

计算社会科学新进展 ——方法与应用

孟天广 副教授

清华大学计算社会科学平台 执行主任

清华大学社会科学学院 院长助理

2019年12月12日

内容提要



大数据+社会科学



计算社会科学：方法论



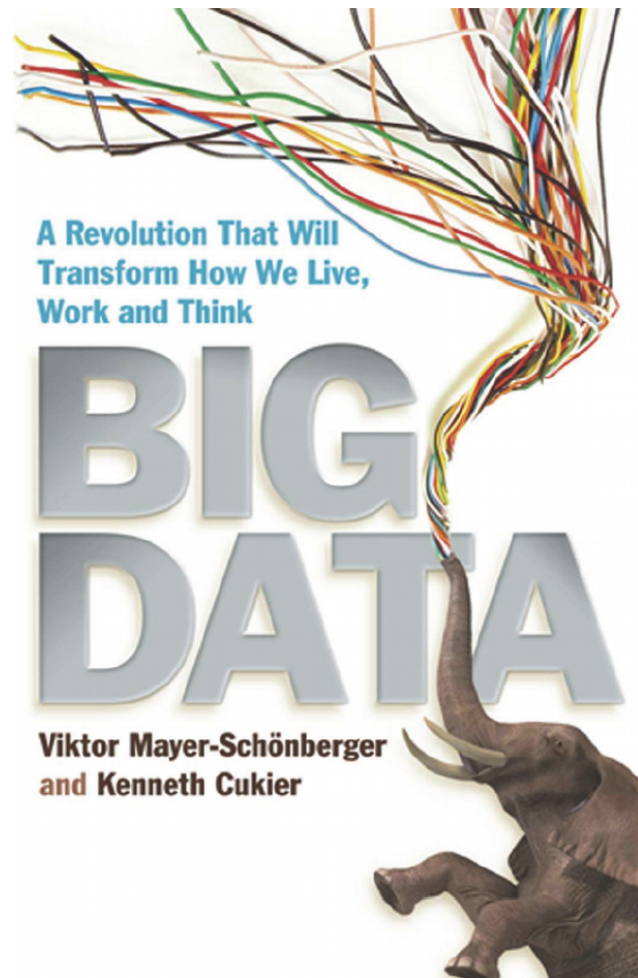
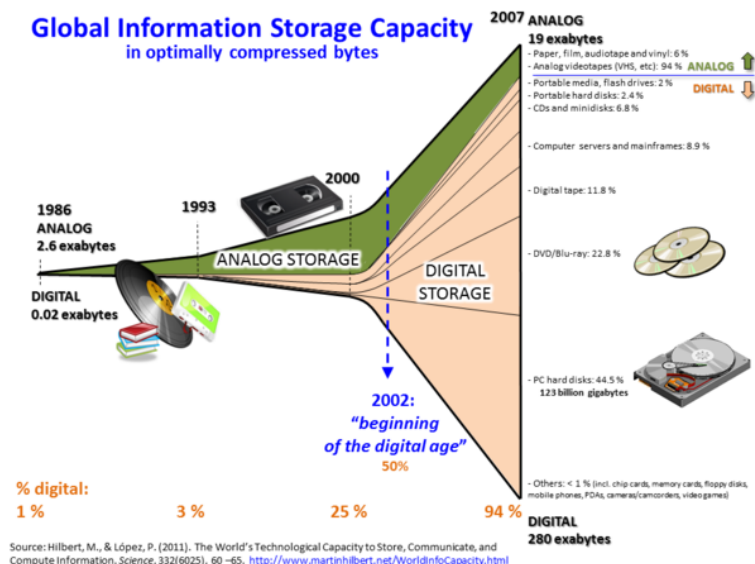
计算社会科学：方法与议题



计算社会科学：前沿应用

大数据时代

- 《科学》2008年提出“大数据”来讨论新信息时代（PB时代）的科学研究；
- 2012年，《纽约时报》刊文宣告“大数据时代已经到来”；



大数据时代

- 大数据是指需要新处理模式才能确保更强的决策力、洞察力和流程优化力的海量、高速增长和多样化的信息财富（Gartner）
- 大数据是以容量大、类型多、存取速度快、应用价值高为主要特征的数据集合，正快速发展为对数量巨大、来源分散、格式多样的数据进行采集、存储和关联分析，从中发现新知识、创造新价值、提升新能力的新一代信息技术和服务业态。（国务院）
- 数据已成为国家基础性战略资源（十三五）

大数据的特征

- Volume(体量大):
 - How much data is really relevant to the problem solution? Cost of processing?
 - *So, can you really afford to store and process all that data?*
- Velocity(增速快):
 - Much data coming in at high speed
 - Need for streaming versus block approach to data analysis
 - *So, how to analyze data in-flight and combine with data at-rest*
- Variety(类型多):
 - A small fraction is structured formats, Relational, XML, etc.
 - A fair amount is semi-structured, as web logs, etc.
 - The rest of the data is unstructured text, photographs, etc.
 - *So, no single data model can currently handle the diversity*
- Veracity(真实性):
 - Accuracy, Precision, Reliability, Integrity
 - *So, what is it that you don't know you don't know about the data?*
- Value(价值高):
 - How much value is created for each unit of data (whatever it is)?
 - *So, what is the contribution of subsets of the data to the problem solution?*

大数据

What is Big Data

- Data
- Analytics
- ✓ Big data is not about Data!
- Industry
- Solution

Why is Big Deal

- Government
- Private Sector
- Science
- Social Science Revolution

大数据+社会科学

- 大数据时代的八大机遇
 - ✓ 海量非结构化数据（信息）
 - ✓ “全量数据”而不是“样本数据”
 - ✓ 丰富、高效的方法工具箱
 - ✓ 机器学习与预测性分析
 - ✓ 强时效性数据
 - ✓ 社会科学知识平民化普及
 - ✓ 良好的组织环境+充分的社会需求

大数据+社会科学

- 大数据+社会科学
 - ✓ 数据驱动
 - ✓ 应用（问题解决）导向
 - ✓ 中观/微观问题
 - ✓ 预测性目标
- 从海量数据中利用机器学习抽离出有价值的信息
 - ✓ 积累海量数据
 - ✓ 利用统计和数学知识+模式识别技术
 - ✓ 发现有意义的新关系、新模式或新趋势

计算社会科学的清华探索

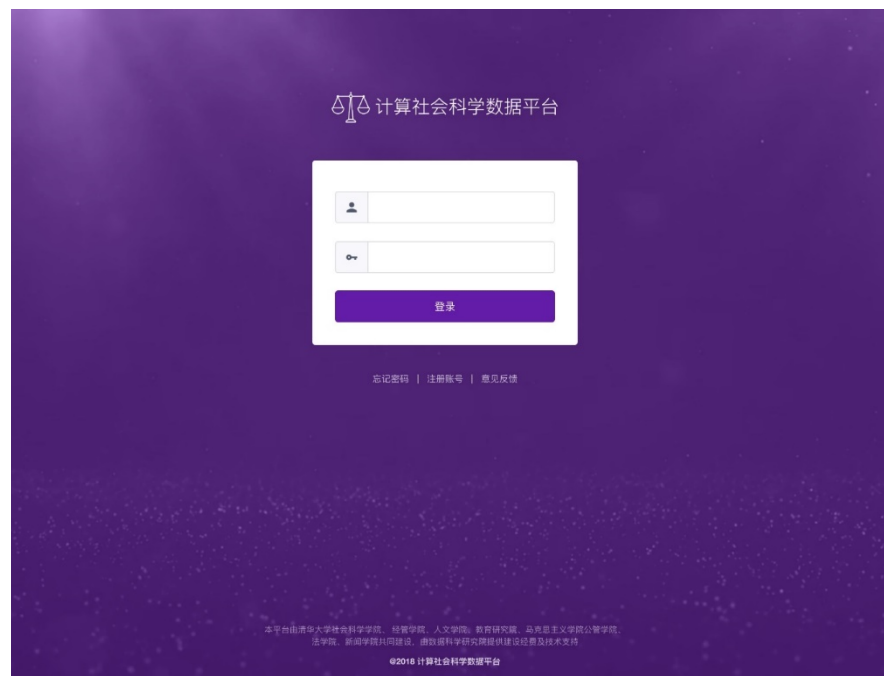
- 清华大学计算社会科学平台
- 平台定位和发展目标
- ✓ 国内首家计算社会科学领域的创新研究机构
- ✓ 营造“创新、共享、开放、合作”的科研环境
- ✓ 立足清华社科、计算科学和数据科学领域的特色优势和交叉研究基础
- ✓ 将社科议题、海量数据、大数据方法相结合，促进跨学科融合，开展创新性研究
- ✓ 集前沿研究、人才培养、科研服务和智库资政于一体
- ✓ 为中国特色计算社会科学的学科建设、清华高水平社科创新成果、促进重大经济社会问题解决提供理论和应用知识

计算社会科学的清华探索

- 清华大学计算社会科学平台

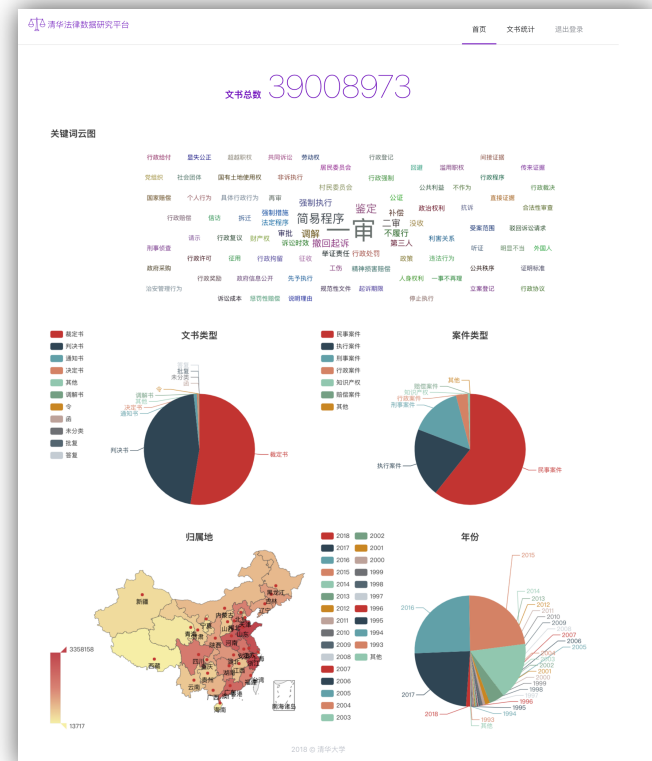
- 法律数据科研条件平台

- 该平台经过一年多的建设，已汇集了8000余万份全国范围内公开的社科内容数据，形成可持续更新的数据库，具备全文检索、分类检索、结构化分析、统计分析、可视化报表等在线服务功能。
- 平台链接：<http://tcd.ids.tsinghua.edu.cn/>



计算社会科学的清华探索

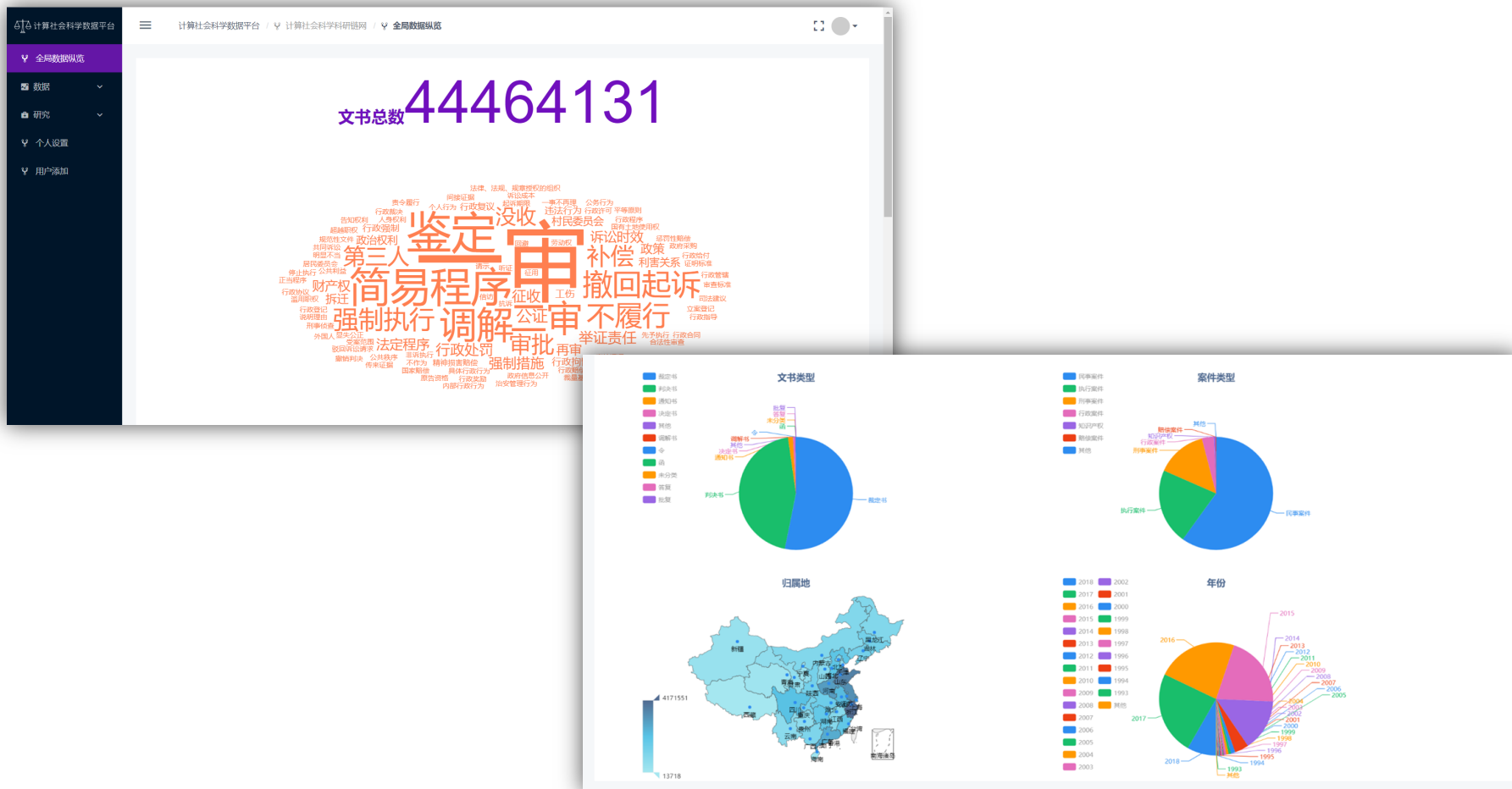
- 2017年11月至今，社科学院与法学院、城管学院、经管学院、人文学院、马克思学院、数据院、计算机系等院系的老师进行过密切的沟通和调研访谈，形成会议记录，明确本平台的建设目的和功能符合师生在社科大数据科研工作需求。研究并确定基于清华邮箱的账号管理模式属于有效的清华身份认证和账号管理模式，符合要求也方便师生使用。



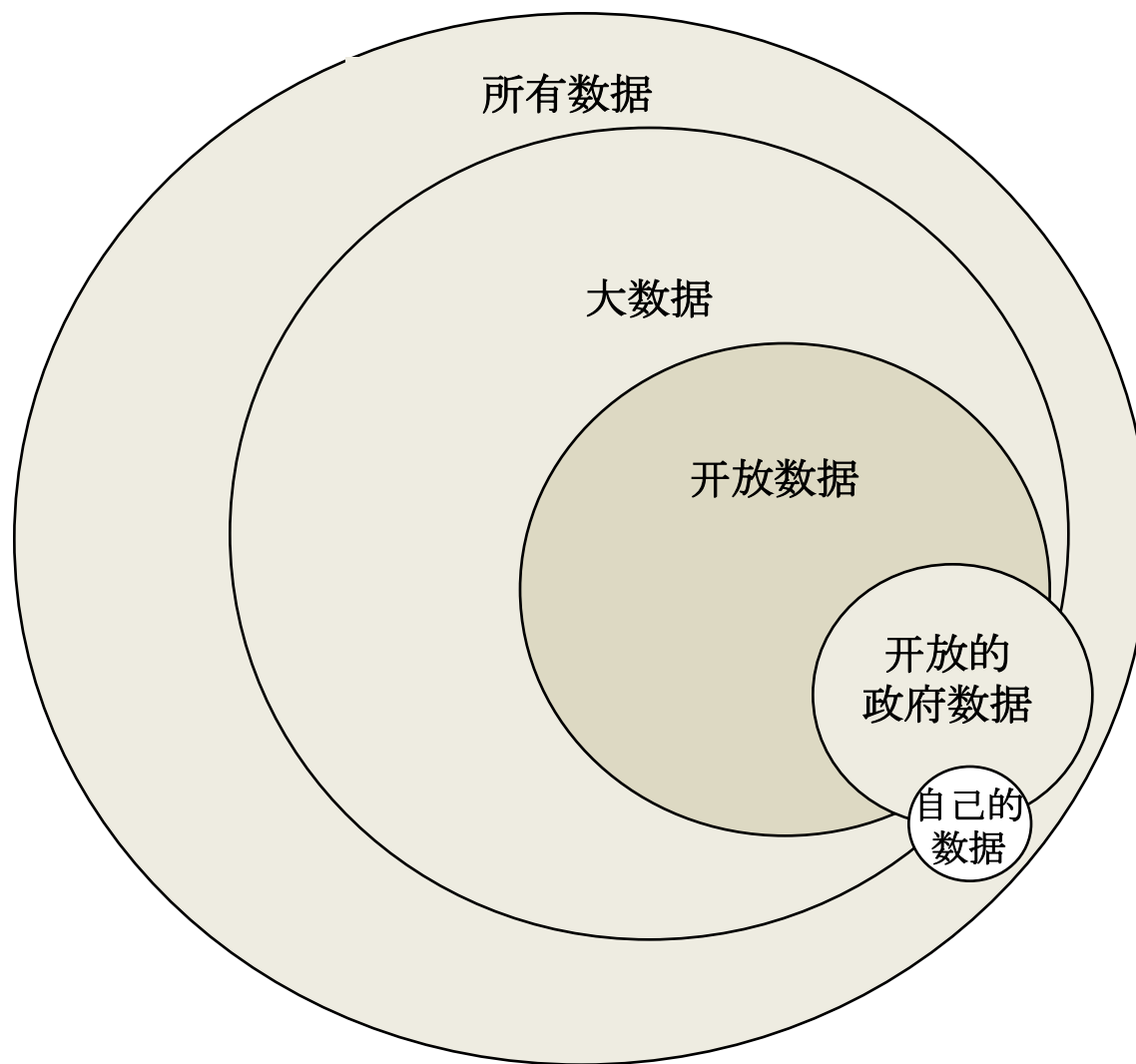


全局数据纵览

- 对系统内文书总数有充分了解，从可视化图形中纵览各类型文书数据占比、案件类型占比、归属地数据占比及统计年间各个年份数据的占比统计数据，以使用户对系统数据有清楚的全局观。

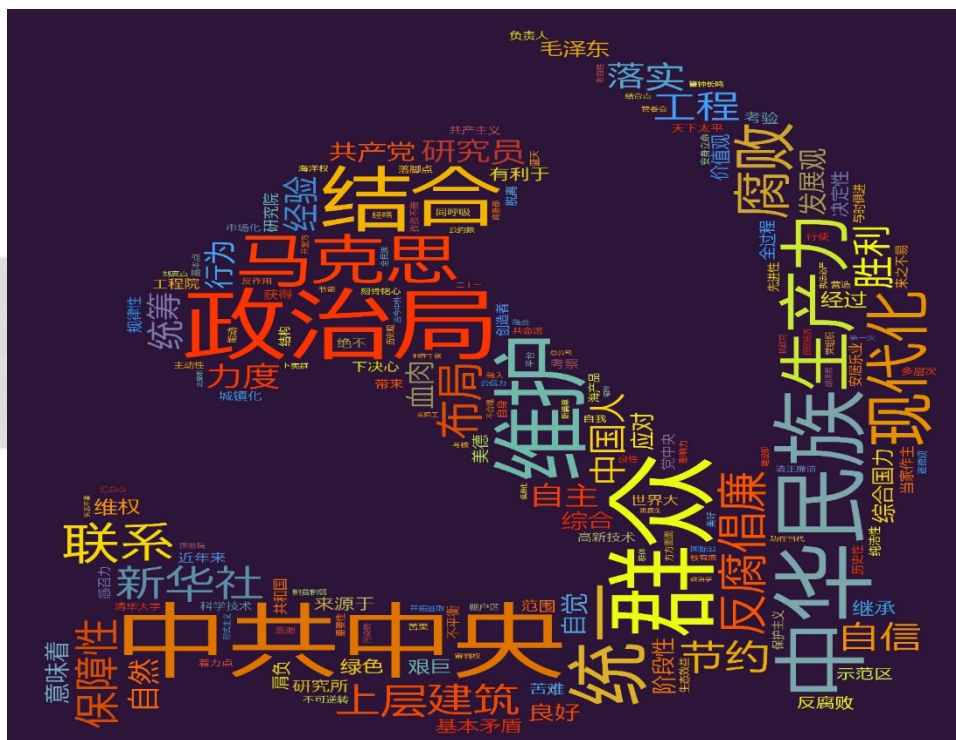
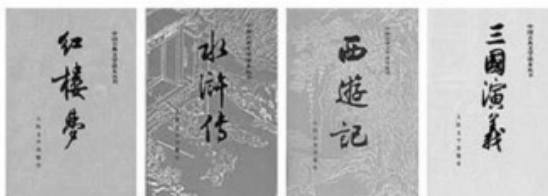
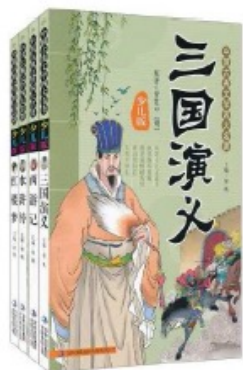


大数据的范畴



文本资料

Text data



多媒体数据（音频、视频、图片）

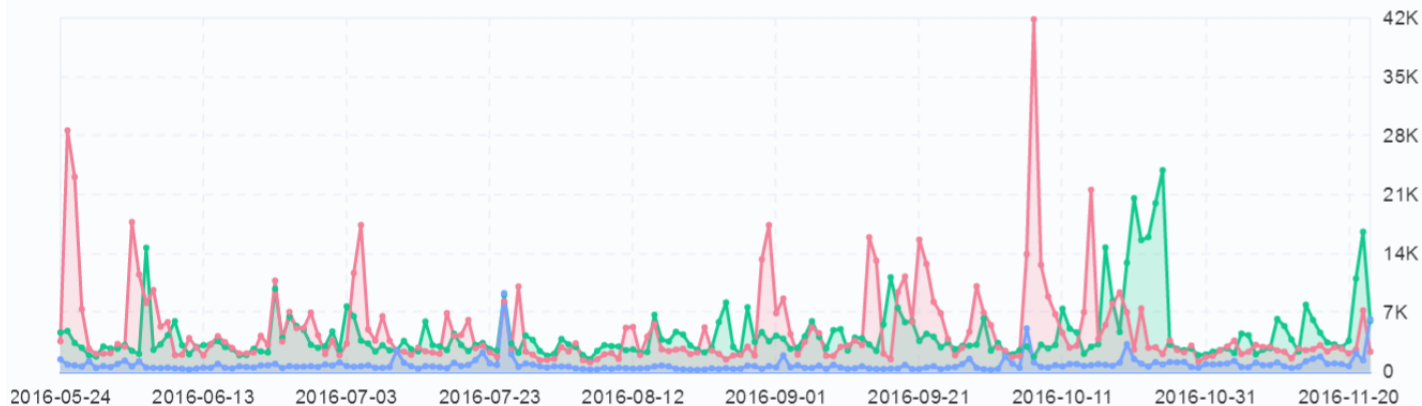
Multimedia data



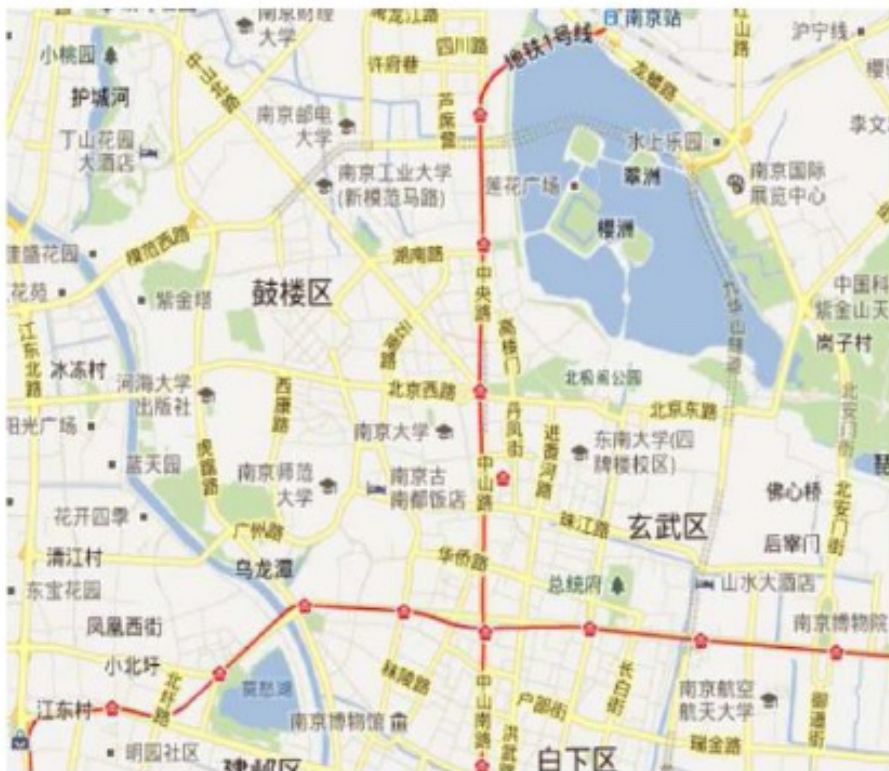
门户网站与新媒体



-O- 中国人民大学 -O- 清华大学 -O- 北京大学

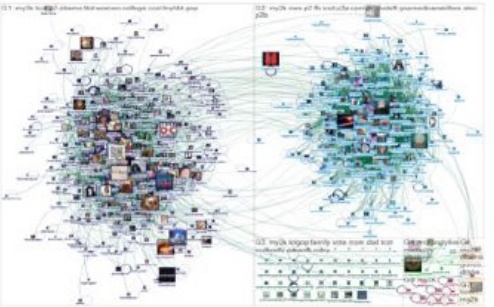
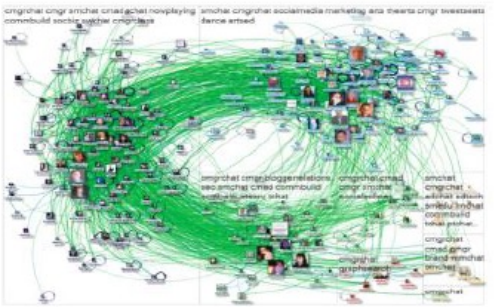
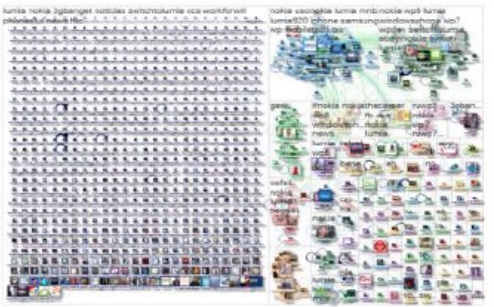
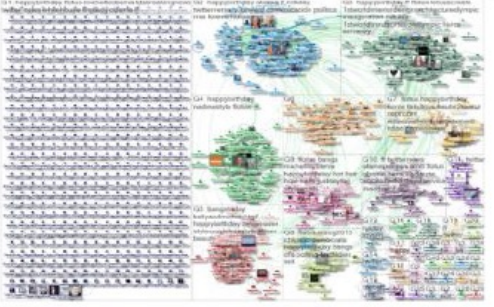
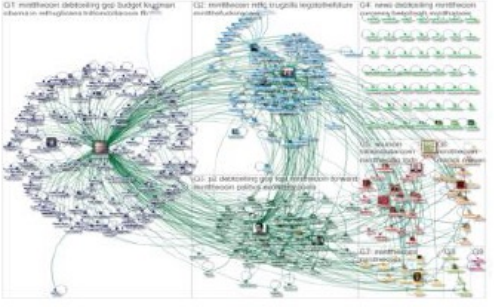



空间数据



网络关系数据

6 kinds of Twitter social media networks

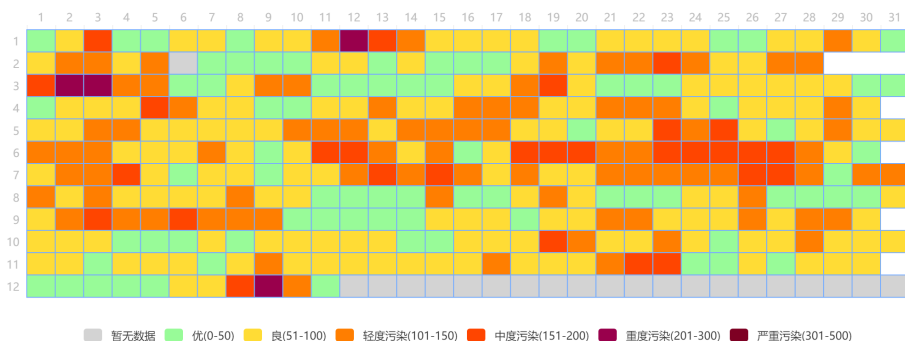
		
<p>Polarized: two dense clusters with little interconnection</p>	<p>In-group: few disconnected isolates, many connections</p>	<p>Brand/Public Topic: many disconnected isolates, some small groups</p>
		
<p>Bazaar: many medium sized groups, some isolates</p>	<p>Broadcast: a hub which is retweeted by many disconnected users</p>	<p>Support: a hub which replies to many disconnected users</p>

行政数据

AQI指数分布

2017 2018 2019

2019年AQI指数全年分布图



计算社会科学：方法论

- 计算社会科学：大数据+社会科学
- ✓ 图灵奖得主J. Gray（2010）：大数据时代将形成数据密集型科学研究“第四范式”。大数据时代的科学研究将不再需要模型和假设，而是利用超级计算能力直接分析海量数据发现相关关系即可获得新知识；
- ✓ 2009年，哈佛大学David Lazer等15位美国学者在《Science》上联合发表了一篇具有里程碑意义的文章“Computational Social Science”；
- ✓ 2014年，哈佛大学Gary King认为大数据方法将终结传统的定量、定性方法分野。

Big Data and Social Analytics certificate course

HARVARD UNIVERSITY HARVARD.EDU

IQSS
The Institute for Quantitative Social Science

CONTACT US OPPORTUNITIES

The IQSS Story

About IQSS ▾ Activities ▾ People News & Events

Helping social scientists understand and solve society's greatest challenges through activities in...

PRINCETON UNIVERSITY

LAZER LAB

HOME PEOPLE PROJECTS PUBLICATIONS DATA EVENTS

TOP NEWS

VOLUNTEER SCIENCE
Volunteer Science!

Volunteer Science offers people around the world to become involved in scientific projects on everything from... more

CENTER FOR
STATISTICS AND
MACHINE LEARNING

WZB
Berlin Social Science Center

Home | Publications | WZB-Mitteilungen

Publications

Big Data in the social sciences



Alexandros Tokhi and Christian Rauh

Big Data has become an ubiquitous buzzword. Social networks, smartphones, and various online applications and websites constantly produce and provide information on an unprecedented scale and level of detail. Data storage is cheap and ever-improving analytical

PENNSYLVANIA STATE UNIVERSITY

BDSS
Big Data Social Science
An Integrative Education and Research Program in Social Data Analytics

Home IGERT SoDA Research News Events People Connect

BIG DATA SOCIAL SCIENCE
Integrative Graduate Education and Research Traineeship at Pennsylvania State University

NEWS More ▾

BDSS IGERT Speaker & Event Series - "Speed Dating and Matchmaking"

BDSS-IGERT Welcomes Its Fifth

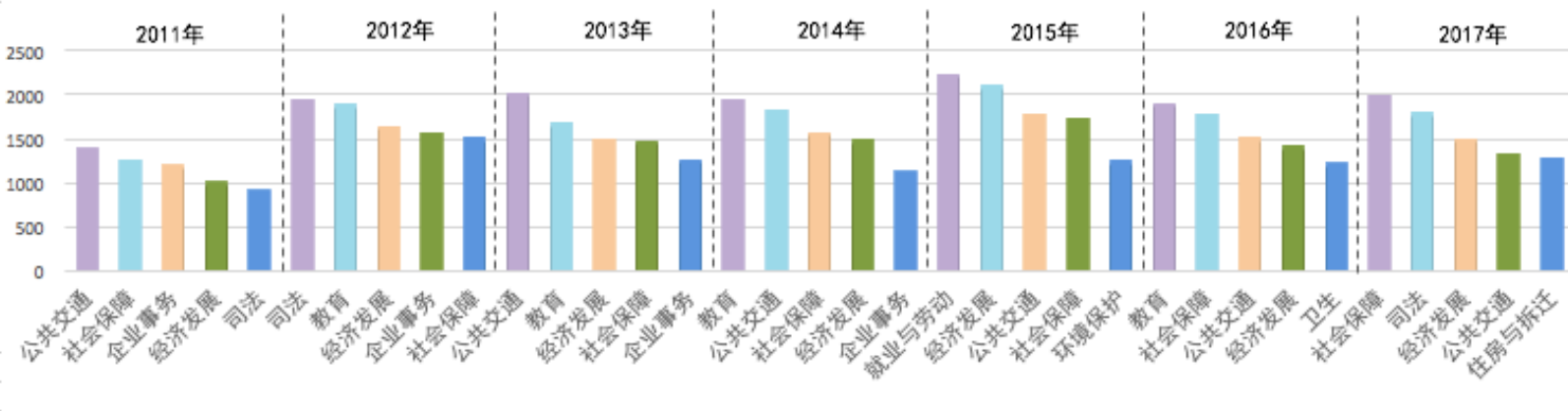
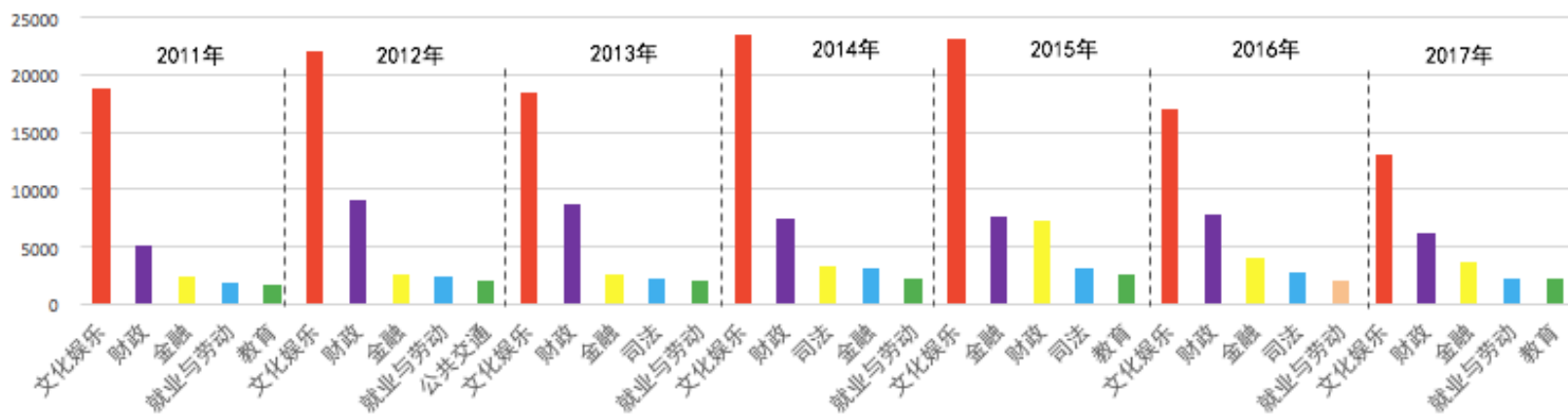
计算社会科学：研究方法

- 网络爬虫
 - 对搜索引擎搜索记录的分析
 - 自动文本分析
 - 视频/图片分析
 - 社会网络分析
 - 空间/时间分析
 - 可视化
-
- 机器学习
 - 自然语言过程
 - 统计分析

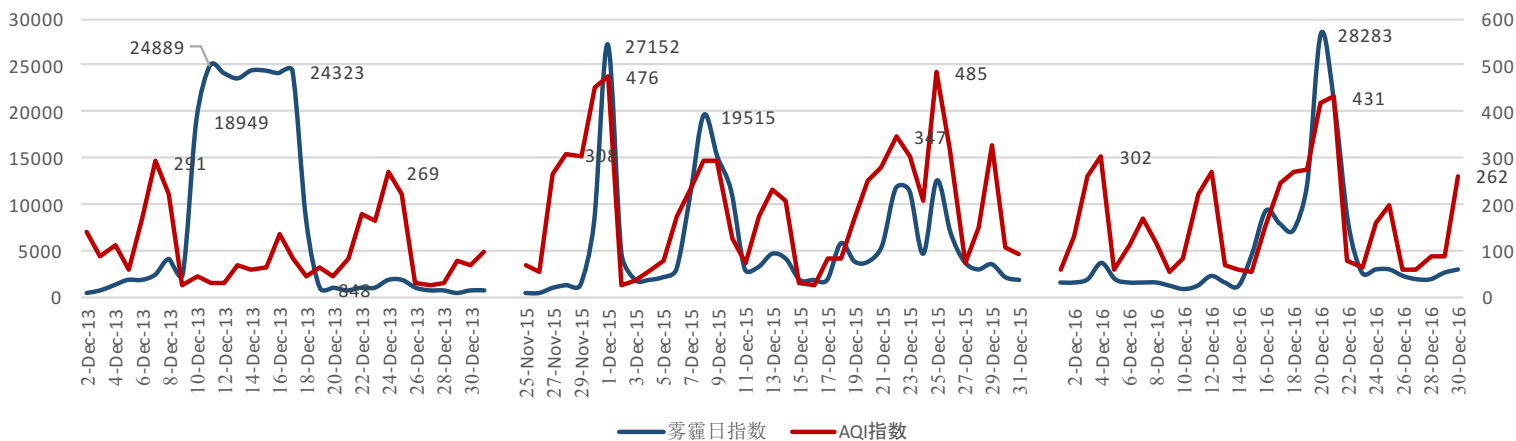
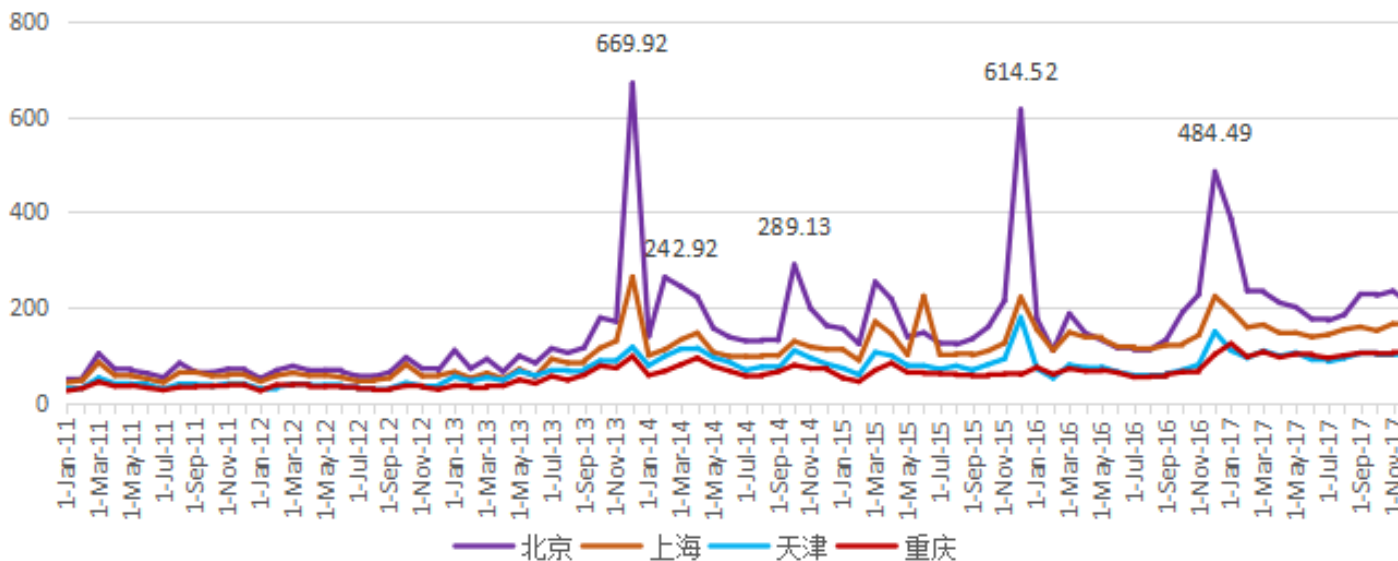
社科研究中的大数据：方法功能

- 作为研究方法的大数据分析
 - 数据采集与管理
 - 测量工具
 - 分类与聚类
 - 关联分析
 - 因果推论（回归分析）
 - 信息呈现（可视化）

搜索指数：Google Trends、百度指数



搜索指数：测量公共关注度



The Local Leader Message Board

The screenshot shows the homepage of the Local Leader Message Board (地方领导留言板) on the People's Net (人民网). The page features a red header with the site's logo and navigation links. A prominent quote from Xi Jinping is displayed in a yellow banner. Below this, there are sections for 'Latest News' (最新动态) and a 'Quick Browse' (快速浏览) area with dropdown menus for selecting provinces and a 'Quick Browse' button. The 'Latest News' section lists several news items related to food subsidies and government transparency. On the right side, there are three QR codes for mobile access and a 'Close' button at the bottom right.

当前位置: 人民网 >> 地方领导留言板

习近平: 网民来自老百姓, 老百姓上了网, 民意也就上了网。群众在哪儿, 我们的领导干部就要到哪儿去。各级党政机关和领导干部要学会通过网络走群众路线, 经常上网看看, 了解群众所思所愿, 收集好想法好建议, 积极回应网民关切, 解疑释惑。

最新动态

安徽网友咨询粮食补贴问题 官方解答

天津网友举报黑幼儿园泛滥 官方进行整治

- 网友给江西省委书记留言获官方回复 共计11条
- 网友给四川省委书记、省长留言获回复 共计35条

网友给河南省委书记、省长留言获答复 共计49条

- 网友给辽宁省省长留言获回复 共计17条
- 网友给辽宁省委书记留言获回复 共计60条

国办印发《通知》: 加强政务公开 做好舆情回应

请选择 > 请选择 > 请选择 快速浏览

北京 | 天津 | 河北 | 山西 | 内蒙古 | 辽宁 | 吉林 | 黑龙江 | 上海 | 江苏 | 浙江 | 安徽 | 福建 | 江西 | 山东 | 河南 | 湖北 | 湖南 | 广东 | 广西 | 海南 | 重庆 | 四川 | 贵州 | 云南 | 西藏 | 陕西 | 甘肃 | 青海 | 宁夏 | 新疆 | 香港 | 澳门 | 台湾 | [点击地区进入留言板留言]

昨日留言 264条 昨日回复 18条

The Interface

河南省 全省历史留言总量: 169463条 历史回复总量: 140792条

Provincial leader message boards

输入搜索内容 关键词 搜索 排行榜

河南省委书记谢伏瞻 Xie Fuzhan, Prov Sec of Henan

查看简历 我要留言

年度总留言量: 9845条 年度公开回复量: 8547条
Annual message: 9845; Total open reply: 8547

Button for "I want to leave a message"

河南省省长陈润儿 Chen Run'er, Governor of Henan

查看简历 我要留言

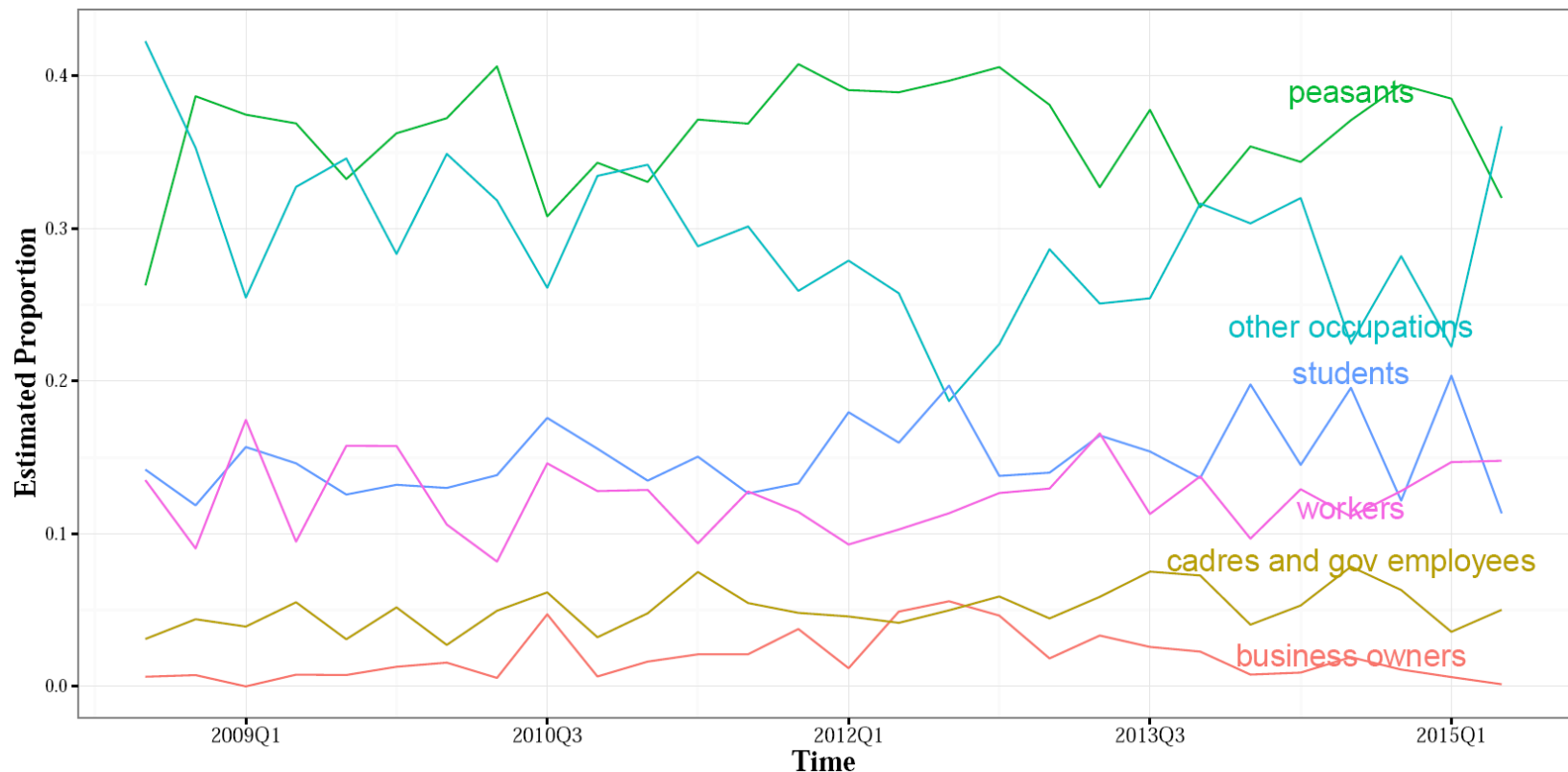
年度总留言量: 3084条 年度公开回复量: 2644条

City leader message boards

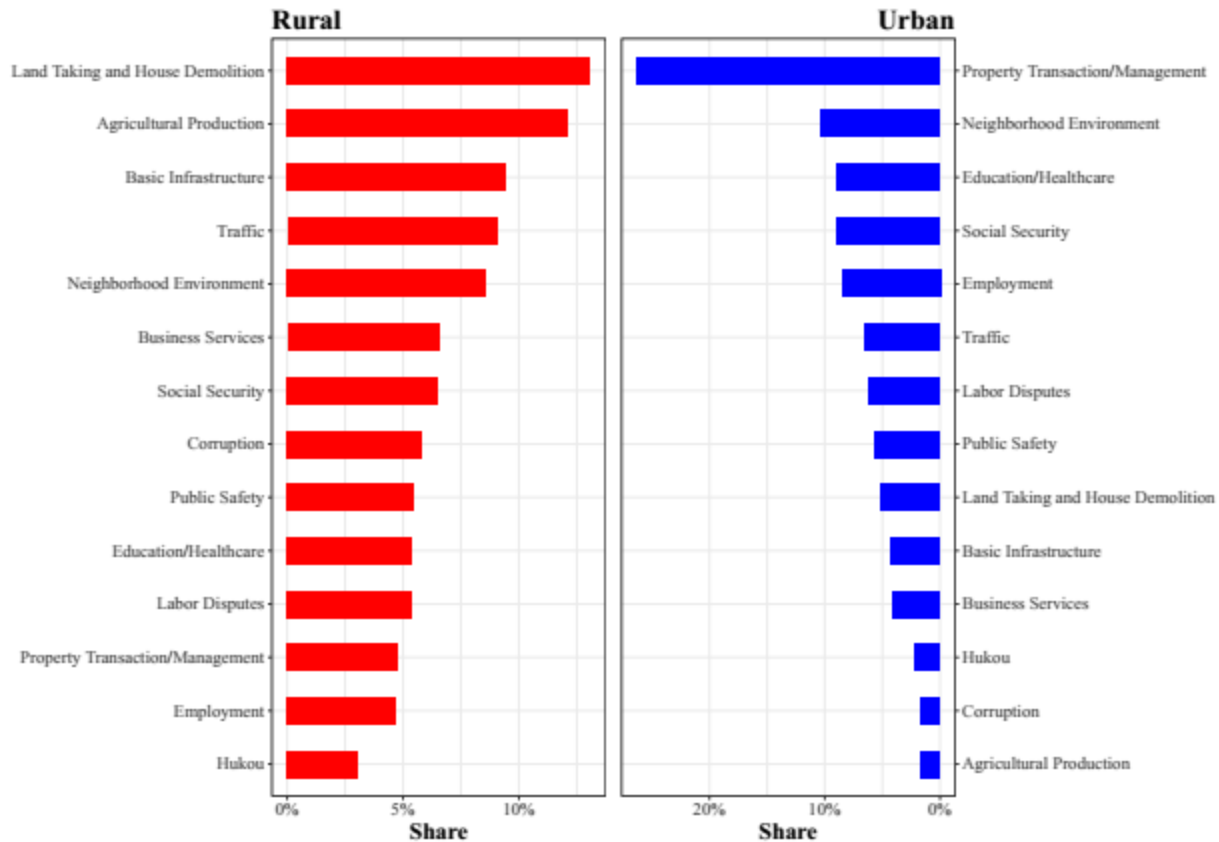
<p>郑州市 Zhengzhou City</p> <p>年度总留言量: 5927条 Annual message: 5927 年度公开回复量: 5267条 Open reply: 5267</p>	<p>开封市</p> <p>年度总留言量: 1433条 年度公开回复量: 799条</p>	<p>洛阳市</p> <p>年度总留言量: 1574条 年度公开回复量: 1430条</p>	<p>平顶山市</p> <p>年度总留言量: 5030条 年度公开回复量: 3808条</p>
<p>安阳市</p> <p>年度总留言量: 913条 年度公开回复量: 690条</p>	<p>鹤壁市</p> <p>年度总留言量: 198条 年度公开回复量: 51条</p>	<p>新乡市</p> <p>年度总留言量: 2357条 年度公开回复量: 1788条</p>	<p>焦作市</p> <p>年度总留言量: 485条 年度公开回复量: 317条</p>
<p>濮阳市</p> <p>年度总留言量: 4564条 年度公开回复量: 4018条</p>	<p>许昌市</p> <p>年度总留言量: 402条 年度公开回复量: 0条</p>	<p>漯河市</p> <p>年度总留言量: 270条 年度公开回复量: 12条</p>	<p>三门峡市</p> <p>年度总留言量: 669条 年度公开回复量: 363条</p>

Who Participate Online?

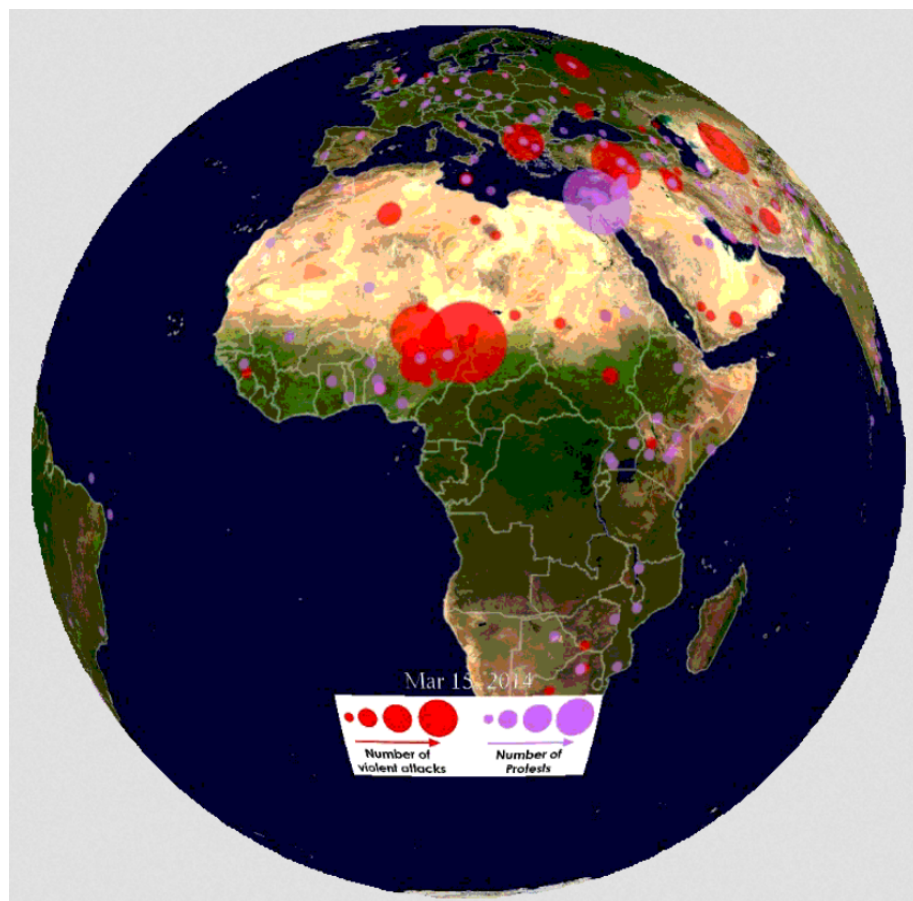
- Supervised Learning with readme (Hopkins and King, 2010)



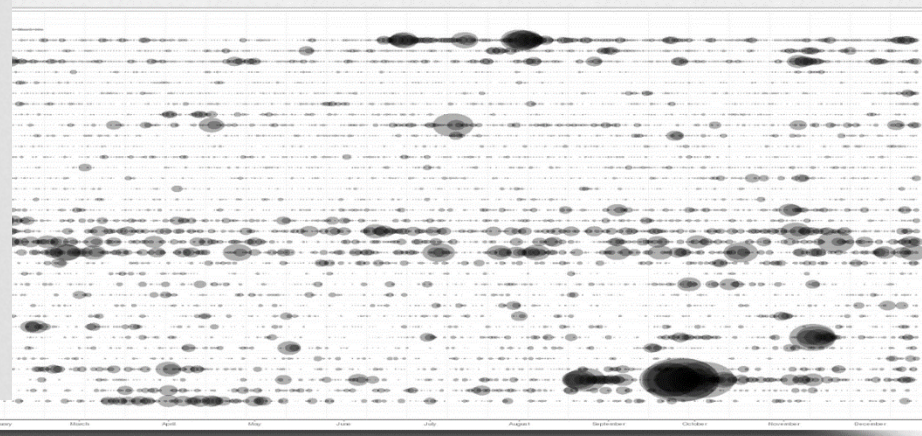
Participate for What?



时间序列分析：全球事件、语调与语言数据库（GDELT）

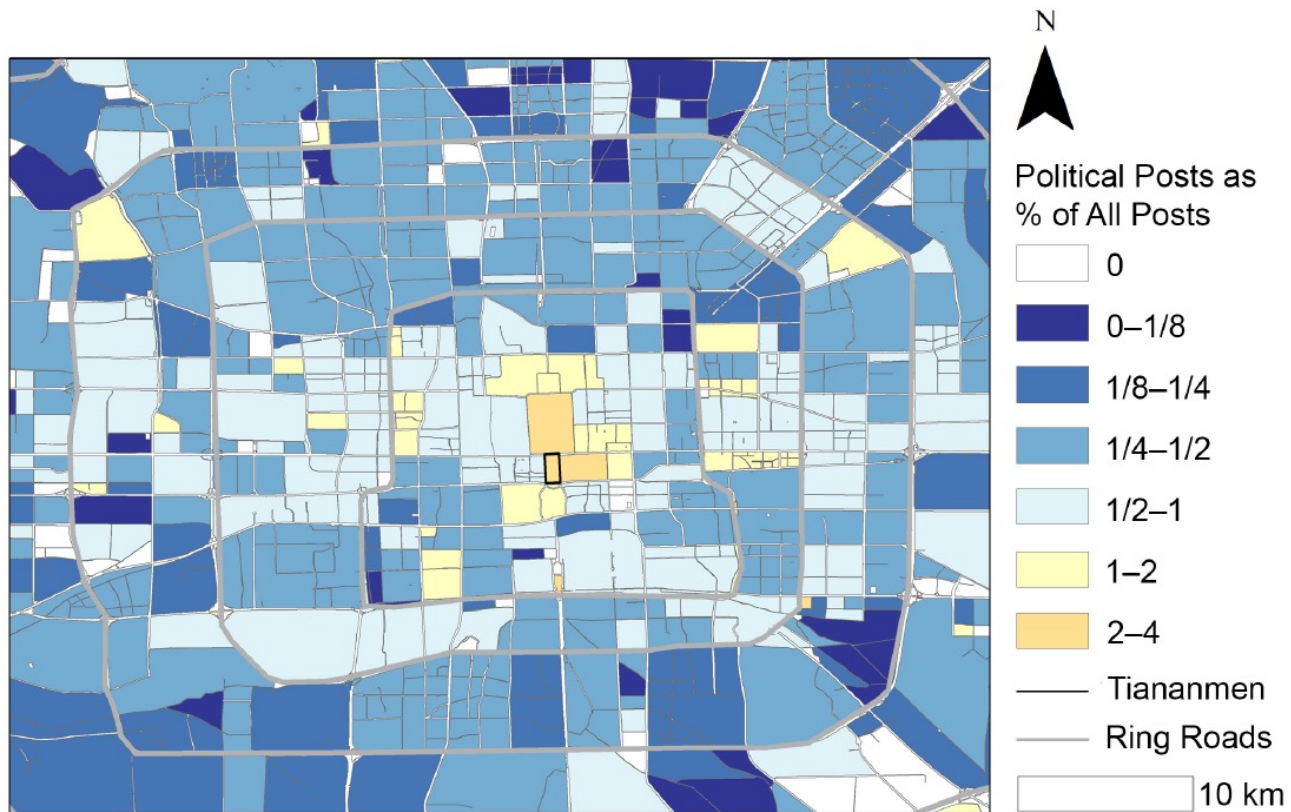


Timeline Visualizer



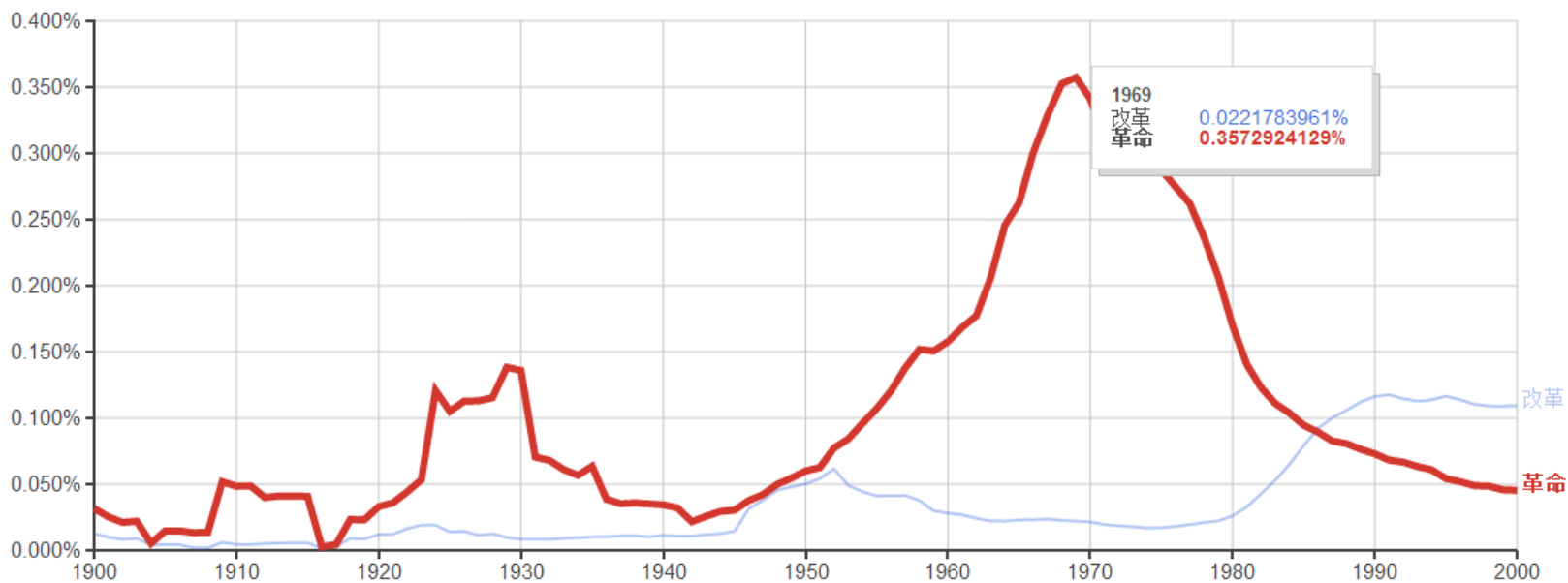
空间分析：空间与政治

Figure 2. Spatial Distribution of Political Posts, Beijing, 1 July 2014–15 June 2015



开源文本分析：Google Ngram

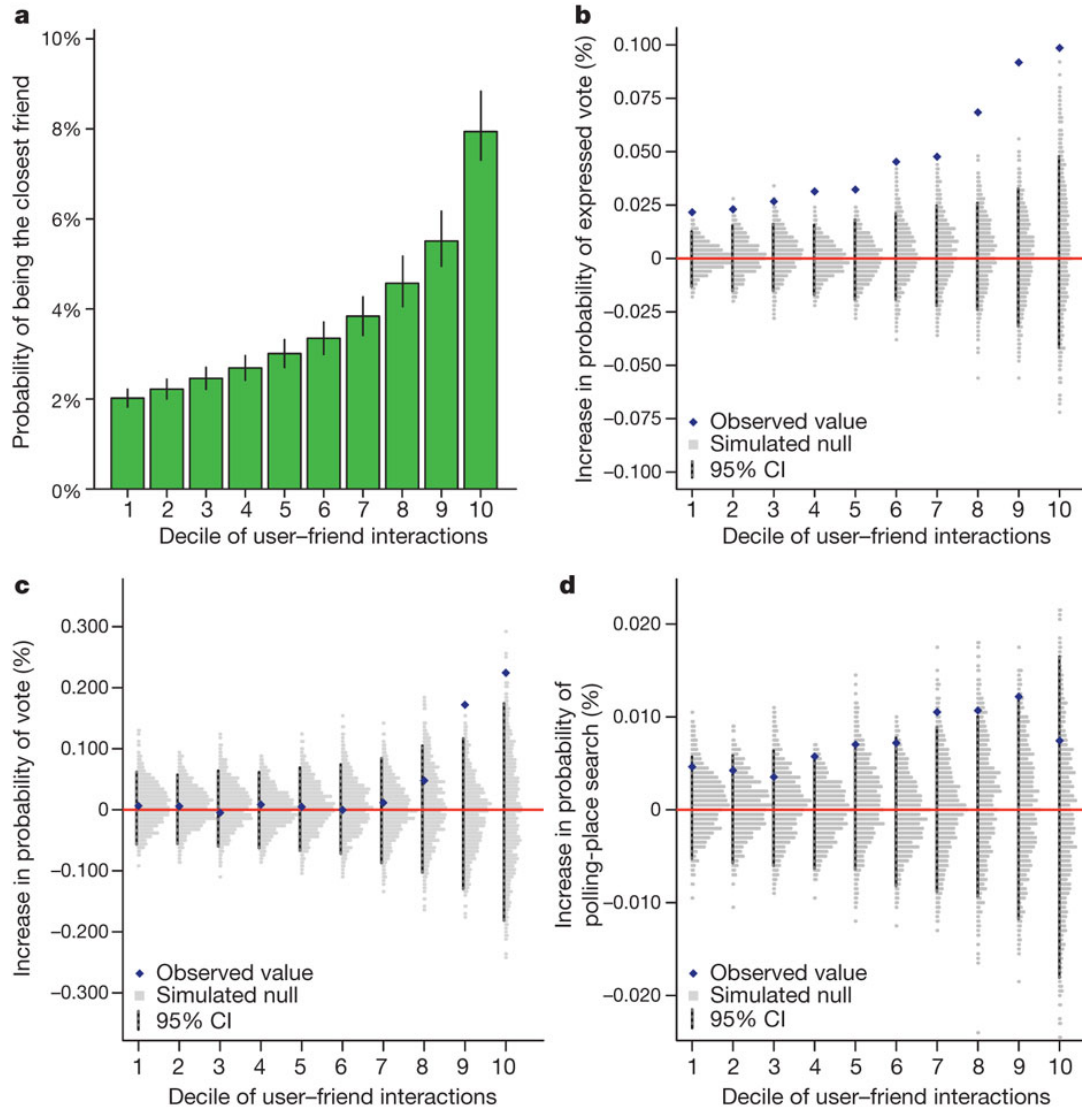
between 1900 and 2000 from the corpus Chinese (simplified) with smoothing of 3 Search lots of books



Robert M. Bond, et al, "A 61-million-person experiment in social influence and political mobilization", *Nature*, Vol.489, No.7415, 2012, pp.295-298.

- Bond等（2012）比较了线上社交网络和面对面社交网络影响政治行为的路径。
- 2010年美国国会大选时对6100万Facebook用户实施发送政治动员消息的随机控制实验；
- 政治动员消息直接影响网民的政治自我表达、信息搜寻和现实投票行为；
- 政治动员消息不仅影响了接受者，还影响了接受者的网友、网友的网友，而这种社会传递效应对投票行为的影响要强于直接动员效应；
- 信息传播更容易发生在具有见面关系的关系密切的朋友中。表明强关系有助于社交网络中对于在线和现实生活中的政治动员。

The effect of mobilization treatment that a friend received on a user's behaviour.



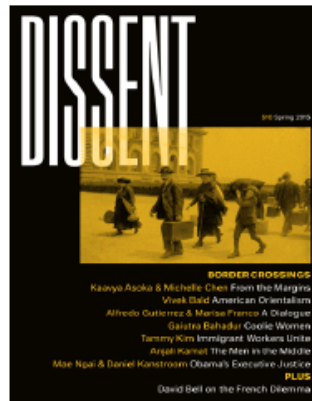
nature

King, G. , Schneer, B. , & White, A. . (2017). How the news media activate public expression and influence national agendas. *Science*, 358(6364), 776-780.

- King等（2017）利用48个社交媒体开展（五年期）田野实验
 - ✓在真实媒体环境中设计并随机化分配媒体资讯（报道）
 - ✓识别媒体报道对个体公共意见（政治知识）的效应（Individual Effect）
 - ✓识别媒体报道对国家政策议程的效应（Collective Effect）

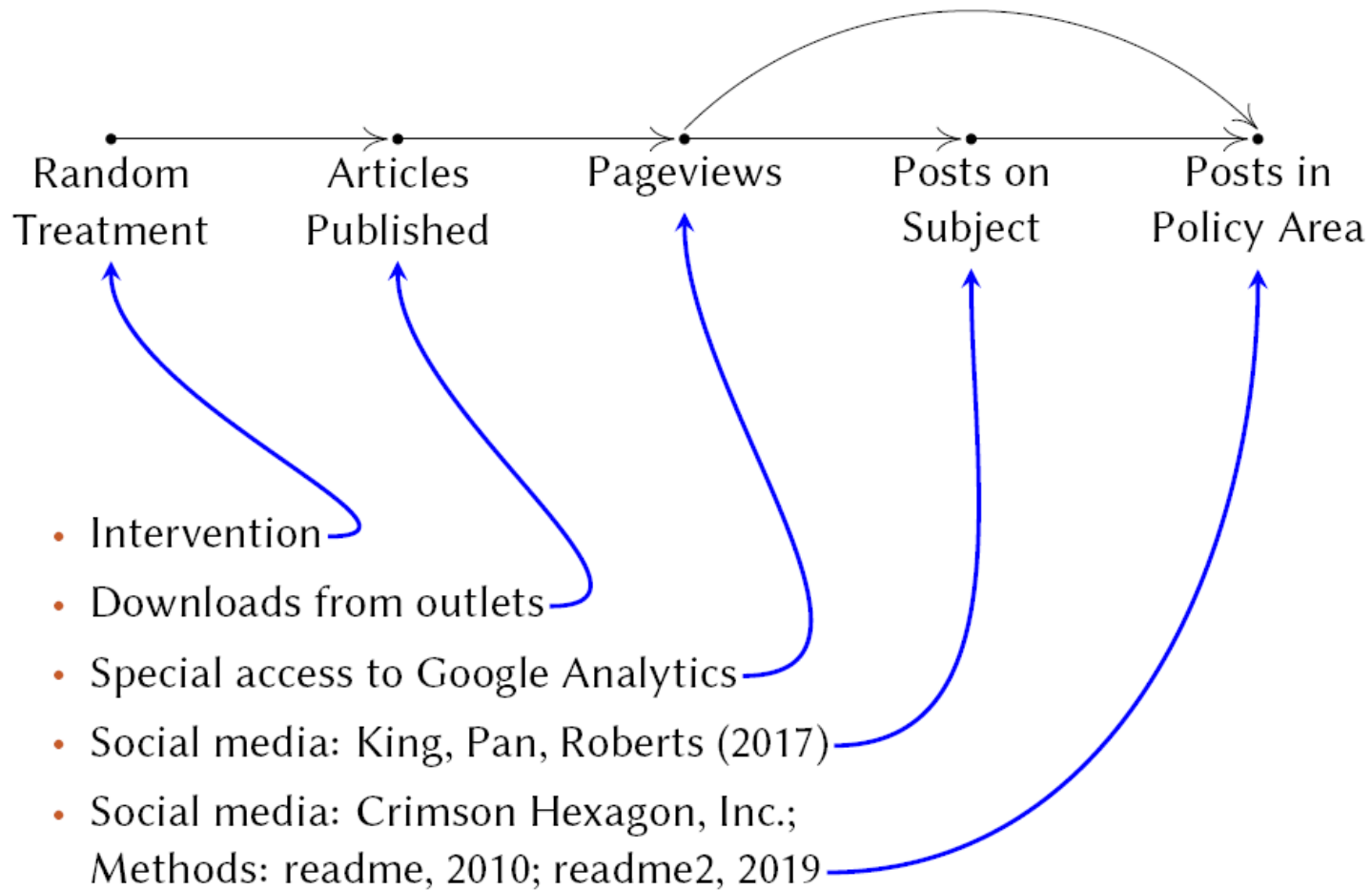
King, G. , Schneer, B. , & White, A. . (2017). How the news media activate public expression and influence national agendas. *Science*, 358(6364), 776-780.

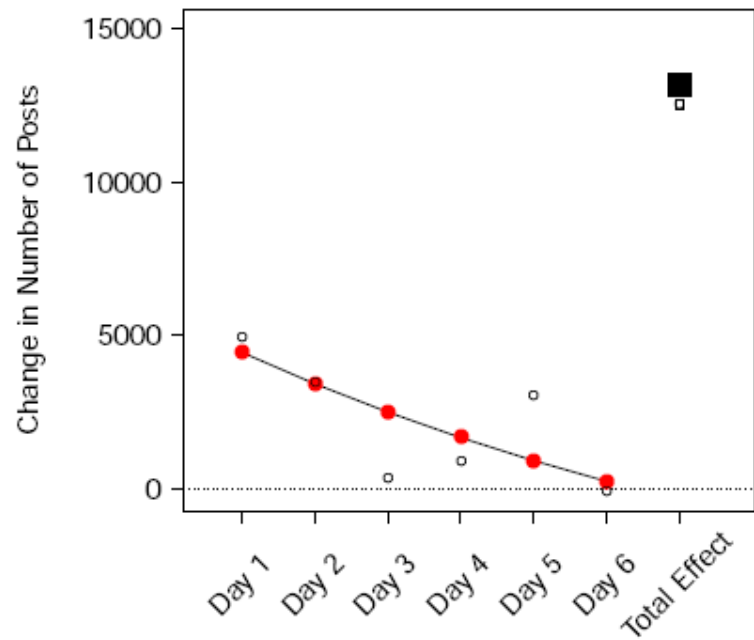
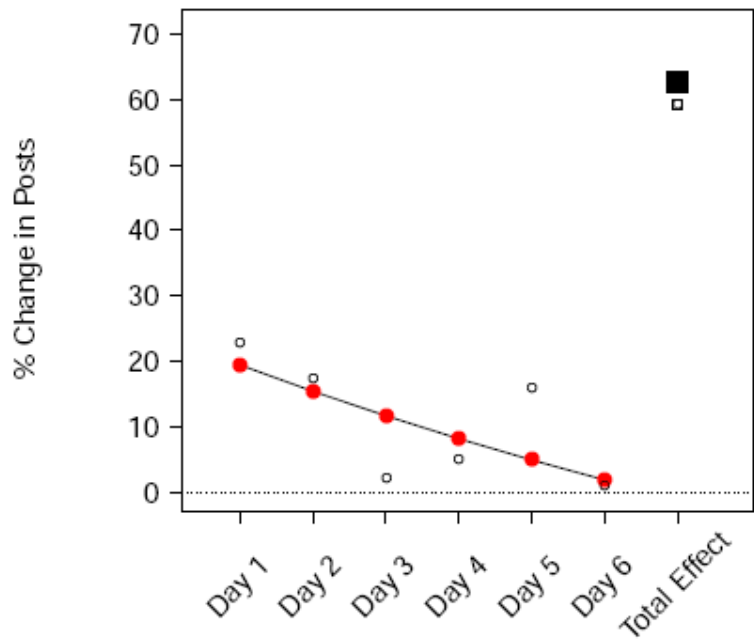
- Signup 48 small media outlets (& > 12 others just for info)
 - 17 for trial runs, 33 in experiment, 2 in both
 - Median size: *The Progressive*, 50,000 subscribers
 - Examples:



- Establish 11 broad *policy areas*
 - Rules: (a) major national importance; (b) interest to outlets
 - race, immigration, jobs, abortion, climate, food policy, water, education policy, refugees, domestic energy production, and reproductive rights

King, G. , Schneer, B. , & White, A. . (2017). How the news media activate public expression and influence national agendas. *Science*, 358(6364), 776-780.

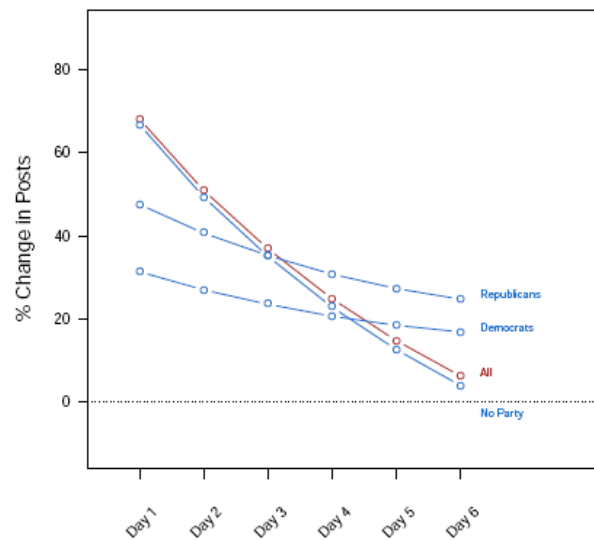
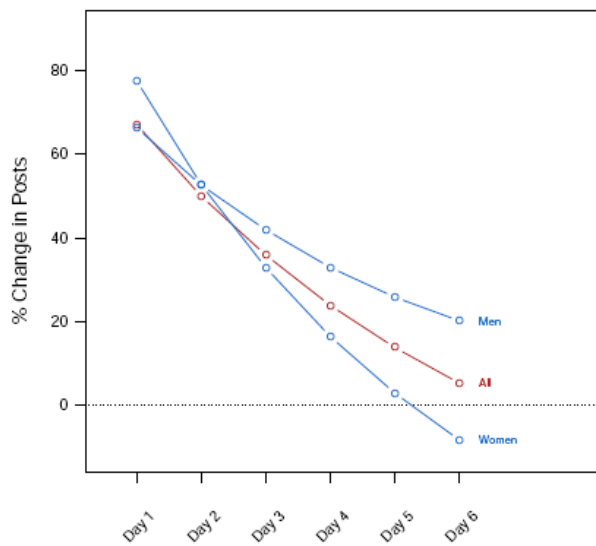
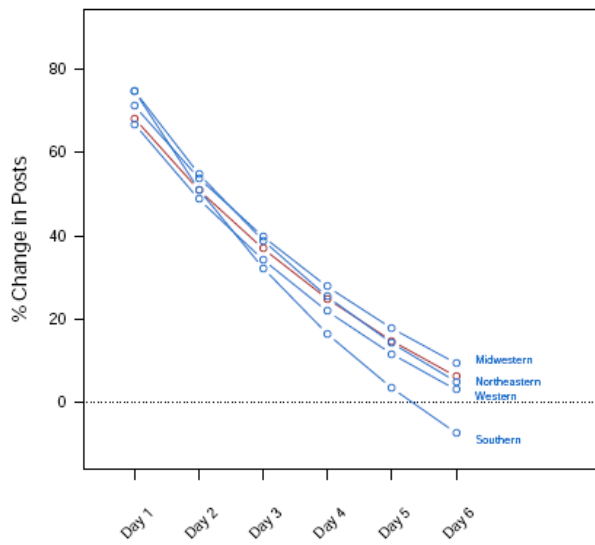




Region

Gender

Party



大数据方法：支持VS批评

- 大数据方法的优势

- 数据：“全量数据”、“消极数据”、“大样本小概率事件”、高维数据
- 方法：机器学习、预测
- 经济性/可行性：低成本、实效性、高效率
- 影响：知识平民化传播

大数据方法的批评

- 大数据方法的局限性

- 数据：“有偏数据”、“分析单位”、“假数据”
- 方法：效度与信度、技术门槛
- 可行性：数据不开放、技术门槛
- 伦理：数据（隐私）权利、社会实验的伦理困境

Q&A