

大数据科学的认知和数学基础引论

WANG Yingxu^{1,2}, 彭军³

1. 清华大学大数据系统软件国家工程重点实验室, 清华大学北京信息科学与技术国家研究中心, 北京 100084
2. 国际认知信息学与认知计算学会; 卡尔加里大学电气与计算机工程系, Schulich 工程学院, Hotchkiss 脑科学研究所, 加拿大卡尔加里 T2N 1N4
3. 重庆科技学院智能技术与工程学院, 重庆 401331

摘要 大数据不仅在科学、工程与计算智能中有着广泛的应用, 而且在人类感知、估计、量化、记忆和推理的认知机制中发挥着基础性作用。通过对大数据科学理论的基础研究, 提出一组大数据系统的一般原理和分析方法。为了从形式上解释大数据的起源和本质, 探讨大数据的认知基础及其数学模型, 严格地引出了根植于科学、工程和社会各个领域的大数据的一般模式。研究发现大数据不再是传统实域上的纯数, 而是一个前所未有的新型数学结构, 称为递归类型化超结构(RTHS)。这一大数据系统的基本拓扑特性揭示了大数据工程的复杂性及其操作与处理的全新认知、理论挑战, 以及可选解决方案。

关键词 大数据科学; 大数据工程; 大数据数学模型; 递归超结构; 认知计算; 计算智能

大数据是人类社会信息化时代的代表现象之一^[1-7]。由于大脑中所表示的认知对象可以根据其抽象层次从下至上分为数据、信息、知识和智能4种形式^[1,6,8-20], 人类活动的几乎所有领域都会产生不断增长的数据。数据在人类认知机制的基本层面如感觉、估计、量化、推理和与现实世界互动等方面均发挥着不可或缺的作用。

数论和数域曾经不断扩展, 经历了从无类型数(\emptyset)、二进制数(\mathbb{B})、自然数(\mathbb{N})、整数(\mathbb{Z})、实数(\mathbb{R})、复数(\mathbb{C})、模糊数(\mathbb{F}), 到超结构数(\mathbb{H} ^[19,21], 简称超数)的演化, 如图1所示^[4-5,13-16,22-24]。图1展示了

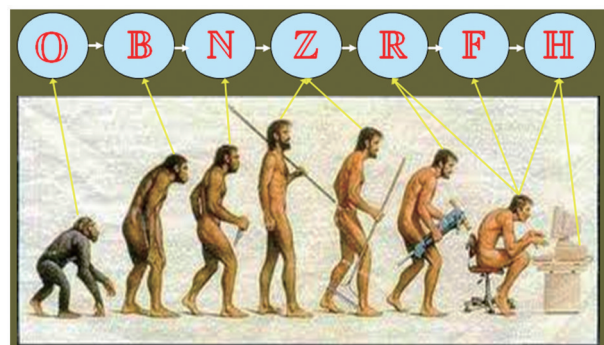


图1 数域与人类数据量化方法的扩展

收稿日期: 2019-11-09; 修回日期: 2020-01-19

基金项目: 国家重点研发计划项目(2016YFB0501504); 国家自然科学基金项目(U1509213)

作者简介: WANG Yingxu, 教授, 研究方向为认知信息学、软件科学、大数据代数和指称数学, 电子信箱: yingxu@ucalgary.ca

引用格式: WANG Yingxu, 彭军. 大数据科学的认知和数学基础引论[J]. 科技导报, 2020, 38(3): 35-46; doi: 10.3981/j.issn.1000-7857.

2020.03.002

人类为了处理现实世界中的实体及其在大脑中的严格表示而对抽象和量化方法的需求的进展过程。大数据科学中的重要发现之一^[19]是数域的需求已经远远超出了传统的 \mathbb{R} 域中的纯数。譬如,近代的模糊数域 \mathbb{F} 的特征是一个二维结构^[25-26],即 $\mathbb{F} = \mathbb{R} \times \mathbb{I} = \{(\mathbb{R}, \mathbb{I})\}$,其中每个成员被表示为在 $[-\infty, +\infty]$ 中的一个变量和在 $(0,1]$ 中的一个关联的隶属度。然而,超数 \mathbb{H} 是数域最新和最一般化的扩展,其可被形式表示为一个类型化 n 元组^[4-5],如定义13所示。

定义1 数据(data)是现实实体或认知对象的量在特定量化尺度下的抽象表示,因而其是有物理量纲和数域类型的数字(digit)。

数据的抽象表达和经典关系模型由 Ullman^[17]和 Codd^[27]提出,其在计算机科学和传统数据库理论中得到了广泛的研究^[3-4,16,22,28-31]。由于数据集的庞大性和复杂性已经超出了传统理论和技术的处理能力,自20世纪90年代以来,科学和工程界就面临着大数据的挑战。“大数据”一词最早是由 Mashey 在1998年提出^[32]。McKinsey 和 Gartner 从工业应用的角度探讨了大数据的高价值,而不仅仅是它的规模^[32]。大数据在科学、工程和社会中无处不在的普遍性在随后的文献中也得到重视^[1-2,32-33]。

定义2 大数据(big data)是在量和质、多样性、存储、检索、提取、计算、语义认知、维护和处理诸方面有别于无类型离散数字的超大规模异构类型量化量所构成的一个超结构(hyperstructure)。

大数据的基本特征是非结构化、异构、单调增长、非描述性、混合/模糊语义,且一致性随时间衰减或熵随时间增加^[23,34]。由于这些固有的复杂性和极大规模的多维超结构对象,大数据科学的理论和大数据工程的方法均面临前所未有的问题。

定义3 大数据工程是用于处理大数据中固有的复杂性和应用需求的系统方法,涵盖大数据系统表示、获取、存储、组织、操作、搜索、分发、标准化、一致性和安全性。

由于大数据呈指数级增长的复杂性和需求,在大数据工程的各个方面和阶段都遇到了前所未有的挑战。这就要求对大数据工程进行理论和形式化研究,从而建立大数据科学。

定义4 大数据科学是关于大数据系统的原理、性质、理论、数学模型和方法论,以及大数据工程的有效组织技术的一门抽象科学。

根据定义1~定义4,基于对大数据本质的研究,大数据科学已经像系统科学、数学科学、软件科学、知识科学、智能科学和认知科学一样,成为当代抽象科学的一门新兴学科^[35]。

本文探讨大数据作为人类推理和工程应用的基本认知对象的哲学、认知、计算和数学基础,以便揭示大数据的认知模型及其计算本质,诸如其类型属性、超结构模型和递归层次操作。

1 大数据科学的认知基础

数据是大脑中最基本的认知对象,因为它通过感觉和量化将现实世界的实体及其属性与抽象认知联系起来^[1,15-16,19,23-24,31,36-42]。数据起源于对实体、属性及其关系的量化和抽象。复杂数据由算术、函数、分析、统计、数值和计算操作生成。数据的演绎分析和归纳综合导致了信息、知识和智能,体现了大数据系统的认知价值。

1.1 数据的认知模型

根据定义1,数据是抽象世界中无形的对象。它们是现实实体的属性或关系及其符号表示的抽象量化。数据的实质曾在概念代数^[19]中研究如下。

定义5 数据的认知模型 C_d ,可被形式化地描述为知识 K 背景下的属性(A)、对象(O)及其与外部概念相对的内部(R^c)、输入(R^i)、输出(R^o)关系的一个超结构,即

$$C_d(\text{数据}: A, O, R^c, R^i, R^o) = \begin{cases} A = \{cognitive_object^*, abstraction, quantity, fact, \\ \quad figure, unit\} \\ O = \{sensory_d, observed_d, experimental_d, measured_d, \\ \quad typed_d, derived_d, statistical_d\} \\ R^c = O \times A \\ R^i \subseteq \times C_d \\ R^o \subseteq C_d \times \end{cases} \quad (1)$$

在上述数据的形式概念模型中,主属性 *cognitive_objects** 表示 C_d (数据)所属的概括概念。 A 中的其他属性表示该概念的蕴含。 O 表示该概念的一组外延或对象。形式概念 C_d (数据)的认知单位 $|C_d|$ 由其内部关系的尺度决定,即 $|C_d| = |R^c| = |O| \times |A| = 7 \times 6 = 42 \text{bir}$, 其中形式概念的认知单位可由二元关系(binary relation, bir)度量^[5]。

二元关系简称为 *bir*(比尔),其作为知识科学

的基本理论^[5]和知识的基本单位被 Wang 在 2018 年发现^[5]。值得注意的是,它与数据和信息的单位 bit (比特)^[34]有着根本的不同。

1.2 数据作为现实实体和抽象属性的量化

数据的普遍性建立在人类认知过程、形式推理、系统量化、数学运算、计算操作和信息处理的普遍性之上。数据如表 1 所示可以分为观察类、推断类、工程类和社会类^[19-20]。

表 1 数据的分类和起源

序号	类别/来源	类型	属性/特性
1.1	观察数据	事实	特征、序数、基数、计数、数量
1.2		状态	存在形式、构造、状态变化
1.3		行为	互动、规范、周期、分布、频率
2.1	推断数据	类比	相似、比较、等价类
2.2		关系	关联、因果、映射、序列、并发
2.3		测量	量化、限定、缩放、标准化、加权、分类
2.4		语义	物理、经验、抽象、数学、规则、性质
2.5		数学	线性、非线性、多项式、解析、微分、积分、细化、函数系统
2.6		统计	概率、范数、偏差、分布、条件、随机过程
2.7		复合	插值、外推、指数、幂函数、阶乘、笛卡尔积、搜索、排序、组合、排列、系统融合
3.1	工程数据	信息系统	信息技术、通信、计算、搜索引擎、多媒体、视频、图像采集
3.2		数字化	事实、统计、信号、媒体、电视、照片、音频、序列、比特流
3.3		数据库	数据中心、公共/工业/个人数据库、云服务器
3.4		知识库	一般、学科、学术、教育、道路、地球、宇宙
3.5		应用领域	工业、搜索引擎、智能系统、媒体、图像、视频、语音、文本、语言
4.1	社会数据	社交网络	商业、个人、学术、兴趣团体、意见收集
4.2		政府	行政、人口普查、投票、法律、图书馆
4.3		数据中心	全球、区域和城市中心、多对多分布中心

大数据的根本来源是人类的集体智慧。产生大数据的典型人类活动有多对多通信、海量数据复制下载、数字图像采集和网络舆情形成等。产生大数据的典型计算是笛卡尔积 ($O(n^2)$)、排序 ($O(n \lg n)$)、搜索(穷举, $O(n^3)$)、知识库更新 ($O(n^2)$), 以及复杂度为 $O(2^n)$ 、 $O(n!)$ 或更高的诸如排列和 NP 问题^[22,43-44]。

大数据的数学性质是对现实世界实体或认知对象的数量的抽象表示,通过量化映射到一定的尺度上。大数据也可能通过极简数学运算产生,比如实域内最大的数据是 $1/0 = \infty$ 。创建大数据的典型数学运算有排列、组合、幂函数、阶乘、发散结构(如

编码器和逻辑门等)。例如:(1) 欧拉数是 $10! = 26525285981219105863630848000000$; (2) 2016 年已知的最大素数^[45] 为 $2^{74207281} - 1$; (3) 有 10000 个输入的与门的状态空间为 2^{10000} 比特,而其语义空间高达 $2^{\lceil I_n \rceil} (I_n - 1) = 2^{\lceil 2^{10000} \rceil} (2^{10000} - 1) \approx 2^{2^{10000} + 10000} \text{bir}$ ^[35]。

定义 6 数据的数学模型 D 是现实实体或认知对象的数量 Q 在度量尺度 S 上的抽象映射,其导出量是一个实数 $D \in \mathbb{R}$, 但其所载有的量纲是 S , 即

$$D \triangleq f_s: Q \rightarrow S \quad (2)$$

定义 7 量化 $\Gamma_S(X)$ 是将未知量 X 映射到给定度量尺度 S 上的一个度量函数及其过程, 即

$$\Gamma_S(X) \triangleq f_s: X \rightarrow S, S \in \mathbb{N}$$

$$= \frac{Q(X)}{[S]} [S] \quad (3)$$

式中, $[S]$ 表示量化单位或量纲。

引理 1 相对于测量标度 S , 数量 X 所对应的数据 $D(X, S)$ 是通过量化 $\Gamma_s(X)$ 产生的实数 $I_x \cdot R_x$ 。在单位 $[S]$ 中, $I_x \cdot R_x$ 由用小数点分隔的两部分组成, 其中 I_x 为整数, R_x 为小于 $[S]$ 的余数。

$$\begin{aligned} D(X, S) \triangleq \Gamma_s(X) &= \frac{Q(X)}{[S]}, S \in \mathbb{N} \\ &= I_x \cdot R_x [S], I_x \in \mathbb{Z}, 0 \leq R_x \in \mathbb{R} < S \quad (4) \\ &= \begin{cases} I_x = \text{mod}(X) \\ R_x = \text{rem}(X) \end{cases} \end{aligned}$$

式中, 数据对象由称之为(整数, 余数)的量化对偶 (I_x, R_x) 表示, 分别由模(mod)和余数(rem)操作生成。

例 1 给定一实体近似长度为 $L \approx 1286 \text{ mm}$ 。根据引理 1 和式(4), 与这一实体相应的数据可以在给定的测量尺度 S 上生成如下:

$$\begin{aligned} S_1 = 1 \text{ m}: D_1 &= \frac{Q(L)}{1} = \frac{Q(L)}{1[\text{m}]} = \frac{1286}{1000} = 1.286 \text{ m} \\ S_2 = 1 \text{ mm}: D_2 &= \frac{Q(L)}{2} = \frac{Q(L)}{1[\text{mm}]} = \frac{1286}{1} = 1286 \text{ mm} \\ S_3 = 1 \text{ }\mu\text{m}: D_3 &= \frac{Q(L)}{3} = \frac{Q(L)}{1[\mu\text{m}]} = \frac{1286}{0.001} = 1286000 \text{ }\mu\text{m} \end{aligned}$$

式中, D_2 表示测量标度 S_2 和量化标度 S^0 相同的特殊情况, 其使量化数据 D_2 等于被测对象 L 。

1.3 数据作为认知和推理的基本对象

在人类认知的层次结构中, 数据、信息、知识和智能之间的关系可以形式化地描述如下。

定义 8 人类大脑中的认知对象 H_{co} 的层次结构是一个四元组, 从下至上分为数据(\mathbb{D})、信息(\mathbb{I})、知识(\mathbb{K})和智能(\mathbb{I}), 即

$$\begin{aligned} H_{co} &\triangleq (\mathbb{D}, \mathbb{I}, \mathbb{K}, \mathbb{I}) \\ &= \begin{cases} \mathbb{D} = f_d: O \rightarrow Q \\ \mathbb{I} = f_i: D \rightarrow S \\ \mathbb{K} = f_k: I \rightarrow C \\ \mathbb{I} = f_i: I \rightarrow B \end{cases} \quad (5) \end{aligned}$$

式中符号分别表示对象(O)、数据(D)、数量(Q)、信息(I)、语义(S)、概念(C)和行为(B)。

基于定义 8, 数据作为人类最基本的认知对象之一具有一系列如文献[24]和[34]中所确定的 20 个基本认知属性, 例如抽象性、概括性、累积性、对认知的依赖性、类型化的元语义、可共享性、无空间尺度、无重量、数据到信息/知识/智能之间的可转换性、多种表示形式、多种承载介质、多种传输形式、对介质的依赖性、对能量的依赖性、无磨损性、时间依赖性、熵守恒、不同于物理属性的信息质量属性、易失真和稀缺性。

2 大数据科学的计算基础

从计算的角度来看, 数据的一个重要抽象属性是它们的类(type)或形态。在类型理论中^[5,16], 类决定了某一形态数据的域、单位和被允许的操作。一个类型系统指定了所有数据对象类的建模和操作规则^[3,16,29,45]。在一给定完备类型系统中, 任何数据对象都可以被指定到一个类型、一个有限类型集, 或它们的组合。每种类都是一个集合, 其中所有数据对象共享一组公共属性、域约束和预定义操作。类可以分为基本和复杂类型。前者是一组最简类型, 其不能被进一步约减; 而后者是由基本类型按照一定的类规则组合而成的混合类型。

定义 9 数据对象 O 是一个有指定类 \mathbb{T} 的变量 $v, v \in V \subset \bar{P}V \subset U$, 其由类型约束 $\mu_v(\mathbb{T})$ 限制以便将此一般类 \mathbb{T} 通过子域裁剪 \mathbb{T}'' 转化为一给定问题的特定类 \mathbb{T}' , 即

$$O \triangleq \langle v: \mathbb{T} \mid \mu_v(\mathbb{T}) = \mathbb{T} \setminus \mathbb{T}'' = \mathbb{T}' \rangle \quad (6)$$

式中, $\mu_v(\mathbb{T}) = \mathbb{T} \setminus \mathbb{T}'' = \mathbb{T}'$ 将给定数域 D 中的 \mathbb{T} 裁剪成问题域 D_p 中的子集 \mathbb{T}' , 即 $v \mid \mathbb{T}' = v \mid \mathbb{T} \setminus v \mid \mathbb{T}''$ 。大数据的论域 U 将在定义 14 中详述。

定义 10 数据对象的基本类型 T_p 是一组完备的且在不丢失逻辑或语义属性的限制下不可再分解的数据类, 如下所示:

$$T_p \triangleq \{N, Z, R, S, L, B, H, P\} \quad (7)$$

式中的符号分别表示诸如自然数(N)、整数(Z)、实数(R)、字符串(S)、逻辑变量(L)、字节(B)、十六进制数(H_x)和指针(P)的类型。

表2不但给定了数据的基本类型,而且确定了它们的表示法、域和属性。表2中的每一个类域都可以用其数学域(\mathbb{T})和问题域(\mathbb{T}')表示,其中问题域 \mathbb{T}' 受给定问题或机器内存(M_{max})空间限制。

表2 数据对象的基本类和域

序号	类型	语法	域	等效性
1	自然数	N	$N \in [0, 65535], N' \in [0, M_{max}-1]$	
2	整数	Z	$Z \in [-32768, +32767]$ $Z' \in [-(M_{max}/2), +(M_{max}/2)-1]$	默认的 算术运
3	实数	R	$R \in [-2147483648, 2147483647]$ $R' \in [-(M_{max}/2), +(M_{max}/2)-1]$	算
4	字符串	S	$S \in [0, 255], S' \in [0, M_{max}-1]$	默认的 字符和 字符串 操作
5	逻辑 变量	L	$L \in \{T, F\}$	布尔常 量 [TIL, FIL]
6	字节	B	$B \in [0, 255], B' \in [0, 1023]$	默认的 二进
7	16 进制	H	$H \in [0, 65535], H' \in [0, M_{max}-1]$	制 操作

虽然十进制系统广泛用于面向人的数据表示,但最精细的数据类是比特(bit)^[13],这是计算机与信息科学融合的基础。复杂的数据类型和系统可以基于比特位进行组合。然而最新发现表明,知识科学与智能科学的基本单位是二进制关系比尔(bir)^[5],这种重要的类为知识科学和智能科学提供了理论基础。

定义 11 数据对象的复杂类型 T_c 包含一组由基本类型派生的组合类型,即

$T_c \triangleq \{\mathbb{T}, SM, TX, A, F, V, T, D, TM, Bir\}$ (8)

式中的符号分别表示诸如任意(\mathbb{T} ,在运行时确定)、结构模型(SM)、文本(TX)、音频(A)、照片(F)、视频(V)、时间(T)、日期(D)、日期-时间(TM),和知识(Bir)的复杂类型。其中Bir是基于单位bir的知识类型。

复杂数据对象的组合类型在表3中给出,用以指定了它们的类、符号、域和属性。在表3中有一些复杂类型是问题相关的,比如{TX, A, F, V}。为方便数据表示和定性,在本文数据对象 x 的组合形式中采用了一种类型后缀约定,其由竖线分隔,即

$x|T$,其中 $|T$ 可以由表2、表3中的任何类型或基于它们的用户定义类型来表示。

表3 数据对象的复杂类和域

序号	类型	符号	域	等效性
1	任意 类型	\mathbb{T}	-	运行时 确定的 任何虚 拟类型
2	结构	SM	$SM \in \tau^* T(Eq. 11) $	默认字 段引用: $x SM.y T$
3	文本	TS= S'	$S' \in [0, M_{max}-1]$	
4	音频	$A=B \times T$	$B' \in [0, 1023],$ $T \in [hh:mm:ss:ms]$	问题相 关
5	照片	$F=B \times B'$	$B' \in [0, 1023]$	
6	视频	$V=$ $B' \times B \times T$	$B' \in [0, 1023], T \in [hh:mm:ss:ms]$	
7	时间	$T=hh:mm:ss:ms$	$hh \in [0, 23], mm \in [0, 59],$ $ss \in [0, 59], ms \in [0, 999]$	
8	日期	$D=YY:MM:DD$	$YY \in [0, 99], MM \in [1, 12],$ $DD \in [1, 31]$	默认的 时间操 作
9	日期/ 时间	$TM=$ $YYYY:MM:DD:hh:mm:ss:ms$	$YYYY \in [0, 9999], MM \in [1, 12],$ $DD \in [1, 31], hh \in [0, 23],$ $mm \in [0, 59], ss \in [0, 59],$ $ms \in [0, 999]$	
10	知识	Bir	$[0, \infty)$	数据的 语义

定义 12 大数据的类型系统 T 是表2、表3所示的17种基本类型和复杂类型的集合,即

$$T \triangleq T_p \cup T_c$$

$$= \{N, Z, R, S, L, B, H, P\} \cup \{\mathbb{T}, SM, TX, A, F, V, T, D, TM, Bir\}$$
 (9)

例 2 大数据应用中任何用户定义的类型都可以通过基本类型、复杂类型及其组合在类型环境 T 中进行形式化描述。例如,多媒体社交网络中的典型大数据对象类型,如文本(电子信件、消息)、语音、照片、视频、知识等均可严格定义如下:

$$\begin{aligned} \text{Text}|TX: & TX \triangleq S \\ \text{Voice}|A: & A \triangleq B \times T \\ \text{Photo}|F: & F \triangleq B \times B \end{aligned}$$
 (10)

VideolV: $V \triangleq B \times B \times T$

KnowledgeBir: $K \triangleq C \times C$

式中, |TX, |A, |F, |V 和 |Bir 是复杂类型或用户定义类型的后缀, 然而 C 是一个如定义 5 所示的认知概念。

大数据科学的另一发现揭示了大数据的形式载体是一种称为类型化元组的超结构^[9]。

定义 13 数据和数量的超结构类型是一个通用的类型化的 n 元组, $\tau^n |T$, 其在满足一定约束条件 \mathbb{T}_i' 下, 将 n 个数据对象封装在异构类 $\mathbb{T}_i (1 \leq i \leq n)$ 的超结构中, 即

$$\begin{aligned} \tau^n |T &\triangleq \left(\overset{n}{R} O_i | SM \right) \\ &= \left(\overset{n}{R} \langle v_i; \mathbb{T}_i' \mid \mathbb{T}_i = \mathbb{T}_i \setminus \mathbb{T}_i' \rangle \right) \end{aligned} \quad (11)$$

式中, $\overset{n}{R} S_i$ 表示一个循环结构或迭代行为^[16,46]; \mathbb{T}_i' 是对标准类型 \mathbb{T}_i 的一个类型约束。

数域的类型约束广泛用于指定用户定义域或特定于问题的类型, 其是大数据建模中用于构造或细化数据类型的一般方法。

3 大数据科学的数学基础

数据是将实体的属性映射到一个度量尺度的量化的结果。大数据是在科学、工程和社会学科中产生的大规模、异构数据的超结构复杂系统。本节将在上述概念模型的基础上, 建立一组大数据的严格数学模型。

3.1 大数据系统的论域

大数据科学的数学结构的整体性质和公理可以通过对大数据系统的抽象模型的归纳而得出。该形式化方法从研究抽象数据和大数据系统的论域开始。

定义 14 大数据与复杂抽象数量的论域 U 是一个六元组:

$$U \triangleq (E, T, Q, R, V, H) \quad (12)$$

式中, E 是一实体和/或其可测量属性集; T 是 E 的类型或性质集; Q 是 E 上的一量化标度集; R 是一

量化关系 ($Q \times E$) 或限定关系 ($E \times T$) 集; V 是一类型化的值集 $Q \times E \times T \rightarrow V | T$; H 是一超结构集 $E \times V \times T$ 。

论域 U 的定义揭示一般大数据比定义 2 中给出的简单数据更加复杂。在 U 的基础上, 将形式化引出大数据系统的数学模型, 分为基本和一般大数据模型。

3.2 大数据系统的数学模型

大数据系统的数学模型可以由引理 1 和定义 14 导出。在简单数据单维结构的基础上, 基本大数据系统可被描述为二维类型化超结构; 而一般大数据系统将被定义为多维类型化超结构。

定义 15 基本大数据系统 (BDS) 的数学模型 $\Theta^2 = \overset{n}{R} \overset{m}{R} d_{ij} | \mathbb{T}_{0j}$ 是论域 U 中的一个二维 $n \times m$ 类型化超结构, 其中数据元素 $d_{ij} | \mathbb{T}_{0j}$ 由类型后缀 $| \mathbb{T}_{0j}$ 指定, 可以是任意一种定义在 T 中的初级或复杂类型, 即

$$\begin{aligned} \Theta^2 &\triangleq \overset{n}{R} \overset{m}{R} d_{ij} | \mathbb{T}_{0j} \\ &= \begin{bmatrix} \theta_0 & e_1 | \mathbb{T}_{01} & e_2 | \mathbb{T}_{02} & \cdots & e_m | \mathbb{T}_{0m} \\ \kappa_1 & d_{11} | \mathbb{T}_{01} & d_{12} | \mathbb{T}_{02} & \cdots & d_{1m} | \mathbb{T}_{0m} \\ \kappa_2 & d_{21} | \mathbb{T}_{01} & d_{22} | \mathbb{T}_{02} & \cdots & d_{2m} | \mathbb{T}_{0m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \kappa_n & d_{n1} | \mathbb{T}_{01} & d_{n2} | \mathbb{T}_{02} & \cdots & d_{nm} | \mathbb{T}_{0m} \end{bmatrix} \end{aligned} \quad (13)$$

式中, $\overset{n}{R} \kappa_i | SM$ 中的每一行表示一个称作为结构模型 ($| SM$) 的类型化元组, 而 $\overset{m}{R} e_j | \mathbb{T}_{0j}$ 中的每一列代表一个具有特定类型 ($| \mathbb{T}_{0j}$) 的域^[47], 其中大 R 算符 $\overset{n}{R} X_i$ 用以表示循环结构或重复行为^[46]。

如定义 15 所示, 大数据 Θ^2 数学模型中的第 0 行 θ_0 是一个特殊的类型化元组, 称为 BDS 的模式。模式 θ_0 指定 BDS 各数据域的结构和约束类。值得注意的是, BDS 数学模型中采用的类后缀 $| \mathbb{T}$ 可用于定义任意类型的大数据, 以便满足对于容纳广泛异构数据的需求。

例 3 给定一个有 100 万用户和 7 个数据域的社交网络大数据系统 BDS_1 。其二维大数据模型的

形式结构 $\Theta^2(BDS_1)$ 可以根据式 (13) 严格描述如下:

$$\theta_0(BDS_1) = (ID|N, UName|S, GName|S, Text|T, Voice|A, Photo|F, Video|V)$$

$$\Theta^2(BDS_1) = \prod_{j=0}^{1000000} R d_{ij} | \mathbb{T}_j$$

$\theta_0(BDS_1)$	ID N	UName S	GName S	Text T	Voice A	Photo F	Video V
κ_1	0000001	John	G_{0001}	$R_t 0000001 T$	$R_a 0000001 A$	$R_p 0000001 F$	$R_v 0000001 V$
κ_2	0000002	Judy	G_{0301}	$R_t 0000002 T$	$R_a 0000002 A$	$R_p 0000002 F$	$R_v 0000002 V$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$\kappa_{1000000}$	1000000	Mike	G_{1806}	$R_t 1000000 T$	$R_a 1000000 A$	$R_p 1000000 F$	$R_v 1000000 V$

(14)

式 (13) 所定义的大数据系统基本模型 Θ^2 可扩展到一般的 n 维类型化超结构大数据系统 Θ^n 。

定义 16 大数据系统 (BDS) 的一般数学模型为 $\Theta^n = \prod_{i_1=0}^{n_1} \prod_{i_2=0}^{n_2} \cdots \prod_{i_q=0}^{n_q} R d_{i_1 i_2 \dots i_q} | \mathbb{T}_{i_1 i_2 \dots i_q}$ 。该模型在 U 中任何 q 维异构数据可被形式化描述成一个递归的类型化超结构 (RTHS), 即

$$\begin{aligned} \Theta^n &\triangleq \prod_{k=q}^1 R \Theta^k (\Theta^{k-1}) \\ &= \prod_{i_1=0}^{n_1} \prod_{i_2=0}^{n_2} \cdots \prod_{i_q=0}^{n_q} R d_{i_1 i_2 \dots i_q} | \mathbb{T}_{i_1 i_2 \dots i_q} \\ &= \prod_{k=q}^1 \begin{bmatrix} \theta_0^k & e_1^k | \mathbb{T}_{01}^k & e_2^k | \mathbb{T}_{02}^k e_2 & \cdots & e_m^k | \mathbb{T}_{0m}^k & \theta_0^{k-1} | P \\ \kappa_1^k & \tau_{11}^k | \mathbb{T}_{01}^k & \tau_{12}^k | \mathbb{T}_{02}^k & \cdots & \tau_{1m}^k | \mathbb{T}_{0m}^k & \kappa_1^{k-1} | P \\ \kappa_2^k & \tau_{21}^k | \mathbb{T}_{01}^k & \tau_{22}^k | \mathbb{T}_{02}^k & \cdots & \tau_{2m}^k | \mathbb{T}_{0m}^k & \kappa_2^{k-1} | P \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \kappa_n^k & \tau_{n1}^k | \mathbb{T}_{01}^k & \tau_{n2}^k | \mathbb{T}_{02}^k & \cdots & \tau_{nm}^k | \mathbb{T}_{0m}^k & \kappa_n^{k-1} | P \end{bmatrix} \end{aligned} \quad (15)$$

$$\begin{aligned} \Theta^3(BDS_2) &= \prod_{k=3}^1 R \Theta^k (\Theta^{k-1}), \quad \Theta^0 = \prod_{i=0}^{1000000} R d_{ij} | \mathbb{T}_{0j} \\ &= \prod_{i_0=0}^{n_0} \prod_{i_1=0}^{n_1} \prod_{i_2=0}^{n_2} R d_{i_0 i_1 i_2} | \mathbb{T}_{i_0 i_1 i_2} \end{aligned}$$

$$= \Theta^3(BDS_2) : \begin{bmatrix} \theta_0^3(BDS_2) & ID|N & UName|S & GName|S & \theta_0^2|P \\ \kappa_1^3 & 0000001 & John & G_{0001} & $\kappa_1^2|P$ \\ \kappa_2^3 & 0000002 & Judy & G_{0301} & $\kappa_2^2|P$ \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \kappa_{1000000}^3 & 1000000 & Mike & G_{1806} & $\kappa_{1000000}^2|P$ \end{bmatrix}$$

式中, 大数据系统模式 $\Theta^2(BDS_1)$ 的 7 个域分别代表识别号(IN)、用户名(IS)、组名(IS)、文本数据(IT)、语音数据(IA)、照片数据(IF)和视频数据(IV)。 $\Theta^2(BDS_1)$ 大数据中如式 (14) 所示的每一记录都受公共模式 $\theta_0(BDS_1)$ 的约束。

式中, $\prod_{j=1}^m R \kappa_j^k | SM = (\tau_1^k | \mathbb{T}_{01}^k, \tau_2^k | \mathbb{T}_{01}^k, \dots, \tau_m^k | \mathbb{T}_{01}^k)$ 表示通过指针 $\theta_0^{k-1} | P$ 递归链接到大数据系统 $\prod_{i=1}^n R \kappa_i^{k-1} | SM$ 中低层结构的一个类型化元组。

通过对比定义 15 和定义 16, 可见二维基本大数据模型 Θ^2 是一般大数据模型 Θ^n 的一个特例。一般大数据模型 $\Theta^n(BDS)$ 中的元素除了终端层外都是一个类型元组, 而基本大数据模型 $\Theta^2(BDS)$ 中的元素是一个终端数据对象。因此, 根据定义 16, 任何前者可以通过有限步转变为后者, 如果 $\Theta^0(BDS)$ 已给定。

例 4 在例 3 的基础上, 建立一个给定社交网络 $\Theta^3(BDS_2)$ 的一般大数据模型, 其中 $\Theta^0 = \prod_{i=0}^{1000000} \prod_{j=0}^4 R d_{ij} | \mathbb{T}_{0j}$ 。按照定义 16, 该问题的解是将 $\Theta^2(BDS_1)$ 分层细化为一个 3 层一般大数据模型, 如式 (16) 所示。

$$\begin{array}{l}
\Downarrow \Theta^2(BDS_2): \\
\begin{array}{|c|c|c|c|c|}
\hline
\theta_0^2(BDS_2) & Rec\#|N & TStamp|TM & Size|B & \theta_0^1|P \\
\hline
\kappa_1^2 & 0000001 & TM_1 & S_1 & \kappa_1^1|P \\
\kappa_2^2 & 0000002 & TM_2 & S_2 & \kappa_2^1|P \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
\kappa_{1000000}^2 & 1000000 & TM_{1000000} & S_{1000000} & \kappa_{1000000}^1|P \\
\hline
\end{array} \\
\Downarrow \Theta^1(BDS_2): \\
\begin{array}{|c|c|c|c|c|}
\hline
\theta_0^1(BDS_2) & Text|T & Voice|A & Photo|F & Video|V \\
\hline
\kappa_1^1 & R_t 0000001|T & R_a 0000001|A & R_p 0000001|F & R_v 0000001|V \\
\kappa_2^1 & R_t 0000002|T & R_a 0000002|A & R_p 0000002|F & R_v 0000002|V \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
\kappa_{1000000}^1 & R_t 1000000|T & R_a 1000000|A & R_p 1000000|F & R_v 1000000|V \\
\hline
\end{array}
\end{array} \quad (16)$$

在 $\Theta^3(BDS_2)$ 中, 每一层自顶向下通过指针 $\theta_0^{k-1}|P$ 链接到下一层, 从而任一 k 层都将被递归地细化 (\Downarrow) 为一个基本二维结构。 $\Theta^3(BDS_2)$ 的顶层设定为 $\theta_0^3(BDS_2)$, 并通过 $\theta_0^2|P$ 将其链接到下一层的二维结构 $\theta_0^2(BDS_2)$ 。第二层对 $\theta_0^2(BDS_2)$ 进行细化, 并通过 $\theta_0^1|P$ 将其链接到下一层的二维结构 $\theta_0^1(BDS_2)$ 。然后, 第一层将 $\Theta^1(BDS_2)$ 表示为在第 0 层给出的一组终端数据对象。

上述定义所给大数据系统的一般数学模型 (Θ^q) 可以用递归层次结构清晰严格地表达任意 q 维大数据结构, 其中每一层都是一个统一的二维基本子结构 (Θ^2)。例 4 显示大数据系统 Θ^q 不再是简单的纯数字一维或二维结构, 而是一个递归类型化的超结构 (RTHS), 其可分解为一组层次化链接的二维超结构 Θ^2 。

大数据系统的一般数学模型 (Θ^q), 不但揭示了大数据的递归超结构的数学实质, 而且为现实世界中任何复杂大数据系统的建模和分析提供了一种严格、通用、灵活和高效的层次细化和综合递归方法。该理论也为大数据在超 R 域中的数学分析 (简称大数据代数^[19]) 奠定了基础。

4 大数据系统的性质

4.1 大数据系统的递归超结构特性

定理 1 大数据系统的一般递归超结构性质可由一个 q 维层次结构 Θ^q 形式化地表述。当最底

层终端结构 Θ^0 已给定为 m 个数据对象的类型化元组, 任一高层结构可以由低层结构递归得到,

$$\Theta^q = \underset{k=1}{\overset{q}{R}} \Theta^k (\Theta^{k-1}), \quad \Theta^0 = \underset{j=0}{\overset{m}{R}} d_j | \mathbb{T}_j \quad (17)$$

证明 定理 1 可以通过定义 15、定义 16 归纳证明如下:

$$\begin{aligned}
\forall \Theta^0 &= \underset{j=0}{\overset{m}{R}} d_j | \mathbb{T}_j, \\
\Theta^q &= \begin{cases} \Theta^1 = \Theta^0 = \left(\underset{j=0}{\overset{m}{R}} d_j | \mathbb{T}_j \right) \\ \Theta^2 = \Theta^1 (\Theta^0) = \Theta^1 \left(\underset{j=0}{\overset{m}{R}} d_j | \mathbb{T}_j \right) \\ \dots \\ \Theta^q \left(\Theta^{q-1} \left(\dots \left(\Theta^1 \left(\underset{j=0}{\overset{m}{R}} d_j | \mathbb{T}_j \right) \right) \dots \right) \right) \\ = \Theta^q \left(\Theta^{q-1} \left(\dots \left(\Theta^1 (\Theta^0) \right) \dots \right) \right) \\ = \underset{k=1}{\overset{q}{R}} \Theta^k (\Theta^{k-1}) \end{cases} \quad (18)
\end{aligned}$$

上述定理是大数据科学基于最新系统科学原理^[35]所得出的另一发现。在定理 1 的基础上, 可以推导出如下 BDS 层次细化和层次抽象的原则。

推论 1 在 U 中任何一个 BDS 都可以通过自上而下的低层结构以不断增加的细节来演绎和分析, 即

$$\Theta^q = \underset{k=q}{\overset{0}{\Downarrow}} \Theta^k = \underset{k=q}{\overset{1}{R}} \Theta^k (\Theta^{k-1}), \quad \Theta^0 = \underset{i=0}{\overset{n}{R}} \underset{j=0}{\overset{m}{R}} d_{ij} | \mathbb{T}_j \quad (19)$$

逆对称地, 推论 1 可以表示为如下层次抽象原则。

推论 2 在 U 中任何一个 BDS 都可以通过自下而上的结构归纳和综合而合成, 即

$$\Theta^q = \prod_{k=0}^q \Theta^k = \overset{q}{R} \Theta^k (\Theta^{k-1}), \quad \Theta^0 = \overset{n}{R} \overset{m}{R} d_{ij} | \mathbb{T}_j \quad (20)$$

该推论可以看作是对被称为“信息隐蔽”的系统建模经验原理的形式化表达。

证明 推论 1 和 2 可引用定理 1 证明。

推论 3 一般大数据库 (BDB) 的数学模型 Θ 是根据定理 1 建立的递归超结构类型化数据库, 其二维和 n 维模型如下:

$$\Theta \triangleq \begin{cases} \Theta^2 = \overset{m_1}{R} \overset{m_2}{R} d_{i_1 i_2} | \mathbb{T}_{i_1 i_2} \\ \Theta^n = \overset{m_1}{R} \overset{m_2}{R} \cdots \overset{m_n}{R} d_{i_1 i_2 \cdots i_n} | \mathbb{T}_{i_1 i_2 \cdots i_n} \end{cases} \quad (21)$$

根据定理 1 和推论 1~推论 3, 任何一般的 n 维大数据系统都可以由低阶的 2 维系统递归而成。反之, 任何基本的二维大数据系统都可以由一般的 n 维系统演绎而成。

4.2 数据、信息、知识和智能之间的关系

虽然大数据的基本属性由递归超结构和类型理论主导, 但信息、知识和智能的基本属性分别由组合机制、形式概念和超函数构成。因此, 在人类认知的层次结构中, 数据、信息、知识和智能之间的关系可以通过以下原理来形式化描述。

定理 2 在认知系统 S_ω 的层次结构中, 认知对象之间的可转换性形成了数据 (\mathbb{D})、信息 (\mathbb{I})、知识 (\mathbb{K})、智能 (\mathbb{I}) 之间的递归结构如下:

$$\begin{aligned} S_\omega &= \overset{1}{R} O^k (O^{k-1}), \quad O^0 = \overset{n}{R} d_i | \mathbb{T}_i \\ &= O^4 (O^3 (O^2 (O^1 (O^0)))) \\ &= \mathbb{I} \left(\mathbb{K} \left(\mathbb{I} \left(\mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \right) \end{aligned} \quad (22)$$

其中层次框架的终端层由具有特定类型的 n 维数据对象 $\overset{n}{R} d_i | \mathbb{T}_i$ 给定。

证明 根据定理 1, 大数据系统的认知转换结构 S_ω 可被递归归纳地证明如下:

$$\begin{aligned} \forall O^0 &= \overset{n}{R} d_i | \mathbb{T}_i \\ O^1 &= O^1(O^0) = O^1 \left(\overset{n}{R} d_i | \mathbb{T}_i \right) = f_d \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \\ &= \mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \end{aligned}$$

$$\begin{aligned} O^2 &= O^2(O^1) = O^2(\mathbb{D}) = f_i \left(f_d \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \\ &= \mathbb{I} \left(\mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \end{aligned} \quad (23)$$

$$\begin{aligned} O^3 &= O^3(O^2) = O^3(\mathbb{I}) = f_k \left(f_i \left(f_d \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \\ &= \mathbb{K} \left(\mathbb{I} \left(\mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \end{aligned}$$

$$\begin{aligned} O^4 &= O^4(O^3) = O^4(\mathbb{K}) = f_l \left(f_k \left(f_i \left(f_d \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \right) \\ &= \mathbb{I} \left(\mathbb{K} \left(\mathbb{I} \left(\mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \right) \end{aligned}$$

$$\begin{aligned} \Rightarrow S_\omega &= O^4(O^3(O^2(O^1(O^0)))) \\ &= \mathbb{I} \left(\mathbb{K} \left(\mathbb{I} \left(\mathbb{D} \left(\overset{n}{R} d_i | \mathbb{T}_i \right) \right) \right) \right) \end{aligned}$$

其中每一给定层次对象都是其直接下层对象的合成。

上述定理 2 揭示了认知对象之间的普遍可转换性和对人类归纳能力的高度依赖, 该定理不但涵盖数据的抽象/量化、信息获取、知识生成和大脑智能/智慧创造, 而且也表明大数据的语义推导和归纳是复杂代数运算, 而不是通常认为的逻辑运算。定理 2 不但解释了低层认知对象是高层认知对象的一个不可或缺的基础和范例, 而且解释了高层认知对象是低层认知对象的一个升华和融合。因此, 对于任一给定层次的认知对象, 其不仅是低层认知对象的语义, 也是高层认知对象的实例。

例 5 一个具有 1000 个输入、 $n=1000$ 的逻辑与门的状态空间所能产生的大数据, $O^0 = \overset{n-1}{R} d_i | \mathbb{T}_i =$

$\overset{999}{R} d_i | \mathbb{B}$, 可根据定理 2 确定为

$$\begin{aligned} O^1 &= \mathbb{D}(O^0) = \mathbb{D} \left(\overset{999}{R} d_i | \mathbb{B} \right) = 2^{1000} | \mathbb{B} \\ O^2 &= \mathbb{I}(\mathbb{D}) = \mathbb{I} \left(\mathbb{D} \left(\overset{999}{R} d_i | \mathbb{B} \right) \right) = 1000 | \mathbb{B} \\ O^3 &= \mathbb{K}(\mathbb{I}) = \mathbb{K}(\mathbb{I}(\mathbb{D})) = \{ \text{rule}_1, \text{rule}_2 \} \end{aligned} \quad (24)$$

$$\begin{cases} \text{rule}_1 | \text{Bir}: \forall \overset{999}{R} d_i | \mathbb{B} = 1 \Rightarrow \text{AND} \left(\overset{999}{R} d_i | \mathbb{B} \right) = 1 \\ \text{rule}_2 | \text{Bir}: \exists \overset{999}{R} d_i | \mathbb{B} = 0 \Rightarrow \text{AND} \left(\overset{999}{R} d_i | \mathbb{B} \right) = 0 \end{cases}$$

$$O^4 = \dot{\mathbb{I}}(\mathbb{K}) = \dot{\mathbb{I}}(\mathbb{K}(\mathbb{I}(\mathbb{D}))) : \forall n, \text{AND}\left(\bigwedge_{i=1}^n d_i | B\right) = \bigwedge_{i=1}^n d_i | B$$

例5表明,在对一个认知对象在不同归纳层次上的量化描述,在数据科学的意义上会有极大的不同。在底部的数据层(O^1),所给逻辑与门表示一个具有 2^{1000} 比特状态空间的超大数据系统。然而,在信息层(O^2),它只表示1000比特的信息,即它的输入尺度。在知识层(O^3),它被归纳为两个Bir的规则,即当所有输入为1时,输出为1;而当任一输入为0时,输出为0。最终,在智能层(O^4),无论输入的尺度(n)是多大,此逻辑与门已升华为与 n 无关的一个通用逻辑与函数。

定理2和例5揭示了人类智能是一种非常高效、快速收敛的知识生成和智慧获取的归纳机制,其中数据仅仅是现实世界几乎无限状态空间中的事实材料和个体实例。归纳的认知过程是大数据转化为知识和智慧的一种典型的、人类智能中最强大的推理能力。

6 结论

介绍了一个从认知、数学和计算3个领域对大数据科学的基础研究。发现了大数据系统是一种新型的超越R域的递归类型化超结构(RTHS),创建了大数据系统(BDS)的基本数学模型($\Theta(BDS)$)和大数据科学的数学方法。研究揭示了大数据系统数据的抽象/量化、信息获取、知识生成和智能/智慧创生,均高度依赖于归纳推理机制,为现实世界中任何复杂大数据系统的建模和分析提供了一种通用和高效的层次细化方法。本文对大数据原理与性质的研究,不仅为解释BDS的本质提供了一个新颖的数学结构,也为大数据科学在大数据工程、机器学习和知识提取等新兴领域的广泛应用提供了一种严格的理论基础。

致谢:本文工作获得加拿大自然科学与工程研究委员会(NSERC)发现基金、清华大学大数据系统软件国家重点实验室及重庆科技学院科研基金的部分资助。陆建华院士和孙家广院士对本文撰写提供了精心指导,王建民

和刘琳对专题文章的组织给予了帮助,魏子麒对文本编辑提供了帮助。

参考文献(References)

- [1] Hassanien A E, Azar A T, Snasel V, et al. Big data in complex systems: Challenges and opportunities[M]. Berlin: Springer, 2015.
- [2] Jacobs A. The pathologies of big data[J]. Queue, 2009, 7(6): 10.
- [3] Mitchell J C. Type systems for programming languages [M]//van Leeuwen J. Handbook of Theoretical Computer Science. Amsterdam: Elsevier, 1990: 365-458.
- [4] Wang Y. Software science: On general mathematical models and formal properties of software[J]. Journal of Advanced Mathematics and Applications, 2014, 3(2): 130-147.
- [5] Wang Y. On cognitive foundations and mathematical theories of knowledge science[J]. International Journal of Cognitive Informatics and Natural Intelligence, 2016, 10(2): 1-24.
- [6] Wang Y. Keynote: On the emergence of abstract sciences and breakthroughs in machine knowledge learning[C]//18th IEEE International Conference on Cognitive Informatics and Cognitive Computing (ICCI*CC 2019). Piscataway N J: IEEE Press, 2009: 5.
- [7] Wang Y. Keynote: The cognitive and mathematical foundations of big data science and blockchain engineering [C]//International Conference on Big Data and Blockchain (ICBDB'19). Piscataway N J: IEEE Press, 2009: 4.
- [8] Bender E A. Mathematical methods in artificial intelligence[M]. Los Alamitos: IEEE CS Press, 1996.
- [9] Berkeley B. Principles of human knowledge[M]. London: Berkeley, 1954.
- [10] Debenham J K. Knowledge systems design[M]. New York: Prentice Hall, 1989.
- [11] McCarthy J, Minsky M L, Rochester N, et al. Proposal for the 1956 dartmouth summer research project on artificial intelligence[R/OL]. [2019-10-31]. <http://www.formal.stanford.edu/jmc/history/dartmouth/dartmouth.html>.
- [12] McCulloch W S. Embodiments of mind[M]. Cambridge: MIT Press, 1965.
- [13] Shannon C E. A mathematical theory of communication [J]. The Bell System Technical Journal, 1948, 27: 379-423, 623-656.
- [14] Turing A M. Computing machinery and intelligence[J].

- Mind, 1950, 59: 433–460.
- [15] von Neumann J. The computer and the brain[M]. New Haven: Yale University Press, 1958.
- [16] Wang Y. Software engineering foundations: A software science perspective[M]. New York: Auerbach Publications, 2007.
- [17] Wang Y. On Abstract Intelligence: Toward a unified theory of natural, artificial, machinable, and computational intelligence[J]. International Journal of Software Science and Computational Intelligence, 2009, 1(1): 1–17.
- [18] Wang Y. Cognitive robots: A reference model towards intelligent authentication[J]. IEEE Robotics and Automation, 2010, 17(4): 54–62.
- [19] Wang Y. Keynote: Big data algebra: A rigorous approach to big data analytics and engineering[C]//17th International Conference on Mathematical and Computational Methods in Science and Engineering (MACMESE'15). Kuala Lumpur, 2015: 2.
- [20] Wang Y. Concept algebra: A denotational mathematics for formal knowledge representation and cognitive robot learning[J]. Journal of Advanced Mathematics and Applications, 2015, 4(1): 62–87.
- [21] Wang Y. In search of denotational mathematics: Novel mathematical means for contemporary intelligence, brain, and knowledge sciences[J]. Journal of Advanced Mathematics and Applications, 2012, 1(1): 4–25.
- [22] Lewis H R, Papadimitriou C H. Elements of the theory of computation[M]. New York: Prentice Hall, 1998.
- [23] Wang Y. On the informatics laws and deductive semantics of software[J]. IEEE Transactions on Systems, Man, and Cybernetics (Part C), 2006, 36(2): 161–171.
- [24] Zadeh L A. Fuzzy logic and approximate reasoning[J]. Synthese, 1975, 30(3/4): 407–428.
- [25] Wang Y. Fuzzy Causal Inferences based on fuzzy semantics of fuzzy concepts in cognitive computing[J]. WSEAS Transactions on Computers, 2014, 13: 430–441.
- [26] Wilson R, Keil F. The MIT encyclopedia of the cognitive sciences[J]. Electronic Resources Review, 2013, 43(4): 282–283.
- [27] Codd E F. A relational model of data for large shared data banks[J]. Communications of the ACM, 1970, 13(6): 377–387.
- [28] Cardelli L, Wegner P. On understanding types, data abstraction, and polymorphism[J]. ACM Computing Surveys, 1985, 17(4): 471–523.
- [29] Guttag J V. Abstract data types and the development of data structures[J]. Communications of the ACM, 1977, 20(6): 396–404.
- [30] Martin-Lof P. An intuitionistic theory of types: Predicative part[J]. Studies in Logic & the Foundations of Mathematics, 1975, 80: 73–118.
- [31] McKinsey B, Gartner D. Big Data means high value, not just volume[M]. Computer Weekly, 2011, 6: 1–2.
- [32] McKinsey B, Gartner D. Big Data means high value, not just volume[J]. Computer Weekly, 2011, 6: 1–2.
- [33] Mashey J R. Big data and the next wave of infrastrass[J]. SGI, 1998: 1–46.
- [34] Wang Y. On mathematical theories and cognitive foundations of information[J]. International Journal of Cognitive Informatics and Natural Intelligence, 2015, 9(3): 41–63.
- [35] Wang Y. Keynote: The emergence of abstract sciences and brain-inspired symbiotic systems[C]//IEEE FDC Workshop on Symbiotic Autonomous Systems in SMC'18.1. Piscataway N J: IEEE, 2018: 3.
- [36] Chapra S C, Canale R P. Numerical methods for engineers with software and programming applications[M]. Boston: McGraw-Hill, 2002.
- [37] Chicurel M. Databasing the brain[J]. Nature, 2000, 406(6798): 822–825.
- [38] Sniijders C, Matzat U, Reips U D. 'Big data': Big gaps of knowledge in the field of internet[J]. International Journal of Internet Science, 2012, 7: 1–5.
- [39] Sternberg R J. In search of the human mind[M]. 2nd ed. New York: Harcourt Brace & Co, 1998.
- [40] Ullman J D, Widom J. A first course in database systems [M]. New York: Prentice Hall, Inc., 1997.
- [41] Wang Y. On cognitive informatics[J]. Brain and Mind, 2003, 4(3): 151–167.
- [42] Wang Y. Keynote: From information revolution to intelligence revolution: big data science vs. intelligence science[C]//Proceedings of 13th IEEE International Conference on Cognitive Informatics and Cognitive Computing (ICCI*CC 2014). London: IEEE CS Press, 2014: 3–5.
- [43] Hartmanis J. On computational complexity and the nature of computer science, 1994 turing award lecture[J]. Communications of the ACM, 1994, 37(10): 37–43.
- [44] Wang Y. On the cognitive complexity of software and its quantification and formal measurement[J]. International Journal of Software Science and Computational Intelligence, 2019, 1(2): 31–53.

- [45] Wikipedia[EB/OL]. [2019- 10- 31]. https://en.wikipedia.org/wiki/Largest_known_prime_number.
- [46] Wang Y. On the big-R notation for describing iterative and recursive behaviors[J]. International Journal of Cognitive Informatics and Natural Intelligence, 2008, 2(1): 17-23.
- [47] Wang Y. Software science: On general mathematical models and formal properties of software[J]. Journal of Advanced Mathematics and Applications, 2014, 3(2): 130-147.

The cognitive and mathematical foundations of big data science

WANG Yingxu^{1,2}, PENG Jun³

1. National Engineering Key Lab for Big Data System Software, School of Software, Beijing National Research Center of Information Science and Technology, Tsinghua University, Beijing 100084, China
2. International Institute of Cognitive Informatics and Cognitive Computing (ICIC), Department of Electrical and Computer Engineering, Schulich School of Engineering and Hotchkiss Brain Institute, University of Calgary, Calgary T2N 1N4, Canada
3. School of Intelligent and Technology Engineering, Chongqing University of Science and Technology, Chongqing 401331, China

Abstract The big data play an indispensable role not only in a wide range of science fields and engineering applications, but also in the cognitive mechanisms of the sensation, the quantification, the qualification, the estimation, the measurement, the memory, and the reasoning of human beings. This paper reviews the basic studies of the theoretical foundations of the big data science, as well as a coherent set of general principles and analytic methodologies for the big data systems. The cognitive foundations of big data are explored in order to formally explain the origin and the nature of the big data. A set of mathematical models of the big data are created to rigorously elicit the general essences and patterns of the big data across pervasive domains in science, engineering, and society. A significant finding about the big data science is that the big data systems in nature are a recursively typed hyperstructure (RTHS) rather than pure numbers. The fundamental topological properties of the big data reveal a set of denotational mathematical solutions for dealing with the inherited complexities and unprecedented challenges in big data engineering.

Keywords big data science; big data engineering; mathematical models; recursively typed hyperstructures (RTHS); cognitive computing; computational intelligence ●



(责任编辑 刘志远)