# SoBigData Research Infrastructure

## Social Mining & Big Data Analytics

# RRI & Data Science
# Anna Monreale

## Dipartimento di Informatica
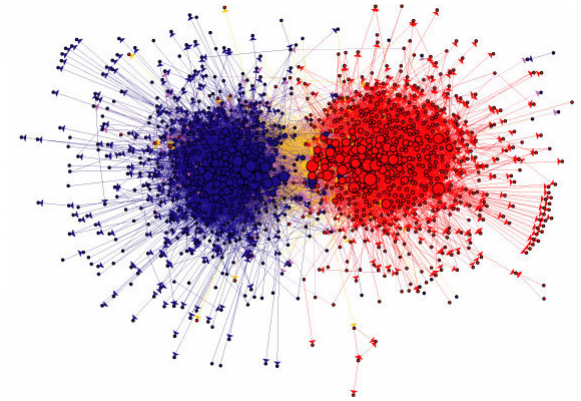## Università di Pisa

# Our digital traces ….

- We produce an unthinkable amount of data while running our daily activities.

- How can we manage all these data? Can we get an added value from them?

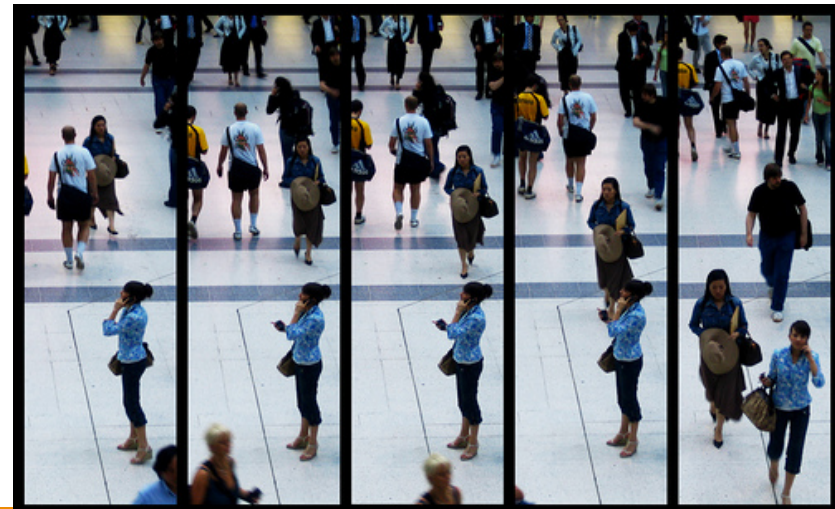# Big data "proxies" of social life

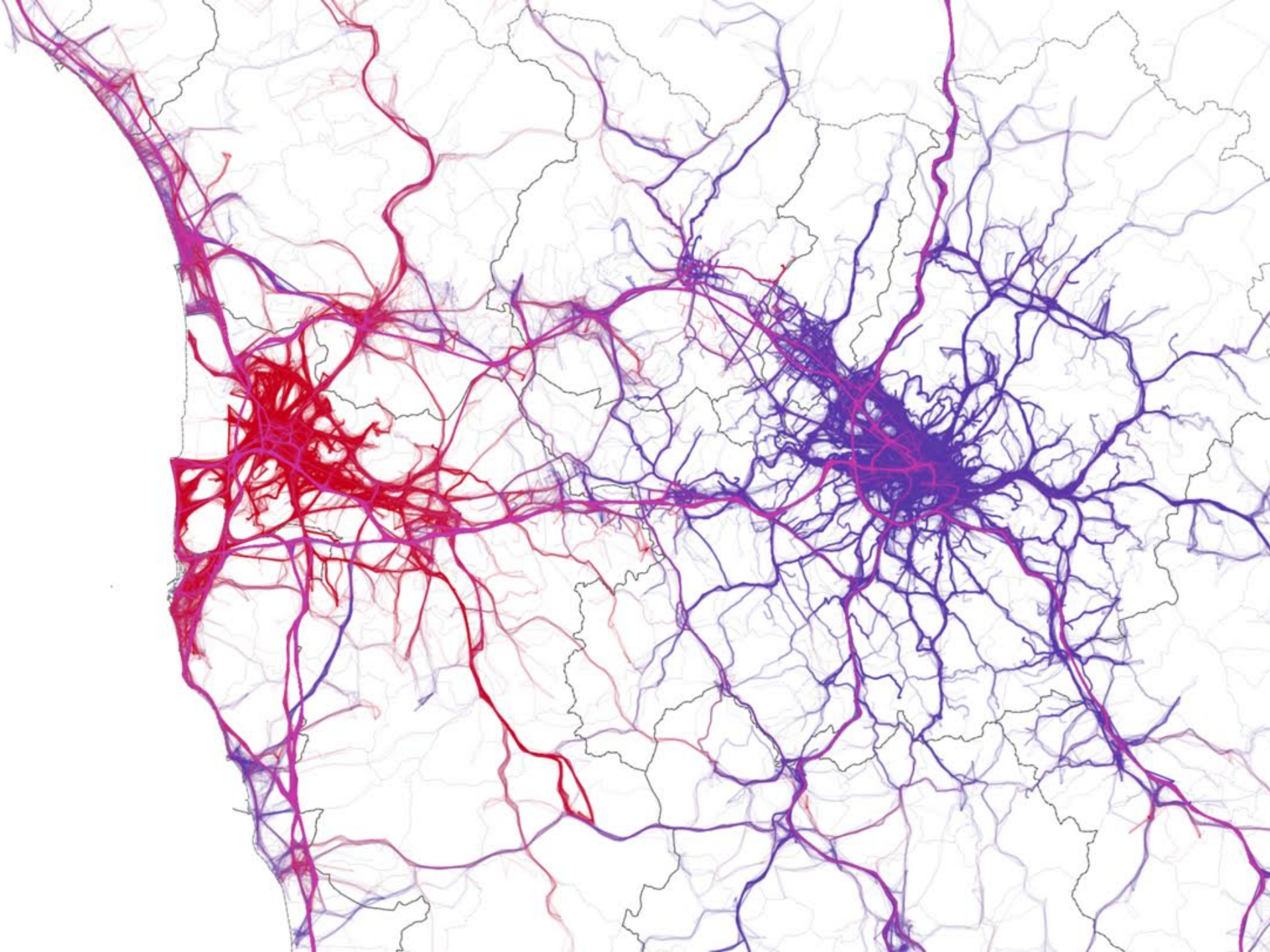**Shopping patterns & lifestyle**

Relationships & social ties

Desires, opinions, sentiments

Movements

2005

Luca Bruno / AP

2013

NBC NEWS

Michael Sohn / AP

# Big Data Analytics & Social Mining

The **main tool** for a **Data Scientist** to measure, understand, **and possibly** predict **human behavior**
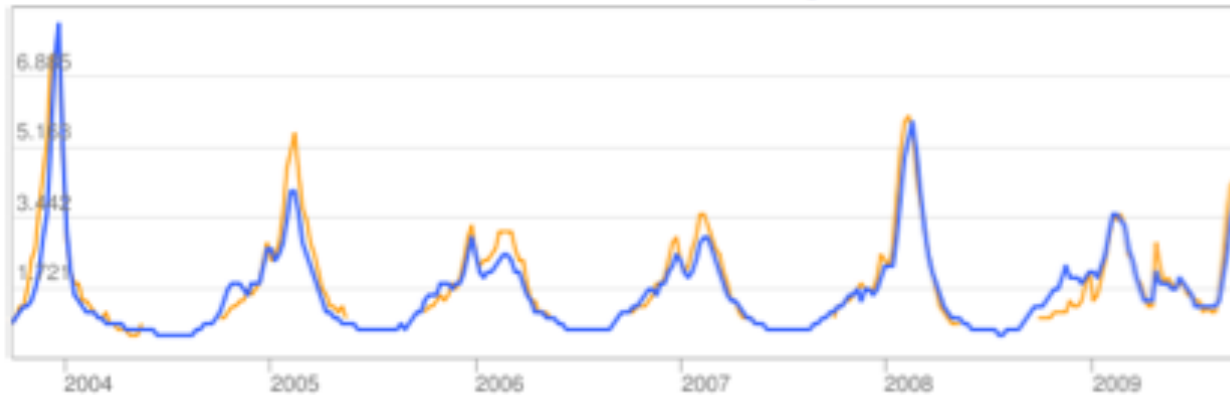
# Nowcasting epidemics



Nature 457, 1012-1014 (2009)

# Mobility Atlas of a City

# Estimation of users city categories by mobile phone data



GSM Calls

Temporal Profile

Profile Map

Unknow 10%

Visitors 20%

Commuter 40%

Residents 20%

Dynamic Residents 10%

# Big Data Analytics & Social Mining

The **main tool** for a **Data Scientist** to measure, understand, **and possibly** predict **human behavior**

**Data Scientist needs to take into account ethical and legal aspects and social impact of data science**

# The perils of big data & data science

- Not only "**privacy**" – (lack of) **protection of personal data**

- Lack of **transparency** on use of data

- Huge **asymmetry** between users' and company's information

- **Access to data**, even own's personal digital traces

- Potential **discrimination** due to profiling

# Legislation Knowledge

- The data scientist should know the **notions, concepts and principles** of the data protection legislation
  - General Data Protection Regulation, European Data Protection Directive
- The data scientist should know the **responsibilities** designed by the laws
  - Data Controller, Data Processor
- The data scientist should act in **compliance with data protection law principles**
  - fair and lawful processing, purpose limitation, **privacy by design** and by default
  - Processing personal data only on a legal basis (e.g. consent, research exceptions etc.)
  - Implement appropriate technical and organizational measures to protect the data and guarantee the privacy right of individuals

**SoBigData** Research Infrastructure

Social Mining & Big Data Analytics

# How to address these issues?

# SoBigData Data Strategies

- Strategies for supporting actors (Data Scientists) in SoBigData RI
  - **Online Training Material & Compliance Self-assessment**
  - **Tool for the privacy risk assessment**

**- On-line Training Material**
**- Self Assessment**

⠿ Apps ▣ Qwertee ⑯ Google ▢ NY ▢ occhiali sole ▢ navigatore ▢ wow ▢ cucina ⑯ Google Traduttore ⱳ Dizionario Italiano-Ing ▢ PhD ▢ nailart ⓫ DEATH NOTE + Speci ◔ Universita' di Pisa We ⓭ Doodle: Book Piazza ⟫ ▢ Other bookmarks

Search News Feed

Go to ⏷    20    ① Francesca Pratesi  ⏷

# SoBigData

| SoBigData.eu | Members | Activity Tracker | Wiki | Twitter Monitor | Stat. Algorithm Importer |

### Statistics

**Your Stats in SoBigData.eu**

ACTIVITY        GOT

↗0  ♡0  💬0     ♡0  💬0

### Top Topics

#wp2
#review

### Recent Documents

📕 Invitation to a Project Rev...
📕 main.pdf
📕 icwsm_main.pdf
📕 icdm_main.pdf
📊 WP9 NDLib.pptx

Show all ...

### Authorisation Options

Personal Token

**About Personal Token**

The personal token has to be used for any programmatic interaction with the services you perform to satisfy your needs.

## Welcome on the On-line Training Material and Self Assessment

Here you can find some basic notion about legal and ethical implication on the use of data and methods available in the SoBigData Research Infrastructure.
If you want to start (it will be necessary about 15 min to complete this educational path), click here

### START!

If you think you already are enough confident and aware about this, you can directly go to the questionaire

### Go to the test

### About

# SoBigData

SoBigData proposes to create the Social Mining & Big Data Ecosystem: a research infrastructure (RI) 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'. Building on several established national infrastructures, SoBigData will ope ...

See more

Other options ...

### VRE Managers and Groups

👥 View Managers

Groups in this VRE:

👥 Production-Support

### Invite Members

e-mail address

Send Invite

Report an issue

https://sobigdata.d4science.org/group/sobigdata.eu

Apps  Qwertee  G Google  NY  occhiali sole  navigatore  wow  cucina  G Google Traduttore  wR Dizionario Italiano-Ing  PhD  nailart  DEATH NOTE + Speci  Universita' di Pisa Wel  Doodle: Book Piazza  »  Other bookmarks

Search News Feed

Go to ▾   20   Francesca Pratesi ▾

# SoBigData

SoBigData.eu | Members | Activity Tracker | Wiki | Twitter Monitor | Stat. Algorithm Importer

**Statistics**

Your Stats in SoBigData.eu

ACTIVITY    GOT

↱0 👍0 💬0    👍0 💬0

**Top Topics**

#wp2
#review

**Recent Documents**

📄 Invitation to a Project Rev...
📄 main.pdf
📄 icwsm_main.pdf
📄 icdm_main.pdf
📄 WP9 NDLib.pptx

Show all ...

**Authorisation Options**

Personal Token

**About Personal Token**

The personal token has to be used for any programmatic interaction with the services you perform to satisfy your needs.

# Ethical Big Data Research

At SoBigData, we aim to promote *ethical big data research*. Big data research is ethical when it attempts to <u>maximise the societal benefits of research, while minimising the harms</u>. This means:
- being clear about the research **purpose**,
- gaining data subjects' **consent**,
- being **transparent** about how data is being processed
- **minimising** potentially sensitive personal data.

We strive to create an environment in which researchers develop, share and improve methodologies for ethical big data research. Our research infrastructure partners you with researchers dealing with the same problems and allows you to share methodologies compliant with legal requirements.

View supplementary material

→

**About**

# SoBigData

SoBigData proposes to create the Social Mining & Big Data Ecosystem: a research infrastructure (RI) 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'. Building on several established national infrastructures, SoBigData will ope ...

See more

Other options ...

**VRE Managers and Groups**

👥 View Managers

Groups in this VRE:

👥 Production-Support

**Invite Members**

e-mail address

Send Invite

Report an issue

## SoBigData

SoBigData.eu | Members | Activity Tracker | Wiki | Twitter Monitor | Stat. Algorithm Importer

**Statistics**

Your Stats in SoBigData.eu

| ACTIVITY | GOT |
|---|---|
| ↗0 👍0 💬0 | 👍0 💬0 |

**Top Topics**

#wp2
#review

**Recent Documents**

📄 Invitation to a Project Rev...
📄 main.pdf
📄 icwsm_main.pdf
📄 icdm_main.pdf
📄 WP9 NDLib.pptx
Show all ...

**Authorisation Options**

Personal Token

**About Personal Token**

The personal token has to be used for any programmatic interaction with the services you perform to satisfy your needs.

# Data Protection Law in the EU

At present, the centerpiece of the European data protection legislation is the Data Protection Directive (DPD) and implementing national laws. From 25 May 2018 a new legal instrument - the General Data Protection Regulation (**GDPR**) – will be directly applicable in all Member States.

| View 95/46/EC Dir. | View GDPR |
|---|---|

Personal data are "any information relating to an identified or identifiable natural person ('data subject') [...]" (GDPR - Article 4(1))
Sensitive data are [...]

← →

**About**

## SoBigData

SoBigData proposes to create the Social Mining & Big Data Ecosystem: a research infrastructure (RI) 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'. Building on several established national infrastructures, SoBigData will ope ...

See more

Other options ...

**VRE Managers and Groups**

👥 View Managers

Groups in this VRE:

👥 Production-Support

**Invite Members**

e-mail address

Send Invite

Search News Feed

Go to ▾     20     Francesca Pratesi ▾

# SoBigData

SoBigData.eu | Members | Activity Tracker | Wiki | Twitter Monitor | Stat. Algorithm Importer

**Statistics**

Your Stats in SoBigData.eu

ACTIVITY          GOT

↪0  👍0  💬0    👍0  💬0

**Top Topics**

#wp2
#review

**Recent Documents**

📄 Invitation to a Project Rev...
📄 main.pdf
📄 icwsm_main.pdf
📄 icdm_main.pdf
📄 WP9 NDLib.pptx

Show all ...

**Authorisation Options**

Personal Token

**About Personal Token**

The personal token has to be used for any programmatic interaction with the services you perform to satisfy your needs.

## Does or will your research involve <u>personal data</u>?

Yes          No          I don't know

## Does this data constitute <u>sensitive data?</u>

Yes          No          I don't know

## Have you <u>pseudonymised</u> your data set?

Yes          No          I don't know

## Have you sent out a notice to the subjects of your data set?

Yes          No          I don't know          Not applicable

**About**

# SoBigData

SoBigData proposes to create the Social Mining & Big Data Ecosystem: a research infrastructure (RI) 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'. Building on several established national infrastructures, SoBigData will ope ...

See more

Other options ...

**VRE Managers and Groups**

👥 View Managers

Groups in this VRE:

👥 Production-Support

**Invite Members**

e-mail address

Send Invite

Report an issue

Search News Feed          Go to ▾   20   Francesca Pratesi ▾

# SoBigData

SoBigData.eu | Members | Activity Tracker | Wiki | Twitter Monitor | Stat. Algorithm Importer

**Statistics**

Your Stats in SoBigData.eu

|  | ACTIVITY | GOT |
|---|---|---|
| | ↗0 👍0 💬0 | 👍0 💬0 |

**Top Topics**

#wp2
#review

**Recent Documents**

📄 Invitation to a Project Rev...
📄 main.pdf
📄 icwsm_main.pdf
📄 icdm_main.pdf
📊 WP9 NDLib.pptx

Show all ...

**Authorisation Options**

Personal Token

**About Personal Token**

The personal token has to be used for any programmatic interaction with the services you perform to satisfy your needs.

---

## Does or will your research involve <u>personal data</u>?

Yes          No          I don't know

## Does this data constitute <u>sensitive data</u>?

Yes          No          I don't know

## Have you <u>pseudonymised</u> your data set?

Yes          No          I don't know

## Have you sent out a notice to the subjects of your data set?

Yes          No          I don't know          Not applicable

Self assessment completed
Even if your research involve personal data, you know the right process to deal with it. Well done!

**Access the RI**

---

**About**

# SoBigData

SoBigData proposes to create the Social Mining & Big Data Ecosystem: a research infrastructure (RI) 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'. Building on several established national infrastructures, SoBigData will ope ...

See more

Other options ...

**VRE Managers and Groups**

👥 View Managers

Groups in this VRE:

👥 Production-Support

**Invite Members**

e-mail address

Send Invite

Report an issue

# Privacy Risk Assessment

# Methodology for PRA

- Service and data format definition

- External information definition

- Simulation of privacy harmful Inferences

- Vulnerability (Risk) quantification

- (Risk mitigation)

# Methods

- Privacy Risk Assessment

    - **Human call habits data (CDR aggregation)**

    - Trajectory data (from GPS observations)

    - Individual shop habits

    - ...

- Privacy Risk Mitigation

# Definition of service
# Classifying user behaviour (Sociometer)

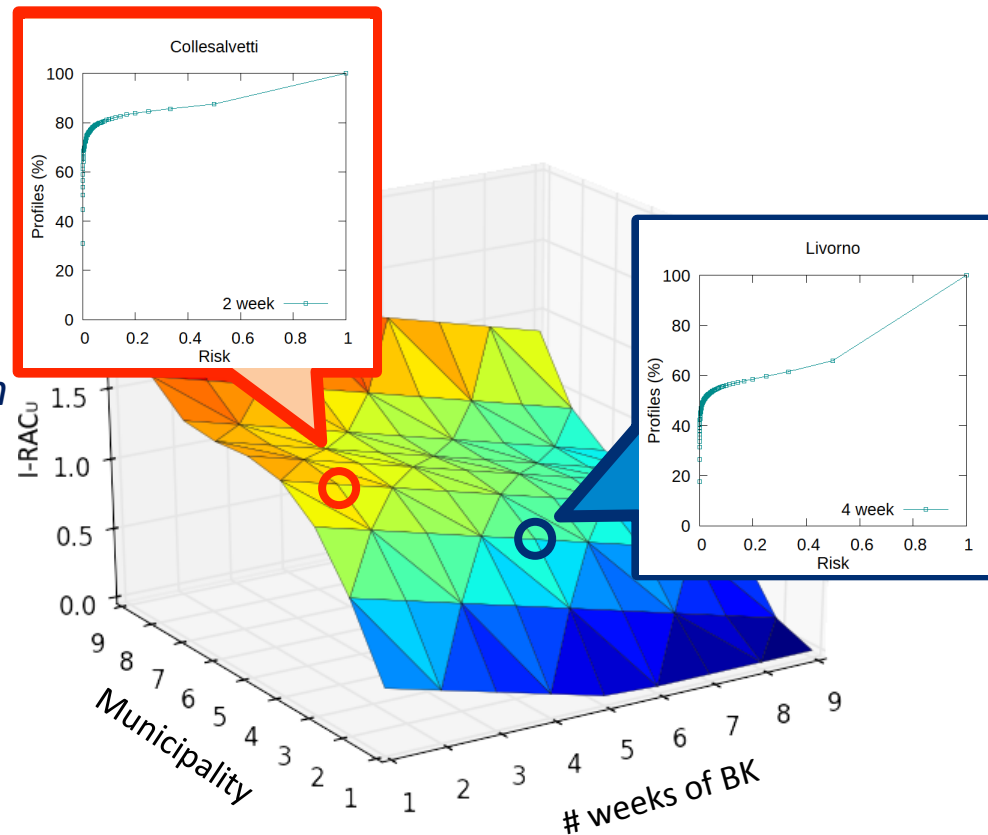# Simulation of privacy harmful Inferences

**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*

# RRI & Data Scientist

- Data Scientist is a professional figure with a mix of competence and knowledge
  - On methods and technologies for managing large amount of data
  - On analytical techniques and modelling of data and data mining
  - On story-telling techniques and data visualization
  - On **ethical and legal aspects** and social impact of data science
  - On appropriate technical and organizational measures for data protection

# Thank you!!!!

**Anna Monreale**

Dipartimento di Informatica

Università di Pisa

anna.monreale@unipi.it

Knowledge Discovery and Delivery Lab
(ISTI-CNR & Univ. Pisa)
www-kdd.isti.cnr.it