

大数据系统和分析技术综述*

程学旗, 靳小龙, 王元卓, 郭嘉丰, 张铁赢, 李国杰

(中国科学院 计算技术研究所 网络数据科学与技术重点实验室, 北京 100190)

通讯作者: 靳小龙, E-mail: jinxiaolong@ict.ac.cn

摘要: 首先根据处理形式的不同,介绍了不同形式数据的特征和各自的典型应用场景以及相应的代表性处理系统,总结了大数据处理系统的三大发展趋势;随后,对系统支撑下的大数据分析技术和应用(包括深度学习、知识计算、社会计算与可视化等)进行了简要综述,总结了各种技术在大数据分析理解过程中的关键作用;最后梳理了大数据处理和分析面临的数据复杂性、计算复杂性和系统复杂性挑战,并逐一提出了可能的应对之策。

关键词: 大数据;数据分析;深度学习;知识计算;社会计算;可视化

中图法分类号: TP301

中文引用格式: 程学旗,靳小龙,王元卓,郭嘉丰,张铁赢,李国杰.大数据系统和分析技术综述.软件学报,2014,25(9):1889-1908.
http://www.jos.org.cn/1000-9825/4674.htm

英文引用格式: Cheng XQ, Jin XL, Wang YZ, GUO JF, Zhang TY, Li GJ. Survey on big data system and analytic technology.
Ruan Jian Xue Bao/Journal of Software, 2014, 25(9): 1240-1252 (in Chinese). http://www.jos.org.cn/1000-9825/4674.htm

Survey on Big Data System and Analytic Technology

CHENG Xue-Qi, JIN Xiao-Long, WANG Yuan-Zhuo, GUO Jia-Feng, ZHANG Tie-Ying, LI Guo-Jie

(Key Laboratory of Network Data Science and Technology, Institute of Computing Technology, The Chinese Academy of Sciences, Beijing 100190, China)

Corresponding author: JIN Xiao-Long, E-mail: jinxiaolong@ict.ac.cn

Abstract: This paper first introduces the key features of big data in different processing modes and their typical application scenarios, as well as corresponding representative processing systems. It then summarizes three development trends of big data processing systems. Next, the paper gives a brief survey on system supported analytic technologies and applications (including deep learning, knowledge computing, social computing, and visualization), and summarizes the key roles of individual technologies in big data analysis and understanding. Finally, the paper lays out three grand challenges of big data processing and analysis, i.e., data complexity, computation complexity, and system complexity. Potential ways for dealing with each complexity are also discussed.

Key words: big data; data analysis; deep learning; knowledge computing; social computing; visualization

近几年,大数据迅速发展成为科技界和企业界甚至世界各国政府关注的热点.《Nature》和《Science》等相继出版专刊专门探讨大数据带来的机遇和挑战.著名管理咨询公司麦肯锡称:“数据已经渗透到当今每一个行业和业务职能领域,成为重要的生产因素.人们对于大数据的挖掘和运用,预示着新一波生产力增长和消费盈余浪潮的到来”^[1].美国政府认为大数据是“未来的新石油”,一个国家拥有数据的规模和运用数据的能力将成为综合国力的重要组成部分,对数据的占有和控制将成为国家间和企业间新的争夺焦点.大数据已成为社会各界关注的新焦点,“大数据时代”已然来临.

什么是大数据,迄今并没有公认的定义.从宏观世界角度来讲,大数据是融合物理世界(physical world)、信息空间和人类社会(human society)三元世界的纽带,因为物理世界通过互联网、物联网等技术有了在信息空间

* 基金项目: 国家重点基础研究发展计划(973)(2014CB340401, 2012CB316303); 国家自然科学基金(61232010, 61100175, 61173008, 61202214); 北京市科技新星计划(Z121101002512063)

收稿时间: 2014-05-09; 定稿时间: 2014-07-01

(cyberspace)中的大数据反映,而人类社会则借助人机界面、脑机界面、移动互联等手段在信息空间中产生自己的大数据映像^[2,3]。从信息产业角度来讲,大数据还是新一代信息技术产业的强劲推动力。所谓新一代信息技术产业本质上是构建在第三代平台上的信息产业,主要是指大数据、云计算、移动互联网(社交网络)等。IDC 预测,到 2020 年第三代信息技术平台的市场规模将达到 5.3 万亿美元,而从 2013 年~2020 年,IT 产业 90%的增长将由第三代信息技术平台驱动。从社会经济角度来讲,大数据是第二经济(second economy^[4])的核心内涵和关键支撑。第二经济的概念是由美国经济学家 Auther 在 2011 年提出的。他指出由处理器、链接器、传感器、执行器以及运行在其上的经济活动形成了人们熟知的物理经济(第一经济)之外的第二经济(不是虚拟经济)。第二经济的本质是为第一经济附着一个“神经层”,使国民经济活动能够变得智能化,这是 100 年前电气化以来最大的变化。Auther 还估算了第二经济的规模,他认为到 2030 年,第二经济的规模将逼近第一经济。而第二经济的主要支撑是大数据,因为大数据是永不枯竭并不断丰富的资源产业。借助于大数据,未来第二经济下的竞争将不再是劳动生产率而是知识生产率的竞争。

相较于传统的数据,人们将大数据的特征总结为 5 个 V,即体量大(volume)、速度快(velocity)、模态多(variety)、难辨识(veracity)和价值大密度低(value)。但大数据的主要难点并不在于数据量大,因为通过对计算机系统的扩展可以在一定程度上缓解数据量大带来的挑战。其实,大数据真正难以对付的挑战来自于数据类型多样(variety)、要求及时响应(velocity)和数据的不确定性(veracity)。因为数据类型多样使得一个应用往往既要处理结构化数据,同时还要处理文本、视频、语音等非结构化数据,这对现有数据库系统来说难以应付;在快速响应方面,在许多应用中时间就是利益。在不确定性方面,数据真伪难辨是大数据应用的最大挑战。追求高数据质量是对大数据的一项重要要求,最好的数据清理方法也难以消除某些数据固有的不可预测性。

为了应对大数据带来的上述困难和挑战,以 Google,Facebook,Linkedin,Microsoft 等为代表的互联网企业近几年推出了各种不同类型的大数据处理系统。借助于新型的处理系统,深度学习、知识计算、可视化等大数据分析技术也得以迅速发展,已逐渐被广泛应用于不同的行业和领域。本文从系统支撑下的大数据分析角度入手,介绍了不同的大数据处理模式与代表性的处理系统,并对深度学习、知识计算等重要的大数据分析技术进行综述,最后指出大数据处理和分析所面临的 3 个核心挑战,并提出可能的解决思路。

1 大数据处理与系统

大数据中蕴含的宝贵价值成为人们存储和处理大数据的驱动力。Mayer-Schönberger 在《大数据时代》一书中指出了大数据时代处理数据理念的三大转变,即要全体不要抽样,要效率不要绝对精确,要相关不要因果^[5]。因此,海量数据的处理对于当前存在的技术来说是一种极大的挑战。目前,人们对大数据的处理形式主要是对静态数据的批量处理,对在线数据的实时处理^[6],以及对图数据的综合处理。其中,在线数据的实时处理又包括对流式数据的处理和实时交互计算两种。本节将详细阐述上述 4 种数据形式的特征和各自的典型应用以及相应的代表性系统。

1.1 批量数据处理系统

利用批量数据挖掘合适的模式,得出具体含义,制定明智的决策,最终做出有效的应对措施实现业务目标是大数据批处理的首要任务。大数据的批量处理系统适用于先存储后计算,实时性要求不高,同时数据的准确性和全面性更为重要的场景。

1.1.1 批量数据的特征与典型应用

(1) 批量数据的特征

批量数据通常具有 3 个特征。第一,数据体量巨大。数据从 TB 级别跃升到 PB 级别。数据是以静态的形式存储在硬盘中,很少进行更新,存储时间长,可以重复利用,然而这样大批量的数据不容易对其进行移动和备份。第二,数据精确度高。批量数据往往是从应用中沉淀下来的数据,因此精度相对较高,是企业资产的一部分宝贵财富。第三,数据价值密度低。以视频批量数据为例,在连续不断的监控过程中,可能有用的数据仅仅有一两秒。因此,需要通过合理的算法才能从批量的数据中抽取有用的价值。此外,批量数据处理往往比较耗时,而且不提供用户

与系统的交互手段,所以当发现处理结果和预期或与以往的结果有很大差别时,会浪费很多时间.因此,批量数据处理适合大型的相对比较成熟的作业.

(2) 典型应用

物联网、云计算、互联网以及车联网等无一不是大数据的重要来源,当前批量数据处理可以解决前述领域的诸多决策问题并发现新的洞察.因此,批量数据处理可以适用于较多的应用场景.本节主要选择互联网领域的应用、安全领域的应用以及公共服务领域的应用这3个典型应用场景加以介绍^[7-12].在互联网领域中,批量数据处理的典型应用场景主要包括:(a) 社交网络:Facebook、新浪微博、微信等以人为核心的社交网络产生了大量的文本、图片、音视频等不同形式的数据.对这些数据的批量处理可以对社交网络进行分析,发现人与人之间隐含的关系或者他们中存在的社区,推荐朋友或者相关的主题,提升用户的体验.(b) 电子商务:电子商务中产生大量的购买历史记录、商品评论、商品网页的访问次数和驻留时间等数据,通过批量分析这些数据,每个商铺可以精准地选择其热卖商品,从而提升商品销量;这些数据还能够分析出用户的消费行为,为客户推荐相关商品,以提升优质客户数量.(c) 搜索引擎:Google等大型互联网搜索引擎与Yahoo!的专门广告分析系统,通过对广告相关数据的批量处理用来改善广告的投放效果以提高用户的点击量.在安全领域中,批量数据主要用于欺诈检测和IT安全.在金融服务机构和情报机构中,欺诈检测一直都是关注的重点.通过对批量数据的处理,可对客户交易和现货异常进行判断,从而对可能存在欺诈行为提前预警.另一方面,企业通过处理机器产生的数据,识别恶意软件和网络攻击模式,从而使其他安全产品判断是否接受来自这些来源的通信.在公共服务领域,批量数据处理的典型应用场景主要包括:(a) 能源:例如,对来自海洋深处地震时产生的数据进行批量的排序和整理,可能发现海底石油的储量;通过对用户能源数据、气象与人口方面的公共及私人数据、历史信息、地理数据等的批量处理,可以提升电力服务,尽量为用户节省在资源方面的投入.(b) 医疗保健:通过对患者以往的生活方式与医疗记录进行批量处理分析,提供语义分析服务,对病人的健康提供医生、护士及其他相关人士的回答,并协助医生更好的为患者进行诊断.当然,大数据的批量处理不只应用到这些领域,还有移动数据分析、图像处理以及基础设施管理等领域.随着人们对数据中蕴含价值的认识,会有更多的领域通过对数据的批量处理挖掘其中的价值来支持决策和发现新的洞察.

1.1.2 代表性的处理系统

由Google公司2003年研发的Google文件系统GFS^[13]和2004年研发的MapReduce编程模型^[14]以其Web环境下批量处理大规模海量数据的特有魅力,在学术界和工业界引起了很大反响.虽然Google没有开源这两项技术的源码,但是基于这两篇开源文档,2006年Nutch项目子项目之一的Hadoop实现了两个强有力的开源产品^[15]:HDFS和MapReduce.Hadoop成为了典型的大数据批量处理架构,由HDFS负责静态数据的存储,并通过MapReduce将计算逻辑分配到各数据节点进行数据计算和价值发现.Hadoop顺应了现代主流IT公司的一致需求,之后以HDFS和MapReduce为基础建立了很多项目,形成了Hadoop生态圈.

MapReduce编程模型之所以受到欢迎并迅速得到应用,在技术上主要有3方面的原因^[16,17].首先,MapReduce采用无共享大规模集群系统.集群系统具有良好的性价比和可伸缩性,这一优势为MapReduce成为大规模海量数据平台的首选创造了条件.其次,MapReduce模型简单、易于理解、易于使用.它不仅用于处理大规模数据,而且能将很多繁琐的细节隐藏起来(比如,自动并行化、负载均衡和灾备管理等),极大地简化了程序员的开发工作.而且,大量数据处理问题,包括很多机器学习和数据挖掘算法,都可以使用MapReduce实现.第三,虽然基本的MapReduce模型只提供一个过程性的编程接口,但在海量数据环境、需要保证可伸缩性的前提下,通过使用合适的查询优化和索引技术,MapReduce仍能够提供很好的数据处理性能.

1.2 流式数据处理系统

Google于2010年推出了Dremel,引领业界向实时数据处理迈进.实时数据处理是针对批量数据处理的性能问题提出的,可分为流式数据处理和交互式数据处理两种模式.在大数据背景下,流式数据处理源于服务器日志的实时采集,交互式数据处理的目标是将PB级数据的处理时间缩短到秒级.

1.2.1 流式数据的特征及典型应用

(1) 流式数据的特征

通俗而言,流式数据是一个无穷的数据序列,序列中的每一个元素来源各异,格式复杂,序列往往包含时序特性,或者有其他有序标签(如 IP 报文中的序号)。从数据库的角度而言,每一个元素可以看作是一个元组,而元素的特性则类比于元组的属性。流式数据在不同的场景下往往体现出不同的特征,如流速大小、元素特性数量、数据格式等,但大部分流式数据都含有共同的特征,这些特征便可用来设计通用的流式数据处理系统。下面简要介绍流式数据共有的特征^[18]。

首先,流式数据的元组通常带有时间标签或其余含序属性。因此,同一流式数据往往是被按序处理的。然而数据的到达顺序是不可预知的,由于时间和环境的动态变化,无法保证重放数据流与之前数据流中数据元素顺序的一致性。这就导致了数据的物理顺序与逻辑顺序不一致。而且,数据源不受接收系统的控制,数据的产生是实时的、不可预知的。此外,数据的流速往往有较大的波动,因此需要系统具有很好的可伸缩性,能够动态适应不确定流入的数据流,具有很强的系统计算能力和大数据流量动态匹配的能力。其次,数据流中的数据格式可以是结构化的、半结构化的甚至是无结构化的。数据流中往往含有错误元素、垃圾信息等。因此流式数据的处理系统要有很好的容错性与异构数据分析能力,能够完成数据的动态清洗、格式处理等。最后,流式数据是活动的(用完即弃),随着时间的推移不断增长,这与传统的数据处理模型(存储→查询)不同,要求系统能够根据局部数据进行计算,保存数据流的动态属性。流式处理系统针对该特性,应当提供流式查询接口,即提交动态的 SQL 语句,实时地返回当前结果。

(2) 典型应用

流式计算的应用场景较多,典型的有两类^[6]:(a) 数据采集应用:数据采集应用通过主动获取海量的实时数据,及时地挖掘出有价值的信息。当前数据采集应用有日志采集、传感器采集、Web 数据采集等。日志采集系统是针对各类平台不断产生的大量日志信息量身订做的处理系统,通过流式挖掘日志信息,达到动态提醒与预警功能。传感器采集系统(物联网)通过采集传感器的信息(通常包含时间、位置、环境和行为等内容),实时分析提供动态的信息展示,目前主要应用于智能交通、环境监控、灾难预警等。Web 数据采集系统是利用网络爬虫程序抓取万维网上的内容,通过清洗、归类、分析并挖掘其数据价值。(b) 金融银行业的应用:在金融银行领域的日常运营过程中会产生大量数据,这些数据的时效性往往较短,不仅有结构化数据,也会有半结构化和非结构化数据。通过对这些大数据的流式计算,发现隐含于其中的内在特征,可帮助金融银行进行实时决策。这与传统的商业智能(BI)分析不同,BI 要求数据是静态的,通过数据挖掘技术,获得数据的价值。然而在瞬息万变的场景下,诸如股票期货市场,数据挖掘技术不能及时地响应需求,就需要借助流式数据处理的帮助。

总之,流式数据的特点是:数据连续不断、来源众多、格式复杂、物理顺序不一、数据的价值密度低。而对应的处理工具则需具备高性能、实时、可扩展等特性。

1.2.2 代表性的处理系统

流式数据处理已经在业界得到广泛的应用,典型的有 Twitter 的 Storm,Facebook 的 Scribe,Linkedin 的 Samza,Cloudera 的 Flume,Apache 的 Nutch。

• Twitter 的 Storm 系统

Storm^[19]是一套分布式、可靠、可容错的用于处理流式数据的系统。其流式处理作业被分发至不同类型的组件,每个组件负责一项简单的、特定的处理任务。Storm 集群的输入流由名为 Spout 的组件负责。Spout 将数据传递给名为 Bolt 的组件,后者将以指定的方式处理这些数据,如持久化或者处理并转发给另外的 Bolt。Storm 集群可以看成一条由 Bolt 组件组成的链(称为一个 Topology)。每个 Bolt 对 Spout 产生出来的数据做某种方式的处理。

Storm 可用来实时处理新数据和更新数据库,兼具容错性和扩展性。Storm 也可被用于连续计算,对数据流做连续查询,在计算时将结果以流的形式输出给用户。它还可被用于分布式 RPC,以并行的方式运行复杂运算。一个 Storm 集群分为 3 类节点:(a) Nimbus 节点,负责提交任务,分发执行代码,为每个工作节点指派任务和监控失

败的任务;(b) Zookeeper 节点,负责 Storm 集群的协同操作;(c) Supervisor 节点,负责启动多个 Worker 进程,执行 Topology 的一部分,这个过程是通过 Zookeeper 节点与 Nimbus 节点通信完成的.因为 Storm 将所有的集群状态保存在 Zookeeper 或者本地磁盘上,Supervisor 节点是无状态的,因此其失败或者重启不会引起全局的重新计算.

Storm 的主要特点是:(a) 简单的编程模型:Storm 提供类似于 MapReduce 的操作,降低了并行批处理与实时处理的复杂性.一个 Storm 作业只需实现一个 Topology 及其所包含的 Spout 与 Bolt.通过指定它们的连接方式,Topology 可以胜任大多数的流式作业需求;(b) 容错性:Storm 利用 Zookeeper 管理工作进程和节点的故障.在工作过程中,如果出现异常,Topology 会失败.但 Storm 将以一致的状态重新启动处理,这样它可以正确地恢复;(c) 水平扩展:Storm 拥有良好的水平扩展能力,其流式计算过程是在多个线程、进程和服务端之间并行进行的.Nimbus 节点将大量的协同工作都交由 Zookeeper 节点负责,使得水平扩展不会产生瓶颈;(d) 快速可靠的消息处理:Storm 利用 ZeroMQ 作为消息队列,极大提高了消息传递的速度,系统的设计也保证了消息能得到快速处理.Storm 保证每个消息至少能得到一次完整处理.任务失败时,它会负责从消息源重试消息.

- LinkedIn 的 Samza 系统

LinkedIn 早期开发了一款名叫 Kafka^[20,21]的消息队列,广受业界的好评,许多流式数据处理系统都使用了 Kafka 作为底层的消息处理模块.Kafka 的工作过程简要分为 4 个步骤,即生产者将消息发往中介(broker),消息被抽象为 Key-Value 对,Broker 将消息按 Topic 划分,消费者向 Broker 拉取感兴趣的 Topic.2013 年,LinkedIn 基于 Kafka 和 YARN 开发了自己的流式处理框架——Samza.Samza 与 Kafka 的关系可以类比 MapReduce 与 HDFS 的关系.Samza 系统由 3 个层次组成,包括流式数据层(Kafka)、执行层(YARN)、处理层(Samza API).一个 Samza 任务的输入与输出均是流.Samza 系统对流的模型有很严格的定义,它并不只是一个消息交换的机制.流在 Samza 的系统中是一系列划分的、可重现的、可多播的、无状态的消息序列,每一个划分都是有序的.流不仅是 Samza 系统的输入与输出,它还充当系统中的缓冲区,能够隔离相互之间的处理过程.Samza 利用 YARN 与 Kafka 提供了分步处理与划分流的处理框架.Samza 客户端向 Yarn 的资源管理器提交作业,生成多个 Task Runner 进程,这些进程执行用户编写的 StreamTasks 代码.该系统的输入与输出来自于 Kafka 的 Broker 进程.

Samza 的主要特性有:(a) 高容错:如果服务器或者处理器出现故障,Samza 将与 YARN 一起重新启动流处理器.(b) 高可靠性:Samza 使用 Kafka 来保证所有消息都会按照写入分区的顺序进行处理,绝对不会丢失任何消息.(c) 可扩展性:Samza 在各个等级进行分割和分布,Kafka 提供一个有序、可分割、可重部署、高容错的系统;YARN 提供了一个分布式环境供 Samza 容器运行.

1.3 交互式数据处理

1.3.1 交互式数据处理的特征与典型应用

(1) 交互式数据处理的特征

与非交互式数据处理相比,交互式数据处理灵活、直观、便于控制.系统与操作人员以人机对话的方式一问一答——操作人员提出请求,数据以对话的方式输入,系统便提供相应的数据或提示信息,引导操作人员逐步完成所需的操作,直至获得最后处理结果.采用这种方式,存储在系统中的数据文件能够被及时处理修改,同时处理结果可以立刻被使用.交互式数据处理具备的这些特征能够保证输入的信息得到及时处理,使交互方式继续进行下去.

(2) 典型应用

在大数据环境下,数据量的急剧膨胀是交互式数据处理系统面临的首要问题.下面主要选择信息处理系统领域和互联网领域做为典型应用场景进行介绍.(a) 在信息处理系统领域中,主要体现了人机间的交互.传统的交互式数据处理系统主要以关系型数据库管理系统(DBMS)为主,面向两类应用,即联机事务处理(OLTP)和联机分析处理(OLAP).OLTP 基于关系型数据库管理系统,广泛用于政府、医疗以及对操作序列有严格要求的工业控制领域;OLAP 基于数据仓库系统(data warehouse)广泛用于数据分析、商业智能(BI)等.最具代表性的处理是数据钻取,如在 BI 中,可以对于数据进行切片和多粒度的聚合,从而通过多维分析技术实现数据的钻取.目前,基

于开源体系架构下的数据仓库系统发展十分迅速,以 Hive^[22]、Pig^[23]等为代表的分布式数据仓库能够支持上千台服务器的规模。(b) 互联网领域.在互联网领域中,主要体现了人际间的交互.随着互联网技术的发展,传统的简单按需响应的人机互动已不能满足用户的需求,用户之间也需要交互,这种需求诞生了互联网中交互式数据处理的各种平台,如搜索引擎、电子邮件、即时通讯工具、社交网络、微博、博客以及电子商务等,用户可以在这些平台上获取或分享各种信息.此外,各种交互式问答平台,如百度的知道、新浪的爱问以及 Yahoo!的知识堂等.由此可见,用户与平台之间的交互变得越来越容易,越来越频繁.这些平台中数据类型的多样性,使得传统的关系数据库不能满足交互式数据处理的实时性需求.目前,各大平台主要使用 NoSQL 类型的数据库系统来处理交互式的数据库,如 HBase^[24]采用多维有续表的列式存储方式;MongoDB^[25]采用 JSON 格式的数据嵌套存储方式.大多 NoSQL 数据库不提供 Join 等关系数据库的操作模式,以增加数据操作的实时性.

1.3.2 代表性的处理系统

交互式数据处理系统的典型代表系统是 Berkeley 的 Spark 系统和 Google 的 Dremel 系统.

• Berkeley 的 Spark 系统

Spark^[26]是一个基于内存计算的可扩展的开源集群计算系统.针对 MapReduce 的不足,即大量的网络传输和磁盘 I/O 使得效率低效,Spark 使用内存进行数据计算以便快速处理查询,实时返回分析结果.Spark 提供比 Hadoop 更高层的 API,同样的算法在 Spark 中的运行速度比 Hadoop 快 10 倍~100 倍^[26].Spark 在技术层面兼容 Hadoop 存储层 API,可访问 HDFS, HBASE, SequenceFile 等.Spark-Shell 可以开启交互式 Spark 命令环境,能够提供交互式查询.

Spark 是为集群计算中的特定类型的工作负载而设计,即在并行操作之间重用工作数据集(比如机器学习算法)的工作负载.Spark 的计算架构具有 3 个特点:(a) Spark 拥有轻量级的集群计算框架.Spark 将 Scala 应用于他的程序架构,而 Scala 这种多范式的编程语言具有并发性、可扩展性以及支持编程范式的特征,与 Spark 紧密结合,能够轻松地操作分布式数据集,并且可以轻易地添加新的语言结构.(b) Spark 包含了大数据领域的数据流计算和交互式计算.Spark 可以与 HDFS 交互取得里面的数据文件,同时 Spark 的迭代、内存计算以及交互式计算为数据挖掘和机器学习提供了很好的框架.(c) Spark 有很好的容错机制.Spark 使用了弹性分布数据集(RDD),RDD 被表示为 Scala 对象分布在一组节点中的只读对象集中,这些集合是弹性的,保证了如果有一部数据集丢失时,可以对丢失的数据集进行重建.

Spark 高效处理分布数据集的特征使其有着很好的应用前景,现在四大 Hadoop 发行商 Cloudera, Pivotal, MapR 以及 Hortonworks 都提供了对 Spark 的支持.

• Google 的 Dremel 系统

Dremel^[27]是 Google 研发的交互式数据分析系统,专注于只读嵌套数据的分析.Dremel 可以组建成规模上千的服务器集群,处理 PB 级数据.传统的 MapReduce 完成一项处理任务,最短需要分钟级的时间,而 Dremel 可以将处理时间缩短到秒级.Dremel 是 MapReduce 的有力补充,可以通过 MapReduce 将数据导入到 Dremel 中,使用 Dremel 来开发数据分析模型,最后在 MapReduce 中运行数据分析模型.

Dremel 作为大数据的交互式处理系统可以与传统的数据分析或商业智能工具在速度和精度上相媲美.Dremel 系统主要有以下 5 个特点:(a) Dremel 是一个大规模系统.在 PB 级数据集上要将任务缩短到秒级,需要进行大规模的并发处理,而磁盘的顺序读速度在 100MB/S 上下,因此在 1s 内处理 1TB 数据就意味着至少需要有 1 万个磁盘的并发读,但是机器越多,出问题概率越大,如此大的集群规模,需要有足够的容错考虑,才能够保证整个分析的速度不被集群中的个别慢(坏)节点影响.(b) Dremel 是对 MapReduce 交互式查询能力不足的有力补充.Dremel 利用 GFS 文件系统作为存储层,常常用它来处理 MapReduce 的结果集或建立分析原型.(c) Dremel 的数据模型是嵌套的.Dremel 类似于 Json,支持一个嵌套的数据模型.对于处理大规模数据,不可避免的有大量的 Join 操作,而传统的关系模型显得力不从心,Dremel 却可以很好地处理相关的查询操作.(d) Dremel 中的数据是用列式存储的.使用列式存储,在进行数据分析的时候,可以只扫描所需要的那部分数据,从而减少 CPU 和磁盘的访问量.同时,列式存储是压缩友好的,通过压缩可以综合 CPU 和磁盘从而发挥最大的效能.(e)

Dremel 结合了 Web 搜索和并行 DBMS 的技术.首先,它借鉴了 Web 搜索中查询树的概念,将一个相对巨大复杂的查询,分割成较小、较简单的查询,分配到并发的大量节点上.其次,与并行 DBMS 类似,Dremel 可以提供了一个 SQL-like 的接口.

1.4 图数据处理系统

图由于自身的结构特征,可以很好地表示事物之间的关系,在近几年已成为各学科研究的热点.图中点和边的强关联性,需要图数据处理系统对图数据进行一系列的操作,包括图数据的存储、图查询、最短路径查询、关键字查询、图模式挖掘以及图数据的分类、聚类等.随着图中节点和边数的增多(达到几千万甚至上亿数),图数据处理的复杂性给图数据处理系统提出了严峻的挑战.下面主要阐述图数据的特征和典型应用以及代表性的图数据处理系统.

1.4.1 图数据的特征及典型应用

(1) 图数据的特征

图数据中主要包括图中的节点以及连接节点的边,通常具有 3 个特征.第一,节点之间的关联性.图中边的数量是节点数量的指数倍,因此,节点和关系信息同等重要,图结构的差异也是由于对边做了限制,在图中,顶点和边实例化构成各种类型的图,如标签图、属性图、语义图以及特征图等.第二,图数据的种类繁多.在许多领域中,使用图来表示该领域的的数据,如生物、化学、计算机视觉、模式识别、信息检索、社会网络、知识发现、动态网络交通、语义网、情报分析等.每个领域对图数据的处理需求不同,因此,没有一个通用的图数据处理系统满足所有领域的需求.第三,图数据计算的强耦合性.在图中,数据之间是相互关联的,因此,对图数据的计算也是相互关联的.这种数据耦合的特性对图的规模日益增大达到上百万甚至上亿节点的大图数据计算提出了巨大的挑战.大图数据是无法使用单台机器进行处理的,但如果对大图数据进行并行处理,对于每一个顶点之间都是连通的图来讲,难以分割成若干完全独立的子图进行独立的并行处理;即使可以分割,也会面临并行机器的协同处理,以及将最后的处理结果进行合并等一系列问题.这需要图数据处理系统选取合适的图分割以及图计算模型来迎接挑战并解决问题.

(2) 典型应用

图能很好地表示各实体之间的关系,因此,在各个领域得到了广泛的应用,如计算机领域、自然科学领域以及交通领域.(a) 互联网领域的应用.随着信息技术和网络技术的发展,以 Web 2.0 技术为基础的社交网络(如 Facebook、人人网)、微博(如 Twitter、新浪微博、腾讯微博)等新兴服务中建立了大量的在线社会网络关系,用图表示人与人之间的关系.在社交网络中,基于图研究社区发现等问题;在微博中,通过图研究信息传播与影响力最大化等问题.除此之外,用图表示如 E-mail 中的人与人之间的通信关系,从而可以研究社会群体关系等问题;在搜索引擎中,可以用图表示网页之间相互的超链接关系,从而计算一个网页的 PageRank 得分等.(b) 自然科学领域的应用.图可以用来在化学分子式中查找分子,在蛋白质网络中查找化合物,在 DNA 中查找特定序列等.(c) 交通领域的应用.图可用来在动态网络交通中查找最短路径,在邮政快递领域进行邮路规划等.当然,图还有一些其他的应用,如疾病爆发路径的预测与科技文献的引用关系等.图数据虽然结构复杂,处理困难,但是它有很好的表现力,因此得到了各领域的广泛应用.随着图数据处理中所面临的各种挑战被不断地解决,图数据处理将有更好的应用前景.

1.4.2 代表性图数据处理系统

现今主要的图数据库有 GraphLab, Giraph(基于 Pregel 克隆), Neo4j, HyperGraphDB, InfiniteGraph, Cassovary, Trinity 以及 Grappa 等.下面介绍 3 个典型的图数据处理系统,包括 Google 的 Pregel 系统, Neo4j 系统和微软的 Trinity 系统.

• Google 的 Pregel 系统

Pregel^[28,29]是 Google 提出的基于 BSP(Bulk synchronous parallel)模型的分布式图计算框架,主要用于图遍历(BFS)、最短路径(SSSP)、PageRank 计算等. BSP 模型是并行计算模型中的经典模型,采用的是“计算-通信-同步”的模式.它将计算分成一系列超步(superstep)的迭代.从纵向上看,它是一个串行模式,而从横向上看,它是

一个并行的模式,每两个超步之间设置一个栅栏,即整体同步点,确定所有并行的计算都完成后再启动下一轮超步.Pregel 的设计思路是以节点为中心计算,节点有两种状态:活跃和不活跃.初始时每个节点都处于活跃状态,完成计算后每个节点主动“Vote to Halt”进入不活跃状态.如果接收到信息,则激活.没有活跃节点和信息时,整个算法结束.

Pregel 架构有 3 个主要特征:(a) 采用主/从(Master/Slave)结构来实现整体功能.一个节点为 Master,负责对整个图结构的任务进行切分,根据节点的 ID 进行散列计算分配到 Slave 机器,Slave 机器进行独立的超步计算,并将结果返回给 Master;(b) 有很好的容错机制.Pregel 通过 Checkpoint 机制实行容错,节点向 Master 汇报心跳维持状态,节点间采用异步消息传输;(c) 使用 GFS 或 BigTable 作为持久性的存储.

Apache 根据 Google 于 2010 年发表的 Pregel 论文开发了高可扩展的迭代的图处理系统 Giraph,现在已经被 Facebook 用于分析社会网络中用户间的关系图中.

• Neo4j 系统

Neo4j^[30]是一个高性能的、完全兼容 ACID 特性的、鲁棒的图数据库.它基于 Java 语言开发,包括社区版和企业版,适用于社会网络和动态网络等场景.Neo4j 在处理复杂的网络数据时表现出很好的性能.数据以一种针对图形网络进行过优化的格式保存在磁盘上.Neo4j 重点解决了拥有大量连接的查询问题,提供了非常快的图算法、推荐系统以及 OLAP 风格的分析,满足了企业的应用、健壮性以及性能的需求,得到了很好的应用.

Neo4j 系统具有以下 5 个特性.(a) 支持数据库的所有特性:Neo4j 的内核是一种极快的图形引擎,支持事物的 ACID 特性、两阶段提交、符合分布式事务以及恢复等;(b) 高可用性:Neo4j 通过联机备份实现它的高可用性;(c) 可扩展性:Neo4j 提供了大规模可扩展性,可以在一台机器上处理数十亿节点/关系/属性的图,也可以扩展到多台机器上并行运行;(d) 灵活性:Neo4j 拥有灵活的数据结构,可以通过 Java-API 直接与图模型进行交互.对于 JRuby/Ruby,Scala,Python 以及 Clojure 等其他语言,也开发了相应的绑定库;(e) 高速遍历:Neo4j 中图遍历执行的速度是常数,与图的规模大小无关.它的读性能可以实现每毫秒遍历 2 000 关系,而且完全是事务性的.Neo4j 以一种延迟风格遍历图,即节点和关系只有在结果迭代器需要访问它们的时候才会被遍历并返回,支持深度搜索和广度搜索两种遍历方式.

• 微软的 Trinity 系统

Trinity^[31,32]是 Microsoft 推出的一款建立在分布式云存储上的计算平台,可以提供高度并行查询处理、事务记录、一致性控制等功能.Trinity 主要使用内存存储,磁盘仅作为备份存储.

Trinity 有以下 4 个特点.(a) 数据模型是超图:超图中,一条边可以连接任意数目的图顶点.此模型中图的边称为超边.基于这种特点,超图比简单图的适用性更强,保留的信息更多;(b) 并发性:Trinity 可以配置在一台或上百台计算机上.Trinity 提供了一个图分割机制,由一个 64 位的唯一标识 UID 确定各结点的位置,利用散列方式映射到相应的机器上,以尽量减少延迟,如图 8 所示.Trinity 可以并发执行 PageRank、最短路径查询、频繁子图挖掘以及随机游走等操作;(c) 具有数据库的一些特点:Trinity 是一个基于内存的图数据库,有丰富的数据库特点,如:在线高度并行查询处理、ACI 交易支持、并发控制以及一致性维护等;(d) 支持批处理:Trinity 支持大型在线查询和离线批处理,并且支持同步和不同步批处理计算.相比之下,Pregel 只支持在线查询处理,批处理必须是严格的同步计算.

微软现在使用 Trinity 作为 Probase 的基础架构,可以从网上自动获得大规模的知识库.Trinity 主要作用是分类建设、数据集成以及查询 Probase.Trinity 也被用于其他的项目中,如 Aether 项目,其功能也在不断的增加中.

1.5 小结

面对大数据,各种处理系统层出不穷,各有特色.总体来说,我们可以总结出 3 种发展趋势:(1) 数据处理引擎专用化:为了降低成本,提高能效,大数据系统需要摆脱传统的通用体系,趋向专用化架构技术.为此,国内外的互联网龙头企业都在基于开源系统开发面向典型应用的大规模、高通量、低成本、强扩展的专用化系统;(2) 数据处理平台多样化:自 2008 年以来克隆了 Google 的 GFS 和 MapReduce 的 Apache Hadoop 逐渐被互联网企业所广泛接纳,并成为大数据处理领域的事实标准.但在全面兼容 Hadoop 的基础上,Spark 通过更多的利用

内存处理大幅提高系统性能。而 Scribe,Flume,Kafka,Storm,Drill,Impala,TEZ/Stinger,Presto,Spark/Shark 等的出现并不是取代 Hadoop,而是扩大了大数据技术的生态环境,促使生态环境向良性化和完整化发展。(3) 数据计算实时化:在大数据背景下,作为批量计算的补充,旨在将 PB 级数据的处理时间缩短到秒级的实时计算受到越来越多的关注。

2 大数据分析

要挖掘大数据的大价值必然要对大数据进行内容上的分析与计算。深度学习和知识计算是大数据分析的基础,而可视化既是数据分析的关键技术也是数据分析结果呈现的关键技术。本节主要介绍深度学习、知识计算和可视化等大数据分析的关键技术,同时也对大数据的典型应用包括社交媒体计算等进行简要综述。

2.1 深度学习

大数据分析的一个核心问题是如何对数据进行有效表达、解释和学习,无论是对图像、声音还是文本数据。传统的研究也有很多数据表达的模型和方法,但通常都是较为简单或浅层的模型,模型的能力有限,而且依赖于数据的表达,不能获得很好的学习效果。大数据的出现提供了使用更加复杂的模型来更有效地表征数据、解释数据的机会。深度学习就是利用层次化的架构学习出对象在不同层次上的表达,这种层次化的表达可以帮助解决更加复杂抽象的问题。在层次化中,高层的概念通常是通过低层的概念来定义的。深度学习通常使用人工神经网络,常见的具有多个隐层的多层感知机(MLP)就是典型的深度架构。

深度学习的起源要追溯到神经网络,20 世纪 80 年代,后向传播(BP)算法的提出使得人们开始尝试训练深层次的神经网络。然而,BP 算法在训练深层网络的时候表现不够好,以至于深层感知机的效果还不如浅层感知机。于是很多人放弃使用神经网络,转而使用凸的更容易得到全局最优解的浅层模型,提出诸如支持向量机、boosting 等浅层方法,以致于此前大部分的机器学习技术都使用浅层架构。转机出现在 2006 年,多伦多大学的 Hinton 等人使用无监督的逐层贪婪的预训练(greedy layer-wise pre-train)方法成功减轻了深度模型优化困难的问题^[33],从而掀起了深度学习的浪潮。Hinton 引入了深度产生式模型 DBN,并提出高效的逐层贪婪的学习算法,使用 DBN 初始化一个深度神经网络(DNN)再对 DNN 进行精调,通常能够产生更好的结果。Bengio 等人^[34]基于自动编码器(auto-encoder)提出了非概率的无监督深度学习模型,也取得了类似的效果。

近几年,深度学习在语音、图像以及自然语言理解等应用领域取得一系列重大进展。从 2009 年开始,微软研究院的 Dahl 等人率先在语音处理中使用深度神经网络(DNN),将语音识别的错误率显著降低,从而使得语音处理成为成功应用深度学习的第 1 个领域^[35]。在图像领域,2012 年,Hinton 等人使用深层次的卷积神经网络(CNN)在 ImageNet 评测上取得巨大突破,将错误率从 26%降低到 15%^[36],重要的是,这个模型中并没有任何手工构造特征的过程,网络的输入就是图像的原始像素值。在此之后,采用类似的模型,通过使用更多的参数和训练数据,ImageNet 评测的结果得到进一步改善,错误率下降至 2013 年的 11.2%^[37]。Facebook 人工智能实验室的 Taigman 等人使用了与文献[36]中类似的神经网络在人脸识别上也取得了很好的效果,将人脸识别的正确率提升至接近人类水平^[38]。此外,图像领域还有一些基于无监督的深度学习研究,比如在 Google Brain 项目中,Le 等人尝试使用完全无标注的图像训练得到人脸特征检测器,使用这些学习到的特征可以在图像分类中取得非常好的效果^[39];Google 的深度学习系统(DistBelief)在获取数百万 YouTube 视频数据后,能够精准地识别出这些视频中的关键元素——猫。在自然语言领域,从 2003 年开始,Bengio 等人使用神经网络并结合分布式表达(distributed representation)的思想训练语言模型并取得很好的效果^[40],不过当时还没有使用到更深层次的模型。2008 年,Collobert 等人训练了包含一个卷积层的深度神经网络,利用学习得到的中间表达同时解决多个 NLP 问题^[41]。尽管这些工作没有取得像图像和语音处理领域如此重大的进展,但也都接近或超过了已有的最好方法。近年来,斯坦福大学的 Socher 等人的一系列工作也值得关注。他们使用递归神经网络(recursive neural network,简称 RNN)在情感分析等问题上取得一系列进展,将现有的准确率从 80%提升到 85%^[42]。在国内,2011 年科大讯飞首次将 DNN 技术运用到语音云平台,并提供给开发者使用,并在讯飞语音输入法和讯飞口讯等产品中得到应用。百度成立了 IDL(深度学习研究院),专门研究深度学习算法,目前已有多项深度学习技术在百度产

品上线.深度学习对百度影响深远,在语音识别、OCR 识别、人脸识别、图像搜索等应用上取得了突出效果.此外,国内其他公司如搜狗、云知声等纷纷开始在产品中使用深度学习技术.

2.2 知识计算

基于大数据的知识计算是大数据分析的基础.知识计算是国内外工业界开发和学术界研究的一个热点.要对数据进行高端分析,就需要从大数据中先抽取出有价值的知识,并把它构建成可支持查询、分析和计算知识库.目前,世界各国各个组织建立的知识库多达 50 余种,相关的应用系统更是达到了上百种.其中,代表性的知识库或应用系统有 KnowItAll^[43,44],TextRunner^[45],NEL^[46],Probase^[47],Satori^[48],PROSPERA^[49],SOFIE^[50]以及一些基于维基百科等在线百科知识构建的知识库,如 DBpedia^[51],YAGO^[52-54],Omega^[55]和 WikiTaxonomy^[56,57].除此之外,一些著名的商业网站、公司和政府也发布了类似的知识搜索和计算平台,如 Evi 公司的 TrueKnowledge 知识搜索平台,美国官方政府网站 Data.gov, Wolfram 的知识计算平台 wolframalpha, Google 的知识图谱 (knowledge graph)、Facebook 推出的类似的实体搜索服务 Graph Search 等.在国内,中文知识图谱的构建与知识计算也有大量的研究和开发工作.代表性工作有中国科学院计算技术研究所的 OpenKN,中国科学院数学研究院陆汝钤院士提出的知件(knowware),上海交通大学最早构建的中文知识图谱平台 zhishi.me,百度推出了中文知识图谱搜索,搜狗推出的知立方平台,复旦大学 GDM 实验室推出的中文知识图谱展示平台等.

支持知识计算的基础是构建知识库,这包括 3 个部分,即知识库的构建、多源知识的融合与知识库的更新.知识库的构建就是要构建几个基本的构成要素,包括抽取概念、实例、属性和关系.从构建方式上,可以分为手工构建和自动构建.手工构建是依靠专家知识编写一定的规则,从不同的来源收集相关的知识信息,构建知识的体系结构^[58].比较典型的例子是知网(Hownet)^[59]、同义词词林^[60]、概念层次网络(HNC)^[61]和中文概念词典(CCD)^[62],OpenCyc^[63]等.自动构建是基于知识工程、机器学习、人工智能等理论自动从互联网上采集并抽取概念、实例、属性和关系^[64,65].比较著名的例子是 Probase^[47],YAGO^[52-54]等.手工构建知识库,需要构建者对知识的领域有一定的了解,才能编写出合适的规则,开发过程中也需要投入大量的人力物力.相反地,自动构建的方法依靠系统自动的学习经过标注的语料来获取规则的,如属性抽取规则,关系抽取规则等,在一定程度上可以减少人工构建的工作量.随着大数据时代的到来,面对大规模网页信息中蕴含的知识,自动构建知识库的方法越来越受到人们的重视和青睐.自动构建知识库的方法主要分为有监督的构建方法和半监督的构建方法两种.有监督的构建方法是指系统通过学习训练数据,获取抽取规则,然后根据这些规则,提取同一类型的网页中的概念、实例、属性和关系.这类方法的缺点是规则缺乏普适性.而且,由于规则是针对特定网页的,当训练网页发生变化,需要重新进行训练来获取规则.半监督的构建方法是系统预先定义一些规则作为种子,然后通过机器学习算法,从标注语料中抽取相应的概念、实例、属性和关系.进一步地,系统根据抽取的结果,发现新的规则,再用来指导抽取相应的概念、实例、属性和关系,从而使抽取过程能够迭代的进行.

多源知识的融合是为了解决知识的复用问题.如前文所述,构建一个知识库的代价是非常大的,为了避免从头开始,需要考虑知识的复用和共享,这就需要对多个来源的知识进行融合,即需要对概念、实例、属性和关系的冲突,重复冗余,不一致进行数据的清理工作,包括对概念、实例进行映射、消歧,对关系进行合并等.这其中概念间关系或分类体系的融合是很关键一部分.按融合方式可以分为手动融合和自动融合.对于规模较小的知识库,手动融合是可行的,但这是一种非常费时而且容易出错的融合方式.相比于手动融合方式,建立在机器学习、人工智能和本体工程等算法上的融合方式具有更好的可扩展性,相关工作包括 YAGO^[52-54],Probase^[47]等.YAGO 知识库将维基百科,WordNet 和 GeoNames 等数据源的知识整合在知识库中.其中,将维基百科的分类体系和 WordNet 的分类体系进行融合是 YAGO 的重要的工作之一.维基百科的分类是一个有向无环图生成的层次结构^[52],这种结构由于仅能反映主题信息,所以容易出错.Probase 提出了一种基于概率化的实体消解(entity resolution)的知识整合技术^[46],将现有结构化数据,如 Freebase,IMDB,Amazon 等整合到 Probase 当中.对多源知识的融合,除了分类体系的融合外,还包括对实体和概念的消解问题,实体和概念的消歧问题等.面对海量知识库时,建立若干个针对不同领域,不同需求的有效的知识融合算法,快速进行多元知识的融合,是亟待进一步解决的问题之一.

大数据时代数据的不断发展与变化带给知识库构建的一个巨大的挑战是知识库的更新问题.知识库的更新分为两个层面,一是新知识的加入;二是已有知识的更改.目前专门针对开放网络知识库的更新工作较少,很多都是从数据库的更新角度展开的,如对数据库数据的增加、删除和修改工作的介绍.虽然对开放网络知识库的更新,与数据库的更新有很多相似之处,但是其本身对更新的实时性要求较高.目前这方面的工作,从更新方式来讲分为两类:一是基于知识库构建人员的更新;二是基于知识库存储的时空信息的更新.前者准确性较高,但是对人力的消耗较大.后者多由知识库自身更新,需要人工干预的较少,但是存在准确率不高的问题.总体上讲,对知识库的更新仍然没有很有效的方法.尤其在面对用户对知识的实时更新需求方面,远远达不到用户的要求.在更新数据的自动化感知方面,缺乏有效的办法自动识别知识的变化,也没有能够动态响应这些变化的更新机制.

2.3 社会计算

以 Facebook、Twitter、新浪微博、微信等为代表的在线社交网络和社会媒体正深刻改变着人们传播信息和获取信息的方式,人和人之间结成的关系网络承载着网络信息的传播,人的互联成为信息互联的载体和信息传播的媒介,社会媒体的强交互性、时效性等特点使其在信息的产生、消费和传播过程中发挥着越来越重要的作用,成为一类重要信息载体.正因如此,当前在线社会计算无论在学术圈和工业界都备受重视,大家关注的问题包括了对在线社会网络结构、信息传播以及信息内容的分析、建模与挖掘等一系列问题.

2.3.1 在线社会网络的结构分析

在线社会网络在微观层面上具有随机化无序的现象,在宏观层面上往往呈现出规则化、有序的现象,为了理清网络具有的这种看似矛盾的不同尺度的结构特性,探索和分析连接微观和宏观的网络中观结构(也称为社区结构)成为了本领域一个重要的研究方向.一般意义上讲,社区结构是指网络节点按照连接关系的紧密程度不同而自然分成若干个内部连接紧密、与外部连接稀疏的节点组,每个节点组相应地被称为社区^[66].社区分析研究目前主要包括社区的定义和度量、社区结构发现和社区结构演化性分析等基本问题^[67].

社区定义或度量大体上分为4类,基于节点的社区定义、基于节点组(社区)的社区定义、基于网络整体的社区定义、基于层次结构的社区定义.目前,社区结构的研究主要集中在基于某种给定社区定义或度量的社区发现上.最具代表性的社区发现算法包括密歇根大学 Newman 等人提出的模块度(modularity)优化方法^[68]、匈牙利科学院 Palla 等人提出的完全子图渗流(clique percolation)方法^[69]、华盛顿大学的 Rosvall 等人提出的基于网络最短编码的 InfoMap 方法^[70]、Airoldi^[71]等人提出的 Mixed Membership Stochastic Block(MMSB)模型,这些社区发现方法在人工构造的测试网络和一些小规模的真实网络上取得了很好的效果.真实世界在线社交网络中的社区结构具有多尺度、重叠等特点,近几年逐步引起研究人员的关注,成为一个研究热点^[72-76].

网络社区的演化性是信息网络的一个基本特性,也是促使大规模信息网络的内容与结构涌现现象及信息大规模传播的基本原因^[77].近几年,在前述社区发现研究的基础上,人们开始研究社区随时间演化的规律^[78].例如,Palla 等人基于完全子图渗流社区发现方法研究社区演化^[79],得到一个有趣结论,小社区的稳定性是保证它存在的前提,大社区的动态性是它存在的基础.Song 等人^[80]考虑了网络结构变化的时间因素,并认为网络演化过程是平滑的,他们使用扩展了的动态贝叶斯网络来建模网络的演化过程,取得了很好的效果.Xing 等人^[81]将网络演化的观点引入到结点的角色分析中.在 MMSB 模型中加入时间因素,他们认为两个相邻的时间片内角色选择方式和角色之间的关系具有一阶马尔可夫性质.此外,社区结构被用于预测网络中潜在存在的边,对于网络演化具有重要意义^[82].

2.3.2 在线社会网络的信息传播模型

在信息传播模型的研究中,最广泛深入研究的是传染病模型^[83,84],除了传染病模型,随机游走模型也是信息传播的基本模型之一^[85],作为最基本的动力学过程之一,随机游走与网络上的许多其他动力学过程(如反应-扩散过程、社团挖掘、路由选择,目标搜索)紧密相关.

近几年,研究人员开始注意到信息传播和传染病传播具有显著不同的特性^[86],包括信息传播的记忆性、社会增强效应、不同传播者的角色不同、消息内容的影响等.Romero 等人提出了 Stickiness 和 Persistence 两个

重要概念^[87],分析不同领域内的 Hashtag 在 Twitter 上的传播过程.Wu 等人分析名人、机构、草根等不同群体之间的消息流向,并分析了不同类型的消息被转发的情况及其生命周期^[88].Lerman 等人从网络动力学角度,通过实际数据分析了 Twitter 中消息传播的特性^[89].Castillo 等人通过特征提取,利用机器学习中分类的方法,对 Twitter 中消息的可信度建模,并预测其中消息的可信性^[90].Phelan 等人提出了一种 Twitter 消息新颖度的度量,并建立了向用户实时推荐新消息的系统^[91].Lerman 等人利用概率方法和先验知识,对 Digg 中的消息建模,预测消息的流行度^[92].当前,对在线社交网络中信息传播的研究主要集中在实证分析和统计建模,对于信息传播机理仍然缺乏深入的理解和有效的建模.

2.3.3 社交媒体中信息检索与数据挖掘

社交媒体的出现对信息检索与数据挖掘的研究提出了新的挑战.不同于传统的 Web 数据,社交媒体中的数据呈现出一些新的特征:(1) 信息碎片化现象明显,文本内容特征越发稀疏;(2) 信息互联被人的互联所取代,社交媒体用户形成的社会关系网络的搜索和挖掘过程中的重要组成部分;(3) 社交媒体的易参与性使得人人具有媒体的特征,呈现出自媒体现象,个人影响力、情感与倾向性掺杂其中.针对这些特点,研究人员在传统信息检索与数据挖掘技术基础上提出了一系列的新模型^[93,94].

鉴于用户所创造的信息往往具有很强的时效性,Yang 等人提出了一种时间序列聚类的方法,从 Twitter 数据中挖掘热门话题发展趋势的规律^[95].因为用户的状态和评论中包含了大众的观点和态度,所以 Bollen 等人通过对 Twitter 中用户的信息进行情感分析,将大众情绪的变化表示为 7 种不同的情绪时间序列,进而发现这些序列能够预测股票市场的走势^[96].此外,基于用户在协作平台上所贡献的内容和标签等信息往往蕴含着丰富的大众知识和智慧这一现象,Hu 等人利用 Wikipedia 中的文章和类别信息来确定用户的查询意图,进而辅助信息检索^[97].社交媒体的检索与挖掘研究在国内也受到了越来越多的重视,包括北京大学、清华大学、哈尔滨工业大学、上海交通大学、浙江大学、复旦大学、中国科学院、微软亚洲研究院等大学和研究机构已经取得了一定的进展,涉及的研究内容包括社会化标签系统中的标签学习和排序^[98,99]、信息抽取和分类^[100]、社会化多媒体检索^[101]、协作搜索和推荐^[102,103]等等.

2.4 可视化

大数据引领着新一波的技术革命,对大数据查询和分析的实用性和实效性对于人们能否及时获得决策信息非常重要,决定着大数据应用的成败.但产业界面对大数据常常显得束手无策.一是因为数据容量巨大,类型多样,数据分析工具面临性能瓶颈.另一原因在于,数据分析工具通常仅为 IT 部门熟练使用,缺少简单易用、让业务人员也能轻松上手实现自助自主分析即时获取商业洞察的工具.因此,数据可视化技术正逐步成为大数据时代的显学.对大数据进行分析以后,为了方便用户理解也需要有效的可视化技术,这其中交互式的展示和超大图的动态化展示值得重点关注.

大数据可视化,不同于传统的信息可视化,面临最大的一个挑战就是规模,如何提出新的可视化方法能够帮助人们分析大规模、高维度、多来源、动态演化的信息,并辅助作出实时的决策,成为了这个领域最大的挑战.为了解决这个问题,我们可以依赖的主要手段是两种,即数据转换和视觉转换.现有研究工作主要聚焦在 4 个方面:(1) 通过对信息流进行压缩或者删除数据中的冗余信息对数据进行简化.其中很多工作主要解决曲面的可视化,使用基本的数据转换方法来对数据进行简化.例如,文献^[104,105]提出通过删除节点以及包含这个节点的三角形进行网络的简化,而 Hoppe 等人^[106,107]则提出了一种渐进网格表达方法,通过有效地删除边及其所属的三角形实现.一些研究人员把上述的这些曲面算法进一步扩展到四面体上^[108,109].上述这些工作存在主要的不同之处在于基本的数据转换步骤中使用不同的错误近似方法.(2) 通过设计多尺度、多层次的方法实现信息在不同的解析度上的展示,从而使用户可自主控制展示解析度.很多已有多尺度算法集中在对地形类数据的渲染上.例如,一些使用固定网格方法的系统建立在直角三角形的层次结构之上^[110,111],而一些不规则三角形网^[112,113]则是通过不把三角形限制在固定网格上的方式来解决这个问题,这两类方法各有利弊.Cignoni 等人^[114,115]则通过利用四叉树纹理层次以及利用三角片面二叉树显示几何形,展示了实时显示大型地形数据集中自适应的几何形和纹理的能力.(3) 利用创新的方法把数据存储在外存,并让用户可以通过交互手段方便地获取相关数据,

这类研究也成为核外算法(out-of-core algorithm),为了应对大规模数据结构无法在内存中存放,而外存访问时间又极大地依赖于外部存储单元的位置,人们设计了一些新的算法与分析工具来解决几何算法^[116,117]以及可视化方法。这类工作重点解决两个问题:a)对算法进行分析得到数据访问的模式,从而重新设计数据结构来最大化访问的局部性;b)在二级存储设备中的数据需要有与算法访问模式匹配的存放布局。Pascucci和Frank^[118]引入了一种新的静态索引体系使得在层次化遍历 n 维规则网格时候,数据的存放布局能够满足上述两个要求。(4)提出新的视觉隐喻方法以全新的方式展示数据,其中,一类典型的方法是“焦点+上下文”方法,它重点对焦点数据进行细节展示,对不重要数据的则简化表示,例如鱼眼视图^[119]。Plaisant提出了空间树(space tree)^[120],一种树形浏览器通过动态调整树枝的尺寸来使其最好地适配显示区域。分层平行坐标方法,作为平行坐标方法的多尺度版本,通过在不同的细节层次使用多的视图来对大规模数据进行表达。

对大数据进行探索和可视化仍然还处在初始阶段,特别是对于动态多维度大数据流的可视化技术还非常匮乏,非常需要扩展现有的可视化算法,研究新的数据转换方法以便能够应对复杂的信息流数据。也需要设计创新的交互方式来对大数据进行可视化交互和辅助决策。

2.5 小结

大数据处理和分析的终极目标是借助对数据的理解辅助人们在不同应用中作出合理的决策。在此过程中,深度学习、知识计算、社会计算和可视化起到了相辅相成的作用。

(1)深度学习提高精度:如前所述,要挖掘大数据的大价值必然要对大数据进行内容上的分析与计算,而传统的数据表达模型和方法通常是简单的浅层模型学习,效果不尽人意。深度学习可以对人类难以理解的底层数据特征进行层层抽象,凝练具有物理意义的特征,从而提高数据学习的精度。因此,深度学习是大数据分析的核心技术;

(2)知识计算挖掘深度:每一种数据来源都有一定的局限性和片面性,只有对各种来源的原始数据进行融合才能反映事物的全貌,事物的本质和规律往往隐藏在各种原始数据的相互关联之中。而借助知识计算可以将碎片化的多源数据整合成反映事物全貌的完整数据,从而增加数据挖掘的深度。因此,基于大数据的知识计算是大数据分析的基础。如何基于大数据实现新知识的感知,知识的增量式演化和自适应学习是其中的重大挑战;

(3)社会计算促进认知:IT技术的发展使得社交媒体成了一类重要的信息载体,承载着对事物的客观或主观描述信息。因此,通过基于社交媒体数据的社会计算可以促进人们对事物的认知。但是,社交媒体大数据往往蕴含着一个体量庞大、关系异质、结构多尺度和动态演化的网络,对它的分析既要有效地计算方法,更需要支持大规模网络结构的图数据存储和管理结构,以及高性能的图计算系统结构和算法;

(4)强可视化辅助决策:对大数据查询和分析的实用性和实效性对于人们能否及时获得决策信息非常重要。而强大的可视化技术,不仅可以对数据分析结果进行更有效的展示,而且可以在大数据分析过程中发挥重要作用。

3 大数据计算面临的挑战与应对之策

尽管大数据是社会各界都高度关注的话题,但时下大数据从底层的处理系统到高层的分析手段都存在许多问题,也面临一系列挑战。这其中有大数据自身的特征导致的,也有当前大数据分析模型与方法引起的,还有大数据处理系统所隐含的。本节对这些问题与挑战进行梳理。

3.1 数据复杂性带来的挑战

大数据的涌现使人们处理计算问题时获得了前所未有的大规模样本,但同时也不得不对更加复杂的对象,如前所述,其典型的特性是类型和模式多样、关联关系繁杂、质量良莠不齐。大数据内在的复杂性(包括类型的复杂、结构的复杂和模式的复杂)使得数据的感知、表达、理解和计算等多个环节面临着巨大的挑战,导致了传统海量数据计算模式下时空维度上计算复杂度的激增,传统的数据分析与挖掘任务如检索、主题发现、语义和情感分析等变得异常困难。然而目前,人们对大数据复杂性的内在机理及其背后的物理意义缺乏理

解,对大数据的分布与协作关联等规律认识不足,对大数据的复杂性和计算复杂性的内在联系缺乏深刻理解,加上缺少面向领域的大数据处理知识,极大地制约了人们对大数据高效计算模型和方法的设计能力。

因此,如何形式化或定量化地描述大数据复杂性的本质特征及其外在度量指标,进而研究数据复杂性的内在机理是个根本问题.通过对大数据复杂性规律的研究有助于理解大数据复杂模式的本质特征和生成机理,简化大数据的表征,获取更好的知识抽象,指导大数据计算模型和算法的设计.为此,需要建立多模态关联关系下的数据分布理论和模型,理清数据复杂度和时空计算复杂度之间的内在联系,通过对数据复杂性内在机理的建模和解析,阐明大数据按需约简、降低复杂度的原理与机制,使其成为大数据计算的理论基石。

3.2 计算复杂性带来的挑战

大数据多源异构、规模巨大、快速多变等特性使得传统的机器学习、信息检索、数据挖掘等计算方法不能有效支持大数据的处理、分析和计算.特别地,大数据计算不能像小样本数据集那样依赖于对全局数据的统计分析和迭代计算,需要突破传统计算对数据的独立同分布和采样充分性的假设.在求解大数据的问题时,需要重新审视和研究它的可计算性、计算复杂性和求解算法.因此,研究面向大数据的新型高效计算范式,改变人们对数据计算的本质看法,提供处理和分析大数据的基本方法,支持价值驱动的特定领域应用,是大数据计算的核心问题.而大数据样本量充分,内在关联关系密切而复杂,价值密度分布极不均衡,这些特征对研究大数据的可计算性及建立新型计算范式提供了机遇,同时也提出了挑战。

因此,需要着眼于大数据的全生命周期,基于大数据复杂性的基本特征及其量化指标,研究大数据下以数据为中心的計算模式,突破传统的数据围绕机器式计算,构建以数据为中心的推送式计算模式,探索弱 CAP 约束的系统架构模型及其代数计算理论,研究分布化、流式计算算法,形成通信、存储、计算融合优化的大数据计算框架;研究适应大数据的非确定性算法理论,突破传统统计学习中的独立同分布假设;也需要探索从足够多 (large enough) 的数据,到刚刚好 (just enough) 的数据,再到有价值 (valuable enough) 的数据的按需约简方法,研究基于自举和采样的局部计算和近似方法,提出不依赖于全量数据的新颖算法理论基础。

3.3 系统复杂性带来的挑战

针对不同数据类型与应用的大数据处理系统是支持大数据科学研究的基础平台.对于规模巨大、结构复杂、价值稀疏的大数据,其处理亦面临计算复杂度高、任务周期长、实时性要求强等难题.大数据及其处理的这些难点不仅对大数据处理系统的系统架构、计算框架、处理方法提出了新的挑战,更对大数据处理系统的运行效率及单位能耗提出了苛刻要求,要求大数据处理系统必须具有高效能的特点.对于以高效能为目标的大数据处理系统的系统架构设计、计算框架设计、处理方法设计和测试基准设计研究,其基础是大数据处理系统的效能评价与优化问题研究.这些问题的解决可奠定大数据处理系统设计、实现、测试与优化的基本准则,是构建能效优化的分布式存储和处理的硬件及软件系统架构的重要依据和基础,因此是大数据分析处理所必须解决的关键问题。

大数据处理系统的效能评价与优化问题具有极大的研究挑战性,其解决不但要求理清大数据的复杂性、可计算性与系统处理效率、能耗间的关系,还要综合度量系统中如系统吞吐率、并行处理能力、作业计算精度、作业单位能耗等多种效能因素,更涉及实际负载情况及资源分散重复情况的考虑.因此,为了解决系统复杂性带来的挑战,人们需要结合大数据的价值稀疏性和访问弱局部性的特点,针对能效优化的大数据分布存储和处理的系统架构,以大数据感知、存储与计算融合为大数据的计算准则,在性能评价体系、分布式系统架构、流式数据计算框架、在线数据处理方法等方面展开基础性研究,并对作为重要验证工具的基准测试程序及系统性能预测方法进行研究,通过设计、实现与验证的迭代完善,最终实现大数据计算系统的数据获取高吞吐、数据存储低能耗和数据计算高效率。

4 结束语

互联网、物联网、云计算技术的快速发展,各类应用的层出不穷引发了数据规模的爆炸式增长,使数据渗

透到了当今每一个行业和业务领域,成为重要的生产因素.大数据因此成为社会各界关注的新焦点,大数据时代已然来临.为了应对不同的业务需求,以 Google,Facebook,Linkedin,Microsoft 等为代表的互联网企业近几年推出了各种大数据处理系统,深度学习、知识计算、可视化等大数据分析技术也得到迅速发展,已被广泛应用于不同的行业和领域.本文根据处理形式的不同,介绍了批量处理数据、流式处理数据、交互处理数据和图数据四种不同形式数据的突出特征和各自的典型应用场景以及相应的代表性处理系统,并总结出引擎专用化、平台多样化、计算实时化是当前大数据处理系统的三大发展趋势.随后,对系统支撑下的深度学习、知识计算、社会计算与可视化四类大数据分析技术和应用进行了简要综述,总结了各种技术在大数据分析理解过程中的关键作用,即深度学习提高精度,知识计算挖掘深度,社会计算促进认知,强可视化辅助决策.本文最后梳理了大数据处理和分析面临的 3 个核心挑战,包括数据复杂性、计算复杂性和系统复杂性,并提出了可能的应对之策.

References:

- [1] Manyika J, Chui M, Brown B, Bughin J, Dobbs R, Roxburgh C, Byers AH. Big data: The next frontier for innovation, competition, and productivity. 2011. http://www.mckinsey.com/insights/business_technology/big_data_the_next_frontier_for_innovation
- [2] Li GJ, Cheng XQ. Research Status and Scientific Thinking of Big Data. *Bulletin of the Chinese Academy of Sciences*, 2012, 27(6): 647–657.
- [3] Wang YZ, Jin XL, Cheng XQ. Network big data: Present and future. *Chinese Journal of Computers*, 2013,36(6):1125–1138.
- [4] Arthur WB. The second economy. 2011. <http://www.images-et-reseaux.com/sites/default/files/medias/blog/2011/12/the-2nd-economy.pdf>
- [5] Mayer-Schönberger V, Cukier K. Big data: A revolution that will transform how we live, work, and think. Houghton Mifflin Harcourt, 2013.
- [6] Sun DW, Zhang GY, Zheng WM. Big data stream computing: Technologies and instances. *Ruan Jian Xue Bao/Journal of Software*, 2014,25(4):839-862.
- [7] Tsourakakis CE. Fast counting of triangles in large real networks without counting: *Algorithms and Laws*, 2008. 608–617. [doi: 10.1109/ICDM.2008.72]
- [8] Chen Y, Alspaugh S, Katz R. Interactive analytical processing in big data systems: A cross-industry study of MapReduce workloads. *Proc. of the VLDB Endowment*, 2012,5(12):1802–1813. [doi: 10.14778/2367502.2367519]
- [9] Stupar A, Michel S, Schenkel R. RankReduce-Processing k -nearest neighbor queries on top of MapReduce. *Large-Scale Distributed Systems for Information Retrieval*, 2010. 13–18.
- [10] Zhou MQ, Zhang R, Xie W, Qian WN, Zhou AY. Security and privacy in cloud computing: A survey. *IEEE*, 2010. 105–112. [doi: 10.1109/SKG.2010.19]
- [11] Feblowitz J. Analytics in oil and gas: The big deal about big data. In: *Proc. of the SPE Digital Energy Conf*. 2013. [doi: 10.2118/163717-MS]
- [12] Yu H, Wang D. Research and implementation of massive health care data management and analysis based on hadoop. *IEEE*, 2012. 514–517. [doi: 10.1109/ICCIS.2012.225]
- [13] Ghemawat S, Gobiuff H, Leung S-T. The Google file system. *ACM*, 2003,37(5):29–43. [doi: 10.1145/1165389.945450]
- [14] Dean J, Ghemawat S. MapReduce: Simplified data processing on large clusters. *Communications of the ACM*, 2008,51(1): 107–113. [doi: 10.1145/1327452.1327492]
- [15] Hadoop. <http://hadoop.apache.org/>
- [16] Dean J, Ghemawat S. MapReduce: A flexible data processing tool. *Communications of the ACM*, 2010,53(1):72–77.
- [17] White T. Hadoop: The definitive guide. O'Reilly Media, Inc., 2012.
- [18] Chakravarthy, Sharma, Jiang Q. Stream data processing: A quality of service perspective: Modeling, scheduling, load shedding, and complex event processing, Springer-Verlag, 2009.
- [19] Storm. <http://storm.incubator.apache.org/>
- [20] Kafka Doc. <http://kafka.apache.org/documentation.html>
- [21] Goodhope K, Koshy J, Kreps J, Narkhede N, Park R, Rao J, Ye VY. Building LinkedIn's real-time activity data pipeline. *IEEE Data Engineering Bulletin*, 2012,35(2):33–45.

- [22] Hive. <https://hive.apache.org/>
- [23] Pig. <https://pig.apache.org/>
- [24] Hbase. <https://hbase.apache.org/>
- [25] MongoDB. <http://www.mongodb.org>
- [26] Zaharia M, Chowdhury M, Franklin M, Shenker S, Stoica I. Spark: Cluster computing with working sets. HotCloud 2010. 2010.
- [27] Melnik S, Gubarev A, Long JJ, Romer G, Shivakumar S, Tolton M, Vassilakis T. Dremel: Interactive analysis of Web-scale datasets. Proc. of the VLDB Endowment, 2010,3(1-2):330–339.
- [28] Malewicz G, Austern MH, Bik AJ, Dehnert JC, Horn I, Leiser N, Czajkowski G. Pregel: A system for large-scale graph processing. In: Proc. of the 2010 ACM SIGMOD Int'l Conf. on Management of Data. ACM, 2010. 135–146. [doi: 10.1145/1327452.1327492]
- [29] Pregel. <http://kowshik.github.io/JPregel/>
- [30] Neo4j. <http://www.neo4j.org/>
- [31] Trinity. <http://research.microsoft.com/trinity>
- [32] Shao B, Wang H, Li Y. Trinity: A distributed graph engine on a memory cloud. In: Proc. of the 2013 Int'l Conf. on Management of Data. ACM, 2013. 505–516. [doi: 10.1145/2463676.2467799]
- [33] Hinton G, Osindero S, Teh Y W. A fast learning algorithm for deep belief nets. Neural Computation, 2006,18(7):1527–1554. [doi: 10.1162/neco.2006.18.7.1527]
- [34] Bengio Y, Lamblin P, Popovici D, Larochelle H. Greedy layer-wise training of deep networks. In: Advances in Neural Information Processing Systems. 2007,19:153.
- [35] Dahl GE, Yu D, Deng L, Acero A. Context-Dependent pre-trained deep neural networks for large-vocabulary speech recognition. IEEE Trans. on Audio, Speech, and Language Processing, 2012,20(1):30–42. [doi: 10.1109/TASL.2011.2134090]
- [36] Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems. 2012. 1097–1105.
- [37] ImageNet. <http://www.image-net.org/challenges/LSVRC/2013/results.php>
- [38] Taigman Y, Yang M, Ranzato M, Wolf L. Deepface: Closing the gap to human-level performance in face verification. In: Proc. of the IEEE CVPR. 2014.
- [39] Le QV. Building high-level features using large scale unsupervised learning. In: Proc. of the 2013 IEEE Int'l Conf. on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2013. 8595–8598. [doi: 10.1109/ICASSP.2013.6639343]
- [40] Bengio Y, Schwenk H, Senécal J S, Morin F, Gauvain JL. Neural probabilistic language models. In: Innovations in Machine Learning. Berlin, Heidelberg: Springer-Verlag, 2006. 137–186. [doi: 10.1007/3-540-33486-6_6]
- [41] Collobert R, Weston J. A unified architecture for natural language processing: Deep neural networks with multitask learning. In: Proc. of the 25th Int'l Conf. on Machine Learning. ACM, 2008. 160–167. [doi: 10.1145/1390156.1390177]
- [42] Socher R, Perelygin A, Wu J Y, Chuang J, Manning CD, Ng AY, Potts C. Recursive deep models for semantic compositionality over a sentiment treebank. In: Proc. of the Conf. on Empirical Methods in Natural Language Processing (EMNLP). 2013. 1631–1642.
- [43] Etzioni O, Cafarella M, Downey D, Popescu AM, Shaked T, Soderland S, Yates A. Unsupervised named-entity extraction from the Web: An experimental study. Artificial Intelligence, 2005,165(1):91–134. [doi: 10.1016/j.artint.2005.03.001]
- [44] Etzioni O, Cafarella M, Downey D, Kok S, Popescu AM, Shaked T, Yates A. Web-Scale information extraction in knowitall: (preliminary results). In: Proc. of the 13th Int'l Conf. on World Wide Web. ACM, 2004. 100–110. [doi: 10.1145/988672.988687]
- [45] Banko M, Cafarella MJ, Soderland S, Broadhead M, Etzioni O. Open information extraction for the Web. In: Proc. of the IJCAI. 2007,7:2670–2676.
- [46] Carlson A, Betteridge J, Kisiel B, Settles B, Hruschka Jr ER, Mitchell TM. Toward an architecture for neverending language learning. In: Proc. of the 24th AAAI Conf. on Artificial Intelligence. Menlo Park: AAAI Press, 2010. 1306–1313.
- [47] Wu W, Li H, Wang H, Zhu KQ. Probbase: A probabilistic taxonomy for text understanding. In: Proc. of the 2012 ACM SIGMOD Int'l Conf. on Management of Data. New York: ACM, 2012. 481–492 [doi: 10.1145/2213836.2213891]
- [48] Gallagher S. How Google and Microsoft taught search to understand the Web. 2012. <http://arstechnica.com/information-technology/2012/06/inside-the-architecture-of-googles-knowledge-graphand-microsofts-satori/>

- [49] Nakashole N, Theobald M, Weikum G. Scalable knowledge harvesting with high precision and high recall. In: Proc. of the 4th ACM Int'l Conf. on Web Search and Data Mining. New York: ACM, 2011. 227–236. [doi: 10.1145/1935826.1935869]
- [50] Suchanek FM, Sozio M, Weikum G. SOFIE: A self-organizing framework for information extraction. In: Proc. of the 18th Int'l Conf. on World Wide Web. New York: ACM, 2009. 631–640. [doi: 10.1145/1526709.1526794]
- [51] Auer S, Bizer C, Kobilarov G, Lehmann J, Cyganiak R, Ives Z. DBpedia: A nucleus for a Web of open data. In: Proc. of the 6th Int'l the Semantic Web and the 2nd Asian Conf. on Asian Semantic Web Conf., ISWC 2007. Piscataway: IEEE, 2007. 722–735. [doi: 10.1007/978-3-540-76298-0_52]
- [52] Biega J, Kuzey E, Suchanek FM. Inside YAGO2s: A transparent information extraction architecture. In: Proc. of the 22th Int'l Conf. on World Wide Web. New York: ACM, 2013. 325–328.
- [53] Hffart J, Suchanek F, Berberich K, Weikum G. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *Artificial Intelligence Journal*, 2013,194(4):28–61. [doi: 10.1016/j.artint.2012.06.001]
- [54] Suchanek F, Kasneci G, Weikum G. YAGO—A core of semantic knowledge. In: Proc. of the 16th Int'l Conf. on World Wide Web. New York: ACM, 2007. 697–706. [doi: 10.1145/1242572.1242667]
- [55] Philpot A, Hovy EH, Pantel P. *Ontology and the Lexicon//The Omega Ontology*. Cambridge: Cambridge University Press, 2008. 35–78.
- [56] Ponzetto S, Navigli R. Large-Scale taxonomy mapping for restructuring and integrating wikipedia. In: Proc. of the 21st Int'l Joint Conf. on Artificial Intelligence, IJCAI 2009. San Francisco: Morgan Kaufmann Publishers, 2009. 2083–2088.
- [57] Ponzetto S, Strube M. Taxonomy induction based on a collaboratively built knowledge repository. *Artificial Intelligence*, 2011, 175(9-10):1737–1756. [doi: 10.1016/j.artint.2011.01.003]
- [58] Shi ZZ. *Knowledge Discovery*. Beijing: Tsinghua University Press, 2002.
- [59] Dong ZD, Dong Q, Hao CL. Theoretical findings of HowNet. *Journal of Chinese Information Processing*, 2007,21(4):3–9.
- [60] Mei LJ, Zhou Q, Cang L, Chen ZS. Merge information in HowNet and TongYiCi CiLin. *Journal of Chinese Information Processing*, 2005,19(1):63–70.
- [61] 黄曾阳.HNC 理论概要. *中文信息学报*,1997,11(4):11–20.
- [62] Yu JS, Yu SW. The structure of Chinese concept dictionary. *Journal of Chinese Information Processing*, 2002,16(4):12–20.
- [63] Xu WY, Liu SY. Logic for knowledgebase systems. *Chinese Journal of Computers*, 2009,32(11):2123–2129.
- [64] Zhong XQ, Liu Z, Dong PP. Construction of knowledge base on hybrid reasoning and its application. *Chinese Journal of Computers*, 2012,35(4):761–766
- [65] Chen LW, Feng YS, Zhao DY. Extracting relations from the Web via weakly supervised learning. *Journal of Computer Research and Development*, 2013,50(9):1825–1835.
- [66] Girvan M, Newman MEJ. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 2002,99(12):7821–7826. [doi: 10.1073/pnas.122653799]
- [67] Fortunato S. Community detection in graphs. *Physics Reports*, 2010,486(3):75–174.
- [68] Newman MEJ. Finding community structure in networks using the eigenvectors of matrices. *Physical Review E*, 2006,74(3): 36–104. [doi: 10.1103/PhysRevE.74.036104]
- [69] Palla G, Derényi I, Farkas I, Vicsek T. Uncovering the overlapping community structure of complex networks in nature and society. *Nature*, 2005,435(7043):814–818. [doi:10.1038/nature03607]
- [70] Rosvall M, Bergstrom CT. Maps of information flow reveal community structure in complex networks. *PNAS*, 2008,105(4): 1118–1123. [doi: 10.1073/pnas.0706851105]
- [71] Airoldi EM, Blei DM, Fienberg SE, Xing EP. Mixed membership stochastic blockmodel. *Journal of Machine Learning Research*, 2008,9:1981–2014.
- [72] Fortunato S, Barthélemy M. Resolution limit in community detection. *Proceedings of the National Academy of Sciences*, 2007,104(1):36–41. [doi: 10.1073/pnas.0605965104]
- [73] Sales-Pardo M, Guimera R, Moreira AA, Amaral LAN. Extracting the hierarchical organization of complex systems. *Proceedings of the National Academy of Sciences*, 2007,104(39):15224–15229. [doi: 10.1073/pnas.0703740104]

- [74] Mucha PJ, Richardson T, Macon K, Porter MA, Onnela JP. Community structure in time-dependent, multiscale, and multiplex networks. *Science*, 2010,328(5980):876–878. [doi: 10.1126/science.1184819]
- [75] Delvenne JC, Yaliraki SN, Barahona M. Stability of graph communities across time scales. *Proceedings of the National Academy of Sciences*, 2010,107(29): 12755–12760. [doi: 10.1073/pnas.0903215107]
- [76] Ahn YY, Bagrow JP, Lehmann S. Link communities reveal multiscale complexity in networks. *Nature*, 2010,466(7307):761–764. [doi:10.1038/nature09182]
- [77] Pastor-Satorras R, Vespignani A. Epidemic spreading in scale-free networks. *Physical Review Letters*, 2001. [doi: 10.1103/PhysRevLett.86.3200]
- [78] Hopcroft J, Khan O, Kulis B, Selman B. Tracking evolving communities in large linked networks. *Proceedings of the National Academy of Sciences*, 2004,6,101(Suppl 1):5249–5253. [doi: 10.1073/pnas.0307750100]
- [79] Palla G, Barabási AL, Vicsek T. Quantifying social group evolution. *Nature*, 2007,446(7136):664–667. [doi:10.1038/nature05670]
- [80] Song L, Kolar M, Xing EP. Time-Varying dynamic bayesian networks. In: *Proc. of the 23rd Neural Information Processing Systems (NIPS 2009)*. 2009.
- [81] Xing EP, Fu W, Song L. A state-space mixed membership blockmodel for dynamic network tomography. *Annals of Applied Statistics*, 2010,4(2):535–566. [doi: 10.1214/09-AOAS311]
- [82] Clauset A, Moore C, Newman MEJ. Hierarchical structure and the prediction of missing links in networks. *Nature*, 2008,453: 98–101. [doi:10.1038/nature06830]
- [83] Anderson RM, May RMC. *Infectious Diseases of Humans: Dynamics and Control*. Oxford Science Publications, 1992.
- [84] Hethcote HW, Van Den Driessche P. Two SIS epidemiologic models with delays. *Journal of Mathematical Biology*, 2000,40:3–26. [doi: 10.1007/s002850050003]
- [85] Noh JD, Rieger H. Random walks on complex networks. *Physical Review Letters*, 2004,92(11):118701.
- [86] Lü L, Chen DB, Zhou T. The small world yields the most effective information spreading. *New Journal of Physics*, 2011,13: 123005. [doi: 10.1088/1367-2630/13/12/123005]
- [87] Romero DM, Meeder B, Kleinberg J. Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on Twitter. In: *Proc. of the 20th Int'l Conf. on World Wide Web (WWW)*. 2011. 695–704. [doi: 10.1145/1963405.1963503]
- [88] Wu S, Hofman JM, Mason WA, Watts DJ. Watts. Who says what to whom on Twitter. In: *Proc. of the 20th Int'l Conf. on World Wide Web*. Hyderabad, 2011. 705–714. [doi: 10.1145/1963405.1963504]
- [89] Lerman K, Ghosh R. Information contagion: An empirical study of the spread of news on Digg and Twitter social network. In: *Proc. of the 4th Int'l Conf. on Weblogs and Social Media (ICWSM)*. Washington, 2010. 90–97.
- [90] Castillo C, Mendoza M, Poblete B. Information credibility on Twitter. In: *Proc. of the 20th Int'l Conf. on World Wide Web*. 2011. 675–684. [doi: 10.1145/1963405.1963500]
- [91] Phelan O, McCarthy K, Smyth B. Using Twitter to recommend real-time topical news. In: *Proc. of the 3rd ACM Conf. on Recommender Systems (RecSys)*. 2009. 385–388. [doi: 10.1145/1639714.1639794]
- [92] Lerman K, Hogg T. Using a model of social dynamics to predict popularity of news. In: *Proc. of the 19th Int'l Conf. on World Wide Web*. 2011. 621–630. [doi: 10.1145/1772690.1772754]
- [93] Cheng XQ, Guo JF, Jin XL. A retrospective of Web information retrieval and mining. *Journal of Chinese Information Processing*, 2011,25(6):111–117
- [94] 沈华伟,靳小龙,任福新,程学旗. 面向社会媒体的舆情分析. *中国计算机学会通讯*, 2012,8(4):32–36.
- [95] Yang J, Leskovec J. Patterns of temporal variation in online media. In: *Proc. of the 4th ACM Int'l Conf. on Web Search and Data Mining (WSDM 2011)*. 2011. 177–186 [doi: 10.1145/1935826.1935863]
- [96] Bollen J, Mao H, Zeng X. Twitter mood predicts the stock market. *Journal of Computational Science*, 2011,2(1):1–8. [doi 10.1016/j.jocs.2010.12.007]
- [97] Hu J, Wang G, Lochovsky F, Sun JT, Sun CZ, Chen Z. Understanding user's query intent with Wikipedia. In: *Proc. of the 18th Int'l Conf. on World Wide Web (WWW 2009)*. 2009. 471–480. [doi: 10.1145/1526709.1526773]

- [98] Wu L, Yang L, Yu N, Hua XS. Learning to tag. In: Proc. of the WWW 2009 MADRID. ACM, 2009. 361–370. [doi: 10.1145/1526709.1526758]
- [99] Liu D, Hua XS, Yang L, Wang M, Zhang HJ. Tag ranking. In: Proc. of the WWW 2009 MADRID. ACM, 2009. 351–360. [doi: 10.1145/1526709.1526757]
- [100] Luo P, Lin F, Xiong Y, Zhao Y, Shi Z. Towards combining Web classification and Web information extraction: A case study. In: Proc. of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2009. 1235–1244. [doi: 10.1145/1557019.1557152]
- [101] Qi GJ, Hua XS, Zhang HJ. Learning semantic distance from community-tagged media collection. *ACM Multimedia*, 2009. 243–252. [doi: 10.1145/1631272.1631307]
- [102] Xue GR, Han J, Yu Y, Yang Q. User language model for collaborative personalized search. *ACM Trans. on Information Systems*, 2009,27(2):1–28. [DOI: 10.1145/1462198.1462203]
- [103] Li B, Yang Q, Xue X. Transfer learning for collaborative filtering via a rating-matrix generative model. In: Proc. of the ICML. 2009. 78. [doi: 10.1145/1553374.1553454]
- [104] Schroeder WJ, Zarge JA, Lorensen WE. Decimation of triangle meshes. *Computer Graphics*, 1992,26(2): 65–70. [doi: 10.1145/133994.134010]
- [105] Renze KJ, Oliver JH. Generalized unstructured decimation. *IEEE Computer Graphics and Applications*, 1996,16(6):24–32. [doi: 10.1109/38.544069]
- [106] Hoppe H. Progressive meshes. In: Proc. of the SIGGRAPH. 1996. 99–108. [doi: 10.1145/237170.237216]
- [107] Hoppe H. View-Dependent refinement of progressive meshes. In: Proc. of the SIGGRAPH. 1997. [doi: 10.1145/258734.258843]
- [108] Chopra P, Meyer J. Tetfusion: An algorithm for rapid tetrahedral mesh simplification. In: Proc. of the Conf. on Visualization 2002. Washington: IEEE Computer Society, 2002. 133–140. [doi: 10.1109/VISUAL.2002.1183767]
- [109] Staadt OG, Gross MH. Progressive tetrahedralizations. In: Ebert D, Hagen H, Rushmeier H, eds. Proc. of the Visualization'98. Los Alamitos: IEEE Computer Society, 1998. 397–402. [doi: 10.1109/VISUAL.1998.745329]
- [110] Evans W, Kirkpatrick D, Townsend G. Right triangular irregular networks. Technical Report, TR97-09, Department of Computer Science, University of Arizona, 1997.
- [111] Mirante A, Weingarten N. The radial sweep algorithm for constructing triangulated irregular networks. *IEEE Computer Graphics and Applications*, 1982,2(3):11–13, 15–21. [doi: 10.1109/MCG.1982.1674214]
- [112] Fowler RJ, Little JJ. Automatic extraction of irregular network digital terrain models. *Computer Graphics (SIGGRAPH'79)*, 1979,13(2):199–207. [doi: 10.1145/965103.807444]
- [113] Silva CT, Mitchell JSB, Kaufman AE. Automatic generation of triangular irregular networks using greedy cuts. In: Proc. of the IEEE Visualization (1995). Los Alamitos: IEEE Computer Society, 1995. 201–208. [doi: 10.1109/VISUAL.1995.480813]
- [114] Cignoni P, Ganovelli F, Gobbetti E, Marton F, Ponchio F, Scopigno R. BDAM: Batched dynamic adaptive meshes for high performance terrain visualization. In: Brunet P, Fellner D, eds. Proc. of the 24th Annual Conf. of the European Association for Computer Graphics (EG 2003), Vol.22. Blackwell: IEEE Computer Society, 2003. 505–514. [doi:10.1111/1467-8659.00698]
- [115] Cignoni P, Ganovelli F, Gobbetti E, Marton F, Ponchio F, Scopigno R. Interactive out-of-core visualization of very large landscapes on commodity graphics platforms. In: Proc. of the ICVS 2003. LNCS, New York: Springer-Verlag, 2003. 21–29. [doi: 10.1007/978-3-540-40014-1_3]
- [116] Goodrich MT, Tsay JJ, Vengroff DE, Vitter JS. Externalmemory computational geometry. In: Proc. of the 34th Annual IEEE Symp. on Foundations of Computer Science (FOCS 1993), 1993. 714–723. [doi: 10.1109/SFCS.1993.366816]
- [117] Matias Y, Segal E, Vitter JS. Efficient bundle sorting. In: Proc. of the 11th Annual ACM-SIAM Symp. on Discrete Algorithms (2000). Society for Industrial and Applied Mathematics, 2000. 839–848.
- [118] Pascucci V, Frank RJ. Global static indexing for real-time exploration of very large regular grids. In: Proc. of the 2001 ACM/IEEE Conf. on Supercomputing (CDROM). New York: ACM, 2001. [doi: 10.1145/582034.582036]
- [119] Plaisant C, Carr D, Shneiderman B. Image-Browser taxonomy and guidelines for designers. *IEEE Software*, 1995,12(2):21–32. [doi: 10.1109/52.368260]

- [120] Plaisant C, Grosjean J, Bederson BB. Spacetrree: Supporting exploration in large node link tree, design evolution and empirical evaluation. In: Proc. of the IEEE Symp. on Information Visualization (InfoVis 2002). Washington: IEEE Computer Society, 2002. 57-64. [doi: 10.1109/INFVIS.2002.1173148]

附中中文参考文献:

- [2] 李国杰,程学旗.大数据研究:未来科技及经济社会发展的重大战略领域.中国科学院院刊,2012,27(6):647-657.
 [3] 王元卓,靳小龙,程学旗.网络大数据:现状与展望.计算机学报,2013,36(6):1125-1138.
 [6] 孙大为,张广艳,郑纬民.大数据流式计算:关键技术及系统实例.软件学报,2014,25(4):839-862.
 [58] 史忠植.知识发现.北京:清华大学出版社,2002.
 [59] 董振东,董强,郝长伶.知网的理论发现.中文信息学报,2007,21(4):3-9.
 [60] 梅立军,周强,臧路,陈祖舜.知网与同义词词林的信息融合研究.中文信息学报,2005,19(1):63-70.
 [62] 于江生,俞士汶.中文概念词典的结构.中文信息学报,2002,16(4):12-20.
 [63] 许文艳,刘三阳.知识库系统的逻辑基础.计算机学报,2009,32(11):2123-2129.
 [64] 钟秀琴,刘忠,丁盘苹.基于混合推理的知识库的构建及其应用研究.计算机学报,2012,35(4):761-766
 [65] 陈立玮,冯岩松,赵东岩.基于弱监督学习的海量网络数据关系提取.计算机研究与发展,2013,50(9):1825-1835.
 [93] 程学旗,郭嘉丰,靳小龙.网络信息的检索与挖掘回顾.中文信息学报,2011,25(6):111-117.



程学旗(1971-),男,安徽望江人,博士,研究员,博士生导师,CCF 高级会员,主要研究领域为网络科学,网络与信息安全,互联网搜索与服务.

E-mail: cxq@ict.ac.cn



靳小龙(1976-),男,博士,副研究员,博士生导师,CCF 会员,主要研究领域为社会计算,网络性能建模与分析,多智能体系统.

E-mail: jinxiaolong@ict.ac.cn



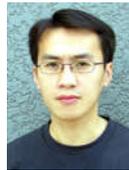
王元卓(1978-),男,博士,副研究员,CCF 高级会员,主要研究领域为社会计算,网络行为分析,信息安全.

E-mail: wangyuanzhuo@ict.ac.cn



郭嘉丰(1980-),男,博士,副研究员,CCF 会员,主要研究领域为信息检索,查询理解,数据挖掘.

E-mail: guojiafeng@ict.ac.cn



张铁赢(1982-),男,博士,助理研究员,CCF 会员,主要研究领域为计算机网络,分布式计算,网络安全.

E-mail: zhangtieying@ict.ac.cn



李国杰(1943-),男,博士,研究员,博士生导师,中国工程院院士,CCF 高级会员,主要研究领域为计算机体系结构,大数据.

E-mail: lig@ict.ac.cn