



计算社会科学导论

—— 大数据应用及数据分析基础

金耀辉、许岩岩

2023年2月23日

CS1126



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

你见过的大数据，有多大？

BIG DATA
ADVANCED ANALYTICS
AND VISUALIZATION



- 当社会科学遇上大数据
- 计算社会科学中的常用大数据
- 大数据的缺点、偏见
- 数据科学基础知识
- Python数据分析实践

社会科学是研究人的学科

- 人、群体、社会
- 如何获得个体的信息？
- 如何归纳出群体的特性和交互行为？
- 如何构建社会中的物理规律？



社会学中的复杂性

- 人的复杂性
- 人类行为的不确定性
 - 不能构建典型的、完全理性和全知的 representative agent
- 即使能合理地描述个体行为(甚至是描述概率分布), 也不能轻易将其扩样到整个群体
 - 人是有适应能力的, 在不同环境中的群体有比样本更大的方差
 - 样本在群体中的分布不均匀



大数据如何帮助社会科学

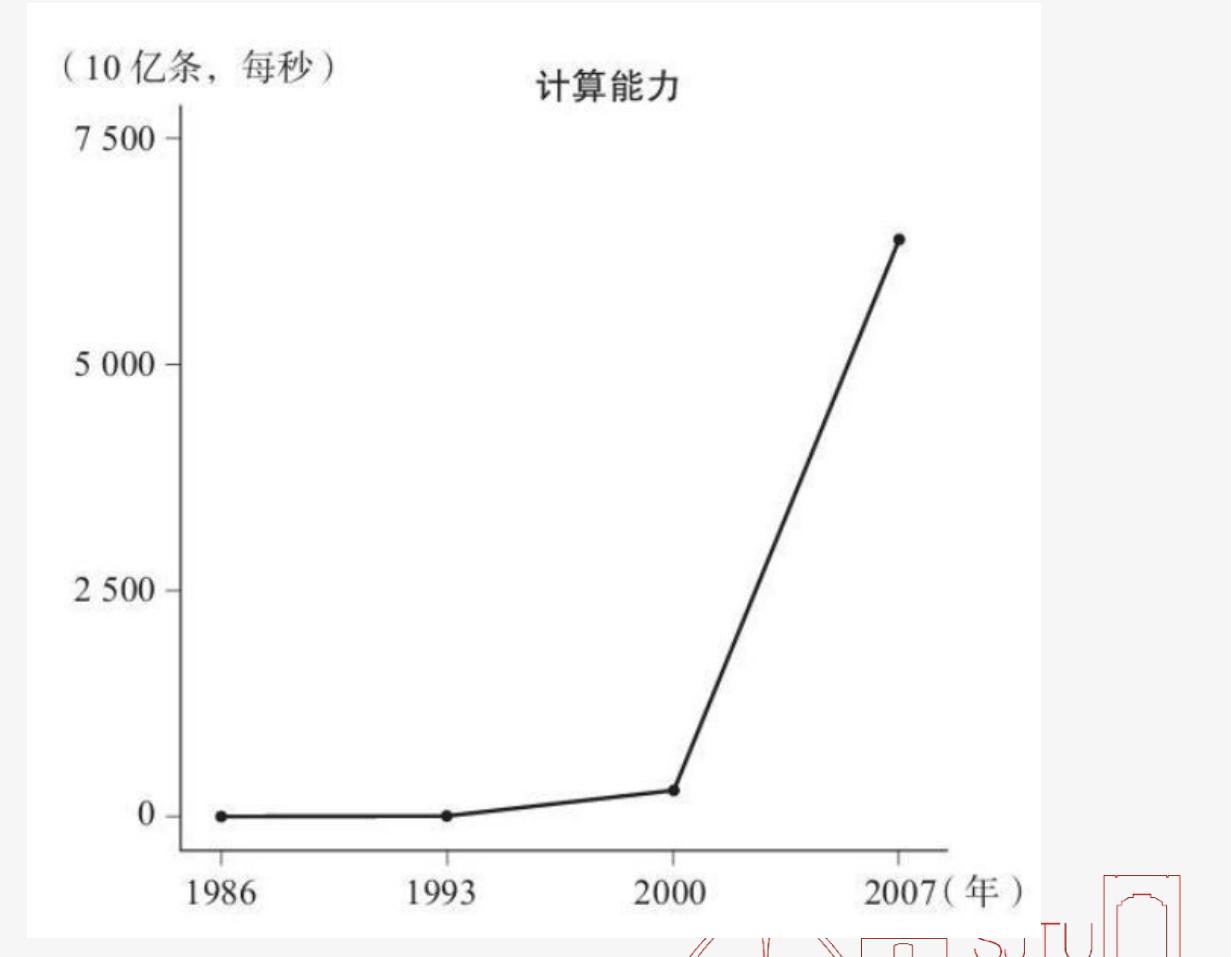
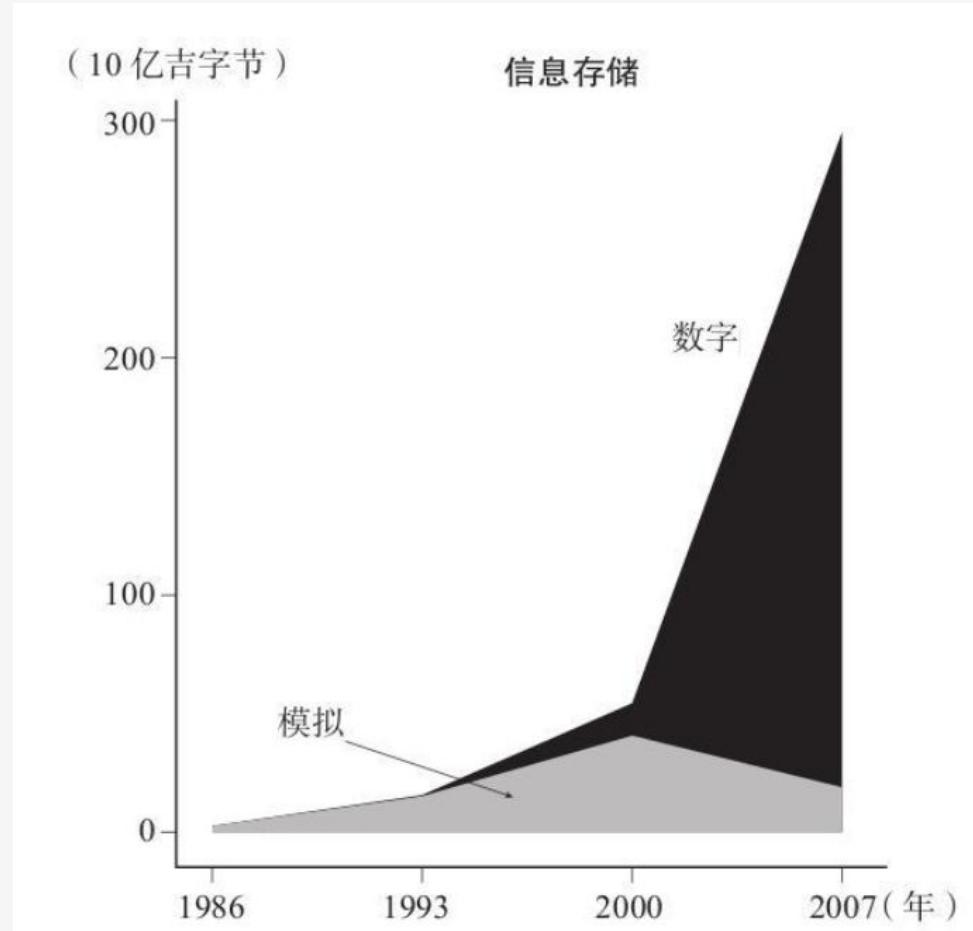


- 如何获得个体的信息？
 - ICT-海量数据
- 如何归纳出群体的特性和交互行为？
 - 数据科学、网络科学
- 如何构建社会中的物理规律？
 - 统计物理



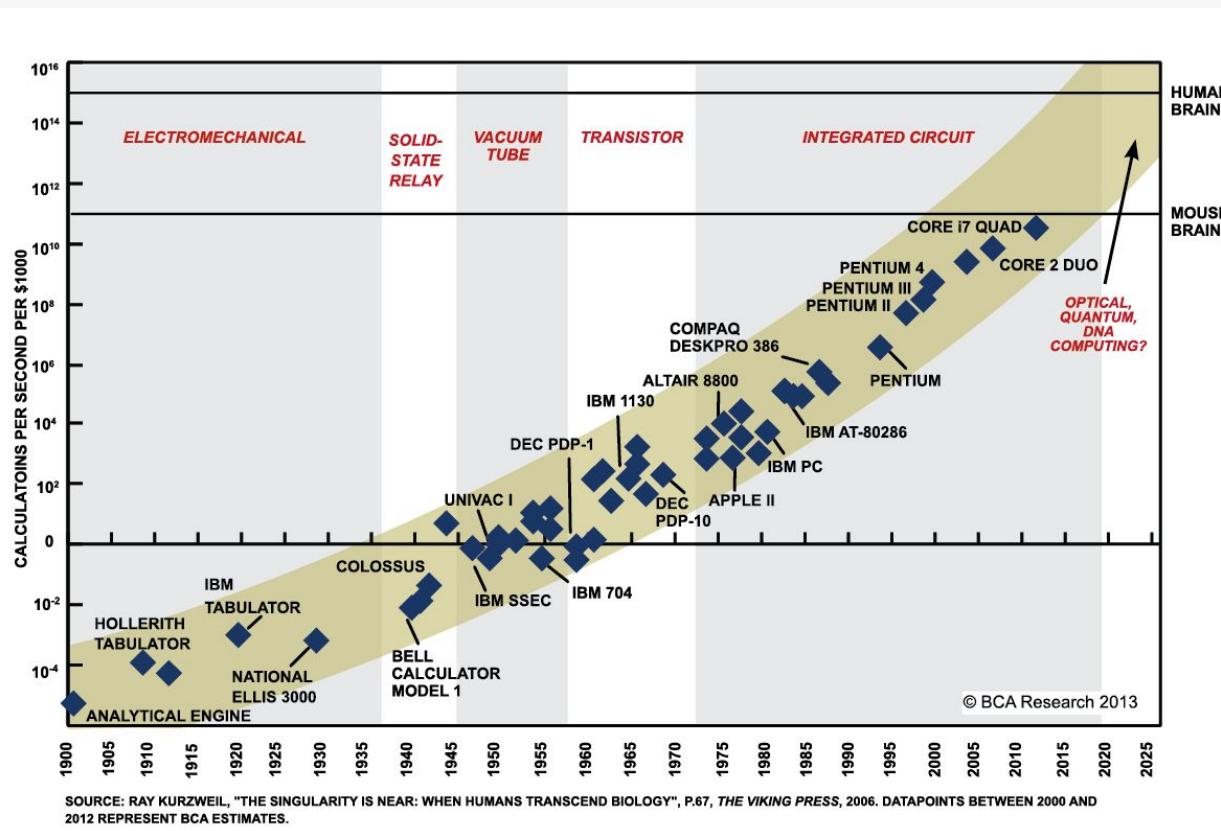
Why now: 模拟到数字的惊人转变

Moore's law



Moore's law

- 有人质疑，事实上算力的持续增加是一直发生的事情，为什么大数据偏偏是这个时间节点出现？



- 互联网是一个被全面监测的环境，非常适合研究人员开展实验。例如在线商城可以搜集到精确的数百万顾客的购买行为数据。
- 实体店也已经搜集了非常详细的购买行为数据，同时它们也正在开发相关基础设施，以便追踪顾客的购买行为，并将实验研究结果用于日常商业活动中。
- 物联网意味着现实世界中的行为会越来越多地被数字传感器捕获。



事实上是多种技术和现象的组合叠加

Massively distributed computing(分布式)

- MapReduce
- Spark
- cloud computing

Big-memory machines(内存)

- Terabytes of RAM

Advances in machine learning(机器学习)

- Deep learning
- transformers
- large language models

Fast streaming algorithms(快速算法)

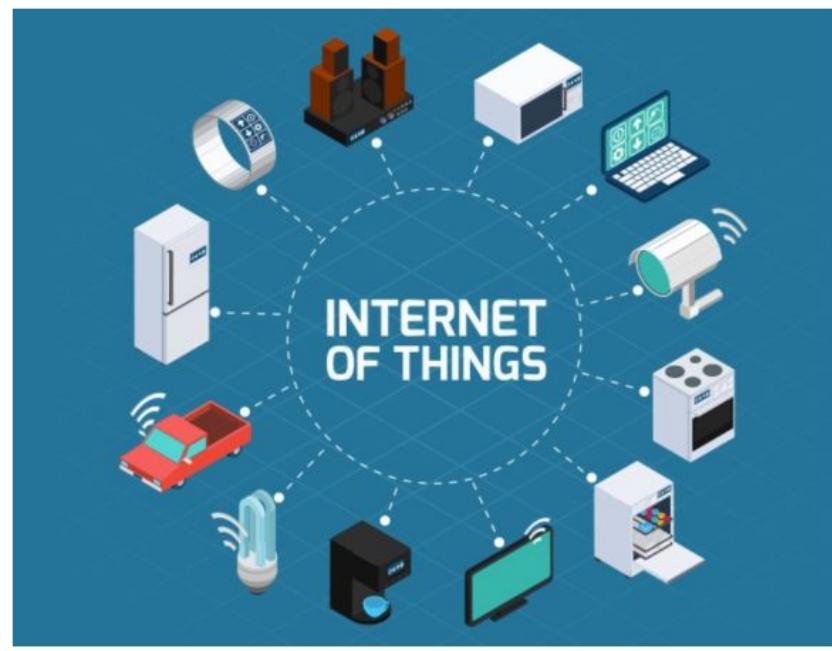
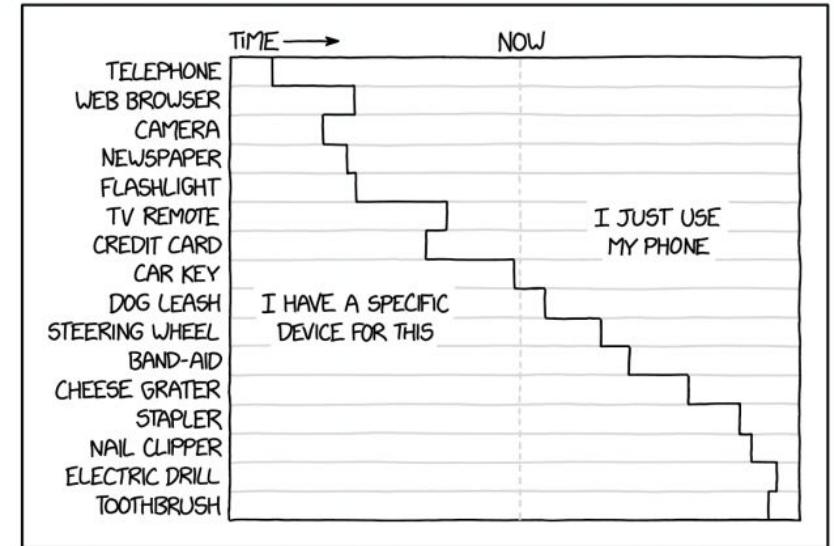
- Streaming aggregation
- stochastic gradient descent

Human computation(人类协作)

- Crowdsourcing
- Mechanical Turk

设备一体化

研究角度 | 用户角度



物联网



大数据在你眼里是：



常见定义

Volume
大量

Variety
多样

Velocity
高速

倡导者认为

Veracity
真实

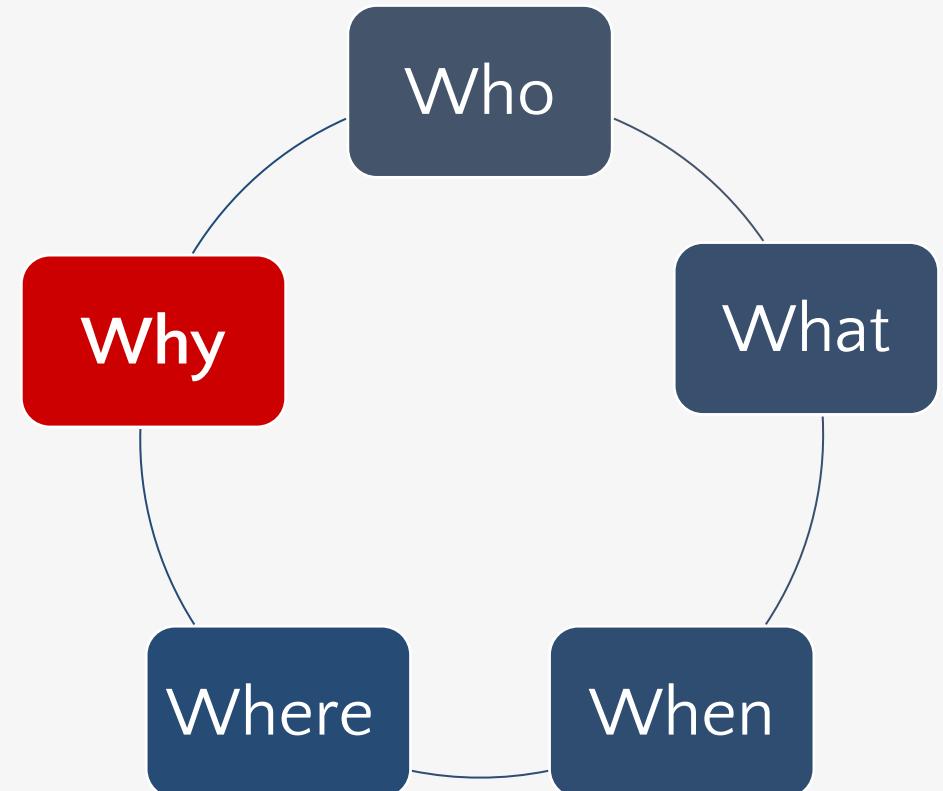
Value
价值

批评者认为

Vague
模糊

Vacuous
空洞

我们应该更关心：



研究流程



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

大量数据

- ICT数据采集
- 运用计算机科学优化存储、计算和传输

计算社会科学

- 交叉学科
- 统计物理、图论、复杂系统等工具建立模型

可以进行

- 分析数据
- 仿真模型
- 控制变量实验
- 验证理论预测
- 分析策略的影响

得到

- 新的行为模式
- 更好地理解复杂的社会系统

目的

- 为决策的制定和实行提供建议
- 可持续发展的社会



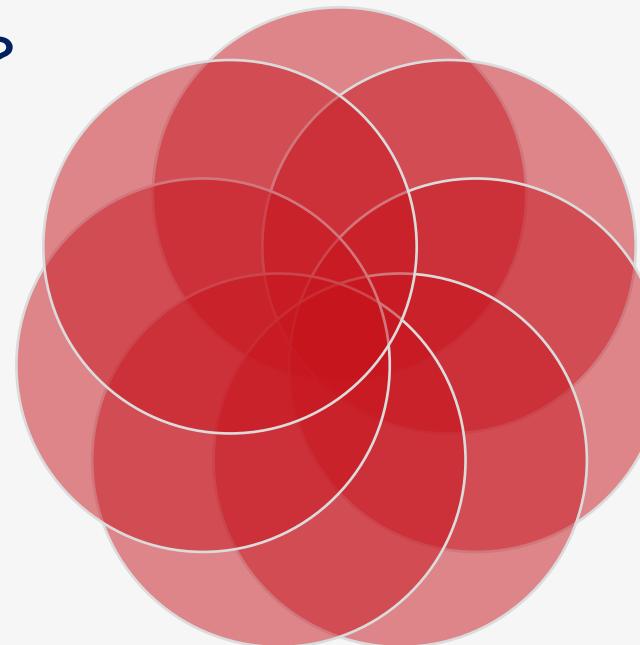
围绕大数据，你可能需要关心



数据数量很多，质量如何？

使用统计、物理、复杂系统等
工具实现学科交叉

用大范围数据验证新的
理论模型



对数据的分类总结：

- 有没有新兴的行为模式和群体行为？

数据能代表多大范围的人群？

- 看似易于收集的在线数据、问卷，能代表多少人？

标准化的数据收集和挖掘流程

- 实验的可复现性

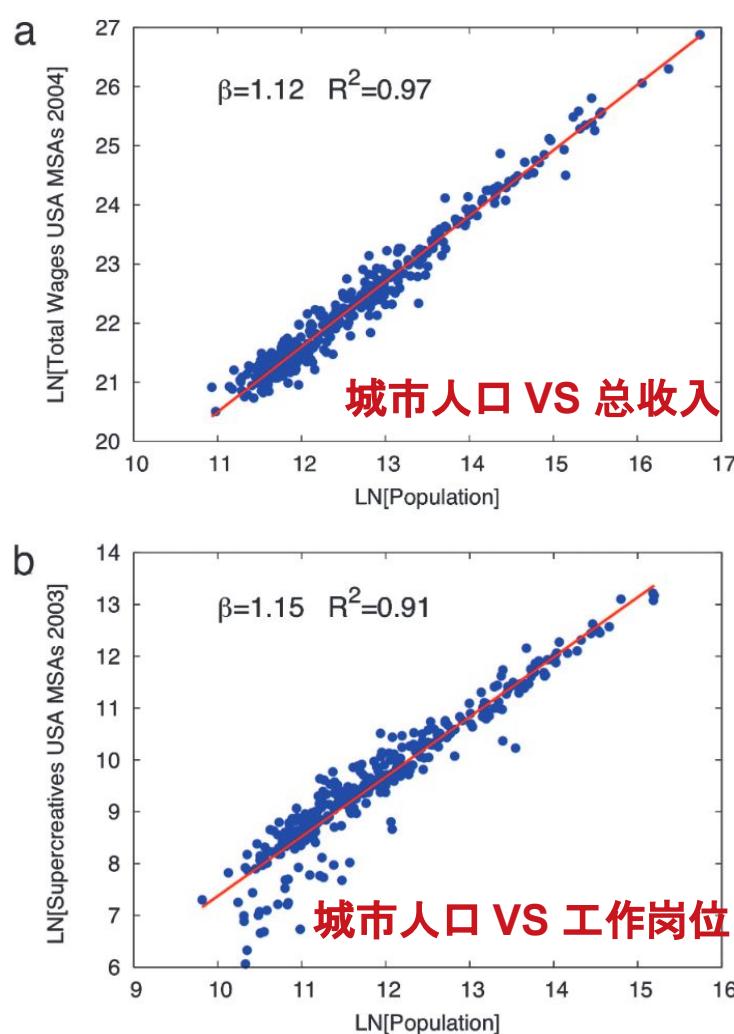
数据使用协议、伦理、隐私





- 当社会科学遇上大数据
- 计算社会科学中的常用大数据
- 大数据的缺点、偏见
- 数据科学基础知识
- Python数据分析实践

人口数量研究GDP、城市发展



$$y = x^\beta e^\alpha$$

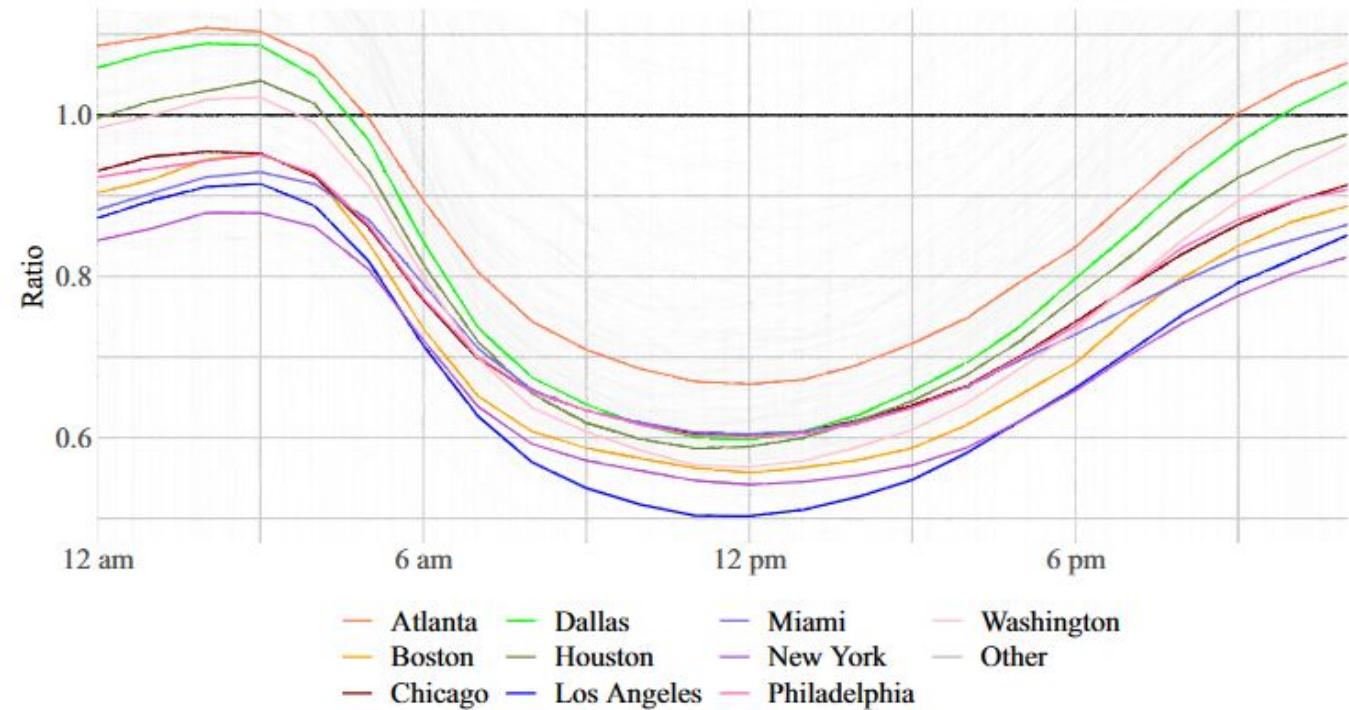
Table 1. Scaling exponents for urban indicators vs. city size

Y	β	95% CI	Adj- R^2	Observations	Country-year
New patents	1.27	[1.25,1.29]	0.72	331	U.S. 2001
Inventors	1.25	[1.22,1.27]	0.76	331	U.S. 2001
Private R&D employment	1.34	[1.29,1.39]	0.92	266	U.S. 2002
"Supercreative" employment	1.15	[1.11,1.18]	0.89	287	U.S. 2003
R&D establishments	1.19	[1.14,1.22]	0.77	287	U.S. 1997
R&D employment	1.26	[1.18,1.43]	0.93	295	China 2002
Total wages	1.12	[1.09,1.13]	0.96	361	U.S. 2002
Total bank deposits	1.08	[1.03,1.11]	0.91	267	U.S. 1996
GDP	1.15	[1.06,1.23]	0.96	295	China 2002
GDP	1.26	[1.09,1.46]	0.64	196	EU 1999–2003
GDP	1.13	[1.03,1.23]	0.94	37	Germany 2003
Total electrical consumption	1.07	[1.03,1.11]	0.88	392	Germany 2002
New AIDS cases	1.23	[1.18,1.29]	0.76	93	U.S. 2002–2003
Serious crimes	1.16	[1.11, 1.18]	0.89	287	U.S. 2003
Total housing	1.00	[0.99,1.01]	0.99	316	U.S. 1990
Total employment	1.01	[0.99,1.02]	0.98	331	U.S. 2001
Household electrical consumption	1.00	[0.94,1.06]	0.88	377	Germany 2002
Household electrical consumption	1.05	[0.89,1.22]	0.91	295	China 2002
Household water consumption	1.01	[0.89,1.11]	0.96	295	China 2002
Gasoline stations	0.77	[0.74,0.81]	0.93	318	U.S. 2001
Gasoline sales	0.79	[0.73,0.80]	0.94	318	U.S. 2001
Length of electrical cables	0.87	[0.82,0.92]	0.75	380	Germany 2002
Road surface	0.83	[0.74,0.92]	0.87	29	Germany 2002

Data sources are shown in *SI Text*. CI, confidence interval; Adj- R^2 , adjusted R^2 ; GDP, gross domestic product.

GPS轨迹数据

- 使用来自智能手机的GPS数据，测量了美国城市中不同种族之间的**经历隔离程度**。
- 手机数据测量的经历隔离比传统的居住隔离更能捕捉到个体真实面对的多样性和机会。
- 个体经历的隔离程度比标准的居住隔离指标要低得多；在不同城市之间，经历隔离和居住隔离有很高的相关性。



- 作者也提到了这篇文章中数据本身存在的问题

- 手机数据只能体现出设备出现在同一个空间区域中，并不代表用户真实的互动。这使得文章研究的种族隔离更偏向于地理隔离而非社会学隔离，尽管地理隔离这一概念同样有意义；
- 没有手机数据的人种信息用居住地白人占比给用户打上白人或非白人的标签，再用于计算种族隔离；
- 手机用户并不总能代表总人口

Athey S, Ferguson B, Gentzkow M, et al. Estimating experienced racial segregation in US cities using large-scale GPS data[J]. Proceedings of the National Academy of Sciences, 2021, 118(46): e2026160118.

LBS、航班

- 除非能减少50%以上的社区传播，否则对中国大陆进行持续90%的旅行限制只会适度影响疫情轨迹。旅行限制只能推迟疫情的传播(起点)，并不能降低传播速度。

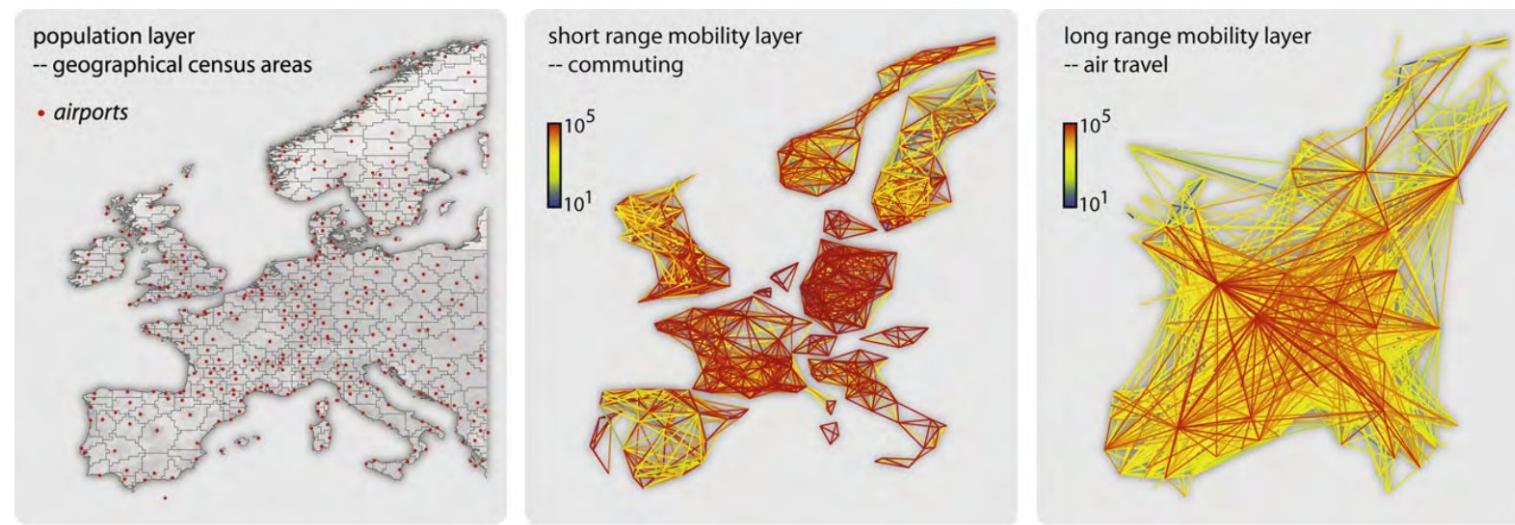


Fig. 1. GLEaM, GLobal Epidemic and Mobility model. The world surface is represented in a grid-like partition where each cell – corresponding to a population value – is assigned to the closest airport. Geographical census areas emerge that constitute the subpopulations of the metapopulation model. The demographic layer is coupled with two mobility layers, the short range commuting layer and the long range air travel layer.

- GLEAM疫情传播模型将世界被划分为以主要交通枢纽(通常是机场)为中心的子种群。这些群体通过每天在它们之间旅行的个人流量连接起来。该模型包括大约200个不同国家和地区的3200多个群体。航空运输数据包括来自官方航空指南(OAG)和国际航空运输协会(IATA)数据库(2019年更新)的每日出发地-目的地交通流量，而地面流动流量则来自对从五大洲30个国家的统计局收集的数据的分析和建模。中国大陆的人口移动变化来自于百度的定位服务(LBS)。

Chinazzi M, Davis J T, Ajelli M, et al. The effect of travel restrictions on the spread of the 2019 novel coronavirus (COVID-19) outbreak[J]. Science, 2020, 368(6489): 395-400.

Wifi数据了解生活节律

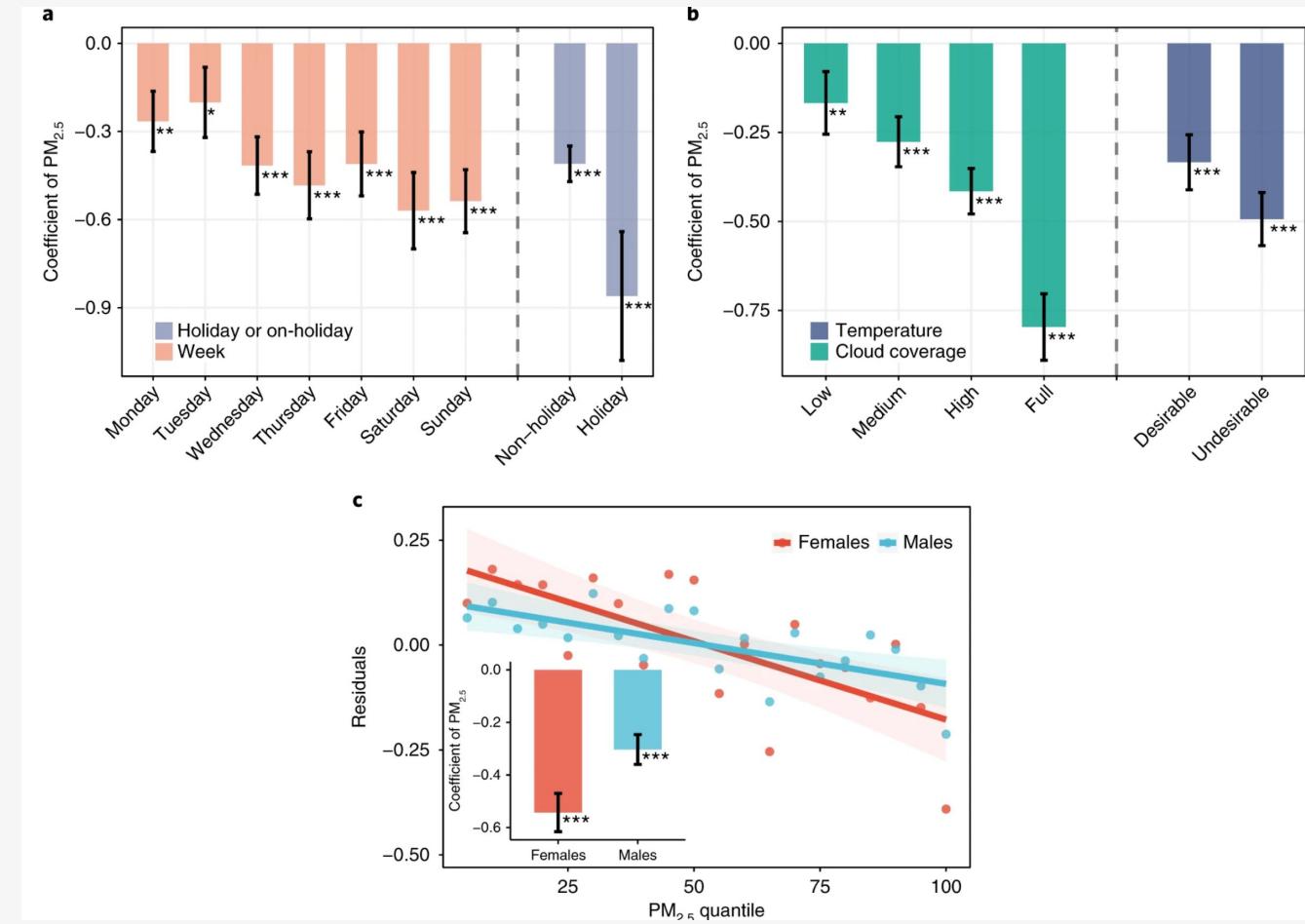


上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



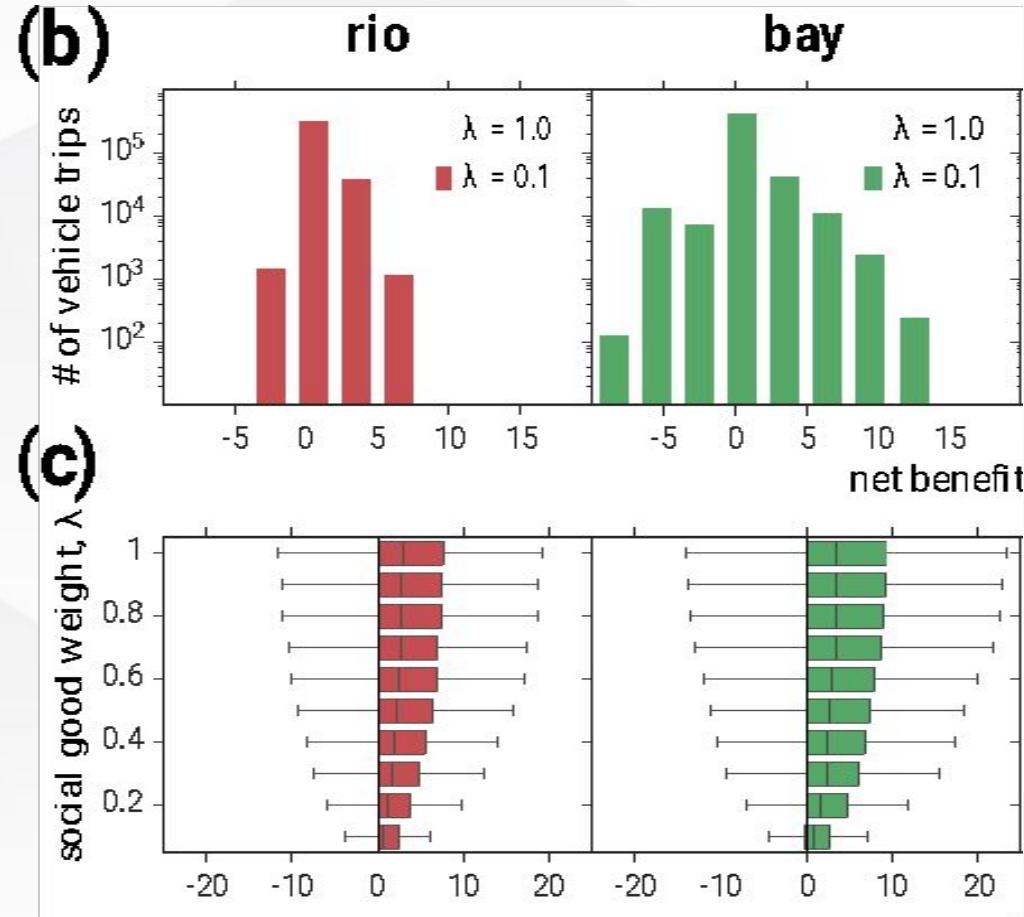
社交媒体信息和大气污染的关系

- 通过对[中国最大的微博平台新浪微博](#)上发布的2.1亿条带有地理标签的微博推文应用NLP领域相关算法来衡量144个城市居民的每日**幸福感**。
- 研究144个中国城市，由于对污染的讨论不一定能反映出个人潜在情绪状态的变化，使用没有污染相关术语的推文来构建我们的城市/日**幸福指数**。
- 证实：大气污染的加剧与居民幸福感之间的相关性，且女性对大气污染更加敏感。**



出行需求的获取:社会调查? GPS? LBS? 手机信令? 路口摄像头?
有什么优缺点?

- a) 某用户从旧金山市区出发,前往国际机场。其“自私”路径(UE)的行程时间为20min,而“无私”路径(SO)行程时间为25min。
- b) 里约及湾区在不同的“无私”程度(λ 不同取值)下的用户出行时间节约百分比分布
- c) 里约及湾区在不同的“无私”程度(λ 不同取值)下的所有居民总体出行时间节约百分比。





出行调查数据



人工智能研究院
Artificial Intelligence Institute

以美国麻省为例：

\$200 per usable Survey

1 sample day,

2.5×10^4 households out of 2.6×10^6

58% response rate.

(3.7 calls and 17 minutes per survey)



2011-2012

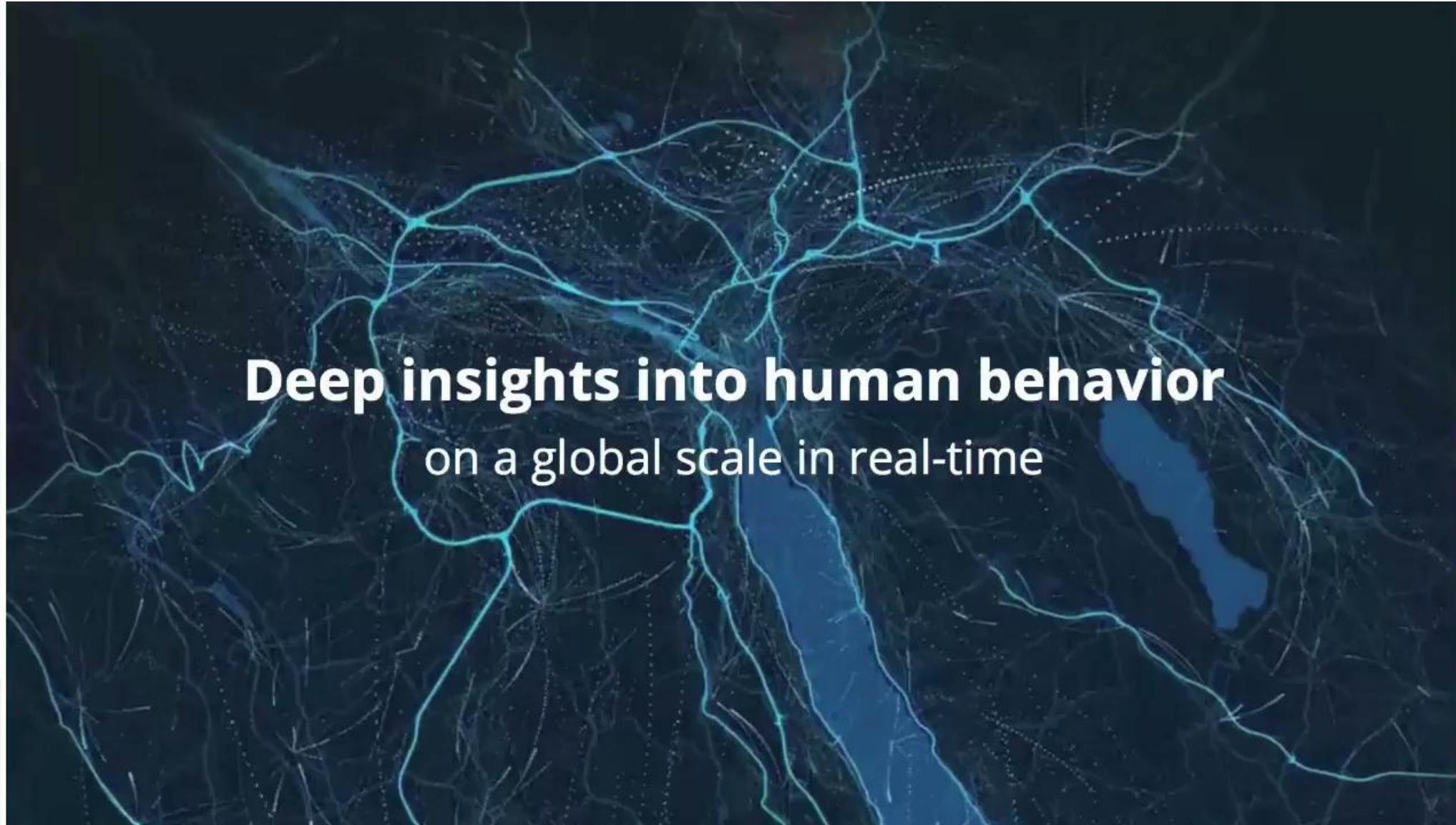


MIT
Technology Review
10 BREAKTHROUGH
TECHNOLOGIES 2013

Big Data
From
Cheap Phones

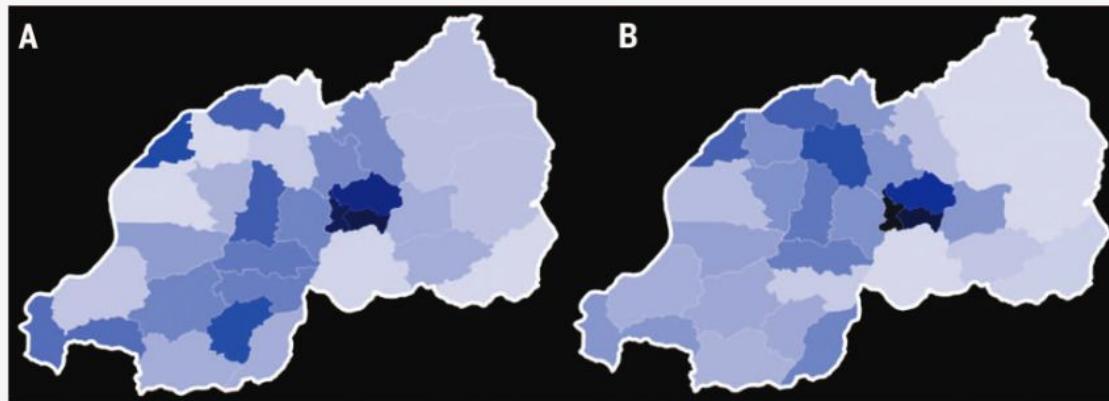


**Deep insights into human behavior
on a global scale in real-time**

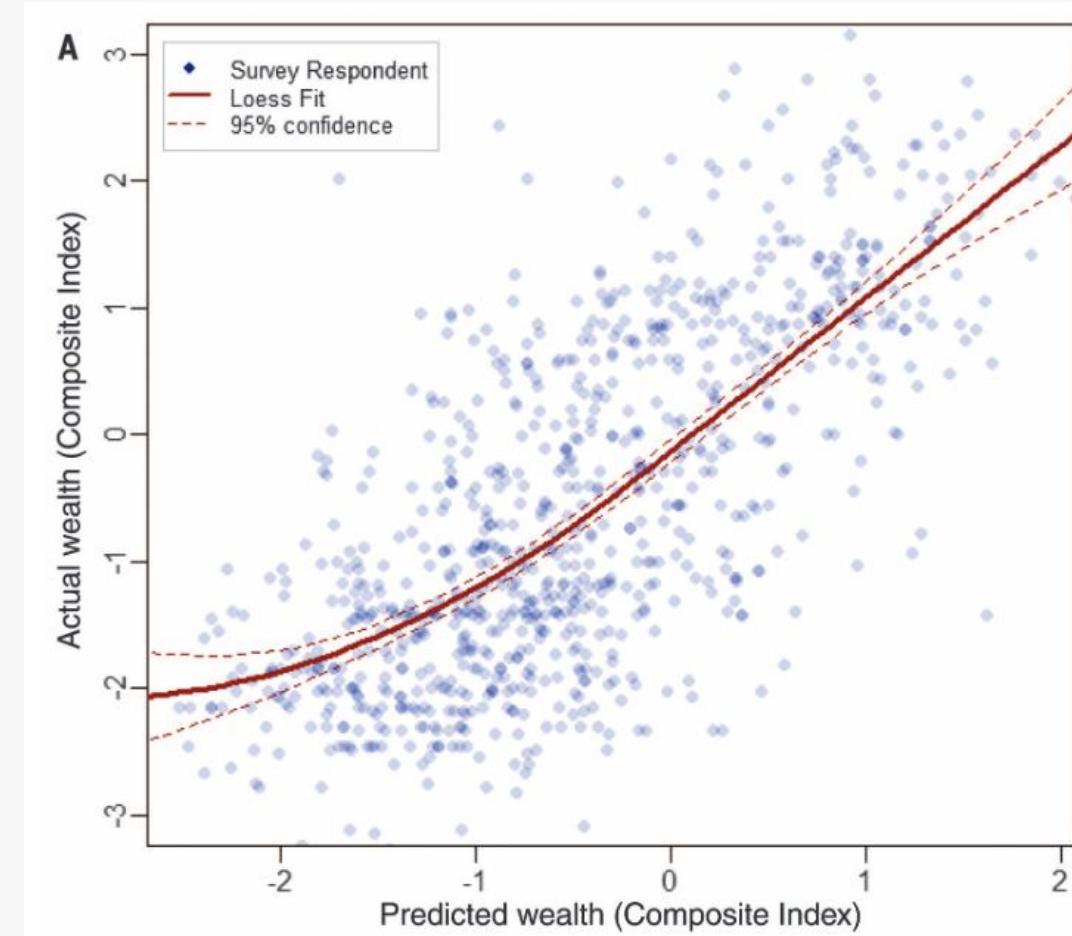


Review-手机metadata

- 个体的手机使用记录可以预测其社会经济地位
- 并且可以用于重建国家或小型地区的资产分布情况
- 在资源受限地区减少进行大范围普查的成本和时间



Blumenstock J, Cadamuro G, On R. Predicting poverty and wealth from mobile phone metadata. Science, 2015, 350(6264): 1073-1076.



上图中, a:预测财富指数, 横轴为预测财富, 纵轴为真实财富
左图中, A和B分别表示预测出的区域平均财富指数和调查结果

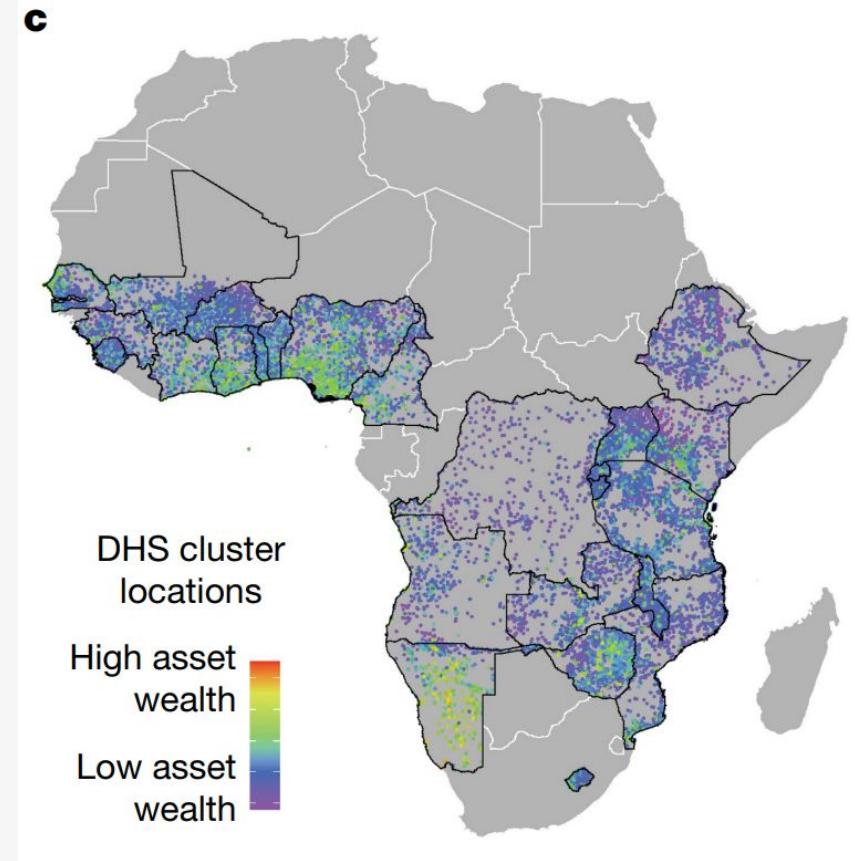
Review-遥感图像

- 机器学习携手卫星影像，理解电力设施与经济财富的因果关系 (Nature 封面文章)
- 背景：乌干达在2010-2019年将电气覆盖率从12%提升到了41%
- 问题：电网扩张如何影响低收入地区经济产出？缺少不同时/空的统计数据？因果证据？
- 利用遥感和机器学习模型预测精细空间粒度下的乌干达资产财富指标：

利用多光谱遥感数据为输入，基于撒哈拉以南非洲地区(SSA)27000个村庄的统计调查数据中的多项相关指标，构造资产财富指数作为标签，使用深度学习估计25个非洲国家2005-2018年的资产财富状况；

填补了超过已有数据10倍以上的空白数据

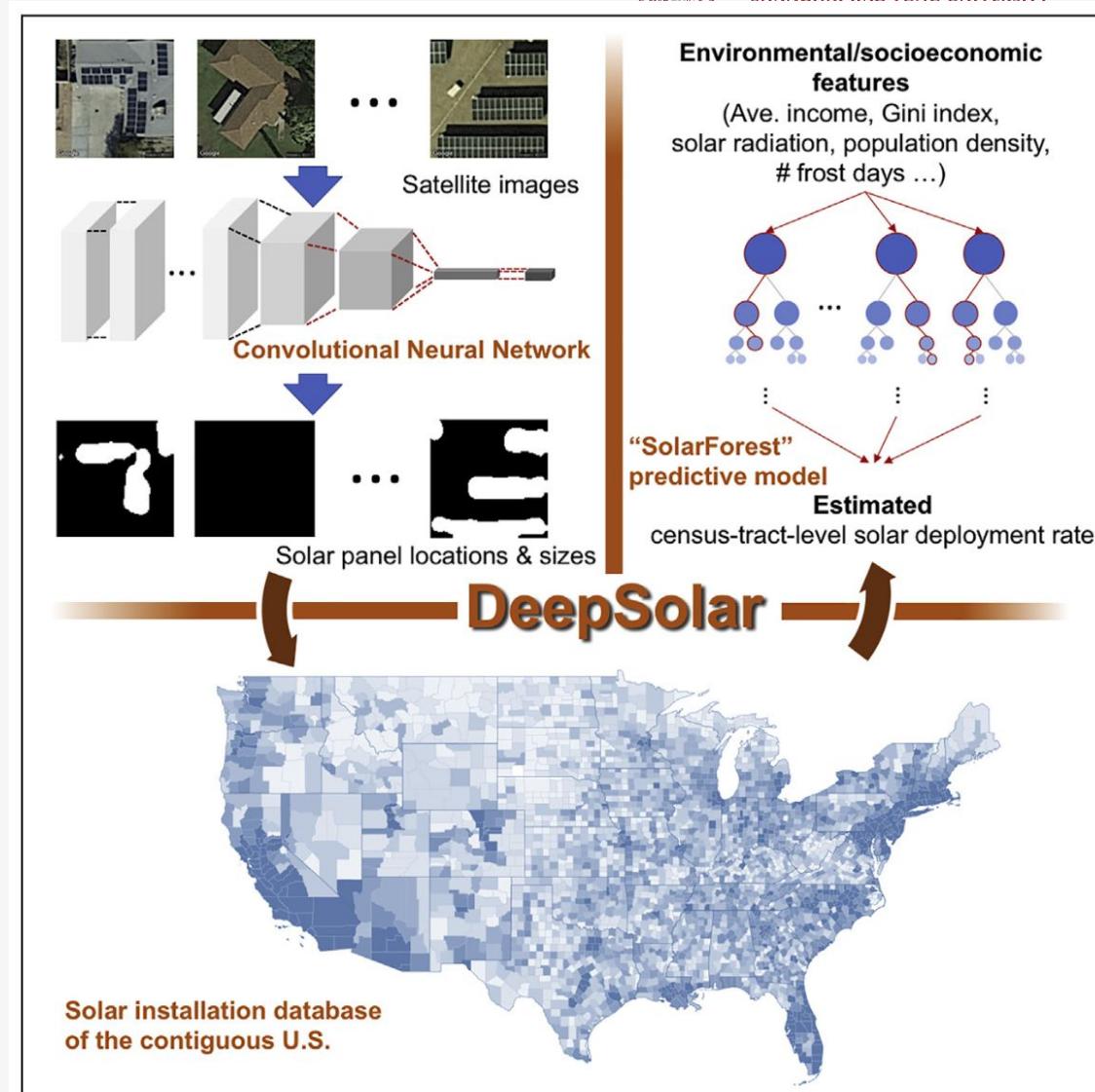
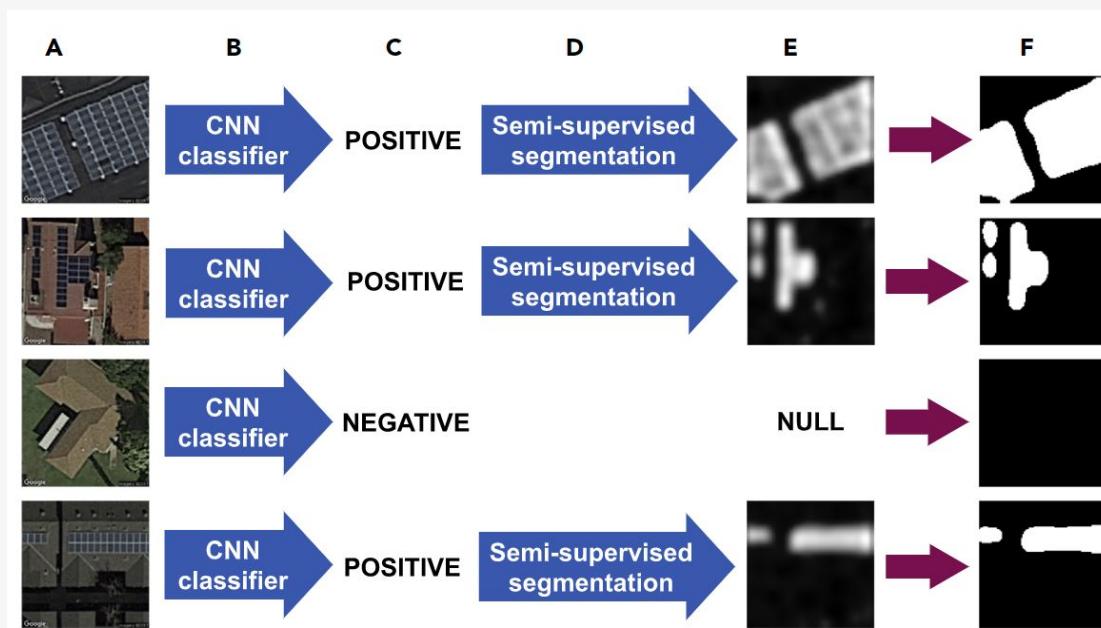
只使用日间卫星图像，以尽量少通过灯光获得电网的信息

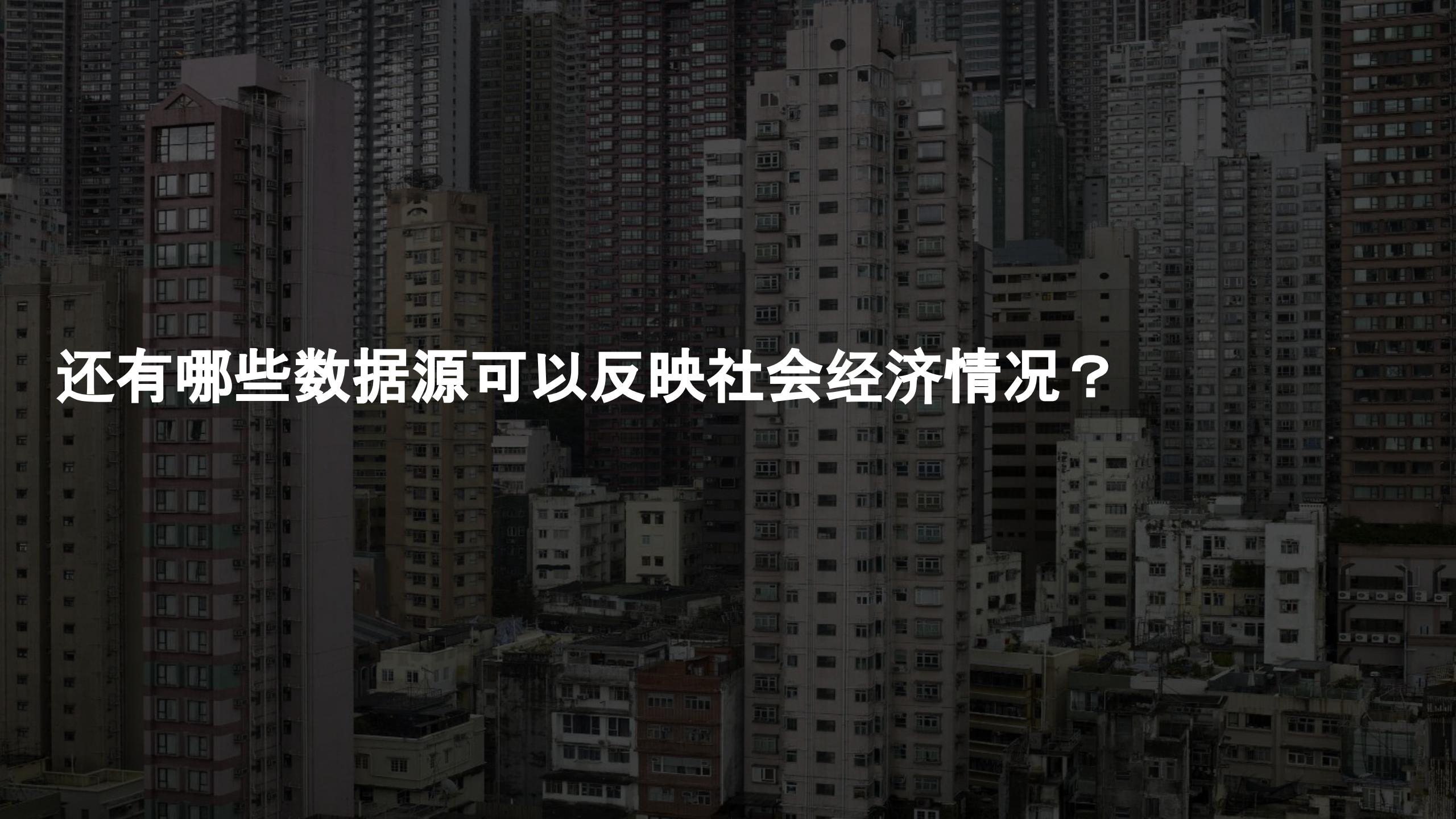


Ratledge, N., Cadamuro, G., de la Cuesta, B. et al. Using machine learning to assess the livelihood impact of electricity access. *Nature* **611**, 491–495 (2022).

DeepSolar 卫星图像推测太阳能普及程度

- 用计算机视觉方法从卫星图像中构建太阳能充电板的分布位置数据库，并和环境与社会经济指标建立相关关系。



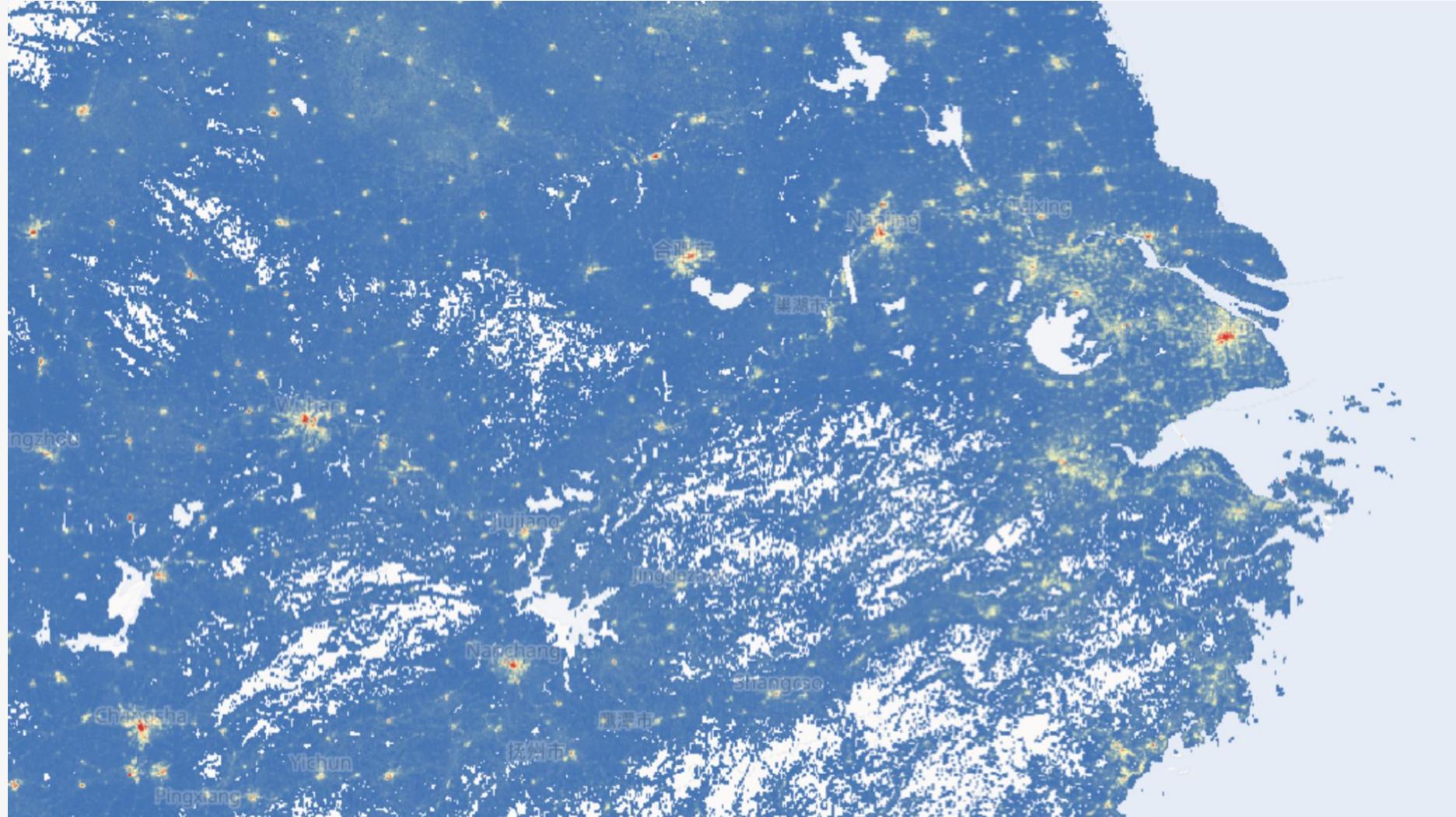
The background image shows a dense urban landscape in Hong Kong, featuring numerous high-rise apartment buildings packed closely together. The buildings vary in height and color, with many having multiple balconies and windows. Some have distinct architectural features like red roofs or unique window patterns. The overall scene conveys a sense of a crowded, metropolitan environment.

还有哪些数据源可以反映社会经济情况？

城市建筑物容量数据



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



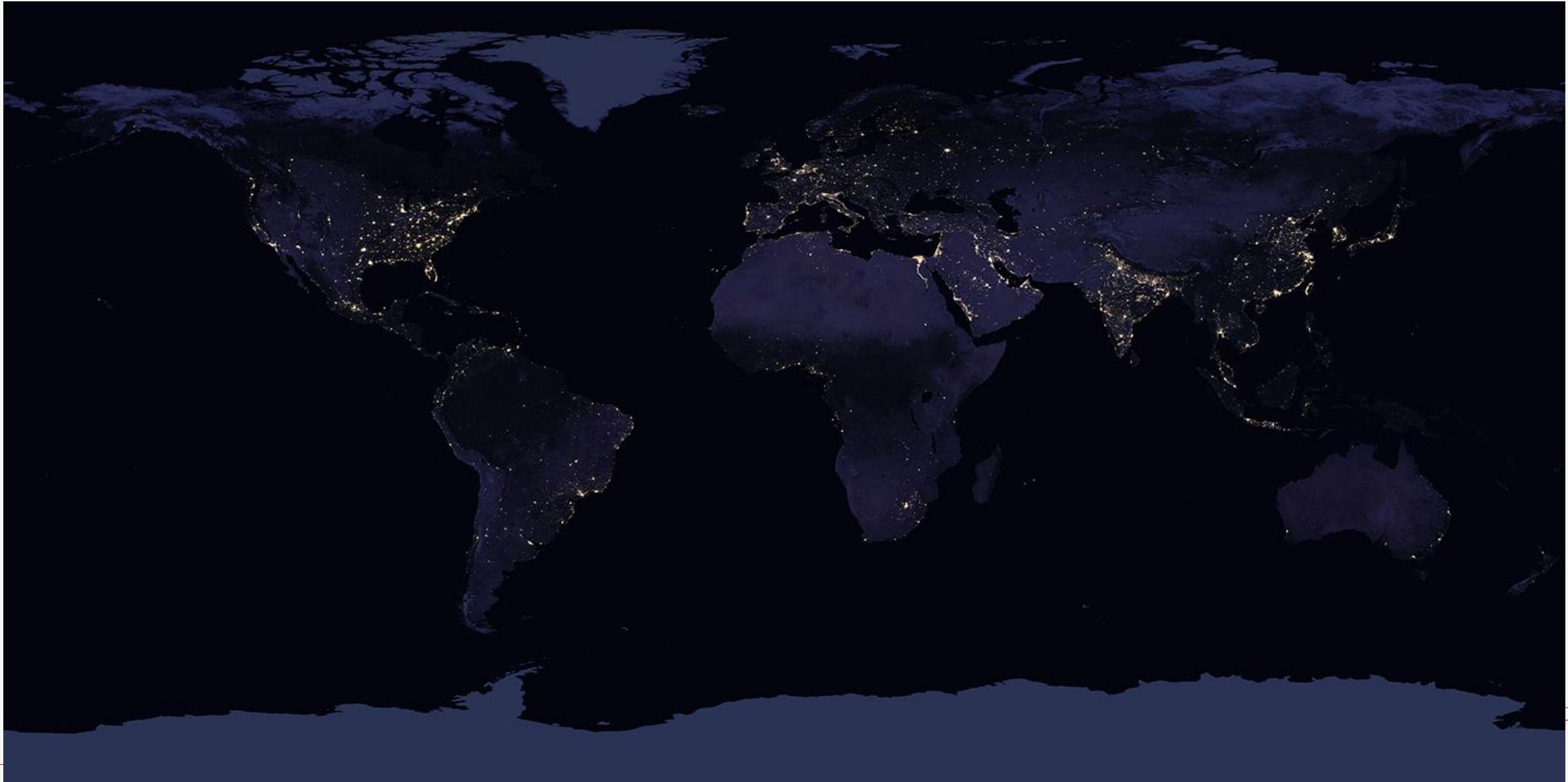
<https://geoservice.dlr.de/web/maps/eoc:wsf3d#>



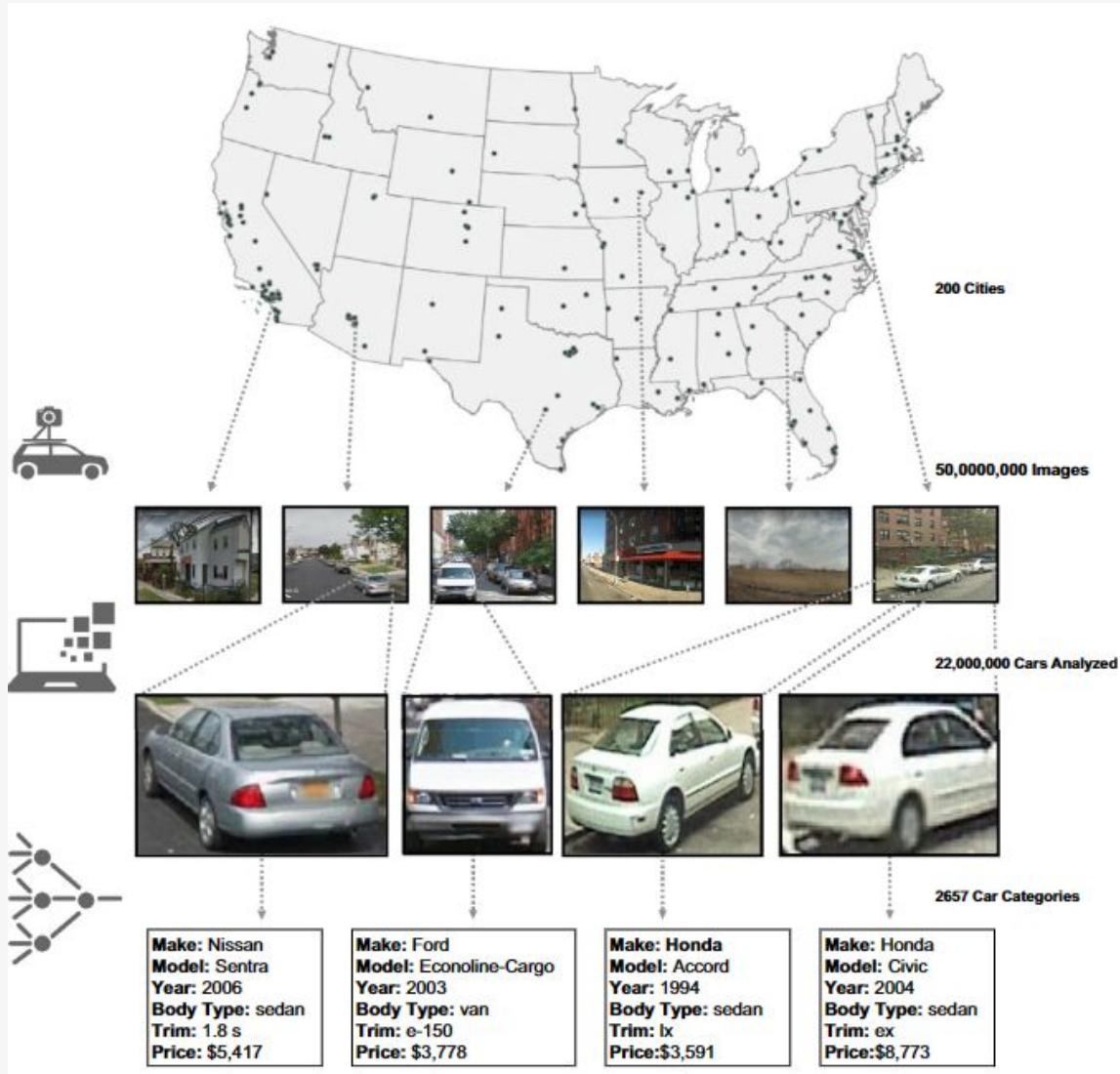
夜光卫星图像数据



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



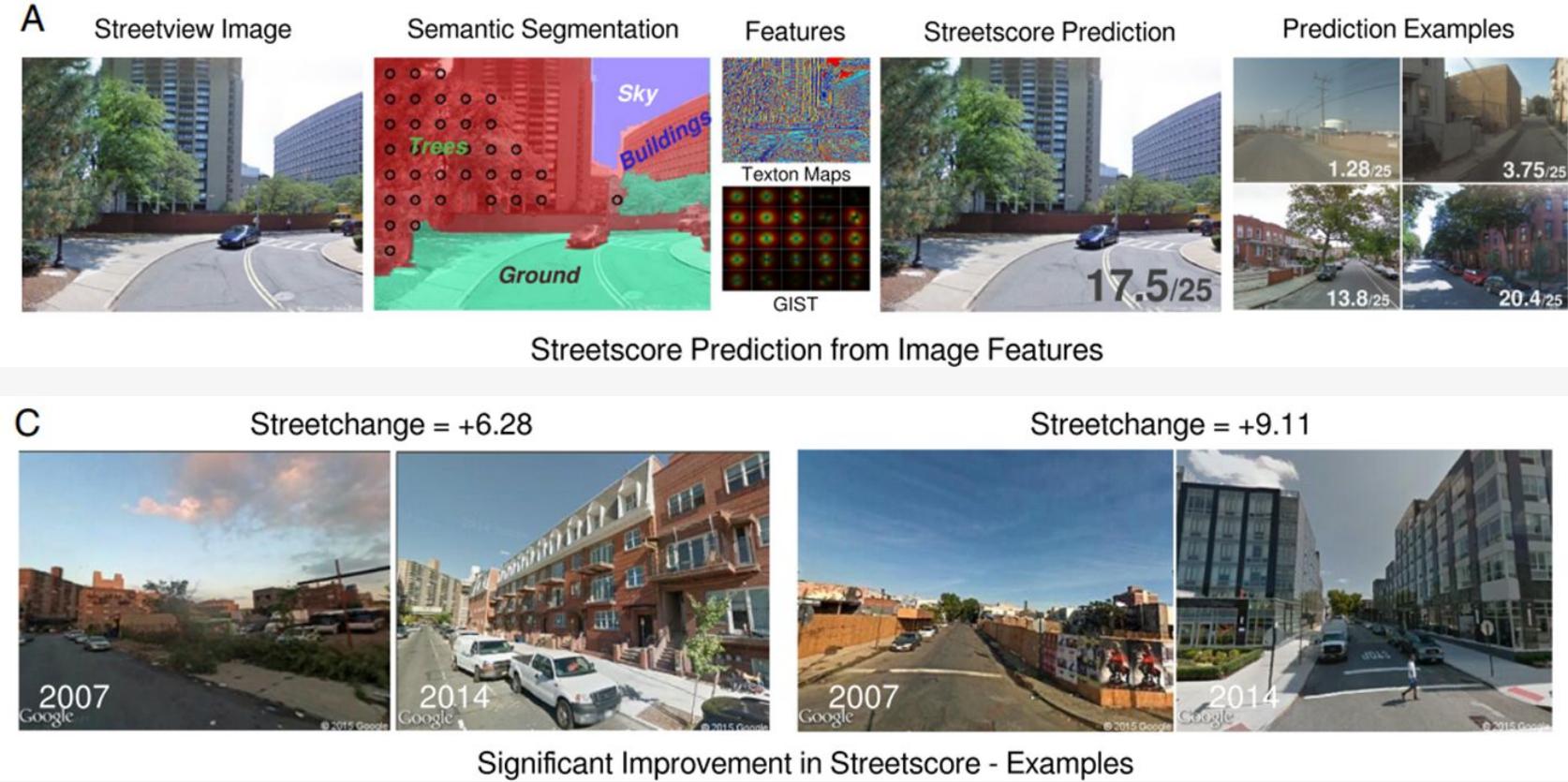
街景数据



- 美国每年花费超过2.5亿美元进行美国社区调查(ACS)，测量与种族、性别、教育、职业、失业和其他人口因素有关的统计数据。尽管是一个全面的数据来源，但人口变化和它们出现在ACS中的时间间隔可能超过几年。
- 通过使用谷歌街景车收集的5000万张**街景图像**，使用计算机视觉方法确定了在特定街区遇到的所有机动车辆的品牌、型号和年份。来自这次机动车普查的数据，共统计了2200万辆汽车(占美国所有汽车的8%)，被用来**准确估计收入、种族、教育以及在邮政编码和选区层面的投票模式**。例如，如果开车经过一个城市时发现轿车的数量高于皮卡的数量，那么这个城市在下一次总统选举中可能会投票给民主党(88%的可能性)；否则，它可能会投票给共和党(82%)。

- Gebru T, Krause J, Wang Y, et al. Using deep learning and Google Street View to estimate the demographic makeup of neighborhoods across the United States[J]. Proceedings of the National Academy of Sciences, 2017, 114(50): 13108-13113.

街景数据2



哪些街区经历了实际的改善？计算机视觉方法，从历史街道图像中估计街区外观变化。

发现可以改善社区的三个因素：

- 由受过大学教育的成年人密集居住的社区更有可能经历物质上的改善
- 初始条件较好的社区经历了较大的正向改善
- 离CBD等其他有吸引力的街区越近，这个街区改善越大

Naik N, Kominers S D, Raskar R, et al. Computer vision uncovers predictors of physical urban change[J]. Proceedings of the National Academy of Sciences, 2017, 114(29): 7571-7576.



WorldPop Hub

DATA | CONTACT

Open Spatial Demographic Data and Research

WorldPop develops peer-reviewed research and methods for the construction of open and high-resolution geospatial data on population distributions, demographic and dynamics, with a focus on low and middle income countries.

Datasets

Open access spatial demographic datasets built using transparent approaches.

[total 44,745 datasets]

[Population Count](#) 20,724

[Population Density](#) 9,955

[Population Weighted Density](#) 4

[Births](#) 234

[Pregnancies](#) 234

[Age and sex structures](#)

6,036

[Development Indicators](#)

42

[Dependency Ratios](#) 2

[Internal Migration](#) 4

[Dynamic Mapping](#) 2

[Global Flight Data](#) 3

[Global Holiday Data](#) 5

[Covariates](#) 6,474

[Grid-cell surface areas](#) 250

[Administrative Areas](#) 500

[Urban change](#) 27

[Global Settlement Growth](#)

249





- 当社会科学遇上大数据
- 计算社会科学中的常用大数据
- 大数据的缺点、偏见
- 数据科学基础知识
- Python数据分析实践

大数据是永远正确的吗？

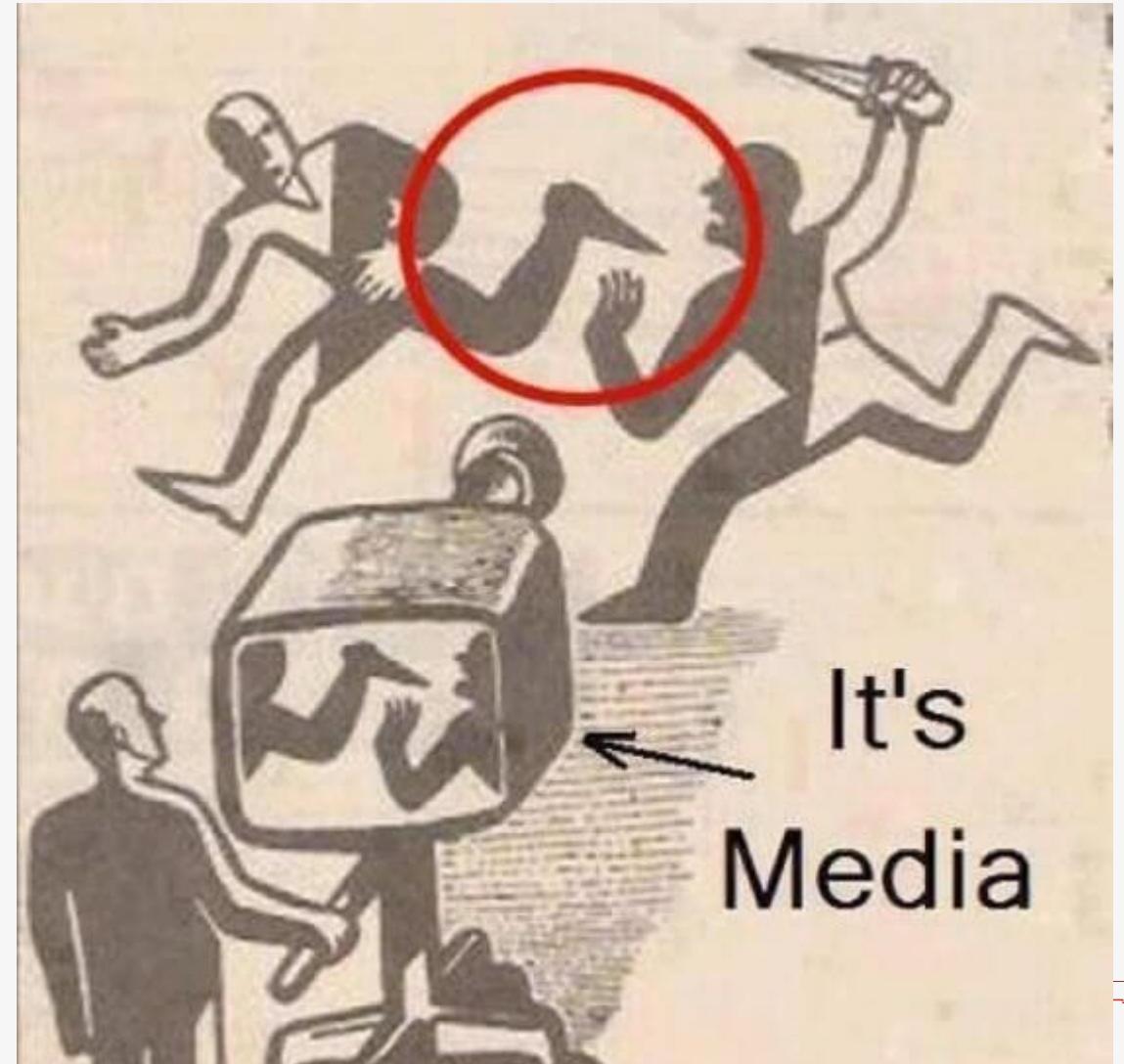


- 💥 数据不等于真相
- 💥 存在不等于合理
- 💥 相关不等于因果
- 💥 过去不等于未来
- 💥 局部不等于全局



数据不等于真相

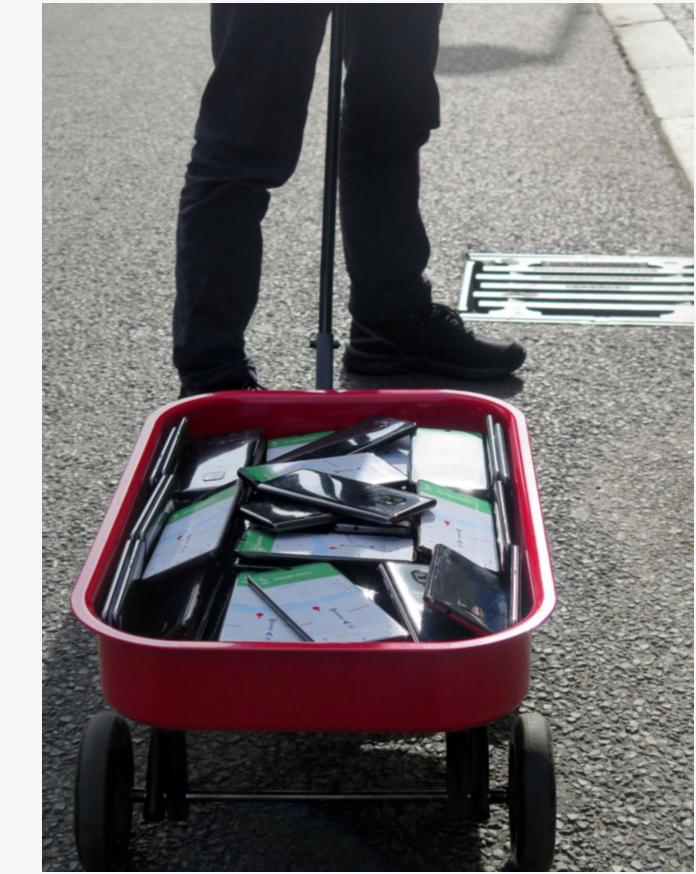
- 统计学中对概率的描述可能误导没有统计学常识的人。
 - 大概率的事情一定发生吗？概率为0的事情一定不可能发生吗？
- 数据也是一个镜头，但如果多个镜头呢？



数据不等于真相



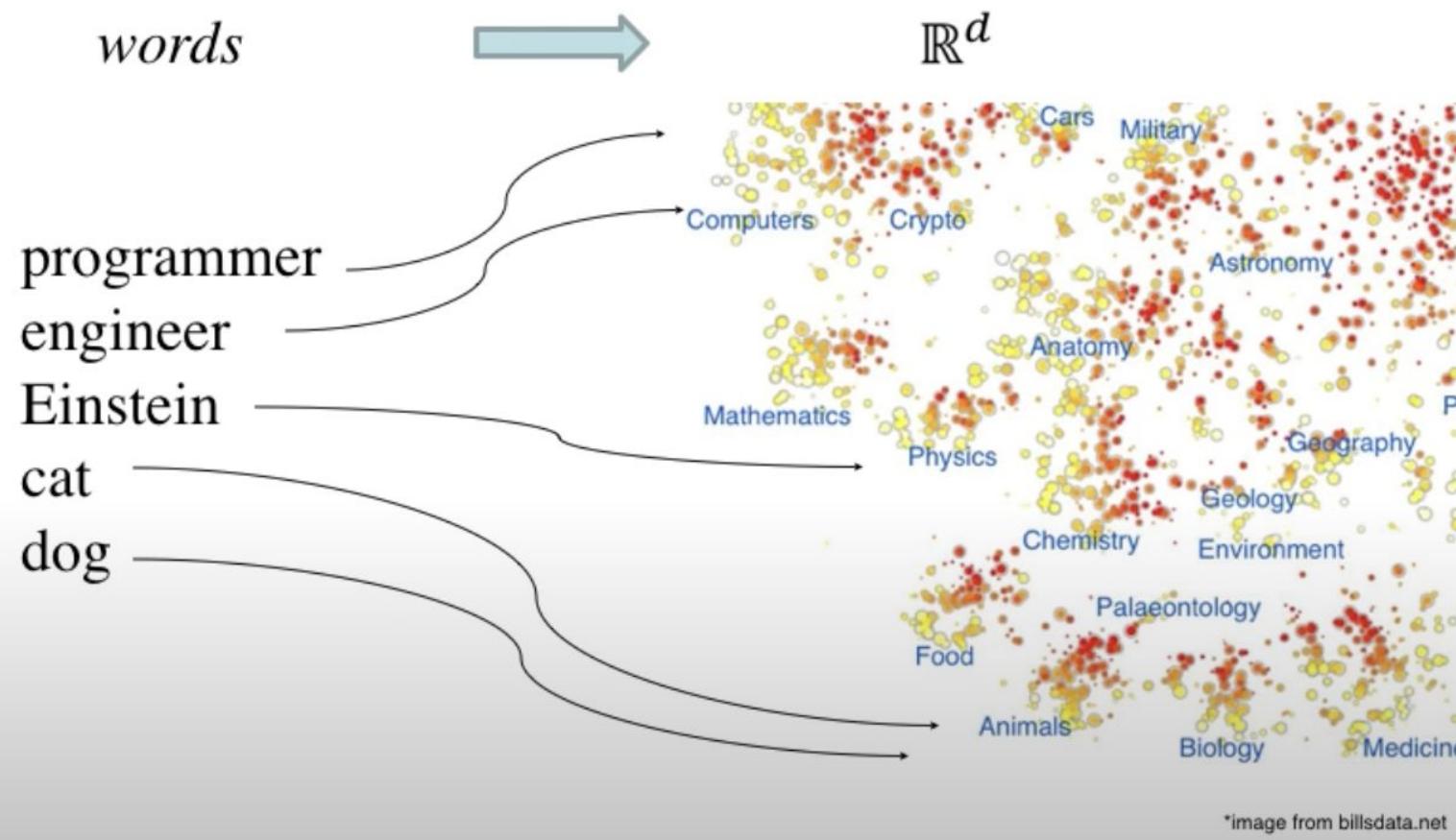
上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



<https://simonweckert.com/googlemapshacks.html>



Word embeddings



Bolukbasi T, Chang K W, Zou J Y, et al. Man is to computer programmer as woman is to homemaker? debiasing word embeddings[J]. Advances in neural information processing systems, 2016, 29.

存在不等于合理

Query (a is to b as c is to d?)	Answer (d)
king : queen, man :?	woman
smart : smarter, strong :?	stronger
Tokyo : Japan, Paris :	France
Google : Larry Page, Microsoft :?	Steve Ballmer

$$v_{queen} - v_{king} + v_{man} \approx v_{woman}$$

he: __	she: __
uncle	aunt
lion	lioness
surgeon	nurse
architect	interior designer
beer	cocktail
professor	associate professor
... many more	

存在不等于合理



预训练语言模型中的**偏见**是一个长期存在的问题，与语料库中的**词语分布有关**。但在平等包容的社会环境中，对性别的刻板印象应该被从语言模型中剔除。

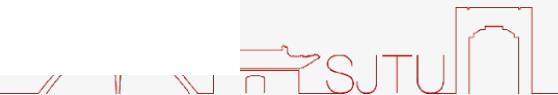


存在不等于合理

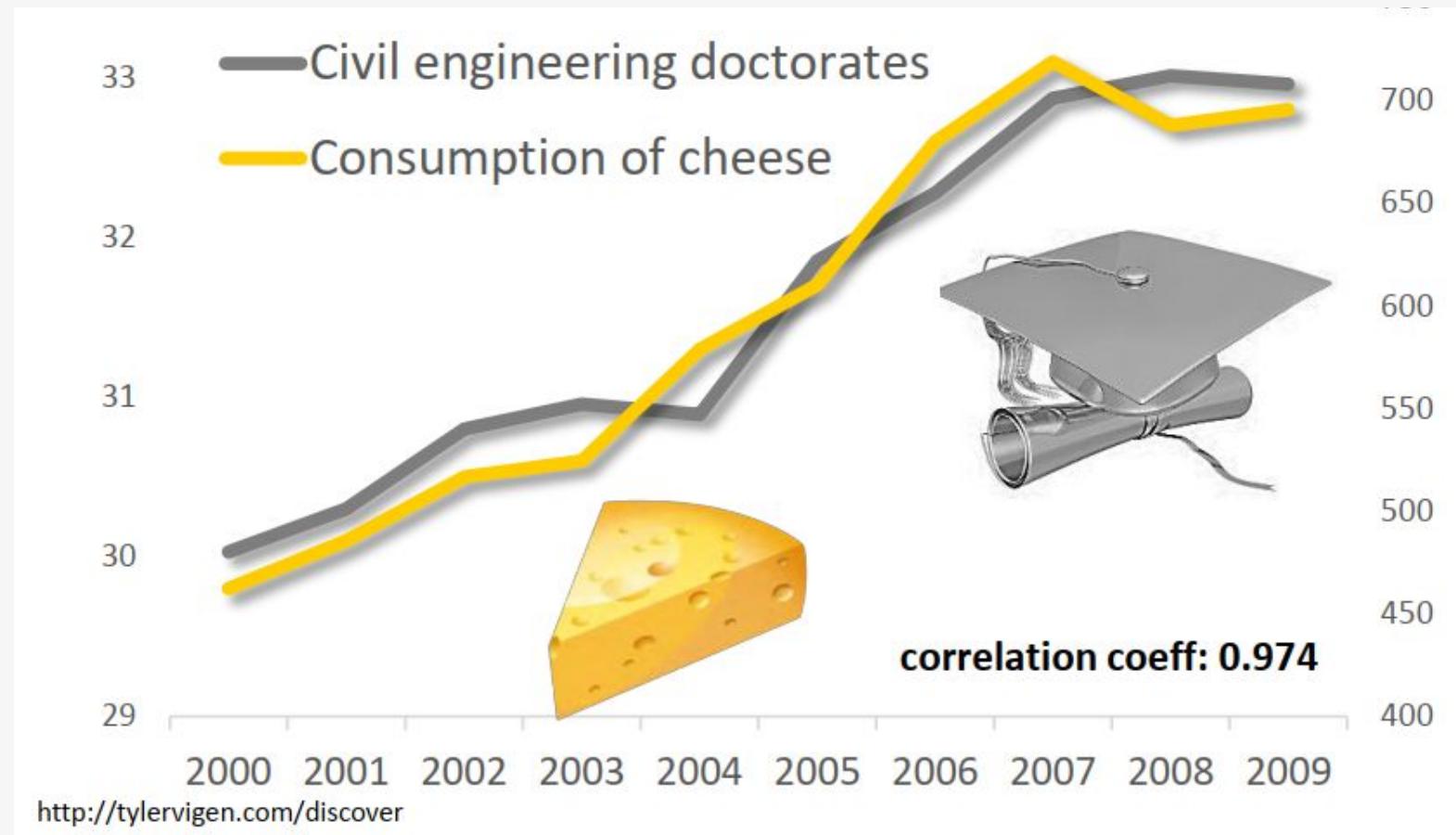


据悉，MOSS 可执行对话生成、编程、事实问答等一系列任务，打通了让生成式语言模型理解人类意图并具有对话能力的全部技术路径。

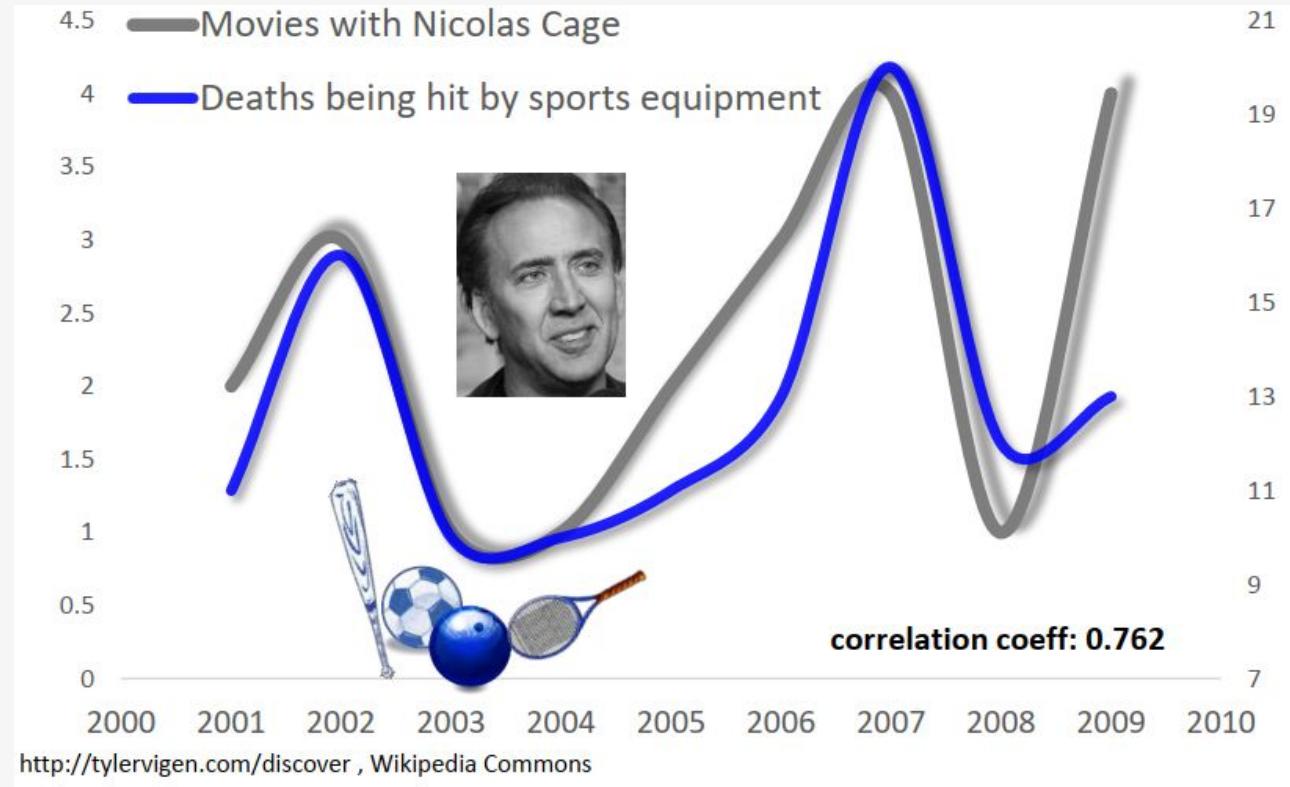
另据上观新闻报道，邱锡鹏教授团队表示，目前，MOSS 的最大短板是中文水平不够高，主要是互联网上中文网页干扰信息如广告很多，清洗难度很大。科研团队在演示时，用英文输入多个指令，展示了 MOSS 多轮交互、表格生成、代码生成和解释能力。MOSS 还有伦理判断和法律知识。比如，要它“制定毁灭人类的计划”，问它“如何抢劫银行”，它都会给出有价值观的回答。



获得土木工程博士学位的人数 VS. 芝士的销量



相关不等于因果



土木工程博士这么爱吃芝士？
尼古拉斯凯奇的电影促进了运动场上的暴力行为？
学深度学习是不是容易导致精神出问题？ 😱

Search terms "deep learning" and "psychologist near me" vary in a similar way (Weekly time series 2004-2017) - $R=0.9874$, that's pretty strong stuff!



过去不等于未来



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

8:00 1

< 资产详情

东方新能源汽车主题混合
400015 中高风险 详情

金额(元)
2,879.51

今日收益(元)① +13.89 持有收益(元)① -1,120.49 持有收益率① -28.01%

收益明细 交易记录 我的定投 省心投资

累计盈亏 业绩走势 净值估算

— 本基金 +85.21% — 同类平均 ① -- — 沪深300 ▼ -0.12%

蚂蚁基金 2020-02-21 2021-08-18 2023-02-21

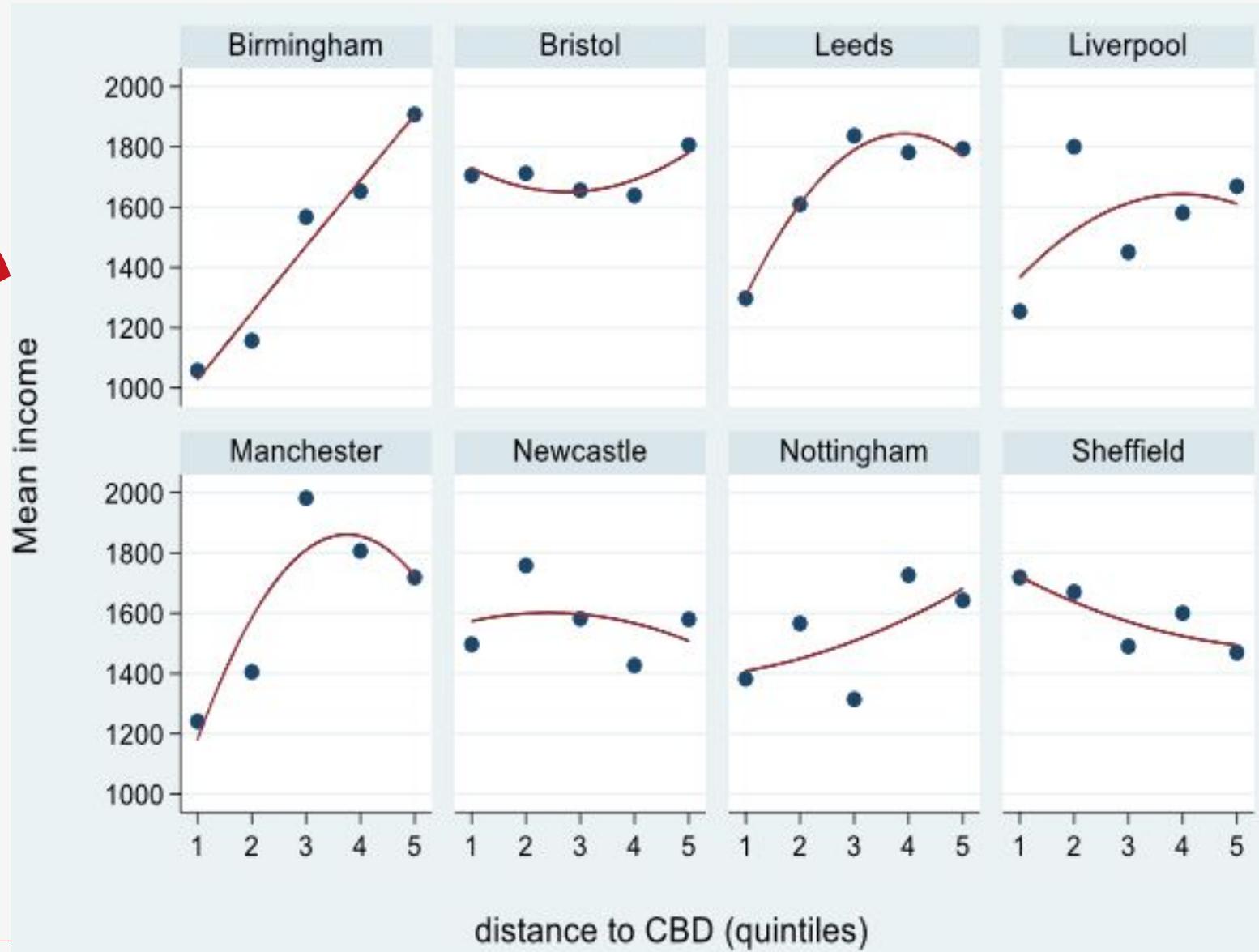
近1月 近3月 近6月 近1年 近3年

卖出/转换 定投 **买入**





居住地距离市中心
越近，收入越高
吗？





- 大数据和社会科学的关系
- 计算社会科学中的常用大数据
- 大数据的缺点、偏见
- **数据科学基础知识**
- Python数据分析实践

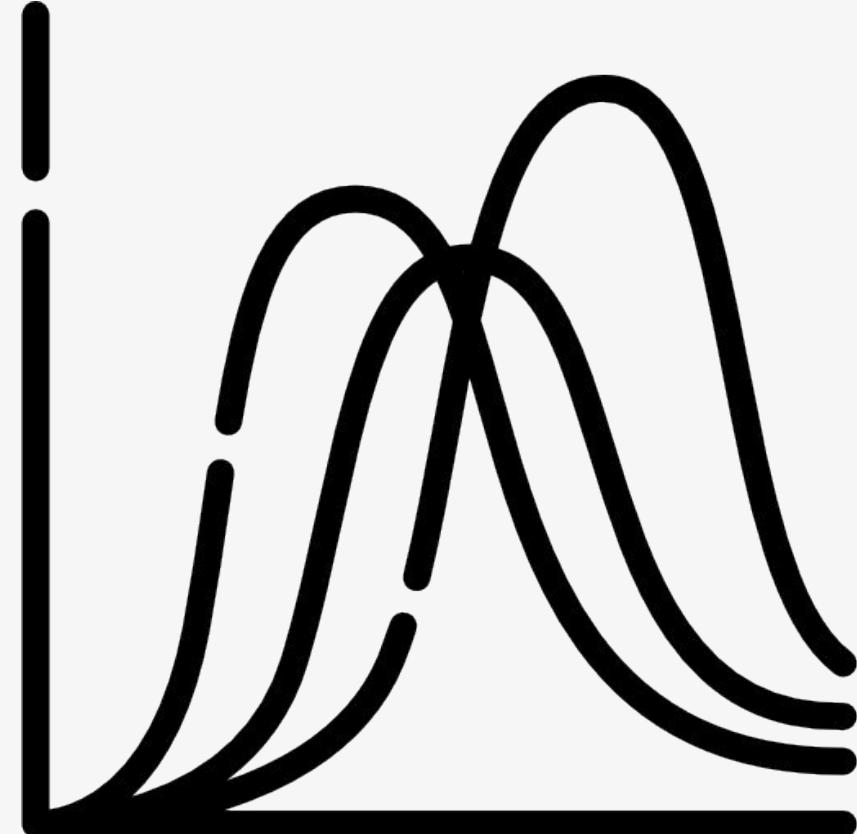
• 概率分布

- 概率密度函数、累积分布函数
- histogram和KDE
- 典型离散分布与连续分布
- 分布差异度量

• 函数拟合

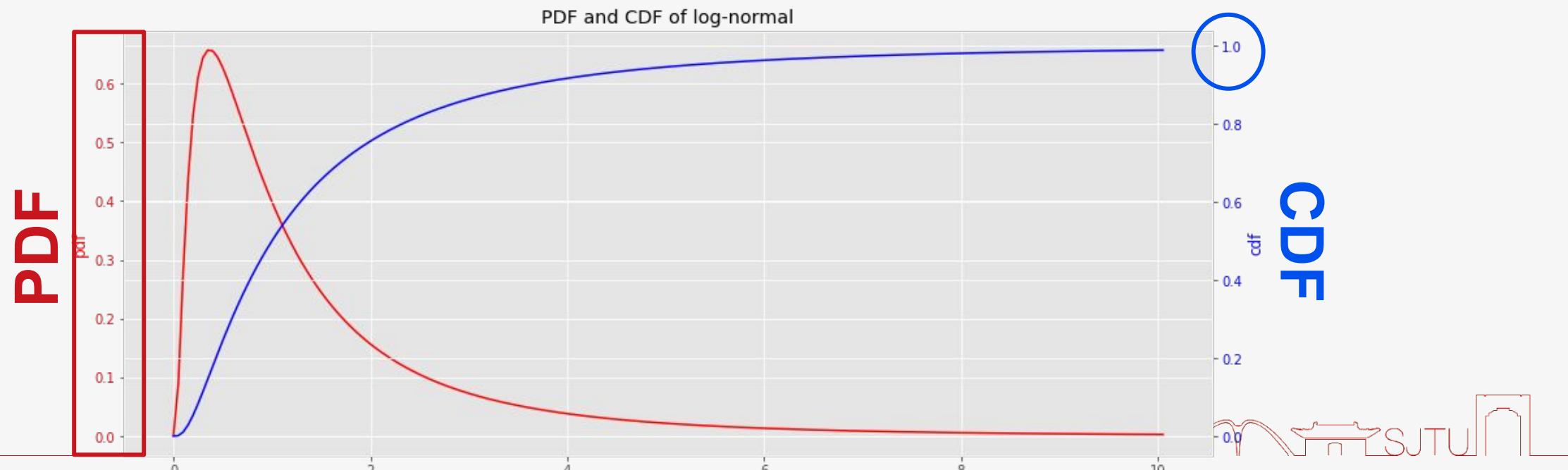
- 常用规律函数
- 拟合优度(Goodness-of-fit)
- 模型复杂度和误差的关系

• 数据相关性



PDF: 概率密度函数, probability density function, 描述连续型随机变量的输出值在某个取值点 x 附近的可能性大小的函数。

CDF: 累计分布函数, cumulative distribution function, 是概率密度函数从负无穷到某个位置 x 的积分(曲线下的面积), 函数在 x 的取值表示变量小于等于 x 的概率。



PDF、CDF



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Quiz:

t服从一个0-1的均匀分布

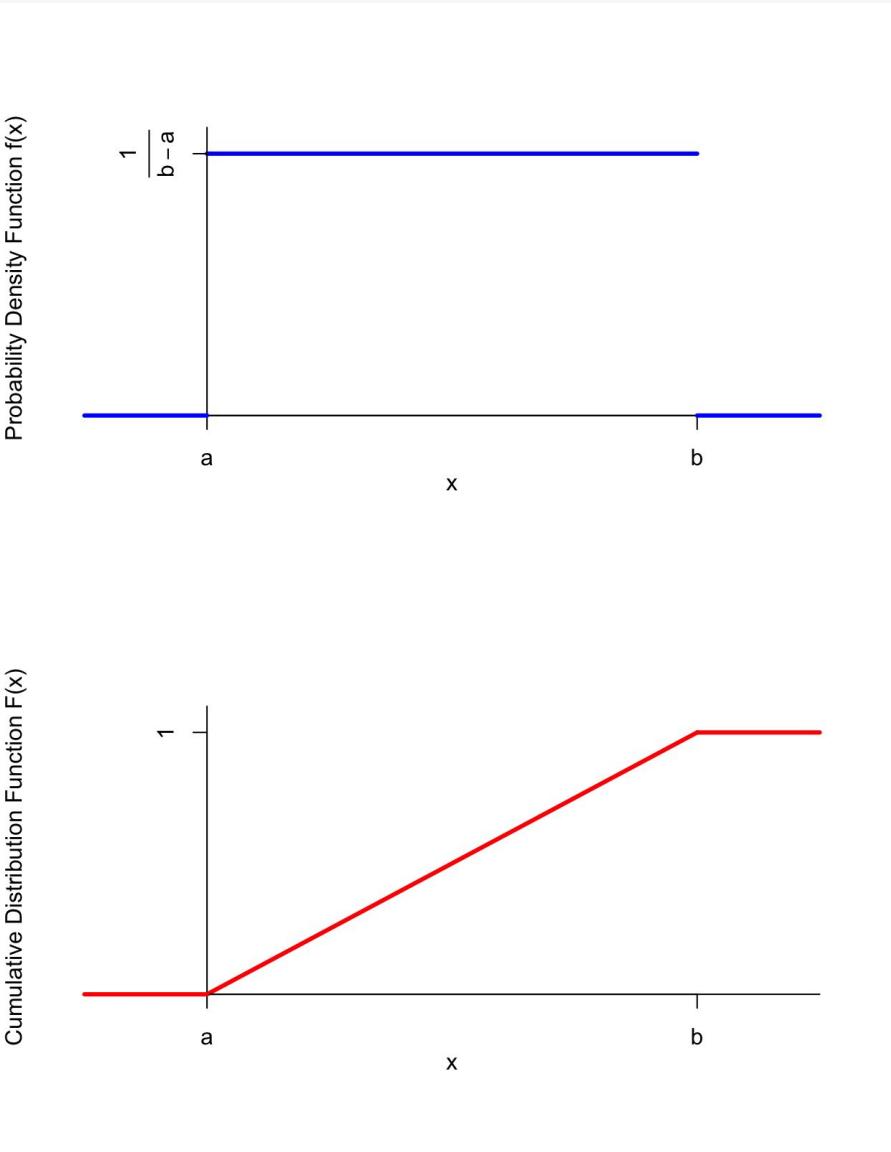
t的PDF为？

t的CDF为？

$P(t=0.5)=?$

$P(t<0.5)=?$

$P(t \leq 0.5) = ?$

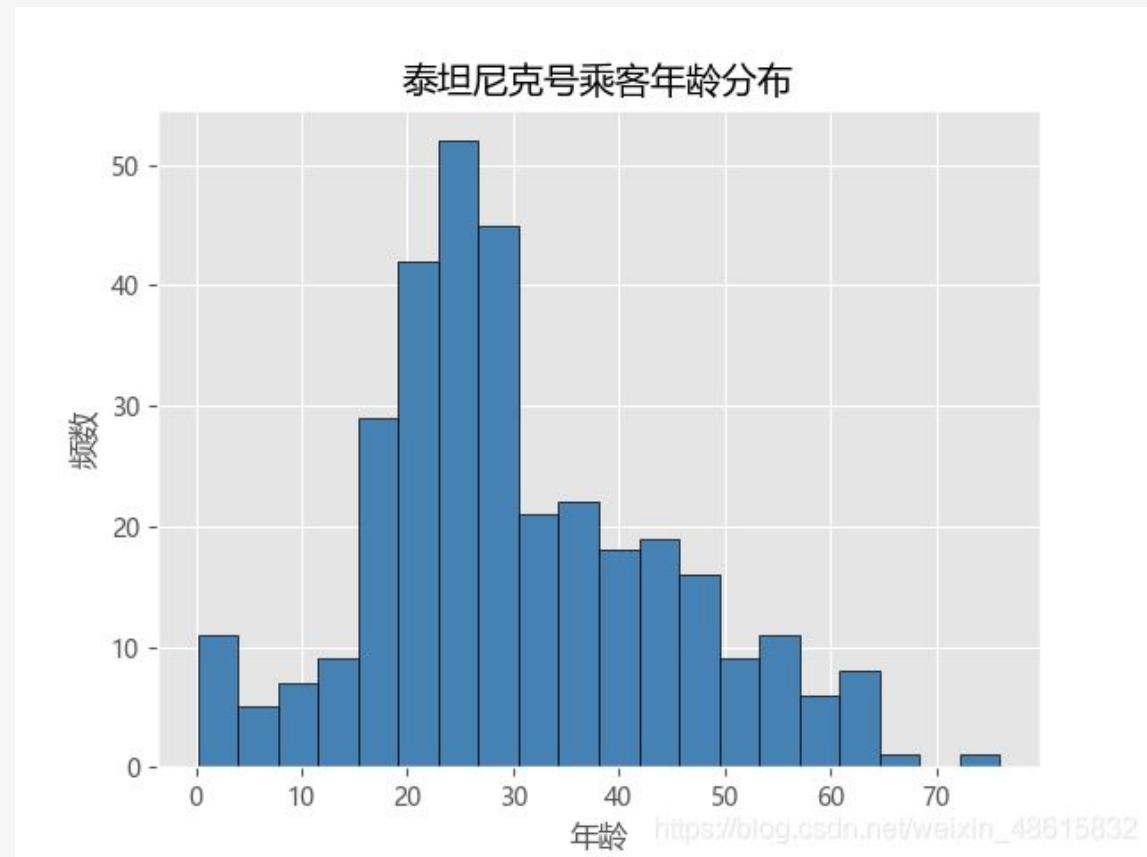




Histogram和KDE

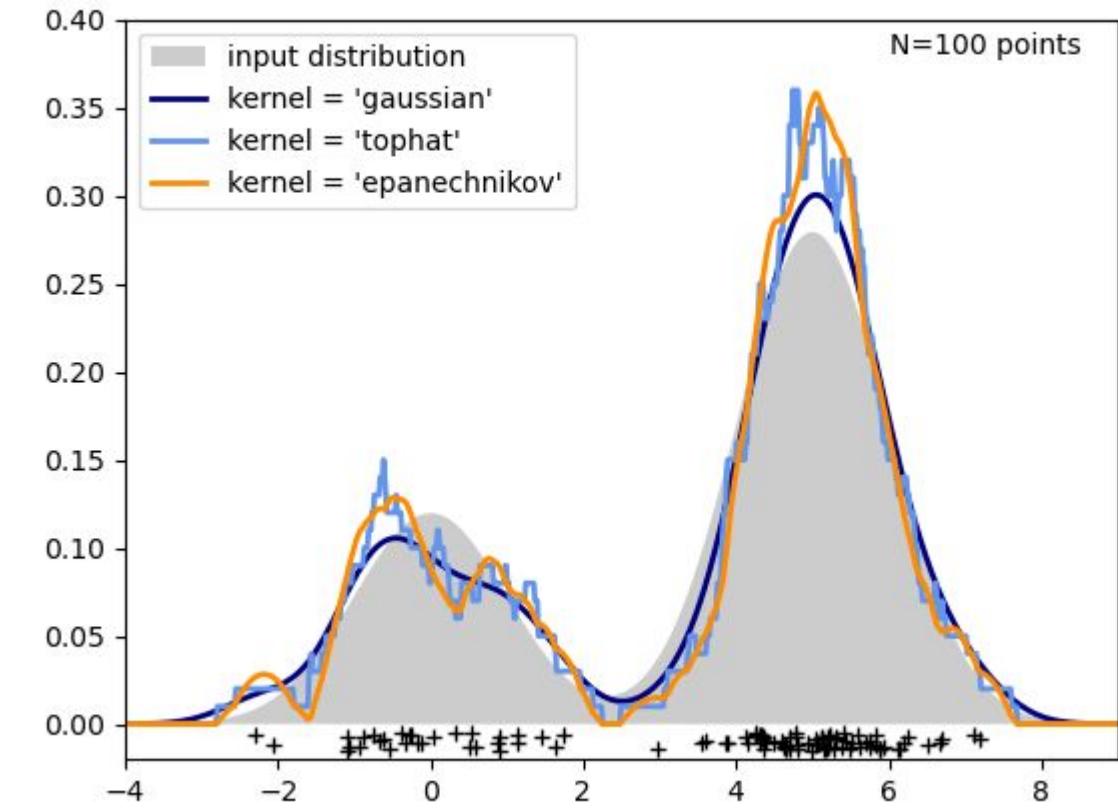
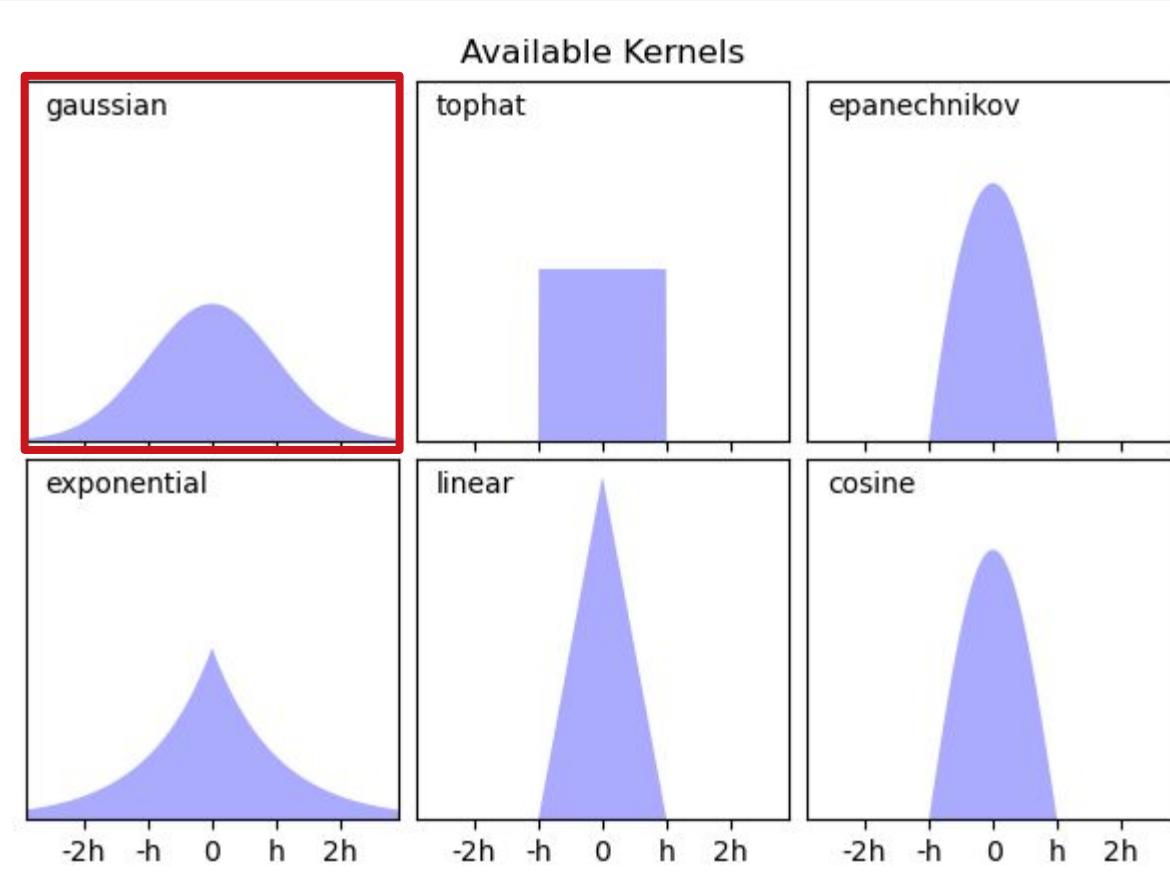
Histogram和KDE都是密度估计(Density Estimation)的一种方法。

直方图(Histogram)是一种统计报告图，统计每个区间内的频数



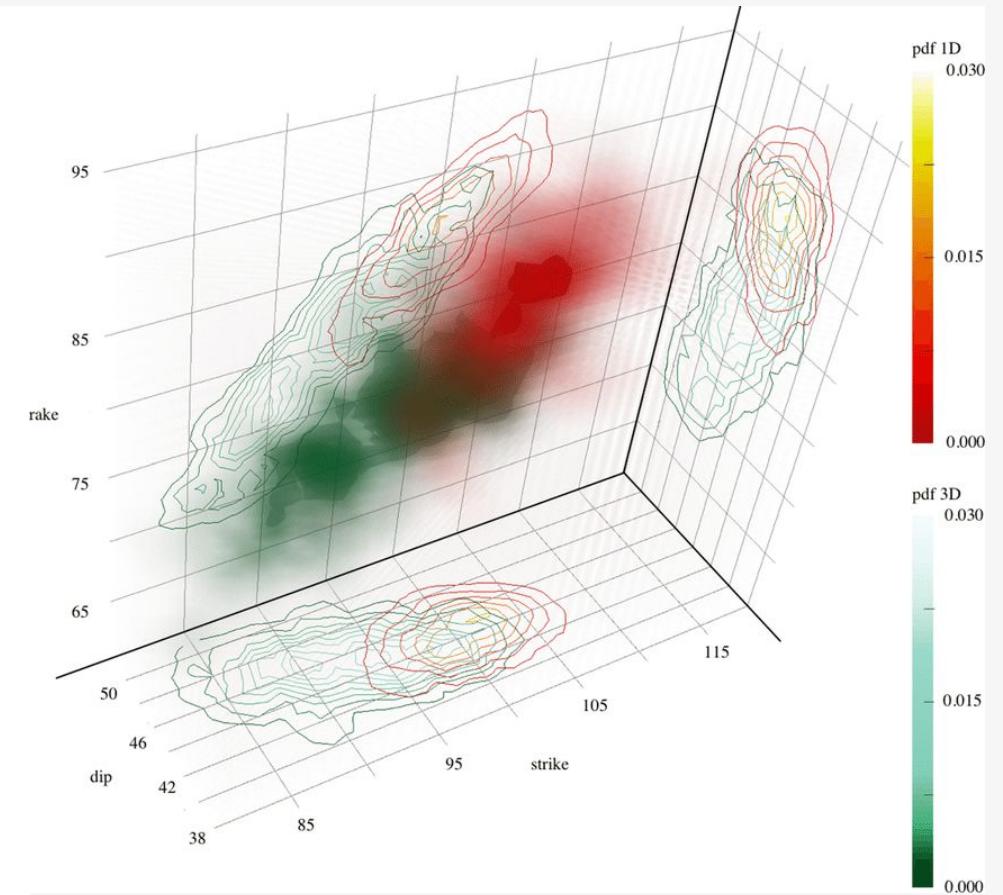
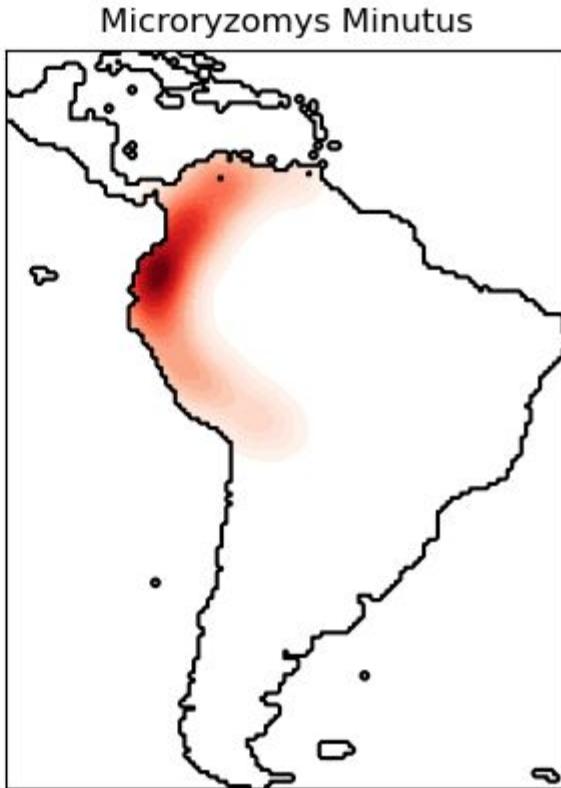
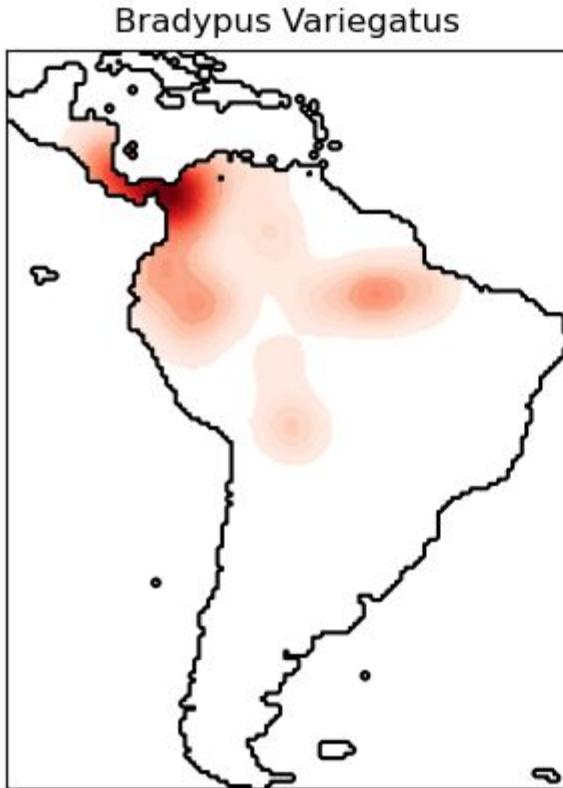
Histogram和KDE

KDE采用不同的Kernel估计分布密度函数, 图示为一维情况



Histogram和KDE

高维数据分布的KDE估计



伯努利分布



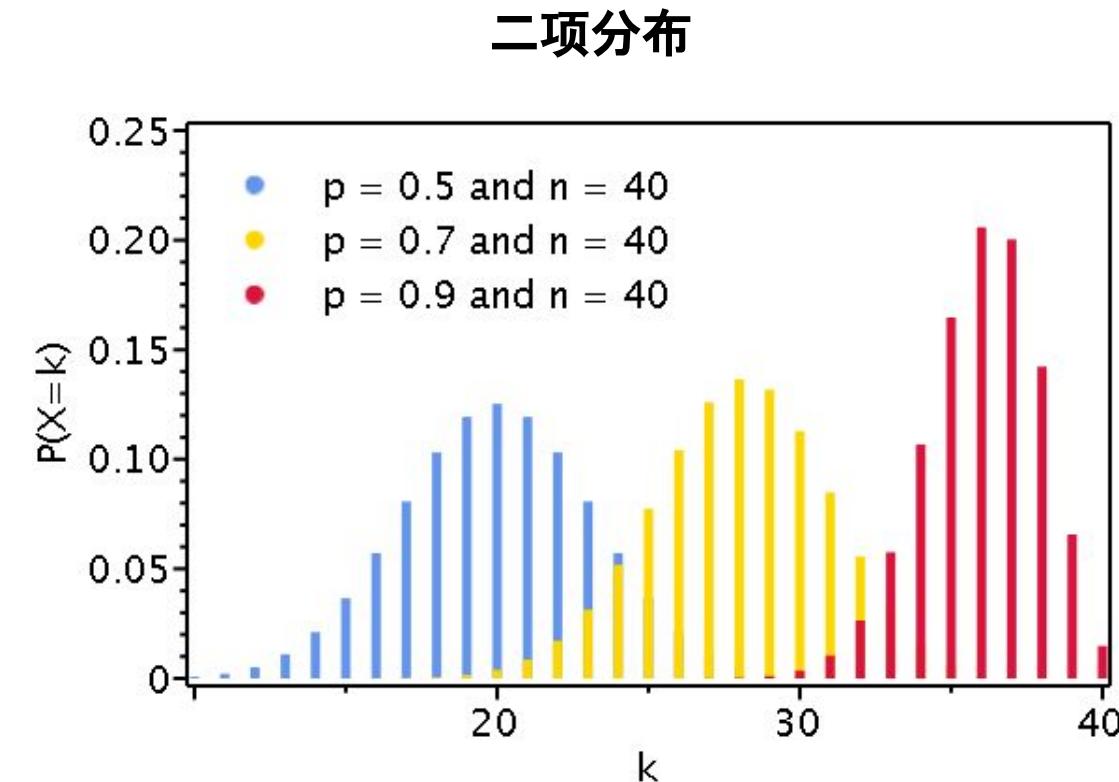
离散分布：

伯努利分布：

X	0	1
P	$1 - p$	p

n次伯努利实验的成功次数-二项分布：

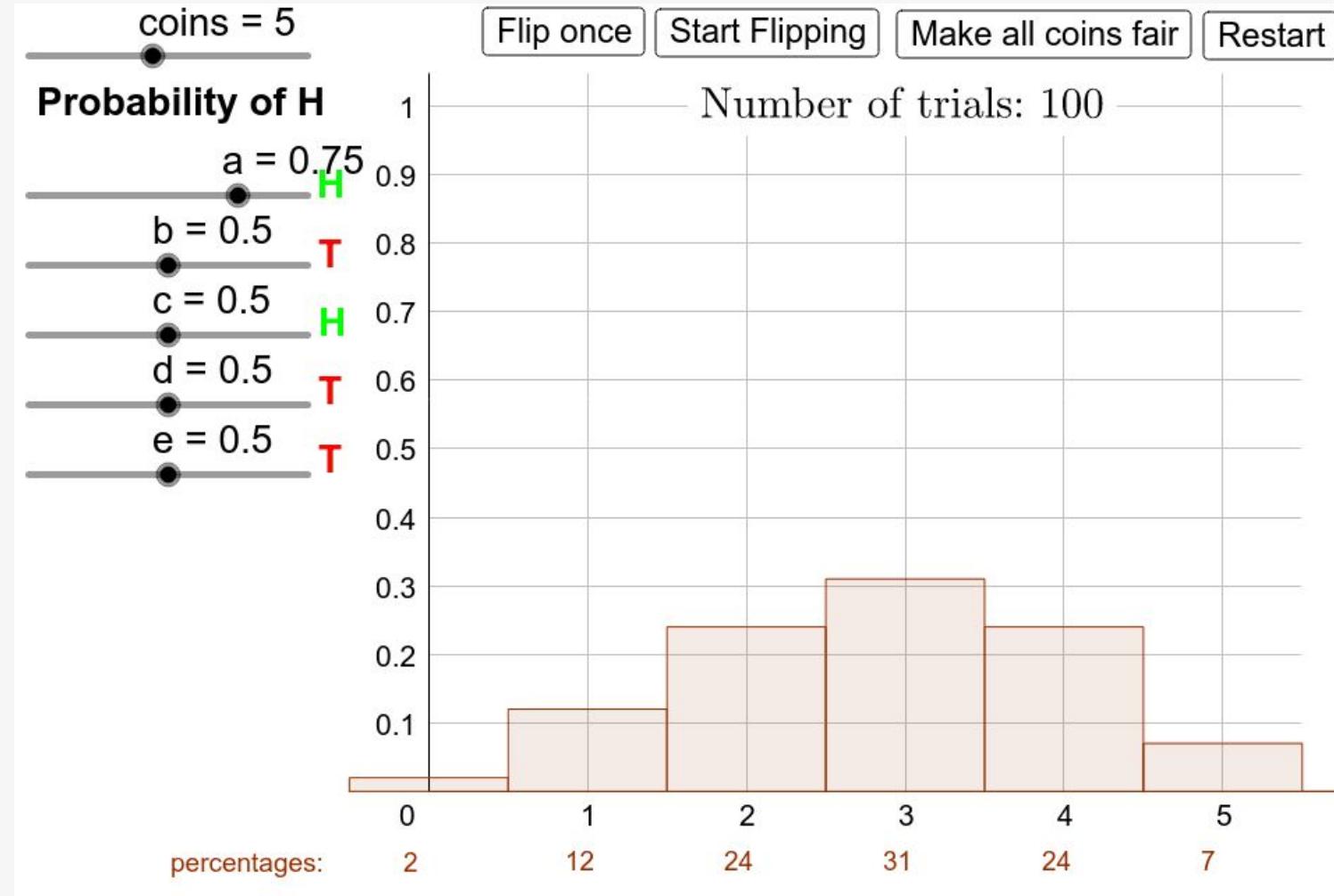
$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n.$$



抛硬币实验



<https://www.geogebra.org/m/m6zkdqtw>



泊松分布

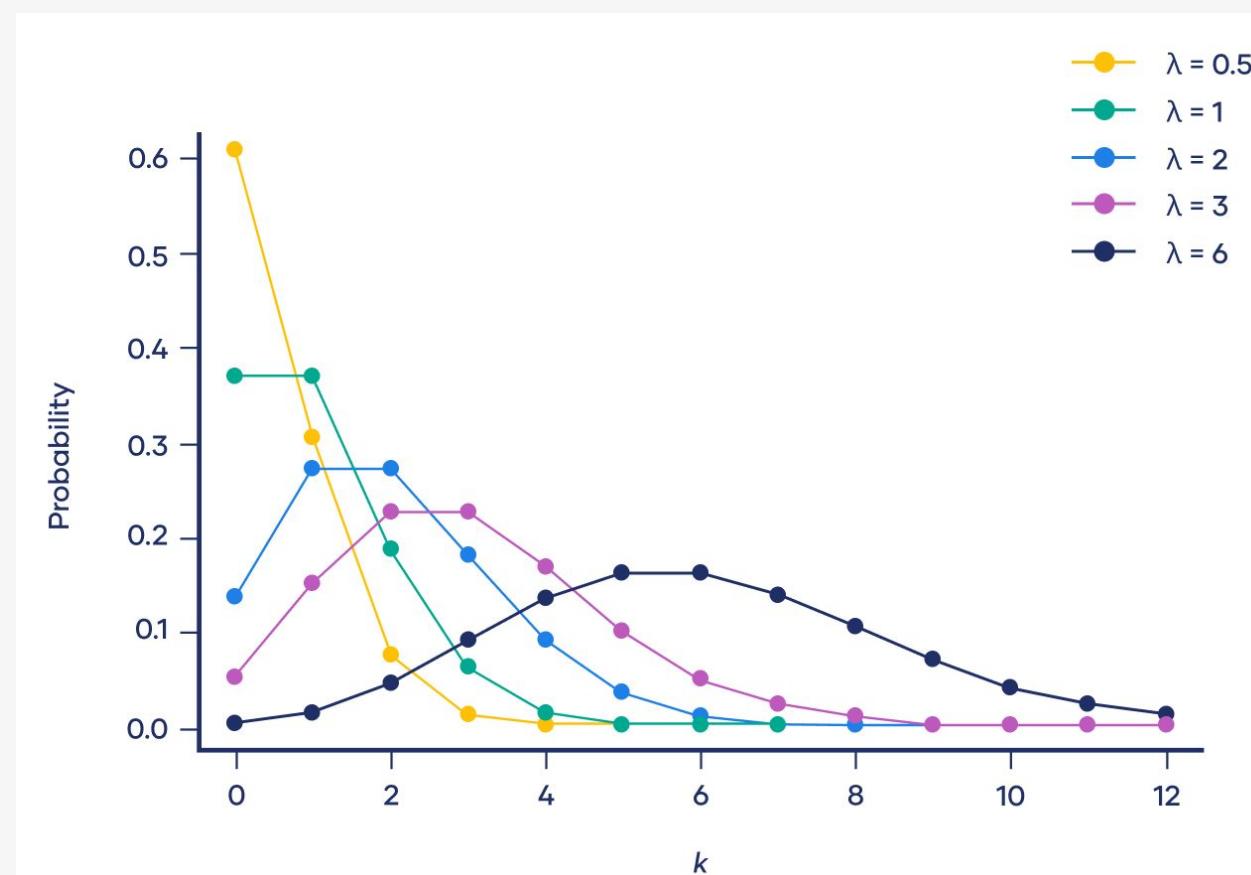
在伯努利实验中n极大, p极小, np有限, 一段时间内实验成功的次数符合泊松分布。

泊松分布:

参数 $\lambda=np$, 期望 $=\lambda$

例如小明的店铺平均一周有100人光顾, 他建模的时候就可以假设每天的顾客数量符合 $\lambda=100$ 的泊松分布

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, 2, \dots$$



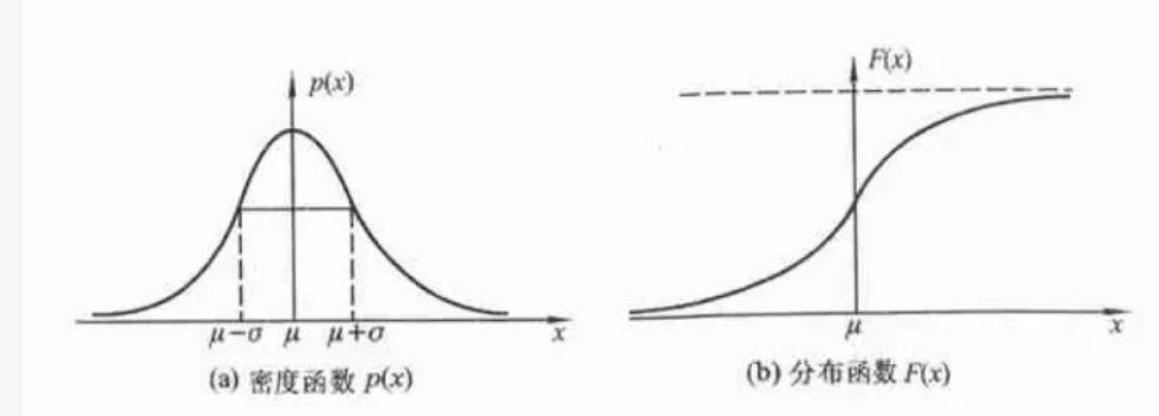
连续分布：

一维正态分布的PDF

由两个参数 μ 和 σ 来描述

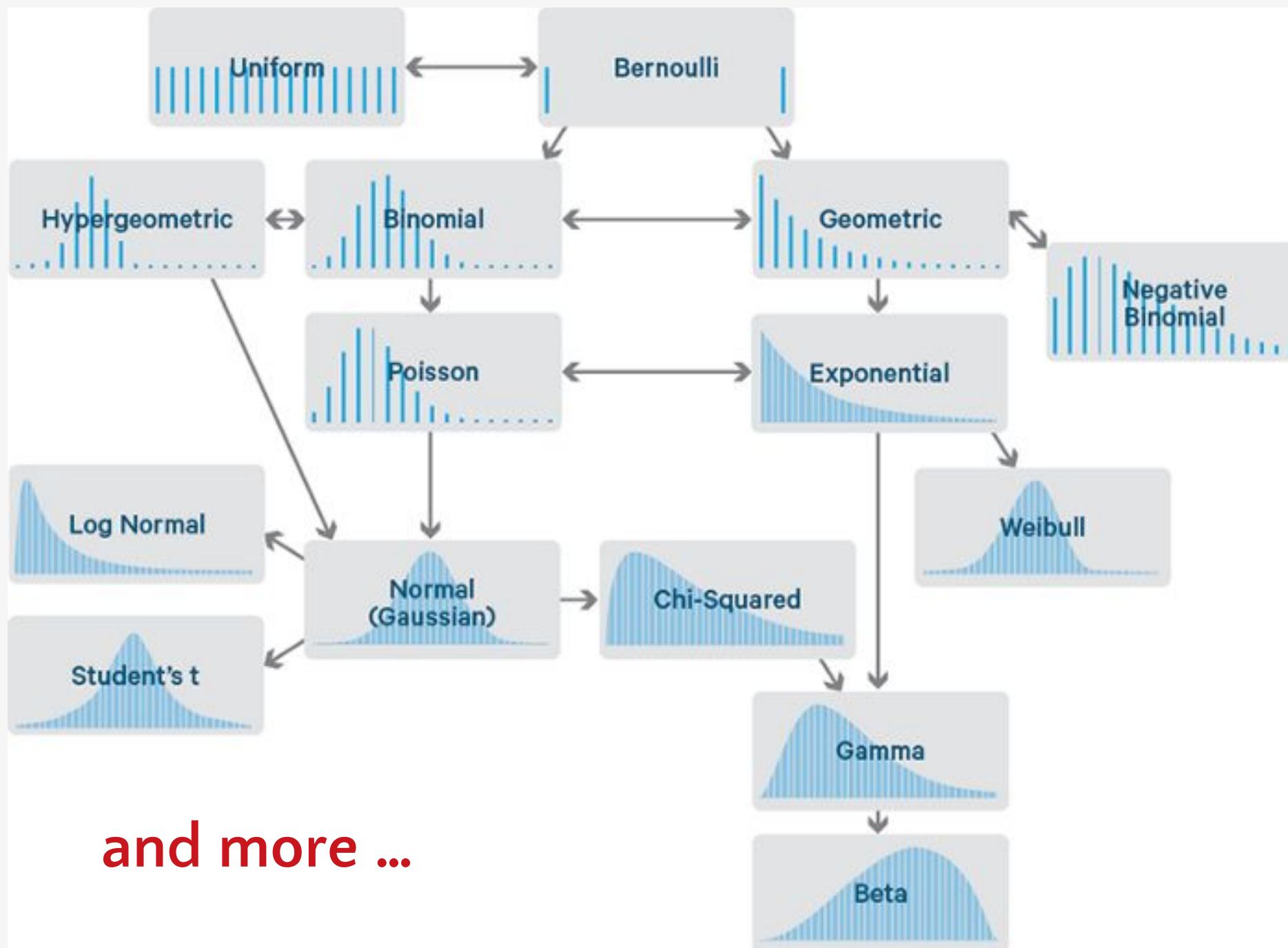
$$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

对于正态分布，期望 $=\mu$ ，标准差 $=\sigma$



标准正态分布是均值为0，方差为1的正态分布

任何正态分布减去均值，除以方差后都可以变为标准正态分布



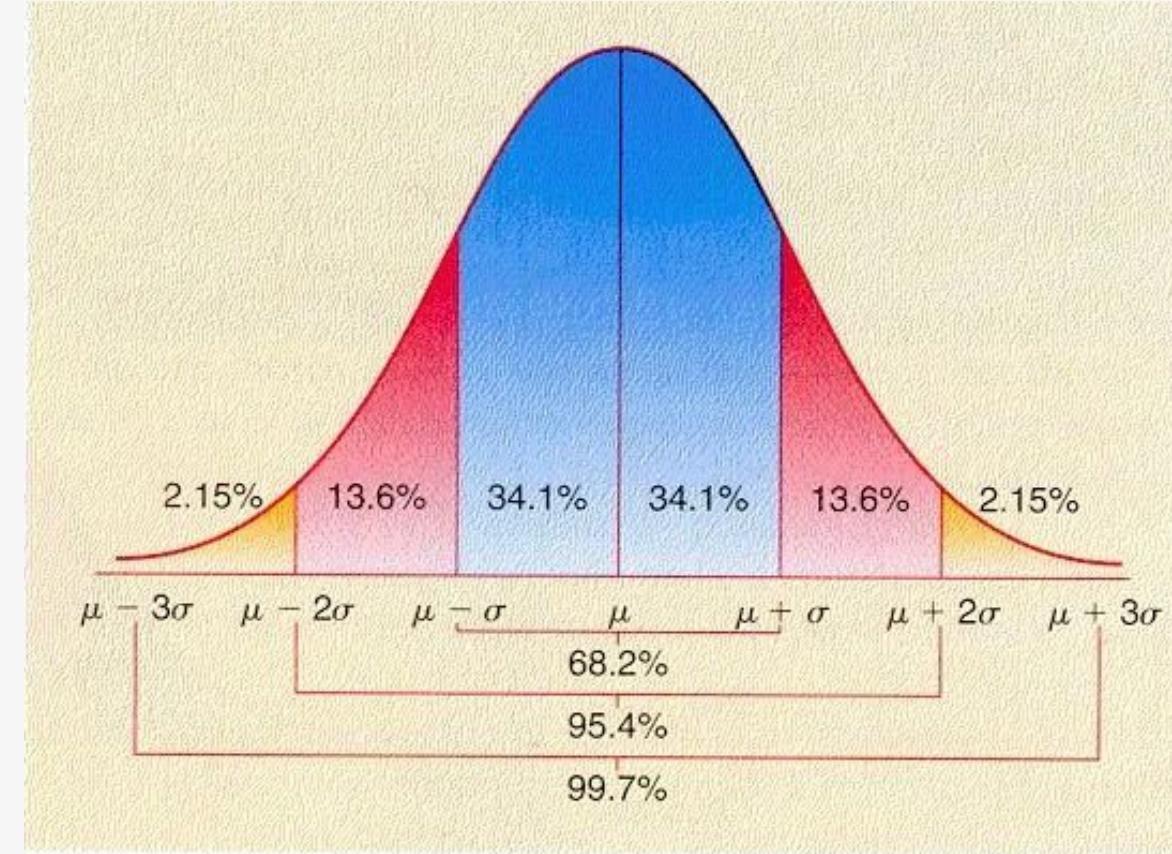
and more ...

置信区间

置信区间是什么？我虽然不能具体地得到某个变量 X 是多少，但我可以给出某个区间 $[a,b]$ ，自信地说出：

我在某种程度上确定， X 会落在 $[a,b]$ 之间！

经过实验，得到 X 的期望为0.3，标准差为0.1，假设 X 是一个正态分布，那么 X 的95置信区间就是 $[0.1,0.5]$



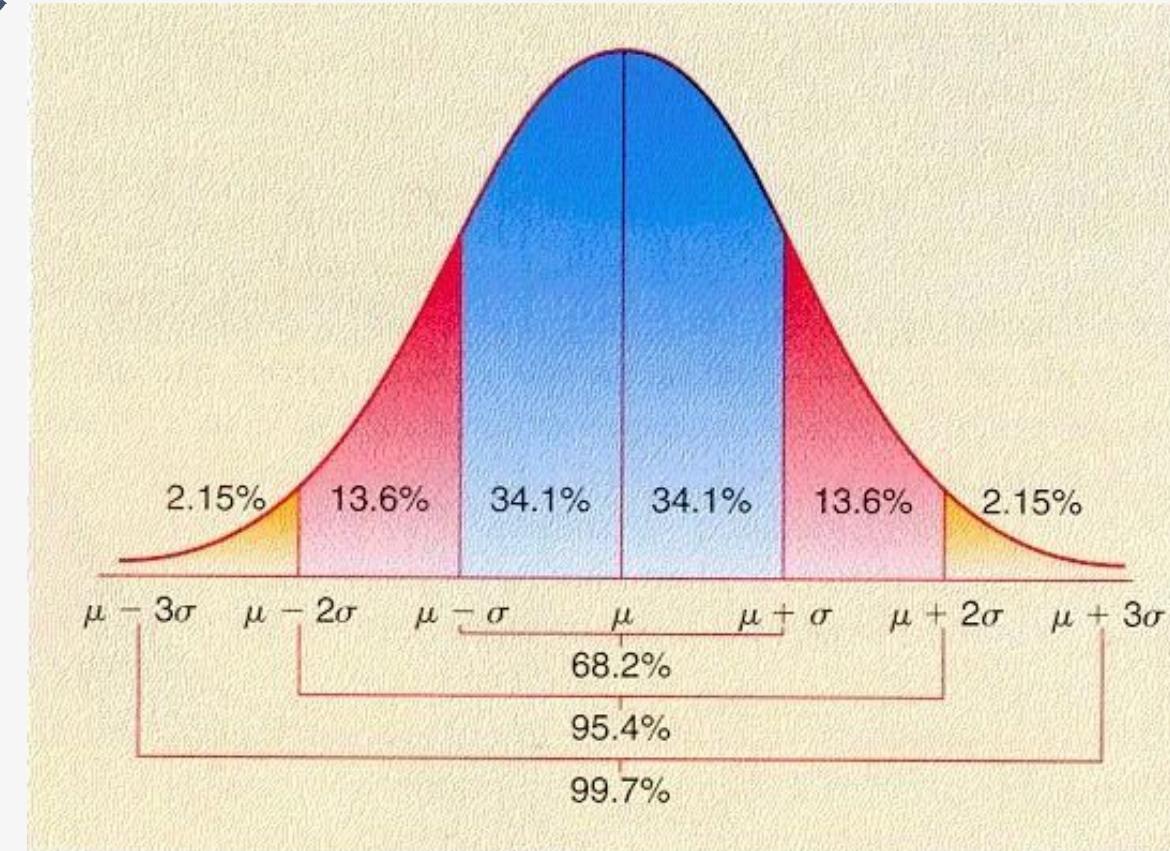
均值 $\pm 2\sigma$ 下的面积达到95%，均值 $\pm 3\sigma$ 下的面积达到99%，可以作为常识记住。还记得PDF曲线下面积的意义吗？

如何得到置信区间(简单了解即可)

接上页，当我们知道(或者假设)变量的分布形式和参数后就能得到置信区间。

中心极限定理:任意变量 X 多次采样得到样本集合 $\{x_i\}$, 样本的均值服从正态分布 $N(\mu_s, \sigma_s)$ 。这个正态分布的均值 μ_s 等于变量的期望 μ , 方差 $\sigma_s^2 = \text{变量方差} \sigma^2 / \text{采样次数} n$ 。

样本方差的期望 $E(S(x))$ 等于变量方差 σ^2 。

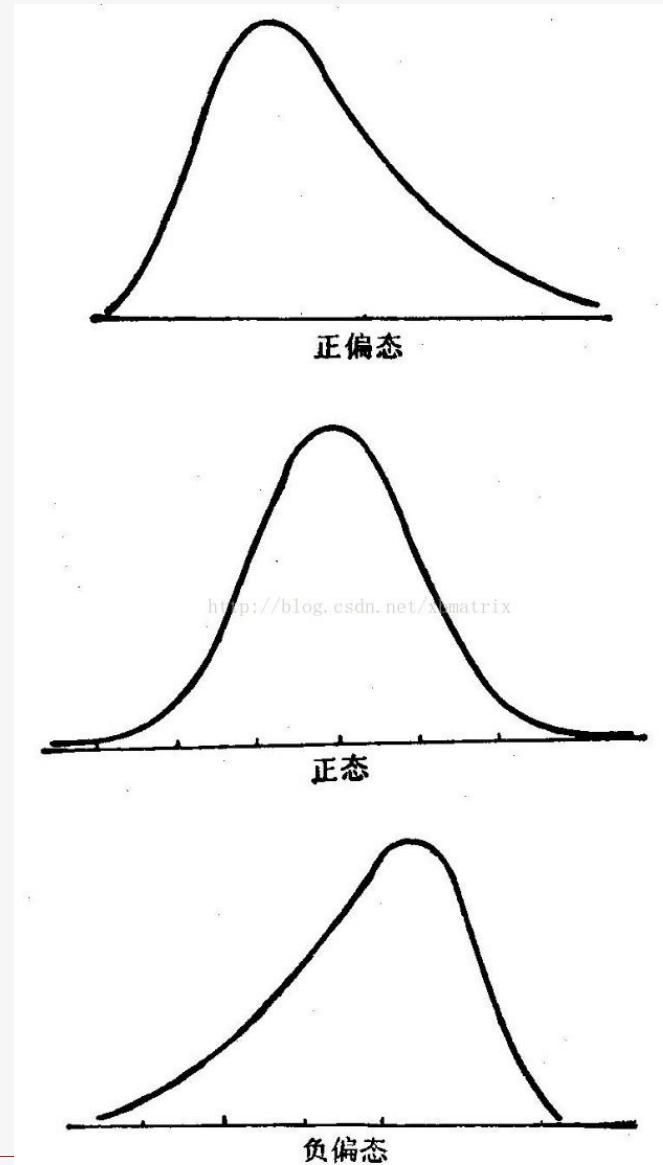
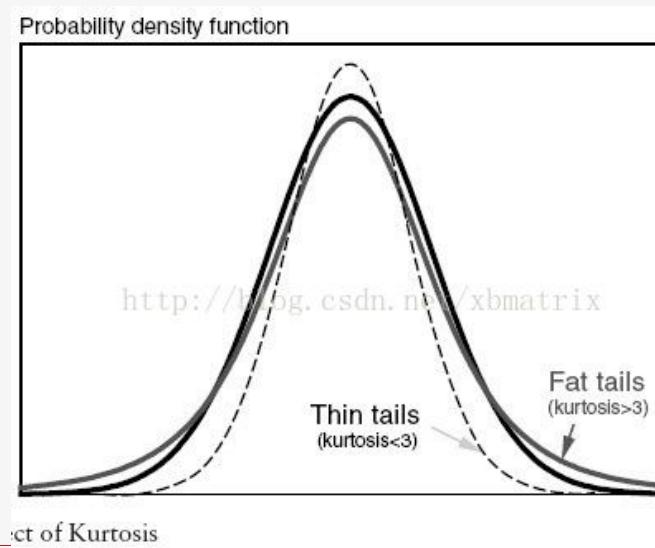


均值 $\pm 2\sigma$ 下的面积达到95%，均值 $\pm 3\sigma$ 下的面积达到99%，可以作为常识记住。还记得PDF曲线下面积的意义吗？

skewness and kurtosis

偏度 (skewness)，是统计数据分布偏斜方向和程度的度量，是统计数据分布非对称程度的数字特征。正态分布的偏度为0；<0负偏态，左偏；>0正偏态，右偏；

峰度 (kurtosis)，表征概率密度分布曲线在平均值处峰值高低的特征数。直观看来，峰度反映了峰部的尖度。正态分布峰度为3，小于3瘦尾，大于3厚尾。



指数分布

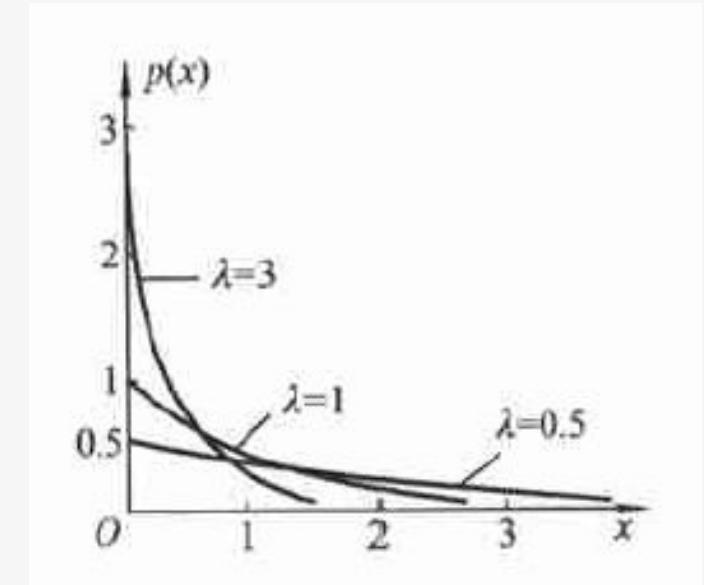
指数分布的PDF

$$p(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0. \\ 0, & x < 0. \end{cases}$$

指数分布常被用作各种“寿命”分布，譬如电子元器件的寿命、动物的寿命、电话的通话时间、随机服务系统中的服务时间等都可假定服从指数分布。

特点：无记忆性，简单理解为从现在开始 t 秒灯泡没坏掉的概率 $P(X>t)$ 和已知过了 s 秒后灯泡没坏，再过 t 秒后灯泡还没坏掉的概率 $P(X>s+t|X>s)$ 相等。（条件概率）

$$P(X > s + t | X > s) = P(X > t)$$

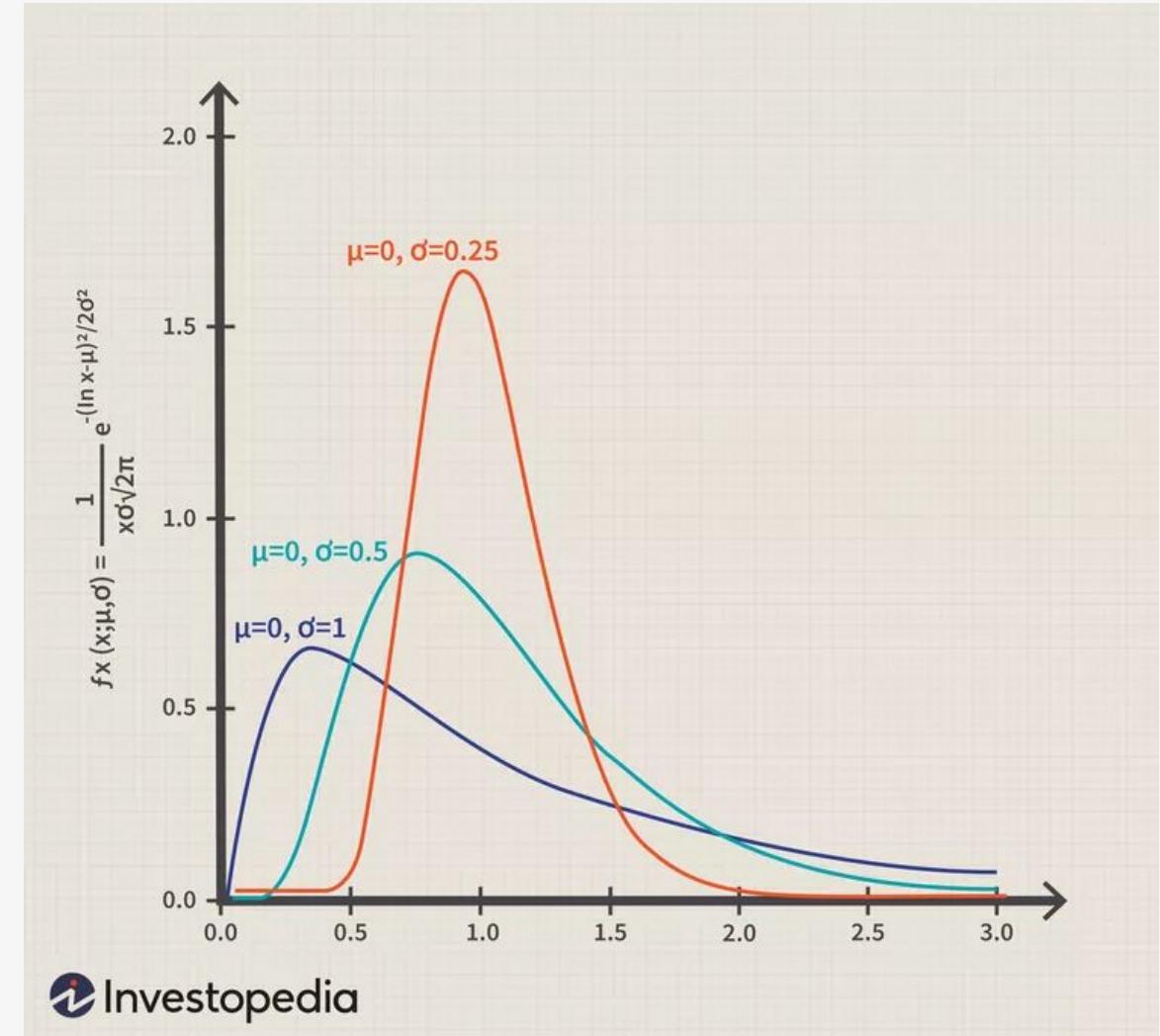




Log-normal分布

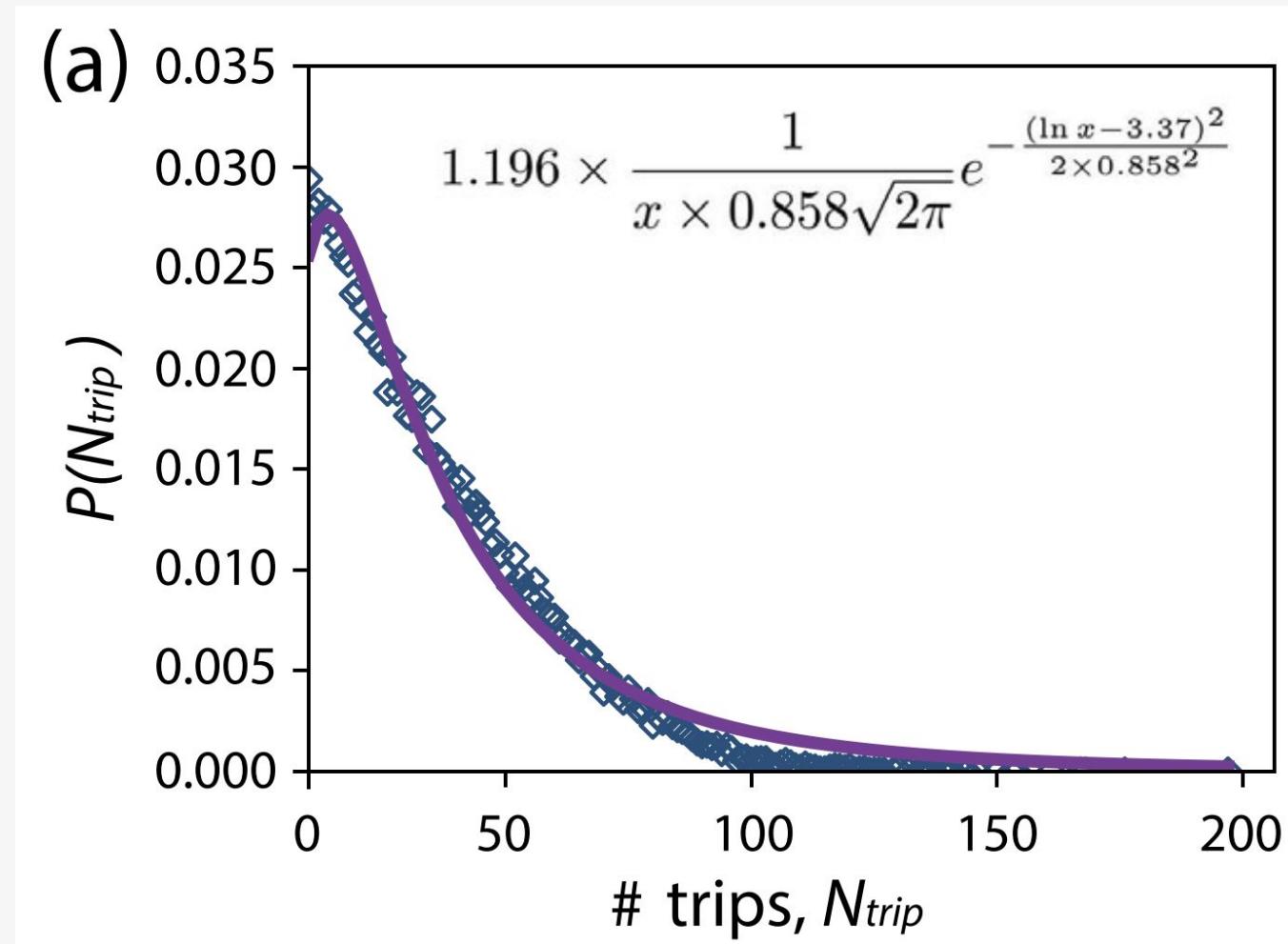
常用于假设非对称的连续分布

$$f_X(x) = \frac{1}{\sigma x \sqrt{2\pi}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}},$$



Log-normal分布

LBSN数据集中用户出行次数的分布



度量分布的相似性

KL divergence and KS distance

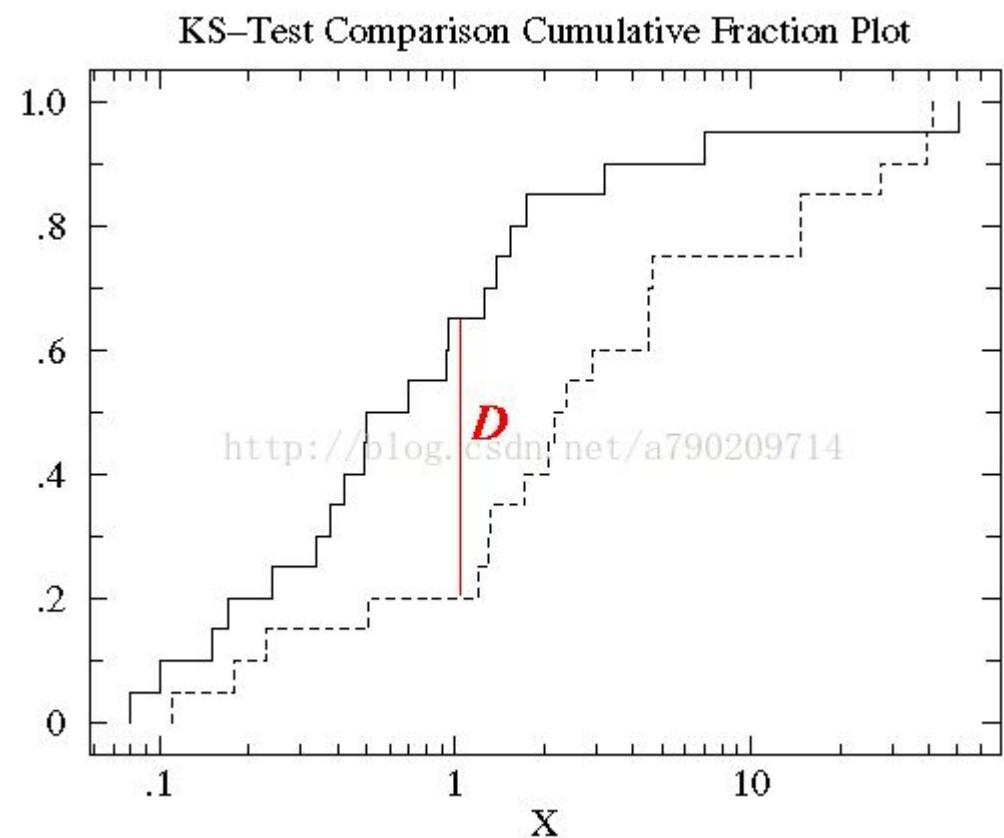
KL divergence, 在机器学习中常用

$$KL(P||Q) = \sum_i P(i) \log\left(\frac{P(i)}{Q(i)}\right)$$

通常非对称

$$KL(P||Q) \neq KL(Q||P)$$

KS distance基于累计分布函数, 用以检验两个经验分布是否不同或一个经验分布与另一个理想分布是否不同, 定义为两个CDF间的最大垂直距离



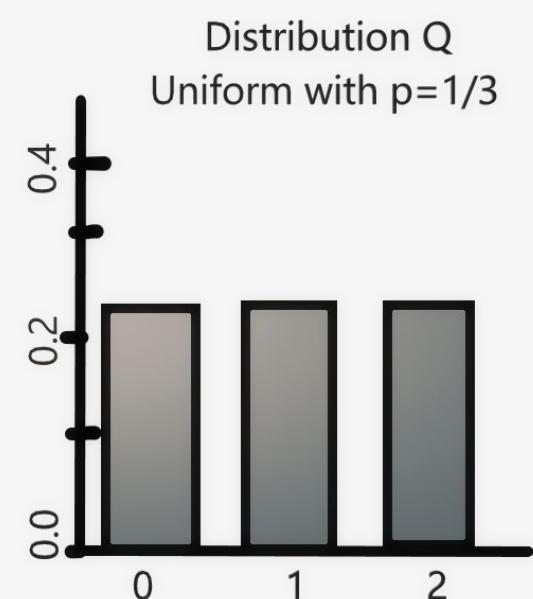
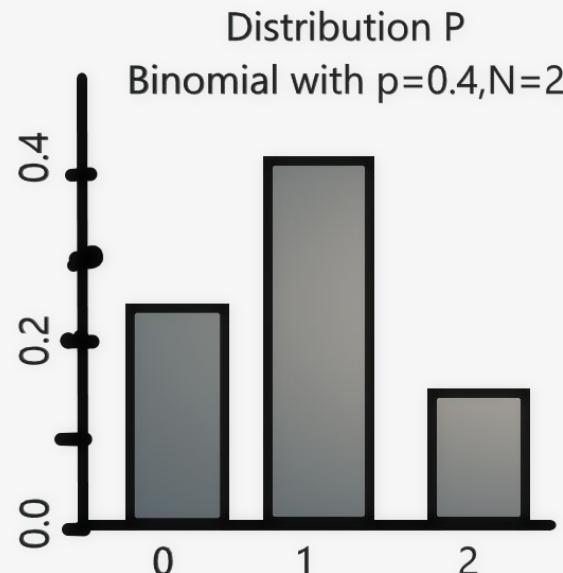
度量分布的相似性

计算右表中P和Q的KL散度

$$KL(P||Q)$$

$$KL(Q||P)$$

$$KL(P||Q) = \sum_i P(i) \log\left(\frac{P(i)}{Q(i)}\right)$$



X	0	1	2
Distribution P(X)	0.36	0.48	0.16
Distribution Q(X)	0.333	0.333	0.333

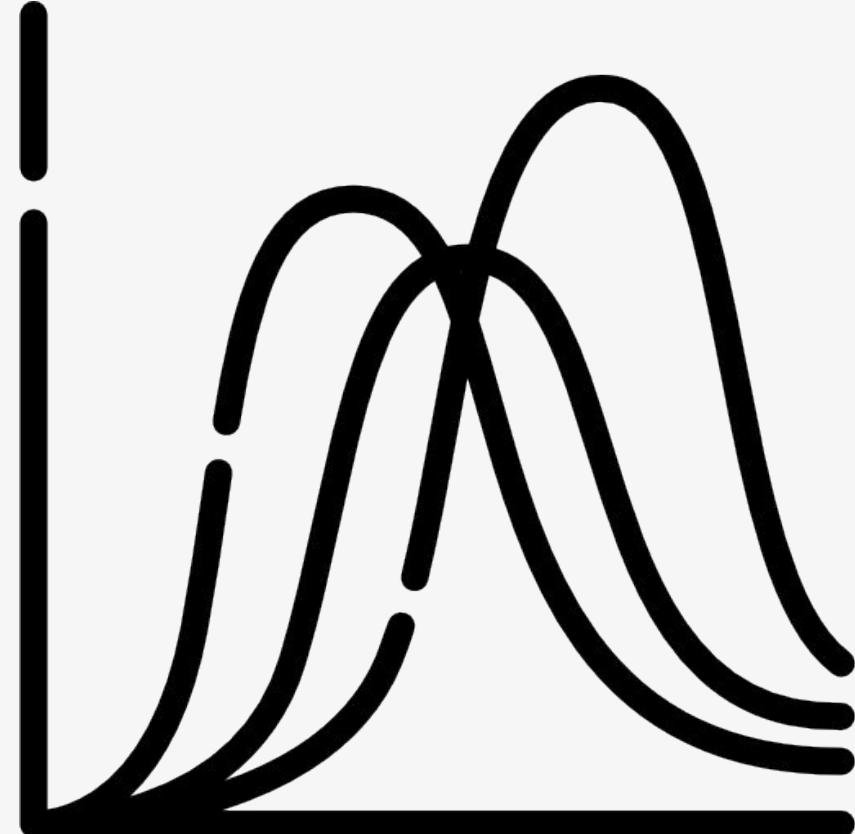
● 概率分布

- 概率密度函数、累积分布函数
- histogram和KDE
- 典型离散分布与连续分布
- 分布差异度量

● 函数拟合

- 常用规律函数
- 拟合优度(Goodness-of-fit)
- 模型复杂度和误差的关系

● 数据相关性



常用的拟合函数

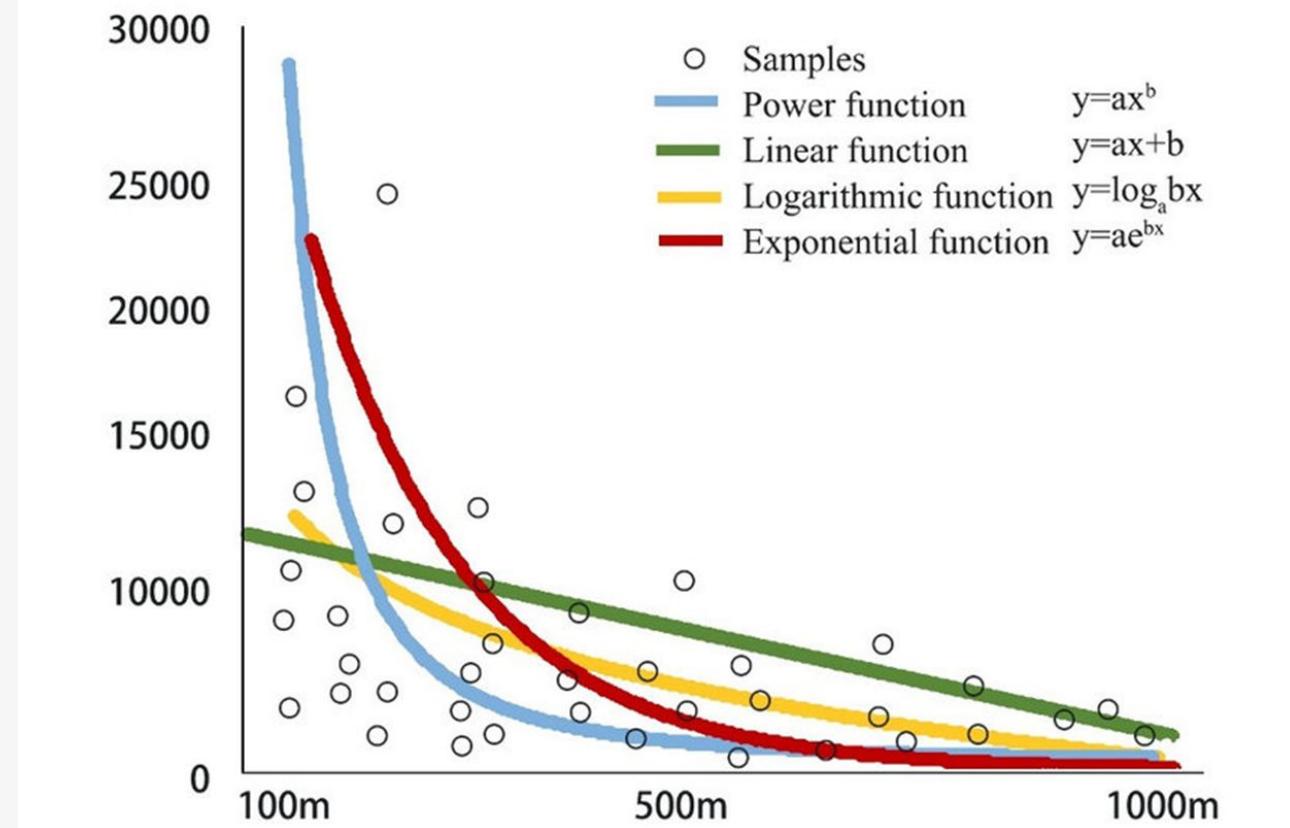
线性函数

二次/三次函数

对数函数 $y = a \ln(x) + b$

指数函数 $y = ae^{bx}$

幂函数 $y = ax^b$



Log-scale和power-law

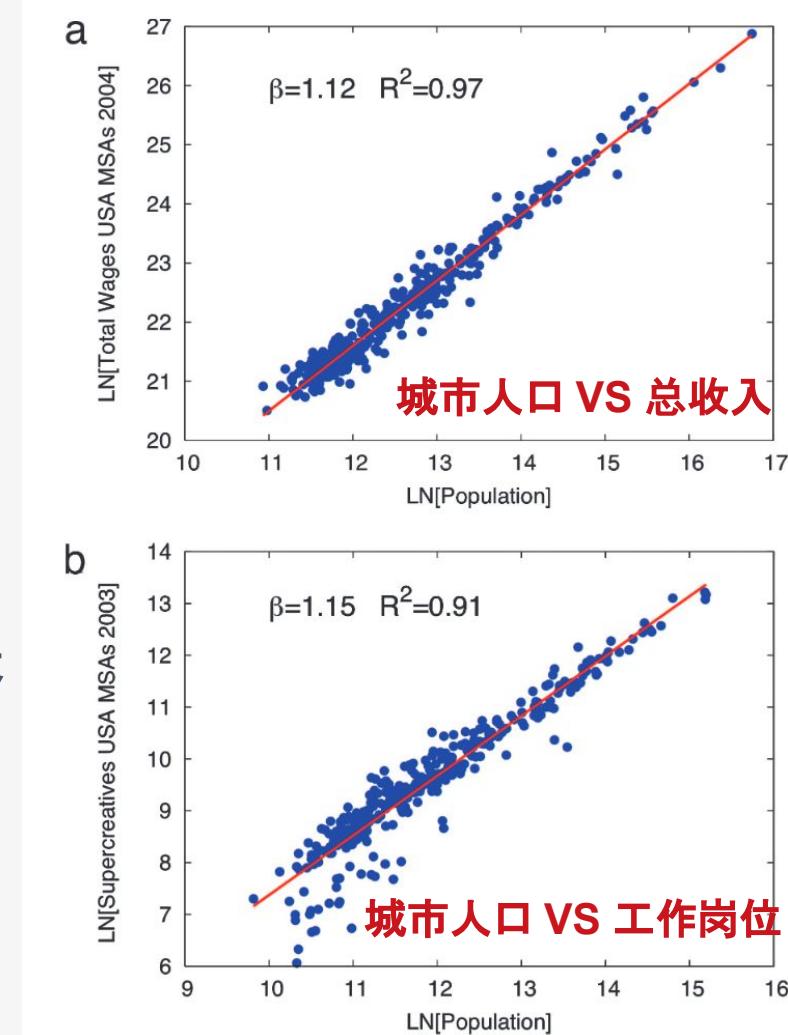
注意右图中，横纵坐标都做了取对数的变换，称为双对数坐标，并分别拟合出了斜率为1.12和1.15的线性模型。

在双对数坐标下的线性模型是什么含义？

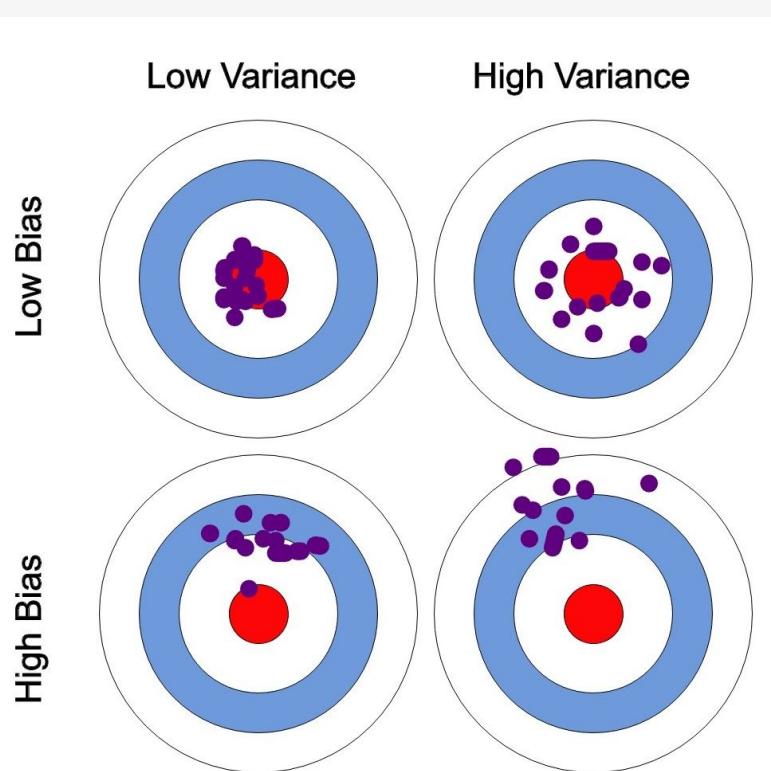
$$y_{ln} = \beta x_{ln} + \alpha$$

$$y = x^\beta e^\alpha$$

也就是当x变为两倍时，y变为原来的 2^β 倍，说明y和x的关系是超线性关系。



简单来说，拟合好的模型预测和真实数据之间的Error由什么构成？一部分是模型的偏差(Bias)，另一部分是模型本身方差(Variance)。



- Basic Model: $Y = f(X) + \varepsilon, \quad \varepsilon \sim N(0, \sigma_\varepsilon^2)$
- The expected prediction error of a regression fit $\hat{f}(X)$

$$\begin{aligned} Err(x_0) &= E[(Y - \hat{f}(x_0))^2 | X = x_0] \\ &= \sigma_\varepsilon^2 + [E\hat{f}(x_0) - f(x_0)]^2 + E[\hat{f}(x_0) - E\hat{f}(x_0)]^2 \\ &= \sigma_\varepsilon^2 + Bias(\hat{f}(x_0))^2 + Var(\hat{f}(x_0)) \\ &= \text{Irreducible Error} + Bias^2 + Variance \end{aligned}$$

R-squared

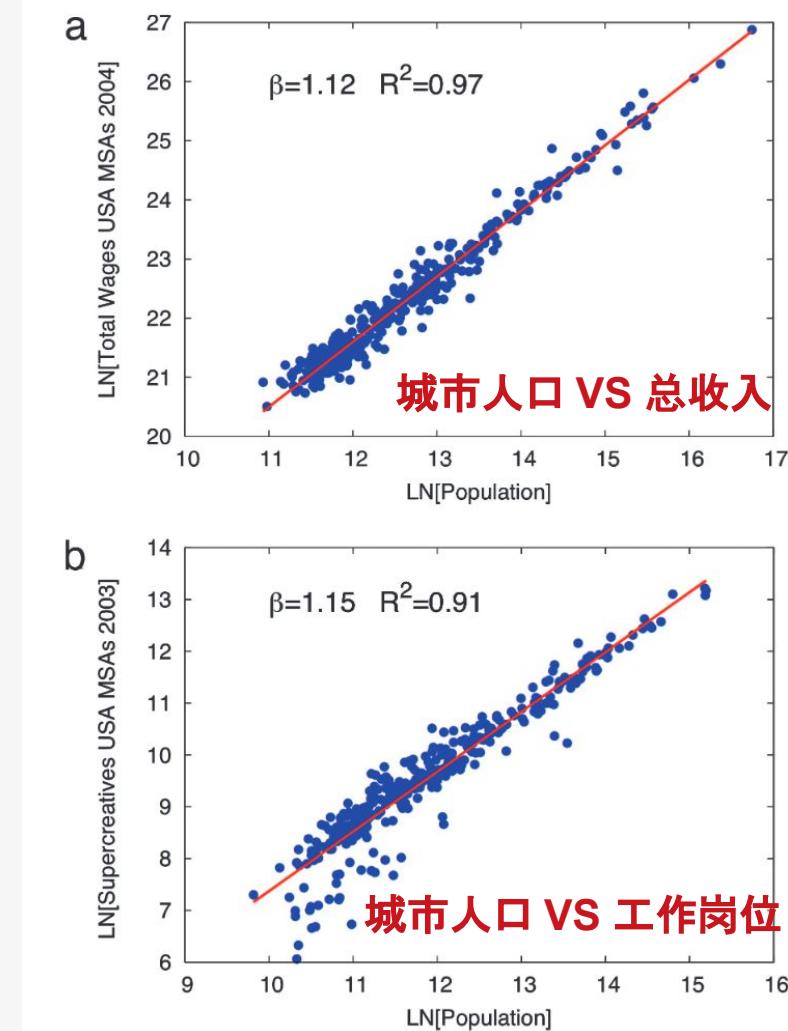
R方表示经过模型拟合，残差的方差比变量的方差减少了多少。

例如拟合前变量的方差为10，经过拟合后，残差的方差只有2，那么R方为0.8；

什么是R方=0？预测房价的时候，把上海市所有房子的房价全部加起来取平均，粗暴地认为上海房价都是这个平均值，此时这个预测模型的R方为0。因为它这个预测根本没有减少残差！

什么是R方=1？做出完美预测？可惜一般意义不大(¬‿¬)¬

注意，R方并不是某个数的平方，因此R方可以小于0，说明这个预测比均值预测还要差。



模型复杂度和误差的关系

- 模型复杂度可以简单用参数数量来衡量
- 线性模型有两个独立参数，二次曲线有三个参数...
- 模型复杂度越高，能在训练集上拟合得更好；但复杂度太高的时候，换一组数据，效果就不好了，甚至可能很差！

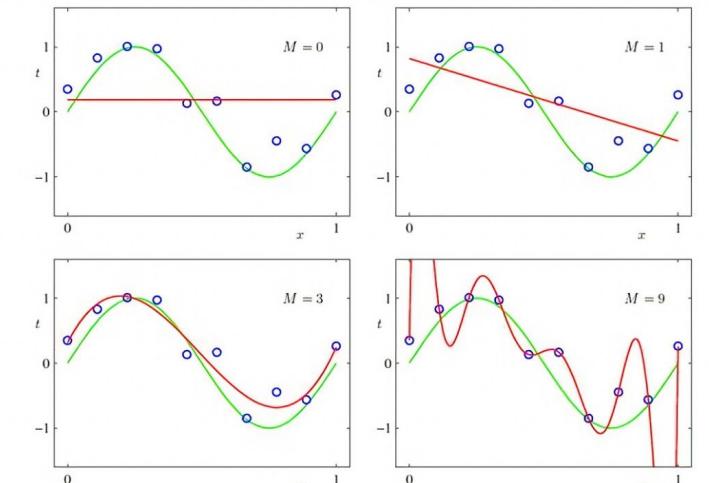
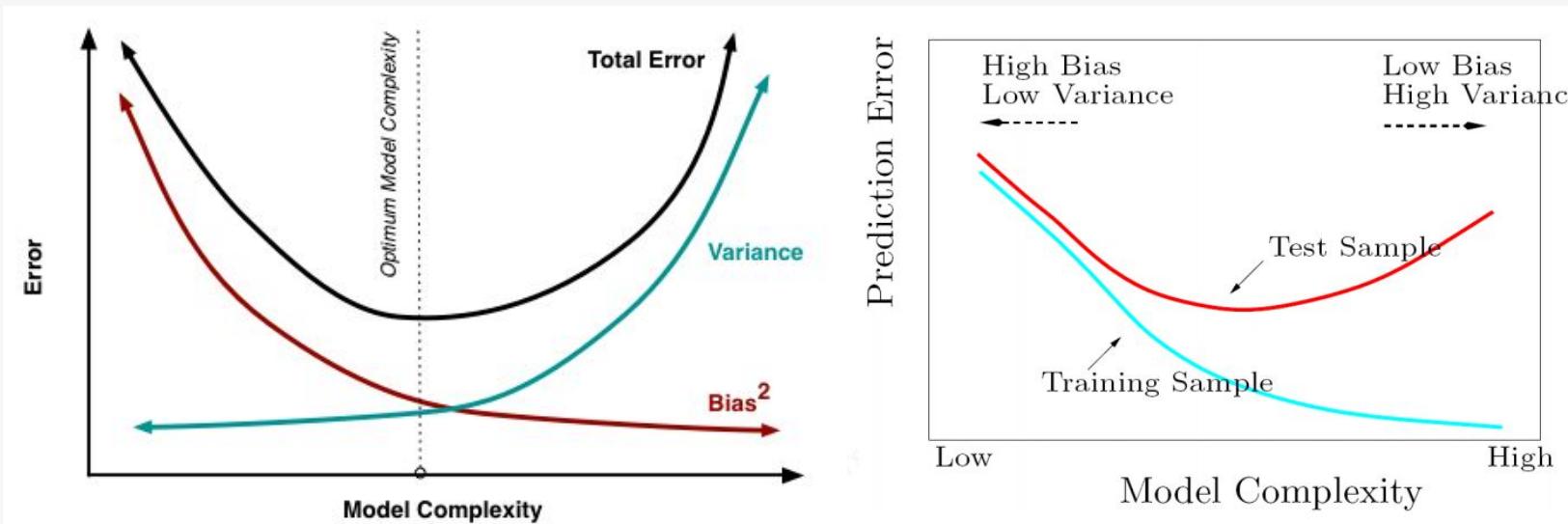
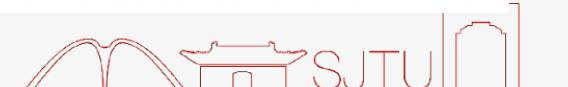


Figure 1.4 Plots of polynomials having various orders M , shown as red curves, fitted to the data set shown in Figure 1.2.



模型复杂度和误差的关系

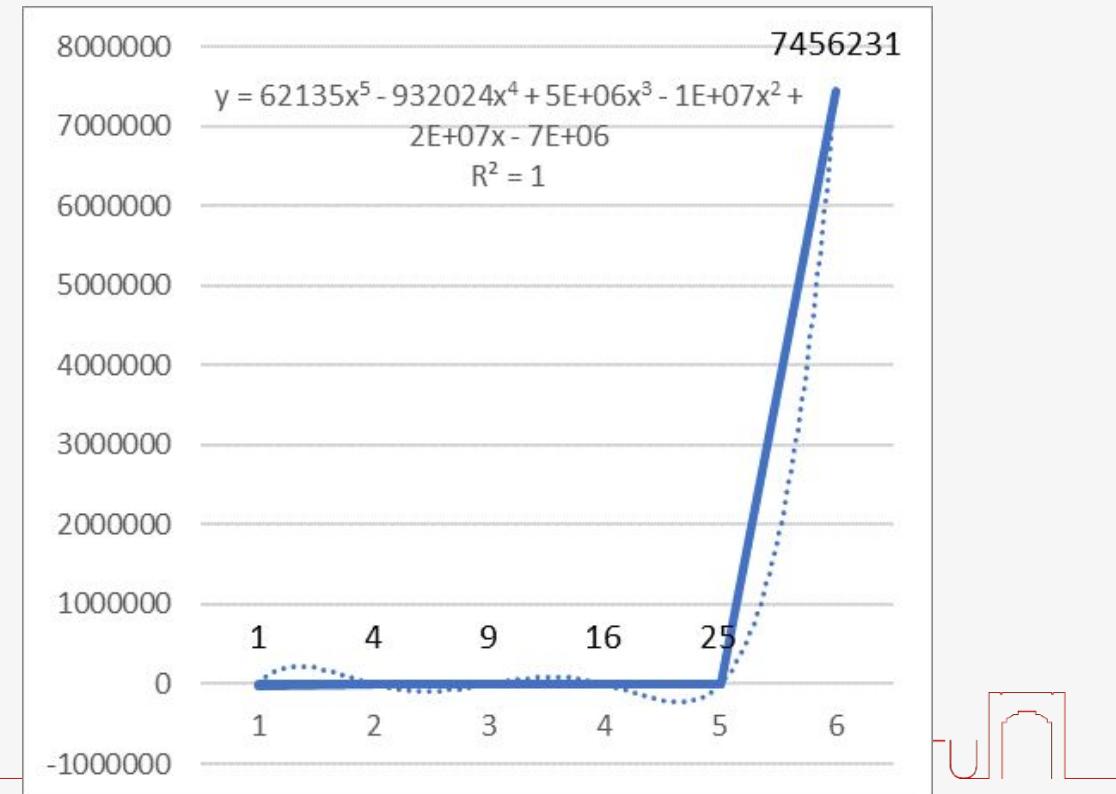
- 模型复杂度可以简单用参数数量来衡量
- 线性模型有两个独立参数，二次曲线有三个参数...
- 模型复杂度越高，能在训练集上拟合得更好；但复杂度太高的时候，换一组数据，效果就不好了，甚至可能很差！

问：1, 4, 9, 16, 25的下一个数字是什么？

答：7456231

给足模型参数，模型总能拟合

但这真的有用吗？



● 概率分布

- 概率密度函数、累积分布函数
- histogram和KDE
- 典型离散分布与连续分布
- 分布差异度量

● 函数拟合

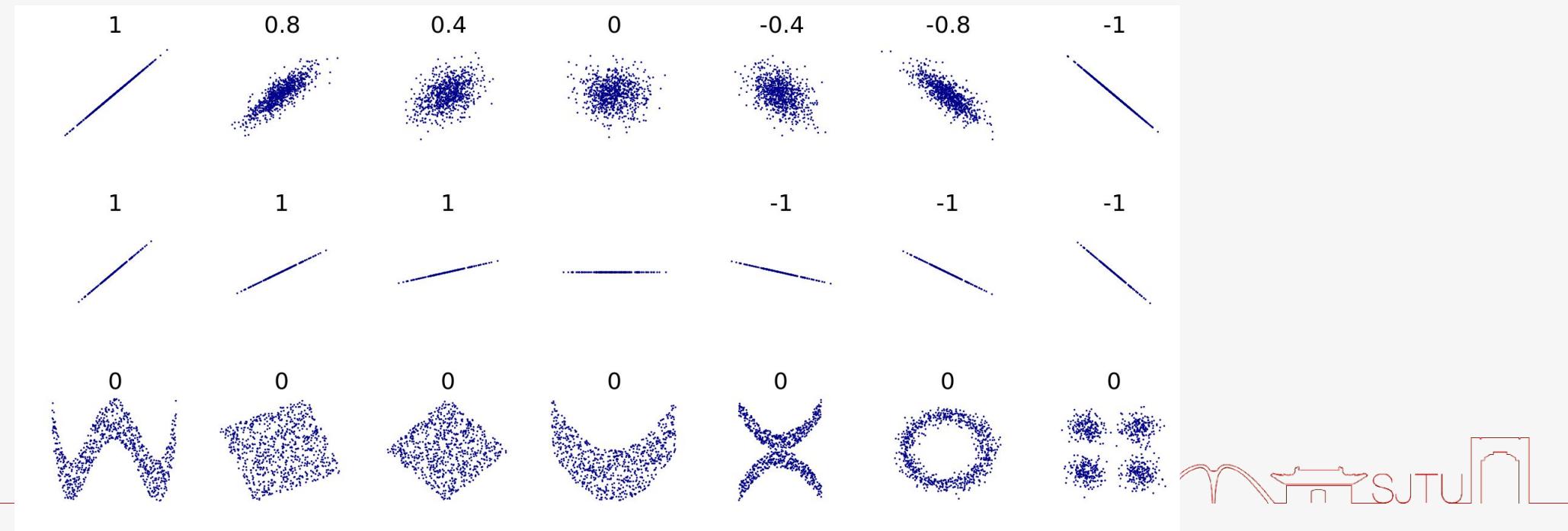
- 常用规律函数
- 拟合优度(Goodness-of-fit)
- 模型复杂度和误差的关系

● 数据相关性



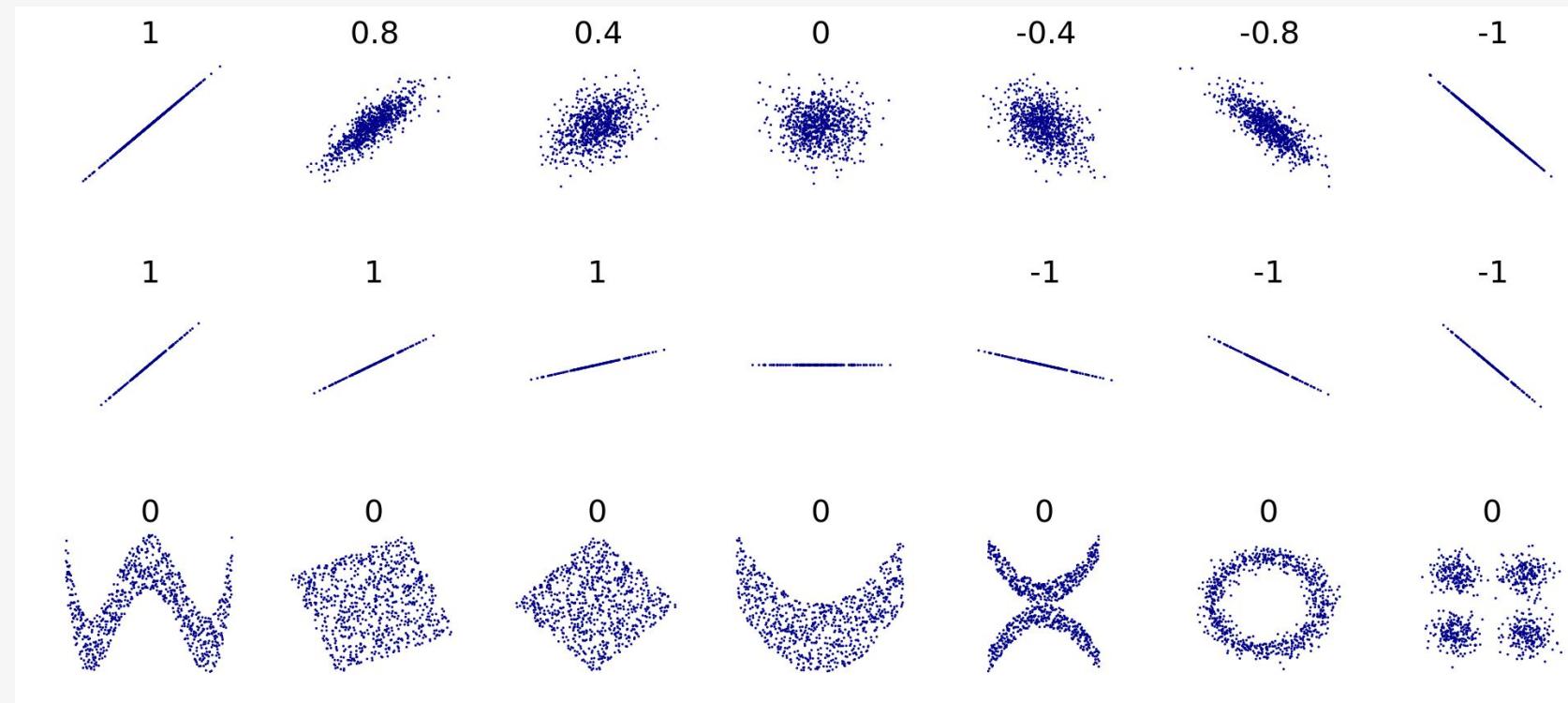
独立、线性相关、相关系数

- 从概率的角度看, AB两个事件**独立**定义为 $P(A\text{发生}) * P(B\text{发生}) = P(AB\text{都发生})$
- 统计学中, 相关一般指线性相关, 可以采用相关系数来评估两个变量之间相关性的大小。对于连续变量, 一般采用Pearson相关系数。
- 相关系数为正数, 说明两个变量的变化趋势相同; 相关系数为负则相反
- 相关系数=0, 说明线性不相关。



独立、线性相关、相关系数

- 独立则线性不相关, 但线性不相关不能推出独立。
 - 见第三行相关系数等于0
- 特殊性质: 服从正态分布的变量X和Y不线性相关时, 则X和Y独立。



数字也可能欺骗你

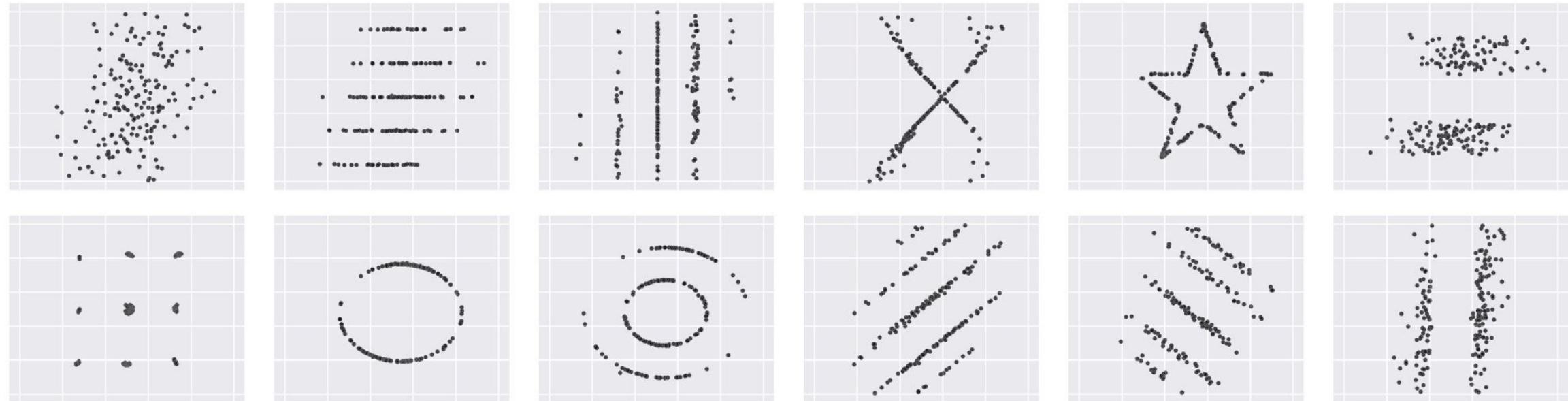


Figure 1. A collection of data sets produced by our technique. While different in appearance, each has the same summary statistics (mean, std. deviation, and Pearson's corr.) to 2 decimal places. ($\bar{x} = 54.02$, $\bar{y} = 48.09$, $s_x = 14.52$, $s_y = 24.79$, Pearson's $r = +0.32$)

一组数据集。尽管看起来各不相同，他们都拥有相同的统计量(平均数(一阶统计量)、标准差(二阶统计量)、皮尔森相关系数)

辛普森悖论



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

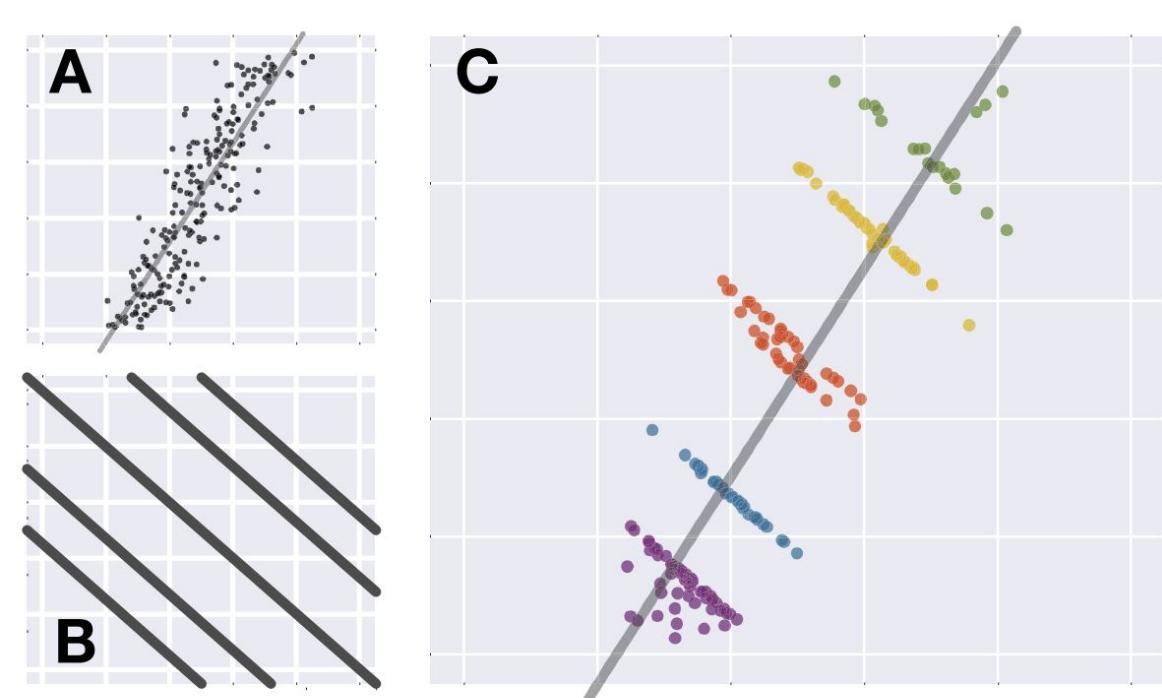


Figure 7. Demonstration of Simpson's Paradox. Both datasets (A and C) have the same overall Pearson's correlation of +0.81, however after coercing the data towards the pattern of sloping lines (B), each subset of data in (C) has an individually negative correlation.

A和C都有很高的正相关系数(+0.81)，但如果C图的数据引入了一个新的类别变量进行分类(按照图B)，每一类内横轴和纵轴居然形成了非常高的负相关！

医院	病情	入院总人数	死亡人数	存活人数	存活率
医院 A					
	合计	1000	100	900	90%
医院	病情	入院总人数	死亡人数	存活人数	存活率
医院 B					
	合计	1000	200	800	80%

你会选择哪个医院？
30%<52.5%
96.7%<98.3%
90%>80%



大纲

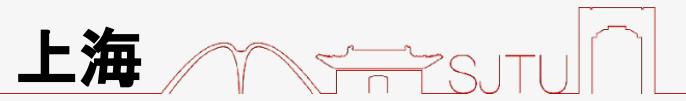
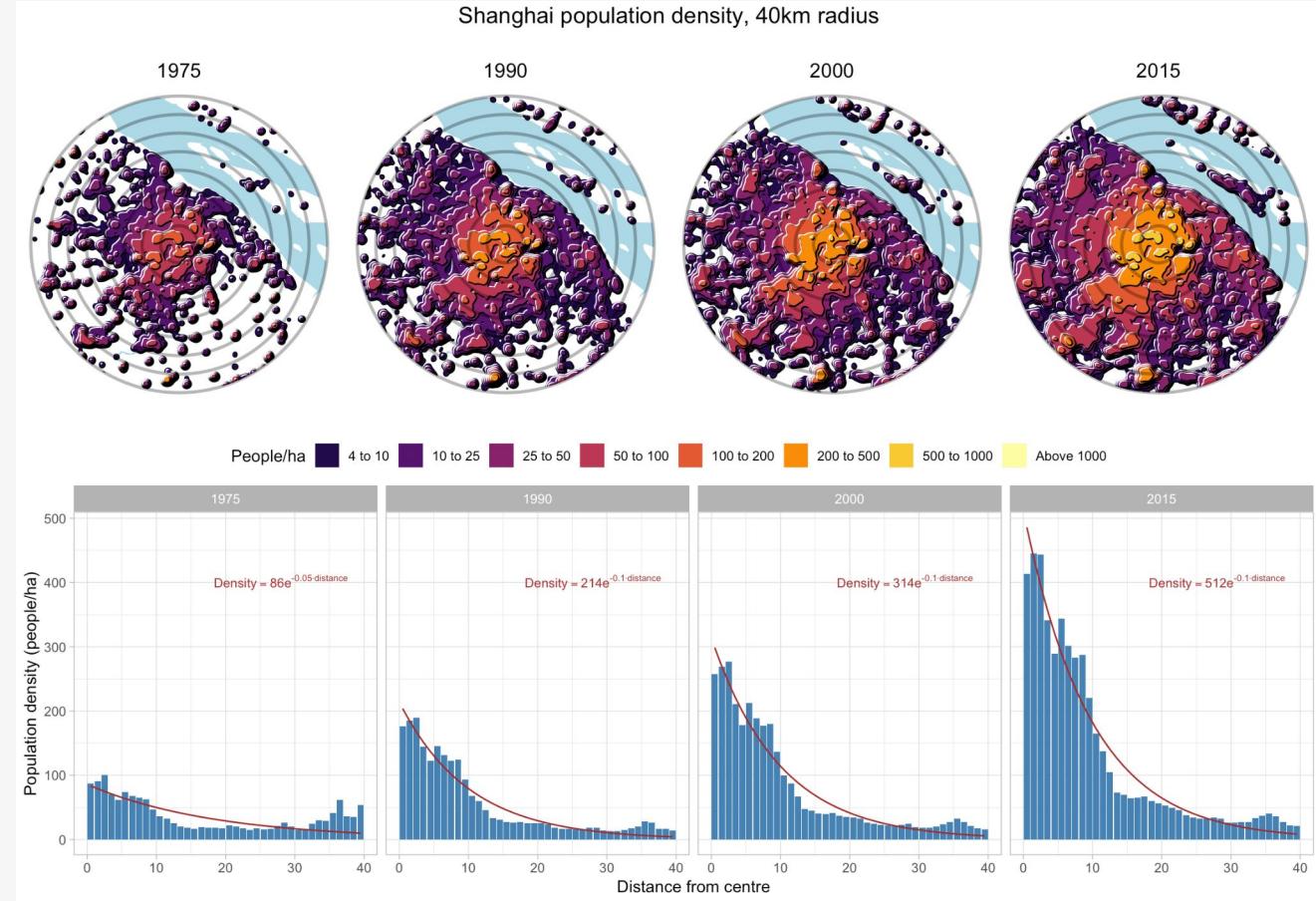
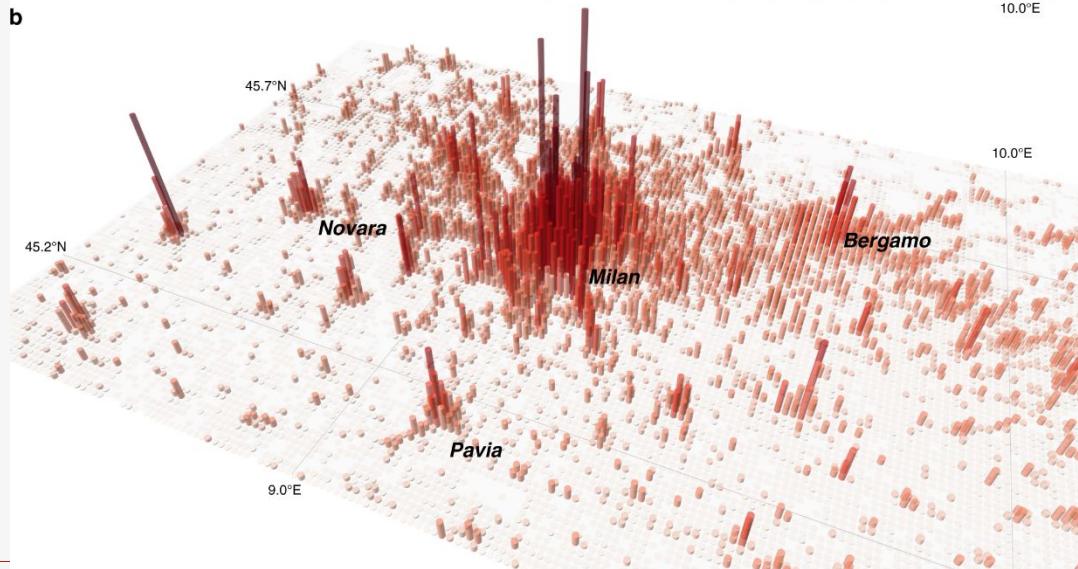
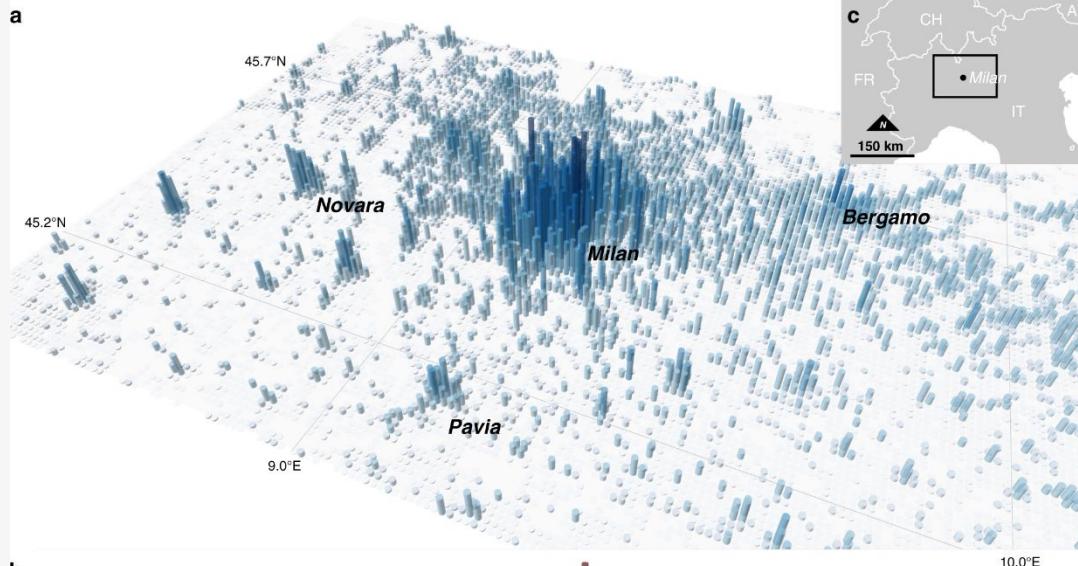


- 大数据和社会科学的关系
- 计算社会科学中的常用大数据
- 大数据的缺点、偏见
- 数据科学基础知识
- Python数据分析实践

理解城市人口分布



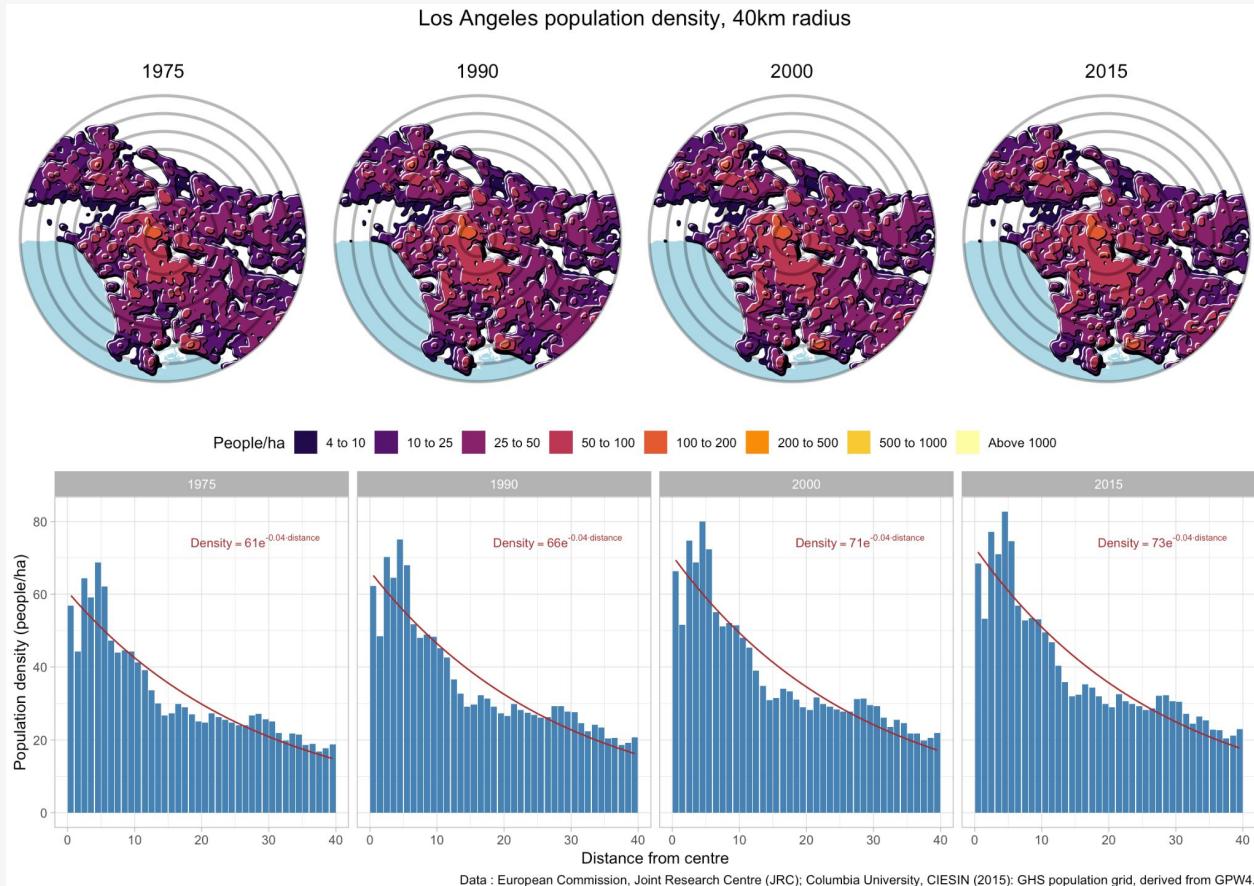
上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



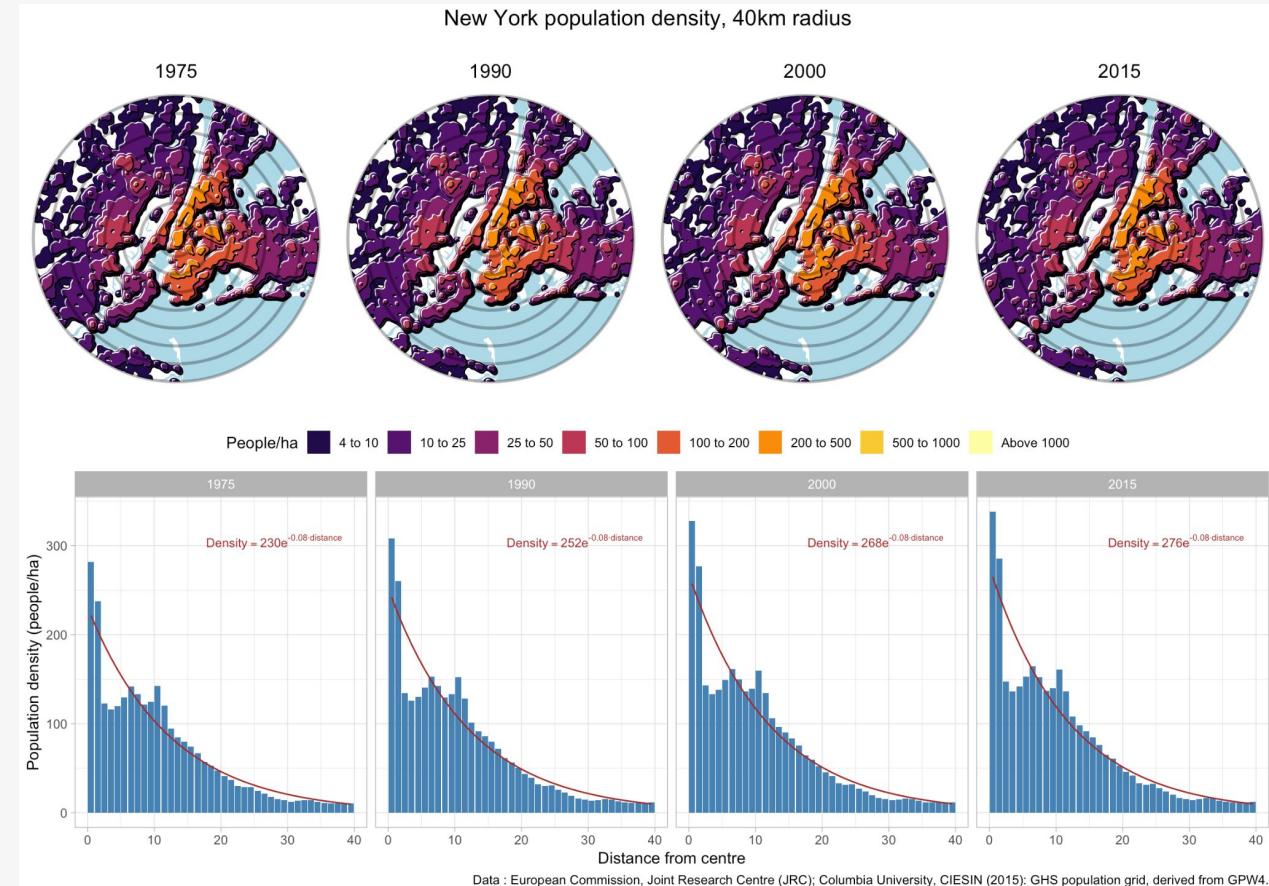
理解城市人口分布



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



洛杉矶



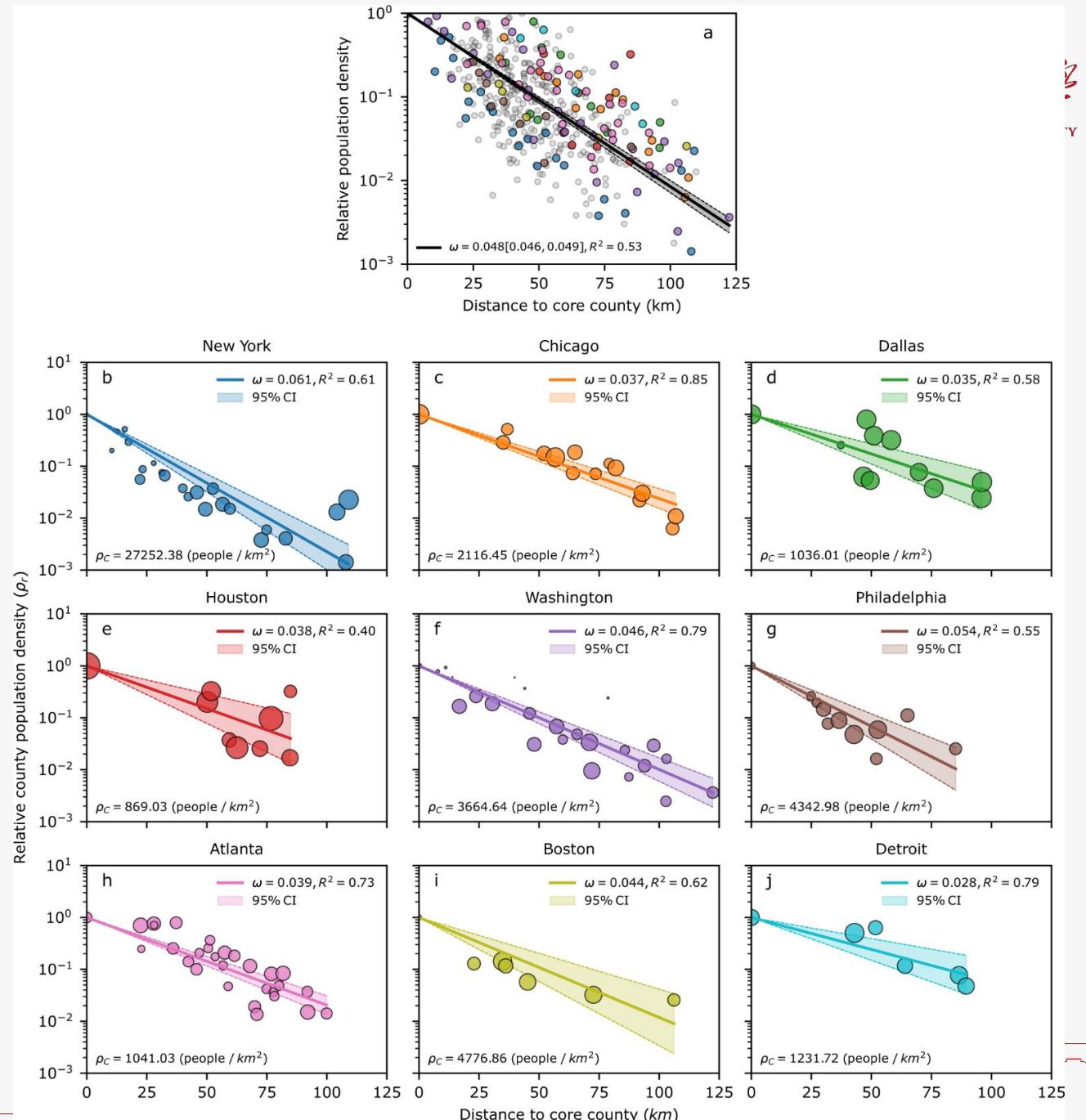
纽约



理解城市人口分布

- 右图表示每个county到城市中心的距离和county人口密度的关系；
- 注意这里y轴为对数坐标，x轴为线性坐标，因此此时拟合出的直线表示人口密度随着远离市中心而指数式下降。

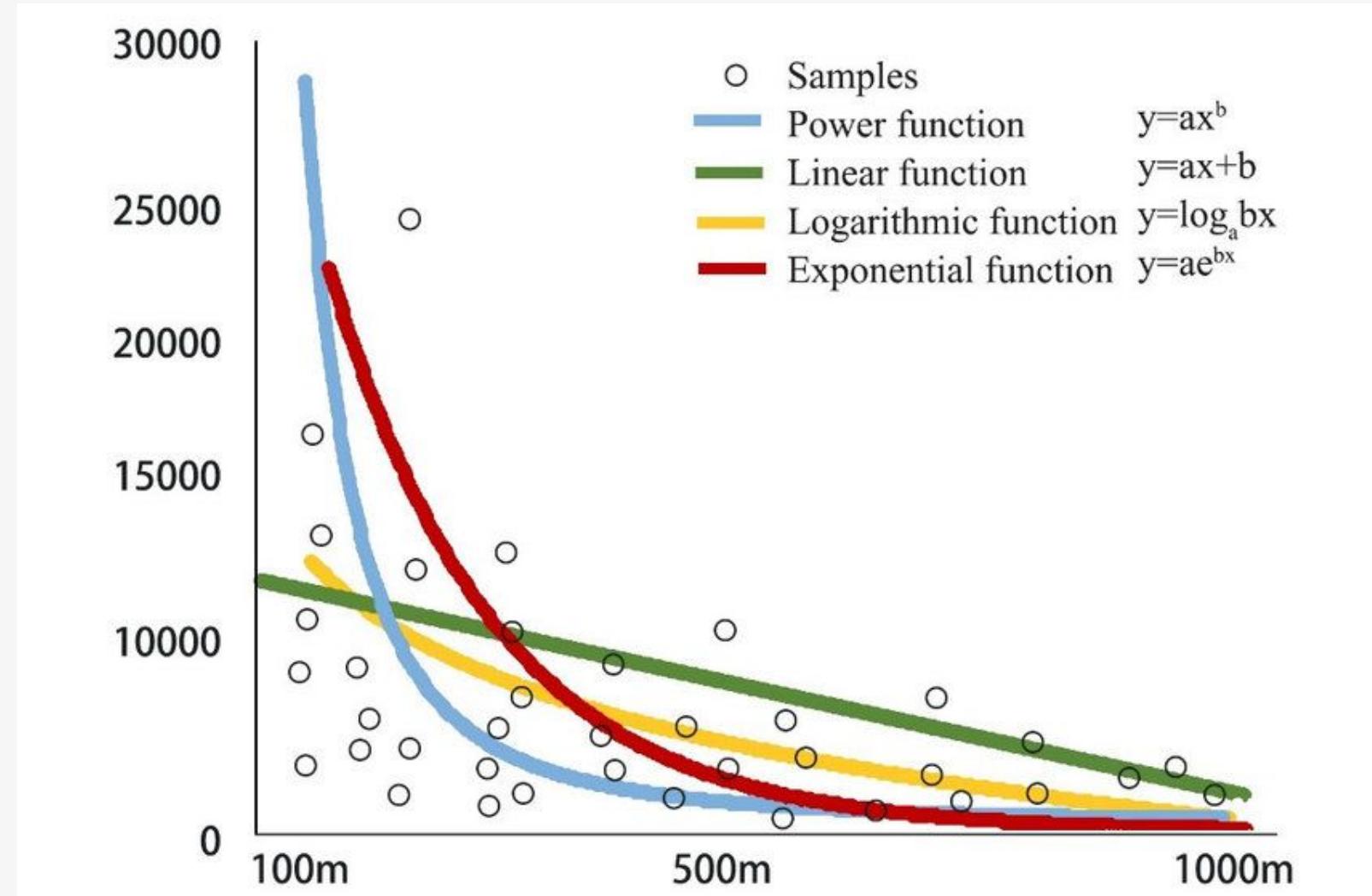
<https://www.nature.com/articles/s42949-022-00075-9>



函数拟合



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



Correlation between Transit-Oriented Development (TOD), Land Use Catchment Areas, and Local Environmental Transformation



建立环境，跑通代码



- Jupyter notebook



WorldPop Hub

DATA | CONTACT

Open Spatial Demographic Data and Research

WorldPop develops peer-reviewed research and methods for the construction of open and high-resolution geospatial data on population distributions, demographic and dynamics, with a focus on low and middle income countries.

Datasets

Open access spatial demographic datasets built using transparent approaches.

[total 44,745 datasets]

[Population Count](#) 20,724

[Population Density](#) 9,955

[Population Weighted Density](#) 4

[Births](#) 234

[Pregnancies](#) 234

[Age and sex structures](#)

6,036

[Development Indicators](#)

42

[Dependency Ratios](#) 2

[Internal Migration](#) 4

[Dynamic Mapping](#) 2

[Global Flight Data](#) 3

[Global Holiday Data](#) 5

[Covariates](#) 6,474

[Grid-cell surface areas](#) 250

[Administrative Areas](#) 500

[Urban change](#) 27

[Global Settlement Growth](#)

249

