



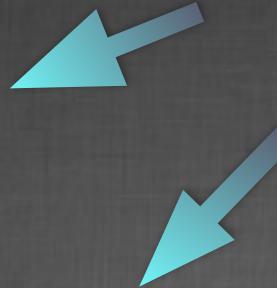
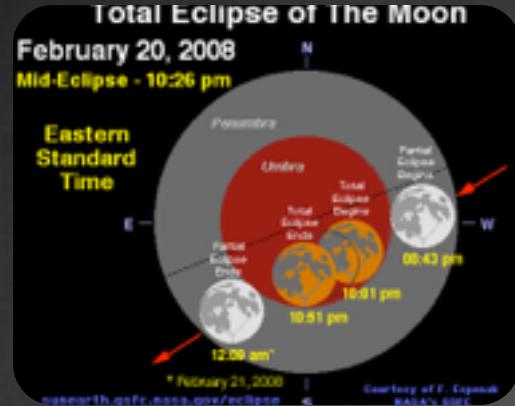
Northeastern

STATISTICS AND MODELING WITH NOVEL DATA STREAMS

Alessandro Vespignani
@alexvespi



MOBS LAB
LABORATORY FOR THE MODELING OF BIOLOGICAL
AND SOCIO-TECHNICAL SYSTEMS



Prediction = is based on what is known (mostly influenced by initial conditions) and provide a single outcome (although it may be probabilistic).

Forecast = given the best current knowledge on the system (knowledge often with uncertainties, model's parameters inferred statistically, etc) there is an ensemble of predictions that are analyzed statistically.

Projection = attempt to describe what would happen under certain assumptions and hypotheses (what if)





THE SOCIAL ATOM

The problem : technosocial systems are unpredictable because include the infinite psychological and cognitive reactions of individuals.
We're too complex and unpredictable.

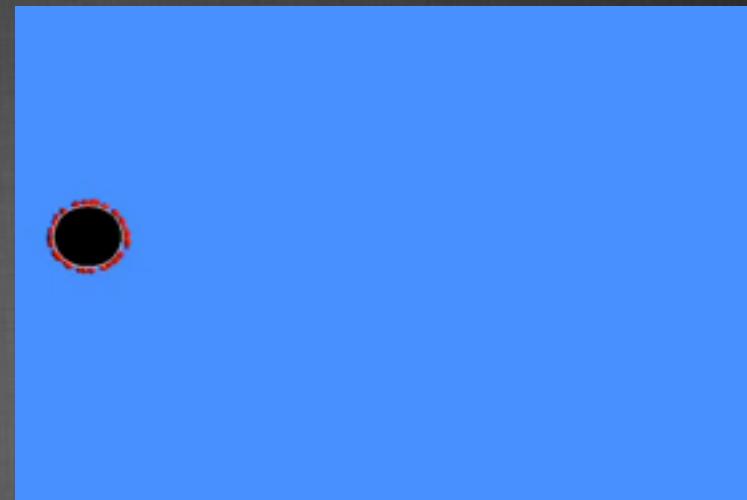


The vision: From the “social atom” or “social molecule” (i.e. small social groups) to the quantitative analysis of “social aggregate states” (Lundberg, Moreno).



THE SOCIAL ATOM

The problem : technosocial systems are unpredictable because include the infinite psychological and cognitive reactions of individuals.
We're too complex and unpredictable.



The vision: From the “social atom” or “social molecule” (i.e. small social groups) to the quantitative analysis of “social aggregate states” (Lundberg, Moreno).



THE SOCIAL ATOM

The problem : technosocial systems are unpredictable because include the infinite psychological and cognitive reactions of individuals.
We're too complex and unpredictable.



The vision: From the “social atom” or “social molecule” (i.e. small social groups) to the quantitative analysis of “social aggregate states” (Lundberg, Moreno).



MATHEMATICAL -> COMPUTATIONAL

Numerical Weather models

- 1920 Richardson integrate manually equations of the atmosphere
- 1950 First numerical weather forecast (24h computation for a 24h forecast)
- 1955 Numerical weather prediction models became operational by the USWB
- 2015 Government and Commercial entities routinely forecast up to three weeks

Numerical Epidemic models

- 1930 Reed-Frost define a simple chain binomial model that they integrate with a “sandbox’ computer
- 1952 First Reed-Frost numerical implementation
- 1980-2000 progress toward the definition of large-scale individual models
- 2005 Large scale agent-based models
- 2105 ???

MATHEMATICAL -> COMPUTATIONAL

Numerical Weather models

- “Primitive” Weather Forecasting Equations
- $p = \rho R T$ Ideal Gas Law (Equation of State)
- $\bar{a}_h = \sum \left(\frac{\bar{F}_h}{m} \right)$ Newton's Second Law of Motion
- $\bar{a}_v = \sum \left(\frac{\bar{F}_v}{m} \right) = (\bar{P} \bar{G} \bar{A})_v - \bar{g}$
- Hydrostatic Law (Obtained from the Equation of Vertical Motion)
- $\Delta T = \Delta q/c_p + (1/\rho)\Delta p$ First Law of Thermodynamics
- $(1/\rho)\Delta \rho/\Delta t = -\text{DIV}$ Conservation of Mass Applied to the Atmosphere (Equation of Continuity)
- $\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} + \omega \left(\frac{\partial T}{\partial p} + \frac{RT}{pc_p} \right) = \frac{J}{c_p}$

Numerical Epidemic models

- a simple hat they box'
 - numerical toward the
 - point-based
 - models
 - 2105 ???
- Zonal wind:
- $$\frac{\partial u}{\partial t} = \eta v - \frac{\partial \Phi}{\partial x} - c_p \theta \frac{\partial \pi}{\partial x} - z \frac{\partial u}{\partial \sigma} - \frac{\partial(\frac{u^2+v^2}{2})}{\partial x}$$
- Meridional wind:
- $$\frac{\partial v}{\partial t} = -\eta \frac{u}{v} - \frac{\partial \Phi}{\partial y} - c_p \theta \frac{\partial \pi}{\partial y} - z \frac{\partial v}{\partial \sigma} - \frac{\partial(\frac{u^2+v^2}{2})}{\partial y}$$
- Temperature:
- $$\frac{\delta T}{\partial t} = \frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} + w \frac{\partial T}{\partial z}$$
- Precipitable water:
- $$\frac{\delta W}{\partial t} = u \frac{\partial W}{\partial x} + v \frac{\partial W}{\partial y} + w \frac{\partial W}{\partial z}$$
- Pressure thickness:
- $$\frac{\partial}{\partial t} \frac{\partial p}{\partial \sigma} = u \frac{\partial}{\partial x} x \frac{\partial p}{\partial \sigma} + v \frac{\partial}{\partial y} y \frac{\partial p}{\partial \sigma} + w \frac{\partial}{\partial z} z \frac{\partial p}{\partial \sigma}$$
- $$\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial \omega}{\partial p} = 0 \quad 0 = -\frac{\partial \phi}{\partial p} - \frac{RT}{p}$$

MATHEMATICAL -> COMPUTATIONAL

Numerical Weather models

- 1920 Richardson integrate manually equations of the atmosphere
- 1950 First numerical weather forecast (24h computation for a 24h forecast)
- 1955 Numerical weather prediction models became operational by the USWB
- 2015 Government and Commercial entities routinely forecast up to three weeks

Numerical Epidemic models

- 1930 Reed-Frost define a simple chain binomial model that they integrate with a “sandbox’ computer
- 1952 First Reed-Frost numerical implementation
- 1980-2000 progress toward the definition of large-scale individual models
- 2005 Large scale agent-based models
- 2105 ???



MATHEMATICAL -> COMPUTATIONAL

Numerical Weather models

- 1920 Richardson integrate manually equations of the atmosphere
- 1950 First numerical weather forecast (24h computation for a 24h forecast)
- 1955 Numerical weather prediction models became operational by the USWB
- 2015 Government and Commercial entities routinely forecast up to three weeks

Numerical Epidemic models

- 1930 Reed-Frost define a simple chain binomial model that they integrate with a “sandbox’ computer
- 1952 First Reed-Frost numerical implementation
- 1980-2000 progress toward the definition of large-scale individual models
- 2005 Large scale agent-based models
- 2105 ???

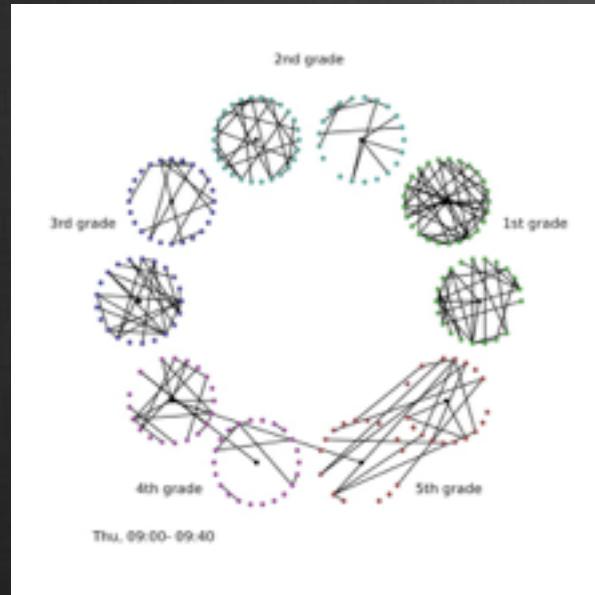
BIG DATA

- Data has fueled spectacular advances in the natural sciences over the last 100 years.
- So what is different now?
- Every 1.2 years, more human-driven socioeconomic data is produced than during all previous history
- Embedded within the data are the raw ingredients for understanding socio-technical and socio-economic systems
- The focus is on understanding these data sets in a statistical sense and more deeply the real world processes which produced the data



DATA ARE CENTRAL IN THE ANALYSIS OF CONTAGION PROCESSES

Human interactions



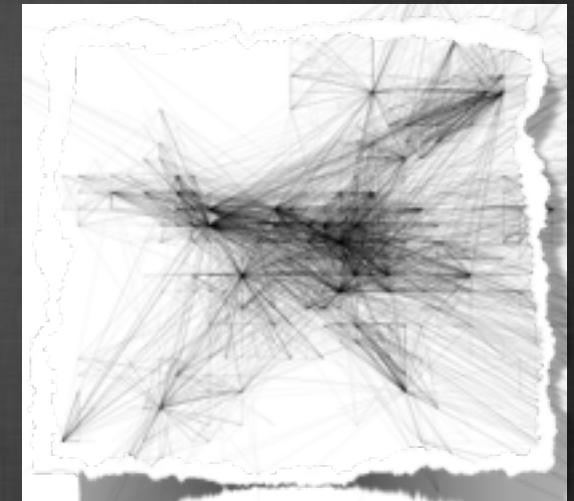
Within school contact
patterns

Human Mobility



Multiscale integration
of mobility networks in
the analysis of
potentially pandemic
pathogens spread.

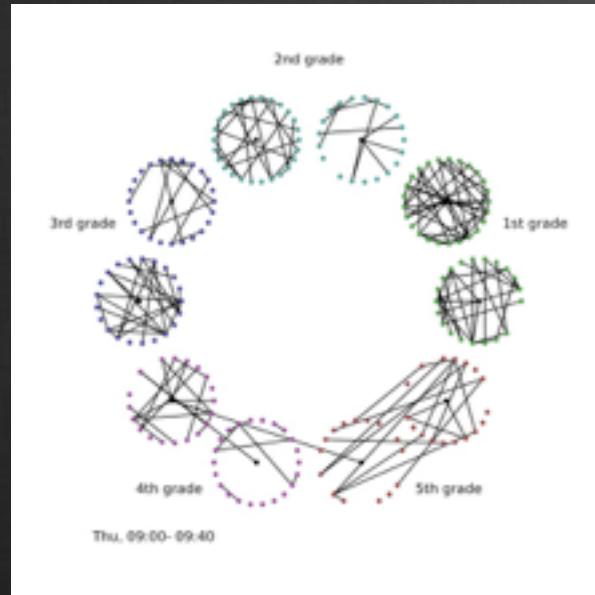
Social networks heterogeneities



Hubs, community,
clustering, heavy
tails, ...

DATA ARE CENTRAL IN THE ANALYSIS OF CONTAGION PROCESSES

Human interactions



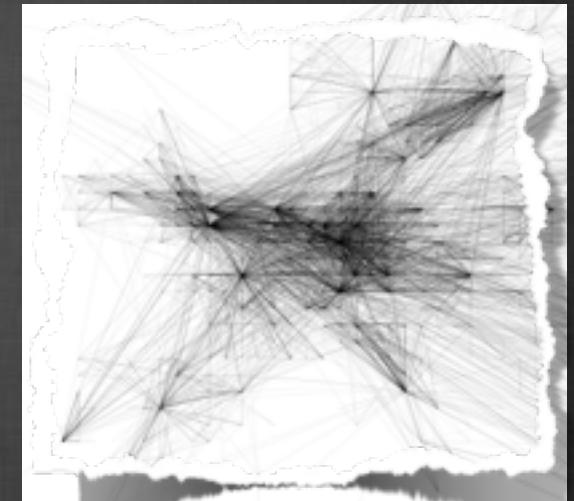
Within school contact
patterns

Human Mobility



Multiscale integration
of mobility networks in
the analysis of
potentially pandemic
pathogens spread.

Social networks heterogeneities



Hubs, community,
clustering, heavy
tails, ...

IN THE BEGINNING WAS DATA....THE OLD DATA. LET THERE BE NOVEL DATA SOURCES: ACTIVE-VS-PASSIVE DATA COLLECTIONS



Participatory platforms



Data harvesting

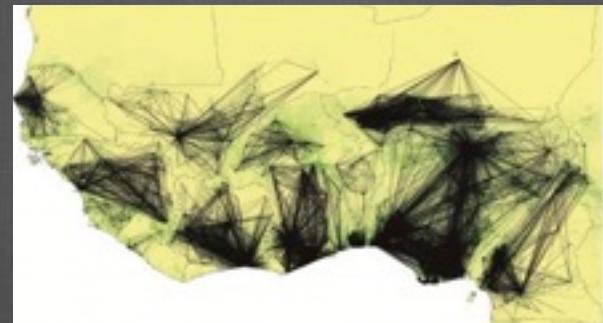
ACTIVE DATA COLLECTION

- Web participatory platforms
- RFID tags, SocioBadges
- Pervasive technologies
- Crowdsourced or manually curated data
-

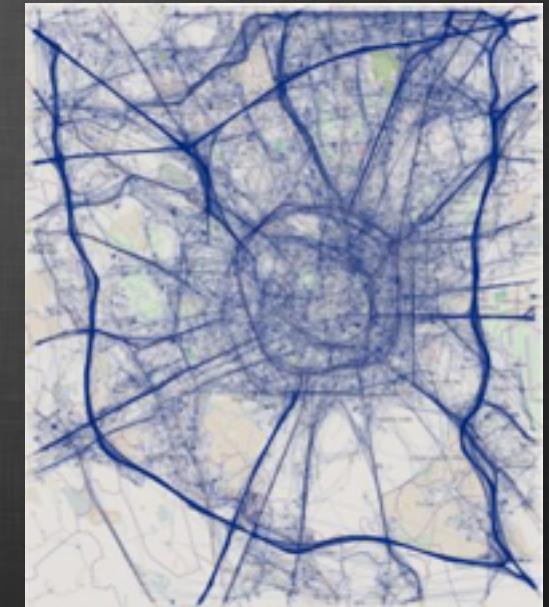
PASSIVE DATA COLLECTION: MOBILITY AT THE “POPULATION” SCALE



Bengtsson et al. Sci. Rep 2015



Wesolowsky et al. PLOS Curr. Out. 2014



Pedreschi et al. et al. PLOS Curr. Out. 2014

PASSIVE DATA ANALYSIS

- Search engine queries (GFT)
- Twitter conversation
- Wikipedia logs
- Over the counter drugs
- Facebook postings
- Restaurant reservation cancellations
- Hospitals parking imagery.

digital traces



digital traces

historical view
temporal horizon
limited reproducibility
limited context
data protection challenges



digital traces

A photograph of a person walking across a vast, light-colored sand dune under a clear blue sky. The person is small in the distance, and their footprints are clearly visible in the sand, forming a path that leads towards them.

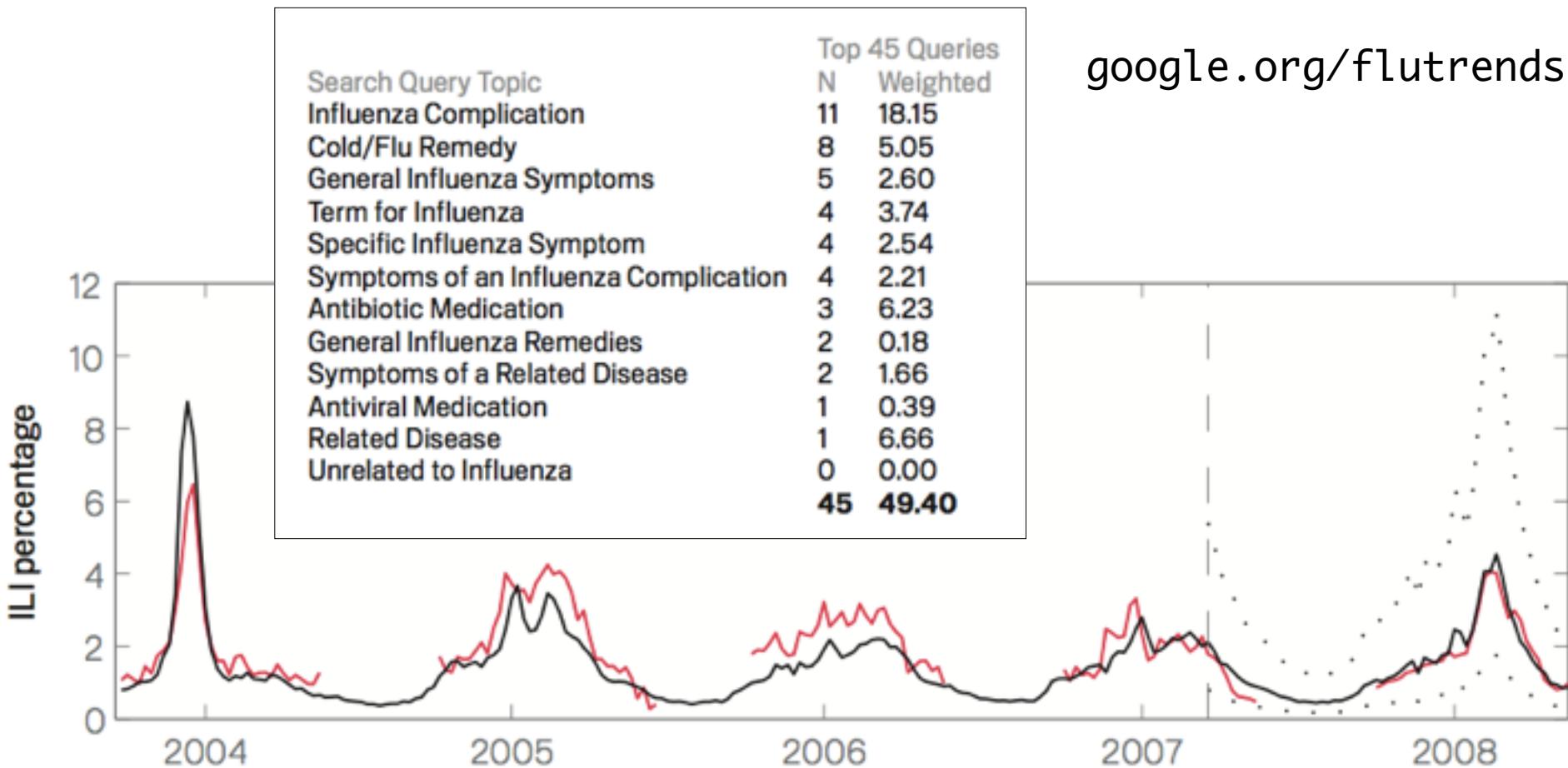
historical view
temporal horizon
limited reproducibility
limited context
data protection challenges

available as a side effect of many activities
machine processable, pattern discovery
high coverage, can work at scale

Detecting influenza epidemics using search engine query data

Jeremy Ginsberg¹, Matthew H. Mohebbi¹, Rajan S. Patel¹, Lynnette Brammer²,
Mark S. Smolinski¹ & Larry Brilliant¹

¹Google Inc. ²Centers for Disease Control and Prevention



NATURE | NEWS

عربي

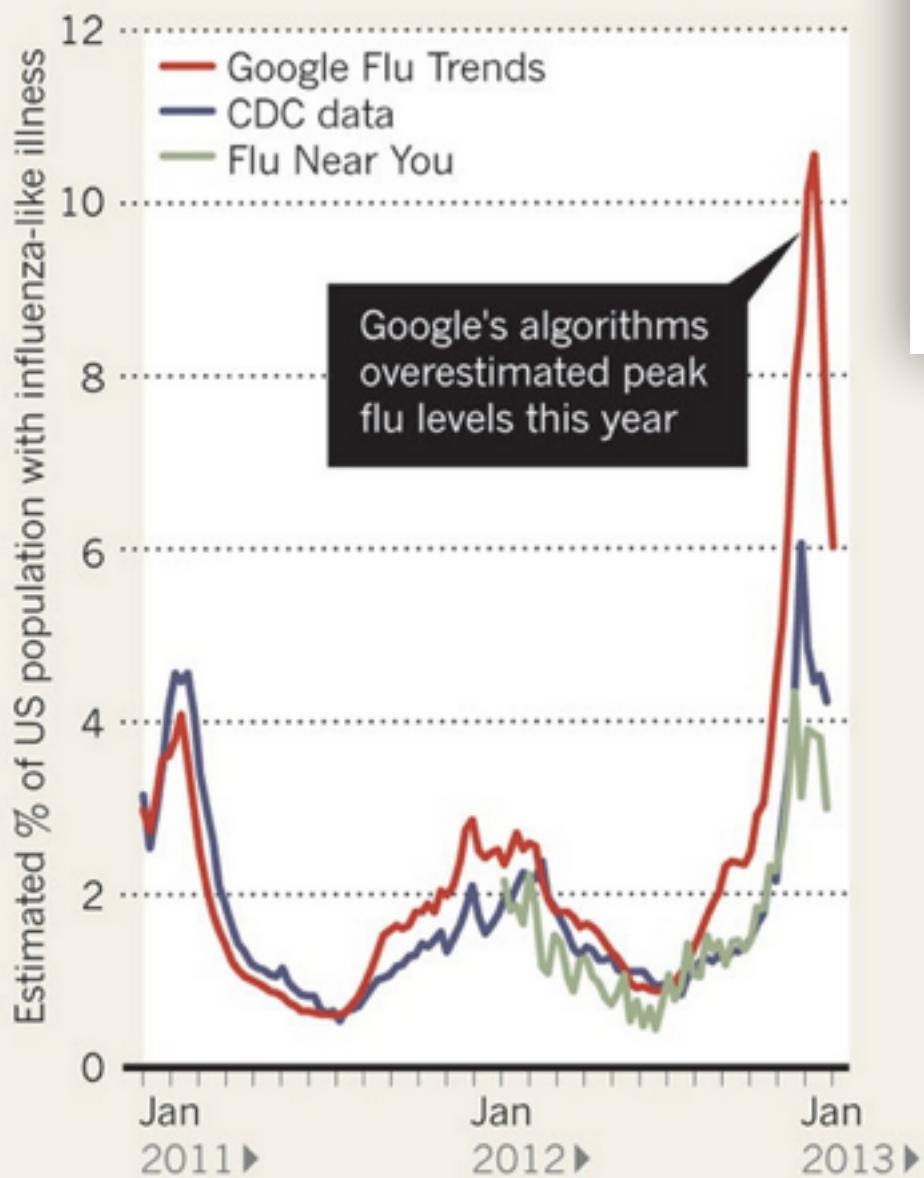
When Google got flu wrong

US outbreak foxes a leading web-based method for tracking seasonal flu.

Declan Butler

FEVER PEAKS

A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.



Google's algorithms
overestimated peak
flu levels this year

NATURE | NEWS

عربي

When Google got flu wrong

Science 14 March 2014:
 Vol. 343 no. 6176 pp. 1203–1205
 DOI: 10.1126/science.1248506

POLICY FORUM

BIG DATA

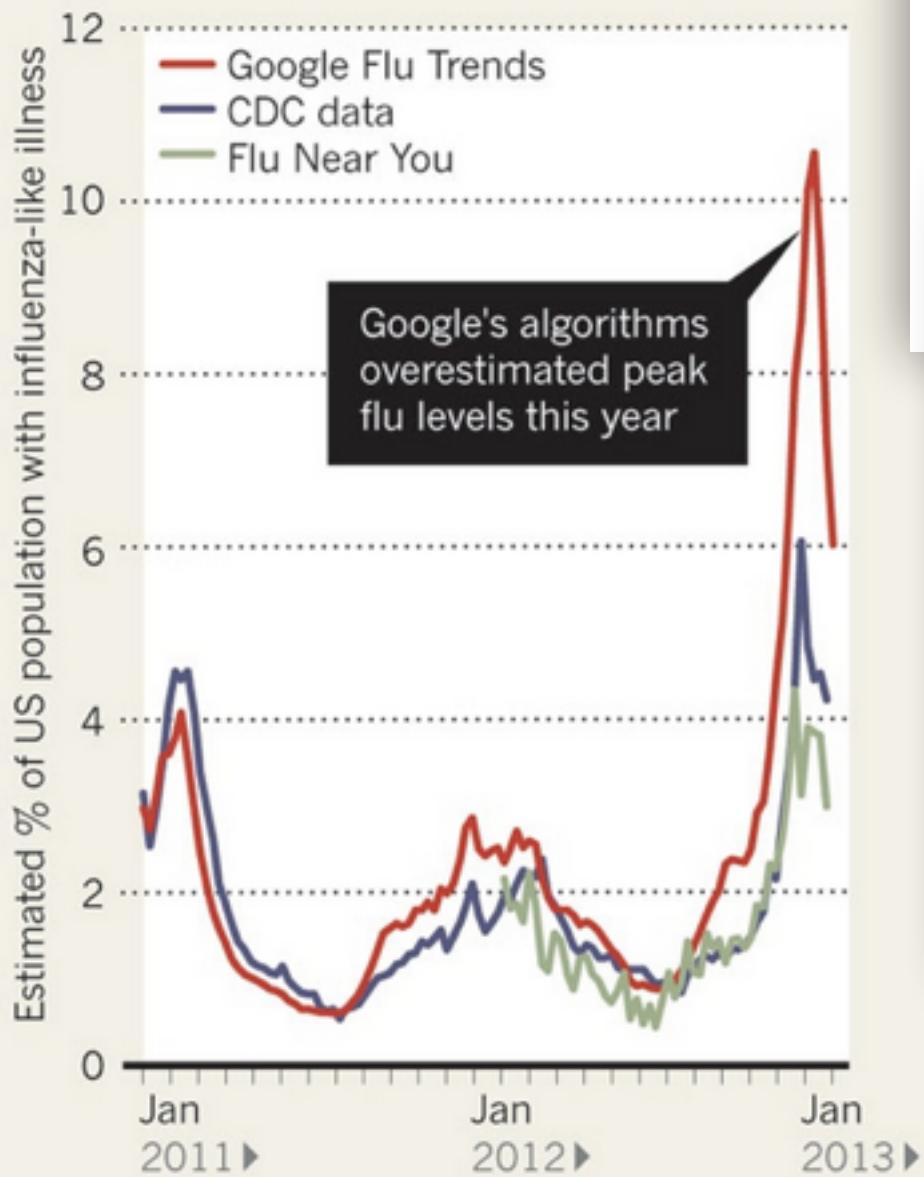
The Parable of Google Flu: Traps in Big Data Analysis

David Lazer^{1,2,*}, Ryan Kennedy^{1,3,4}, Gary King³, Alessandro Vespignani^{5,6,3}[Author Affiliations](#)[*Corresponding author. E-mail: \[d.lazer@neu.edu\]\(mailto:d.lazer@neu.edu\).](#)

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason the executives or the creators of the flu tracking system would have hoped. *Nature* reports predicting more than double the proportion of doctor visits for influenza-like illness than the Centers for Disease Control and Prevention (CDC), which bases its estimates on samples from laboratories across the United States (1, 2). This happened despite the fact that GFT often outperforms CDC reports. Given that GFT is often held up as an exemplary use of big data, what lessons can we draw from this error?

FEVER PEAKS

A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.



NATURE | NEWS

عربي

When Google got flu wrong

Science 14 March 2014:
Vol. 343 no. 6176 pp. 1203–1205
DOI: 10.1126/science.1248506

< Prev | T

Read

POLICY FORUM

BIG DATA

The Parable of Google Flu: Traps in Big Data Analysis

David Lazer^{1,2,*}, Ryan Kennedy^{1,3,4}, Gary King³, Alessandro Vespignani^{5,6,3}

FT Magazine

Home World Companies Markets Global Economy
Arts Magazine Food & Drink House & Home Lunch with the FT Style Books

March 28, 2014 11:38 am

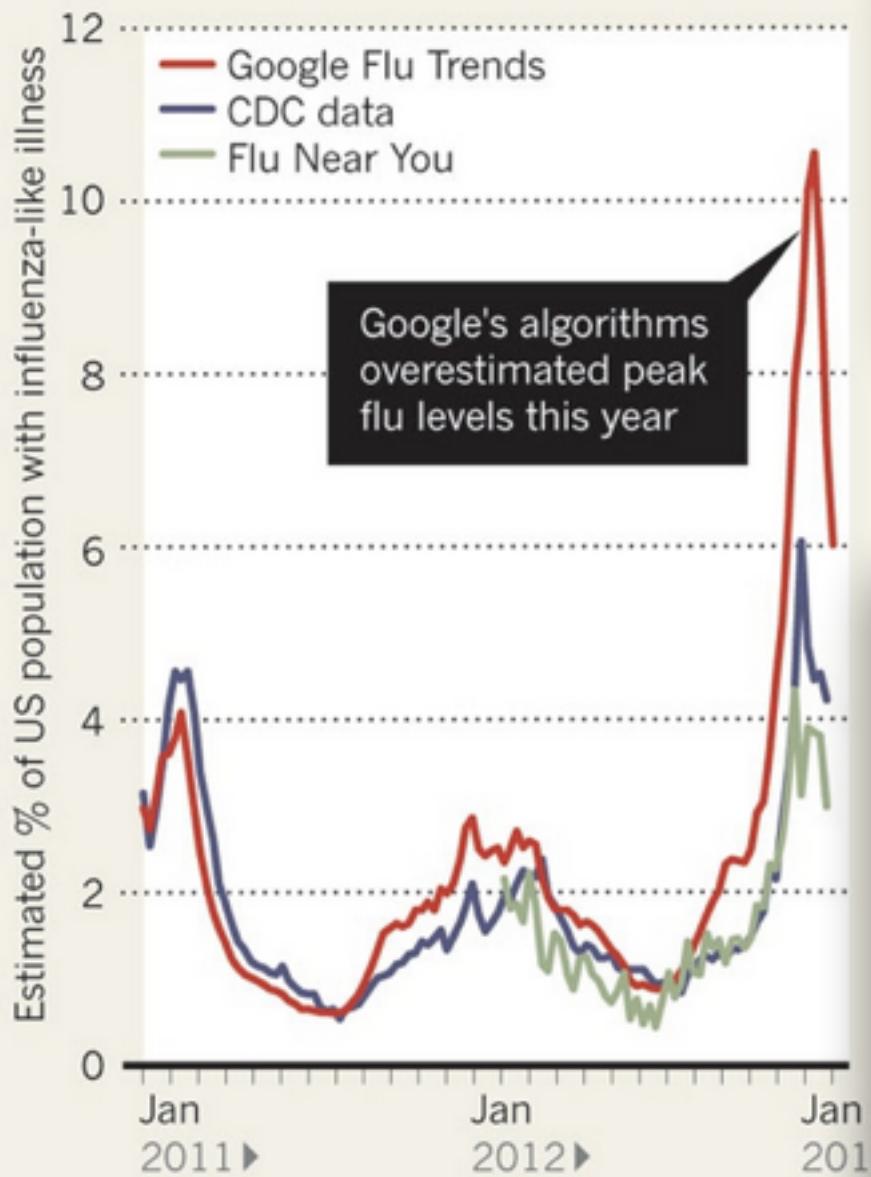
Big data: are we making a big mistake?

By Tim Harford

Big data is a vague term for a massive phenomenon that has rapidly become an obsession with entrepreneurs, scientists, governments and the media

FEVER PEAKS

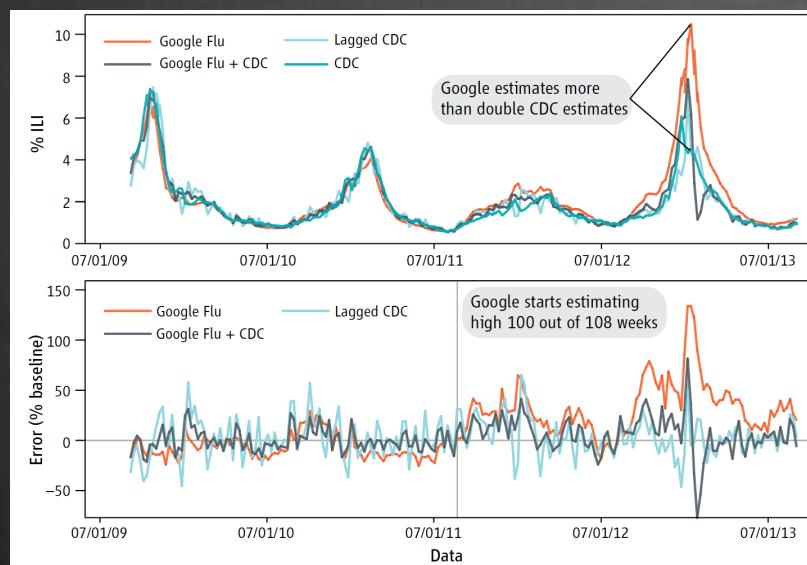
A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.



DIFFERENT FORECAST PHILOSOPHIES

Use surrogate signal in algorithm trained on historical data (generally CDC time series) to achieve lead time (real time data collection, time-series extrapolation)

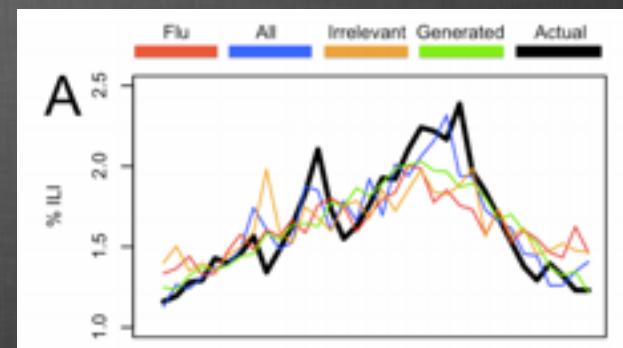
Google Flu Trend (paradigm)



Case study on GFT and other non-generative models simple Lagged regression can be 90% "good" (Lazer et al. Science 2014).

Twitter based approaches

- Salathe'; Culotta; Dredze etc. Etc. (since the first paper by Signorini et al.);
- Word selection
- Linear regression, Multiple linear regression; SVM Regresssion; EFS
- High-level geographical resolution
- Full natural language processing



Statistical biases, "zombies" etc

- Well discussed in the literature
- Similarities & difference with GFT.

THIS IS ~~FORECAST~~ NOWCAST !!!

We miss the microscopic generative foundations and models



MICRO/INDIVIDUAL BASED LARGE-SCALE COMPUTATIONAL MODELS FOR INFECTIOUS DISEASE SPREADING

Individual based,structured metapopulation models

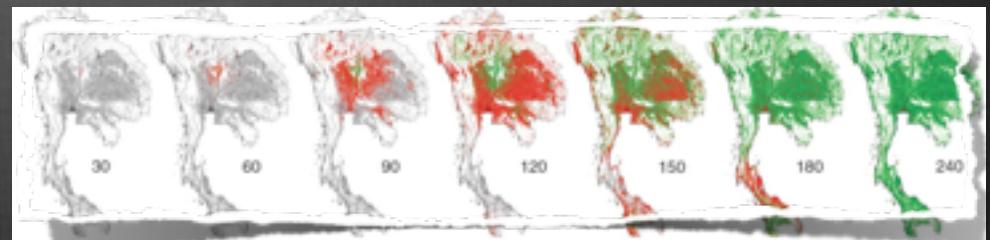
- Collection of sub-populations coupled by traveling individuals. [Ravchev, Longini. Mathematical Biosciences (1985) 50 urban areas worldwide; Viboud et al., Valleron et al., Brockmann et al., Colizza et al., Balcan et al.]

Agent based models

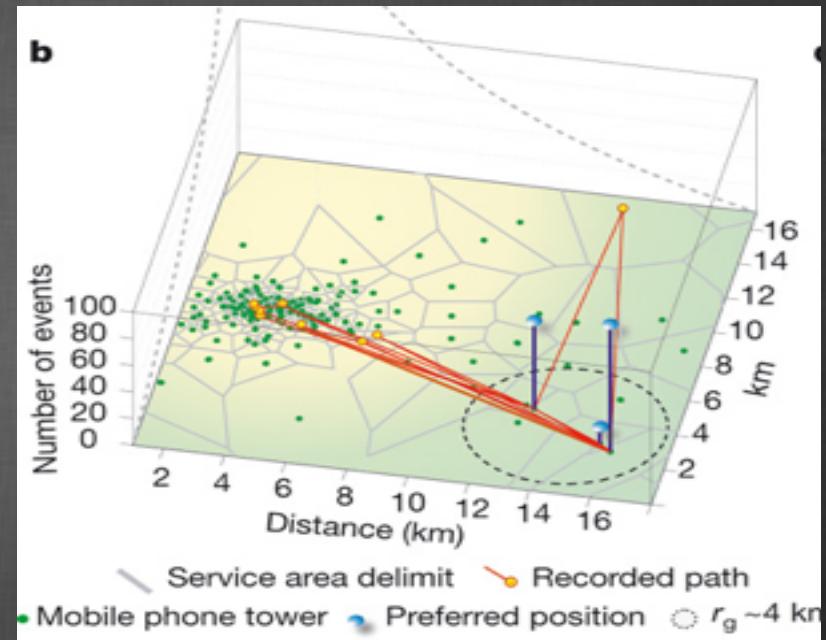
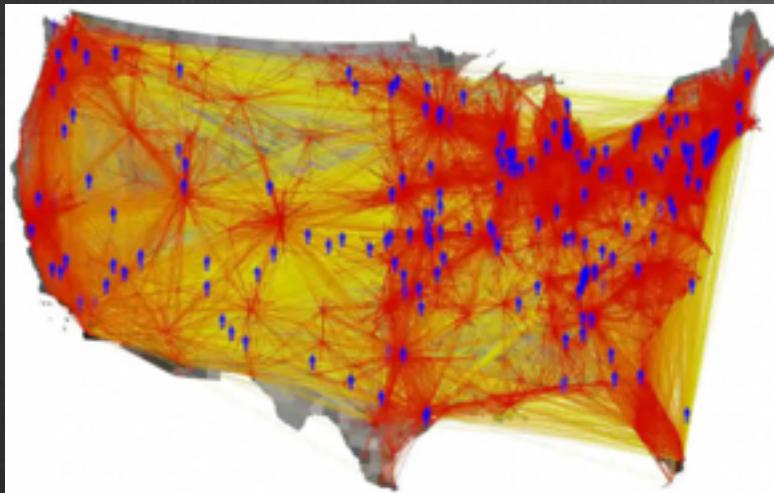
- (highly detailed individual based simulation– household+work+mobility detailed representation)
 - Eubank&Barrett, Longini, Halloran, Ferguson, Burke, Merler and the FBK

Challenges

- Data hungry models,
- Lack of information=assumption
- Set of parameters usually not available initially (vaccination rate, pre-exposure immunity etc.)

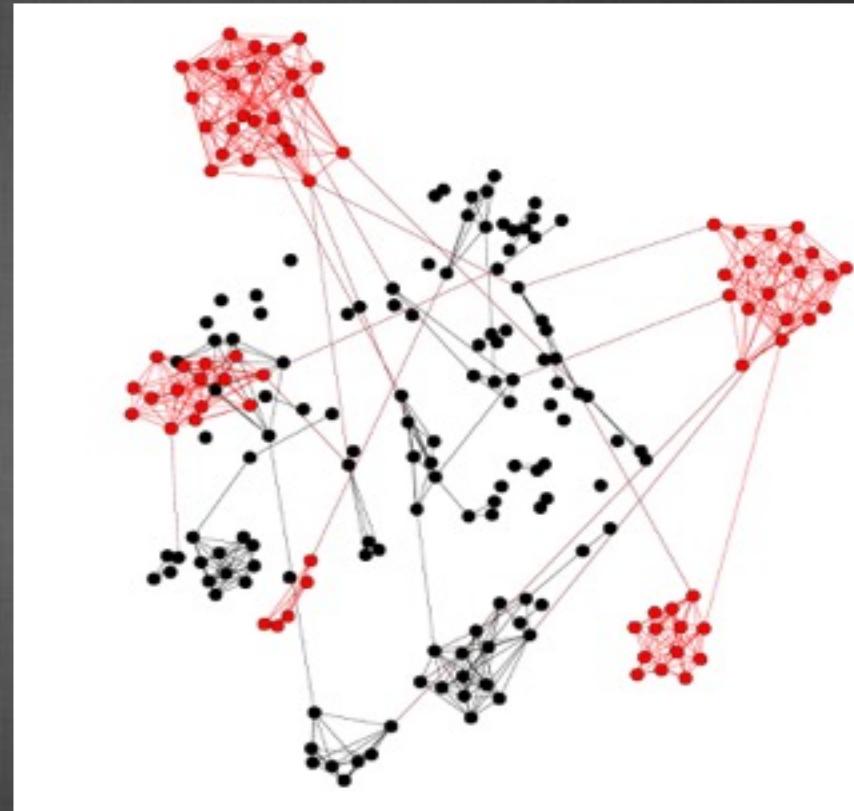
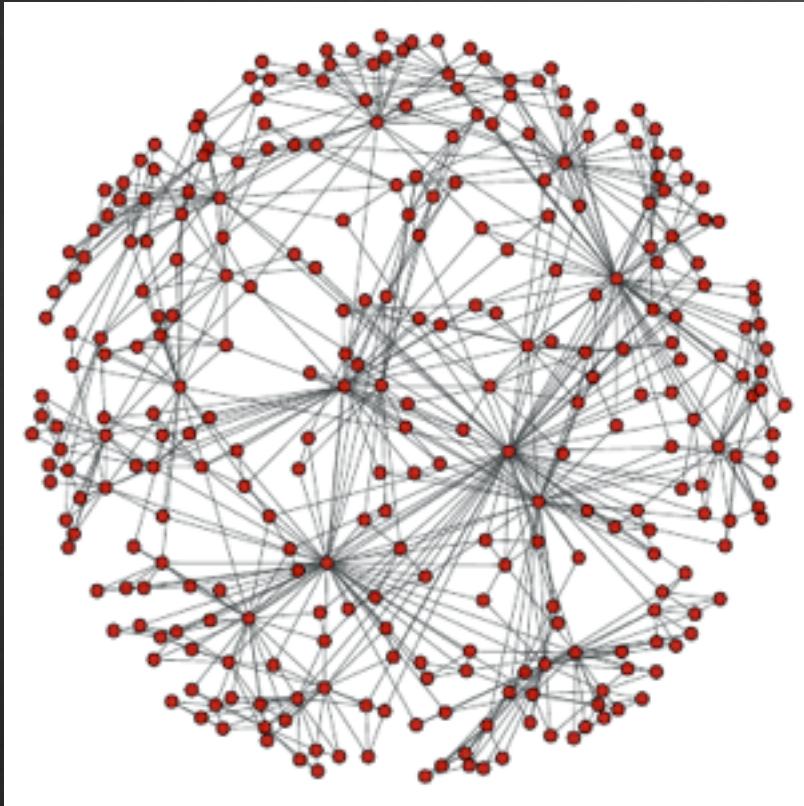


FROM GEOGRAPHY TO SOCIAL SPACE



Geographical areas/census
Mobility

FROM GEOGRAPHY TO SOCIAL SPACE



**Structured communities in
the abstract social space
define by knowledge and
information**

TWITTER NETWORK ON MAY 15TH PROTEST MOVEMENT IN SPAIN

Moreno et al. (BIFI, Universidad de Zaragoza)



TWITTER NETWORK ON MAY 15TH PROTEST MOVEMENT IN SPAIN

Moreno et al. (BIFI, Universidad de Zaragoza)



TRANSFORMING THE WAY WE APPROACH SOCIO-TECHNICAL SYSTEMS

- Social computational science
- Computational epidemiology
- Science of Science
- Information systems and data science
 - General classes of tools methods and models generally applicable to complex techno-social systems:
 - New techniques for the generation of models, social analytics, real time empirical data generation.
 - Scenario analysis: systemic risk, economic impact etc.