

大数据时代的解决方案

# 块数据

# 4.0

人工智能时代的  
激活数据学

大数据战略重点实验室◎著

# 块数据 4.0

人工智能时代的  
激活数据学

大数据战略重点实验室◎著

图书在版编目 ( CIP ) 数据

块数据 4.0 : 人工智能时代的激活数据学 / 大数据  
战略重点实验室著. -- 北京 : 中信出版社, 2018.5

ISBN 978-7-5086-8886-2

I. ①块… II. ①大… III. ①经济管理-数据管理-  
通俗读物 IV. ①F2-39

中国版本图书馆 CIP 数据核字 ( 2018 ) 第 073607 号

块数据 4.0——人工智能时代的激活数据学

著 者：大数据战略重点实验室

出版发行：中信出版集团股份有限公司

( 北京市朝阳区惠新东街甲 4 号富盛大厦 2 座 邮编 100029 )

承 印 者：北京通州皇家印刷厂

开 本：880mm×1230mm 1/32

印 张：11.5 字 数：230 千字

版 次：2018 年 5 月第 1 版

印 次：2018 年 5 月第 1 次印刷

广告经营许可证：京朝工商广字第 8087 号

书 号：ISBN 978-7-5086-8886-2

定 价：59.00 元

版权所有·侵权必究

如有印刷、装订问题，本公司负责调换。

服务热线：400-600-8099

投稿邮箱：author@citicpub.com

## 编撰委员会

总 顾 问 陈 刚 闫傲霜 李再勇

编 委 会 主 任 李再勇

编委会常务副主任 许 强 陈 晏

编 委 会 副 主 任 聂雪松 徐 昊 连玉明

主 编 连玉明

副 主 编 朱颖慧 武建忠

执 行 副 主 编 宋希贤 宋 青 胡海荣

主 要 撰 稿 人 连玉明 朱颖慧 武建忠 宋 青

胡海荣 宋希贤 张俊立 张龙翔

范贤昱 龙荣远 黄 倩 邹 涛

翟 斌 郑 婷 陈 威 何 露

姜 璠 陈 鹏 胡亚男 田翠梅

学 术 秘 书 李瑞香 江 岸

眼看乾坤一局棋，满枰黑白子离离。

铿然一子成何劫，唯有苍苍妙手知。

这是被称为“波斯李白”的诗人奥马尔·海亚姆（1048—1122）的《鲁拜集》中的诗句。作为一名精通天文和数学的大学者，奥马尔认为，宇宙的规律是可以探知的，并可以用严密而美妙的数学方式表示出来。前定与随机，必然与偶然，向来是人文科学中长期争论不休的命题。自然科学理论始终受实验和观测的检验，而它的每一个重大发现又都会反馈到文化和社会的层面，对人的哲学和历史观有所启示。

决定性和概率性一直被当作数学、物理等学科对自然界的描述方式。在牛顿创立古典力学之后的250年间，直至20世纪20年代，决定论长期处于主导地位，基于概率论的统计描述或者说数据的描述，则一直属于不得已情况下所采用的辅助手段。决

定性的牛顿力学从计算和预测的观点来看，实际上也具有内秉随机性，这就是微观层次上的混沌运动。大量隐藏在暗数据背后的某些看似简单原因所导致的复杂后果，则渐渐成为混沌研究的重要信息。混沌不是无序和混乱。与人们习以为常的周期排列或对称形状的数据相比，大自然和人类社会中的很多数据其实就是一种没有周期性次序的混沌。在理想模型中，它可能包含着无穷的内在层次，层次之间存在自相似性或不尽相似。在观察手段和技术的分辨率不高时，只能看到每一个层次或某一种类型的结构。但技术条件改变或提高后，在远离不能识别之处就会出现更小尺度上的结构。零维的点、一维的线、二维的面、三维的体和四维的时空，是人们现在所能认知的数据空间。如果在不远的将来，我们真的进入一个超数据时代，现有的技术和描述手段也许就无法对这种高度无序数据的混沌运动进行分形，而关于相变和临界现象理论的框架也需要一个新的重构。

这时我们不由得想到那个著名的洛伦兹“蝴蝶效应”理论，其实和这个理论相联系的还有一个被称为“湍流发生机制”的观点，认为向湍流的转变由少数自由度决定，经过两三次突变，运动就到了维数不高的奇怪吸引子上。这里所谓的吸引子是指运动轨迹长时间之后的终极形态，它可能是稳定的平衡点或周期性的轨道，也可能是继续不断变化、没有明显规则或次序的许多回转曲线。无论是蝴蝶效应还是湍流发生机制，其实都是对我们现在正在研究的激活数据学的一种理论上的关照和呼应。事实上，大数据乃至超数据时代的数据运动，就是这样一种处于混沌和分形

之下的对数据运动轨迹及其规律的研究。许多看起来杂乱无章、随机起伏的数据变化或时空穿越，可能造成的就是类似亚马逊级别的数据风暴。如果说上述蝴蝶效应粉碎的是本就无法实现的长期天气预报的幻梦，那么紧接着的奇怪吸引子告诉我们的是，人类对于天气的实际预报能力并没有因那只蝴蝶的翅膀而受到任何影响，相反，却因对于更加混沌的数据的研究而提高了。激活数据学就是一种基于复杂理论及混沌研究的关于未来大数据乃至超数据时代的理论假说，就像上面讲的天气预报，但它所关心的并不是下个星期的晴雨冷热，而可能是未来10年耕种季节的平均降水量和平均气温。激活数据学研究使以往根据统计原则所做的预报上升为数据动力学的预报，也就是应用了似是随机现象的内在规律，从而提高了预测单个轨道近期行为的精确度，并丰富了长期预报的办法。

同样，我们还可以考察一个似静实动的模型。让沙子从一个漏斗孔中缓缓落到桌面上，形成渐渐变大的沙堆，总有最后新添加的某一粒沙子会在整个沙堆勉强维持平衡的锥面上导致一次“雪崩”，使一撮沙子滑到堆底，雪崩留下的小洼地会被后续的沙流填平，直到下一次更大的雪崩。在我们收集所有这些雪崩的数据后，可以发现它们的大小和间隔遵循某些数据动力规律，而沙堆模型无疑也启发了我们对于数据激活状态中的相变和突变的研究。无论数据的平衡态的相变或非平衡的临界多么不确定，可以确定的是，在搜索、融合、激活和碰撞等一系列状态下，数据在某一个临界点附近的扰动必然会导致某种全局性后果。当然，是

否存在可以被“激活”的“数据蝴蝶”或“数据吸引子”，还需要我们进一步探寻，但这并不否定我们的所愿，而仅仅需要我们从实际数据的研究和挖掘中进一步加以发现。

作为一种理论假说，激活数据学就像一座朝向深邃的大数据宇宙的“天眼”。它是未来人类进入云脑时代的预报，是关于混沌的数据世界的跳出决定论和概率论的非此即彼、亦此亦彼的复杂理论的大数据思维范式的革命。从一定意义上来说，大数据就是面向未来社会人类需要破译的“基因”。正如因发明一种DNA（脱氧核糖核酸）快速测序方法而获得1980年诺贝尔化学奖的吉尔伯特针对生物学研究范式的变化指出的，“正在兴起的新的范式在于，所有的基因将被知晓，今后生物学研究项目的起点将是理论的。一位科学家将从理论的假设开始，然后才转向实验室去检验该假设”。是的，借助日渐深入的人工智能的发展，大数据的理论研究正在激发人类的新的假想和猜测。正是这种假想和猜测，让我们以某种“对称破缺”的方式去探知深邃未知的数据海洋，发现诸多社会发展法则背后产生影响甚至支配的物质和数字的力量。

人生是一种快变量，语言是一种慢变量，而数据将是一种突变量。虽然“未知”依然是现实的一部分，但是身处海量数据大爆发时代，人们坚信，未来已来！从“块数据1.0”到现在的“块数据4.0”，我们一直在持续探讨这个已来的未来，尤其是基于对“以人为原点的数据社会学的范式革命”的认知。事实上，从一开始，我们就没有把大数据仅仅看作所谓的“大”的数据，而是把大数据看作一种“活”的数据，因为只有激活，大数据才有生命，



才有社会属性，才能成为未来世界人们赖以生存与发展的土壤和空气。最后，套用《爆裂》一书中关于现代世界生存的九大原则中“系统优于个体”的表述：真正具有竞争性的是一个系统，而非一个特别强大的个体；是一套能够保证不断成功的制度，而不是一个天才个人的行为。同样，激活数据学就是这样一个思想的系统，就是要为我们身处的这个大数据时代找到一个解决方案，这个方案可以构建一个融合数据、计算和场景的系统，让我们在大数据的时空中真正“思考和行动”起来。世界正处于根本结构性变革中，我们必须具备这样一种能力，即下意识地适应和发现因不适应我们的旧习惯而被忽视的事情。

连玉明

大数据战略重点实验室主任

2018年4月3日于北京

绪 论	大数据时代的解决方案	001
<b>第一章</b>	<b>超数据时代的数据拥堵</b>	
第一节	小数据时代、大数据时代和超数据时代	012
	(一) 小数据时代	012
	(二) 大数据时代	016
	(三) 超数据时代	021
第二节	奇点来临：数据大爆炸	026
	(一) 数据连接型社会：数据量化世界	026
	(二) 数据大爆炸：海量、复杂与失控	032
	(三) 数据失真、数据依赖与数据安全	035
第三节	数据拥堵与数据治理	038
	(一) 数据拥堵的由来	038
	(二) 从生命周期视角思考数据拥堵	040

- (三) 数据拥堵的治理范式 044

## 第二章 激活数据学：基于块数据理论的解决方案

### 第一节 复杂理论与块数据 052

- (一) 复杂性的涌现 052
- (二) 块数据的数据观 055
- (三) 数据学与数据科学 058

### 第二节 激活数据学的提出 061

- (一) 激活数据学的由来 061
- (二) 激活数据学的理论框架 062
- (三) 激活数据学的时代价值 068

### 第三节 激活数据学与数据激活机理 071

- (一) 数据搜索：智能感知 071
- (二) 关联融合：智能聚合 072
- (三) 自激活：智能决策 074
- (四) 热点减量化：智能筛选 076
- (五) 群体智能：智能碰撞 077

## 第三章 数据搜索：智能感知

### 第一节 智能感知与交互 082

- (一) 生物感知 082
- (二) 机器感知 085
- (三) 交互识别 088

### 第二节 搜索引擎：连接人与信息 091

- (一) 从“寻物”到“搜数” 091

	(二) 谷歌搜索: 让流动的信息产生智能	095
	(三) 搜索引擎的工作原理	098
第三节	搜索引擎到人工智能的终极演进	103
	(一) 全局化范围搜索	103
	(二) 智能化目标识别	106
	(三) 无界化协同感知	110
<b>第四章</b>	<b>关联融合: 智能聚合</b>	
第一节	人脑信息的处理与融合	118
	(一) 对象感知	118
	(二) 情景关联	120
	(三) 信息融合	122
第二节	智能数据处理	124
	(一) 大数据融合处理模式	124
	(二) 数据融合处理局限	127
	(三) 基于人脑模式的数据关联融合	129
第三节	数据融合: 构建新型数据关系	134
	(一) 降维去噪	135
	(二) 关联识别	140
	(三) 融合重构	144
<b>第五章</b>	<b>自激活: 智能决策</b>	
第一节	脑认知与类脑计算	154
	(一) 神经元与神经网络	154

- (二) 从学习到决策 157
- (三) 人脑智能决策对机器学习的启示 163

## 第二节 让机器像人一样思考 166

- (一) 从“深蓝”到“阿尔法元” 166
- (二) 构造人工神经网络 169
- (三) 深度学习驱动机器智能决策 171

## 第三节 智能判断与决策 175

- (一) 提取特征 175
- (二) 构建模型 177
- (三) 决策输出 181

# 第六章 热点减量化：智能筛选

## 第一节 遗忘，是为了更好的记忆 187

- (一) 人脑的记忆存储极限 187
- (二) 记忆的选择性封存 192
- (三) 遗忘也是一种学习 196

## 第二节 删除，数据取舍之道 199

- (一) 数字记忆是生物记忆的延伸 199
- (二) 全面数字存储下的信息失控 203
- (三) 数字记忆与信息取舍 206

## 第三节 筛选，选择最优决策 210

- (一) 数据匹配与简约 210
- (二) 优化算力配置 214
- (三) 选择最优算法 218

## 第七章 群体智能：智能碰撞

### 第一节 头脑风暴：发现好想法和做出好决策 226

- (一) 创造力是发现好想法的源泉 226
- (二) 群体合作与互动 228
- (三) 群体决策与判断 230

### 第二节 群体学习：从个体智能到群体智能 234

- (一) 个体智能的局限 234
- (二) 从生物群体到机器人群体 237
- (三) 群体机器人的行为协作 242

### 第三节 群体空间：人脑智慧和机器智能的交互 244

- (一) 人机优势互补 244
- (二) 机器智能进阶 246
- (三) 人机社会化协作 250

## 第八章 激活数据学的应用场景

### 第一节 激活数据学下的自动驾驶 258

- (一) 智能驾驶引领新一轮工业革命 258
- (二) 激活数据学在无人驾驶中的应用场景 261
- (三) 激活数据学为智能驾驶提供理论依据 264

### 第二节 激活数据学下的城市大脑 267

- (一) 城市大脑：城市的数据智能中枢 267
- (二) 激活数据学优化城市大脑的系统应用 269
- (三) 激活数据学让城市大脑更智慧 273

第三节 激活数据学下的医疗影像 274

- (一) 人工智能赋能医疗影像 274
- (二) 激活数据学在医疗影像中的应用策略 276
- (三) 激活数据学提升医疗影像价值 280

第四节 激活数据学下的智能语音 282

- (一) 智能语音交互：进阶的交互模式 282
- (二) 智能语音技术提升的路径选择 284
- (三) 激活数据学开启语音交互新时代 287

第九章 云脑时代：开启数字文明新纪元

第一节 驱动云脑时代的“三驾马车” 294

- (一) 数据驱动 294
- (二) 计算驱动 297
- (三) 场景驱动 301

第二节 区块链：人工智能任性发展的“保险阀” 303

- (一) 哲学视域下的人工智能风险 303
- (二) 区块链与秩序互联网 309
- (三) 区块链重塑人工智能时代新生态 312

第三节 数权法与数字文明新时代 315

- (一) 云脑时代的制度安排与法律规制 315
- (二) 数权法构建数字文明新秩序 319
- (三) 构建网络空间人类命运共同体 322

参考文献 327

术语索引 341

后记 349

## 大数据时代的解决方案

本书探讨的主题是大数据时代激活数据学的提出、运行机理及场景应用。激活数据学是以充分发挥人机群体智能为核心，综合运用数据科学、生命科学和社会科学提出的海量数据存储、处理的解决方案。激活数据学将确立一个新的观察人类智能和机器智能的视角，引导人们重新审视数据无限膨胀可能造成的人类认知障碍，重新思考维持一个健康、安全和有效的数字社会的根本办法，建立与人类智能复杂性同步的人工智能系统，开启用复杂性系统思维认识未来世界和改造未来世界之旅。

## 大数据时代面临的问题与挑战

在人类文明的伊始，人与人的第一声交流即意味着“连接”的开始。语言使人与人连接，并促使用于记载事物的文字、数字



符号产生，这样的“连接”便产生了“数据”，并演化为人类文明最初的信息与知识。在漫长的农耕文明时代，“连接”主要以语言沟通和书面文字沟通的形式存在。进入工业文明时代，“连接”开始通过无线电台、电报、电视的形式存在，但这样的“连接”产生的信息往往是单向性且缺乏互动的。互联网时代，人和人开始通过网络进行复杂交错的互动连接。社交网站、电子邮件、搜索引擎、聊天工具……人类建立连接的方式趋于多样化、多维化，人类社会产生的数据也因而大量积累。与此同时，数据的价值越来越受到人们的重视。数据深刻作用于政治、经济、文化等领域，带来更多的创新机会，从生产、生活到科研，一个大数据时代正在开启。

在美丽的贵州省黔南州布依族苗族自治州平塘县，被称为“中国天眼”的世界上最大的单口径射电望远镜——FAST（500米口径球面射电望远镜）已于2016年9月25日落成启用。FAST的计算速度需达到每秒200万亿次以上，存储容量需达到10PB<sup>①</sup>以上。这一世界级的工程将帮助人们捕捉到更多来自宇宙的信息，它的背后是“天文级”的海量数据存储和复杂的计算。

随着时间的推移、科学任务的深入，以及数据的大量采集，未来对计算速度和存储容量的需求将爆炸式增长，数据量和计算量都将“大得惊人”。

数据是没有边际的，而计算力、存储力始终存在物理极限。

---

<sup>①</sup> 1PB=2<sup>50</sup>B。