# Data Science for Historical Inquiries



**Donato Malerba**

*Computer Science Dept. – University of Bari*
*Big Data Lab - CINI*

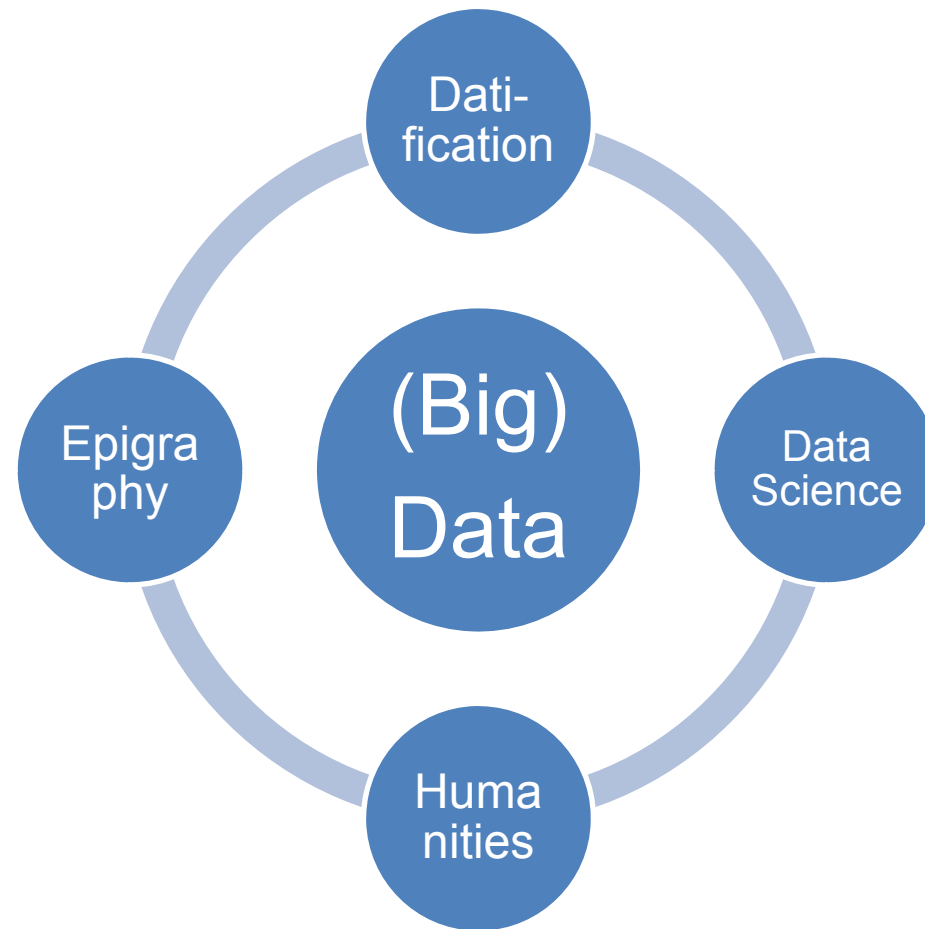**Off the Beaten track**
Bari, 25th September 2015
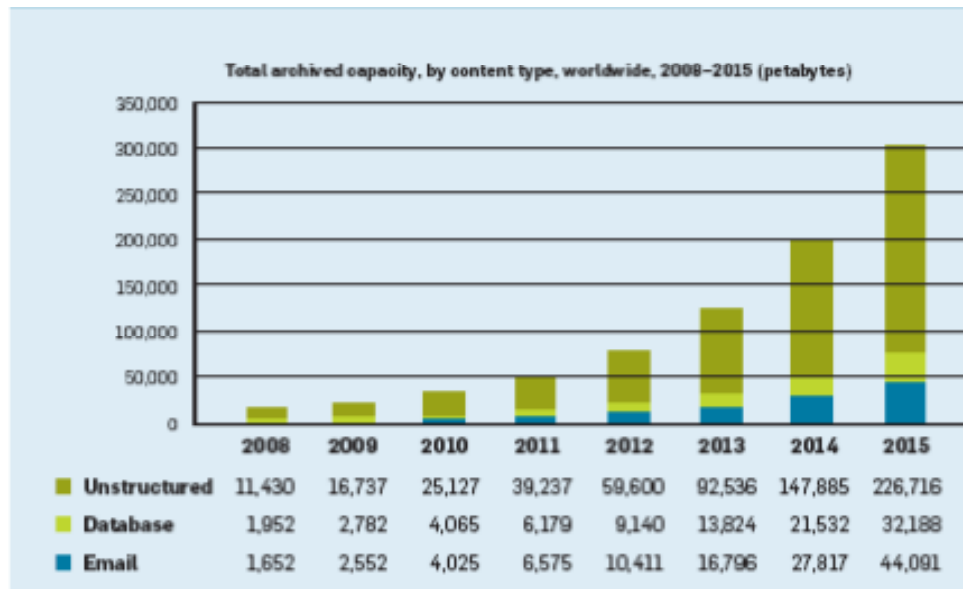
UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini
consorzio
interuniversitario
nazionale
per l'informatica

# An estimated growth ...

- Yearly increase of data volume: 30-40%
- The data volume doubles every 2.5 years



Total archived capacity, by content type, worldwide, 2008–2015 (petabytes)

| | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |
|---|---|---|---|---|---|---|---|---|
| Unstructured | 11,430 | 16,737 | 25,127 | 39,237 | 59,600 | 92,536 | 147,885 | 226,716 |
| Database | 1,952 | 2,782 | 4,065 | 6,179 | 9,140 | 13,824 | 21,532 | 32,188 |
| Email | 1,652 | 2,552 | 4,025 | 6,575 | 10,411 | 16,796 | 27,817 | 44,091 |

Exponential growth
- 2,7 ZB ($10^{21}$ bytes) in 2012
- 35 ZB in 2020

# *Datification*

- Neologism referring to the conversion into digital format (data) of:

- Movies, music, books, etc. (previously on films, paper, vinyl, etc.)

- Phone conversations, mail, radio and TV programs

# *Datification*

- *Facebook*: our social networks are data,

- *Twitter*: our sentiments are data,

- *LinkedIn*: our professional experiences are data

# *Devices generate data …*

- **Internet of Things** (IoT): each thing is uniquely identifiable and is able to interoperate within the existing Internet infrastructure. Things have an active role and exhibits "intelligent" behaviors.

# Devices generate data ...
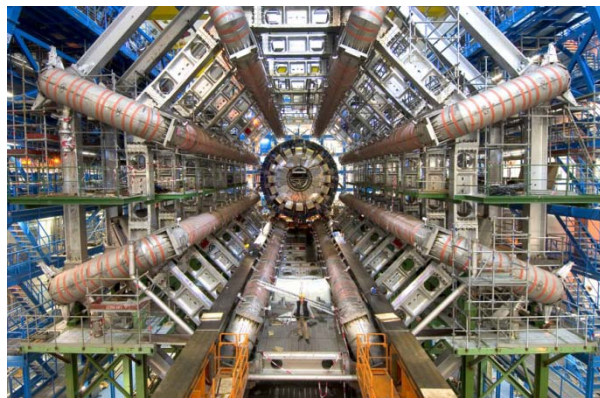
- Due to our symbiosys with digital technologies we are becoming "living sensors"

- 7 billion people and 6,8 billion mobile phones

- «*We are like a digital Tom Thumb, who leaves (digital) crumbs along the way*».

# *Science generates data ...*

- Large volumes of data are now generated in:
  - **Genomics** ➜ next generation sequencing
  - **Astrophysics** ➜ astronomical observatories
  - **Particle Physics** ➜ Large Hadron Collider
  - **Neuroscience** ➜ Human Brain project
  - ...

# Companies generate data ...

- Nowaday every big business is a digital business:
  - **Alibaba** the largest shop in the world without warehouse.
  - **Uber** the largest company in the world without cars.
  - **Airbnb** the largest network for lodging without a single hotel.

# *Companies generate data ...*

- Purchase orders, invoices, dispatches, ...

- Data collected in information systems are considered an *(intangible) asset.*

- *Facebook:* (tangible) *assets* for 6,3 billion$ but at its stock market flotation its value was 104 billion$.

- Data are an intangible asset, but only **5‰** of business data are actually processed

# The data deluge

February 2010

# Big Data

- Data collection has gotten to the point where the volume, the velocity and the variety of incoming data demands from **non conventional tools** to extract and process information in (near) rela time.

# Big Data: a revolution?

- *The true revolution is not in technologies we use to process data, but in the way we use data.*

- *The larger scale of processed data paves the way to new potential applications which would not be possible otherwise*

# *Data Science*

- Data science is a body of principles and techniques for applying data-intensive analysis to investigate phenomena, acquire new knowledge, and correct and integrate previous knowledge with measures of correctness, completeness, and efficiency of the derived results.

# *Big Data vs. Data Science*

- Data Science vs. Big Data
  - Data Science does not always need Big Data, however the steady increase of data makes Big Data an important issue for Data Science.

# Big Data: *an example*

- Spring 2009: pandemic influenza caused by the virus A (H1N1) emerged, beginning in Mexico and quickly spreading to the United States and around the world

  *USA Centers for Disease Control and Prevention*:

  state and local health departments collected data on influenza-like cases + epidemiological models
  ➔ a lag of two weeks

# Big Data: *An example*

**Google Flu Trends**: uses anonymized, aggregated internet search activity to provide near-real time estimates of influenza activity

**Stime storiche**                                    Visualizza dati per: | Stati Uniti |

*Detecting influenza epidemics using search engine query data.*
**Nature 457, 1012-1014 (19 February 2009)**

Attività influenzale Stati Uniti

Stima sull'influenza          ● Stima di Google Trend influenzali ● Dati Stati Uniti

8.866

6.650

4.433

2.217

2004   2005   2006   2007   2008   2009   2010   2011   2012   2013

Stati Uniti: dati ILI (Influenza-Like Illness) forniti pubblicamente dagli U.S. Centers for Disease Control.

UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini
consorzio interuniversitario nazionale per l'informatica

# Big Data: *A change of perspective*

- Three Big changes in the way data are analyzed:

  1. Process all available data (the obsolescence of sampling)

  2. Accept increased measurement error in return for more data

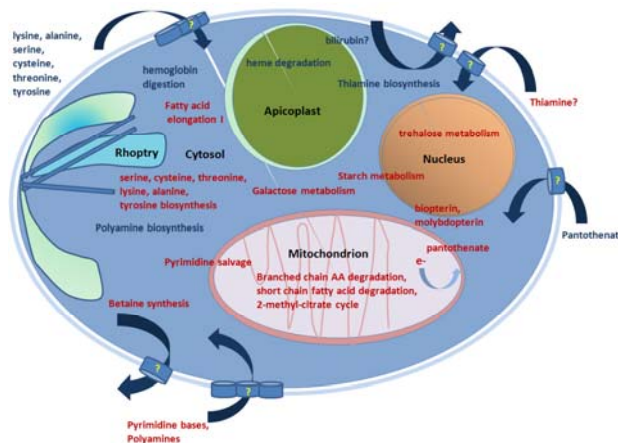  3. Move away from the age-old search for causality (focus shift, from causation to correlation)

# Big Data for a *data-driven science*

- Big Data are causing a radical change in scientific paradigms.

- Thomas Kuhn: *a scientific revolution is a paradigm shift*

- Two traditional paradigms in science:
  - Experimentation
  - Theory

# Two new paradigms

- Computational simulation (or third paradigm): Ken Wilson, Nobel laureate in physics, 1982.

- At the base of the *computational sciences*
  - *computational biology*, simulate the behaviour of biological systems, *metabolic pathways* or a cell or the way a protein is produced.

# Two new paradigms

- Data intensive knowledge discovery (or fourth paradigm): Jim Gray, computer scientist.

- At the base of the *science informatics*

  - *bioinformatics*, an interdisciplinary field that develops methods and software tools for understanding biological data.

# Data-intensive Knowledge Discovery

Steps:

- Data capture

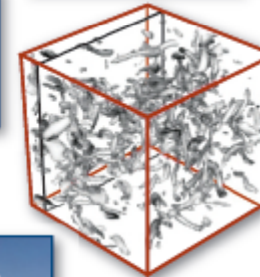- Data curation

- Data analysis

- Result publishing

# The Science Paradigms

## Science Paradigms

- Thousand years ago:
  science was **empirical**
  *describing natural phenomena*

- Last few hundred years:
  **theoretical** branch
  *using models, generalizations*

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi Gp}{3} - K\frac{c^2}{a^2}$$

- Last few decades:
  a **computational** branch
  *simulating complex phenomena*

- Today: **data exploration** (eScience)
  *unify theory, experiment, and simulation*

  – Data captured by instruments
    or generated by simulator

  – Processed by software

  – Information/knowledge stored in computer

  – Scientist analyzes database/files
    using data management and statistics

UN
DEGI
ALDO MORO

consorzio
interuniversitario
nazionale
per l'informatica

# What's the impact of these trends on Humanities?

- Large-scale digitization projects have vastly increased the quantity of cultural heritage material in several humanistic areas.

- Humanities scholars are increasingly incorporating computational tools and methods in all phases of their research.

- Large investments in digital infrastructures supported by funding agencies.

UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini

consorzio
interuniversitario
nazionale
per l'informatica

# What's the impact of these trends on Humanities?

- Most researchers in humanities explore data manually, using their knowledge and expertise to extract the information they deem relevant.

# What's the impact of these trends on Humanities?

Example:

Els Witte (2006) studied the image of the nation in the Belgian Revolution (1828-1847) by manually browsing six newspapers and selecting 350 articles that expressed an opinion.

- This is a nice example of a new perspective on the study of "public opinion" or collective "mentalities"

# What's the impact of these trends on Humanities?

- However, there are corpora for historical research that are simply too large to be examined in their entirety and to be inspected manually

- Sampling reduce the amount of data to manageable proportions, but the manual inspection strongly limits the subsequent analysis.

# What's the impact of Data Science on Humanities?

- One of the promise of the <span style="color:red">convergence of data science and humanities</span> is that it will enable us to investigate much larger quantities of data.

- We are now entering a new phase in which historians are able to analyze massive volumes of data in various formats (records, texts, images, …)

# What's the impact of these trends on humanities?

- New techniques of large-scale data analysis allow historians to manage big data sets that were impossible to manage earlier.

- Data science can reduce the effort required of humanities researchers to obtain useful information from large repositories of digitized cultural heritage (e.g., medieval manuscripts)

# What's the impact of these trends on humanities?

Advantages of adapting Data Science methodologies to humanities:

- **Reproducibility** of results

- Promotion of **collaborative** work (in contrast to current research which is predominantly individualistic)

- **New research questions**

# What's the impact of Data Science on Epigraphy?

Several epigraphic repositories currently contain large corpora of pictures and the textual document representation thereof, which have been stored and annotated on several levels of interest.

http://www.eagle-network.eu,

http://edh-www.adw.uni-heidelberg.de,

http://www.edb.uniba.it,

http://eda-bea.es,

http://www.epigraphik.uni-hamburg.de,

http://usepigraphy.brown.edu,

http://www.edr-edr.it

UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini
consorzio
interuniversitario
nazionale
per l'informatica

# What's the impact of Data Science on Epigraphy?

A vital question:

*how, if at all, should the work of epigraphists adapt to the presence of orders of magnitude more potential source material?*

# What's the impact of Data Science on Epigraphy?

From the perspective of traditional research, <span style="color:red">little has changed</span>: the fact that most of recorded human intellectual output is now accessible does not increase the ability of an epigraphist to read it.

# What's the impact of Data Science on Epigraphy?

Evidently, any fundamental advances must come from the fact that this material is now available for computational processing.

# What's the impact of Data Science on Epigraphy?

We argue that <span style="color:red">a new and still unexplored frontier of digital epigraphy projects is that of enabling automatic analysis of information</span> currently stored in epigraphic repositories in order to extract implicit, previously unknown and potentially useful knowledge from them.

# What's the impact of Data Science on Epigraphy?

The application of Data Science practices may help to reveal interesting relationships among

- – linguistic style,

- – positioning,

- – dating of inscriptions,

- – ...

thus creating new links between different pieces of data.

# What's the impact of Data Science on Epigraphy?

In general, this approach can help to organize large collections of inscriptions, introduce younger scholars to the field of epigraphy, and identify anomalies that can be later explored using more traditional methods, as already done in computational historiography (Mimno, 2012).

# A concrete example: EDB

45.000 Christian inscriptions of Rome, including inscriptions published in the Inscriptiones Christianae Vrbis Romae septimo saeculo antiquiores, nova series (ICVR) editions.



www.edb.uniba.it

# A concrete example: EDB

## Data Sources



Stores the **text** of the epigraphs and a set of **metadata** about dating, original context, current location, related literature, etc.

# A concrete example: EDB

## Data Sources



Stores **photos of the epigraphs.**

This repository can be either internal or external, as well as an integration of internally produced and external resources.

# A concrete example: EDB

## Data Sources



Stores **metadata about scientific papers** in which epigraphs have been studied.

Each epigraph stored in the Epigraph Database can be associated to one or more papers stored in this database.

# A concrete example: EDB

## Data Sources



Stores data about **geographic positions**.

One or more geographic positions can be associated to each epigraph
(e.g. locations where the epigraph was found and/or the position where the epigraph is currently located).

# A concrete example: EDB

## Data Sources



**Supports retrieval tasks**.

It gives the possibility to enhance and/or perform query expansion to better match data stored in the epigraph database.

# A concrete example: EDB

## Data Sources



Stores the **knowledge** extracted by the Data Mining Module, i.e. patterns of interest in the form of rules, clusters, etc.

# A concrete example: EDB

## The Data Mining Module

Allows data analysts to **execute data mining algorithms** on the available data, in order to discover valuable knowledge, which is then stored in the Knowledge Base.



UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini consorzio interuniversitario nazionale per l'informatica

# A concrete example: EDB

- <span style="color:red">Change Detection</span>

- Since epigraphs can usually be associated to an historical period (dating), it could be interesting to identify how social and cultural changes over time have affected the epigraphs.

- Possible aspects to study...

Orthography

Used Materials

Phonetics

Executing Techniques

UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini consorzio interuniversitario nazionale per l'informatica

# A concrete example: EDB

- The data analysis task:

identification of emerging patterns, i.e. patterns which show relevant changes (in frequency) over time. The temporal dimension has a central role.

The emerging patterns can be ranked according to some measures of relevance, such as the **growth rate**, which represents the relative variation of the **support** of the pattern in the considered time intervals.

# A concrete example: EDB

- **Target objects**: epigraphs
- **Task-relevant objects**: materials, executing techniques, kinds of writing, etc.
- **Proper time intervals** on the basis of a discretization strategy (e.g. equal width, equal frequency, clustering-based discretization).
- **Features of interest** among all the available ones, in order to focus the algorithm only on the relevant data.

# A concrete example: EDB

Table I.  Dataset Obtained after the Application of EW Discretization (left) and EF Discretization (right)

| Block | Interval | All<br>Epigraphs | Precise Dating<br>Epigraphs | Block | All<br>Interval | Epigraphs | Precise Dating<br>Interval | Epigraphs |
|-------|----------|----------|----------|-------|----------|----------|----------|----------|
| 1 | 200–249 | 1,671 | 5 | 1 | 200–287 | 2,479 | 200–359 | 199 |
| 2 | 250–299 | 929 | 29 | 2 | 288–325 | 2,479 | 360–375 | 189 |
| 3 | 300–349 | 7,418 | 108 | 3 | 326–348 | 2,479 | 376–388 | 186 |
| 4 | 350–399 | 2,841 | 596 | 4 | 349–349 | 2,479 | 389–401 | 200 |
| 5 | 400–449 | 1,853 | 245 | 5 | 350–383 | 2,479 | 402–439 | 183 |
| 6 | 450–500 | 162 | 103 | 6 | 384–500 | 2,479 | 440–500 | 129 |

- Information considered: text, material, support, …  for 14,874 epigraphs

UNIVERSITÀ DEGLI STUDI DI BARI ALDO MORO

cini consorzio interuniversitario nazionale per l'informatica

# A concrete example: EDB

$-epigraph(E), \neg opisthographic(E), material(E, M), support(M, 'Tabula Marmorea'),$
$engraving\_technique(E, T), name(T, 'Insculptus')$
$[\textbf{200} - \textbf{299}] : [0.33..0.52] \nearrow [\textbf{300} - \textbf{349}] : 0.70 \quad \textbf{GR} = 1.65$

A moderate increase of single-sided epigraphs engraved with the insculptus technique on tabula marmorea in the time interval [300–349] with respect to the interval [200–299]

# A concrete example: EDB

$—epigraph(E), \neg opisthographic(E), material(E, M), support(M, \text{'Tabula Marmorea'}),$
$engraving\_technique(E, T), name(T, \text{'Insculptus'})$
$[\mathbf{200 - 299}] : [0.33..0.52] \nearrow [\mathbf{300 - 349}] : 0.70 \quad \mathbf{GR} = 1.65$

This may be due to the greater tolerance for Christians under the Emperor Constantine the Great (306–337 AD) and the consequent diffusion of "official" marble epigraphs in public places.

This is a significant change, since the first Christian inscriptions were usually written on tiles and bricks.

UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO

cini
consorzio
interuniversitario
nazionale
per l'informatica

# A bibliographical reference



Discovering Novelty Patterns from the Ancient Christian Inscriptions of Rome

GIANVITO PIO, FABIO FUMAROLA, ANTONIO E. FELLE, DONATO MALERBA, and MICHELANGELO CECI, University of Bari Aldo Moro

# Conclusions

- Digital technologies are opening up a great transformation in humanities

- Will these new technologies and approaches change the nature of historical inquire?

- Will we see a gradual change, with tools or techniques increasingly being added to established practice, or are we facing a revolution?

# Conclusions

- How historical knowledge is validated in the XXI century? In the era of Big Data, are data-driven approaches appropriate for the purpose of historical knowledge validation?

- When is history manipulated? How can epigraphic resources be used to disclose these manipulations? Are data mining methods useful to reveal these manipulations?

# Conclusions

- Is it possible to algorithmically discover in the data (inscriptions) the various, equally plausible, storytellings of the past?

- What large outstanding questions can epigraphists hope to address by the convergence of data science and epigraphy?

- Should we expect that data science methods will set agendas for research in epigraphy?

# Conclusions

- There is an urgent need for a critical reflection within the epigraphic community, and, more in general, among historians, on the <span style="color:red">epistemological implications</span> of the current <span style="color:red">data revolution</span>.

- Some preliminary studies (Pio et al, 2014) have barely begun to tackle the problem, despite the rapid changes in research practices presently taking place.

# Conclusions

What's the room for the epigraphist in a data-driven world?

The role of the epigraphist in these studies is particularly important.

Their pose research questions to data scientists, they give feedback on results, thus allowing an iterative refinement of analysis algorithms and the development of a user-friendly digital tool.

# Thank you