



第8章 数据可视化

(PPT版本号: 2022--V1)





CONTENTS

- 8.1 可视化概述
- 8.2 可视化图标
- 8.3 可视化工具
- 8.4 可视化典型案例





8.1

可视化概述

8.1.1 什么是数据可视化

8.1.2 可视化的发展历程

8.1.3 可视化的重要作用



8.1.1 什么是数据可视化

- 数据可视化是指将大型数据集中的数据以图形图像形式表示，并利用数据分析和开发工具发现其中未知信息的处理过程
- 数据可视化技术的基本思想是将数据库中每一个数据项作为单个图元素表示，大量的数据集构成数据图像，同时将数据的各个属性值以多维数据的形式表示，可以从不同的维度观察数据，从而对数据进行更深入的分析

8.1.2 可视化的发展历程

霍乱地图分析了霍乱患者分布与水井分布之间的关系，发现在有一口井的供水范围内患者明显偏多，据此找到了霍乱爆发的根源是一个被污染的水泵



图 反映霍乱患者分布与水井分布的地图

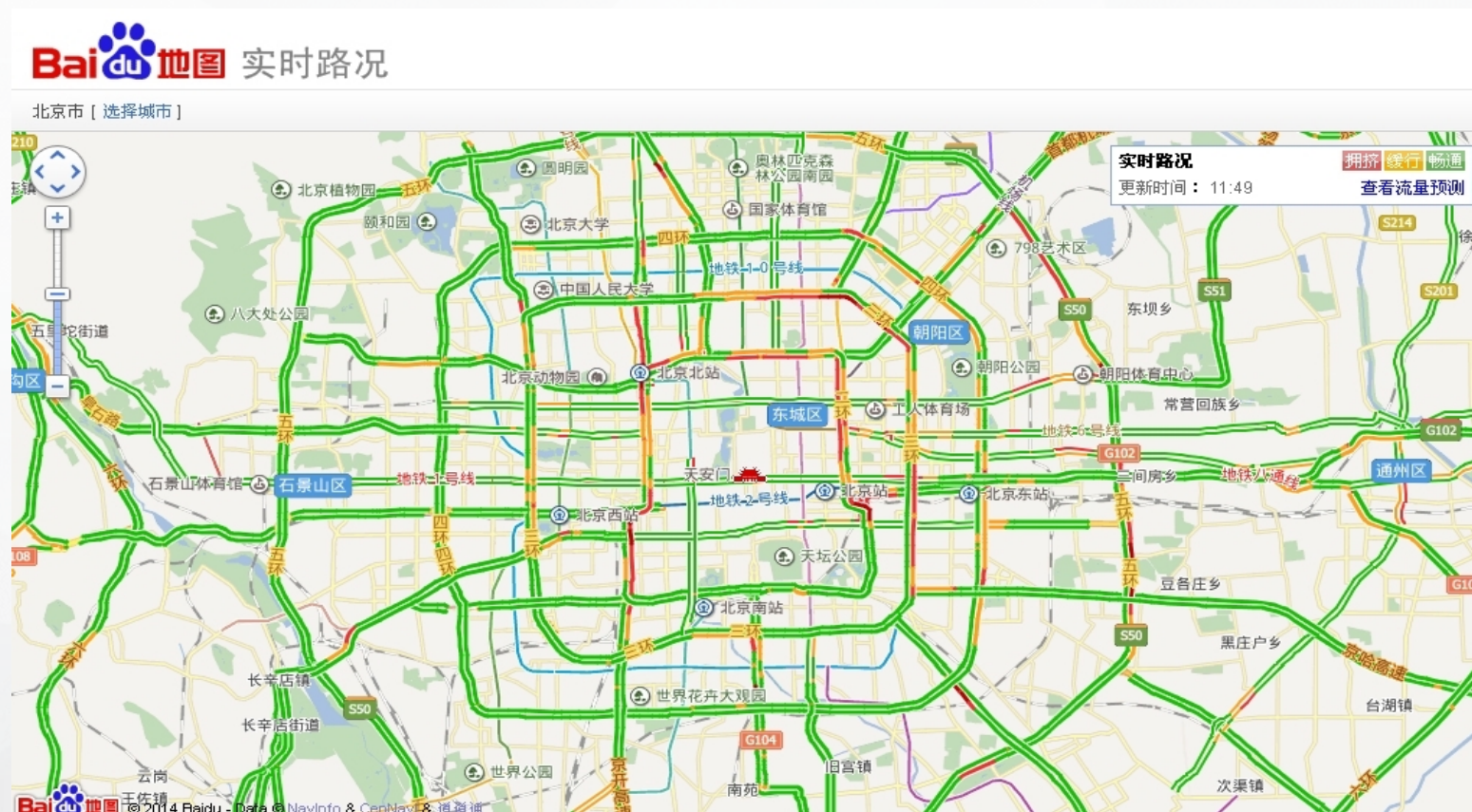
8.1.2 可视化的发展历程

- 20世纪50年代，随着计算机的出现和计算机图形学的发展，人们可以利用计算机技术在电脑屏幕上绘制出各种图形图表，可视化技术开启了全新的发展阶段。最初，可视化技术被大量应用于统计学领域，用来绘制统计图表，比如圆环图、柱状图和饼图、直方图、时间序列图、等高线图、散点图等，后来，又逐步应用于地理信息系统、数据挖掘分析、商务智能工具等，有效促进了人类对不同类型数据的分析与理解
- 随着大数据时代的到来，每时每刻都有海量数据在不断生成，需要我们对数据进行及时、全面、快速、准确的分析，呈现数据背后的价值，这就更需要可视化技术协助我们更好地理解和分析数据，可视化成为大数据分析最后的一环和对用户而言最重要的一环

8.1.2 可视化的发展历程

在大数据时代，可视化技术可以支持实现多种不同的目标：

(1) 观测、跟踪数据



8.1.2 可视化的发展历程

(2) 分析数据

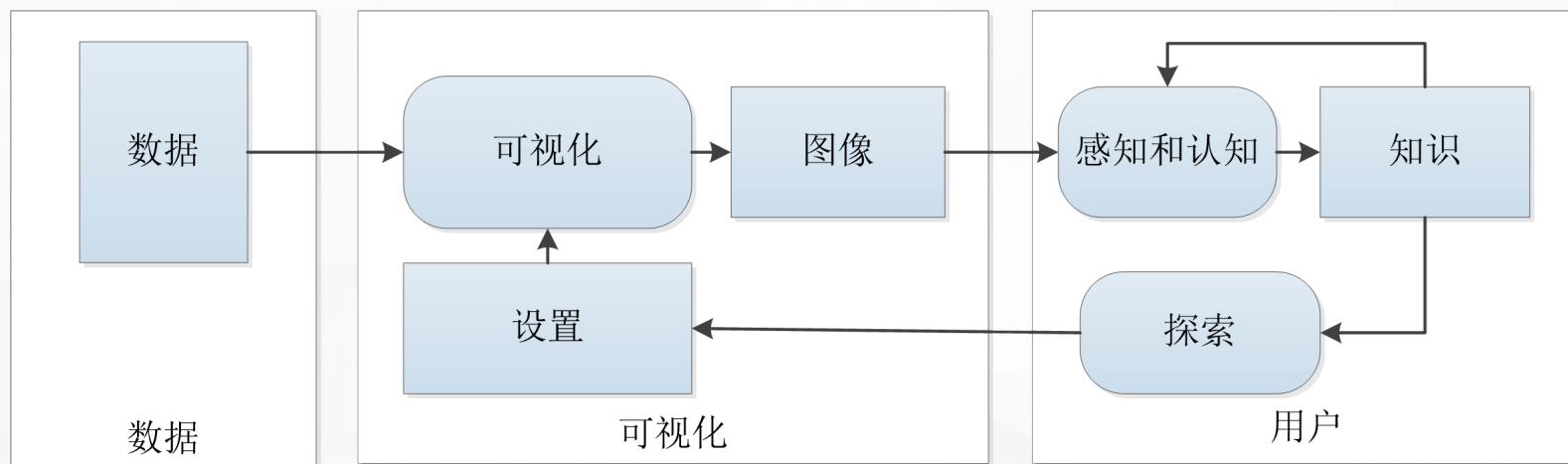


图 用户参与的可视化分析过程

8.1.2 可视化的发展历程

(3) 辅助理解数据



图 微软“人立方”展示的人物关系图

8.1.2 可视化的发展历程

(4) 增强数据吸引力



图 一个可视化的图表新闻实例



8.2

可视化图表



8.2可视化图表

表 最常用的统计图表类型及其应用场景

图表	维度	应用场景
柱状图	二维	指定一个分析轴进行数据大小的比较，只需比较其中一维
折线图	二维	按照时间序列分析数据的变化趋势，适用于较大的数据集
饼图	二维	指定一个分析轴进行所占比例的比较，只适用于反映部分与整体的关系
散点图	二维或三维	有两个维度需要比较
气泡图	三维或四维	其中只有两维能够精确辨识
雷达图	四维以上	数据点不超过6个

8.2可视化图表

除了上述常见的图表以外，数据可视化还可以使用其他图表，具体如下：

(1) 漏斗图。漏斗图适用于业务流程比较规范、周期长、环节多的流程分析，通过漏斗各环节业务数据的比较，能够直观地发现和说明问题所在。

(2) 树图。树图是一种流行的、利用包含关系表达层次化数据的可视化方法，它能将事物或现象分解成树枝状，因此又称“树型图”或“系统图”。树图就是把要实现的目的与需要采取的措施或手段，系统地展开，并绘制成图，以明确问题的重点，寻找最佳手段或措施。

(3) 热力图。以特殊高亮的形式显示访客热衷的页面区域和访客所在的地理区域的图示，它基于GIS坐标，用于显示人或物品的相对密度。

(4) 关系图。基于3D空间中的点线组合，再加以颜色、粗细等维度的修饰，适用于表征各节点之间的关系。

(5) 词云。通过形成“关键词云层”或“关键词渲染”，对网络文本中出现频率较高的“关键词”给予视觉上的突出。

(6) 桑基图。也被称为“桑基能量分流图”或“桑基能量平衡图”，它是一种特定类型的流程图，图中延伸的分支的宽度对应数据流量的大小，通常应用于能源、材料成分、金融等数据的可视化分析。

(7) 日历图。以日历为基本维度的、对单元格加以修饰的图表。



8.3 可视化工具

8.3.1入门级工具

8.3.2信息图表工具

8.3.3地图工具

8.3.4时间线工具

8.3.5高级分析工具



8.3.1 入门级工具

- Excel是微软公司的办公软件Office家族的系列软件之一，可以进行各种数据的处理、统计分析和辅助决策操作，已经广泛地应用于管理、统计、金融等领域

8.3.2 信息图表工具

信息图表是信息、数据、知识等的视觉化表达，它利用人脑对于图形信息相对于文字信息更容易理解的特点，更高效、直观、清晰地传递信息，在计算机科学、数学以及统计学领域有着广泛的应用。

1. Google Chart API

谷歌公司的制图服务接口Google Chart API，可以用来为统计数据并自动生成图片，该工具使用非常简单，不需要安装任何软件，可以通过浏览器在线查看统计图表。

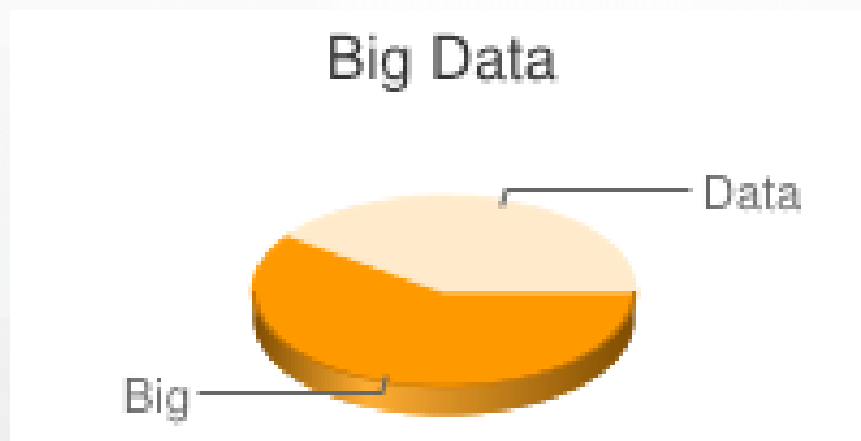


图 通过浏览器在线查看Google Chart统计图表

8.3.2 信息图表工具

2. ECharts

ECharts是由百度公司前端数据可视化团队研发的图表库，可以流畅地运行在PC和移动设备上，兼容当前绝大部分浏览器（IE8/9/10/11、Chrome、Firefox、Safari等），底层依赖轻量级的Canvas类库ZRender，可以提供直观、生动、可交互、可高度个性化定制的数据可视化图表。ECharts提供了非常丰富的图表类型，包括常规的折线图、柱状图、散点图、饼图、K线图，用于统计的盒形图，用于地理数据可视化的地图、热力图、线图，用于关系数据可视化的关系图、treemap，用于多维数据可视化的平行坐标，以及用于BI的漏斗图、仪表盘，并且支持图与图之间的混搭，能够满足用户绝大部分分析数据时的图表制作需求。

8.3.2 信息图表工具

3. D3

D3是最流行的可视化库之一，是一个用于网页作图、生成互动图形的JavaScript函数库，提供了一个D3对象，所有方法都通过这个对象调用。D3能够提供大量线性图和条形图之外的复杂图表样式，例如Voronoi图、树形图、圆形集群和单词云等（如图10-8所示）。



图 D3提供的可视化图表

8.3.2 信息图表工具

4. Tableau

Tableau是桌面系统中最简单的商业智能工具软件，更适合企业和部门进行日常数据报表和数据可视化分析工作。Tableau实现了数据运算与美观的图表的完美结合，用户只要将大量数据拖放到数字“画布”上，转眼间就能创建好各种图表。

5. 大数据魔镜

大数据魔镜是一款优秀的国产数据分析软件，它丰富的数据公式和算法可以让用户真正理解探索分析数据，用户只要通过一个直观的拖放界面就可创造交互式的图表和数据挖掘模型。

8.3.3 地图工具

- 地图工具在数据可视化中较为常见，它在展现数据基于空间或地理分布上有很强的表现力，可以直观地展现各分析指标的分布、区域等特征。当指标数据要表达的主题跟地域有关联时，就可以选择以地图作为大背景，从而帮助用户更加直观地了解整体的数据情况，同时也可以根据地理位置快速地定位到某一地区来查看详细数据。
- 右图就是以数据地图形式呈现的2008年世界各国GDP数据，图中，颜色越深的国家，其GDP越高。

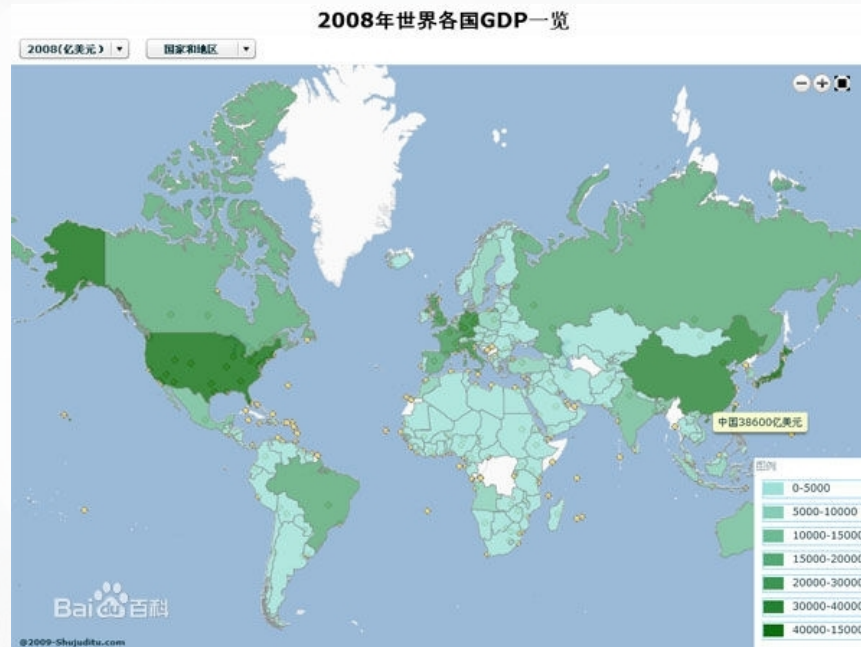


图 2008年世界各国GDP数据地图

8.3.3 地图工具

- **1. Google Fusion Tables**

Google Fusion Tables让一般使用者也可以轻松制作出专业的统计地图。该工具可以让数据表呈现为图表、图形和地图，从而帮助发现一些隐藏在数据背后的模式和趋势。

- **2. Modest Maps**

Modest Maps是一个小型、可扩展、交互式的免费库，提供了一套查看卫星地图的API，只有10KB大小，是目前最小的可用地图库，它也是一个开源项目，有强大的社区支持，是在网站中整合地图应用的理想选择。

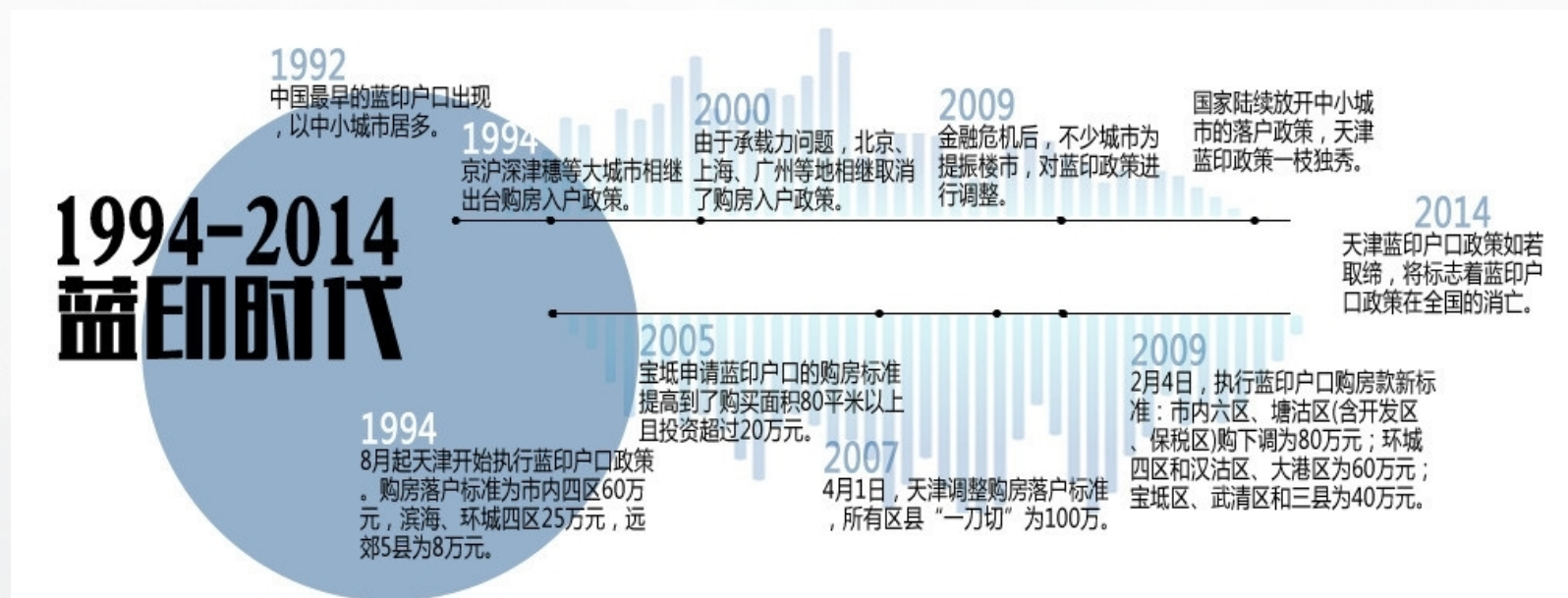
- **3. Leaflet**

Leaflet是一个小型化的地图框架，通过小型化和轻量化来满足移动网页的需要。

8.3.4 时间线工具

时间线是表现数据在时间维度的演变的有效方式，它通过互联网技术，依据时间顺序，把一方面或多方面的事件串联起来，形成相对完整的记录体系，再运用图文的形式呈现给用户。时间线可以运用于不同领域，最大的作用就是把过去的事物系统化、完整化、精确化。自2012年Facebook在F8大会上发布了以时间线格式组织内容的功能后，时间线工具在国内外社交网站中开始大面积流行。

图10-10显示了我国户籍制度在1994年到2014年间随时间的演变情况，它采用了时间线表示方法。



8.3.4 时间线工具

- **1. Timetoast**

Timetoast是在线创作基于时间轴事件记载服务的网站，提供个性化的时间线服务，可以用不同的时间线来记录你某个方面的发展历程、心理路程、进度过程等等。

Timetoast基于 flash 平台，可以在类似 flash时间轴上任意加入事件，定义每个事件的时间、名称、图像、描述，最终在时间轴上显示事件在时间序列上的发展，事件显示和切换十分流畅，随着鼠标点击可显示相关事件，操作简单。

- **2. Xtimeline**

Xtimeline 是一个免费的绘制时间线的在线工具网站，操作简便，用户通过添加事件日志的形式构建时间表，同时也可给日志配上相应的图表。不同于Timetoast的是，Xtimeline是一个社区类型的时间轴网站，其中加入了组群功能和更多的社会化因素，除了可以分享和评论时间轴外，还可以建立组群讨论所制作的时间轴。

8.3.5 高级分析工具

- **1. R**

R是属于GNU系统的一个自由、免费、源代码开放的软件，它是一个用于统计计算和统计制图的优秀工具，使用难度较高。R的功能包括数据存储和处理系统、数组运算工具（具有强大的向量、矩阵运算功能）、完整连贯的统计分析工具、优秀的统计制图功能、简便而强大的编程语言，可操纵数据的输入和输出，实现分支、循环以及用户可自定义功能等，通常用于大数据集的统计与分析。

- **2. Weka**

Weka是一款免费的、基于Java环境的、开源的机器学习以及数据挖掘软件，不但可以进行数据分析，还可以生成一些简单图表。

- **3. Gephi**

Gephi是一款比较特殊也很复杂的软件，主要用于社交图谱数据可视化分析，可以生成非常酷炫的可视化图形。

8.3.5 高级分析工具

4.Python

Python是一种面向对象的解释型计算机程序设计语言，由荷兰人吉多·范罗苏姆（Guido van Rossum）于1989年发明。Python是纯粹的自由软件，源代码和解释器CPython遵循GPL（GNU General Public License）协议。Python具有丰富和强大的库。它常被称为“胶水语言”，能够把用其他语言制作的各种模块（尤其是C/C++）很轻松地连接在一起。Python也是一种很好的可视化工具，可以开发出各种可视化效果图，Python可视化库可以大致分为：基于matplotlib的可视化库、基于JavaScript的可视化库、基于上述两者或其他组合功能的库。



8.4

可视化典型案例

8.4.1 全球黑客活动

8.4.2 互联网地图

8.4.3 编程语言之间的影响力关系图

8.4.4 百度迁徙

8.4.5 世界国家健康与财富之间的关系

8.4.6 3D可视化互联网地图APP



8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具

8.3.4 时间线工具



1.2

大数据时代

1.2.1 第三次信息化浪潮

1.2.2 信息科技为大数据时代提供技术支撑

1.2.3 数据产生方式的变革促成大数据时代的来临



1.2.1 第三次信息化浪潮

- 根据IBM前首席执行官郭士纳的观点，IT领域每隔十五年就会迎来一次重大变革

表1-1 三次信息化浪潮

信息化浪潮	发生时间	标志	解决问题	代表企业
第一次浪潮	1980年前后	个人计算机	信息处理	Intel、AMD、IBM、苹果、微软、联想、戴尔、惠普等
第二次浪潮	1995年前后	互联网	信息传输	雅虎、谷歌、阿里巴巴、百度、腾讯等
第三次浪潮	2010年前后	物联网、云计算和大数据	信息爆炸	将涌现出一批新的市场标杆企业

1.2.2 信息科技为大数据时代提供技术支撑

1. 存储设备容量不断增加

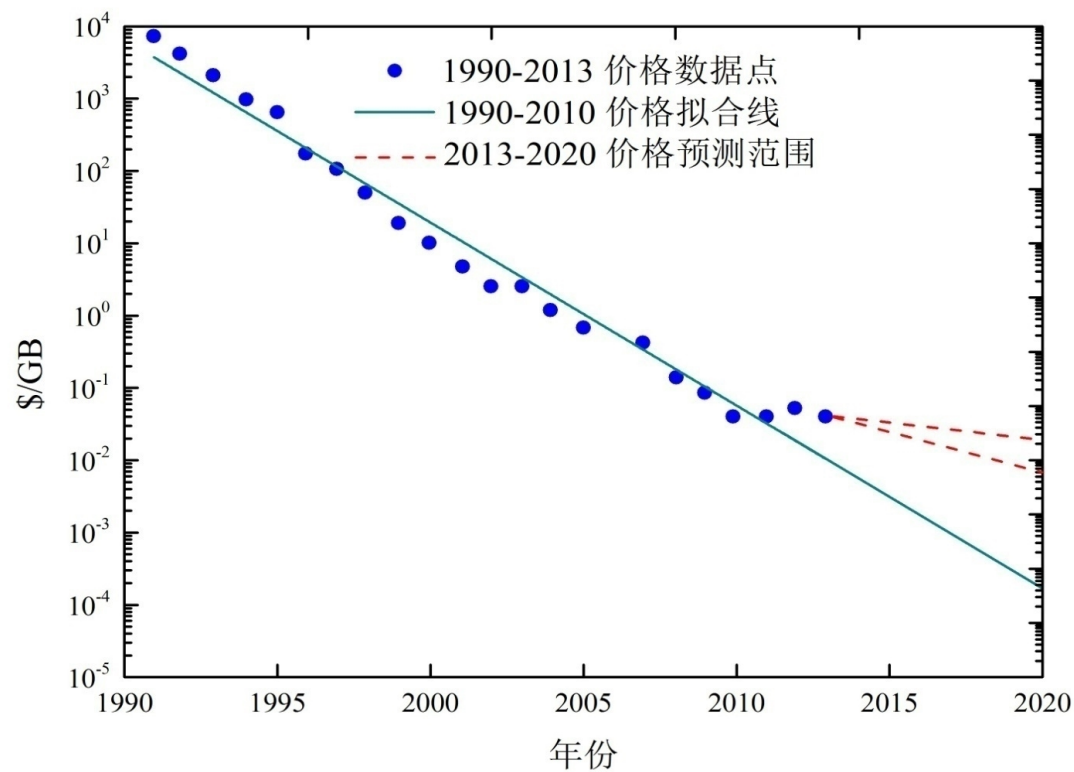


图 存储价格随时间变化情况

1.2.2 信息科技为大数据时代提供技术支撑

2. CPU处理能力大幅提升

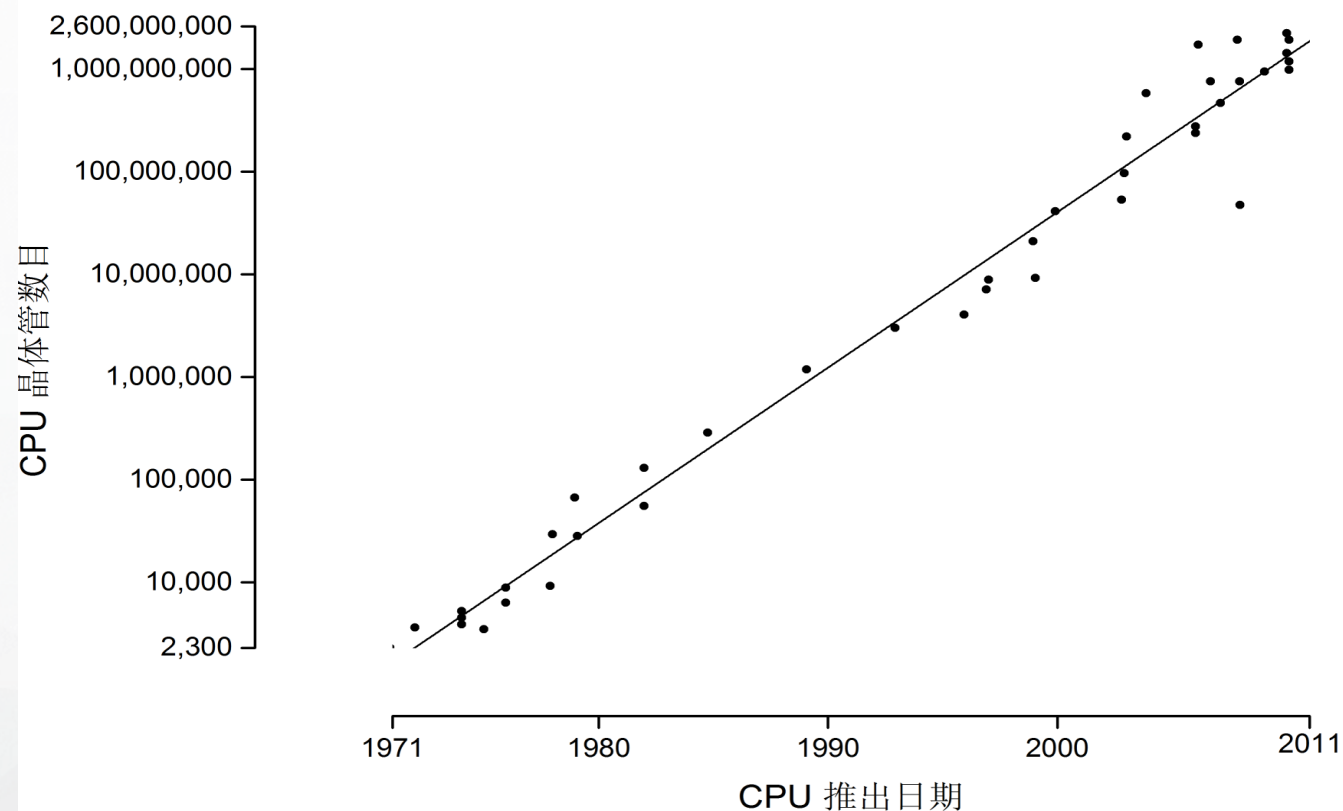


图 CPU晶体管数目随时间变化情况

1.2.2 信息科技为大数据时代提供技术支撑

- 在信息化基础设施方面，据工业和信息化部官网消息，截至2019年12月底，我国互联网宽带接入端口数量达9.16亿个，其中，光纤接入端口占互联网接入端口的比重达91.3%；光缆线路总长度已达4750万公里，相当于在京沪高铁线上往返1.8万余次。同时，近五年来固定宽带和移动宽带资费平均下降90%，速率提升6倍。目前，我国已基本实现“城市光纤到楼入户，农村宽带进乡入村”。
- 据中国信息通信研究院（简称中国信通院）数据，截至2020年2月底，全国建设开通5G基站达16.4万个，5G网络建设基础不断夯实。2020年中国将建设60万~80万个5G基站。

1.2.2 信息科技为大数据时代提供技术支撑

3. 网络带宽不断增加

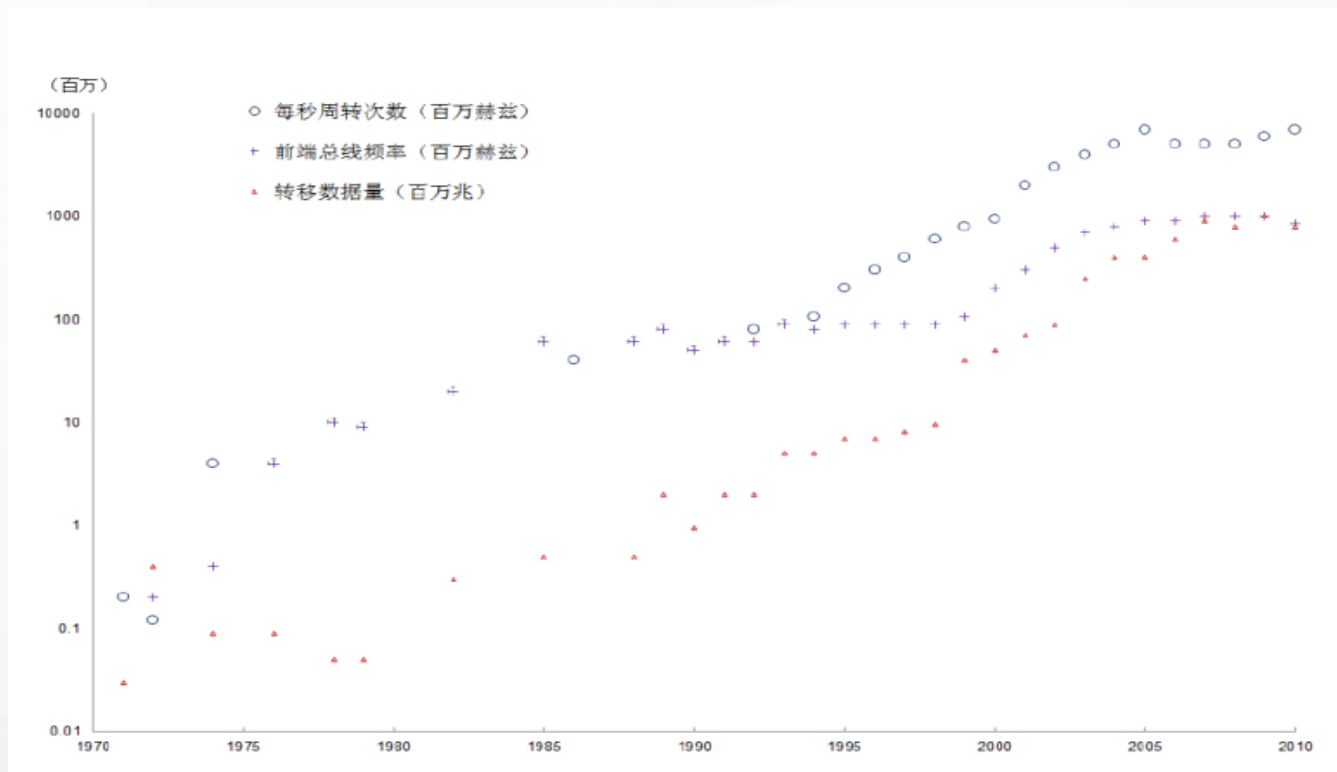


图 网络带宽随时间变化情况

1.2.3 数据产生方式的变革促成大数据时代的来临



图 数据产生方式的变革



1.3

大数据的发展历程



1.3 大数据的发展历程

表 大数据发展的三个阶段

阶段	时间	内容
第一阶段：萌芽期	上世纪90年代至本世纪初	随着数据挖掘理论和数据库技术的逐步成熟，一批商业智能工具和知识管理技术开始被应用，如数据仓库、专家系统、知识管理系统等。
第二阶段：成熟期	本世纪前十年	Web2.0应用迅猛发展，非结构化数据大量产生，传统处理方法难以应对，带动了大数据技术的快速突破，大数据解决方案逐渐走向成熟，形成了并行计算与分布式系统两大核心技术，谷歌的GFS和MapReduce等大数据技术受到追捧，Hadoop平台开始大行其道
第三阶段：大规模应用期	2010年以后	大数据应用渗透各行各业，数据驱动决策，信息社会智能化程度大幅提高



1.4

世界各国的大数据发展战略

1.4.1 美国

1.4.2 英国

1.4.3 法国

1.4.4 韩国

1.4.5 日本

1.4.6 中国



1.4 世界各国的大数据发展战略

国家	战略
美国	稳步实施“三步走”战略，打造面向未来的大数据创新生态
英国	紧抓大数据产业机遇，应对脱欧后的经济挑战
法国	通过发展创新性解决方案并应用于实践来促进大数据发展
韩国	以大数据等技术为核心应对第四次工业革命
日本	开放公共数据，夯实应用开发
中国	实施国家大数据战略，加快建设数字中国

1.4.1 美国

- 美国是率先将大数据从商业概念上升至国家战略的国家，通过稳步实施“三步走”战略，在大数据技术研发、商业应用以及保障国家安全等方面已全面构筑起全球领先优势。
- 第一步是快速部署大数据核心技术研究，并在部分领域积极开发大数据应用。
- 第二步是调整政策框架与法律规章，积极应对大数据发展带来的隐私保护等问题。
- 第三步是强化数据驱动的体系和能力建设，为提升国家整体竞争力提供长远保障。

» 1.4.2 英国

- 英国政府于2010上线政府数据网站Data.gov.uk，同美国的Data.gov平台功能类似，但主要侧重于大数据信息挖掘和获取能力的提升
- 在2012年发布了新的政府数字化战略，实现大数据驱动的社会经济增长
- 2013年英国政府加大了对大数据领域研究的资金支持

1.4.3 法国

- 2011年7月，法国启动了开放数据项目，通过实现公共数据在移动终端上的使用，最大限度地挖掘数据的应用价值。项目内容涉及交通、文化、旅游和环境等领域。
- 2013年12月，法国政府发布《数字化路线图》，明确了大数据是未来要大力支持的战略性高新技术。
- 此外，法国中小企业、创新和数字经济部推出大数据规划，在2013年至2018年在法国巴黎等地创建大数据孵化器

1.4.4 韩国

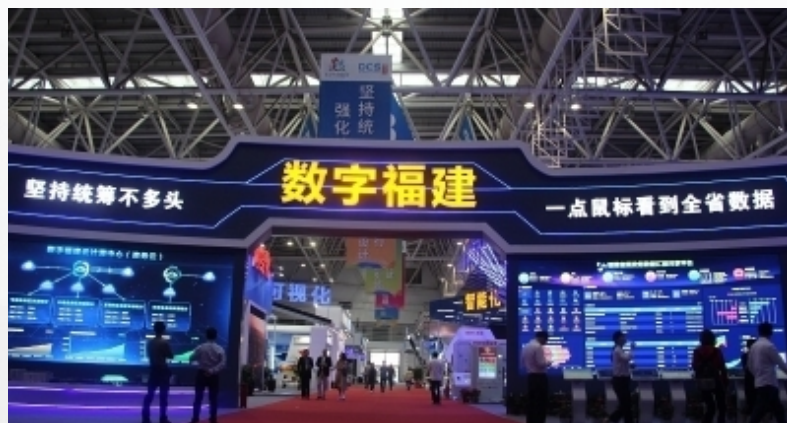
- 韩国的智能终端普及率以及移动互联网接入速度一直位居世界前列，这使得其数据产出量也达到了世界先进水平
- 在朴槿惠政府倡导的“创意经济”国家发展方针指导下，韩国多个部门提出了具体的大数据发展计划
- 2016年年底，韩国发布以大数据等技术为基础的《智能信息社会中长期综合对策》，积极应对第四次工业革命的挑战

1.4.5 日本

- 2010年5月，日本发达信息通信网络社会推进战略本部发布了以实现国民本位的电子政府、加强地区间的互助关系等为目的的《信息通信技术新战略》。
- 2012年6月，日本IT战略本部发布电子政务开放数据战略草案
- 2012年7月，日本政府推出了《面向2020年的ICT综合战略》，大数据成为发展的重点
- 2013年6月，日本公布新IT战略——创新最尖端IT国家宣言，明确了2013-2020年期间以发展开放公共数据为核心的日本新IT国家战略

1.4.6 中国

- 2015年8月，国务院印发了《促进大数据发展行动纲要》。党的十八届五中全会将大数据上升为国家战略。在党的十九大报告中，习近平总书记明确指出：“推动互联网、大数据、人工智能和实体经济深度融合”。
- 2018年4月22日-24日，首届“数字中国”建设峰会在福建省福州市举行。





1.5

大数据的概念

1.5.1 数据量大

1.5.2 数据类型繁多

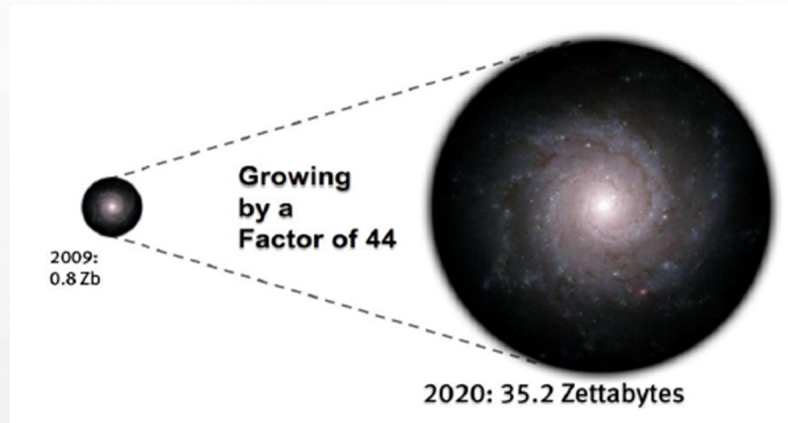
1.5.3 处理速度快

1.5.4 价值密度低



1.5.1 数据量大

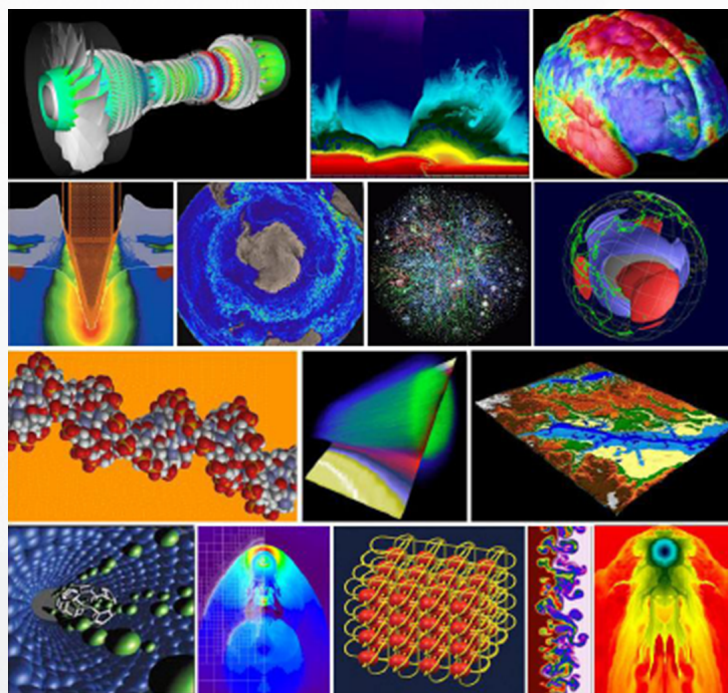
- 根据IDC作出的估测，数据一直都在以每年50%的速度增长，也就是说每两年就增长一倍（大数据摩尔定律）
- 人类在最近两年产生的数据量相当于之前产生的全部数据量
- 预计到2020年，全球将总共拥有35ZB的数据量，相较于2010年，数据量将增长近30倍



TERABYTE	10 的 12 次方	一块 1TB 硬盘		200,000 照片或 mp3 歌曲
PETABYTE	10 的 15 次方	两个数据中心机柜		16 个 Blackblaze pod 存储单元
EXABYTE	10 的 18 次方	2,000 个机柜		占据一个街区的 4 层数据中心
ZETTABYTE	10 的 21 次方	1000 个数据中心		纽约曼哈顿的 1/5 区域
YOTTABYTE	10 的 24 次方	一百万个数据中心		特拉华州和罗德岛州

1.5.2 数据类型繁多

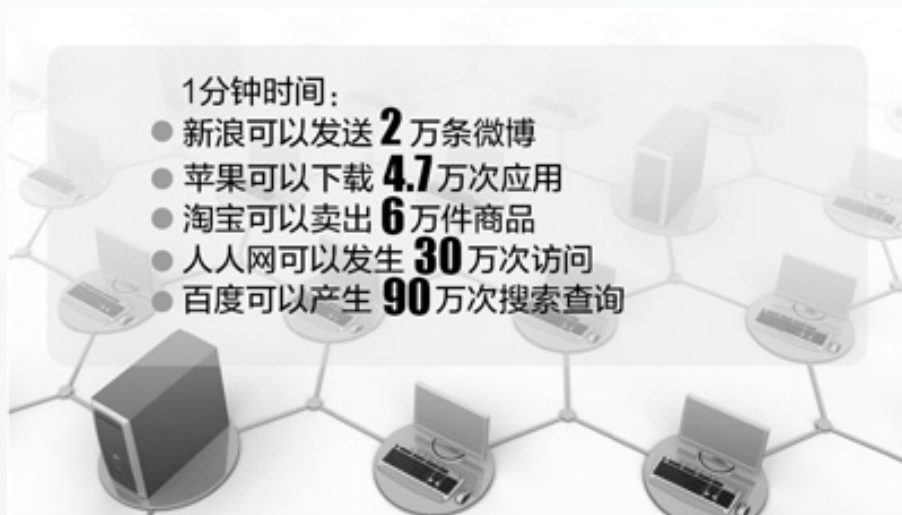
- 大数据是由结构化和非结构化数据组成的
 - 10%的结构化数据，存储在数据库中
 - 90%的非结构化数据，它们与人类信息密切相关



- 科学研究
 - 基因组
 - LHC 加速器
 - 地球与空间探测
- 企业应用
 - Email、文档、文件
 - 应用日志
 - 交易记录
- Web 1.0数据
 - 文本
 - 图像
 - 视频
- Web 2.0数据
 - 查询日志/点击流
 - Twitter/ Blog / SNS
 - Wiki

1.5.3处理速度快

- 从数据的生成到消耗，时间窗口非常小，可用于生成决策的时间非常少
- 1秒定律：这一点也是和传统的数据挖掘技术有着本质的不同



1.5.4 价值密度低

价值密度低，商业价值高

以视频为例，连续不间断监控过程中，可能有用的数据仅仅有一两秒，但是具有很高的商业价值





1.6

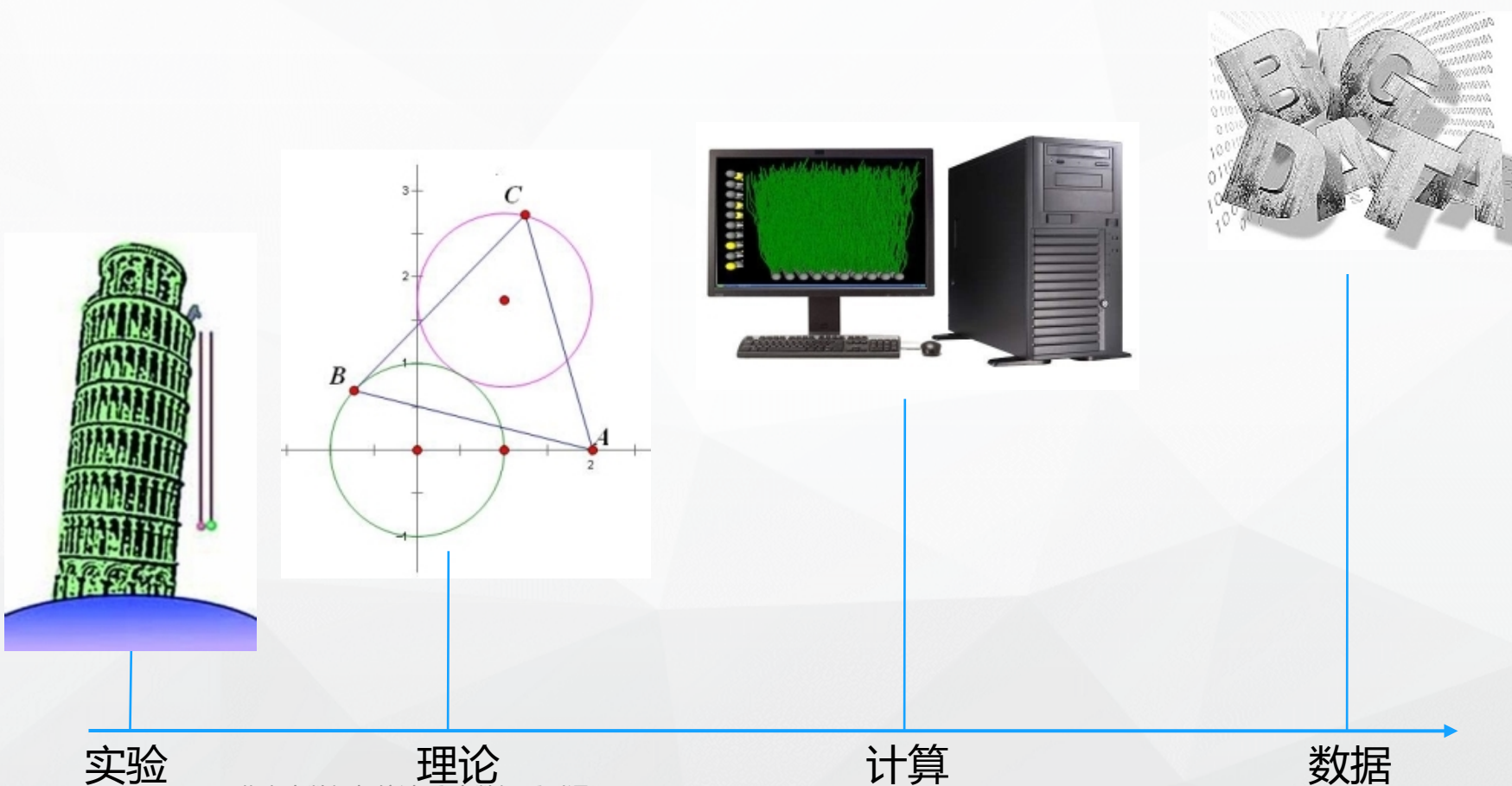
大数据的影响

- 1.6.1 大数据对科学研究的影响
- 1.6.2 大数据对社会发展的影响
- 1.6.3 大数据对就业市场的影响
- 1.6.4 大数据对人才培养的影响



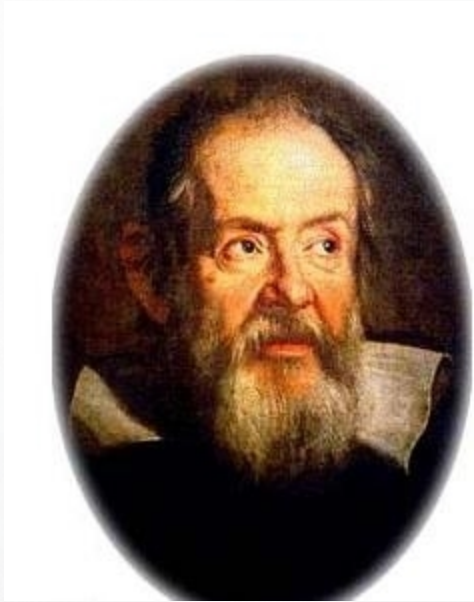
1.6.1 大数据对科学研究的影响

图灵奖获得者、著名数据库专家Jim Gray 博士观察并总结人类自古以来，在科学研究上，先后历经了实验、理论、计算和数据四种范式



1.6.1 大数据对科学研究的影响

科学研究第一种范式：实验



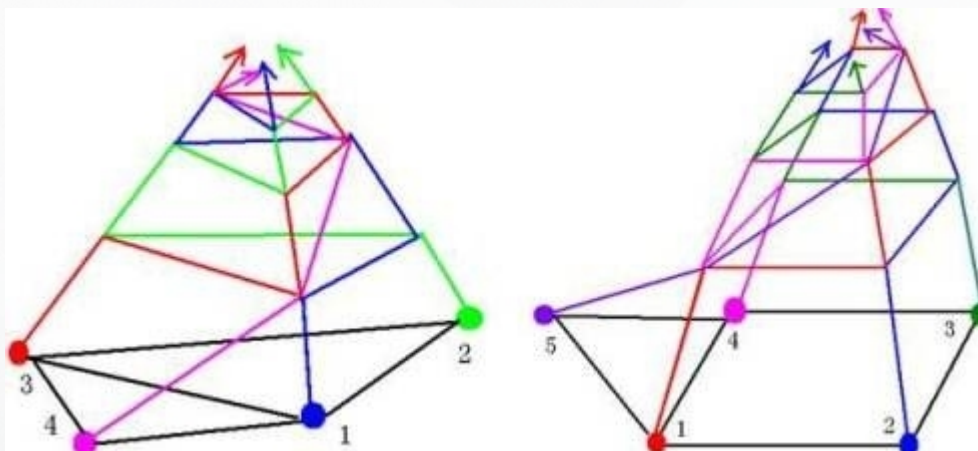
伽利略



伽利略在比萨斜塔做两个铁球同时落地实验

1.6.1 大数据对科学研究的影响

科学研究第二种范式：理论



几何理论



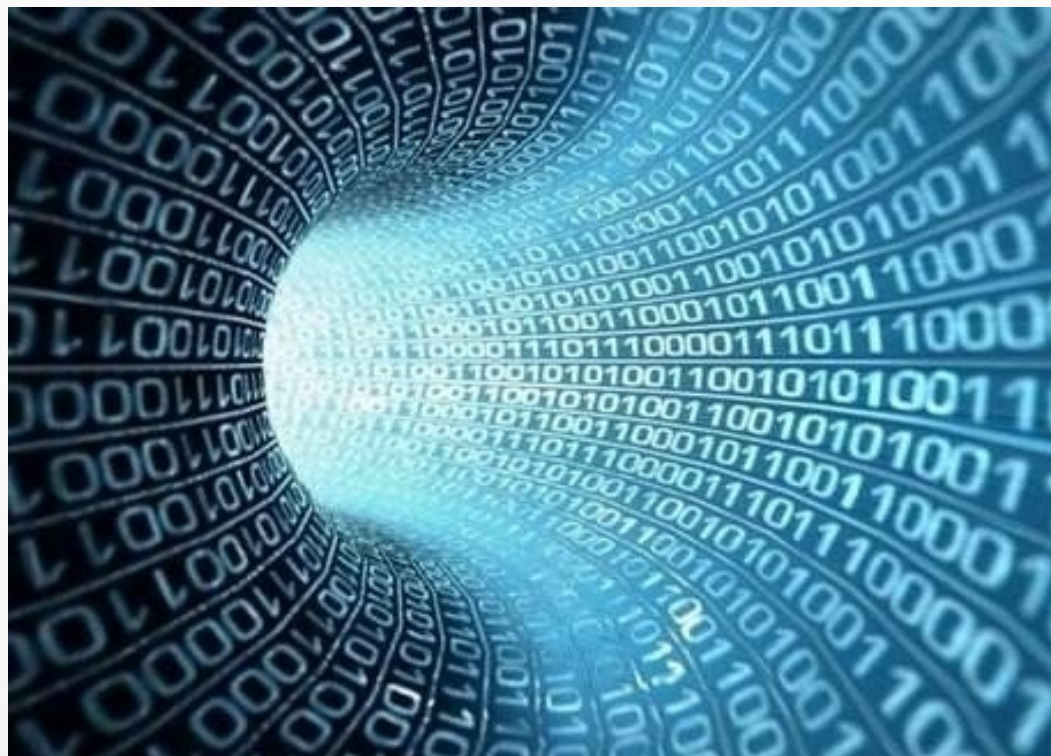
牛顿三大定律

1.6.1 大数据对科学研究的影响

科学研究第三种范式：计算



1.6.1 大数据对科学研究的影响



大数据时代，以数据为中心

➤ 1.6.2 大数据对社会发展的影响

- 大数据决策逐渐成为一种新的决策方式
- 大数据成为提升国家治理能力的新途径
- 大数据应用有力促进了信息技术与各行业的深度融合
- 大数据开发大大推动了新技术和新应用的不断涌现

1.6.3 大数据对就业市场的影响

大数据的兴起使得数据科学家成为热门职业



- 麦肯锡报告，到2018年，在“具有深入分析能力的人才”方面，美国面临着14万到19万的缺口，“可以利用大数据分析来做出有效决策的经理和分析师”缺口则会达到150万
- 国内有大数据专家估算过，5年内国内的大数据人才缺口会达到130万，以大数据应用较多的互联网金融为例，这一行业每年增速达到4倍，届时，仅互联网金融需要的大数据人才就是现在需求的4倍以上
- 根据第四届中国贵州人才博览会发布《全国大数据人才需求指数报告》，2016年2月份，贵阳大数据人才月薪已逼近8000元

1.6.4 大数据对人才培养的影响

- 大数据时代到底需要什么样的人才？
- 一是计算机技术相关人才，包括平台搭建和应用开发
- 二是统计学相关人才，包括数学、建模、算法
- 三是业务人才，就是要有一定的专业领域知识，只有明白目标领域知识的人才能了解数据的意义以及指导数据分析的方向并判断数据分析结果的可信性





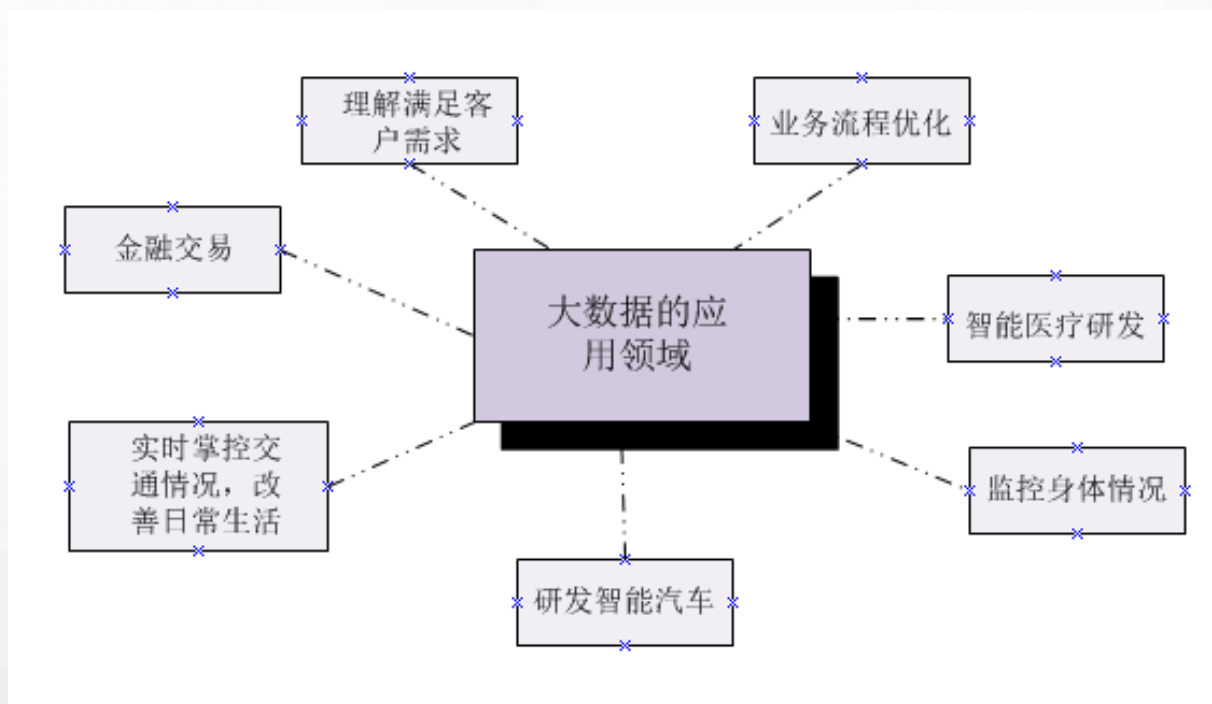
1.7

大数据的应用



1.7 大数据的应用

- 大数据无处不在，包括金融、汽车、零售、餐饮、电信、能源、政务、医疗、体育、娱乐等在内的社会各行各业都已经融入了大数据的印迹



1.7 大数据的应用

- 就企业而言，对大数据的掌握程度可以转化为经济价值的源泉
- 就政府而言，大数据的发展将会提高政府科学决策水平，改变政府传统“拍脑袋”式决策，变为用数据说话，利用大数据分析社会、经济、人文生活等规律，从而为国家宏观调控、战略决策、产业布局等夯实根基
- 在医疗领域，大数据也有不俗表现
- 大数据也悄然地影响着绿茵场上强弱的较量



1.8

大数据产业



1.8大数据产业

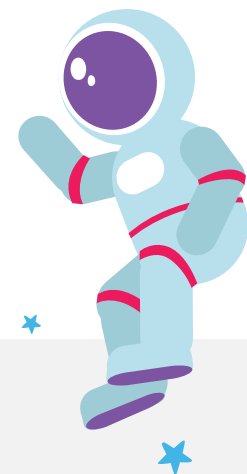
- 大数据产业是指一切与支撑大数据组织管理和价值发现相关的企业经济活动的集合

产业链环节	包含内容
IT基础设施层	包括提供硬件、软件、网络等基础设施以及提供咨询、规划和系统集成服务的企业，比如，提供数据中心解决方案的IBM、惠普和戴尔等，提供存储解决方案的EMC，提供虚拟化管理软件的微软、思杰、SUN、Redhat等
数据源层	大数据生态圈里的数据提供者，是生物大数据（生物信息学领域的各类研究机构）、交通大数据（交通主管部门）、医疗大数据（各大医院、体检机构）、政务大数据（政府部门）、电商大数据（淘宝、天猫、苏宁云商、京东等电商）、社交网络大数据（微博、微信、人人网等）、搜索引擎大数据（百度、谷歌等）等各种数据的来源
数据管理层	包括数据抽取、转换、存储和管理等服务的各类企业或产品，比如分布式文件系统（如Hadoop的HDFS和谷歌的GFS）、ETL工具（Informatica、Datastage、Kettle等）、数据库和数据仓库（Oracle、MySQL、SQL Server、HBase、GreenPlum等）
数据分析层	包括提供分布式计算、数据挖掘、统计分析等服务的各类企业或产品，比如，分布式计算框架MapReduce、统计分析软件SPSS和SAS、数据挖掘工具Weka、数据可视化工具Tableau、BI工具（MicroStrategy、Cognos、BO）等等
数据平台层	包括提供数据分享平台、数据分析平台、数据租售平台等服务的企业或产品，比如阿里巴巴、谷歌、中国电信、百度等
数据应用层	提供智能交通、智慧医疗、智能物流、智能电网等行业应用的企业、机构或政府部门，比如交通主管部门、各大医疗机构、菜鸟网络、国家电网等



大数据与统计系

大数据系列课程



大数据



大数据与统计系数字教师网

<http://hssj.wxcgi.com/>