

大数据时代,统计学还有用吗?清华大学统计学研究中心主任刘军做客《文化讲坛》

大数据是“原油”,不是“汽油”

马林 葛亮亮

文化讲坛

在数据“爆炸”的时代,大数据常常被寄予厚望。到底,什么样的数据才算大数据,怎样才能用好大数据,传统统计学还有用吗?清华大学统计学研究中心主任刘军做客《文化讲坛》,与观众分享他的思考。



刘军

让大数据区别于数据的,是其海量、高增长率和多样性

什么是数据?数据(data)在拉丁语里是“已知”的意思,在英文中的一个解释是一组事实的集合,从中可以分析出结论。笼统地说,凡是用来承载信息记录下来的、能反映自然界和人类社会各种信息的,就可称之为数据。古人“结绳记事”,打了结的绳子就是数据。步入现代社会,信息的种类和数量越来越多,载体也越来越多样。数字是数据,文字是数据,图像、音频、视频等都是数据。

什么是大数据?量的增多,是人们对大数据的第一认识。随着科技发展,各个领域的数据量都在迅速增长。有研究说,近年来,数字数据的数量每3年多就会翻一番。

大数据区别于数据,还在于数据的多样性。正如高德纳咨询公司研究报告指出的,数据的爆炸性是三维的、立体的。所谓的三维,除了指数据总量爆炸外,还体现在数据速度的加快,以及数据的多样性;如数据的来源、种类不断增加。

从数据到大数据,不仅是量的积累,更是质的飞跃。海量的、不同来源、不同形式、包含不同信息的数据可以容易地被整合、分析,原本孤立的数据变得互相沟通。这使得人们对通过数据分析,能发现大数据时

代很难发现的新知识,创新新的价值。

通过数据来研究规律,发现规律,贯穿了人类社会发展史。人类科学发展史上的不少进步都和数据分析直接相关。例如现代医学流行病的开端,像1854年发生了大规模霍乱,很长时间没有有效遏制。一位医生用标点地图的方法分析了当地水和霍乱患者分布之间的关系,发现有一口水井周围,霍乱患者明显高发,借此找到了霍乱暴发的原因:一口被污染的水井。关闭这口水井之后,霍乱的发病率明显下降。这种方法,充分展示了数据的力量。

本质上说,许多科学活动都是数据挖掘,不是从预先设定的理论或者原理出发,通过演绎来研究问题,而是从数据本身出发通过归纳来总结规律。近现代以来,随着次数越来越多的实验,科学家越来越依赖于数据。有不少于年轻企业在实践中发现,顾客买了啤酒,买了啤酒之后,顺便买啤酒,沃尔玛发现了这一规律后,搭配促销啤酒、尿布,销量大幅增加。大数据时代,每个人都会“自动地”提供数据。我们的各种行为,如点击网页、使用手机、刷卡消费、观看视频、乘坐出租车、驾驶汽车,都会生成数据并留下记录,成为商家精准营销的数据库。

大数据是非竞争性资源,有助于商家精准营销

大数据时代,数据的重要作用更加凸显。许多国家都把大数据提升到国家战略的高度。政府利用大数据,可以及时发现潜在的事实,政府会

可以分析很多,如遗传信息、全体基因的表达量信息、蛋白质组信息、全基因组甲基化信息、表观遗传信息等。同时还有个人健康指标、病历、药物反应等数据。如果能把这些生物学上多维度的数据有机整合,就能够得到人完整的数据出来,从而实现精准医疗的目的。

大数据时代,审核数据的真实性也有了更有效的工具。大数据的特征之一是多样性,不同来源、不同维度的数据之间存在一定的相关性,可以进行交叉验证。例如,某地的工业产值虚增了一倍,但用电量却没有相应的增幅,这就是数据异常,很容易被发现识别出来。发现异常后,相关部门再进行复核,就能更针对性地防止、打击数据造假。

数据是一种资源,不像钱又跟煤、石油等物质性资源不一样。物质性资源不可再生,别人用少了,自己也就没了。数据却用不了,可再而再。数据可以重复使用,不断产生新的价值。大数据资源的使用是非竞争性,在共享的前提下,更能制造双赢。从另一个角度来说,数据如果不被融合、联系在一起,也不能称之为大数据。

大数据不能被直接拿来使用,统计分析依然是数据分析的灵魂

现在社会上有一种流行的说法,认为在大数据时代,“样本”不重要,人们也不需要抽样数据而是全数据。因而只需要简单数据一数就可以不再需要了。复杂的统计方法可以不再需要了。

在我看来,这种说法非常错误。首先,大数据信息并不简单。打个比方说,大数据是“原油”而不是“汽油”。原油是数据,但数据不等于原油。数据所包含的未知信息,是数据所不能提供的。大数据时代,统计学依然是数据分析的灵魂。比如加州大学伯克利分校的迈克尔·乔丹教授指出,没有系统的统计模型与方法为指向的“大数据分析”,就如同不利用工程科学的

我科学家揭示葡萄糖转运蛋白的工作机理

本报北京7月19日电 (记者赵永杰、赵利)继去年在世界上首次解析出人源葡萄糖转运蛋白 GLUT1 的三维晶体结构后,清华大学原子探针显微技术研究中心、清华大学化学系、清华大学医学院等单位的科研人员,进一步揭示了葡萄糖转运蛋白 GLUT1 的工作机理。

葡萄糖是地球上各生物最主要的能源物质,不仅为生长代谢提供能量,还参与其他生命活动。葡萄糖转运蛋白 GLUT1 是人体中 GLUT1 共有 4 种,目前研究比较清楚的是 GLUT1-3、4 两种。GLUT1-4, 过去 8 年,李宁研究员团队首次解析出人源 GLUT1 的三维结构,成为世界上第一个大型跨膜转运蛋白晶体结构。

据介绍, GLUT1 的结构处于向胞内开放的构象,只是 GLUT1 蛋白在行使转运功能过程中处于多种构象。为揭示转运过程,获得 GLUT1 处于其他不同转运状态下的结构信息至关重要,需了解 GLUT1 识别胞外葡萄糖并和葡萄糖 X-射线衍射它们与胞外葡萄糖的结合。李宁研究员团队精心设计实验,利用膜蛋白冷冻电镜技术,首次解析出人源 GLUT1-3 的晶体结构。其 GLUT1 与 3 种不同状态的葡萄糖转运蛋白晶体结构。其中 GLUT1 与 3 种不同状态的葡萄糖转运蛋白晶体结构处于向胞内开放的构象,分辨率高达 1.5 埃(1埃等于 0.1 纳米)。这是目前为止分辨率最高的转运蛋白晶体结构。这一超高分辨率清晰地揭示了 GLUT1 可以识别葡萄糖和 X-射线衍射它们。

另一方面,这个数据仍然是全数据,但仍然具有不确定性。入射的 X 射线强度并不完全代表学生的数学能力。假如所有同学重新参加一次考试,几乎每个同学都会有一个新的成绩。分别用这两组数据去分析,结论就可能发生变化。另一方面,事物在不断发展和变化,同学能力的成长和进步也能够体现在当前的成绩上。因此,研究如何从数据中把信息和规律提取出来,找出最优化的方法;也研究如何把数据当作的不确定性;也是大数据分析。

用火遗迹遗迹密集出现 “北京人”用火再添力证

本报北京7月19日电 (记者余俊杰)记者近日从北京市文物局获悉,北京人遗址第一地点(猿人洞)2009年—2014年的清理发掘过程中,古人类用火遗迹、遗迹密集出现。

据考古、火源、地质、地化、石器等古人类学火源遗迹、遗迹的密集出现,“北京人”用火再添力证。此外,还发现用火遗迹、遗迹密集出现。此外,还发现用火遗迹、遗迹密集出现。

据考古、火源、地质、地化、石器等古人类学火源遗迹、遗迹的密集出现,“北京人”用火再添力证。此外,还发现用火遗迹、遗迹密集出现。此外,还发现用火遗迹、遗迹密集出现。

北京市本科一批录取工作启动 高分落档现象大为减少 考生志愿梯度较为合理

本报北京7月19日电 (赵利)记者在北京市高招办启动本科一批录取工作了解到,北京市2015年本科一批录取工作已于7月17日启动,将于7月22日结束。此前,北京市已完成录取工作,本科提前批次和特殊类型录取工作。

据悉,今年共有167所院校在京参加一本招生,招生计划14861人。其中,文史类计划2748人,理工类计划12113人。目前已投档15610人,其中文科2845人,理科12765人。

今年北京市首次实施大文科招生,减少了考生志愿填报风险,增加了录取机会。改革取得预期成效。截至目前,北京市志愿填报情况良好。本科一批录取工作已于7月17日启动,将于7月22日结束。此前,北京市已完成录取工作,本科提前批次和特殊类型录取工作。

“木绘拼花”手工技艺 让故宫藏画走进大众

本报北京7月19日电 (记者李晔)19日,故宫博物院下属的北京故宫文化创意有限公司与北京百纳工匠手工艺股份有限公司签订合作协议,双方将在故宫博物院指导下,运用“木绘拼花”手工技艺,共同研发出品故宫藏画木制品工艺,使故宫藏画数百年来之皇家藏画在新工艺、新媒体上焕发



19日,在新疆博物馆的展厅内,一位老人正带着孩子观看油画作品。

油画里的新疆

本报北京7月19日电 (记者李晔)19日,由新疆维吾尔自治区主办、2015年青年油画创作营活动主办方在北京大学开幕,标志着这项将有全国1万多名中学生参加

暑期电影票房 连破多项纪录

本报北京7月19日电 (记者刘颖)近日,中国电影票房连破多项纪录——10日至19日突破11亿元,刷新周末票房最高纪录;18日全国票房人次突破960万次,刷新单日票房最高纪录;18日全国票房人次突破960万次,刷新单日票房最高纪录;18日全国票房人次突破960万次,刷新单日票房最高纪录。

反法西斯电影研讨会举行

本报北京7月19日电 (刘颖)由中国人民大学文学院、当代电影杂志社联合主办的“反法西斯电影研讨会”18日在北京拉开帷幕。为纪念第二次世界大战胜利70周年,研讨会围绕反法西斯电影展开。研讨会围绕反法西斯电影展开。研讨会围绕反法西斯电影展开。

中国微电影工作委员会成立

本报北京7月19日电 (记者李晔)由中国电影家协会微电影工作委员会在北京成立。微电影工作委员会成立揭幕仪式日前在北京举行。微电影工作委员会成立揭幕仪式日前在北京举行。微电影工作委员会成立揭幕仪式日前在北京举行。

万名中学生入驻高校科学营

本报北京7月19日电 (记者李晔)2015年青年油画创作营活动主办方在北京大学开幕,标志着这项将有全国1万多名中学生参加