

# 科技大数据

## BIG DATA 因你而改变

Big Data of Science and Technology: the Driving Force to Make Great Change

戴国强 赵志耘 / 主编



科学技术文献出版社

SCIENTIFIC AND TECHNICAL DOCUMENTATION PRESS

# 科技大数据

## BIG DATA 因你而改变

Big Data of Science and Technology: the Driving Force to Make Great Change

戴国强 赵志耘 / 主编



科学技术文献出版社  
SCIENTIFIC AND TECHNICAL DOCUMENTATION PRESS

· 北京 ·

## 图书在版编目 (CIP) 数据

科技大数据：因你而改变 / 戴国强，赵志耘主编. —北京：科学技术文献出版社，2018. 8

ISBN 978-7-5189-4725-6

I. ①科… II. ①戴… ②赵… III. ①科学技术—数据处理—研究 IV. ①G203

中国版本图书馆 CIP 数据核字 (2018) 第 174970 号

## 科技大数据：因你而改变

策划编辑：李蕊 责任编辑：李晴 张红 杨瑞萍 责任校对：文浩 责任出版：张志平

出版者 科学技术文献出版社

地址 北京市复兴路15号 邮编 100038

编务部 (010) 58882938, 58882087 (传真)

发行部 (010) 58882868, 58882870 (传真)

邮购部 (010) 58882873

官方网址 [www.stdp.com.cn](http://www.stdp.com.cn)

发行者 科学技术文献出版社发行 全国各地新华书店经销

印刷者 北京时尚印佳彩色印刷有限公司

版次 2018年8月第1版 2018年8月第1次印刷

开本 710×1000 1/16

字数 362千

印张 28.5

书号 ISBN 978-7-5189-4725-6

定价 128.00元



版权所有 违法必究

购买本社图书，凡字迹不清、缺页、倒页、脱页者，本社发行部负责调换

## 《科技大数据：因你而改变》

### 编写组

主 编	戴国强	赵志耘			
编写人员	戴国强	赵志耘	袁 伟	张英杰	吴 思
	常 春	刘 伟	张闪闪	杨代庆	王立学
	雷 雪	陈 亮	贾 佳	张均胜	杨 岩
	张兆锋	石崇德	何彦青	望俊成	王 政
	张玄玄	张爱霞	赵 辉	邢晓昭	许 燕
	汪芸辉	张 静	宋培彦	刘 蔚	张志娟
	徐 峰	王 勇	李维波	傅俊英	郑 佳
	雷孝平	张海超	曹 燕	付鑫金	

# 序 言

钱塘江潮般的大数据浪潮汹涌来袭，紧随其后的却是天际边渐入眼帘的一幅泼墨山水，精彩的画面也徐徐铺展开来。

大数据，因为万事万物皆可量化而精彩。

当我们踏入 21 世纪，扑面而来的社会变革、生产力的急速发展、新技术的迅速产业化，都不及大数据年均 50% 的增长速度那么突出。在信息化 1.0 时代，微电子技术的快速发展实现了模拟声音和图像的数字化；而在信息化 2.0 时代，方兴未艾的新一代信息技术革命则正在加速撬动全球数字化、数据化进程。越来越多原先无法测量的东西能被度量，原先无法跟踪的行为可以被跟踪，原先无法记录的活动被数字化记录。互联网的搜索流水线正源源不断地生产着巨量数据，而物联网又悄然将量化的触角进一步延展到可以追索的极限时间和广袤空间，万物皆可量化的理想之光正在逐步照进现实。

大数据，因为万众徜徉数据海洋而精彩。

受交通大数据分享极大地便利和丰富了大众出行的启发，兼有数据生产者和数据使用者双重身份的普罗大众，对开放共享大数据之意义的认识更加深刻，对开发利用大数据之价值的需求更加迫切，对享受进而陶醉于大数据之便利的实践更加频繁。在这种较为宽松适宜的社会氛围下，数据资源就像深埋于地底、彼此隔绝的万千暗河，快速涌流出地面聚成数据之河，汇成数据海洋，又在旺盛需求的光热作用下升华成诸如政府数据云、消费数据云、医疗数据云、科学数据云、资源数据云等类型多样、伸手可及的共享数据彩云，数据资源的这一华丽转身为全社会

谋事创业提供了丰沛甘霖，孕育了无限可能。数据共享应用又源源不断创造出的新数据资源，并不断注入到数据海洋，进一步实现从 Share（享有）向 Partake（分享）【Part（参与）+Take（获取）】的转变，最终打造出一个可持续的数据循环生态系统。

大数据，因为催生创新雨后春笋而精彩。

由于数字化、网络化、智能化、泛在化和可视化等硬件技术的持续进步，大数据资源类型不断拓展、存储中心分布更加灵活、数据传输能力越来越大，越来越快，结构化和非结构化数据加快融合统一，开发利用数据的资源和技术约束不断缩小，让创新者将更多精力和智慧投入多维度、宽视角、广领域的数据分析利用和价值创造中。创新越来越多地与大数据紧紧捆绑在一起，全社会研究大数据、挖掘大数据、存储大数据、利用大数据、服务大数据已成蔚然之势，并在科学研究、技术开发、市场应用、产业发展等创新全链条上蓬勃开展。基于大数据的创新范式让创新思维变得更加宽阔，让创新舞台变得空前宽广，让创新的道路变得更加多元。

大数据，因为颠覆传统生活模式而精彩。

曾几何时，随时找到便捷的出行工具，随身携带手机买遍天下，随心所欲地与远方亲朋好友“面对面”聊天，随地买到去往天涯海角的火车票，还是国人心中遥不可及的梦想。如今，这些梦想不仅成真，而且还在不断地颠覆着我们吃、穿、行、游、购、娱等生活的传统模式。在它们鲜亮的外表背后，大数据与大数据技术发挥着默默无闻却又居功至伟的作用。同时，企业、个人思维模式也在发生着巨大改变，消费行为、日常行为等从不可预测而变成大概率可预测，信息、产品获取变得更具个人定制化特征，因果思维更多地被相关性思维所取代。因为大数据传统生活模式、思维模式正在被快速颠覆，我们对世界的认知也得以不断延伸。

大数据的精彩还远未充分呈现，其精彩程度甚至远远超出我们目前

的想象，尚隐藏在我们未掌握、未探寻到的所在。问题的关键在于是不是需要 we 继续深入挖掘？这还得从当下的世界发展困境说起。在全球化深入发展的今天，相当多的领域都已处于充分竞争状态，传统发展模式下的生产力潜力几乎已发挥殆尽，未来发展之路仍“浮云遮望眼”。尽管“创新者生，守旧者亡”“唯改革者进，唯创新者强，唯改革创新者胜”的道理早已为世人所认同，但谁来驱动创新呢？面对新一轮科技革命和产业变革加快涌起的浪潮，看似“条条大路通罗马”，其实因“九曲回肠”而难以“望尽天涯路”！未来潜于史，知古而鉴今。回顾历次工业革命的历程不难发现，每次生产力的巨大跃升，都是由于适应和满足了在当时即将到来的巨大市场需求变化，无论是以珍妮纺织机为代表的标准化制造和机械化、以福特汽车工厂为代表的流水线规模化制造和电气化，还是以工业控制和全球化分工生产为代表的模块化制造和信息化，莫不如此。

未来已走来，唯变为不变。当前，第四次工业革命晨曦乍现，以人工智能、量子信息、移动通信、物联网、大数据与云计算为代表的新一代信息经济，以基因编辑、生物合成、精准医疗、再生医学为代表的生物经济，以3D打印、自主控制系统、机器人、碳纳米管、石墨烯为代表的未来制造经济纷至沓来，新技术、新产品、新服务令人眼花缭乱，新产业、新模式和新业态更是层出不穷，不断冲击着我们对发展的认知。面对战略性新技术引擎“群雄并起”的局面，哪些“新动能”真正能挑起第四次工业革命的大梁呢？若能准确把握未来社会的主流市场需求，问题就能迎刃而解。

未来之动力，大数据为钥。当前，我们的社会正阔步迈向未来。其中，发展动力向更多靠内在动力转型，生产商品向定制化、服务化、绿色化升级，生产模式向网络化、智能化跨越，生产组织向扁平化、虚拟化、融合化推进，创新模式向以大数据驱动为核心的第四范式更新，发展空间加速向虚拟疆域和现实疆域并举转变。这些转变中蕴藏着的精致需求都直接或间接与以大数据为中心的新一代信息技术有着密切关系，

为此，习近平总书记将大数据谓之信息化发展的新阶段，大数据是把握并开启未来社会发展的一把钥匙。

相对于变化，找到“不变”更为关键，大数据则使其成为可能。第四次工业革命仍将依靠技术跃迁，人仍将作为支撑“核心”，而支撑引导智能化生产的是丰富多样、可共享和利用的“大数据”，它将成为第四次工业革命的标志。把大数据的潜在能力有效转化为产业需求的数据，赋能于所有参与生产活动、生产组织的新型生产力参与者，我们才有可能在全球新一轮变革与发展的竞争中继续高歌猛奏。

大数据时代，世界既变得更加透明，也越来越神秘。如果数据本身有误或不够全面，即使经过再严谨的逻辑推演、再繁复的模型计算、再高深的知识挖掘，也很难得出一个正确的结论，陷入“Garbage in, Garbage out”的循环。由于大数据的强大赋能，任何决策失误都将带来指数型风险。

以生物医药领域的热点——脑疾病治疗为例。随着人类寿命的延长，越来越多的脑疾病严重影响了人们的生活质量并给他们的家庭带来极大负担。尽管生命科学家和医学专家们通过有氧训练或者认知训练刺激神经元的新生来预防阿尔茨海默病，通过药物来刺激神经元生成从而改善抑郁症状，通过新生神经元修补神经损伤……但长期的努力却并没有取得预期的效果。

问题也因此而产生，以上治疗是否都与新生神经元有关？没有取得进展的科学问题能否归纳为“成人脑是否存在新生神经元”？

关于这一问题的研究，2018年两家著名刊物——《细胞—干细胞》第四期杂志发表的博尔德里尼团队研究结论与3月《自然》在线发表刊载的索雷尔斯团队论文观点却截然相反。索雷尔斯团队研究发现胎儿和婴儿的大脑中有大量神经元母细胞和新生神经元，但这些细胞的数量自一岁开始急剧减少，至13岁仍能观察到新生神经元，而所有成年个体的脑标本却未发现一例新生神经元。博尔德里尼团队的研究则发现，海马

区域的神经元母细胞数量随年龄的增长越来越少，由这些母细胞分化出的神经元却并没有随年龄增长而减少。在研究样本中，即使是老年人，在死亡时海马区域仍然有数以千计的新生神经元。

尽管两个研究团队都以人脑作为研究对象，都是对不同成熟阶段的神经元进行免疫荧光染色并用 DCX 和 PSA-NCAM 标记新生神经元，然而得到的结论却完全相反！

是博尔德里尼团队标记新生神经元时发生污染？还是索雷尔斯团队由于使用传统化学手段保存样本，技术上落后于博尔德里尼团队的大脑急速冰冻法？是 DCX 的表达并不恒定，对处理时间过于敏感？还是某个团队没有对海马区域进行完整检查？如果以上皆非，还有何原因导致两项研究结果迥异？其重要判别点又是什么？

如果以上治疗与新生神经元无关的话，那又是什么影响着治疗的有效性？已有的研究数据积累量显然无法回答这一问题，犹如盲人摸象，由于每个人只摸到了一小部分，导致对这个问题还难以达到认知上的统一。进一步设想一下，如果索雷尔斯团队的研究结果反映了真实情况，即新生神经元在成年人中几乎不存在，那之前的治疗研究、巨量投入、学科建设、人才培养等，是不是都有点南辕北辙？

以上讨论表明，数据样本量不足已成为制约科学研究、技术开发和创新的重大瓶颈因素，从思维、哲学和方法论层面思考科学技术发展的基本问题无可避免。事实上，科学家和工程师们执着探寻的因果关系与相关关系的研究对象，早已由原来的人类世界与物理世界的两极相互作用，扩大到人类世界、物理世界与网络世界的三极交互作用，影响经济和社会发展的因素显著增加。在这种情况下，如果不能首先探寻清楚基本问题，就无法看清事物发展脉络并总结其发展规律，更不要奢谈开展预测和科学决策。而大数据深度分析无疑为研究基本问题并提供解决方案等后续工作提供了一个“金刚钻”和“导航灯”，例如，通过更全面的分析而摸清上述不同研究产生矛盾结果的原因所在，并据此提出新的

研究思路。

如果说过去 20 年是 DNA 时代，也就是基于数据 (Data)、网络 (Net) 和自动化 (Automatic) 的组合，通过发挥“基因”的遗传效应实现全球的高速发展，那么未来将升级到 RNA 时代，不仅继续发挥遗传效应，而且还将在适者生存和不确定性的前提下复制优秀品质并引导良性变异，推动经济社会继续保持协调、高速、有序发展。这里的 RNA 实质上概括了大数据时代的三大变化特征。

大数据“赋能”数据资源 (Resource)。在大数据时代，数据快速从成果化向资源化、资产化、资本化三位一体转变，即通过深入挖掘使用大数据的潜在价值，实现大数据的保值、增值、赋值能力，即为数据“赋能”。

大数据向正新信息科技转变 (Next Generation of Information Science)。在摩尔定律达到极限的今天 (晶体管价格在连续 50 年下降后逆向上涨，若非市场操作，则可判断“摩尔定律”极限将至)，传统信息技术所支撑的巨量存储和计算模式将被迫提升到科学存储和智能计算的新篇章，并据此将会源源不断地产生新的科学发现、技术创造、产品创新。

大数据与人工智能 (AI) 紧密结合。今天人类社会面临的挑战是必须能够完成传统机器和人类自身不愿做、不敢做和做不到的事情。资源和环境的约束要求我们必须冗余数据和极小数据两种极限情况下也能做出正确处理，这意味着我们必须适应未来的极简需求，呈现于外部的极简，蕴藏于极其复杂的处理与计算，这绝非人类智慧的核心优势所在。多维度关联数据、不断优化的系统算法及不断升级的计算机软硬件将帮助我们更加深入地认识世界、改造世界。本书力求通过知识组织和数据挖掘等工具方法将纷繁的数据进行全面梳理和组织加工，为科学寻踪和发现提供条件；通过结构和非结构化数据存储的科学方法，在大数据的丛林中依靠可靠的路径和定位知识关联；通过机器智能帮助用户按语义提取和处理数据；利用大数据开发出最新的翻译软件，使我们获取

最新的信息更加便利；利用机器智能解决信息过载问题，实现“人找信息”到“信息找人”的巨大转变；利用大数据公共服务平台及后台技术支撑手段，提供从政府服务、产业服务、企业服务到个人服务的完整服务谱系。数据处理、数据流动、关联与分析、结果动态展示……我们希望借助这些科学的方法和手段展现大数据的潜在精彩。

在过去 40 年里，中国科学技术信息研究所围绕大数据持续耕耘不辍，本书内容就是长期努力的重要呈现，是对全所大数据研究实践工作的阶段性总结。写作本书的更重要目的是力图增进读者对大数据技术的发展和应用有更全面的了解，对依托大数据开展研究和创新创业的重要性有更深入的理解。由于大数据所涉及领域众多，技术本身仍在迅猛发展，我们的研究也在不断深化中，不当之处敬请读者指正。

在第四次工业“革命”中人是最为关键的要素，未来跨越式发展要依靠大数据，更要依靠人类智慧！让我们回归“人定胜天”（人自身安定胜过老天的帮助）的本意，认识并按照规律踏实勤勉工作而不要奢求老天的帮助。大数据的精彩需要我们共同去了解、理解、发现、实现。

因为，有你、有我，大数据才会真的精彩！

戴国强

2018 年 9 月 28 日

于中国科学技术信息研究所

# 目 录

<b>第一章 大数据思维与科技大数据</b> .....	<b>001</b>
1.1 概念辨析：科技大数据的内涵与外延 .....	001
1.1.1 科技大数据基础概念 .....	001
1.1.2 科技大数据基础理论 .....	004
1.2 技术视角：“数”领风骚的流派之争 .....	009
1.2.1 大数据技术发展路径 .....	010
1.2.2 大数据技术架构 .....	015
1.3 数据视角：数据今生“岭”与“峰” .....	022
1.4 管理视角：“治理”而非“管理”的能力 .....	024
1.4.1 科技大数据治理现状 .....	024
1.4.2 科技大数据能力需求 .....	026
1.5 服务视角：大数据服务于“四链” .....	030
1.5.1 服务创新：科技大数据 + 创新链 .....	031
1.5.2 服务产业：科技大数据 + 产业链 .....	033
1.5.3 服务资金：科技大数据 + 资金链 .....	036
1.5.4 服务人才：科技大数据 + 人才链 .....	037
<b>第二章 主题词表：规范科技大数据的基石</b> .....	<b>042</b>
2.1 知识散沙：科技大数据需要结构化 .....	042
2.1.1 科技大数据需要知识单元结构化 .....	042

2.1.2	科技大数据需要语义关联化 .....	043
2.2	什么是主题词表 .....	044
2.2.1	科技数据组织主题法 .....	044
2.2.2	主题词表及其家族 .....	045
2.2.3	科技术语 .....	047
2.2.4	科技概念 .....	049
2.2.5	主题词表结构 .....	050
2.2.6	主题词表构建方法 .....	051
2.2.7	《汉表》今与昔 .....	056
2.3	主题词表应用：以新型《汉表》为例 .....	057
2.3.1	使用新型《汉表》进行文本分析 .....	058
2.3.2	智能化的检索：跨语言检索、扩检与缩检 .....	061
2.3.3	知识组织：科技术语聚合与关联 .....	064
2.4	新时代下主题词表的挑战与展望 .....	066
2.4.1	主题词表面临与存在的挑战 .....	066
2.4.2	主题词表未来的发展方向 .....	068
<b>第三章</b>	<b>元数据：科技大数据知识关联之源 .....</b>	<b>072</b>
3.1	元数据概览 .....	072
3.1.1	元数据发展脉络 .....	073
3.1.2	元数据的“困境” .....	074
3.2	兼容并蓄：多源异构元数据的融合 .....	076
3.2.1	元数据资源掘金 .....	077
3.2.2	海量元数据集成整合 .....	078
3.2.3	元数据语义关联构建 .....	080
3.2.4	案例：元数据集成管理系统 .....	082

3.3	包罗万象：元数据关联挖掘及应用 .....	087
3.3.1	元数据关联 .....	087
3.3.2	关联数据：科技大数据及关联的描述 .....	090
3.3.3	知识链接：科研实体的关联网络 .....	092
3.3.4	科技大数据知识链接服务 .....	095
3.4	元数据发展展望 .....	101
<b>第四章</b>	<b>语义信息抽取：计算机理解科技大数据的钥匙 .....</b>	<b>105</b>
4.1	计算机为何要理解科技大数据 .....	105
4.1.1	海量数据去粗取精的需要 .....	105
4.1.2	内容深度挖掘的需要 .....	106
4.1.3	知识服务的需要 .....	107
4.2	计算机怎样理解科技大数据 .....	108
4.2.1	技术框架：信息抽取 .....	108
4.2.2	关键技术一：命名实体识别 .....	109
4.2.3	关键技术二：实体消歧 .....	112
4.2.4	关键技术三：语义关系抽取 .....	113
4.3	信息抽取实践：以美国专利局数据为例 .....	118
4.3.1	专利信息概念模型的建立 .....	118
4.3.2	专利数据集的形成和预处理 .....	119
4.3.3	专利信息的抽取 .....	120
4.3.4	实体语义关系的应用：专利诉讼案中证据专利的 识别 .....	127
4.4	信息抽取面临的不足和挑战 .....	130

## 第五章 大数据时代的伯乐：从科技团队创新能力评价看科技评价

### 工作 ..... 133

5.1 科技评价：不评价，无优化.....	133
5.1.1 科技评价的概念、特点及功能 .....	134
5.1.2 科技评价工作的分类.....	136
5.1.3 常用的科技评价方法.....	139
5.2 基于文献计量的科技团队创新能力评价 .....	141
5.2.1 科技团队的概念和类型.....	141
5.2.2 科技团队的生命周期.....	142
5.2.3 科技团队创新能力评价模型 .....	146
5.2.4 科技团队评价指标体系设计 .....	149
5.2.5 科技评价案例.....	152
5.3 评价工作：评过去，期未来.....	158
5.3.1 评价方法的适用性和局限性 .....	158
5.3.2 评价指标的选择原则和依据 .....	159
5.3.3 评价结果的影响和使用.....	160
5.4 本章小结.....	161

## 第六章 创新图谱：基于数据的科技创新决策可视化平台..... 163

6.1 科技决策的趋势与需求.....	163
6.1.1 科技决策日益需要科技大数据支撑 .....	163
6.1.2 信息可视化及分析有助于科技决策 .....	165
6.2 科技创新图谱关键技术.....	168
6.2.1 数据集成与组织.....	168
6.2.2 数据可视化分析.....	170
6.3 科技创新图谱系统平台.....	176
6.3.1 数据要素 .....	176

6.3.2	功能架构 .....	177
6.3.3	城市创新 .....	179
6.4	科技创新图谱总结与展望 .....	198
<b>第七章</b>	<b>机器翻译：多语言科技大数据智能服务的桥梁 .....</b>	<b>201</b>
7.1	科技大数据离不开机器翻译 .....	202
7.1.1	跨语言知识挖掘的需求 .....	202
7.1.2	机器翻译发展态势 .....	203
7.1.3	机器翻译的主流方法 .....	205
7.2	神经机器翻译让机器翻译不再“神经” .....	208
7.2.1	神经网络 .....	208
7.2.2	神经机器翻译模型 .....	211
7.2.3	神经机器翻译前沿进展 .....	215
7.3	让科技交流没有语言障碍 .....	216
7.3.1	利用知识组织解决科技机器翻译领域自适应 .....	217
7.3.2	“科信智译”翻译服务 .....	221
7.4	本章小结 .....	225
<b>第八章</b>	<b>智能问答：塑造大数据时代的新型服务 .....</b>	<b>230</b>
8.1	智能问答系统发展史 .....	230
8.2	几种智能问答系统的服务实践 .....	232
8.2.1	Siri 的崭露头角 .....	232
8.2.2	一鸣惊人的沃森 .....	232
8.2.3	基于检索服务开发的智能问答系统 .....	233
8.2.4	京东 JIMI 在零售领域的探索 .....	233
8.2.5	阿尔法小蛋机器人在教育领域的探索 .....	234
8.2.6	小 i 机器人在客服领域的探索 .....	234

8.3	智能问答系统的实现方法.....	235
8.3.1	基于规则的问答系统.....	235
8.3.2	基于知识图谱的问答系统.....	235
8.3.3	基于端到端算法的问答系统.....	236
8.4	关键技术.....	237
8.4.1	自然语言处理.....	238
8.4.2	知识图谱.....	240
8.4.3	信息检索.....	241
8.4.4	深度学习.....	242
8.5	“中信所小科”智能问答机器人.....	244
8.5.1	“中信所小科”智能问答机器人是什么.....	244
8.5.2	“中信所小科”智能问答服务的主要架构.....	245
8.5.3	“中信所小科”智能问答服务实践分享.....	249
8.6	展望.....	250
<b>第九章 资源平台：知识聚集和共享的枢纽.....</b>		<b>252</b>
9.1	科技报告：详细记录科研过程和结果的大文本.....	252
9.1.1	什么是科技报告.....	253
9.1.2	国家科技报告大数据管理.....	257
9.1.3	国家科技报告大数据共享服务.....	265
9.2	科技成果：为大数据时代产学研结合助力.....	271
9.2.1	现状及挑战：科技成果去哪了.....	271
9.2.2	国家科技成果数据库：让成果大数据安身立命.....	272
9.2.3	科技成果成熟度评价：大浪淘沙现珍珠.....	277
9.3	专利：创新者们的机会和陷阱.....	285
9.3.1	专利分析能带来什么.....	285