

# 大数据分析 与 数据挖掘

简祯富 许嘉裕 / 编著

DATA

清华大学出版社

# 大数据分析 与 数据挖掘

简祯富 许嘉裕 / 编著

清华大学出版社  
北京

## 内 容 简 介

随着移动通信和行动装置普及、物联网和网络发展,以及云端技术的不断进步,现今数据产生、搜集和储存方式比以往更为方便。数据挖掘与大数据分析可以从海量数据中,找到值得参考的样型或规则,转换成有价值的信息、洞察或知识,创造更多新价值。

本书主要介绍数据挖掘与大数据分析的理论方法与实践应用,并加入丰富的实务案例介绍,具体说明如何应用数据挖掘与大数据分析技术以解决真实问题,深入浅出地剖析从数据中掏金的秘诀。全书共分为13章,内容涵盖数据挖掘基本概念与数据准备、数据挖掘的方法与实证、数据挖掘的进阶运用;书中也提供R语言与编程实例辅以说明,使读者更能融会贯通地应用数据挖掘方法,进而提升大数据分析和数字决策能力。

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

大数据分析数据挖掘/简祯富,许嘉裕编著. --北京:清华大学出版社,2016  
ISBN 978-7-302-42425-3

I. ①大… II. ①简… ②许… III. ① 统计数据—统计分析 ②数据采集 IV. ①O212.1 ②TP274

中国版本图书馆CIP数据核字(2015)第306758号

责任编辑:冯 昕  
封面设计:张京京  
责任校对:王淑云  
责任印制:宋 林

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)

印装者:北京鑫海金澳胶印有限公司

经 销:全国新华书店

开 本:185mm×260mm 印 张:23

字 数:560千字

版 次:2016年3月第1版

印 次:2016年3月第1次印刷

印 数:1~2000

定 价:49.00元

产品编号:065663-01

### “为大于其细，行远必自迩！”

1992年我到美国威斯康星大学麦迪逊分校(UW-Madison)攻读决策科学与作业研究博士时,发现我在新竹“清华大学”念的概率、统计、实验设计和统计方法等课程的教科书作者竟然都是麦迪逊的教授,所以选择统计作为副修;另一方面,我又在麦迪逊的医疗系统研究分析中心(Center for Health Systems Research and Analysis, CHSRA)担任研究助理,参与由Gustafson教授领导的大型研究团队发展的“综合医疗促进支持系统”(Comprehensive Health Enhancement Support System, CHESS),计划的目的是借着提供信息(information)、转介服务(referral to service providers)、决策支持(decision support)和社会援助(social support)等方式,帮助面对疾病和健康危机的人(如癌症和艾滋病患者)及其亲友取得相关信息、寻求可利用的资源、分析决策,以及社群服务和互相扶持等。我的主要工作是分析系统所搜集的使用数据和用户填写的问卷调查数据等,并在每周研究团队的定期会议上进行汇报,通过各种可能的分析和数据探索,以证明CHESS的效益。因为我的指导教师当时只是团队中的助理教授,所以我特别卖力分析,生怕工作不保就没有奖学金了。有一天,研究团队的一位成员在会议后告诉我说,我做的工作好像“数据挖掘”(data mining),他认为数据挖掘的方法将来可能会超越统计,虽然当时我觉得怎么可能有一种最近才发展的方法,可以超越已有几百年根基的统计学,但也让我注意到数据挖掘这个研究领域。

1996年我回到新竹“清华大学”任教,即成立“决策分析研究室”(Decision Analysis Laboratory, DALab),和研究伙伴与学生们包括本书共同作者许嘉裕博士一起投入决策分析、数据挖掘和优化的研究和实践工作,并通过产学合作计划作研究,然而却苦无合适的教材训练学生,特别是结合实际案例的课本,因此就持续借着整理产学合作研究成果、撰写期刊论文和指导学生论文之机,准备撰写教科书的基础材料。数据挖掘和大数据分析是方法论,也是实证推导模式(empirically derived model),因此必须结合方法发展与实证研究以检验研究效度。决策分析研究室研究团队与台积电、旺宏、合达电、联发科、广达电脑、创意电子、晶元光电、采钰、关东鑫林、茂迪、普生、力晶、世界先进等公司建立双赢的产学合作机制,做到学术研究贡献能够接连获奖,而实际效益能够达到合作厂商产业化的要求,作为更深一层理论研究的基础;更有幸从2005年借调台积电三年,实际应用所发展的分析方法在企业营运中,领导研究室的学生们和工业工程处同仁们一起推动台积电“IE十大建设”并发展相关的分析技术和数字决策系统,提供数字化系统化之决策依据,而从中得到产业导师宝贵的指导和回馈,也累积实战的经验和心得;进而执行台湾“科技部”“IC产业同盟”(Semiconductor Technologies Empowerment Partners Consortium, STEP Consortium)暨

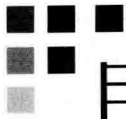
深耕工业基础技术计划,并成立“清华-台积电卓越制造中心”(NTHU-TSMC Center for Manufacturing Excellence),把累积多年的实证及大数据分析技术,推广到半导体供应链上、下游和其他高科技产业,借此提升产业的决策分析和智能制造能力;并通过主办“清华IC学堂”“半导体大数据分析竞赛”及产学合作成果发表研讨会等活动,培养具备跨界创新、团队合作能力的“资料科学家”。因此,本书在编撰过程中一再修改更新,希望一方面能深入介绍数据挖掘与大数据分析的基础方法和工具,另一方面则通过跨领域的实际案例和范例程序,以具体培养结合理论与实务的决策科学家。

非常感谢新竹“清华大学”和元智大学的良好学术研究环境和科学园区的地利人和,使我们可以结合理论与实务,从产业大数据和具体问题的实证中发展适用的方法、检验所学,再进而导向更深一层的研究。随着问题的广度和复杂度以及合作伙伴的阶层和领域而不断成长,这一路走来,虽然整个研究团队一直秉持自强不息、行胜于言的精神努力提升,但也得力于产业先进和合作伙伴们的提携协助和计划执行过程中的指导,因此要感谢的人非常多,希望借着本书的出版能使更多读者从中得到启发和实际的帮助,以造福社会和产业,也算是间接回报所有关心和帮助我们的人。尽管本书经过长期的准备,但完稿阶段所花费的心力远远超过预期,特别感谢专任助理梁婉玲编辑汇总的工作和与出版社的联络,减少本书错误的可能,以及决策分析研究室同学们一起打拼完成各项研究计划,这也是本书各案例的论文均引用完整作者名单的原因;也感谢在“数据挖掘”课程教学中每位互动的学生,让我们得到教学相长和调整教材的回馈建议。本书自2014年在台湾出版以来,引发学术界和产业界的广泛回响,成为多所大学和各大企业的指定教材。感谢北京清华大学出版社理工分社张秋玲社长和冯昕主任的支持,将全书重新编辑改版,去芜存菁,并增添一章全新章节,使内容更加丰富完整。然而,本书疏漏之处在所难免,盼诸位领导和前辈,不吝赐教,以提升大数据分析 and 数字决策能力。

简祯富 许嘉裕 谨识

IC产业同盟,2015冬





## 第 1 篇 大数据分析 with 数据挖掘导论

<b>第 1 章 大数据分析 with 数据挖掘概论</b> .....	<b>3</b>
1.1 前言 .....	3
1.2 大数据分析的应用 .....	6
1.3 数据挖掘与数字决策 .....	8
1.4 数据挖掘和大数据分析架构与步骤 .....	9
1.4.1 问题定义与架构 .....	10
1.4.2 数据准备 .....	11
1.4.3 建立挖掘模式 .....	11
1.4.4 结果解释与评估 .....	12
1.5 数据挖掘的问题类型.....	13
1.5.1 分类 .....	13
1.5.2 预测 .....	13
1.5.3 聚类 .....	14
1.5.4 关联规则 .....	14
1.6 数据挖掘模式.....	14
1.7 结论.....	15
1.8 本书架构.....	17
问题与讨论 .....	17
<b>第 2 章 数据与数据准备</b> .....	<b>19</b>
2.1 数据取得.....	20
2.2 大数据分析的基础: Hadoop .....	22
2.2.1 Hadoop 架构 .....	22
2.2.2 Hadoop 分布式文件系统 .....	23
2.2.3 MapReduce .....	24
2.3 数据类型.....	25
2.4 数据尺度.....	26
2.5 数据检查.....	28
2.6 数据探索与可视化.....	29



2.7	数据整合与清理	32
2.8	数据转换	36
2.8.1	数据数值转换	36
2.8.2	数据属性转换	37
2.9	数据归约	38
2.9.1	数据维度归约	38
2.9.2	数据数值归约	44
2.10	数据分割	46
2.11	应用实例——半导体厂制造技术员人力资源管理质量提升	47
2.11.1	案例背景	47
2.11.2	数据准备	47
2.12	结论	50
	问题与讨论	51

## 第 2 篇 数据挖掘方法与实证

第 3 章	关联规则	55
3.1	关联规则的定义与说明	55
3.2	关联规则的衡量指针	57
3.3	关联规则的类型	59
3.4	关联规则算法	60
3.4.1	Apriori 算法	62
3.4.2	Partition 算法	65
3.4.3	DHP 算法	66
3.4.4	MSApriori 算法	68
3.4.5	FP-Growth 算法	70
3.5	多维度关联规则	75
3.6	多阶层关联规则	76
3.7	关联规则的应用	79
3.8	R 语言与关联规则分析	79
3.9	应用实例——电力公司配电事故定位的研究	83
3.9.1	案例背景	83
3.9.2	数据准备	84
3.9.3	关联规则推导	85
3.10	结论	88
	问题与讨论	88
第 4 章	决策树分析	93
4.1	决策树的建构	93



4.1.1	数据准备 .....	94
4.1.2	决策树的分支准则 .....	96
4.1.3	决策树修剪 .....	104
4.1.4	规则提取 .....	106
4.2	决策树的算法 .....	107
4.2.1	CART .....	108
4.2.2	C4.5/C5.0 .....	108
4.2.3	CHAID .....	109
4.3	决策树分类模型评估 .....	110
4.4	R语言与决策树分析 .....	112
4.4.1	CART决策树分析 .....	112
4.4.2	C5.0决策树分析 .....	114
4.4.3	CHAID决策树分析 .....	115
4.5	应用实例——建构cDNA生物芯片的数据挖掘模式 .....	117
4.5.1	案例背景 .....	117
4.5.2	数据准备 .....	117
4.5.3	生物芯片数据的决策树构建 .....	118
4.5.4	规则解释与评估 .....	119
4.6	结论 .....	120
	问题与讨论 .....	120
<b>第5章</b>	<b>人工神经网络 .....</b>	<b>127</b>
5.1	人工神经网络的基本结构 .....	130
5.2	网络学习法 .....	132
5.3	反向传播人工神经网络 .....	134
5.3.1	网络架构 .....	134
5.3.2	学习算法 .....	136
5.3.3	反向传播人工神经网络步骤 .....	137
5.3.4	反向传播人工神经网络范例 .....	138
5.4	自组织映射网络 .....	139
5.4.1	网络架构 .....	140
5.4.2	学习算法 .....	142
5.4.3	SOM人工神经网络步骤 .....	143
5.4.4	自组织映射图网络范例 .....	143
5.5	自适应共振理论人工神经网络 .....	146
5.5.1	网络架构 .....	147
5.5.2	ART1网络算法 .....	148
5.5.3	适应性共振网络范例 .....	150
5.6	R语言与人工神经网络 .....	152





5.6.1	反向传播人工神经网络	152
5.6.2	自组织映射网络	154
5.6.3	自适应共振理论人工神经网络	155
5.7	应用实例——半导体生产周期时间预测与管控	158
5.7.1	案例简介	158
5.7.2	数据分群	159
5.7.3	数据配适与预测	160
5.7.4	信息整合与敏感度分析	161
5.7.5	案例小结	162
5.8	结论	163
	问题与讨论	163
<b>第6章</b>	<b>聚类分析</b>	<b>165</b>
6.1	聚类分析法简介	165
6.1.1	聚类分析的阶段	166
6.1.2	相似度的衡量	166
6.1.3	聚类分析方法	169
6.2	层次聚类分析法	170
6.3	划分聚类分析法	174
6.3.1	K 平均法	174
6.3.2	K 中心点法	176
6.4	以密度为基础的分群算法	179
6.5	以模式为基础的分群算法	181
6.5.1	期望最大化算法	181
6.5.2	自组织映射图网络	182
6.6	R 语言与聚类分析	182
6.7	应用实例——黄光机台聚类分析	184
6.7.1	案例简介	184
6.7.2	验证两阶段分群算法	185
6.7.3	案例小结	187
6.8	结论	187
	问题与讨论	188
<b>第7章</b>	<b>朴素贝叶斯分类法与贝叶斯网络</b>	<b>190</b>
7.1	贝叶斯定理	190
7.2	朴素贝叶斯分类法	192
7.3	贝叶斯网络	196
7.3.1	贝叶斯网络的理论基础	196
7.3.2	贝叶斯网络的不一致性修正	201

7.4	R 语言与贝叶斯分类 .....	203
7.5	应用实例——电力公司馈线事故定位系统 .....	207
7.5.1	案例简介与问题架构 .....	207
7.5.2	数据整理与贝叶斯网络图构建 .....	208
7.5.3	给定贝叶斯推理网络的参数 .....	209
7.5.4	验证贝叶斯推理网络 .....	210
7.5.5	案例小结 .....	210
7.6	结论 .....	211
	问题与讨论 .....	211
<b>第 8 章</b>	<b>粗糙集理论 .....</b>	<b>215</b>
8.1	粗糙集理论 .....	215
8.2	粗糙集理论基本概念 .....	215
8.2.1	信息系统与决策表 .....	216
8.2.2	等价关系 .....	216
8.2.3	近似空间 .....	217
8.2.4	近似集合的准确率 .....	218
8.2.5	分类的准确率与属性相依程度 .....	219
8.2.6	简化 .....	219
8.3	粗糙集理论产生分类规则 .....	222
8.4	粗糙集理论与其他分类方法的比较 .....	223
8.5	R 语言与粗糙集理论 .....	224
8.5.1	决策表与等价关系 .....	225
8.5.2	近似空间 .....	225
8.5.3	简化与规则推演 .....	226
8.6	应用实例——TFT-LCD 数组事故诊断 .....	227
8.6.1	案例简介 .....	227
8.6.2	分析过程 .....	227
8.6.3	案例小结 .....	230
8.7	结论 .....	231
	问题与讨论 .....	231
<b>第 9 章</b>	<b>预测与时间数据分析 .....</b>	<b>234</b>
9.1	回归分析 .....	234
9.1.1	回归分析基本介绍 .....	234
9.1.2	参数估计 .....	237
9.1.3	回归模型解释与评估 .....	237
9.1.4	多重回归分析 .....	239
9.1.5	共线性 .....	239



9.2	逻辑回归 .....	240
9.2.1	概率与胜算 .....	240
9.2.2	逻辑回归模式 .....	240
9.3	时间序列分析 .....	242
9.4	时间数据的分析步骤 .....	243
9.5	模式选择与建立 .....	244
9.5.1	时间序列平滑法 .....	246
9.5.2	平稳型时间序列 .....	247
9.5.3	无定向型时间序列 .....	251
9.5.4	趋势型、季节型与介入事件型时间序列 .....	252
9.6	阶次选取与参数估计 .....	254
9.7	模式评估 .....	255
9.7.1	拟合优度检定 .....	255
9.7.2	预测误差衡量 .....	256
9.8	R语言与时间数据分析 .....	257
9.9	应用实例——半导体光罩需求预测 .....	261
9.9.1	案例简介与问题架构 .....	261
9.9.2	数据准备与数据处理 .....	261
9.9.3	需求波动侦测分析过程 .....	262
9.9.4	案例小结 .....	263
9.10	结论 .....	264
	问题与讨论 .....	265
<b>第 10 章</b>	<b>集成学习与支持向量机 .....</b>	<b>268</b>
10.1	集成学习 .....	268
10.1.1	Bagging .....	268
10.1.2	Boosting .....	269
10.2	支持向量机 .....	272
10.2.1	可区分情况(separable case) .....	272
10.2.2	不可分状况(non-separable case) .....	274
10.2.3	非线性分类 .....	275
10.3	R语言与随机森林集成学习模型 .....	276
10.3.1	利用随机森林进行分类 .....	276
10.3.2	利用随机森林评估变量重要性 .....	277
10.4	结论 .....	278
	问题与讨论 .....	278

**第3篇 数据挖掘进阶运用**

<b>第 11 章 商业智能</b> .....	<b>281</b>
11.1 商业智能概述 .....	281
11.2 应用实例——交通信息预测 .....	283
11.3 个案研究——人力资源数据挖掘 .....	283
11.3.1 案例说明 .....	283
11.3.2 分析过程 .....	284
11.3.3 案例小结 .....	291
11.4 应用实例——机票价格预测 .....	292
11.5 个案研究——产品需求预测 .....	292
11.5.1 半导体产品需求预测架构 .....	292
11.5.2 分析过程 .....	297
11.5.3 案例小结 .....	303
11.6 结论 .....	303
问题与讨论 .....	304
<b>第 12 章 制造智能</b> .....	<b>305</b>
12.1 序言 .....	305
12.2 WAT 参数特征提取与关联分析 .....	307
12.2.1 案例说明 .....	307
12.2.2 分析过程 .....	308
12.2.3 案例小结 .....	312
12.3 半导体 CP 测试数据挖掘与晶圆图样型分类 .....	312
12.3.1 案例背景 .....	312
12.3.2 分析过程 .....	313
12.3.3 案例小结 .....	318
12.4 低良率事故诊断与制程关联分析 .....	318
12.4.1 案例说明 .....	318
12.4.2 分析过程 .....	319
12.4.3 案例小结 .....	323
12.5 半导体制造管理的数据挖掘 .....	324
12.5.1 案例背景 .....	324
12.5.2 分析过程 .....	324
12.5.3 案例小结 .....	329
12.6 结论 .....	330
问题与讨论 .....	331



<b>第 13 章 数字决策及商业分析与优化</b> .....	<b>332</b>
13.1 决策信息系统 .....	332
13.1.1 决策信息系统 .....	332
13.1.2 决策信息系统的架构 .....	333
13.1.3 应用实例——电性测试机台维修的决策支持系统 .....	334
13.2 商业分析与优化 .....	339
13.2.1 商业分析与优化 .....	339
13.2.2 商业分析与优化的基本要素 .....	340
13.2.3 商业分析与优化的应用 .....	341
13.3 数字决策 .....	342
13.4 结论 .....	343
问题与讨论 .....	344
<b>参考文献</b> .....	<b>345</b>



第 1 篇

---

# 大数据分析 with 数据挖掘导论



## 第 1 章

# 大数据分析 with 数据挖掘概论

## 1.1 前言

随着信息科技的进步和网络的发达、计算机运算能力的增强以及数据搜集与储存技术持续改进的影响,大幅改变数据的分析和应用方式,“大数据分析”(big data analytics)和**数据挖掘(data mining)**可以发掘先前未知且潜在有用的信息样型(patterns)或规则(rules),进而转化为有价值的信息或知识,帮助决策者迅速做出适当的决策,是现代企业重要的竞争优势。

由于自动化的生产环境、智能手机的普及、电子商务的发展、物联网的建立以及社交网络的发达,现在多数人都可以不受时空地点限制地上网,浏览社交网络,在网络上聊天、购物,以及实时收看与查询最新的新闻报道与文章等,也可以用来管理远程的生产和服务系统。当你在微博上打卡点赞、收发电子邮件、到便利商店购买零食、搭乘大众交通工具、经过停车场利用信用卡缴费时,这些日常生活中的习惯与动作,随时随地正透过网络记录,快速累积成巨量数据或大数据。过去对商品的评价主要是通过口口相传,而现在则是借由在线文章发表,由社交网络快速扩散,这意味着网络经营的重要性已开始逐渐大过实体经营,大数据分析正引领着数字决策并带来新商机。

“数据”在经济学中属于非竞争性的商品,其与物质性的东西(例如食物、车等)不同,并不会因为使用次数增加而降低价值或造成耗损。因此,零售业者累积的事务数据可以一再使用,根据不同目的提取不同的数据,或运用于不同的目标对象上(Mayer-Schonberger & Cukier, 2013)。除了传统的统计分析和数据挖掘外,大数据分析技术和应用正改变我们的生产方式、服务系统和生活形态。

每一秒,一间大型医院会增加 12 万笔健康相关的生理数据;每一分钟,YouTube 网站会接收到民众上传总长达 72 小时的视频;每一天,一家银行的信用卡交易次数达 500 万笔。时间分秒走过的同时,大量数据也随时都在快速累积,如图 1.1 所示。而在全世界数兆个传感器、超过五亿部智能手机、十亿台计算机上,每一天不断运作所产生的数据量估计高达 25 亿 GB(胡世忠, 2013)。科技研究公司 IDC 更预估,到 2020 年全球数据量将累积达 40 000 ZB(Gantz & Reinsel, 2012),数据储存单位如表 1.1。



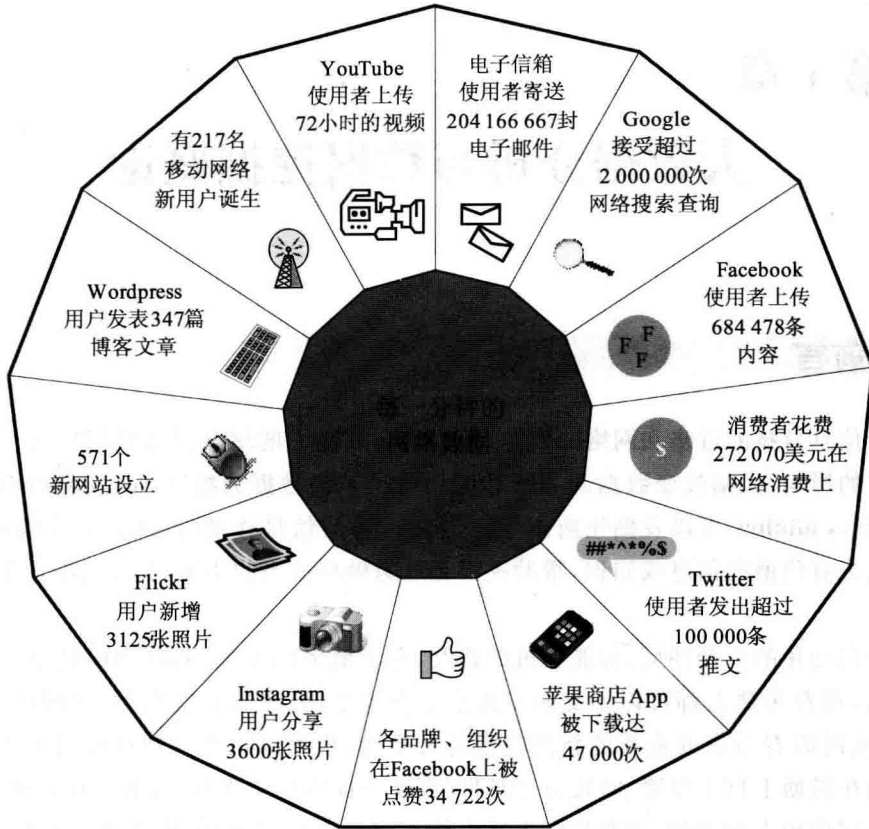


图 1.1 持续增加的大数据(胡世忠, 2013)

表 1.1 数据的储存单位

储存单位/B	文件储存单位
Kilobyte (KB)	1 KB=1024 B=2 <sup>10</sup> B
Megabyte (MB)	1 MB=1024 KB=2 <sup>20</sup> B
Gigabyte (GB)	1 GB=1024 MB=2 <sup>30</sup> B
Terabyte (TB)	1 TB=1024 GB=2 <sup>40</sup> B
Perabyte (PB)	1 PB=1024 TB=2 <sup>50</sup> B
Exabyte (EB)	1 EB=1024 PB=2 <sup>60</sup> B
Zettebyte (ZB)	1 ZB=1024 EB=2 <sup>70</sup> B
Yottabyte (YB)	1 YB=1024 ZB=2 <sup>80</sup> B

大量的传感器与电子卷标置入到日常生活的电子设备中,例如手机、监控摄影机、环境温度传感器、水电天然气表等,随时感测人们的生活动态。例如,电力公司为了节省能源,开发的智能电表和智能电网即装置了大量的传感器,24小时不间断地测量与传输终端顾客的电力使用信息。对终端顾客而言,智能电表能实时显示家中的用电量,协助用户调整用电习惯。对电力公司而言,则可透过实时用电量的监控,掌握电网供电状态,当耗电量可能超过