



SoBigData

Research Infrastructure



Data Science
Opportunities, Risks, Capabilities

Salvatore Rinzivillo




ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"





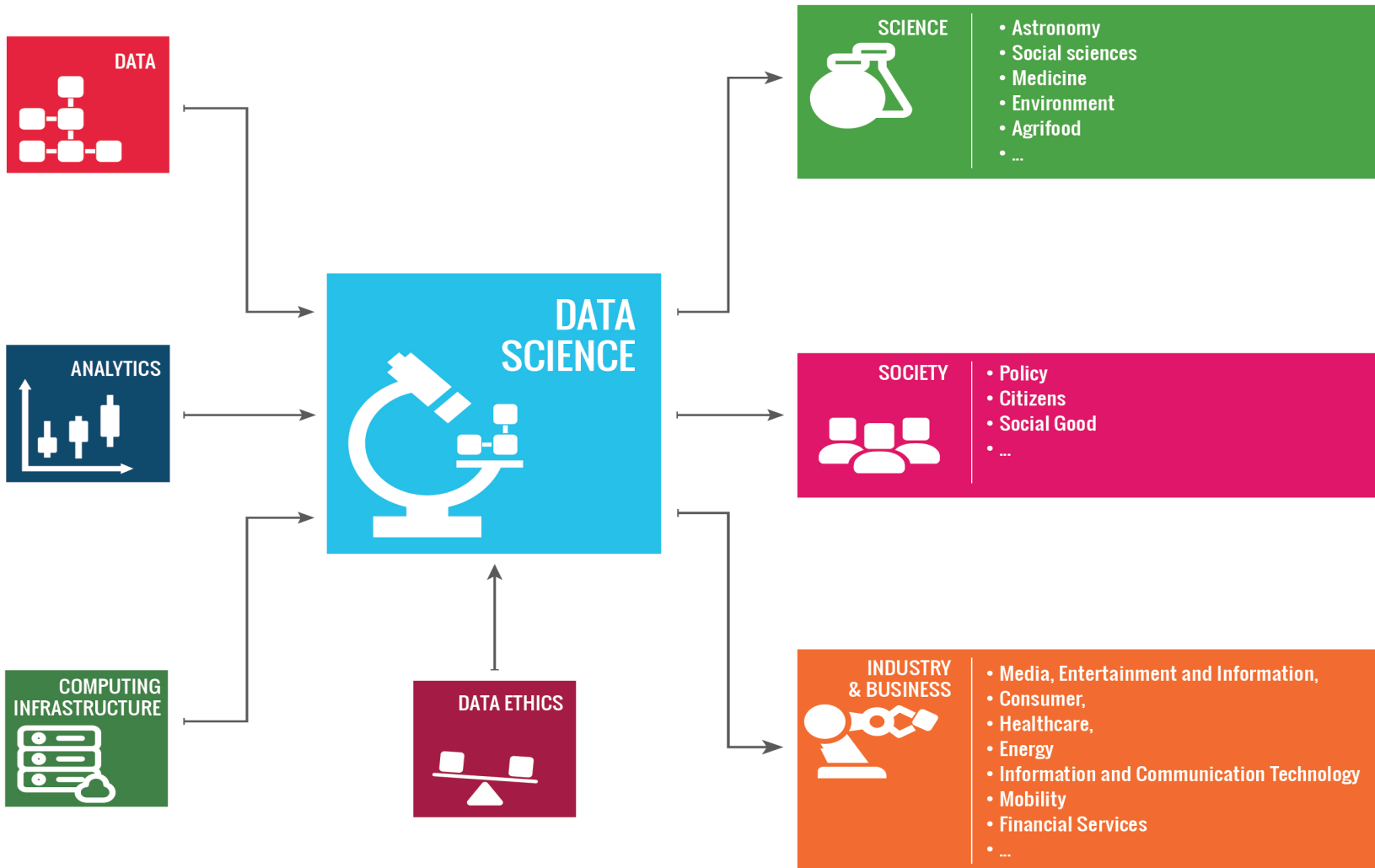
Data Science and BigData: a Game-changer for Science and Innovation

Document for G7 Academy, March 2017, authored by

- Fabio Beltram: Scuola Normale, Pisa.
 - Fosca Giannotti: Istituto Scienza e Tecnologie dell'Informazione, CNR, Pisa.
 - Dino Pedreschi: Dipartimento di Informatica, Univ. Pisa, Pisa
- 

What is data science?

data availability, sophisticated analysis techniques, and scalable infrastructures brought what we call today “Data Science”

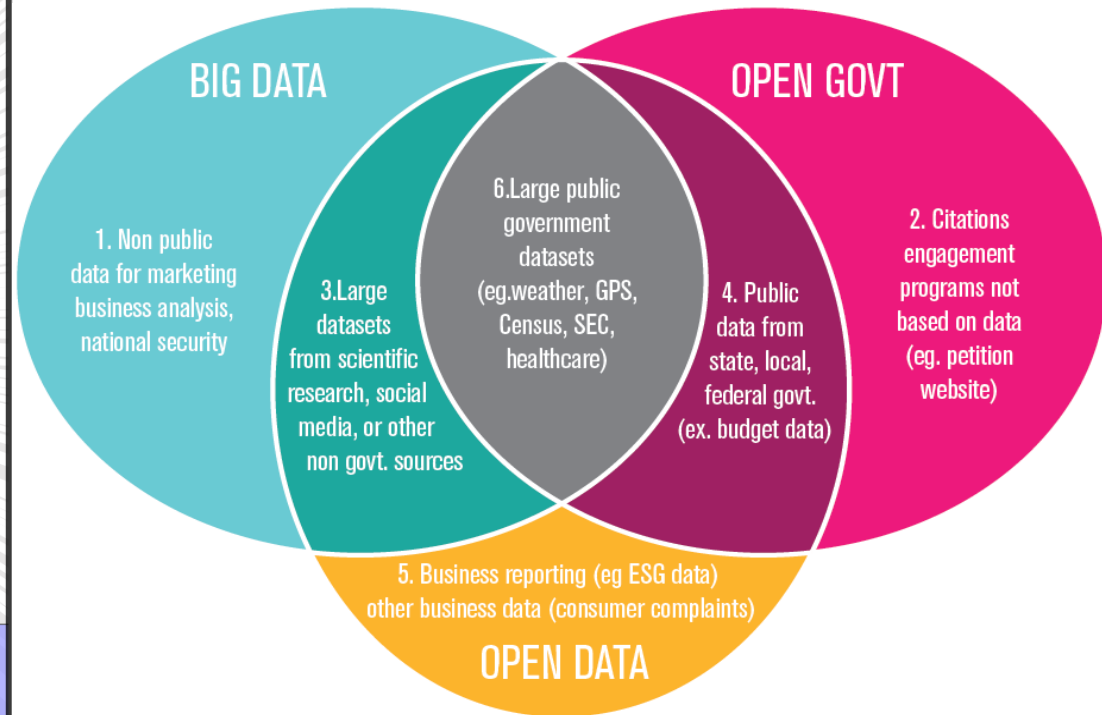
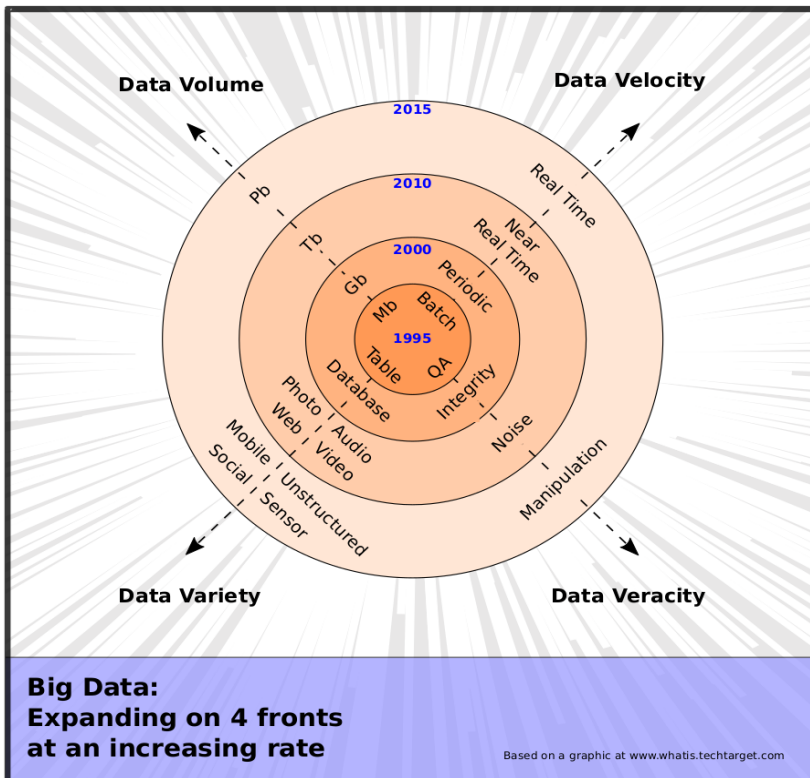


What is data science?

“Data science includes data extraction, data preparation, data exploration, data transformation, storage and retrieval, computing infrastructures, various types of data mining and statistical learning, presentation of explanations and predictions, and the exploitation of results taking into account ethical, social, legal, and business aspects.”

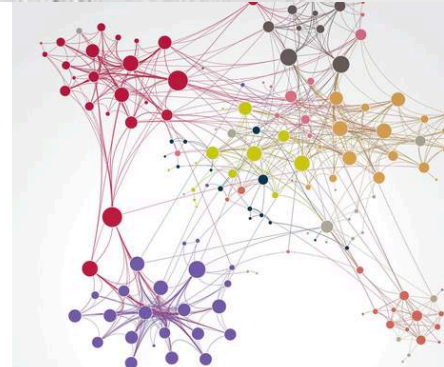
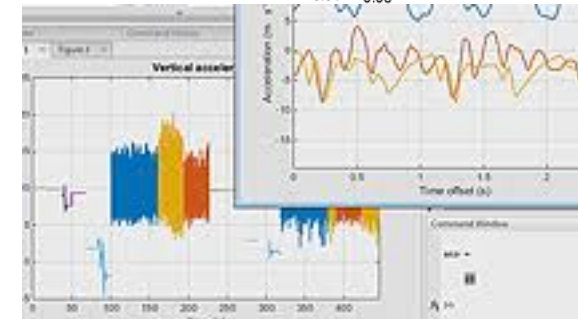
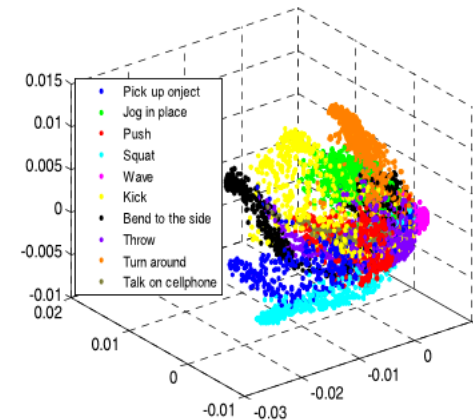
The data.

Data may be structured or unstructured, big or small, static or streaming.



The analytics

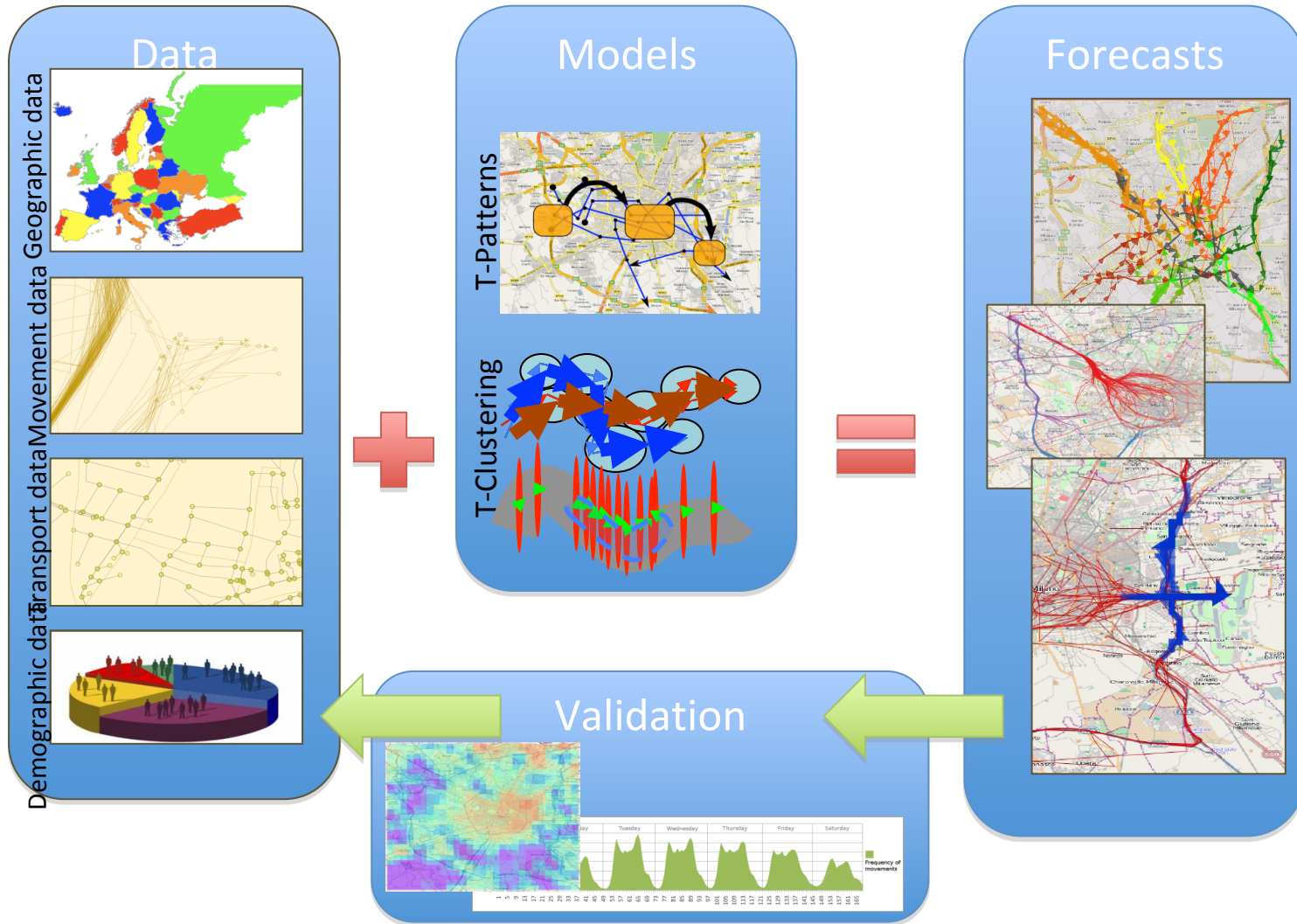
- **Data mining algorithms** for automated pattern discovery highlight the structure hidden in massive datasets.
- **Machine learning - nowadays “deep learning”** - methods exploit large “training” datasets of examples to learn general rules and models to classify data and predict outcomes,
- **Network science** has unveiled the magic of shifting from the statistics of populations to the statistics of interlinked entities, connected by the ties of their mutual interactions;

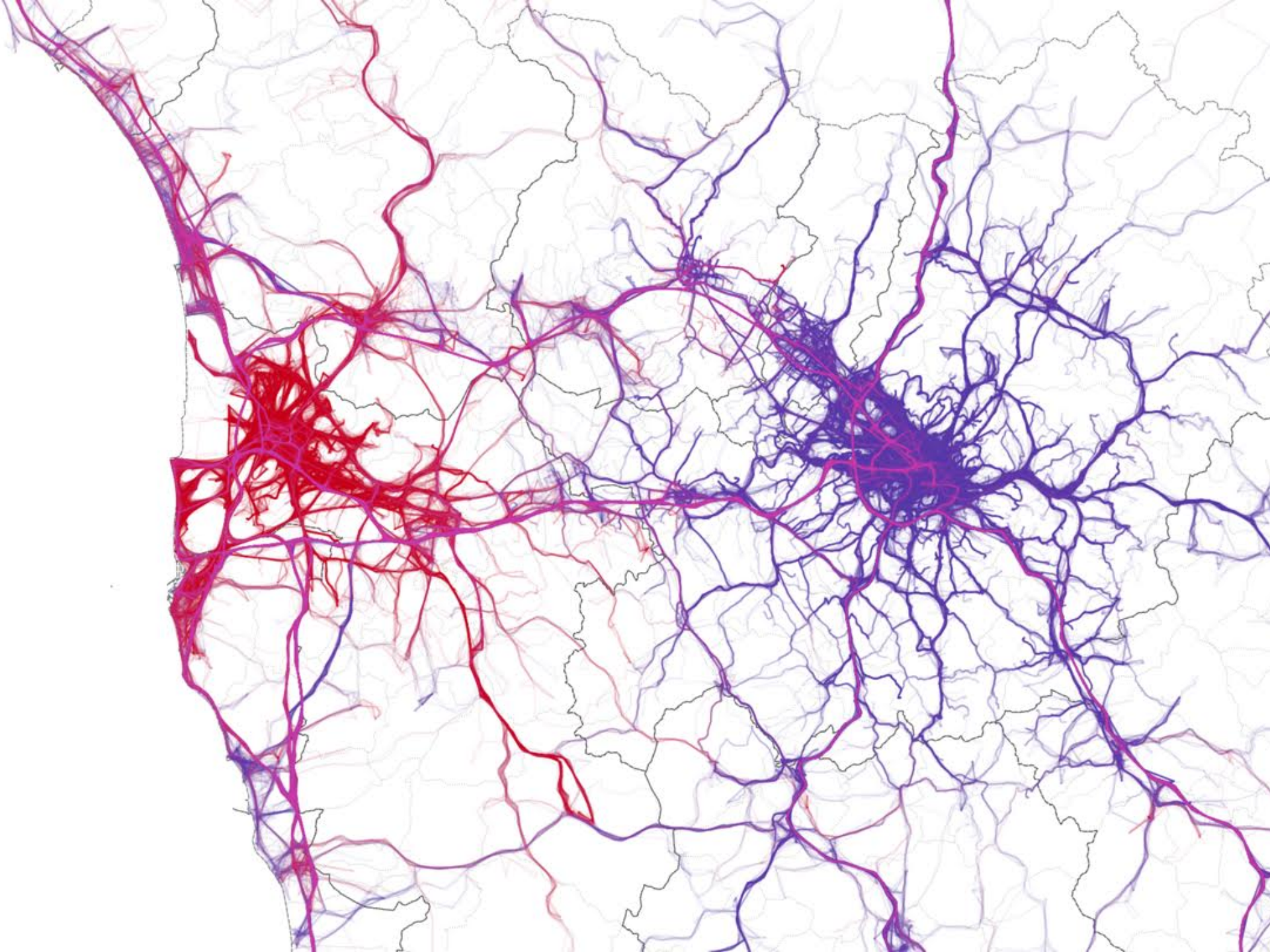


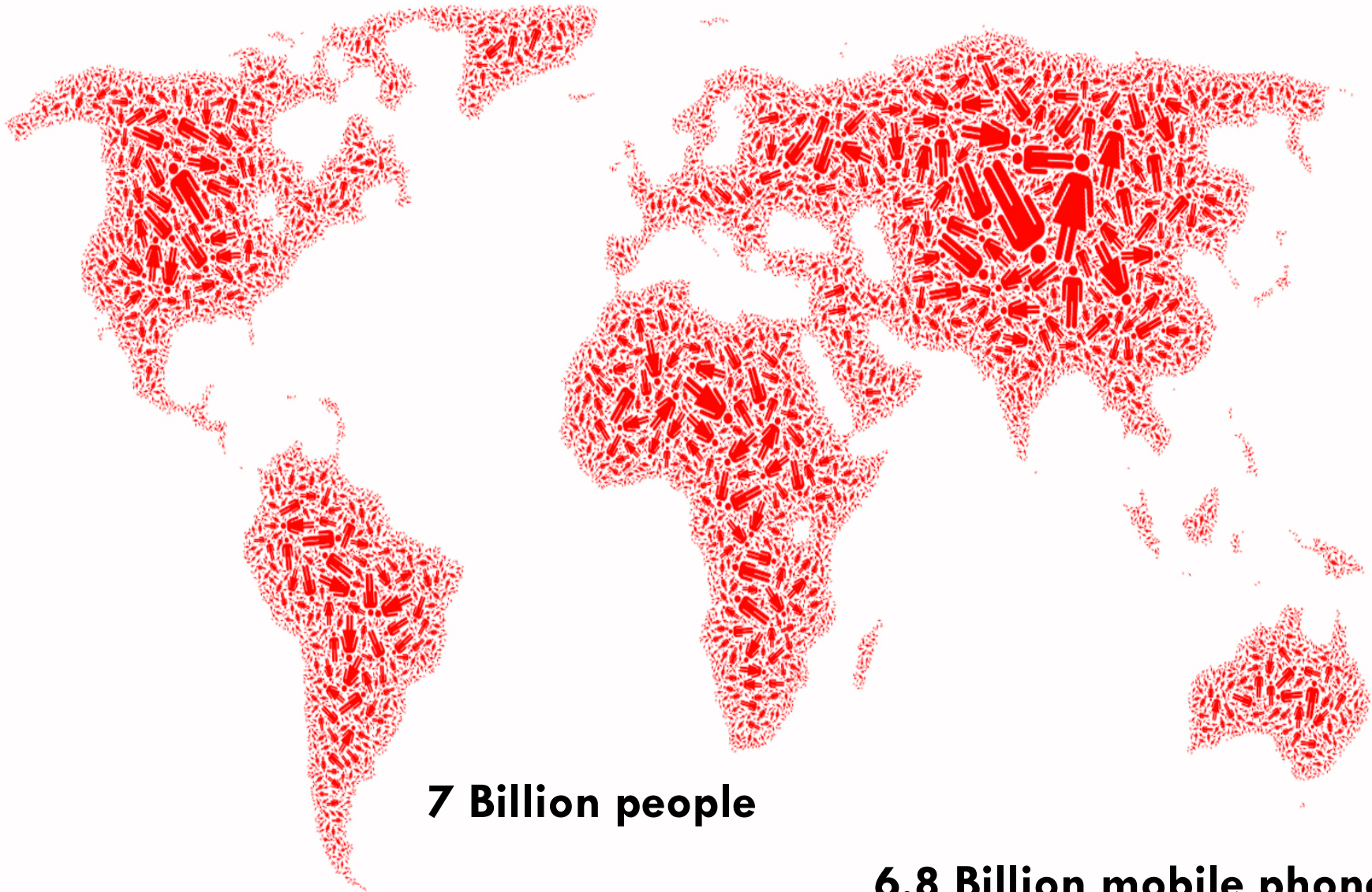
Albert-László Barabási

**NETWORK
SCIENCE**

From DATA to KNOWLEDGE







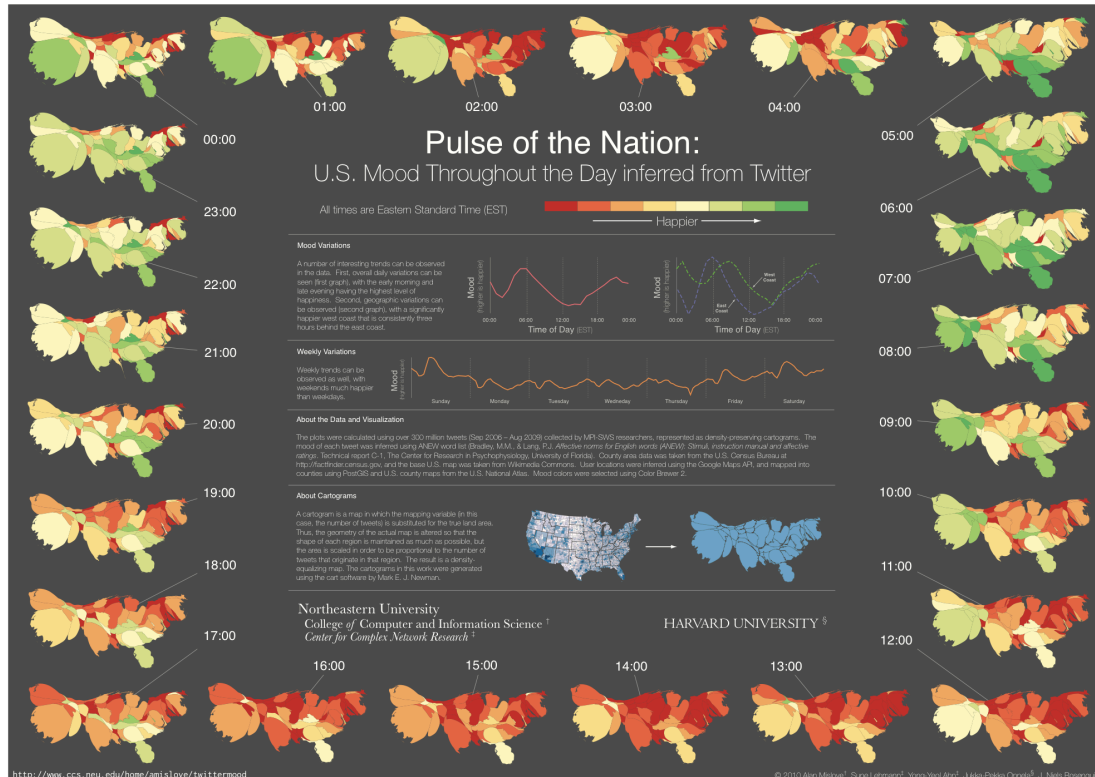
7 Billion people

6.8 Billion mobile phones

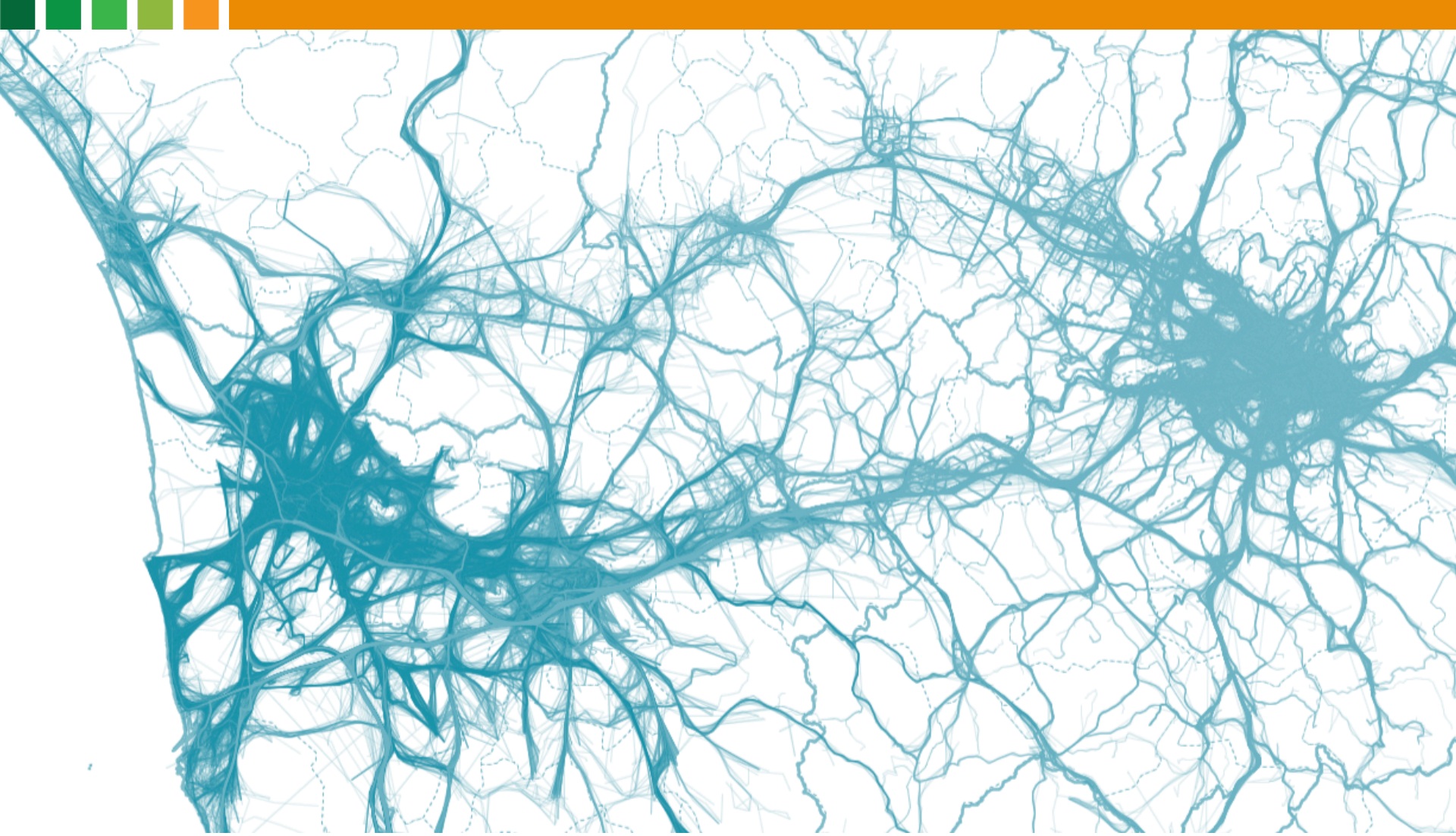




Measuring happiness with twitter



- **Computational social science** “*is now using digital tools to analyze the rich and interactive lives we lead to answer questions that were previously impossible to investigate*”. (Mann. PNAS January 19, 2016, vol. 113 no. 3)



URBAN MOBILITY ATLAS

Pisa Pisa



NUMBER OF VEHICLES

5,615

Source: OCTO Telemati...



RESIDENTS

89,694

Source: census 2001



DAYS

31

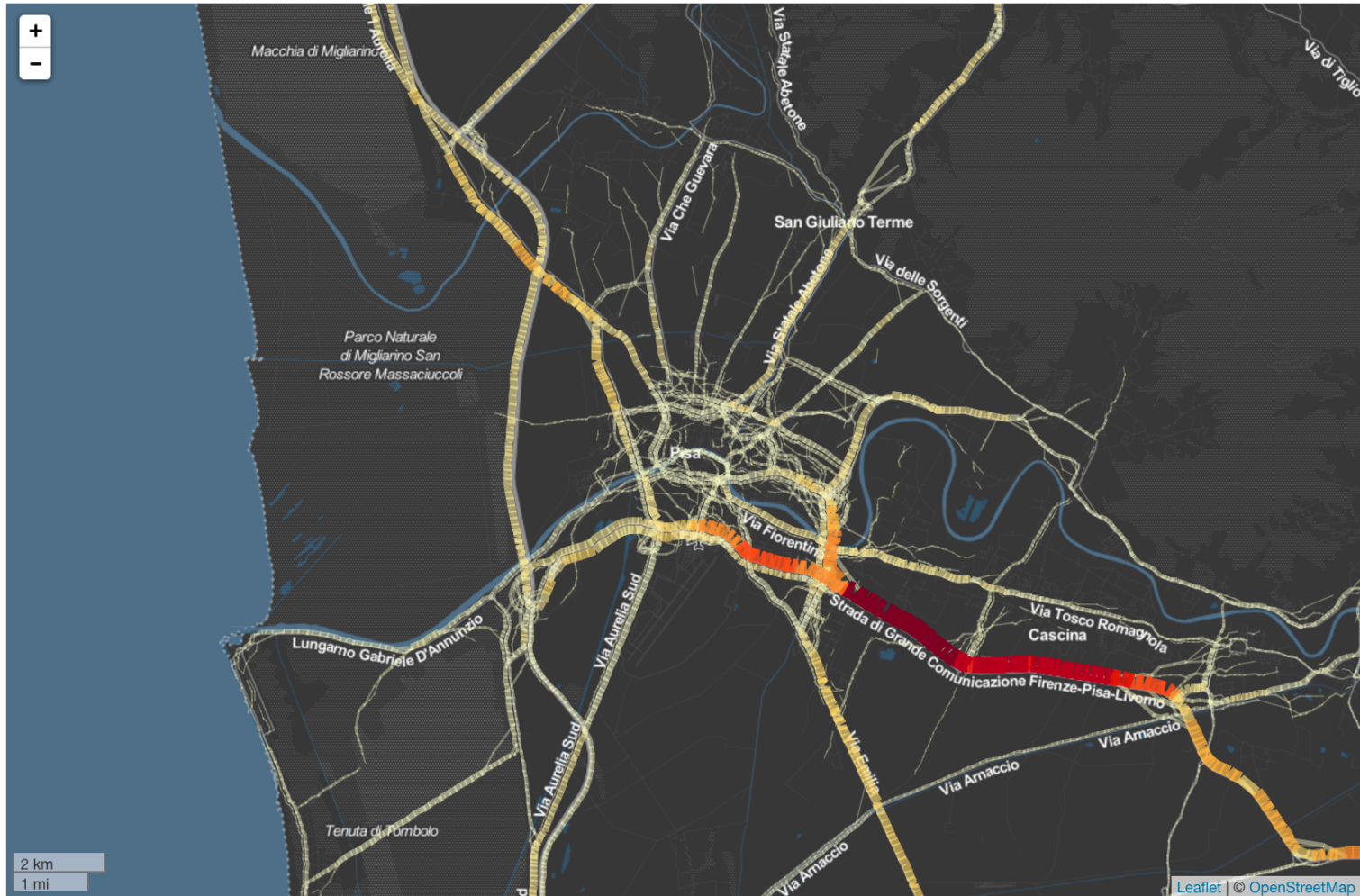
May 1-31, 2011



SYSTEMATICS TRAJS

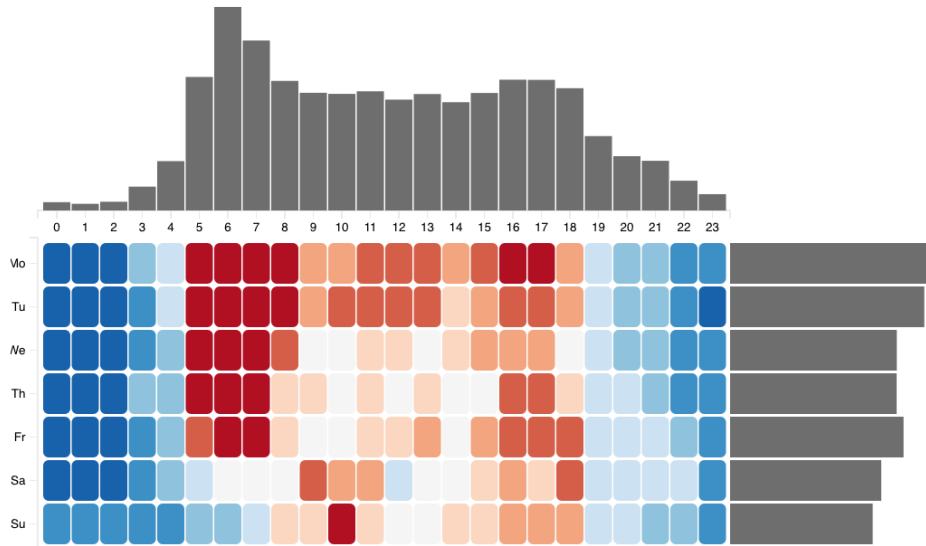
38.88%

Flows and Traffic



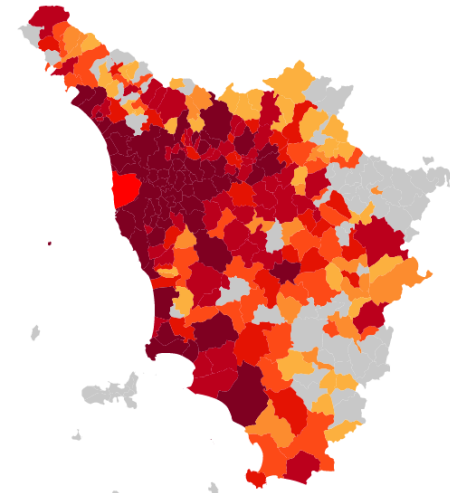
Flows of traffic exiting from the city.

Temporal Matrix



Time distribution of trips entering the city during a typical week. Trips can be filtered by occasional or systematic.

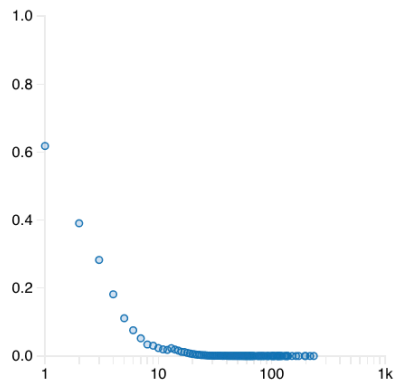
OD Map



Map of origins and destinations

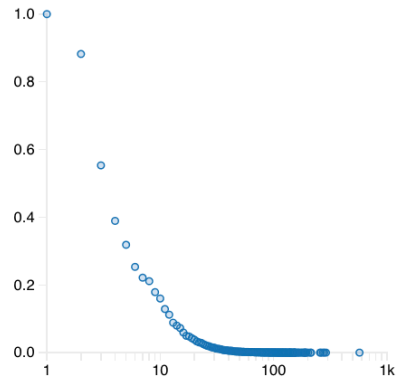
Lengths of Trajs

set of buttons here



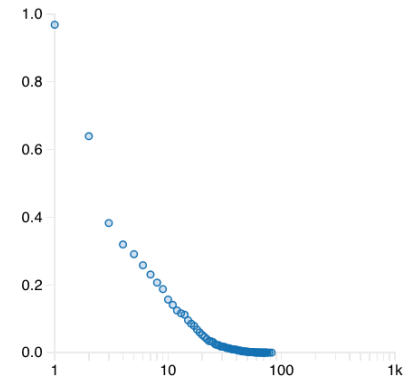
Lengths of Trajs

Duration of Trajs

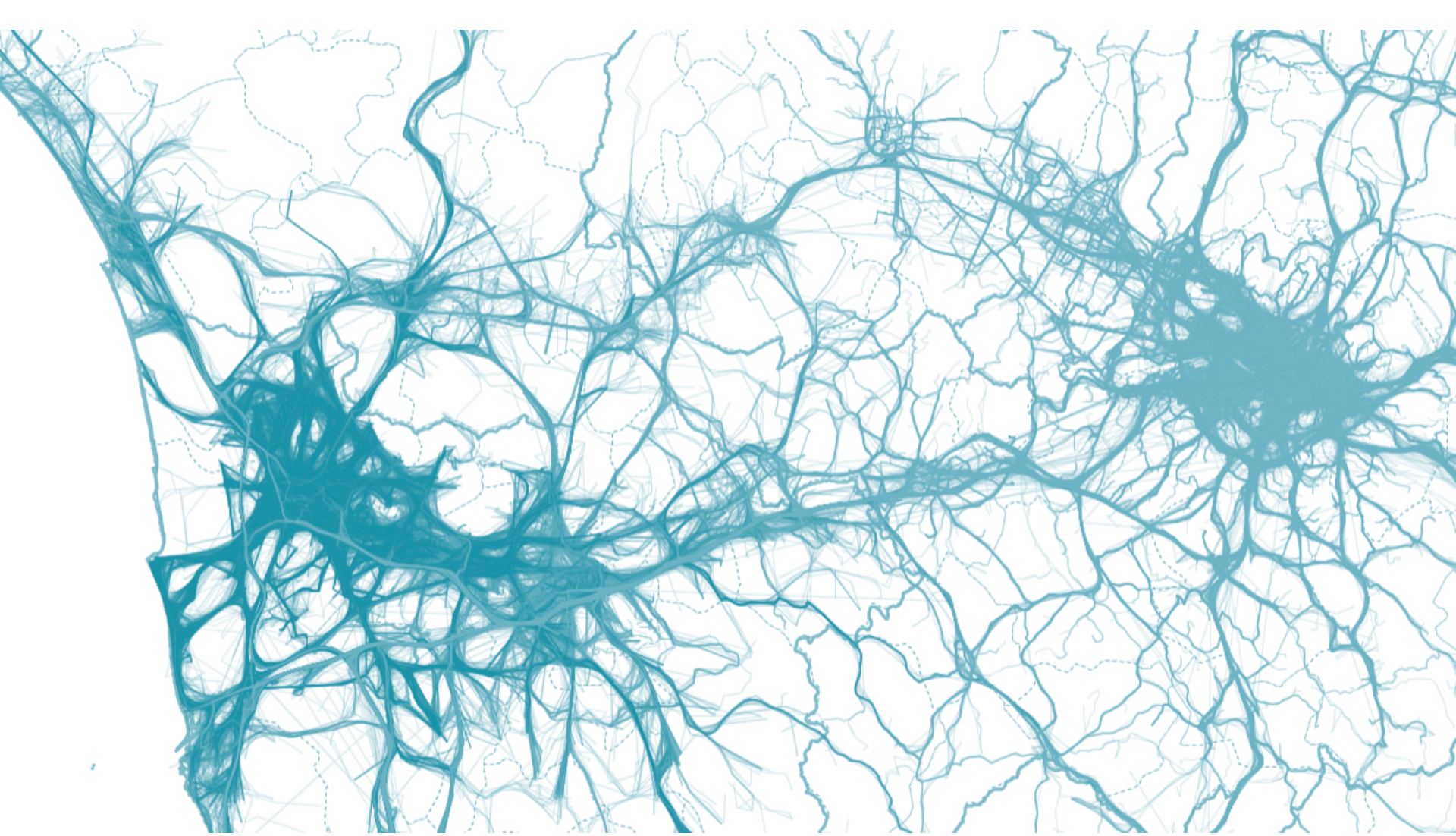


Distribution of durations

Speeds of Trajs

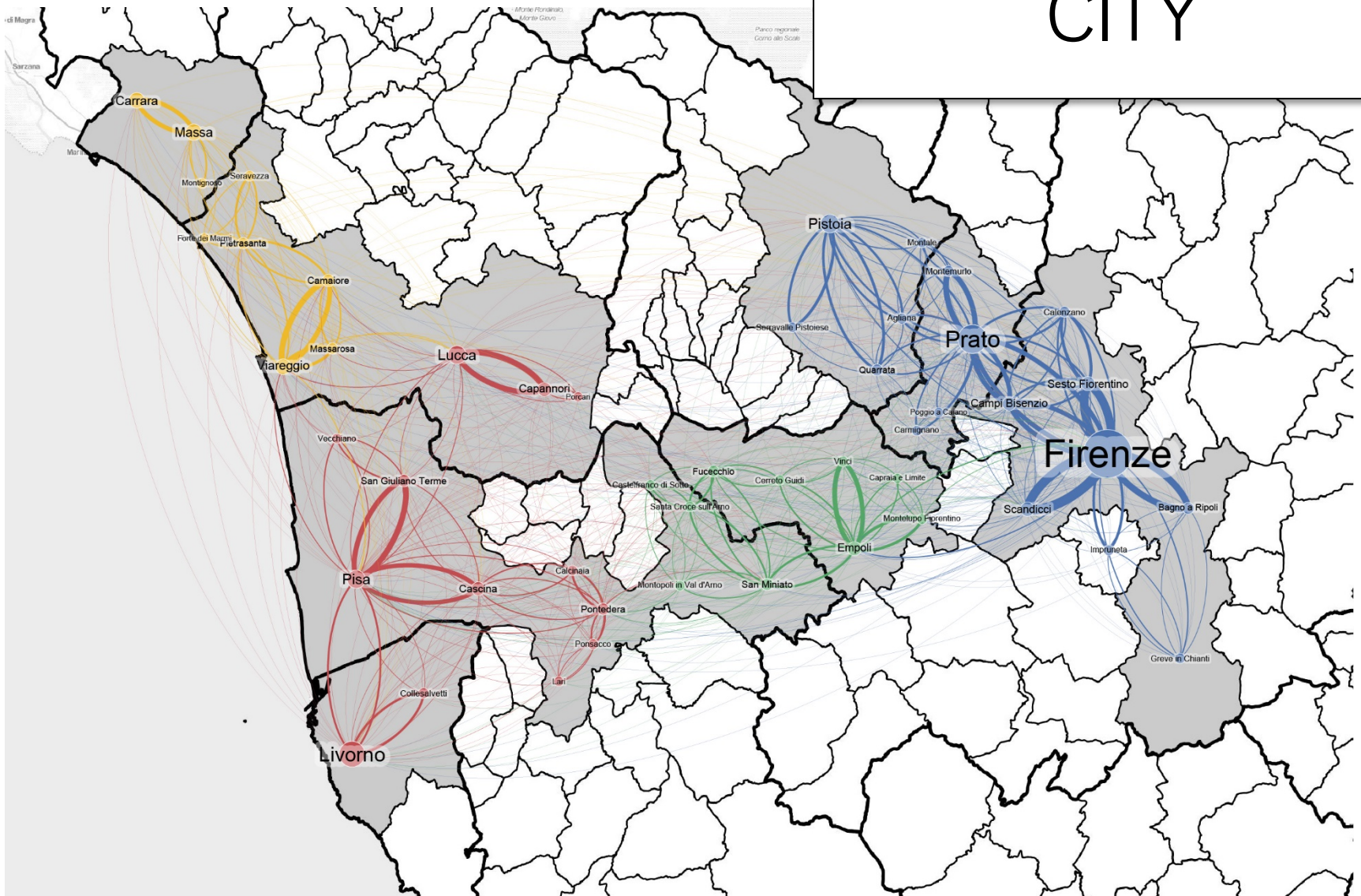


Distribution of speeds



THE POLYCENTRIC CITY

POLYCENTRIC CITY

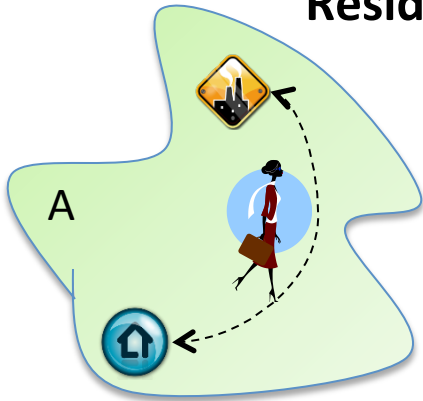




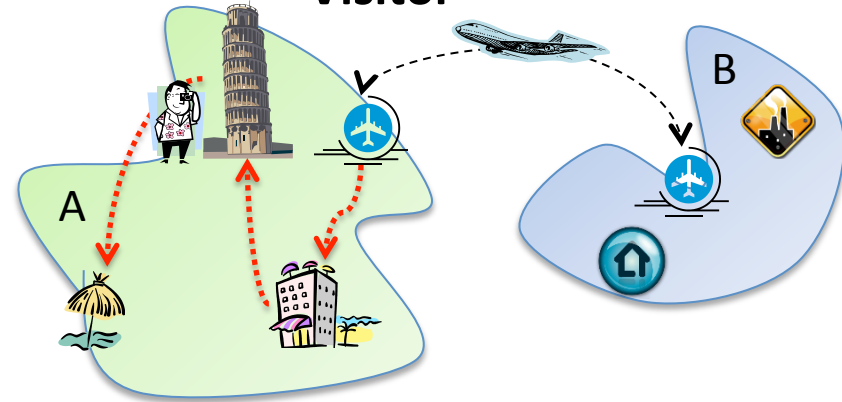
REAL TIME DEMOGRAPHY

Sociometer: Estimating User Category from mobile phone

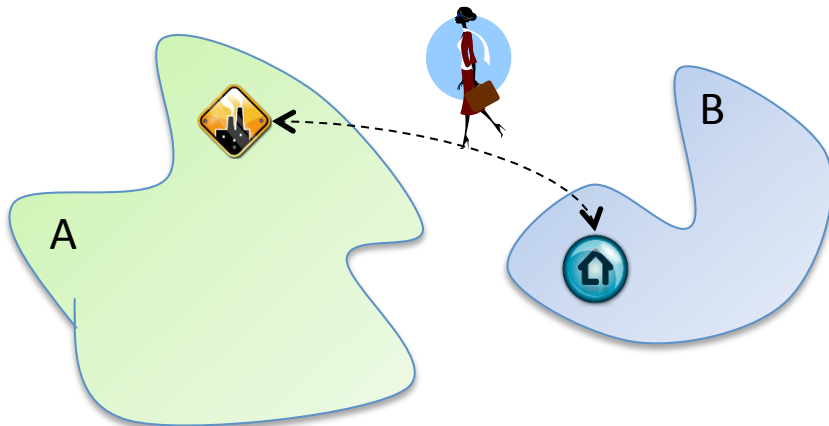
Resident



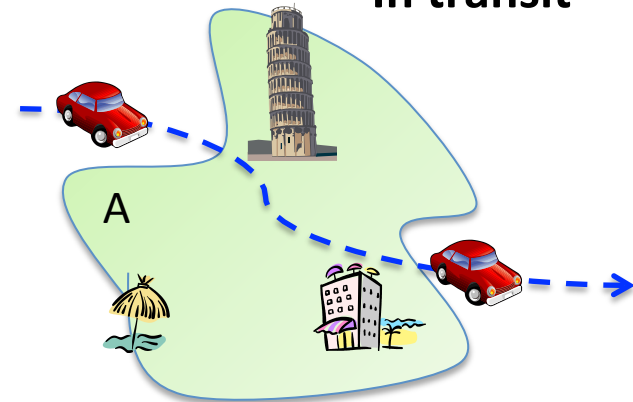
Visitor

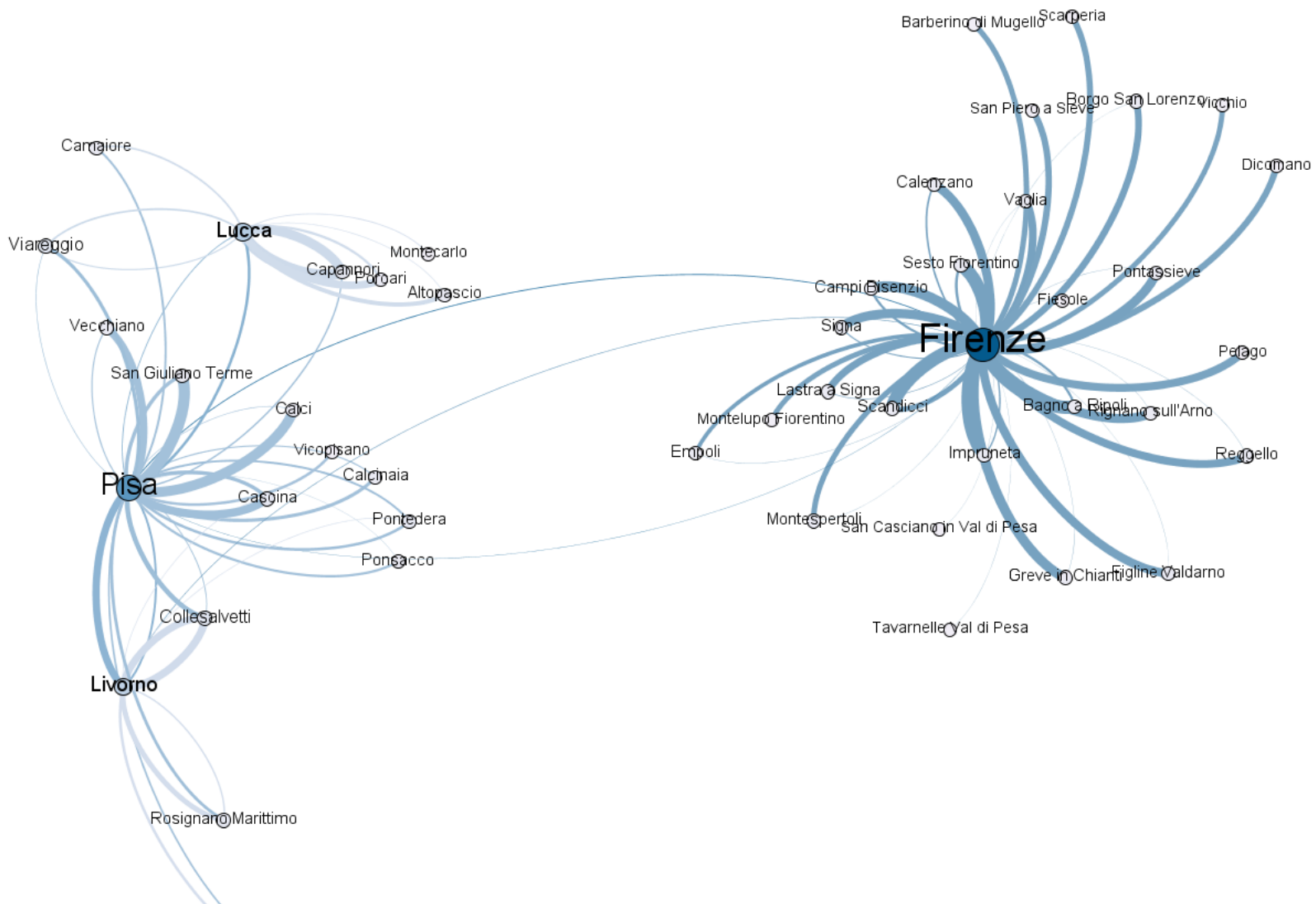


Commuter



In transit

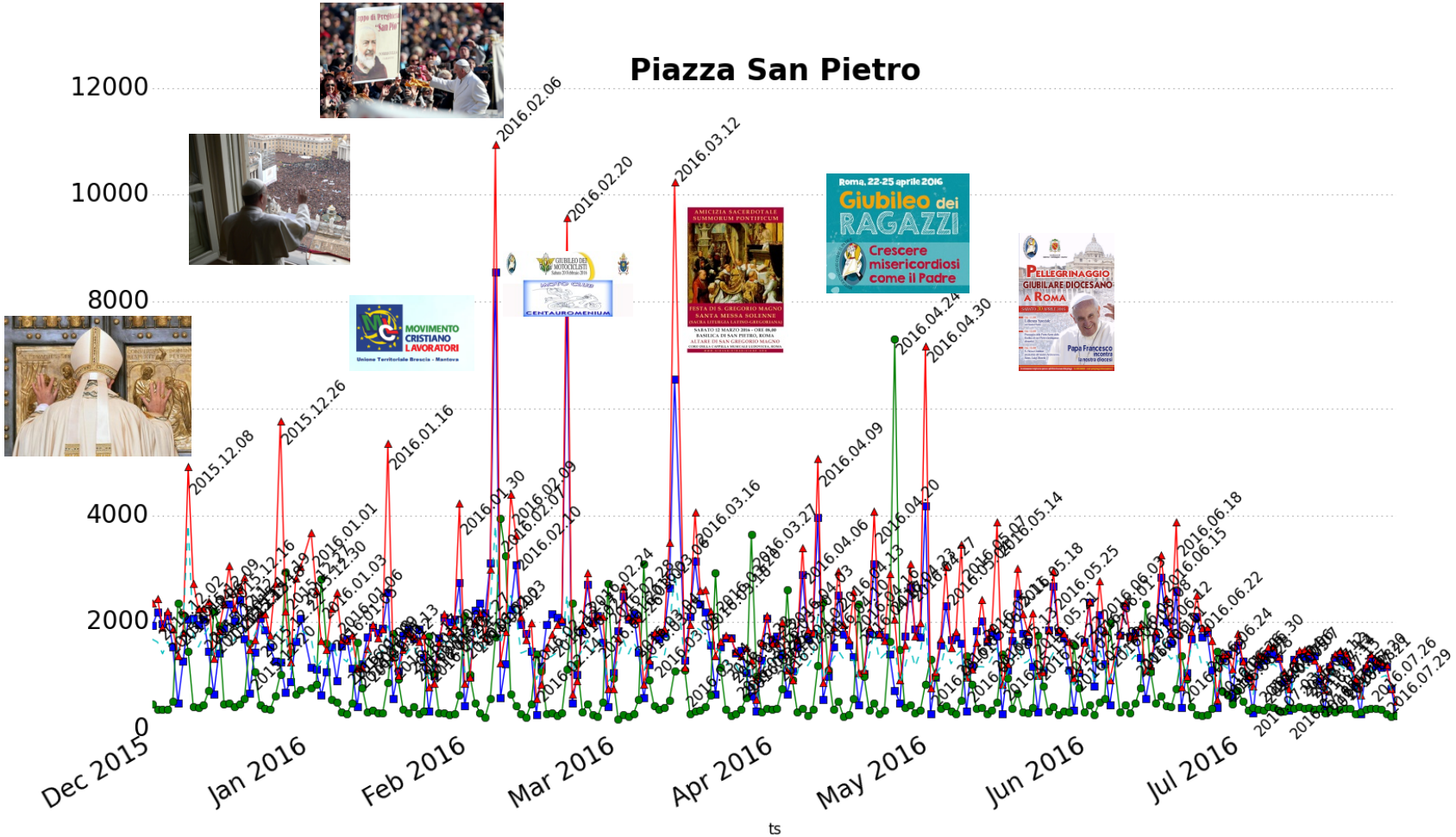


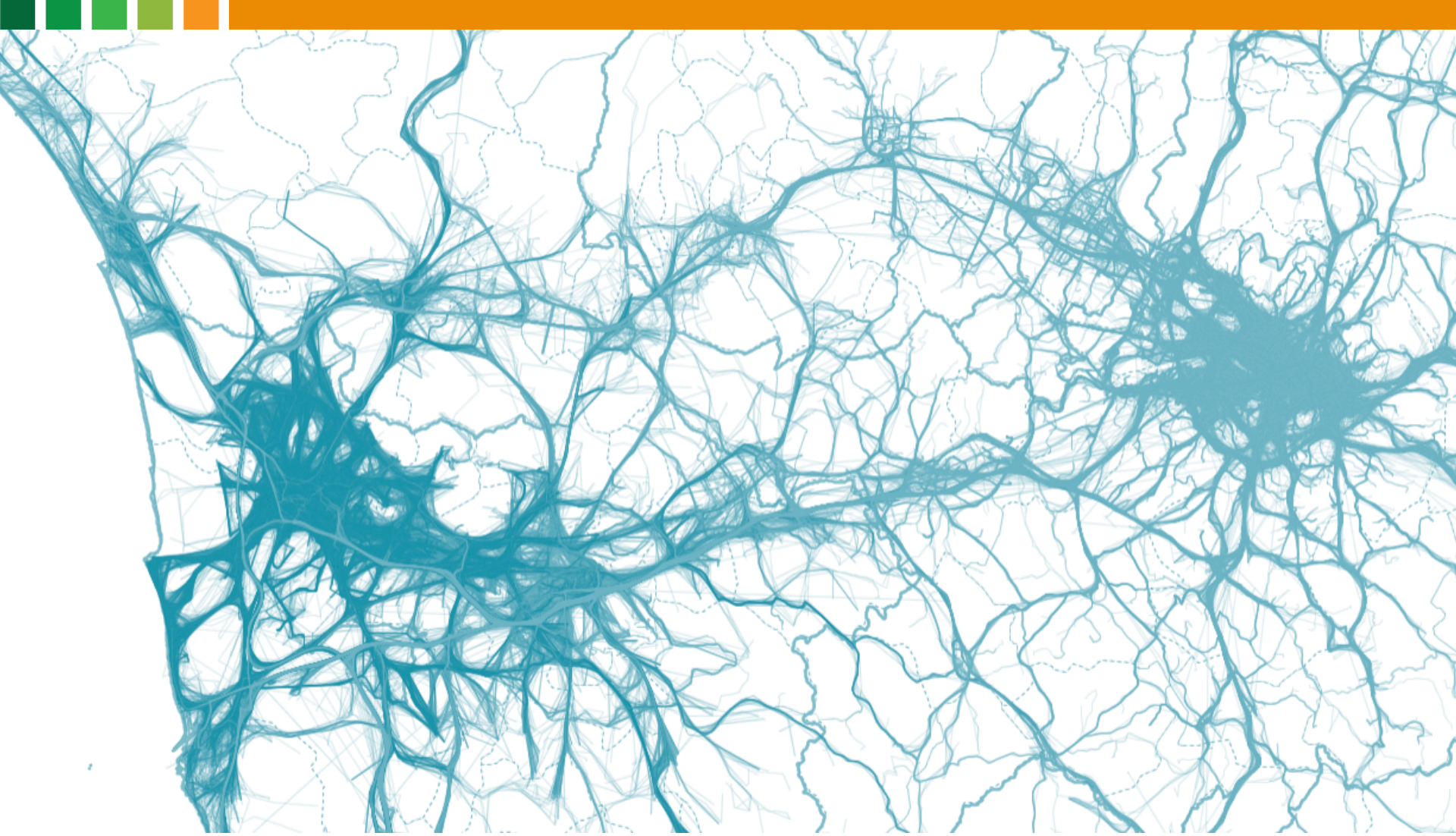


San Pietro Square

- commuter
- ▲ passingby
- visitor
- dynamic_resident

Piazza San Pietro

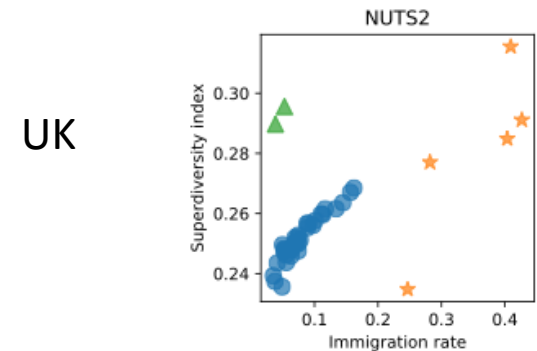
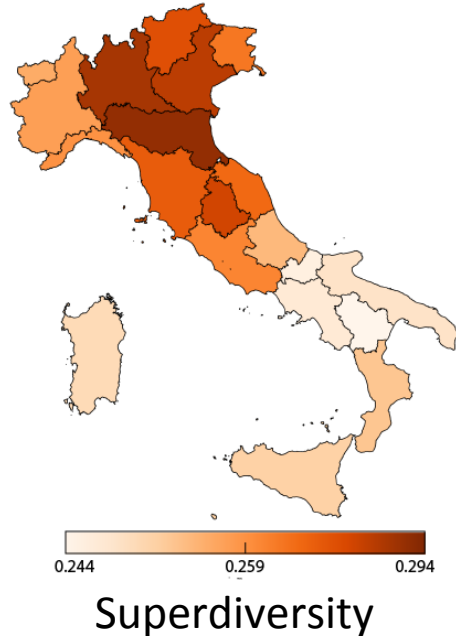
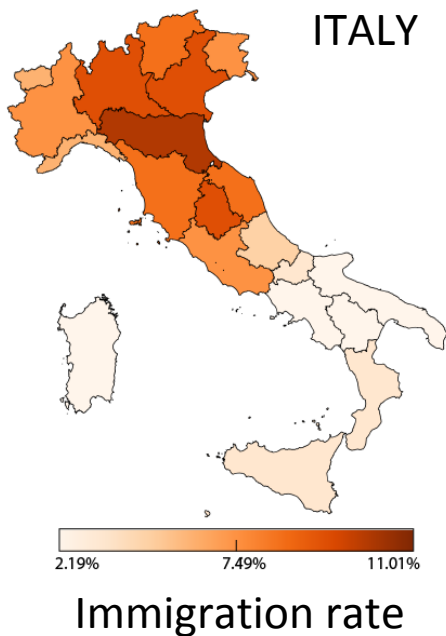
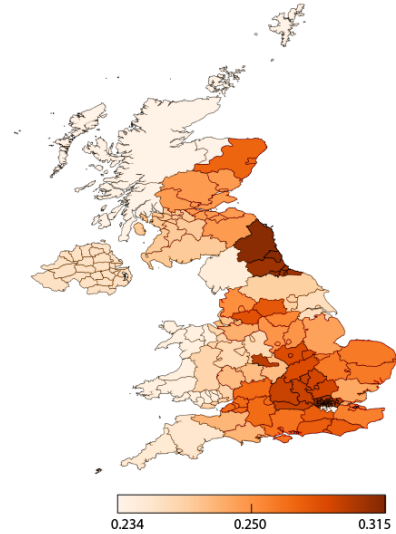
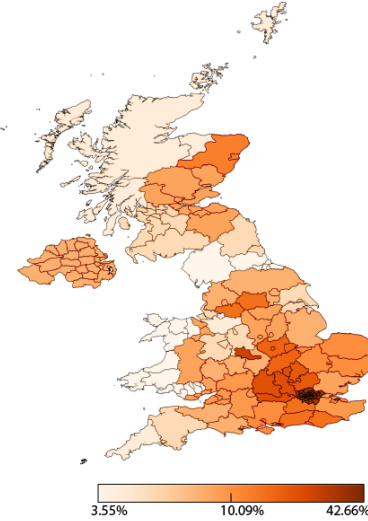
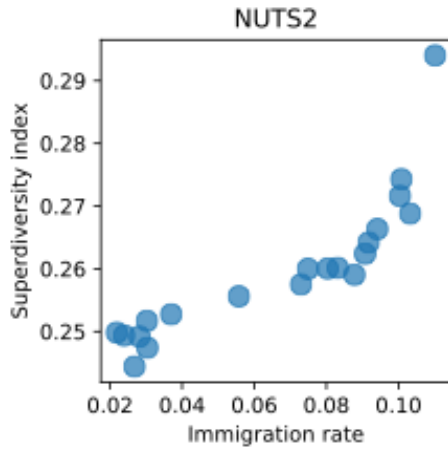




DIVERSITY & WELLBEING



Migration and Superdiversity on Twitter



- Immigration rates: JRC D4R data challenge
- Superdiversity: distance between the emotional content of words in the standard language and on Twitter

Societal Debates

- By analysing millions of datasets of public debates on social media and newspaper articles, it is possible to understand which are the most discussed topics, how they emerge and evolve in time and space and how opinions polarize.

3 Million Brexit Tweets Reveal Leave Voters Talked About Immigration More Than Anything Else

Groundbreaking analysis shows immigration, not sovereignty or the NHS, dominated the conversation – and making British judges responsible for British law was a key theme for supporters.



James Ball
BuzzFeed Special
Correspondent



Chris Applegate
Editorial Developer, UK

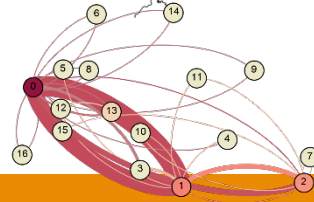
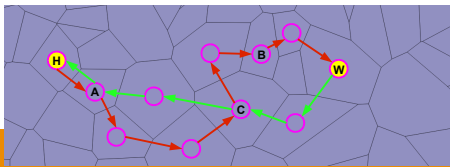
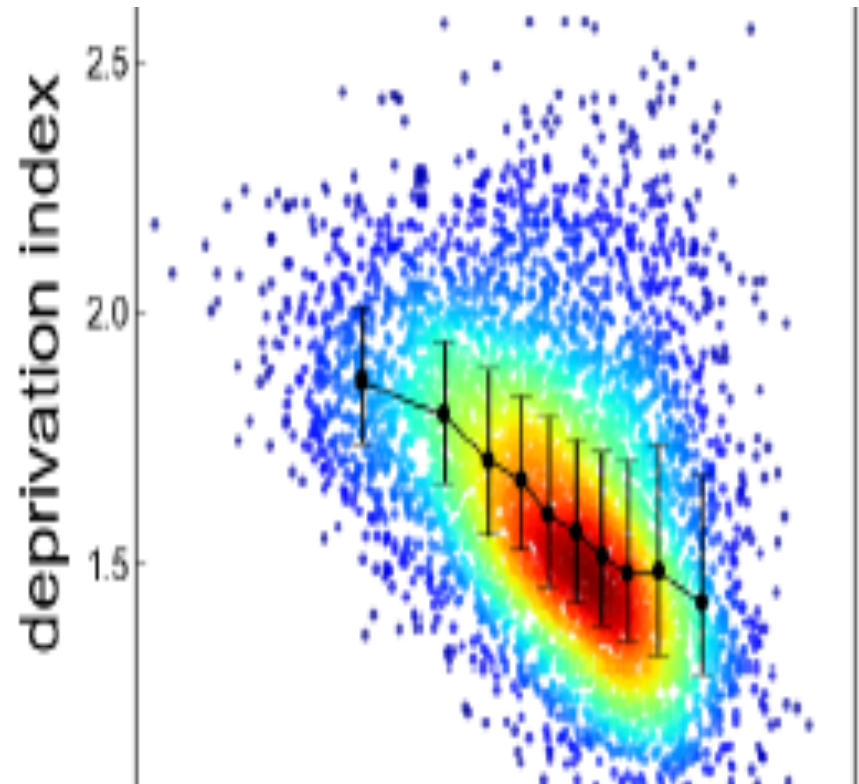
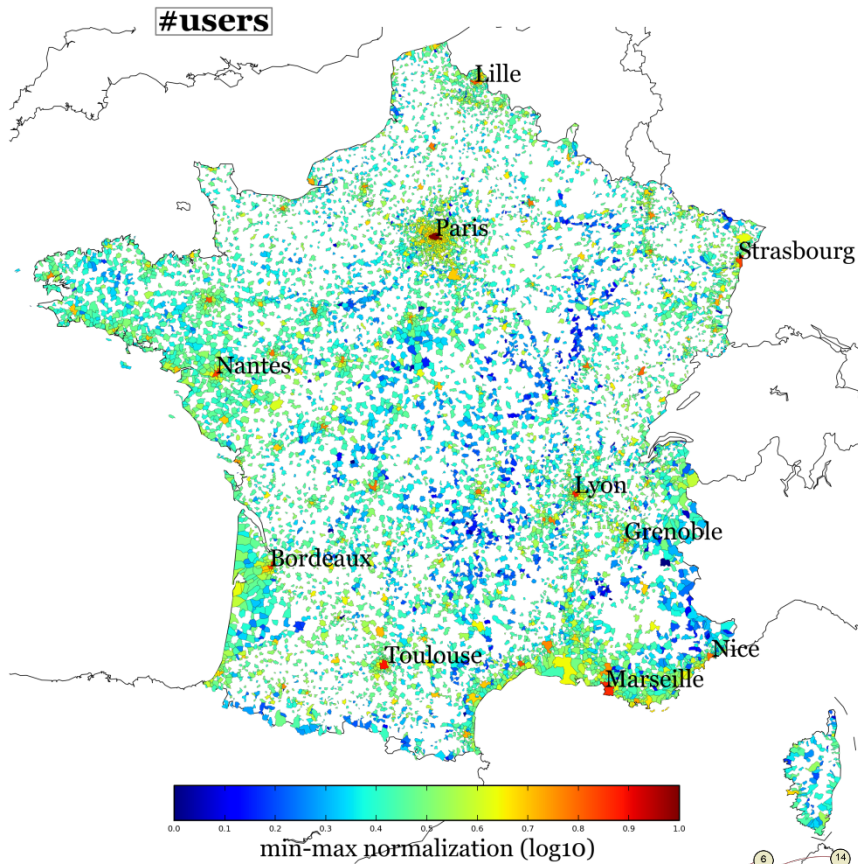
posted on Dec. 9, 2016, at 2:03 p.m.



Big Data: Diversity and economic development



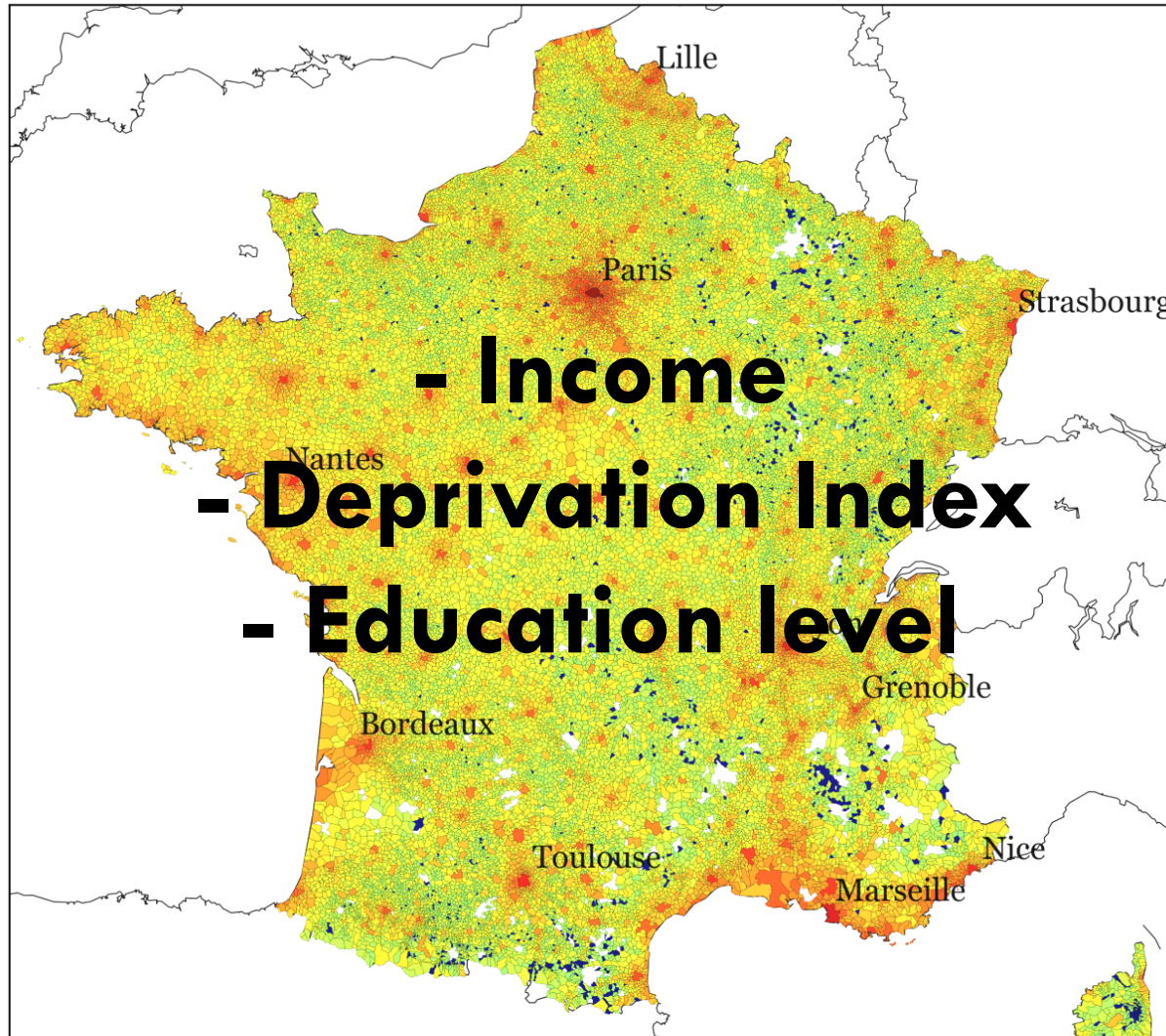
Mobility Diversity and Wellbeing



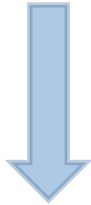
$$-\sum_{T'_i \subset T_i} P(T'_i) \log_2 [P(T'_i)]$$

Economic Measures

7,000 French cities



20 million users
200 million calls



hadoop
user filtering



6 million users
mobility trajs
social network



Four individual measures

- Radius of gyration
- Social degree

volume

- Mobility entropy
- Social diversity

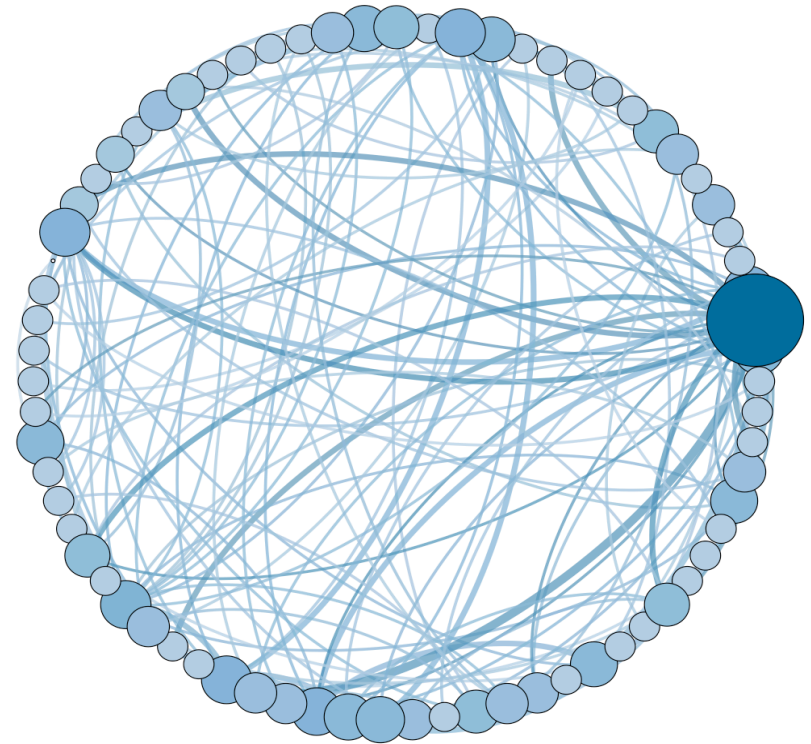
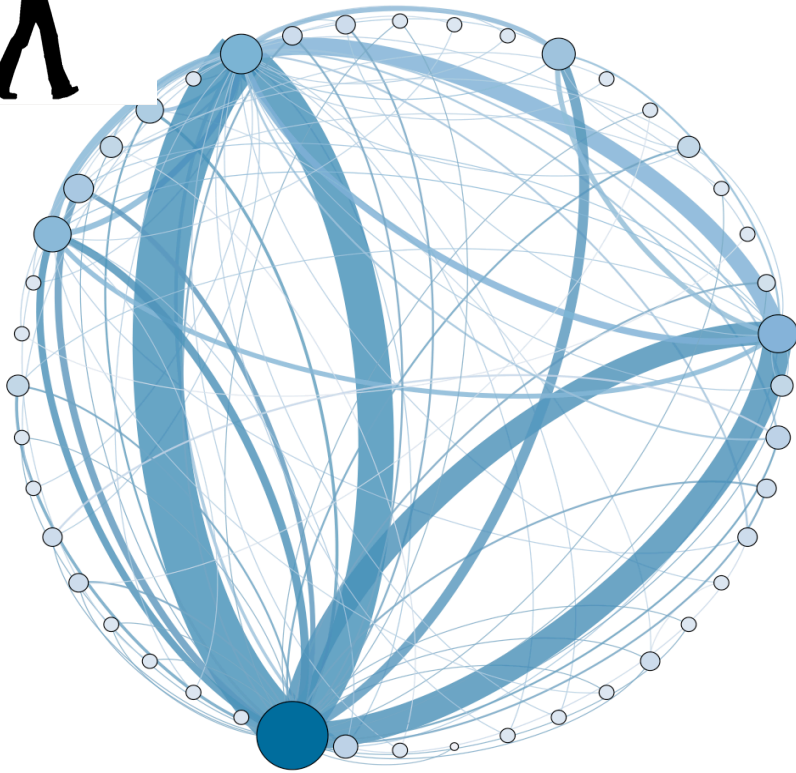
diversity

$$-\sum_{T'_i \subset T_i} P(T'_i) \log_2 [P(T'_i)]$$

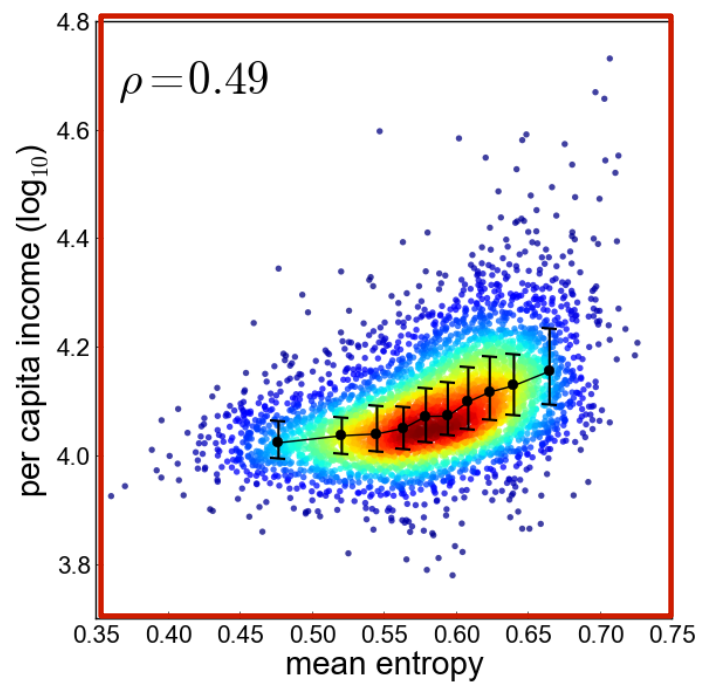
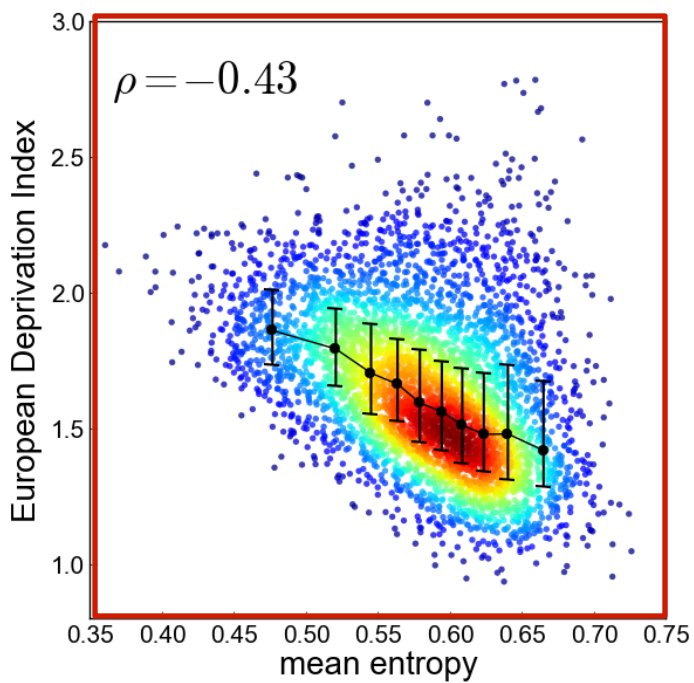
Diversity of individual mobility networks



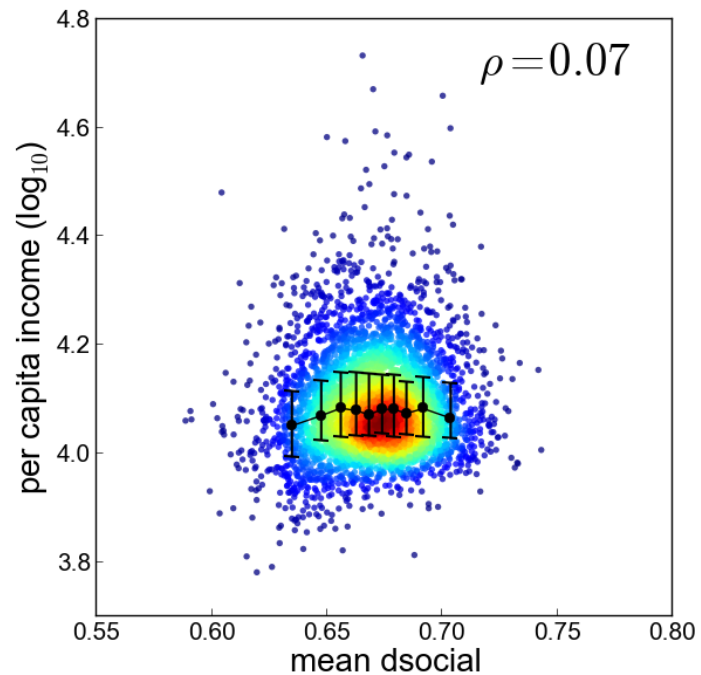
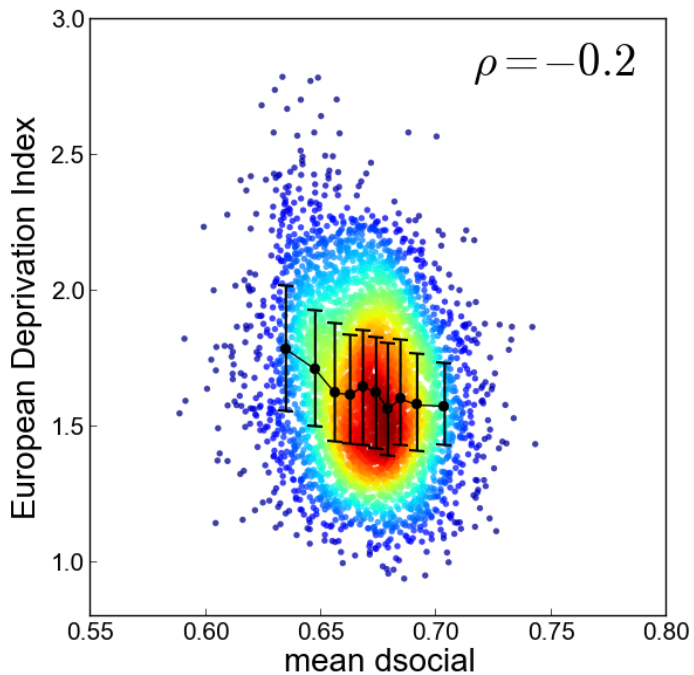
dreamstime.com



mobility

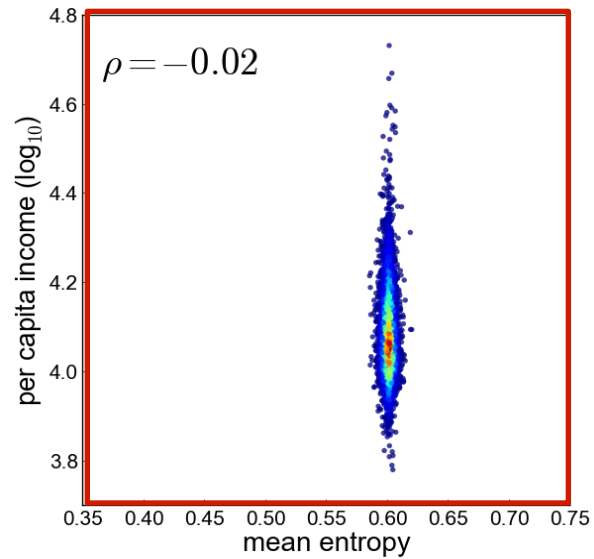
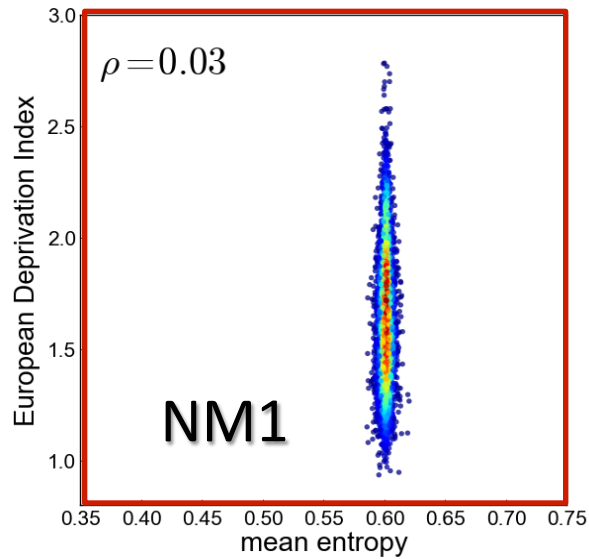


sociality

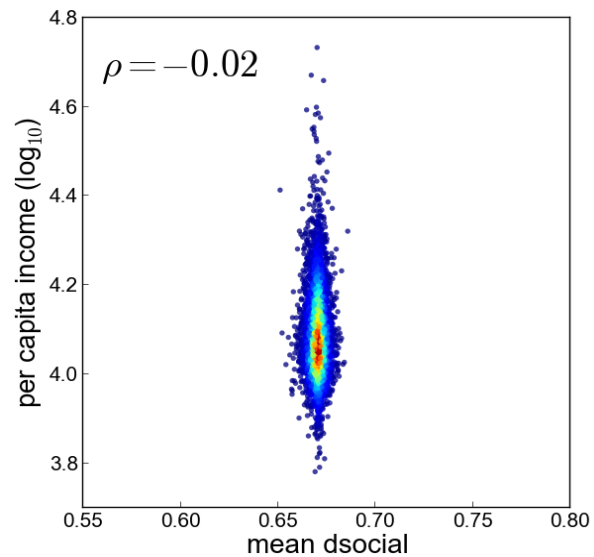
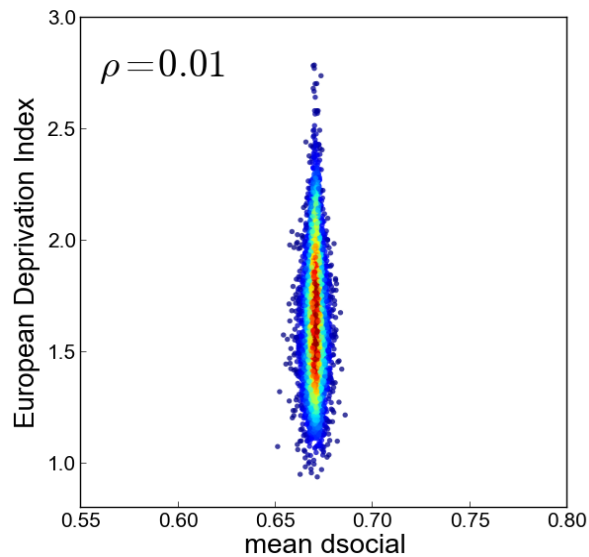


What on a null model: randomizing people

mobility



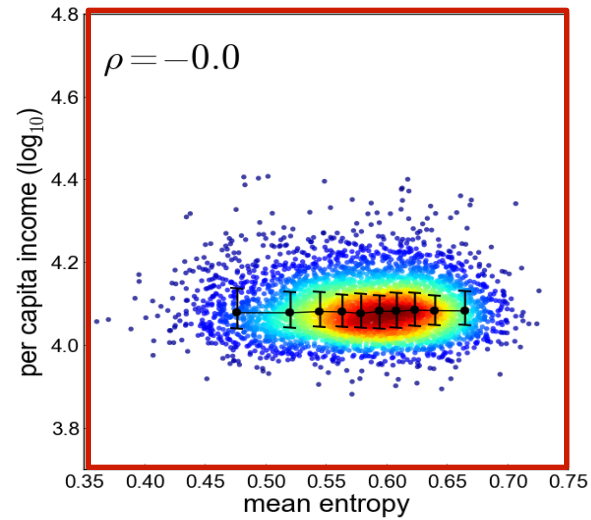
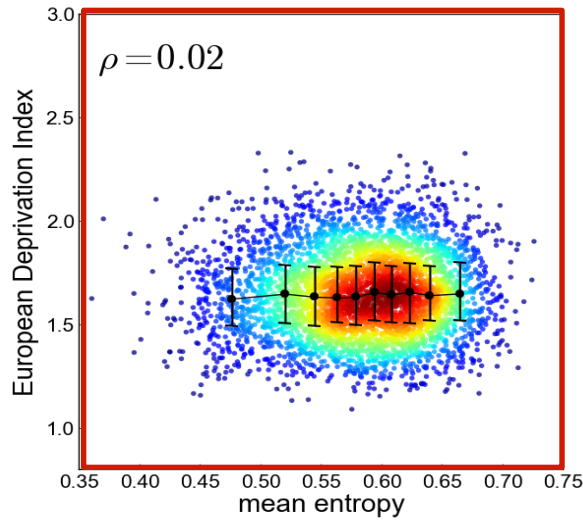
sociality



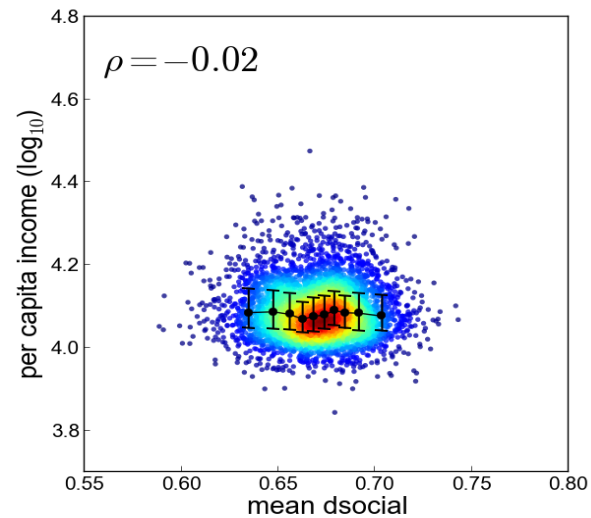
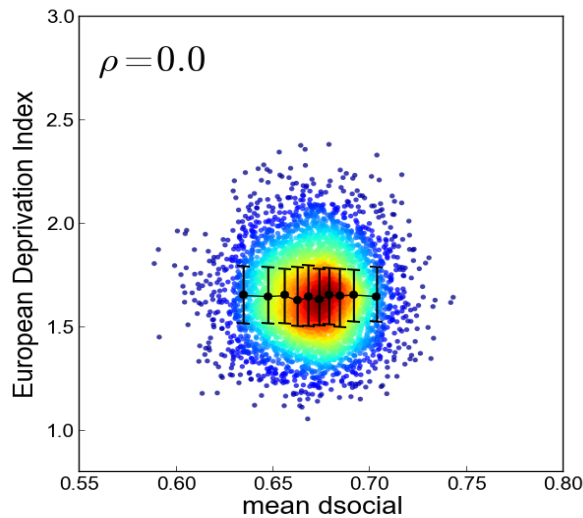
What on a null model: randomizing EDI

EDI

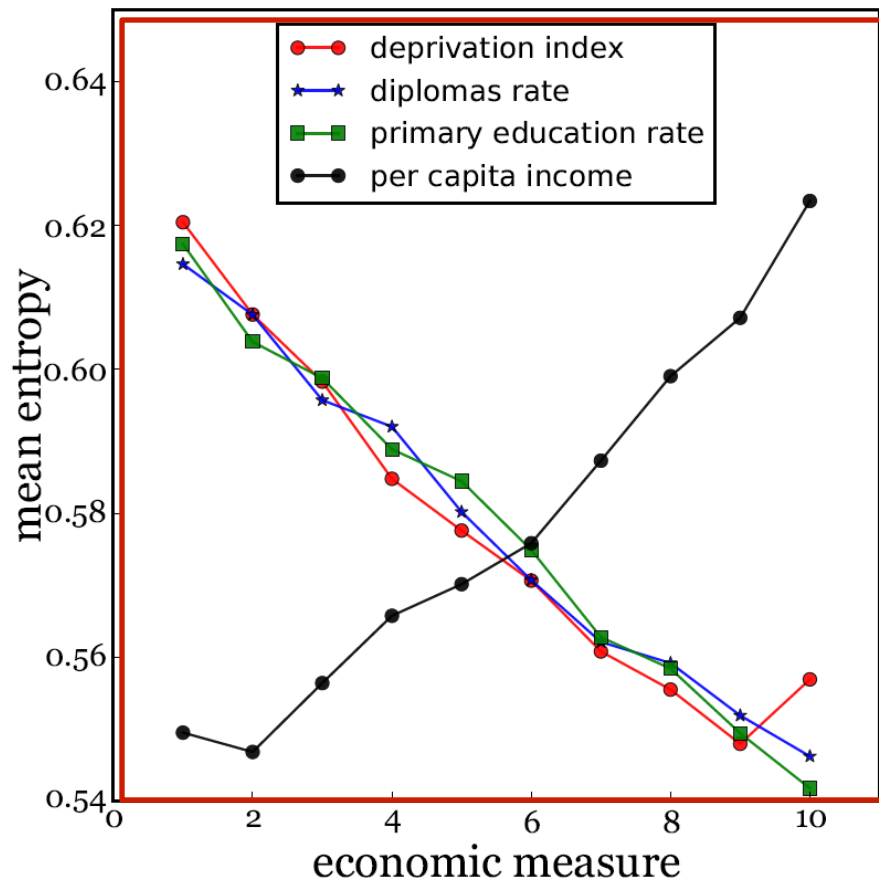
mobility



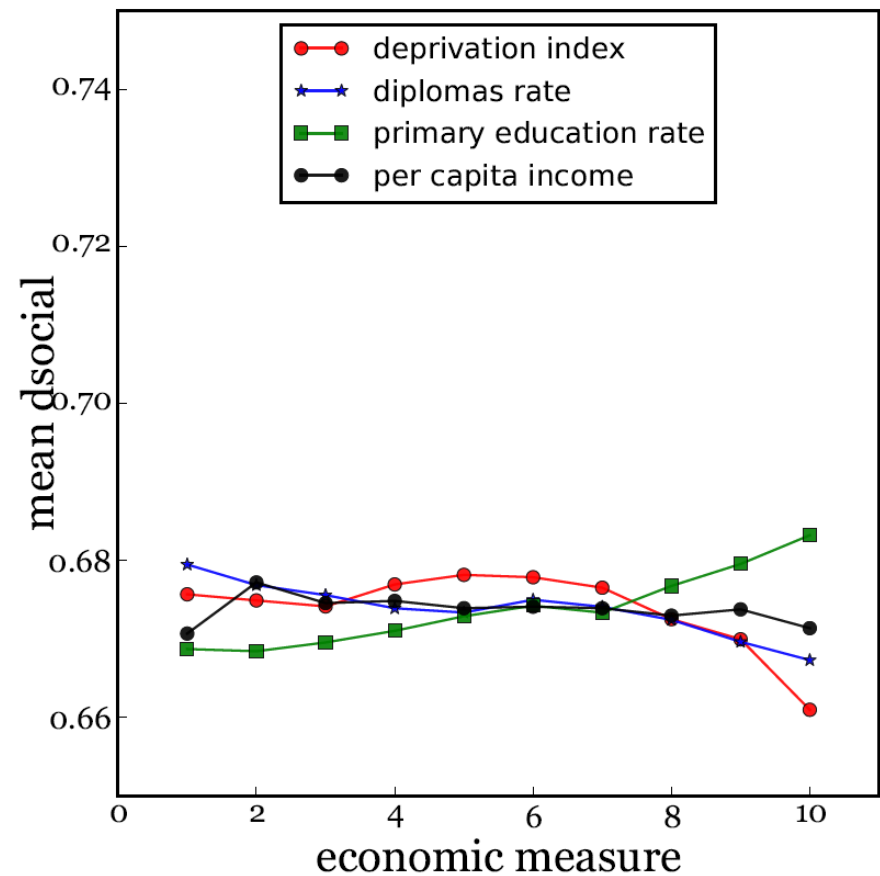
sociality



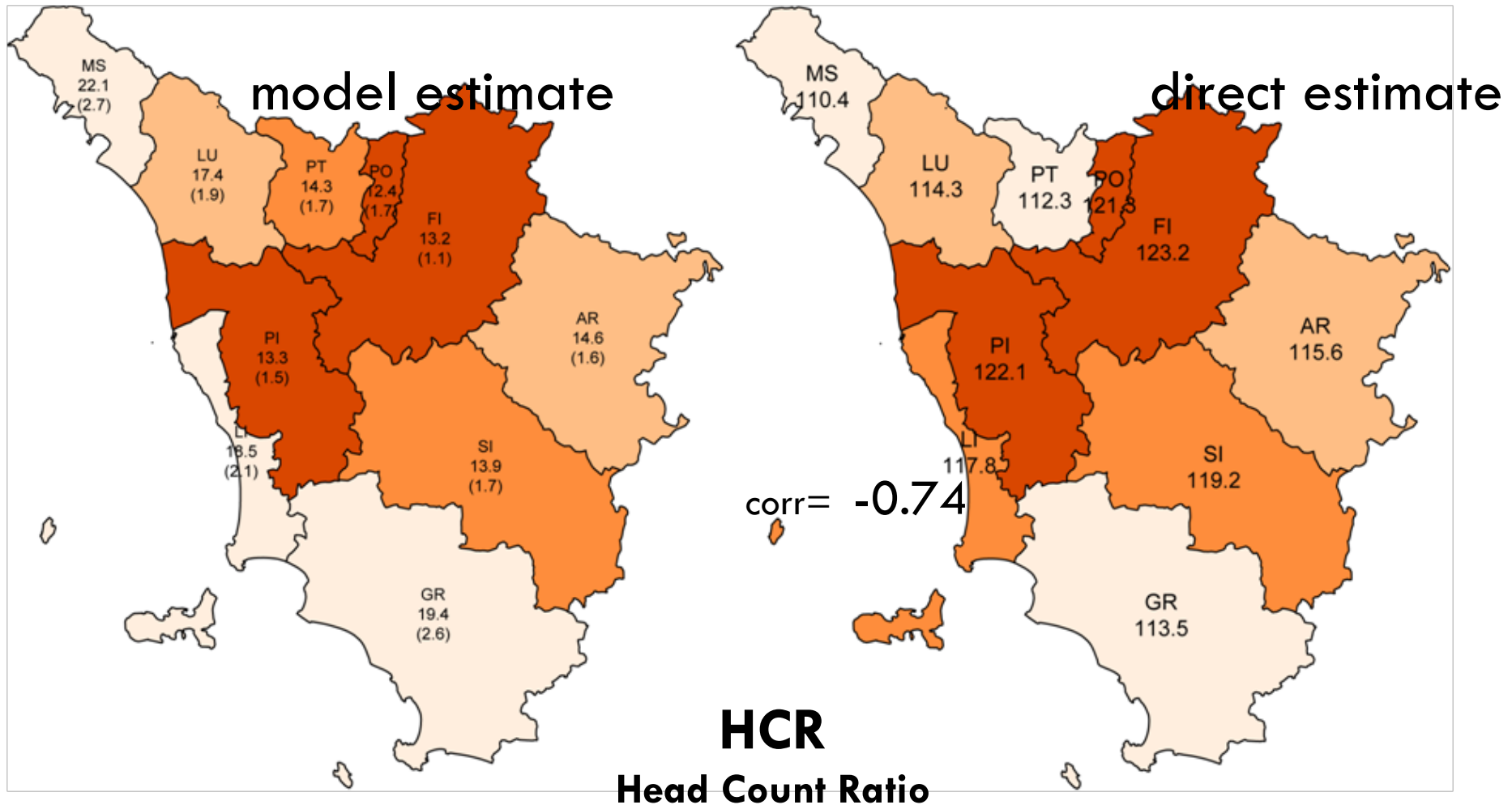
mobility



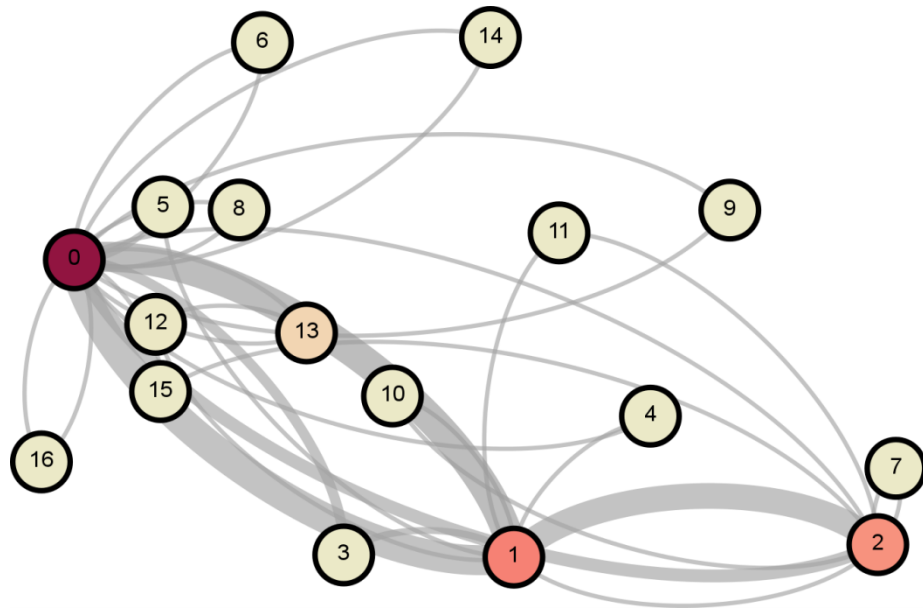
sociality



...also in Tuscany



Behind the scene: Individual Mobility Networks



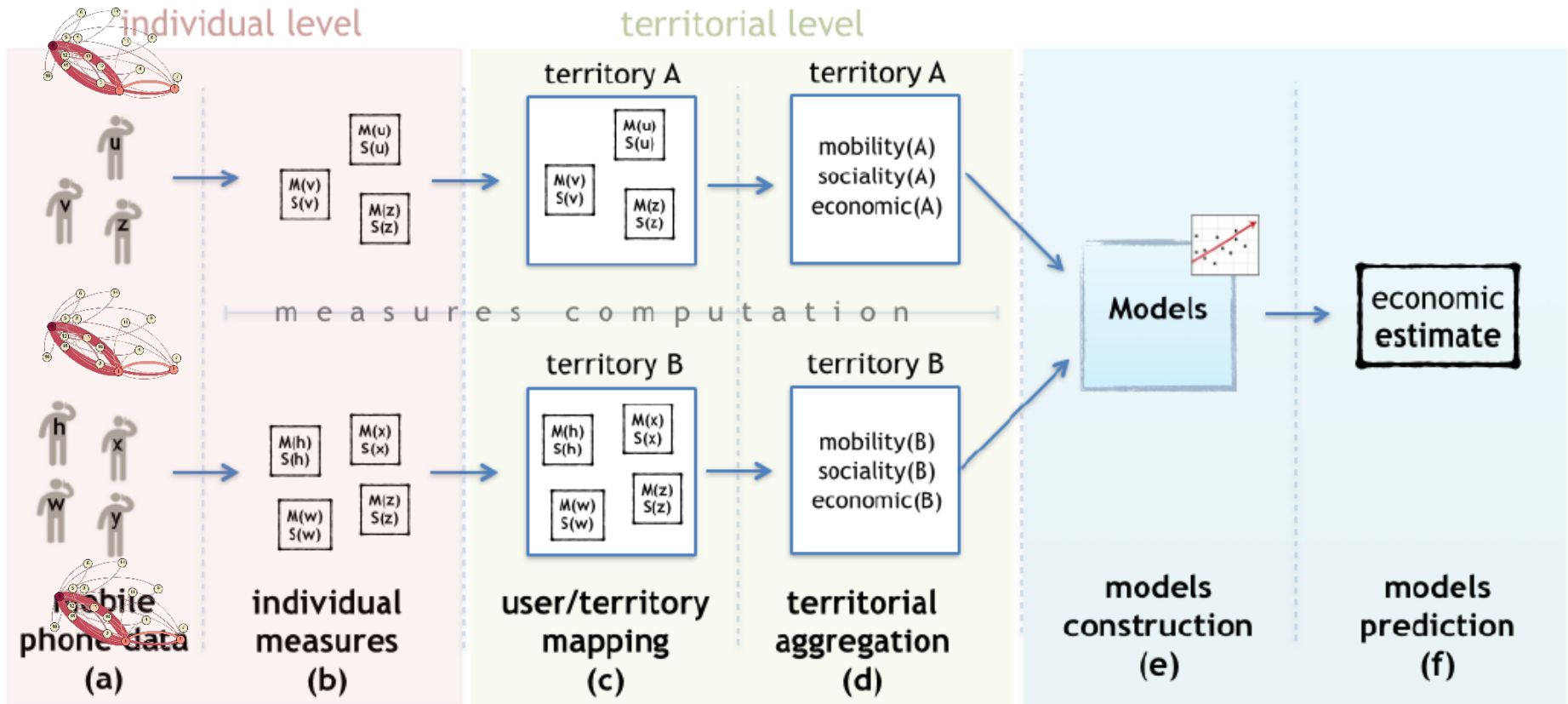
Trip Features

Length
Duration
Time Interval
Average Speed

Network Features


<i>centrality</i>	clustering coefficient average path length
<i>predictability</i>	entropy
<i>hubbiness</i>	degree betweenness
<i>volume</i>	edge weight flow per location

Behind the scene





Discussion

1. **Mobility** diversity is linked to wellbeing
 2. **Entropy** is stable across age/gender but varies with wellbeing
 3. **Geography** matters
 4. Big Data as a pillar for official statistics
- 

Soccer Player Ratings

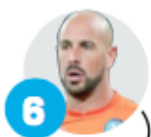


The New York Times

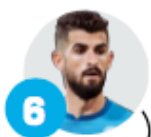
“I don’t ask Messi for more to get an 8 or 9,” he said. “The thing is, Messi tends to play better more frequently, so he usually gets a good rating. You see what you see, and you try to be honest. It’s all you can do.”

Rating performances is a
complex task,
can we reproduce it?

ZIELINSKI, EREDE DI QUALITÀ



REINA
Per i fantacalcisti e perché sull'unico pallone rischia la salute andando nella pozzanghera.



HYSAJ
Eppure il campo non gli manca (non gli mancherebbe) ma le energie forse un pochino sì.



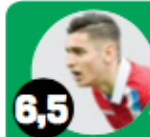
ALBIOL
Con le ciabatte, in stile salotto, lasciando che la Spal gli vada a battere addosso.



KOULIBALY
Il solito «energumeno»: di forza, di prepotenza e con autorevolezza ritrovata.



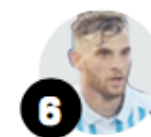
MARIO RUI
Rischia il giallo (e la squalifica) e quindi poi si contiene, limitandosi.



MERET
E' bravo, reattivo, istintivo e frena Insigne ma soprattutto Callejon.



SALOMON
Non sceglie: aspetta o attacca Insigne e rischia di finire a gambe all'aria.



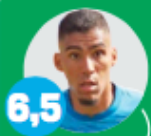
VICARI
Sta là dietro e oppone il corpo e la posizione alle rare verticalizzazioni.



FELIPE
Si stacca troppo, aprendo la corsia centrale per Allan, perché Callejon lo distrae.



LAZZARI
Gli mancano le coperture e poi dà un senso di anarchia tagliando sempre, troppo.



ALLAN
Il gol che riconsegna il primato in classifica, prima di correre per sé e per gli altri.



JORGINHO
Geometrie apprezzabili, però senza avere intorno uomini che pedalino come si dovrebbe.



HAMSIK
Il pallido capitano rimane dietro i suoi standard e l'ammonizione gli fa male.



CALLEJON
Apre per Allan e lo manda in porta e poi (sembra) governa i carichi di fatica.



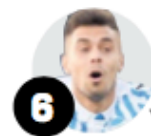
MERTENS
E' la prima sponda nell'1-0 ma è anche un po' vago, quasi distante dalla partita.



SCHIATTARELLA
Si ritrova con Hamsik, lo contiene e persino lo costringe a stargli dietro.



VIVIANI
Gli viene meno il gusto di osare e palleggia con paura addosso che diventa nemica.



GRASSI
Perde lo scatto di Allan, poi dà movimento e pure eleganza ad un centrocampio piatto.



DRAMÈ
Quasi si isola e lascia che da quelle parti, ma senza esagerare, il Napoli vada.



KURTIC
L'unica preoccupazione è Jorginho e spreca non l'occasione ma il suo tempo.



INSIGNE
Insegue il gol, e si vede, però Meret e il palo lo costringono a soffrire ancora.



ZIELINSKI
(25' st)
E' di impatto ma anche di talento (e che ruleta!). Hamsik ha un erede di qualità assoluta.



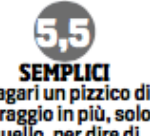
ROG
(41' st)
Va a coprire il campo, per restringerlo, nel finale da domare con intelligenza.



DIAWARA
(45' st)
L'ultimo argine per il recupero che diventa ampio e comunque pericoloso.



SARRI
Piccole tracce di Napoli, qualcosina all'avvio, poi una gestione eccessiva.



SEMPlici
Magari un pizzico di coraggio in più, solo quello, per dire di averci provato.



ANTENUCCI
Non gli arriva uno straccio di pallone, ma non ne va neanche a inseguire.



COSTA
(16' st)
In un contesto blando a cui può solo garantire di fungere da cerniera.



FLOCCARI
(30' st)
E' il jolly che si va a cercare: magari una palla sporca. Ma bisognerebbe arrivare a lui.



PALOSCHI
(37' st)
Aggiunge spiccioli di minutaggio ad una gara in cui l'attacco non esiste.



GAVILUCCI
Già non averla complicata, semplice com'era, sa di buon senso. Comoda così eh

wyscout

La Gazzetta dello Sport

Corriere dello Sport

TUTTOSPORT

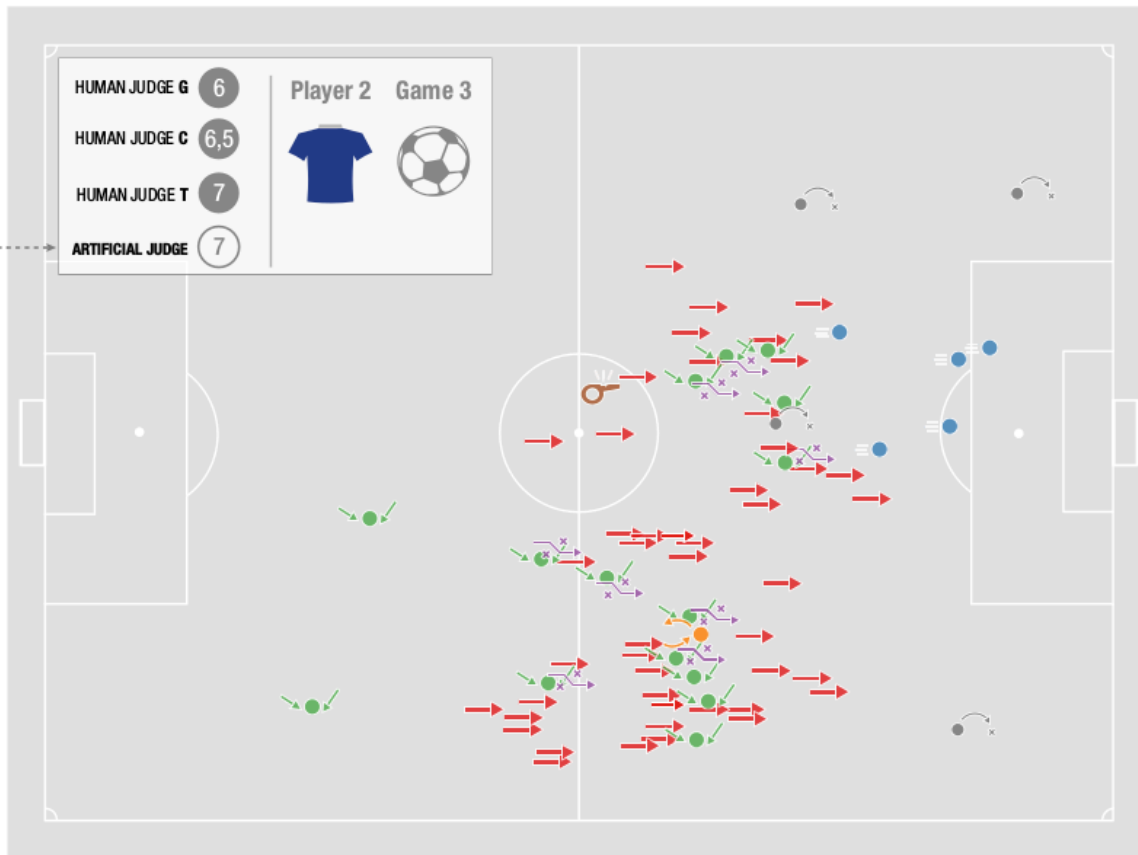
700 players

760 games

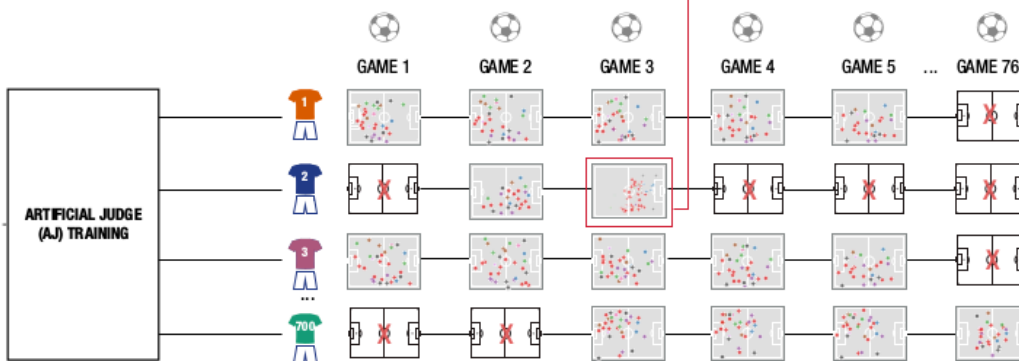
1M events



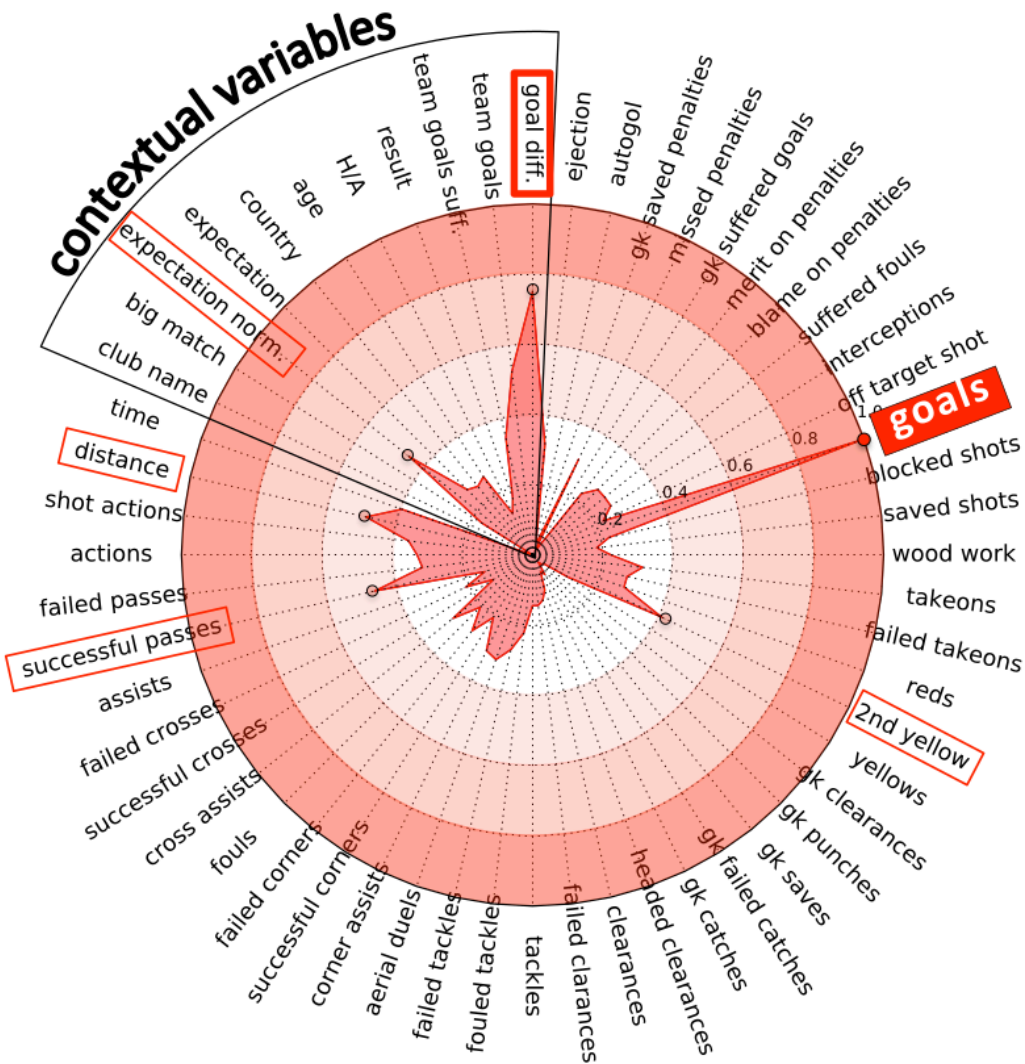
(a)



(b)



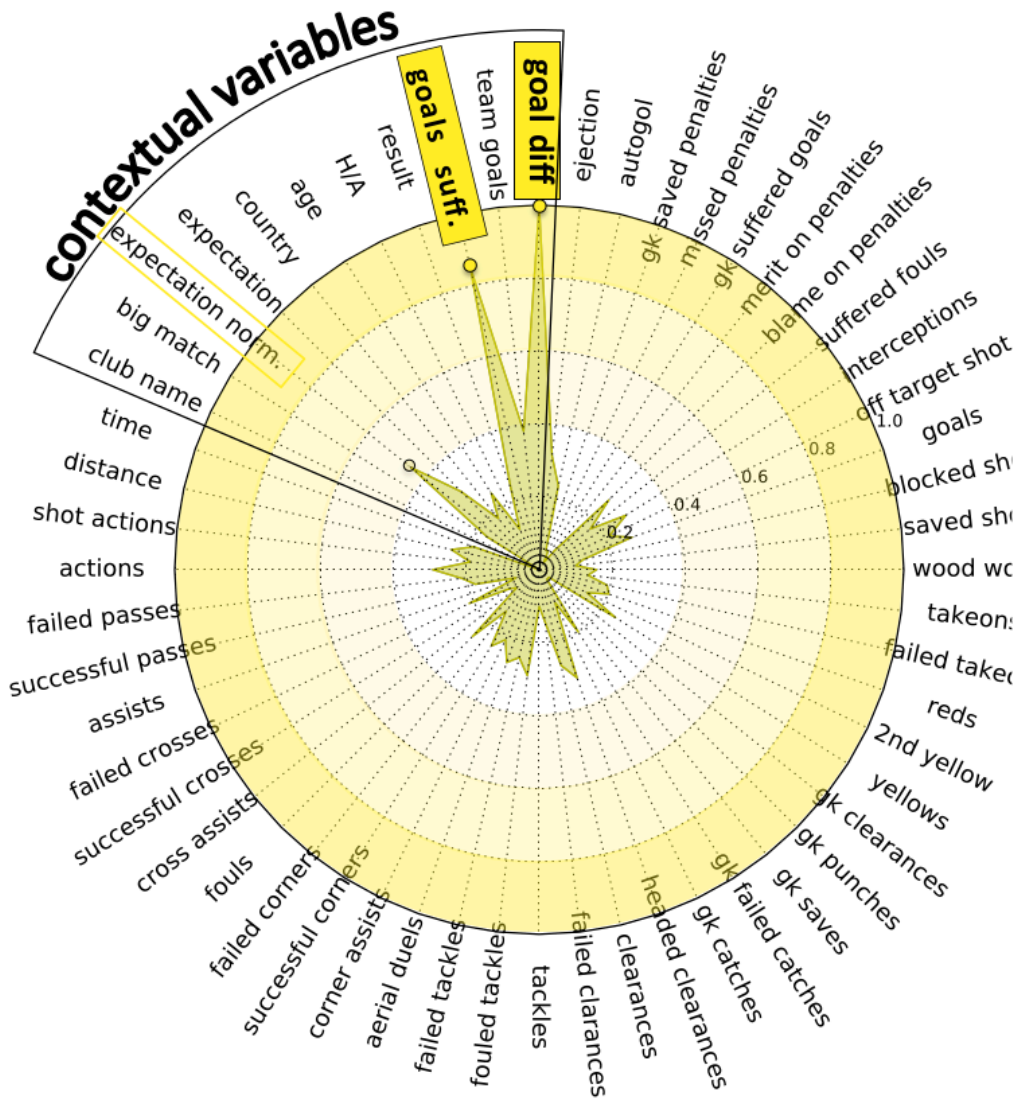
Forwards



Features that matter

- 1) Just a subset of the features matter (20)
- 2) Contextual features are highly important
- 3) >90% of the features have negligible importance

Defenders



Features that matter

- 1) Just a subset of the features matter (20)
- 2) Contextual features are highly important
- 3) >90% of the features have negligible importance
- 4) the same features has different importance in different roles

A Multidisciplinary piece of art

Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dessiné par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite Paris, le 20 Novembre 1869

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en traits de ces zones. Le rouge désigne les hommes qui ont péri en Russie, le noir ceux qui ont survécu. Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Fezardac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Tous mieux faits jugés à l'œil la diminution de l'armée; j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilow n'ont rejoint nos Ouches en Wilhelk, avaient toujours marché avec l'armée.

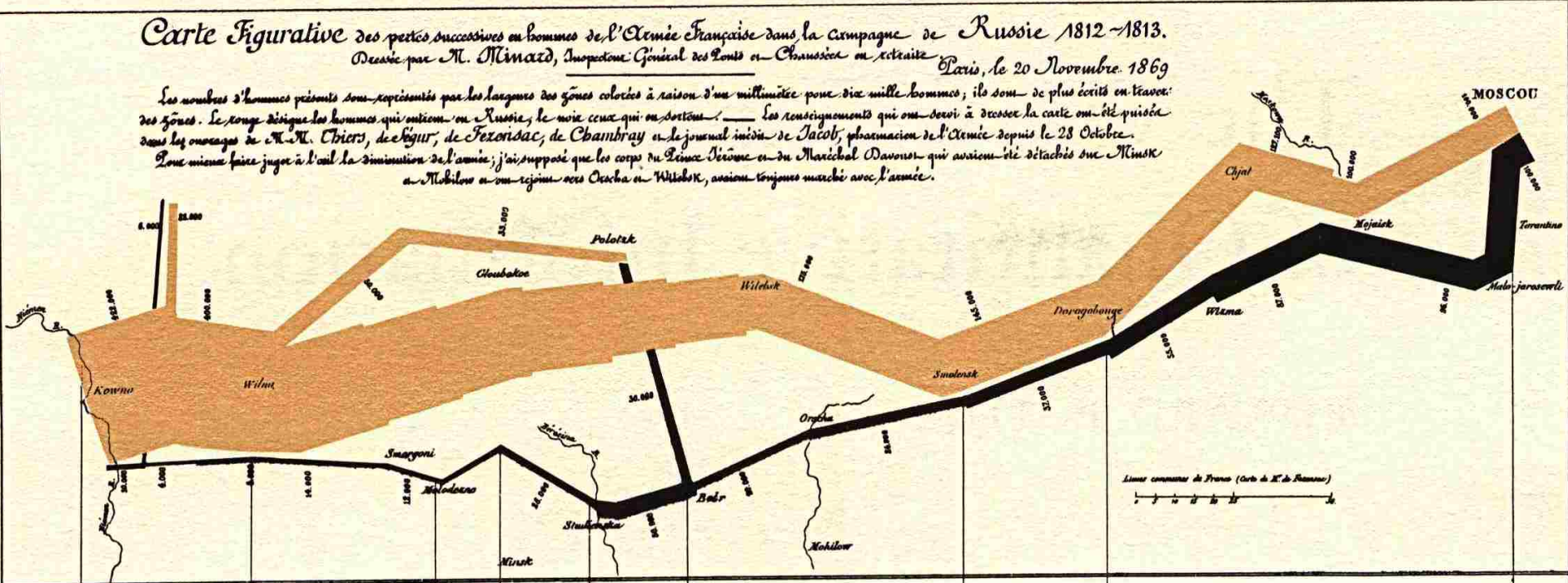
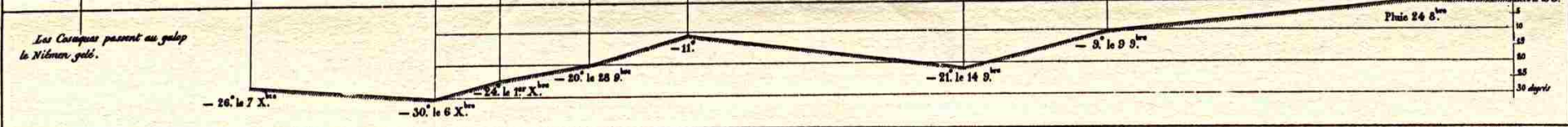


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

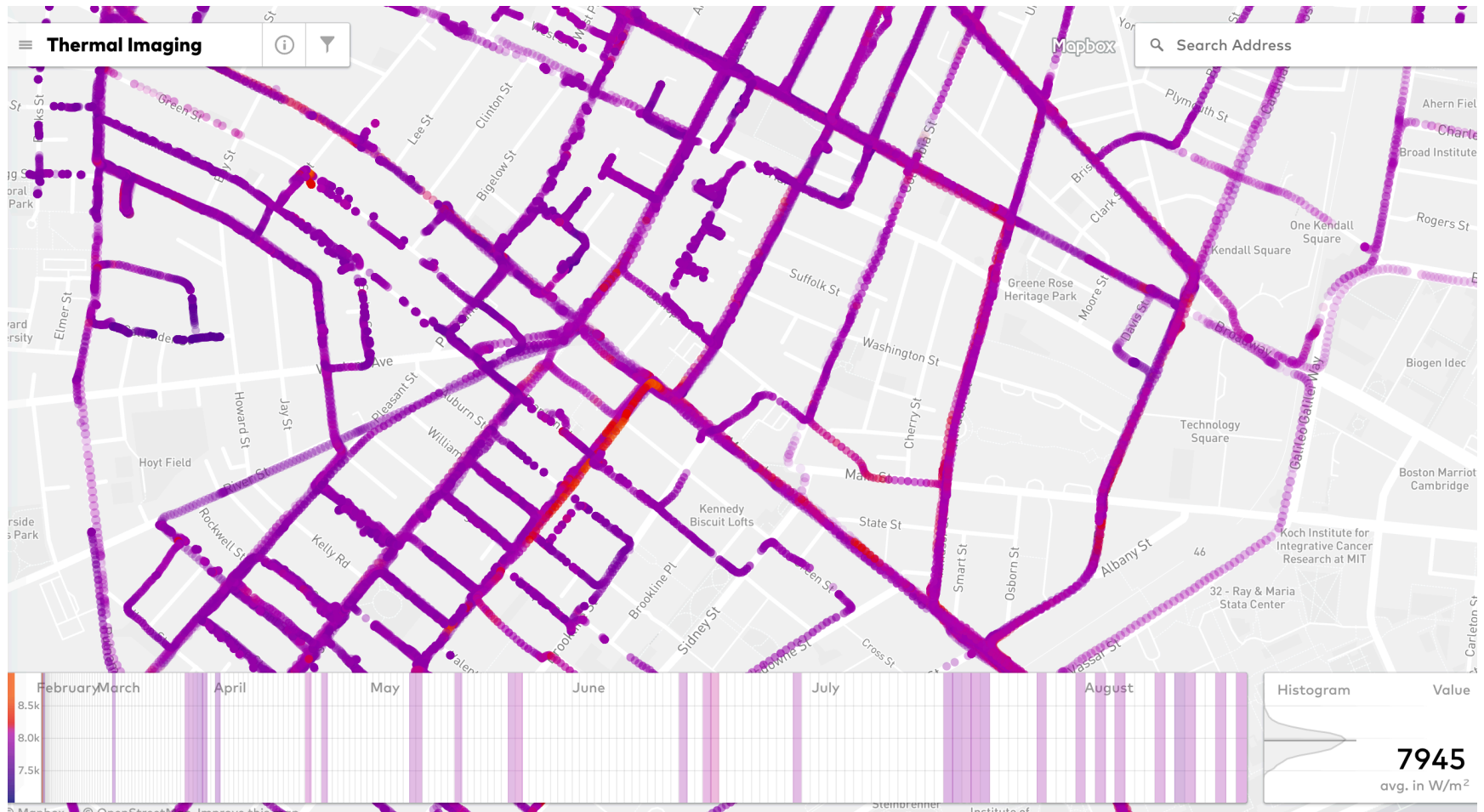


Auq. par Raynier, à Par. J^{de} Maria J^{de} O^{de} à Paris.

Imp. Lith. Raynier et Desvres.

Charles Minard. "Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813", 1869.

City Scanners – Senseable City Lab (MIT)

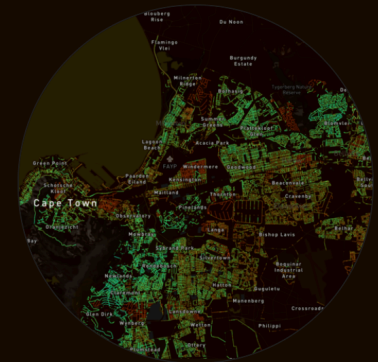
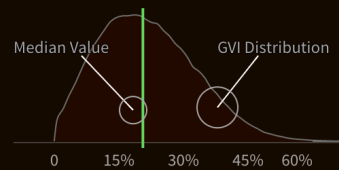


<https://youtu.be/Y9wTuLQkLzc>

<http://senseable.mit.edu/cityscanner/>

Treepedia – Senseable City Lab (MIT)

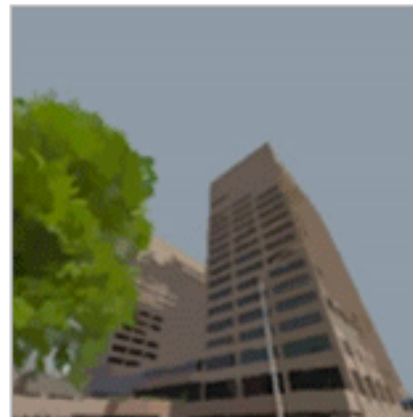
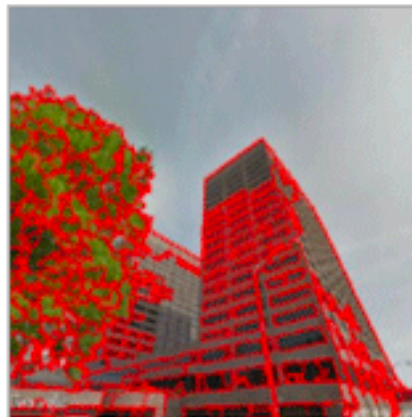
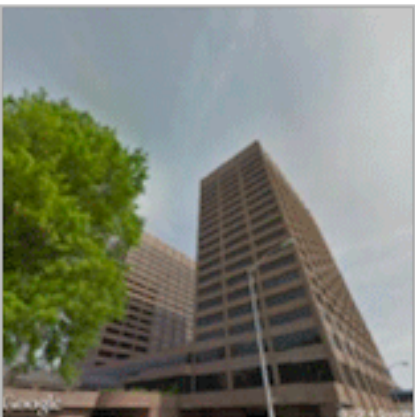
Compare
Different Cities



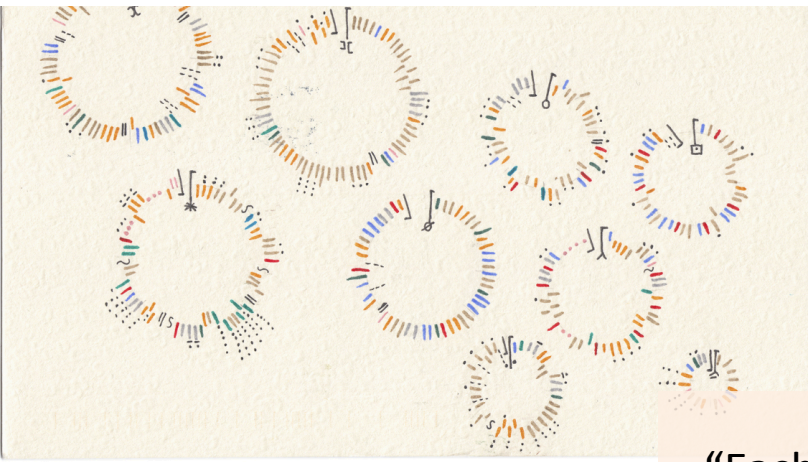
<http://senseable.mit.edu/treepedia>

Green View Index

- Use Google Street View images to estimate green canopy coverage of a city
- Focus on individual point of view (instead of satellite imagery)



Dear Data



66 DEAR DATA
WEEK 08: PHONE ADDICTION!

HOW TO READ IT: Every circle represents a PLACE or SITUATION where I checked my phone, somehow ordered from left to right according to how many times I did it in that place. Every single LINE is a SINGLE TIME I interacted with my phone, ordered chronologically per each place.

PLACES/SIT.:
 x while walking
 * while working
 II while waiting for sth or s. body
 o in the Bathroom
 o on the couch
 □ on the bed
 ^ other places at home
 % cafe/restaurants shops....

NEW YORK FROM: 100 GEORGIA LUPU 05 NOV 2 11209 BROOKLYN - NY - USA

SEND TO: STEFANIE POSAVEN LONDON - UK - ENGLAND

ATTRIBUTES:
 → OUTSIDE = I picked it PURPOSELY
 → INSIDE = Because of an alert
 = turned the phone facing the table not to see it
 --- didn't pick it because I didn't want to report

“Each week, and for a year, we collected and measured a particular type of data about our lives, used this data to make a drawing on a postcard-sized sheet of paper, and then dropped the postcard in an English “postbox” (Stefanie) or an American “mailbox” (Giorgia)!”

<http://www.dear-data.com>

FROM: LONDON 97

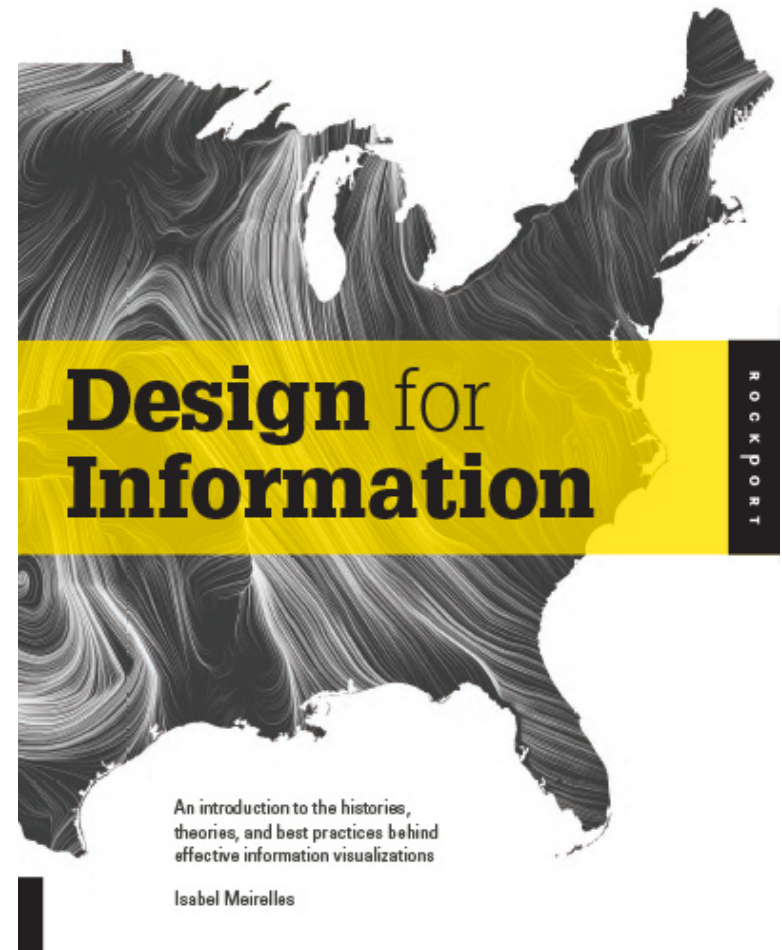
REASONS FOR PICKING UP PHONE:
 CHECK THE TIME
 PLAY MUSIC
 I WANT TO ALARM
 CHECK/SEND EMAIL
 FEELS OF HABIT, NO REAL REASON
 CHECK SOCIAL MEDIA

NOTES:
 I SPENT LOADS OF TIME + WASTED MULTIPLE CARDS
 AND I STILL AM NOT TOTALLY HAPPY W/ MY DRAWING! OH WELL, YOU WIN!

BROOKLYN, NY 11209
USA

BY AIR MAIL par avion Royal Mail®

Suggested Readings





SoBigData

Research Infrastructure



Social Mining &
Big Data Ecosystem
H2020 - www.sobigdata.eu
September 2015- August 2019

GARR Conference – 16th November 2017

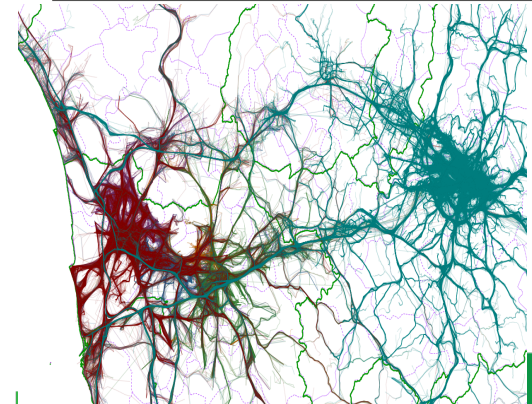
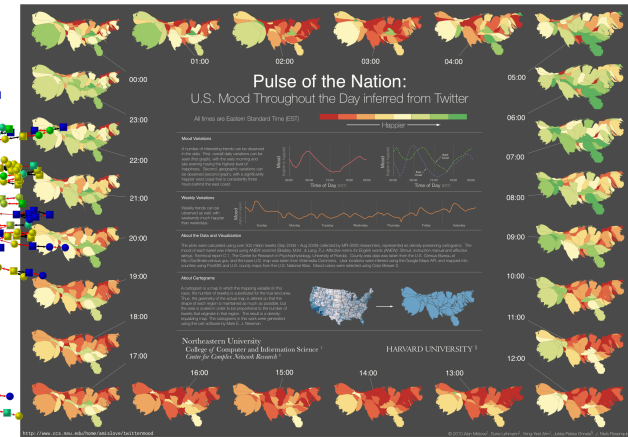
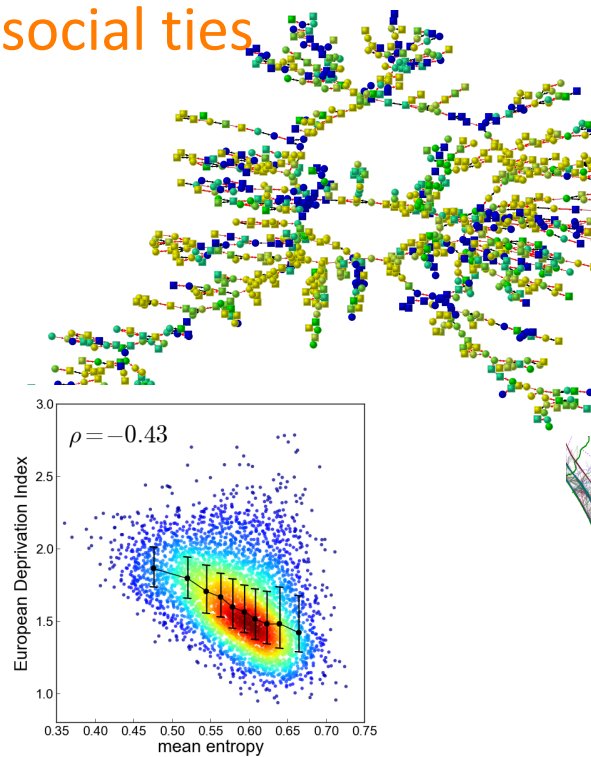
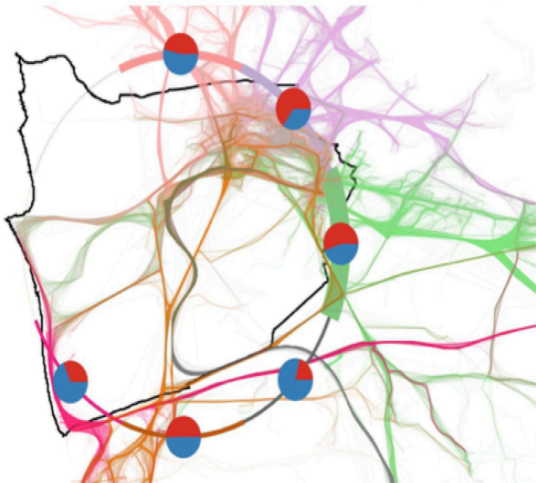
@SoBigData (<https://twitter.com/SoBigData>)

<https://www.facebook.com/SoBigData>

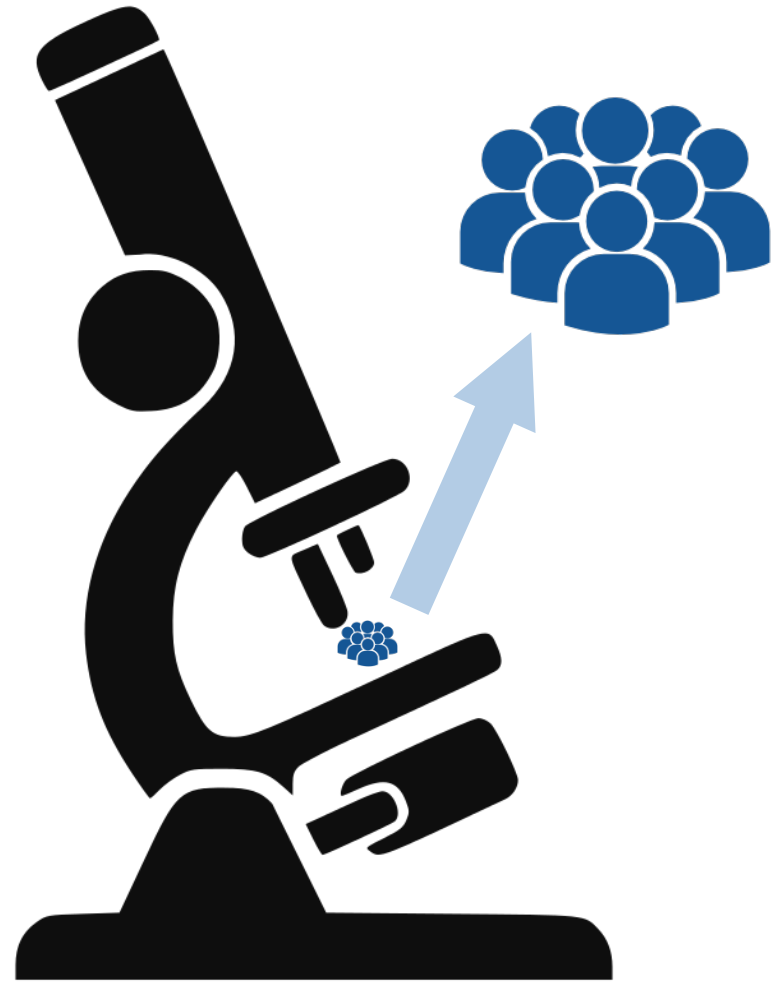


What is Social Mining

- Automated discovering patterns and models of human behaviour across the various social dimensions that have big data “proxies”
 - desires and opinions
 - relationships and social ties
 - life-styles
 - mobility



**Social mining:
making sense
of big data to
understand
society**

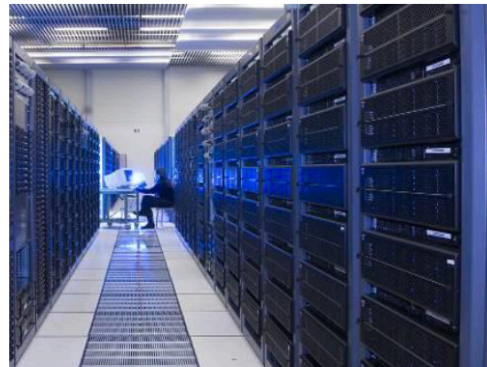
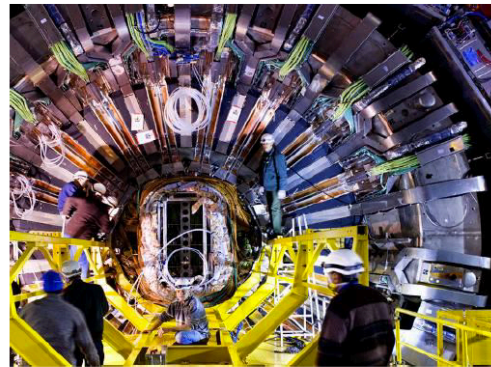


Research Infrastructures

Research infrastructures are *facilities, resources and services* used by the research communities to conduct research and foster innovation.

Major scientific equipments

Knowledge-based resources



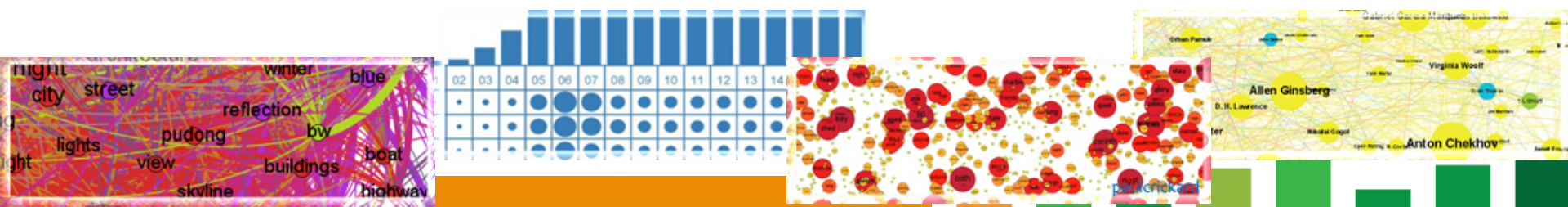
e-infrastructures



SoBigData GOAL is...



TO CONSTRUCT THE **Multidisciplinary European Infrastructure on Big Data and Social Data Mining (the Social Mining CERN)** providing an integrated ecosystem for **ethic-sensitive scientific discoveries** and advanced applications of social data mining on the various dimensions of social life, as recorded by “big data”.





The pillars for reaching the goal

- **an ever-growing, distributed data ecosystem for procurement, access and curation of big social data, within an ethic-sensitive context, based on**
 - innovative strategies for acquiring social big data for research purposes,
 - using both opportunistic means offered by social sensing technologies and
 - participatory means based on user involvement as prosumers of social data and knowledge.



The pillars for reaching the goal

- **an ever-growing, distributed platform of interoperable, social data mining methods and associated skills:**
 - tools, methodologies and services for mining, analysing, and visualising complex and massive datasets,
 - harnessing the techno-legal barriers to the ethically safe deployment of big data for social mining.



The pillars for reaching the goal

- **Building the Social Mining community of scientific, industrial, and other stakeholders (e.g. policy makers),**



The path to achieve the goals

- **Integrate European national infrastructures and centres of excellence in big data analytics, social mining and data science**
 - 1. Text and Social Media Mining (TSMM)***
 - 2. Social Network Analysis (SNA)***
 - 3. Human Mobility Analytics (HMA)***
 - 4. Web Analytics (WA)***
 - 5. Visual Analytics (VA)***
 - 6. Social Data (SD)***

Integrating national research Infrastructures



GATE⁰¹¹
general architecture
for text engineering

LIVINGARCHIVE
Lectives

SoBigData
Euro Lab on Big Data Analytics
& Social Mining

 **TARTU ÜLIKOOL**
1632

 **L3S Research Center**

 **Fraunhofer**



The Consortium





The path to achieve the goals

- **Grant access (both virtual and trans-national on-site)** to the SoBigData RI to multi-disciplinary scientists, innovators, public bodies, citizen organizations, SMEs, as well as data science students at any level of education.
- **joint research, and extensive networking and innovation actions**

Big Data Ecosystem

- Open Data
- Restricted Data
- Virtual Collections

Social Mining

- Text & Social Media Mining
- Social Network Analysis
- Human Mobility Analytics
- Web Analytics
- Visual Analytics
- Social Data

Ethical and Legal Framework



Virtual Access

E-infrastructure



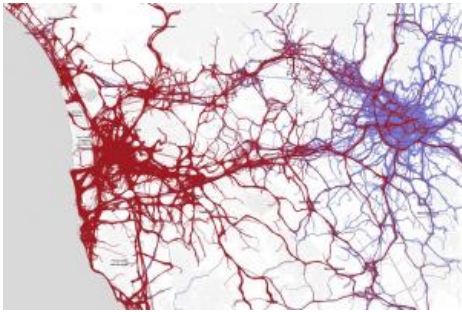
Transnational Access

Open calls
Exploratory projects



Networking

Training
Dissemination
Innovation Accelerator

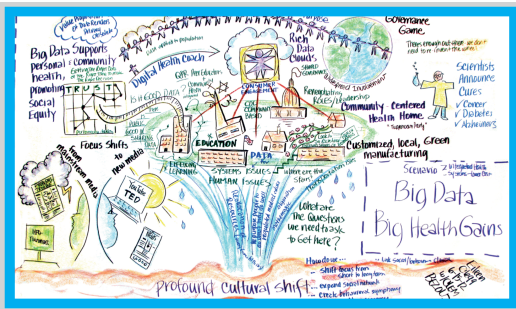
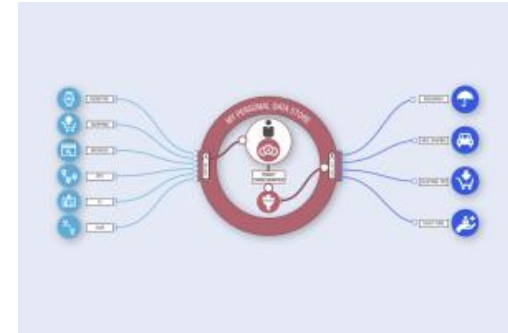


City of Citizens

This exploratory tells stories about cities and people living in it. We describe those territories by means of data, statistics and models.

Well-being & Economic Performance

Can Big Data help us to understand relationships between economy and daily life habits? We use data of purchases in supermarkets and investigate people's behavior.



Societal Debates

We study public debates on social media and newspaper. We can identify themes, following the discussions around them and tracking them through time and space.

Migration Studies

Could Big Data help to understand the migration phenomenon? We try to answer to some questions about migrations in Europe and in the world.



An aerial, high-angle photograph of a large, diverse crowd of people scattered across a vast, green, textured field. The people are seen from above, appearing as small, colorful figures. They are engaged in various activities, some standing in small groups, others walking or sitting. The overall scene conveys a sense of a large public gathering or event. A white rectangular box with orange text is overlaid on the center of the image.

Data Scientists have an obligation to take into account the ethical and legal aspects and the social impact of Data Science

Ethics and Security



Legal and Ethical framework

Define and implement the legal and ethical framework of the SoBigData RI, in accordance with the European and national legislations

Monitor of research

Monitor the compliance of experiments and research protocols with the framework

Privacy-by-design

The development of big data analytics and social mining tools with Value-Sensitive Design and privacy-by-design methodologies

The GDPR

- It entered into force on 25 May 2018
- Introduces important novelties
 - New Obligations
 - New Rights



EUROPEAN DATA PROTECTION SUPERVISOR

Opinion 7/2015

Meeting the challenges of big data

*A call for transparency, user control, data
protection by design and accountability*

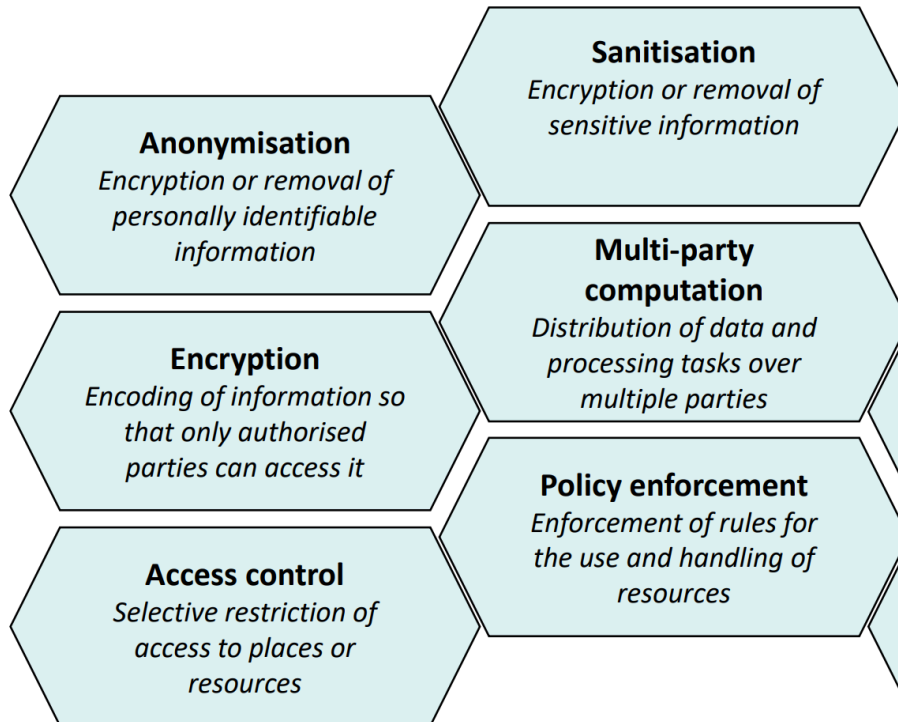


New Elements in the EU GDPR

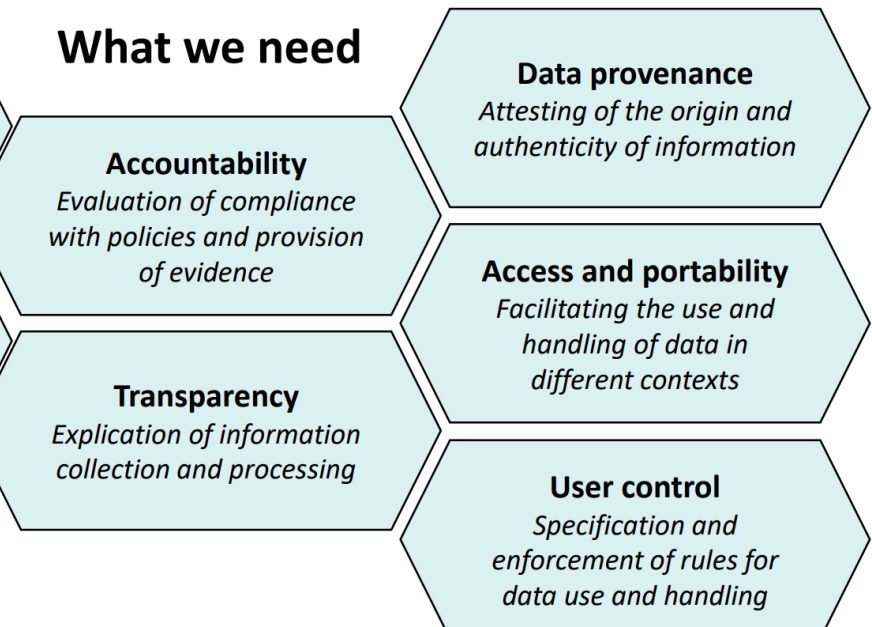
- New Obligations for Data Processors
- GDPR Outside EU
- Accountability Principle
- **Privacy by Design**
- Principle of Transparency
- Data Portability
- Right of Oblivion
- Profiling
- **The right of explanation**
- Research Data & GDPR

PET technology

What is mainly done



What we need



What is coming up



PRIVACY BY DESIGN

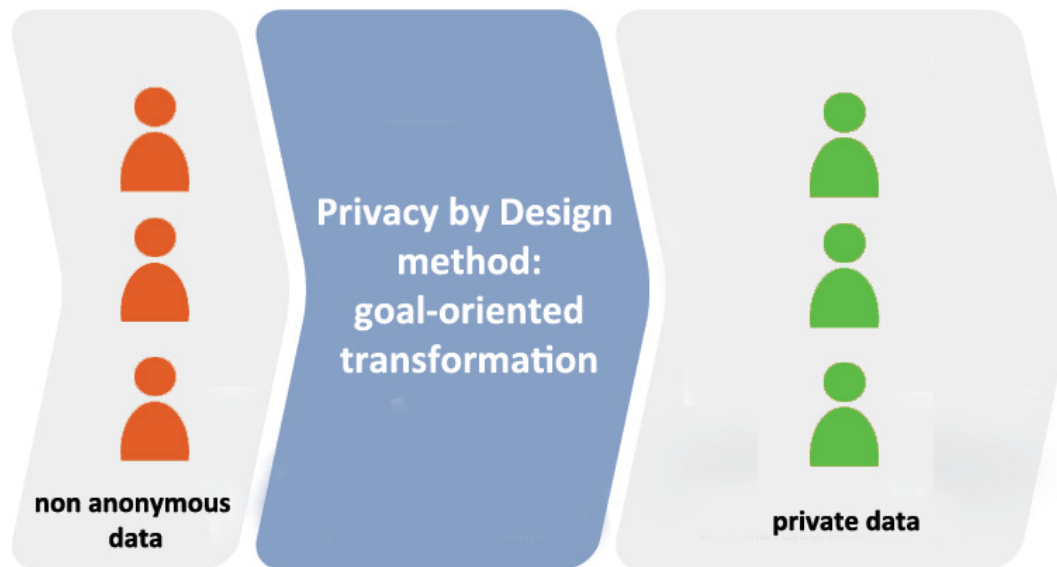


Privacy-by-Design

1. **Proactive** not reactive; preventative not remedial
2. Privacy as the **default** setting
3. Privacy **embedded** into design
4. **Full functionality** – positive-sum, not zero-sum
5. End-to-end security – **full lifecycle protection**
6. Visibility and **transparency** – keep it open
7. Respect for user privacy – keep it **user-centric**

Privacy by design big data analytics

- Design analytical process that implement the **privacy-by-design & by-default** principle



- Consider privacy at every stage of the service implementation
- Integrate privacy requirements “by design” into business models.



PRUDEnce: a System for Assessing Privacy Risk vs Utility in Data Sharing Ecosystems



Knowledge Discovery and Delivery Lab
(ISTI-CNR & Univ. Pisa)

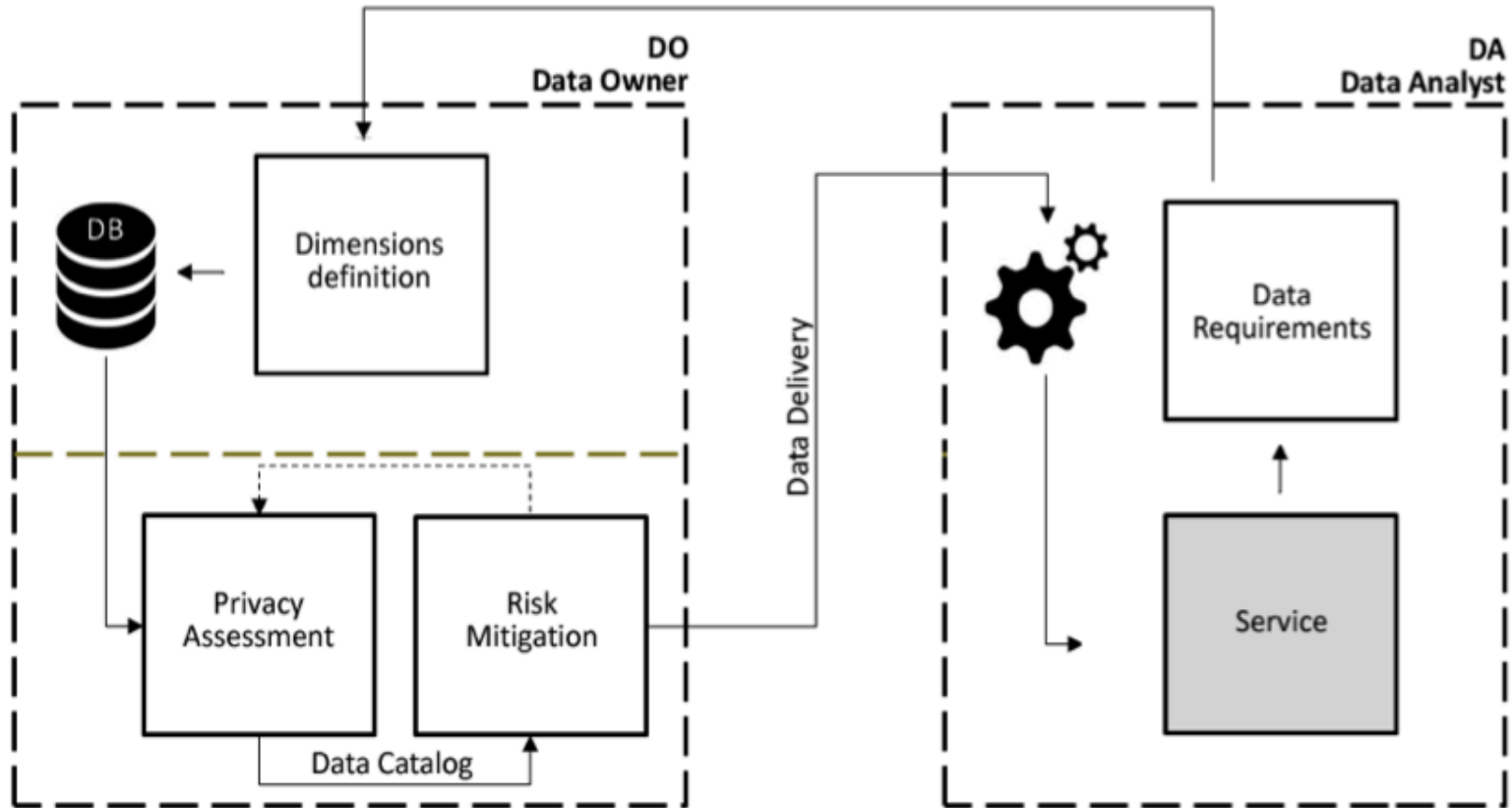
www-kdd.isti.cnr.it



Privacy by Design Methodology in PRUDENCE

- *PRUDENCE* is designed with assumptions about
 - The **sensitive data** that are the subject of the analysis
 - The **attack model**, i.e., the knowledge and purpose of a malicious party that wants to discover the sensitive data
 - The **target analytical questions** that are to be answered with the data
- *PRUDENCE is capable* to
 - transform the data into an anonymous version with a **quantifiable privacy guarantee**
 - guarantee that the analytical questions can be answered correctly, within a **quantifiable** approximation that specifies the **data utility**

Privacy Risk Assessment Framework



PRIVACY-AWARE FRAMEWORK FOR DATA SHARING

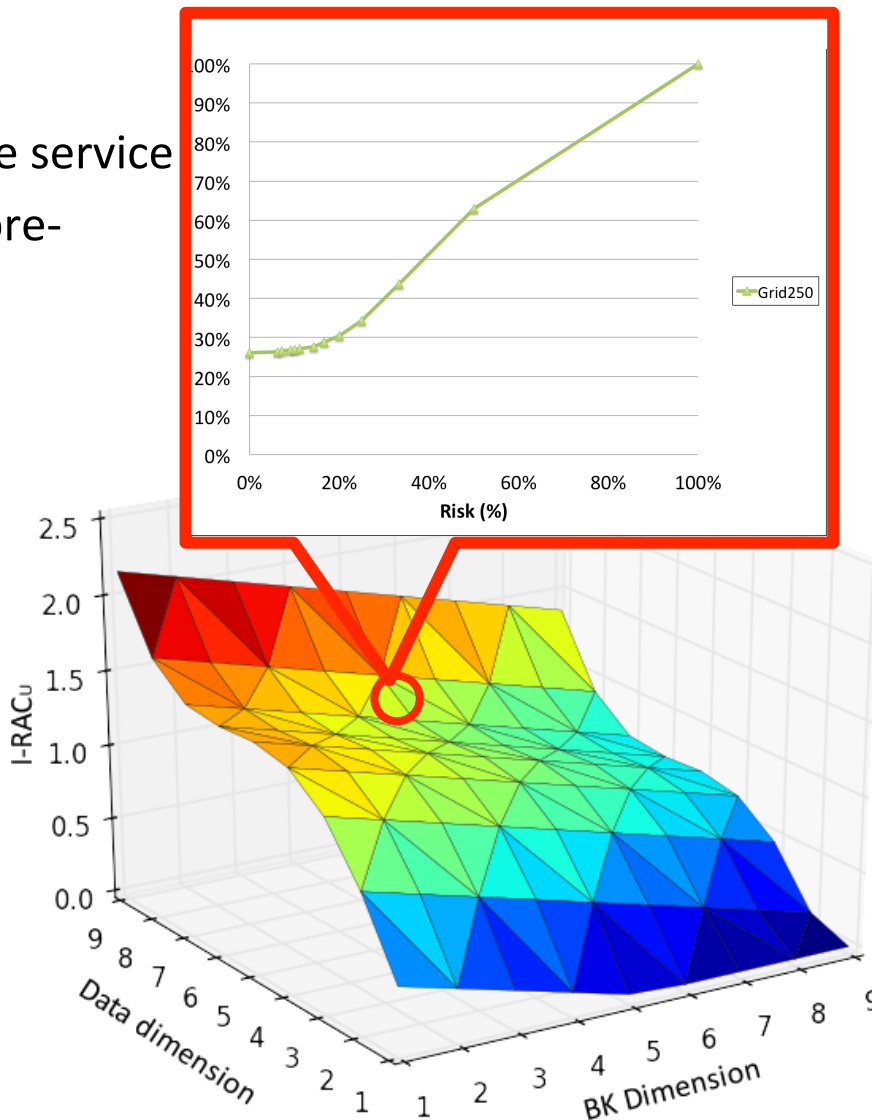
Data Catalog

For each:

- **Data Format**, i.e., the data needed for the service
- **Risk Assessment Setting**, i.e., the set of pre-processing and privacy attacks

The Data Catalog provides:

- **Quantification of Privacy Risk**, i.e., the evaluation of the real risk of re-identification
- **Quantification of Data Quality**, i.e., the quality level we can achieve with private data, compared with the data quality of original data.



Simulation of privacy harmful Inferences

Data dimension:

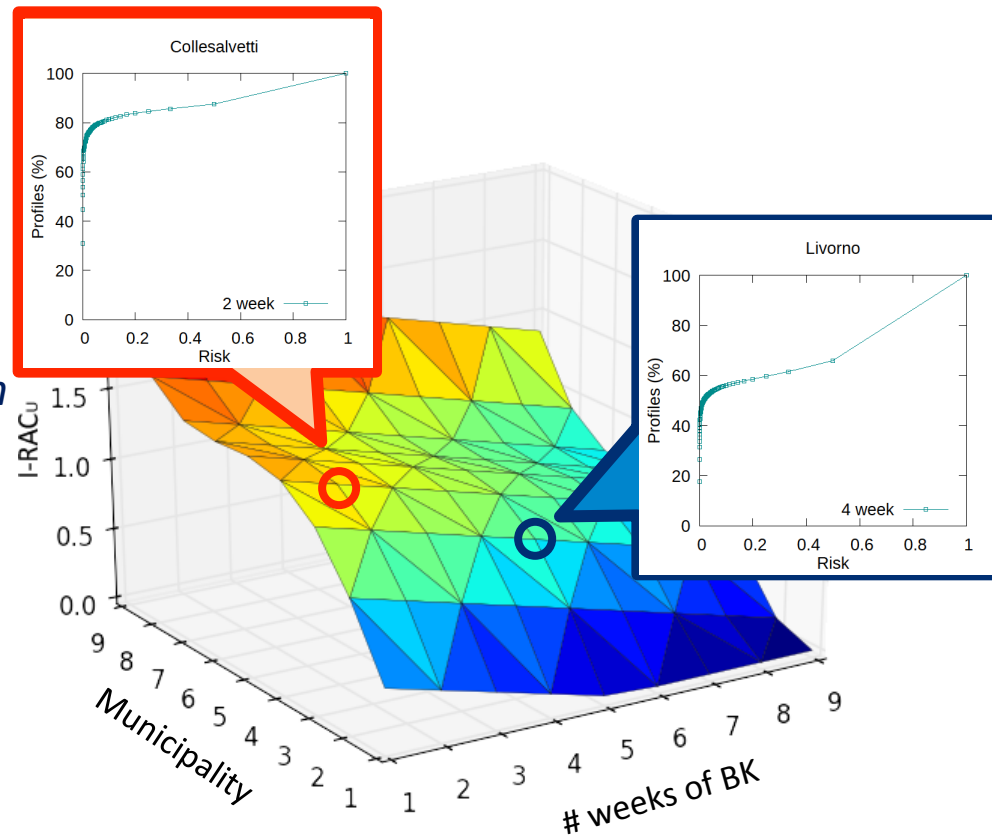
The spatial area in which the analysis is performed.

Background Knowledge dimension:

The temporal window (in weeks) in which the attacker recorded the user activity.

I-RACu:

An indicator of the risk of re-identification of the users





How is it possible to define services GDPR compliant based on GPS data?

SOME PRACTICAL EXAMPLES (1)



Services that need for GPS data

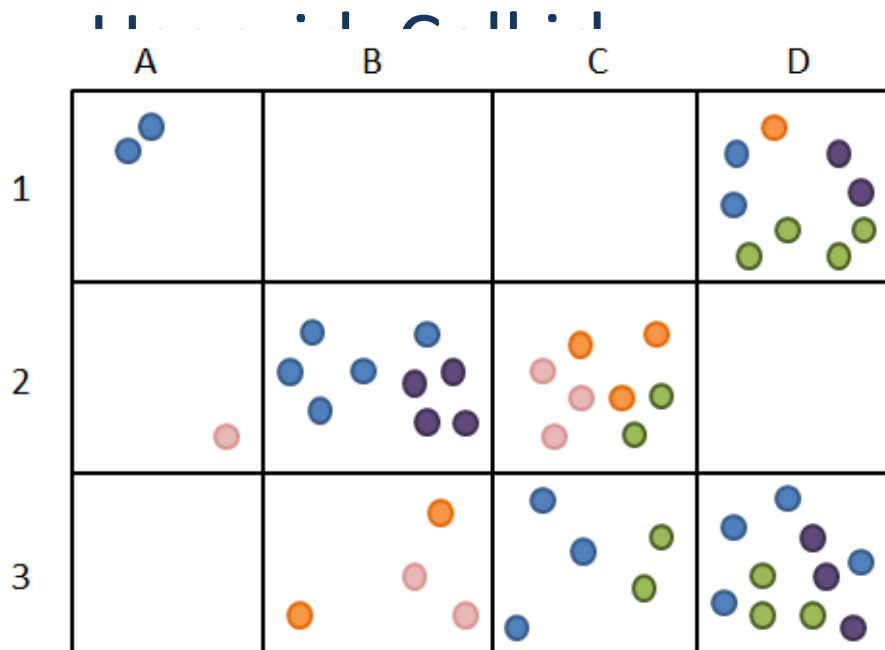
- Parking Assistance
- Geolocalized Marketing Advices
- Traffic jam analysis and prevention
- Navigation systems development
- Route/destination prediction
- Selection of the best location where to open a new facility
 - franchise store
 - fuel station
 - shopping mall
- [...]

Example1: Individual Presences

- Possible services:
 - Developing Parking Assistance
 - Geolocalized Marketing Advices
- These services do not need for all individual trajectories: specific movements are not necessary
- The only information needed is the last position of an individual (and maybe the time)

Data description

For each user, list of locations (grid cells) that the user has frequently visited ($\#visit > \text{threshold}$)



Blue: $\langle B2,5 \rangle, \langle D3,4 \rangle, \langle C3,3 \rangle, \langle A1,2 \rangle, \langle D1,2 \rangle$

Green: $\langle D1,4 \rangle, \langle D3,3 \rangle, \langle C2,2 \rangle, \langle C3,2 \rangle$

Orange: $\langle C2,3 \rangle, \langle B3,2 \rangle$

Purple: $\langle B2,4 \rangle, \langle D3,3 \rangle, \langle D1,2 \rangle$

Pink: $\langle C2,3 \rangle, \langle B3,2 \rangle$

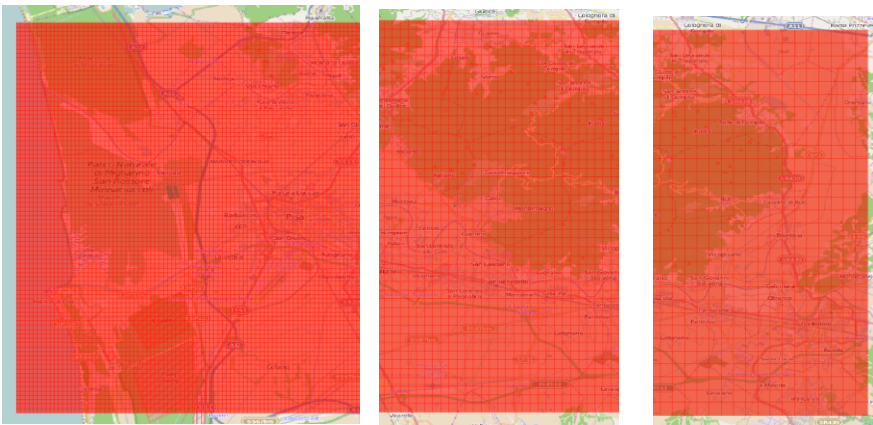
Data Dimensions

Grid size: defines the granularity of the spatial information released about each user

Frequency threshold: defines a filter on the data DO can distribute

Spatial granularity used:

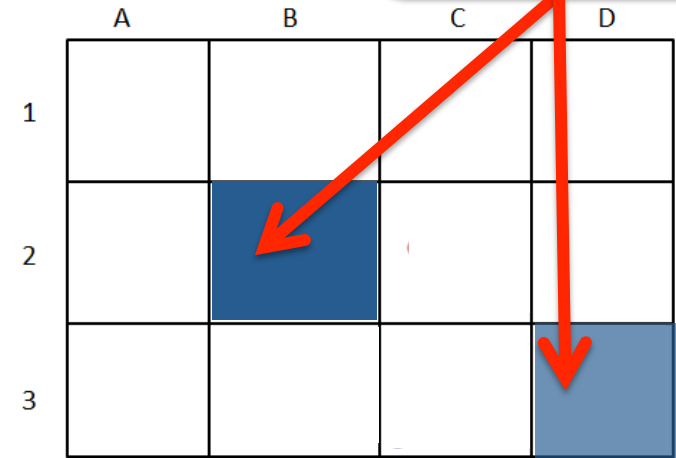
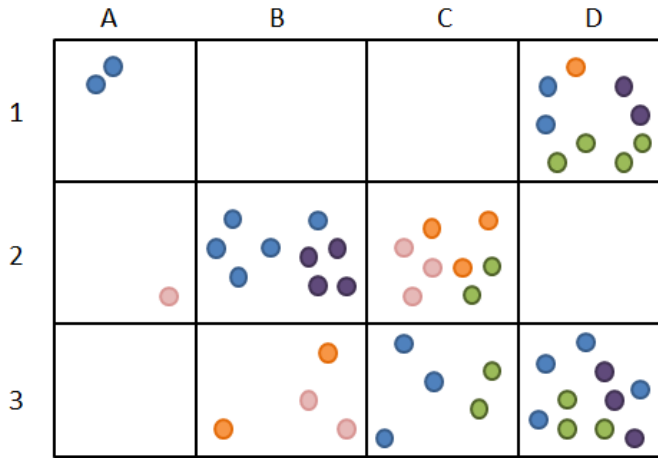
Grids (cell side): 250, 500 and 750 meters



Frequency threshold: 1, 4, 7, 10, 13

Attack 1: Top-k places

Background Knowledge:
Top-k places



The attacker knows the first k location(s) of his target

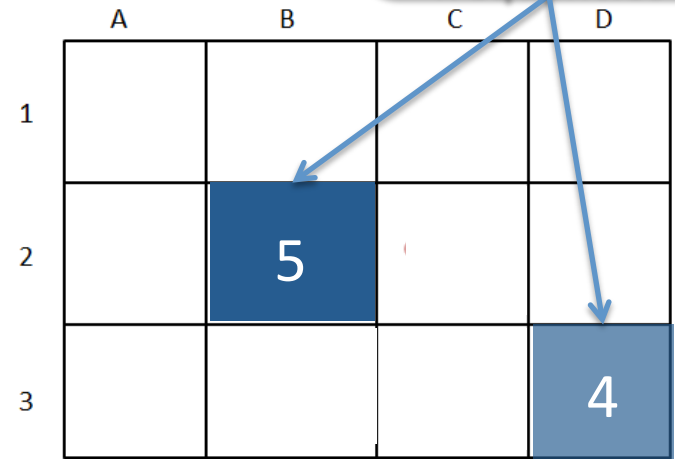
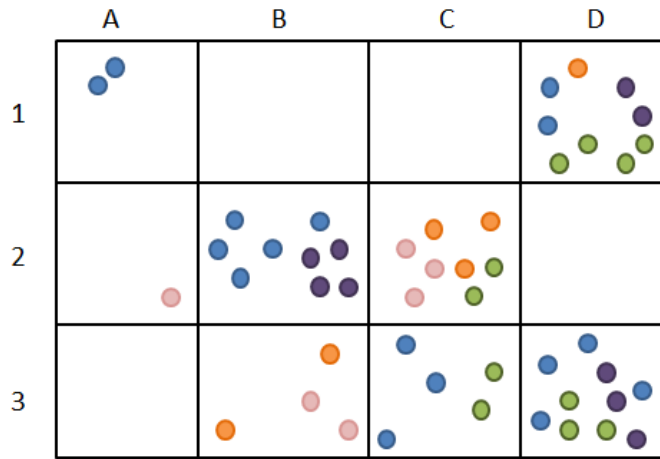
Background Knowledge Dimension:

- Number of locations known ($h = 1, 2, 3$)

E.g., Mr. Smith lives in B2 and works in D3

Attack 2: Add frequencie

Background Knowledge:
Top-k places
and their
frequencies



The attacker knows the first k location(s) of his target, and also the exact frequency

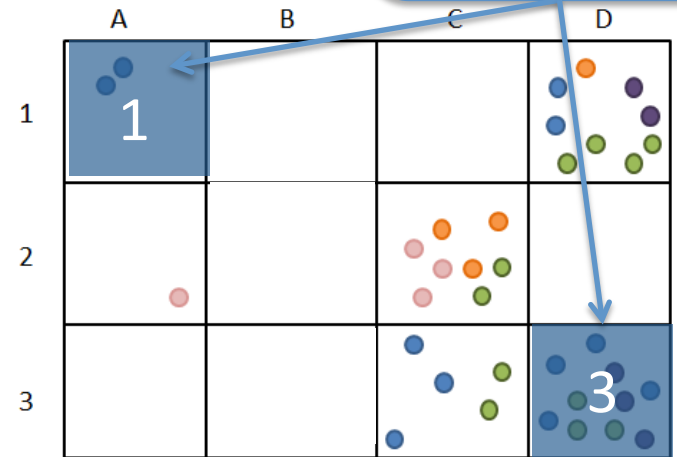
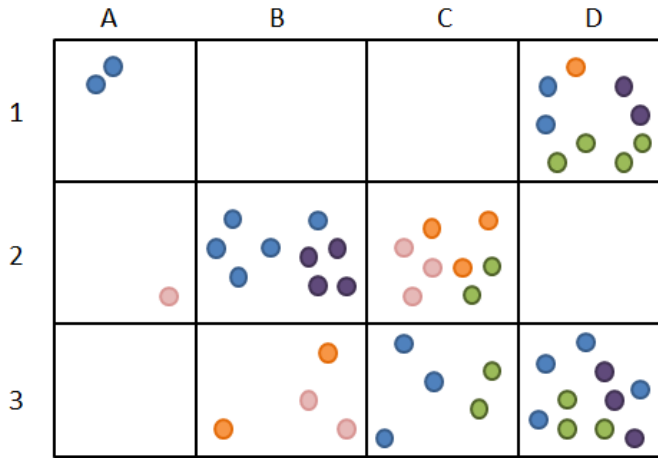
Background Knowledge Dimension:

- Number of locations known ($h = 1, 2, 3$)

E.g., Mr. Smith lives in B2 (and he parked there 5 times) and works in D3 (and he went to work 4 times)

Attack 3: Casual observati

Background Knowledge: some places and lower bounds to their frequencies



The attacker knows some location(s) with minimum frequencies

Background Knowledge Dimensions:

- Number of locations known ($h = 1, 2, 3$)
- Minimum frequency associate to the known locations (100% of original freq, 50% of original freq, only presence)

E.g., Mr. Smith was seen once in A1 and 3 times in D3

Simulation of Attack

- We simulate the chosen attack (or all of them)
- At the end we obtain a list of individuals with their own probability of re-identification

Pseudo ID	Probability
100	1/3
101	1/10
102	1/50
203	1/30
205	1/25
...	...
452	1/30

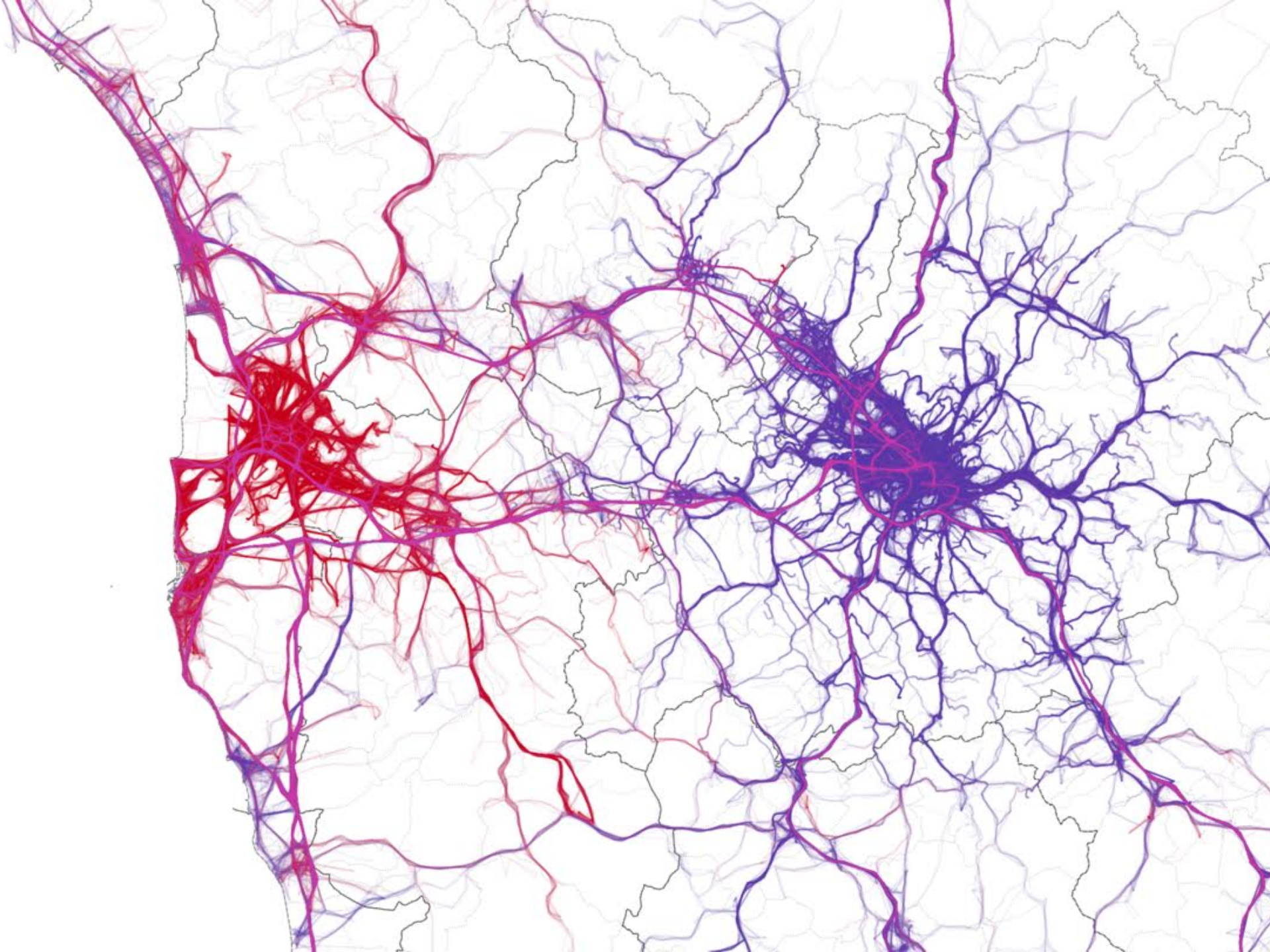
What next?



- Having in mind a privacy threshold (e.g., $1/20$)
- We see that many of our individual are already safe
- We can act (anonymizing) only on the other ones (e.g., 100&101)

Pseudo ID	Probability
100	$1/3$
101	$1/10$
102	$1/50$
203	$1/30$
205	$1/25$
...	...
452	$1/30$

But we need to go further!

- A city cannot be managed centrally, from a control room.
- Our cities are complex networks of interactions
 - the outcome for everybody depends not only on individual choices but it is conditioned by everybody else's choices.



- 
- A granular capability of citizens to self-organize, collaborate and coordinate their actions from the bottom-up is more efficient and resilient
 - But requires to align individual interests and goals with those of the collectivity in the system.
 - We humans have a limited perception of ourselves as a social, collective living being
- 

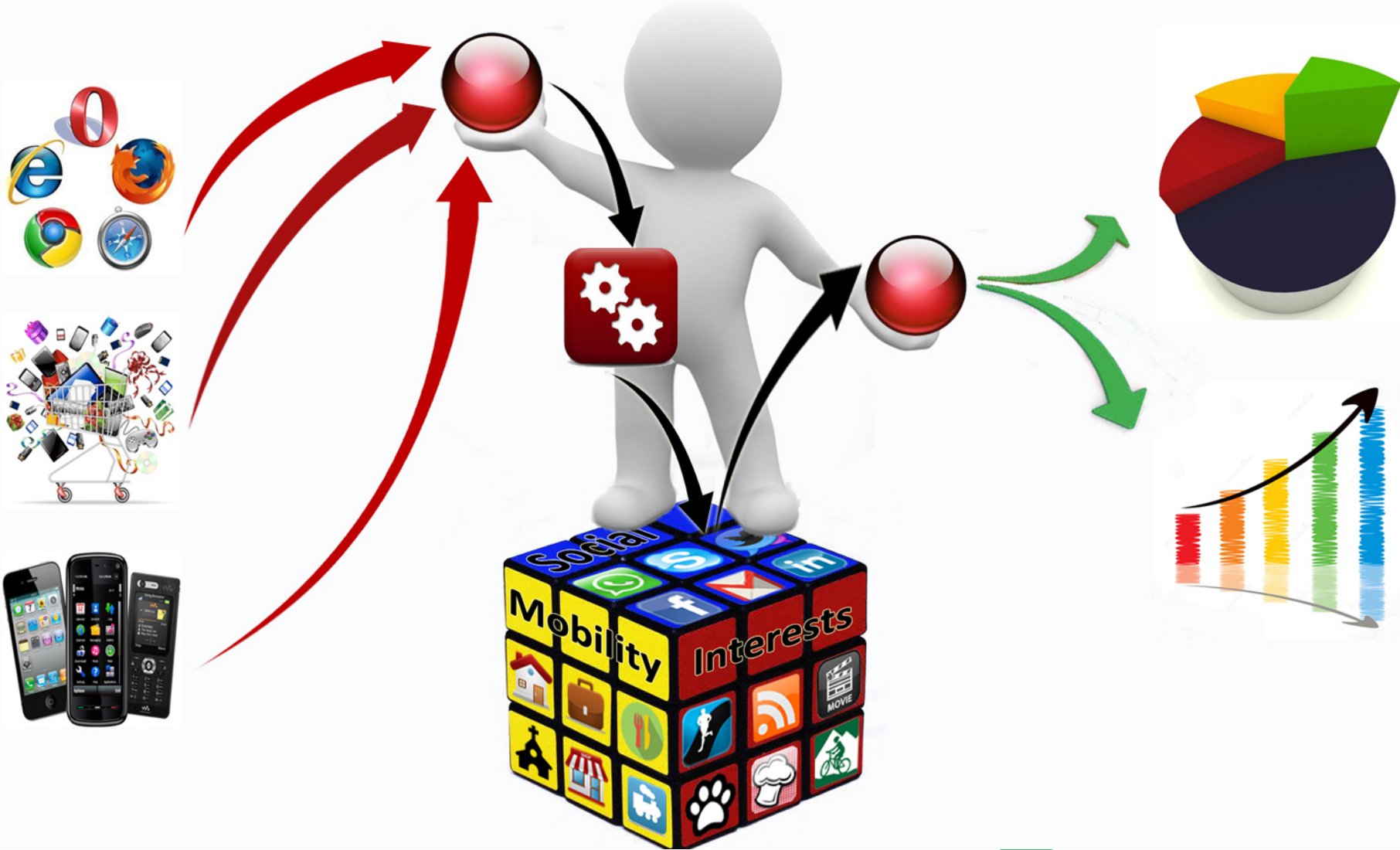


TOWARDS A PERSONAL DATA ECOSYSTEM





A user-centric ecosystem for personal big data



Personal Data Ecosystem



Where am I? Comparison with the community

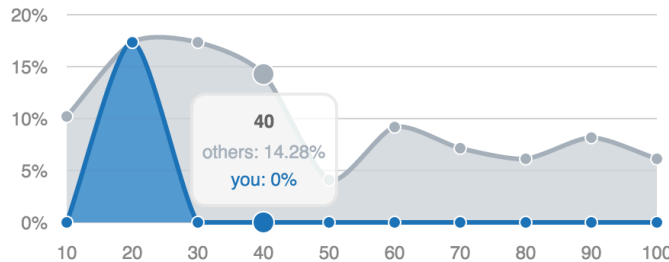
MyRoutine Mario Rossi

mariorossi ▾

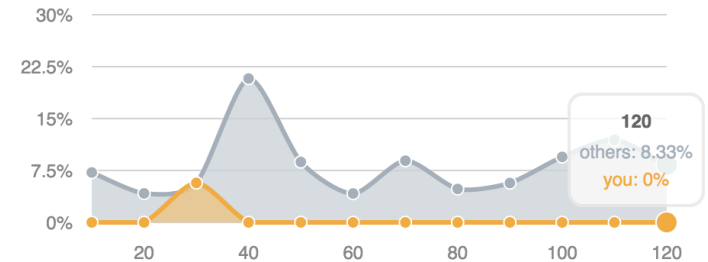
- Home
- Mobility Network
- Shopping Profile
- Where I Am?
- Statistics



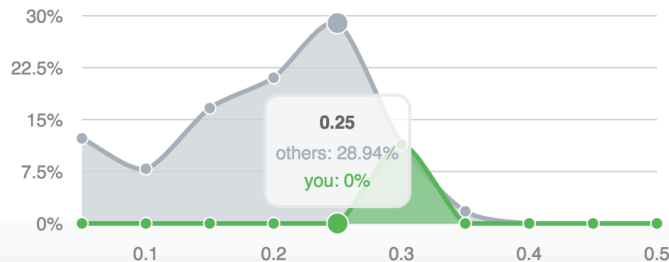
Radius of Gyration



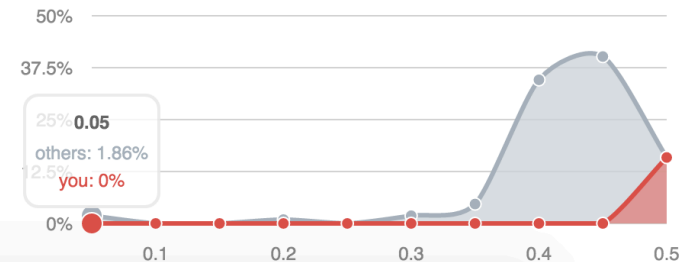
Travel Time



Basket Predictability




Time and Space Predictability





Towards a User centric data market

- We need a Personal Data Ecosystem
 - to acquire, integrate and make sense of our own data
 - and to connect with our peers and the surrounding urban community and infrastructure
 - to the purpose of developing the **collective awareness** needed to face our grand challenges
- 

A smart city is a city of participating, aware citizens



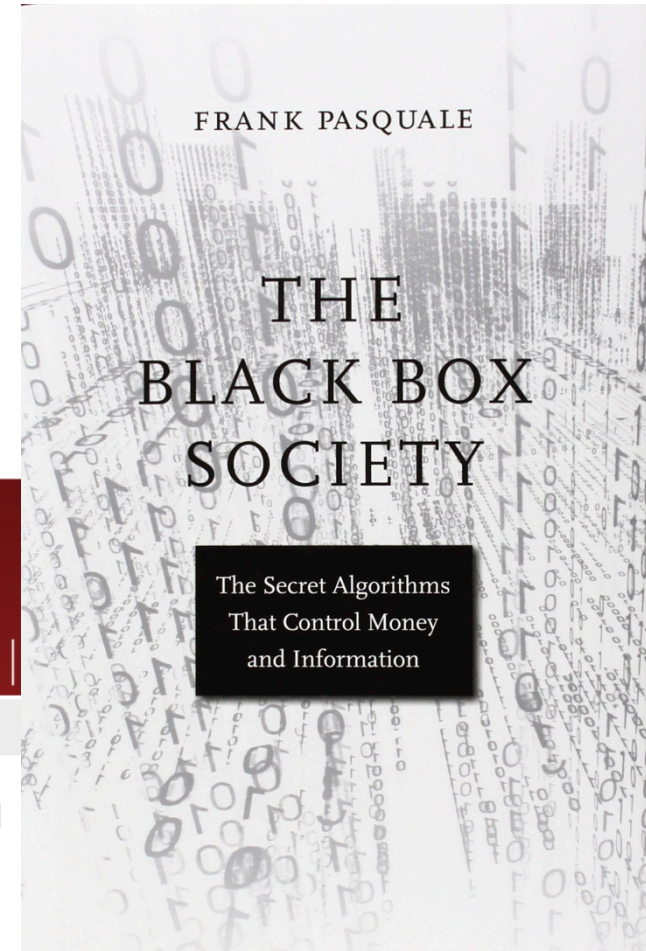


RIGHT TO EXPLANATION



Transparent algorithms to build trust

- **Systems that recommend humans making a decision should explain why**



nature International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video

Archive > Volume 537 > Issue 7621 > Editorial > Article

NATURE | EDITORIAL



More accountability for big-data algorithms

To avoid bias and improve transparency, algorithm designers must make data sources and profiles public.

21 September 2016

Big Data, Big Risks

- **Big data is algorithmic, therefore it cannot be biased!**
And yet...
- All traditional evils of social discrimination, and many new ones, exhibit themselves in the big data ecosystem
- Because of its tremendous **power**, massive data analysis must be used **responsibly**
- Technology alone won't do: also need **policy**, **user involvement** and **education** efforts



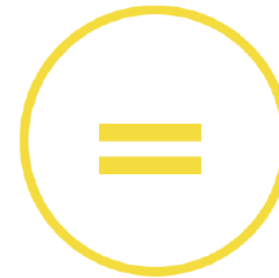
Fairness





Diversity



Transparency




Neutrality

- 
- By 2018, 50% of business ethics violations will occur through improper use of big data analytics
 - [source: Gartner, 2016]
- 



The danger of black boxes

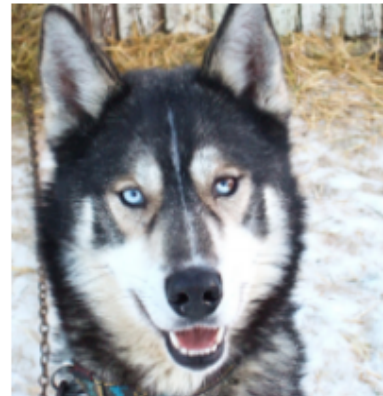
- The COMPAS score (Correctional Offender Management Profiling for Alternative Sanctions)
 - A 137-questions questionnaire and a predictive model for “risk of crime recidivism.” The model is a proprietary secret of Northpointe, Inc.
 - The data journalists at propublica.org have shown that the model has a strong ethnic bias
 - blacks who did not reoffend are classified as high risk twice as much as whites who did not reoffend
 - whites who did reoffend were classified as low risk twice as much as blacks who did reoffend.
- 

The danger of black boxes

- An accurate but untrustworthy classifier may result from an accidental bias in the training data.
- In a task of discriminating wolves from huskies in a dataset of images, the resulting deep learning model is shown to classify a wolf in a picture based solely on ...

The danger of black boxes

- An accurate but untrustworthy classifier may result from an accidental bias in the training data.
- In a task of discriminating wolves from huskies in a dataset of images, the resulting deep learning model is shown to classify a wolf in a picture based solely on ... **the presence of snow in the background!**

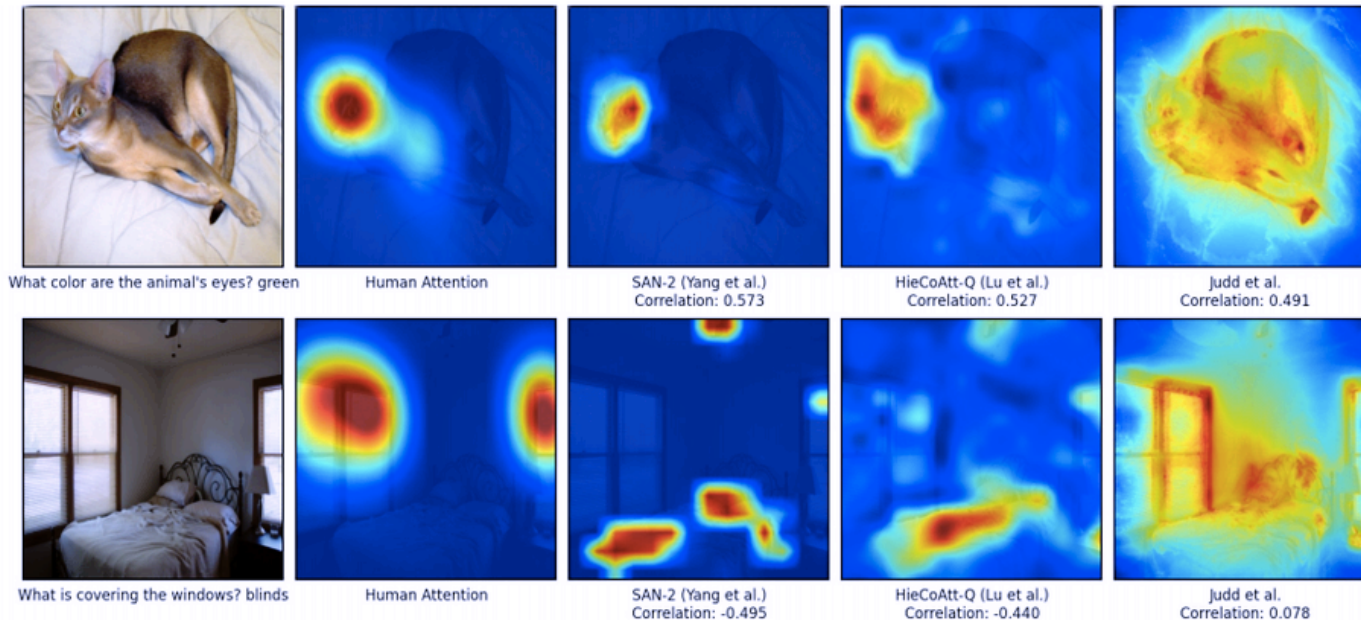


(a) Husky classified as wolf



(b) Explanation

Deep learning is creating computer systems we don't fully understand



"THEY'RE PICKING [ANSWERS] BASED ON BIASES IN THE DATA SETS, RATHER THAN FROM FACTS ABOUT THE WORLD."



TOWARDS EXPLANABLE AI



Goal

- Develop a logical/statistical framework consisting of a family of *algebras of rules* of adequate expressiveness, designed to tackle two tasks in a systematic way:
- the *explanation by design* – **XbD** – problem:
 - given a dataset of **training decision records**, how to develop a machine learning decision model *together with its explanation*;
- the *black box explanation* – **BBX** – problem:
 - given the decision records produced by an inscrutable black box decision model, how to *reconstruct an explanation* for it.

Discrimination-aware Data Mining

Dino Pedreschi Salvatore Ruggieri Franco Turini

Dipartimento di Informatica, Università di Pisa
L.go B. Pontecorvo 3, 56127 Pisa, Italy
{pedre,ruggieri,turini}@di.unipi.it

KDD'08, August 24–27, 2008, Las Vegas, Nevada, USA.
Copyright 2008 ACM 978-1-60558-193-4/08/08 ...\$5.00.



FROM DISCRIMINATION DISCOVERY TO BLACK BOX EXPLANATION



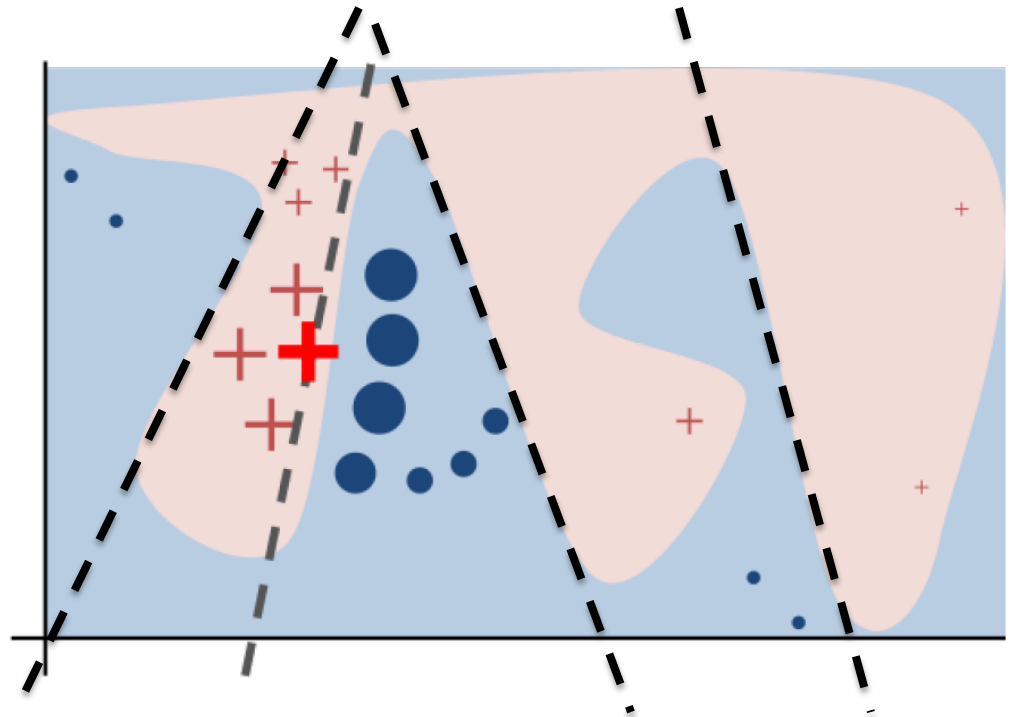
Research direction

- Rule-based discrimination discovery can be generalized to systematically tackle the broader problem of **explaining black box decision making systems** in an **agnostic way**
 - without taking into account the internals of the decision model, either **algorithmic, human, or combination thereof**.
- An **explanation** is a comprehensible representation of a decision model, acting as an interface between the model and the human.

From local rules to global explanations

Issues:

- Rule language – expressiveness vs comprehensibility
- Coverage – find enough rules to capture the whole data/decision space
- Simplification – find optimally simple set of rules by reasoning on/manipulating rules
- Fidelity proxying of the black-box behavior



Original picture from
Ribeiro et al., KDD 2016

- **A Survey Of Methods For Explaining Black Box Models**
- [Riccardo Guidotti, Anna Monreale, Franco Turini, Dino Pedreschi, Fosca Giannotti](#)
- <https://arxiv.org/abs/1802.01933>
- Submitted, 2018

Ethically Aligned Design

ETHICS+IN+ACTION >

The IEEE Global Initiative on Ethics of
Autonomous and Intelligent Systems



CHARACTER

MORALITY

ETHOS

CONVENTION

GOODNESS

READ *Ethically Aligned Design,*

The most comprehensive, crowd-sourced
global treatise regarding the Ethics of Autonomous
and Intelligent Systems available today.

FAIR

First Aid for
Responsible
data scientist

TRY IT ON:

[fair.sobigdata.eu/
moodle](https://fair.sobigdata.eu/moodle)



The SoBigData online course developed to ensure that people are familiar with the basic elements about ethics, data protection, and intellectual property law

SoBigData

FAIR

FAIR – First Aid for Responsible Data Scientists by SoBigData

You are not logged in. (Log in)

FIRST AID FOR DATA SCIENTIST

This course rises within the EU SoBigData project: we developed this online course in order to make sure that all users are familiar with the basic elements about: ethics, data protection, and intellectual property law.

Access to the platform

Username

ML

Password

.....

Log in

Forgotten your username or password?

New account



A SoBigData initiative

The European Research Infrastructure on Big Data



For Researchers

who want to become more aware of ethical issues



For Companies

which want to contribute to our community



For Students

who are curious about ethics, privacy and law