书契。”即在一根绳子上打结，用以记事。上古时期的中国及秘鲁印第安人皆有此习惯。到了近代，一些没有文字的民族，仍然采用结绳记事来作为数据记录方式传递信息。古人采取的结绳方法，据古书记载为：“事大，大结其绳；事小，小结其绳，结之多少，随物众寡。”

图 1-1 所示的是古代印加人采用的一种结绳记事的方法，用来计数或者记录历史。事大，大结其绳；事小，小结其绳。不过，这种记事的方法已经失传，目前还没有人能够了解其全部含义。

随着数字和文字的出现，古人开始以更加明确的形式来记录数据。古埃及人创造了莎草纸，用来进行记录。埃及博物馆中陈列的各种莎草纸文书、图画表明，莎草纸是人类历史上最早、最便利的书写材料之一，是记录古埃及历史的主要载体（图 1-2）。

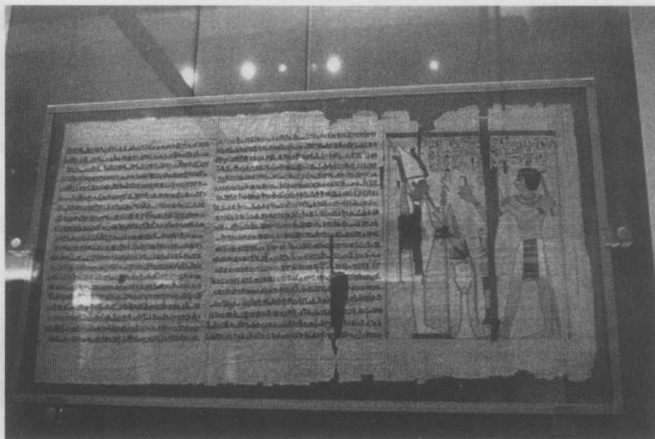


图 1-2 古埃及人用莎草记录数据

我们的祖先在汉代发明了造纸术（图 1-3），这种数据记录的方式一直延续到现在。这里一定要强调，本书中所提到的数据，不光是涉及数字形式的数据，以文本、图像、语音等方式所记录的数据都是数据科学所研究的内容。造纸术的发明和改进，让文本形式的数据记录方式开始盛行起来。

最早的录音机，也叫留声机，诞生于 1877 年，是发明大王——爱迪生所制造的。爱

迪生发现了电话传话器里的模板随着说话声而震动的现象，于是他拿短针做了试验，从中得到了很大的启发。说话的快慢高低能使短针产生相应的颤动。那么，反过来，这种颤动也一定能发出原先的说话声音。于是，他开始研究声音重发的问题。

在1877年8月15日，爱迪生让助手克瑞西按图样制造出一台由大圆筒、曲柄、受话机和模板组成的怪机器。爱迪生指着这台怪机器对助手说：“这是一台会说话的机器。”他取出一张锡箔，卷在刻有螺旋槽纹的金属圆筒上，让针的一头轻擦着锡箔转动，另一头和受话机连接。爱迪生摇动曲柄，对着受话机唱起了“玛丽有只小羊羔，雪球儿似一身毛……”。唱完后，把针又放回原处，再轻悠悠地摇动曲柄。接着，机器不紧不慢、一圈又一圈地转动着，唱起了“玛丽有只小羊羔……”，与刚才爱迪生唱的一模一样。在他身旁的助手们，见到一架会说话的机器，都惊讶得说不出话来。

“会说话的机器”诞生的消息，轰动了全世界。1877年12月，爱迪生公开表演了留声机（图1-4），外界舆论马上把他誉为“科学界之拿破仑”，留声机是19世纪最让人振奋的三大发明之一。即将开幕的巴黎世界博览会立即把它作为时新展品展出。就连当时美国总统海斯也在留声机旁转了2个多小时。



图 1-3 造纸流程

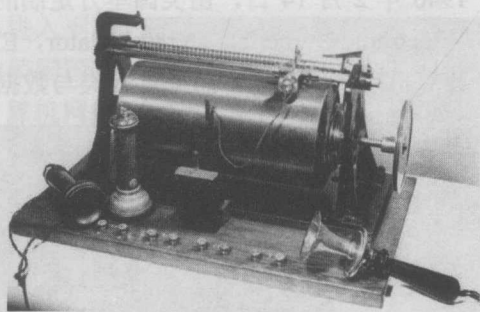


图 1-4 最早的留声机

10年后，爱迪生又把留声机上的大圆筒和小曲柄改成类似时钟发条的装置，由发动机带动一个薄薄的蜡制大圆盘转动，从此以后留声机才广为普及。留声机的发明，让音频数据的记录成为可能。

在公元前400年前，墨子所著《墨经》中已有针孔成像的记载；在13世纪，欧洲也出现了利用针孔成像原理制成的映像暗箱，人们可以走进暗箱观赏映像或描绘景物。

但直到1822年，法国的涅普斯才在感光材料上制出了世界上第一张照片，不过当时成像不太清晰，而且需要8个小时的曝光。1826年，他又在涂有感光性沥青的锡基底版上，通过暗箱拍摄了一张照片。

1839年，法国的达盖尔制成了第一台实用的银版照相机，它是由两个木箱组成，把

一个木箱插入另一个木箱中进行调焦，用镜头盖作为快门，来控制长达三十分钟的曝光时间，从而拍摄出清晰的图像，最终实现了静止图像数据的记录（图 1-5）。

1874 年，法国的朱尔·让桑发明了一种摄影机。他将感光胶片卷绕在带齿的供片盘上，在一个钟摆机构的控制下，供片盘在圆形供片盒内做间歇供片运动，同时钟摆机构带动快门旋转，每当胶片停下时，快门开启曝光。让桑将这种相机与一架望远镜相接，能以每秒一张的速度拍下行星运动的一组照片。让桑将其命名为摄影枪，这就是现代电影摄影机+Y6R 的始祖。摄影机（图 1-6）的发明，使得运动图像数据的记录成为可能。

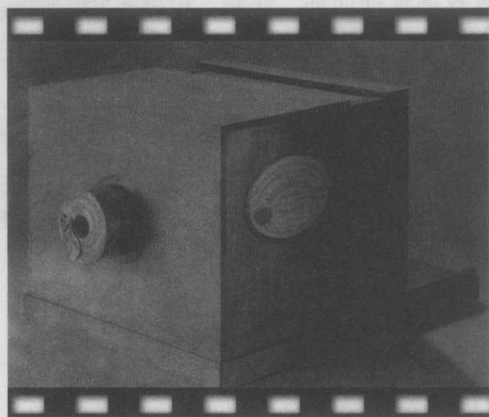


图 1-5 最古老的照相机

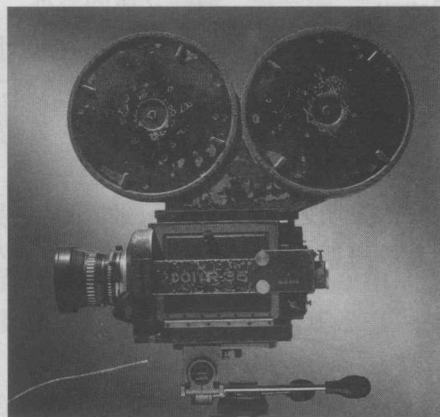


图 1-6 早期的摄影机

1946 年 2 月 14 日，由美国军方定制的世界第一台电子计算机“电子数字积分计算机”（Electronic Numerical and Calculator, ENIAC）在美国宾夕法尼亚大学问世，这表明电子计算机时代的到来。从此，人类与数据的关系进入了第二个时代（图 1-7）。

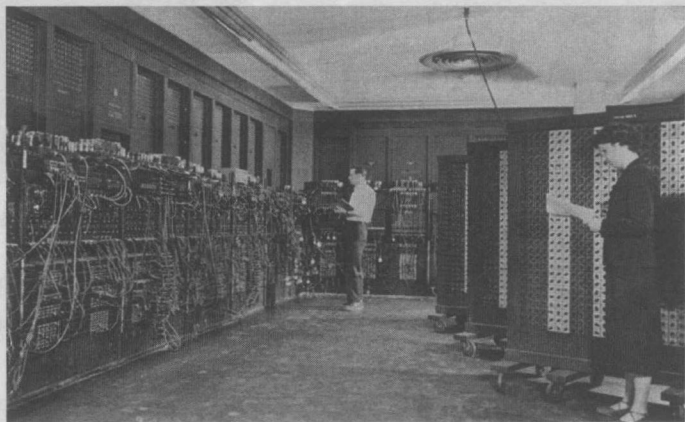


图 1-7 世界上第一台电子计算机

计算机的飞速发展，为数据的存储和处理模式带来了巨大的变革。以往的数据需要存储在纸张、胶片、磁带等介质上，一方面，数据的存储无法进行压缩，另一方面，介质的存储需要占用大量的空间。而计算机的发明，从本质上改变了这一点。数据可以通

过多种算法进行压缩。而且随着半导体工业的发展,存储能力不断增强,数据所需要的存储实体空间也在不断缩小。如今,一块小小的优盘就可以存储 GB 量级的数据,为我们节约了大量的数据存储空间。

随着计算机技术的发展,数据的处理能力也在不断提升。在计算机发明以前,数据都是通过人工的方式来进行处理。而有了计算机的帮助,通过各种各样的计算方式和统计软件,我们可以快速地处理数据。根据最新的统计,目前世界上最快的计算机——中国制造的天河2号(图1-8),处理速率已经达到了每秒钟进行22.86千万亿次浮点操作的水平。



图1-8 当前世界上最快的计算机——天河2号

互联网的出现,使人类与数据之间的关系进入到第三个时代的标志(图1-9)。最早的网络,是由美国国防部高级研究计划局(ARPA)建立的。现代计算机网络中的很多概念和方法,如分组交换技术都来自于ARPANET。ARPANET不仅进行了租用线互联的分组交换技术研究,而且做了无线、卫星网的分组交换技术研究,其结果就是加速了TCP/IP的问世。

1977—1979年,ARPANET推出了TCP/IP体系结构和协议。1980年前后,ARPANET上的所有计算机开始了TCP/IP协议的转换工作,并以ARPANET为主干网建立了初期的Internet。到1983年时,ARPANET的全部计算机完成了向TCP/IP的转换,并在UNIX(BSD 4.1)上实现了TCP/IP。到1984年时,美国国家科学基金会NSF规划建立了13个国家超级计算中心及国家教育科技网,随之替代了ARPANET的骨干地位。1988年,Internet开始对外开放。到了1991年6月,在联通Internet的计算机中,商业用户首次超过了学术界用户,这是Internet发展史上的一个里程碑,从此Internet的成长速度一发不可收拾。

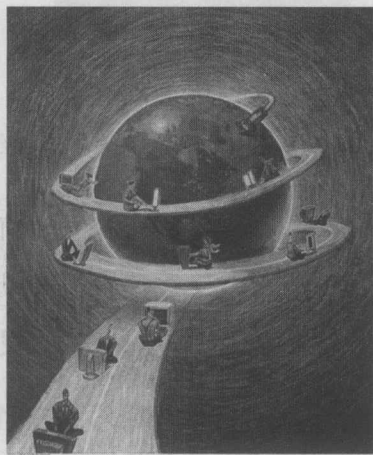


图1-9 网络与数据

互联网的精神就在于“开放、分享、平等、合作”。网络的出现,让人与人之间的距

离变得越来越短，地球村的概念也随之产生。通过网络，我们可以越洋对话，可以浏览海量的数据，可以实时地关注国际上最新的事件。网络让数据的产生和共享进入了一个崭新的时代。

网络时代的来临，造就了数据的大爆炸。据统计，2012年年底，有超过6 000万用户，通过社交网站 Facebook 发布了超过300亿条的新内容；游戏商 Zynga 每天要处理超过1 PB 容量的玩家数据；每天通过视频网站 YouTube 被浏览的视频量大约为20亿次；每个月通过微博 Twitter 所进行的搜索量会达到320亿次。让我们感受一下，这是多么庞大的数据量。

通过传感器网络搜集的数据又是另一大来源。所谓传感器网络，就是由大量部署在作用区域内的、具有无线通信与计算能力的微小传感器节点，通过自组织的方式所构成的，能根据环境自主完成指定任务的分布式智能化网络系统。

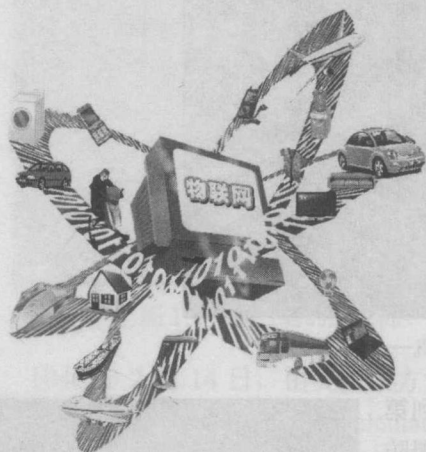


图 1-10 物联网

传感器网络综合了多种先进技术，如传感器技术、嵌入式计算技术、现代网络及无线通信技术、分布式信息处理技术，等等。它能够通过各种集成化的微型传感器协作，来实时监测、感知和采集各种环境或监测对象的信息，并可通过嵌入式系统对信息进行处理，并且通过无线网络将所感知到的信息传送到用户终端。

利用传感器网络，通过感知识别技术，让物品“开口说话、发布信息”，融合物理世界和信息世界，便可以建立物联网（图 1-10）。物联网的“触手”是位于感知识别层的大量信息生成设备，包

括 RFID、传感器网络、定位系统等。传感器网络所感知的数据正是物联网海量信息的重要来源之一。

互联网和物联网，正是我们数据收集来源的两大重要渠道，推动了大数据时代的来临。

1.2.3 数据、信息与知识

数据、信息与知识，这三个概念，在后面的学习中会多次出现。在使用这三者时，往往会存在一些概念上的交叠，容易混淆，在这里先做一下区分。

这三者之间最主要的区别是所考虑的抽象层次不同。数据是最低层次的抽象，信息次之，知识则是最高层次的抽象。数据是原始的、零散的，数据本身是没有意义的，数据经过了处理依然是数据，只有经过解释和理解才有意义。从数据抽象到信息的过程，就是对数据解读和释义的过程。

我们对数据进行解释和理解之后，才可以从数据中提取出有用的信息。对信息进行整合和呈现，则能够获得知识。例如，世界第一高峰珠穆朗玛峰的高度 8 844.43 m，可