

Responsible AI

Anna Monreale
Università di Pisa



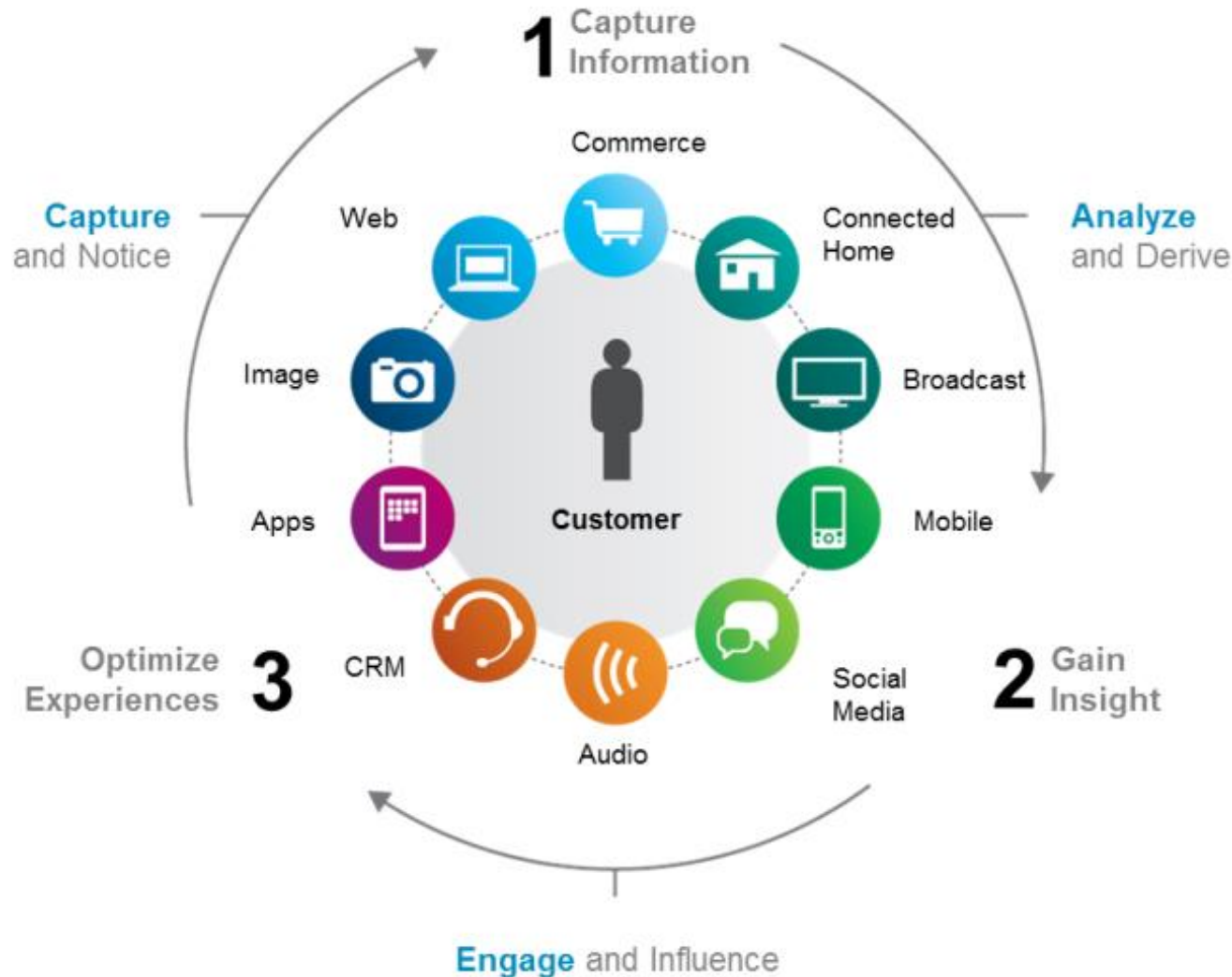
Knowledge Discovery and Delivery Lab
(ISTI-CNR & Univ. Pisa)
www-kdd.isti.cnr.it

Our digital traces

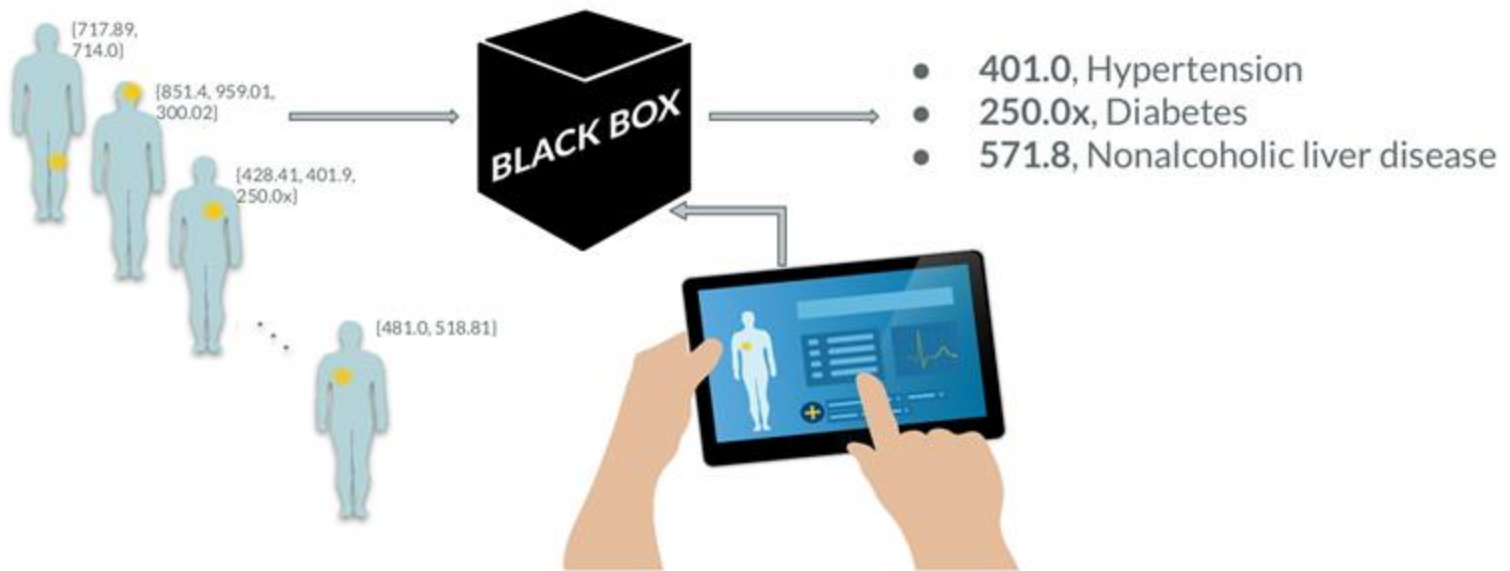
- We produce an unthinkable amount of data while running our daily activities.
- How can we manage all these data? Can we get an added value from them?



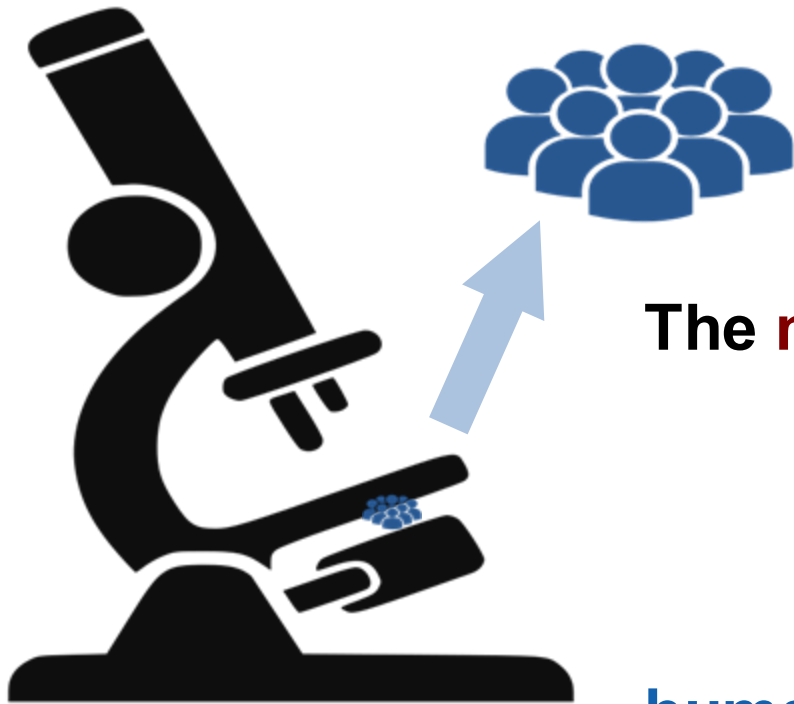
Big Data: new, more carefully targeted services



AI in healthcare



AI, Big Data Analytics & Social Mining



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

An aerial photograph of a large crowd of people scattered across a green field. The people are small, colorful figures from this high angle. A white rectangular box with a thin black border is centered horizontally and contains the text.

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

EU Ethics Guidelines for AI – (2019)

Human-centric approach: AI as a means, not an end

Trustworthy AI as our foundational ambition, with three components

Lawful AI

complying with all applicable laws and regulations

Ethical AI

ensuring adherence to ethical principles and values

Robust AI

perform in a **safe, secure** and **reliable** manner, both from technical and a social perspective, with safeguards to foresee and prevent unintentional harm

Requirements

1. Human agency and oversight

- Fundamental rights
- Human agency
- Human oversight

2. Technical robustness

- Resilience to attack and security
- Safety
- Accuracy
- Reliability and reproducibility

3. Privacy and data governance

- Privacy and data protection
- Quality and integrity of data
- Access to data

4. Transparency

- Traceability
- Explainability



Requirements

5. Diversity, non-discrimination and fairness

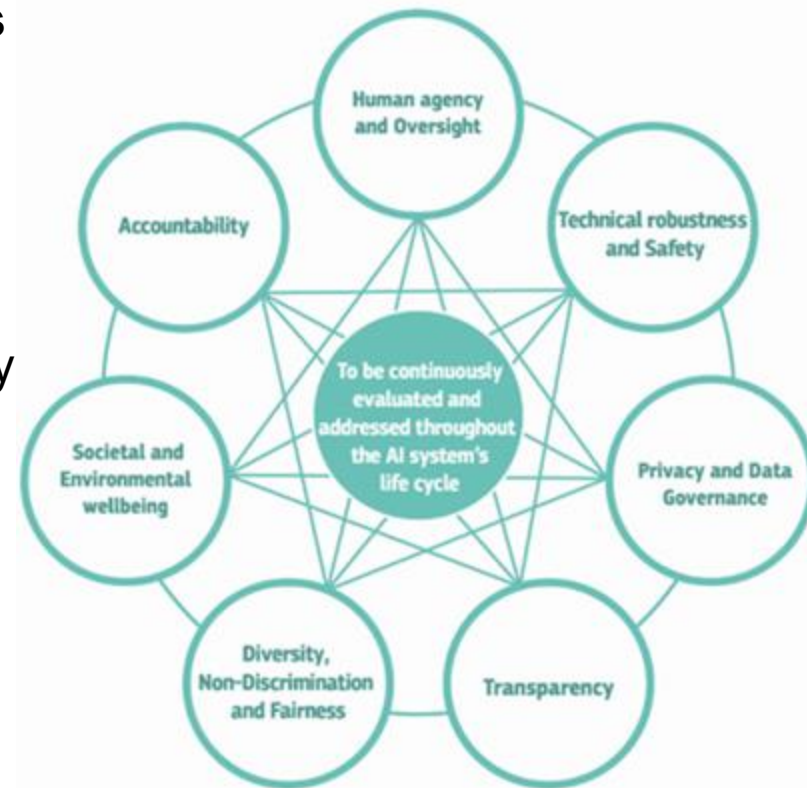
- Avoidance of unfair bias
- Accessibility and universal design
- Stakeholder Participation

6. Societal and environmental well-being

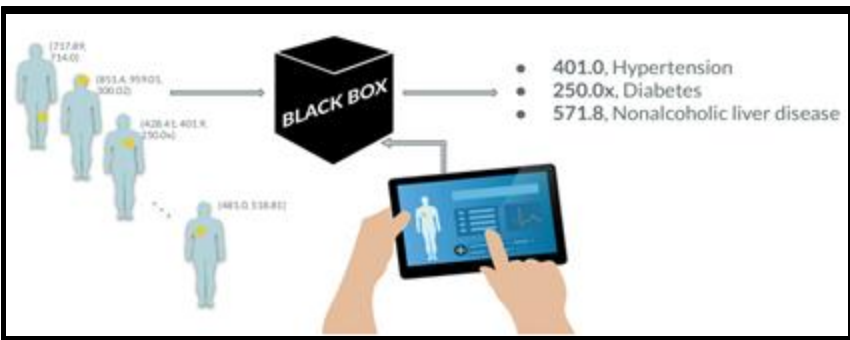
- Sustainable and environmentally friendly AI
- Social impact
- Society and Democracy

7. Accountability

- Minimisation and reporting of negative impacts
- Auditability
- Minimisation and reporting of negative impacts
- Trade-offs



Big Data & AI





WHAT IS A.I.?

A practical definition of AI

‘artificial intelligence system’ (AI system)

means a system that

1. receives machine and/or **human-based data** and inputs
2. infers how to achieve a given set of human-defined objectives using learning, reasoning or modelling implemented with the techniques and approaches listed in **Annex I**
3. generates outputs in the form of content (generative AI systems), predictions, recommendations or decisions, which influence the environments it interacts with.

Machine Learning

Deep Learning

Other statistical approaches

AI Act, TITLE I, Article 3





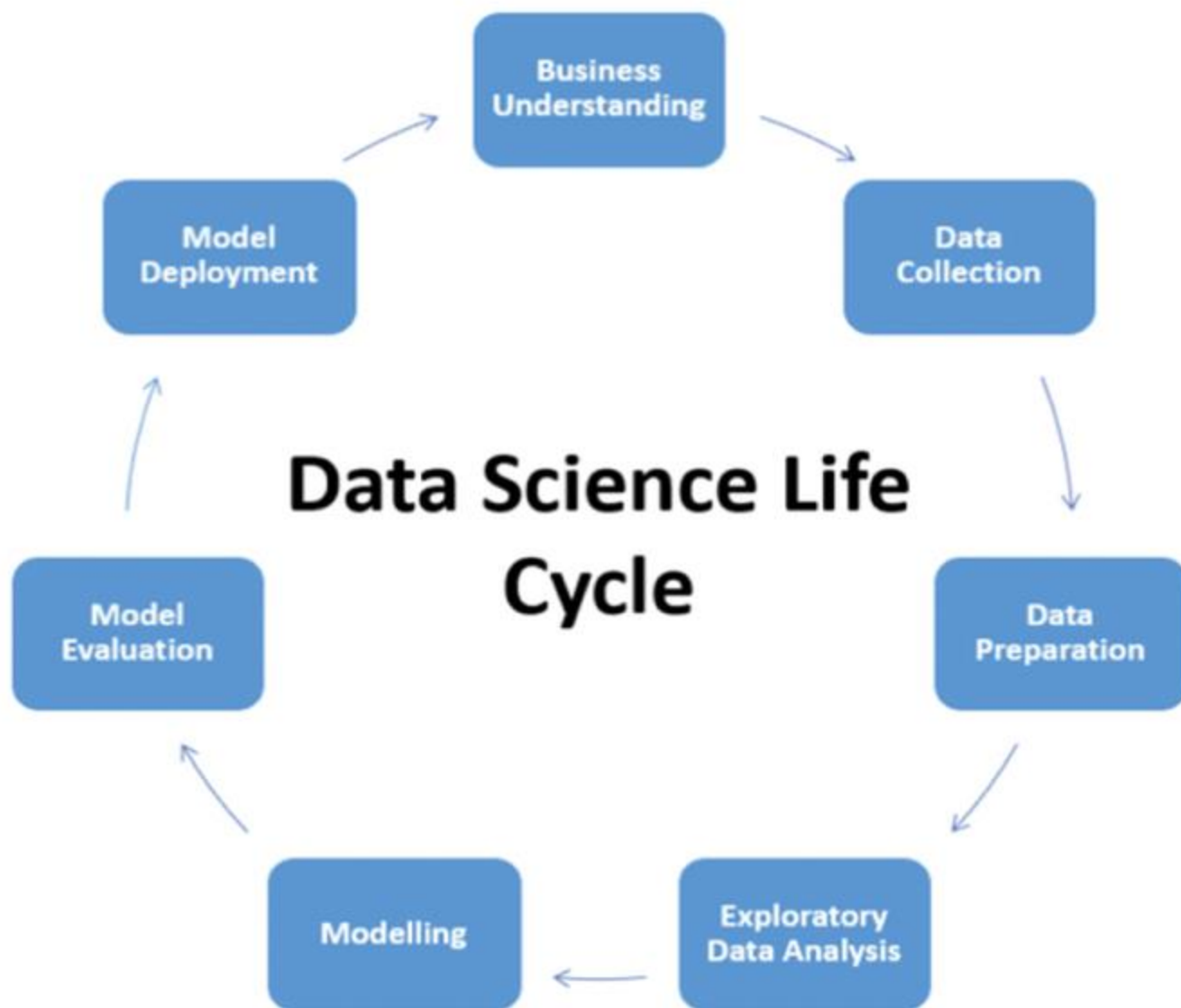


Balancing AI Innovation and Human rights

GDPR



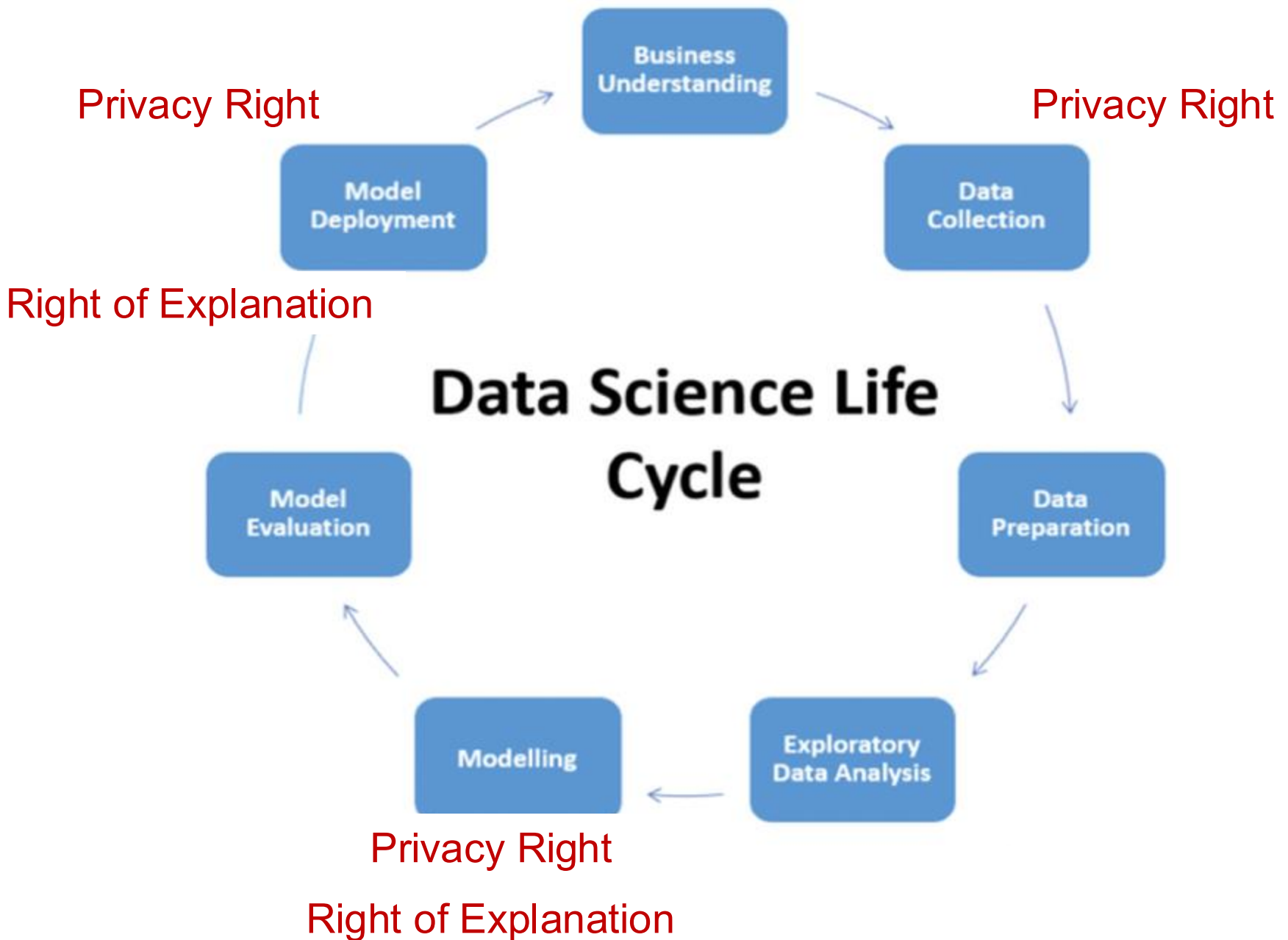
AI Act



An aerial photograph of a large crowd of people scattered across a green field. The people are seen from above, appearing as small figures in various colors. The crowd is distributed across the entire frame, with some denser clusters and some more sparse areas. Two white rectangular boxes with red text are overlaid on the image, one on the left and one on the right.

Privacy Right

Right of Explanation



PRIVACY & DATA PROTECTION

EU Legislation for protection of personal data

- European directives:
 - Data protection directive (95/46/EC)
 - ePrivacy directive (2002/58/EC) and its revision (2009/136/EC)
 - General Data Protection Regulation (May 2018)

<http://eur-lex.europa.eu/legal-content/EN/TXT/HTML/?uri=CELEX:32016R0679&from=IT>



General Data Protection Regulation

Data Protection & Right of Explanation

Personal Data



Personal data is defined as **any information** relating to an identity or identifiable natural person.



An **identifiable person** is one who can be identified, **directly or indirectly**, in particular by reference to an identification number or to one or more factors specific to his physical, physiological, mental, economic, cultural or social identity.

TYPES OF PERSONAL DATA IN GDPR

PERSONAL DATA

Names, emails, phone numbers



Display pictures, social media IDs and profile URLs

Website logs like IP addresses, user agents and device IDs



Cookies and Radiofrequency identification tag (RFID) tags



Audio and video recordings of users

Payment details like bank account number and credit card information



Geolocation data

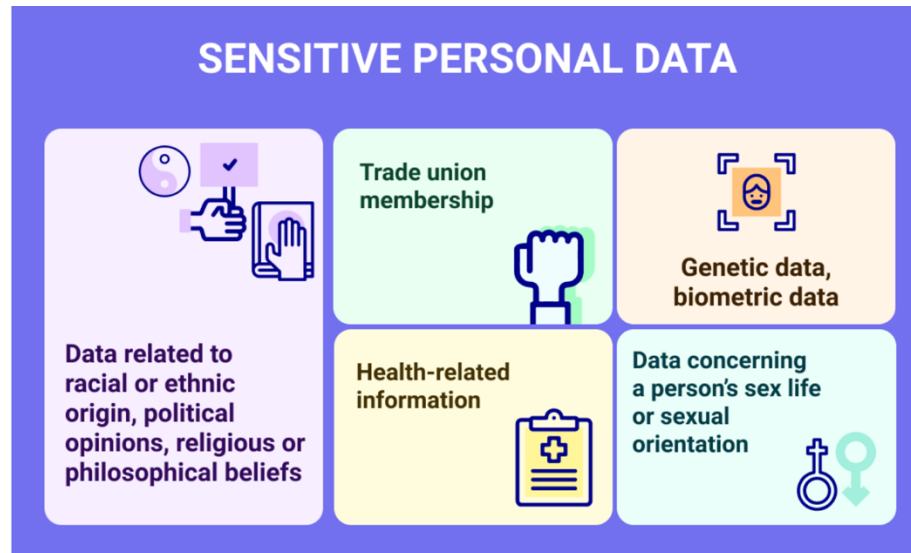


Current or previous employee data

Physical, physiological, genetic, mental, economic, cultural, or social identifiers



Special Categories of Personal Data



Sensitive personal data is a specific set of “**special categories**” that must be treated with extra security

EU Directive (95/46/EC) and GDPR

- **GOALS:**
 - protection protection of individuals with regard to the **processing** of personal data
 - the free movement of such data
 - User control on personal data
- The term “process” covers anything that is done to or with personal data:
 - collecting
 - recording
 - organizing, structuring, storing
 - adapting, altering, retrieving, consulting, using
 - disclosing by transmission, disseminating or making available, aligning or combining, restricting, erasing, or destroying data.

Anonymity according to EU Law



- The principles of protection must apply to any information concerning an identified or identifiable person
- To determine whether a person is identifiable, account should be taken of all the means likely reasonably to be used either by the controller or by any other person to identify the said person
- The principles of protection shall not apply to data rendered anonymous in such a way that the data subject is no longer identifiable

Privacy by Design Principle



Privacy by design is an approach **to protect privacy by inscribing it into the design specifications of information technologies**, accountable business practices, and networked infrastructures, from the very start



Developed by Ontario's Information and Privacy Commissioner, Dr. Ann Cavoukian, in the 1990s



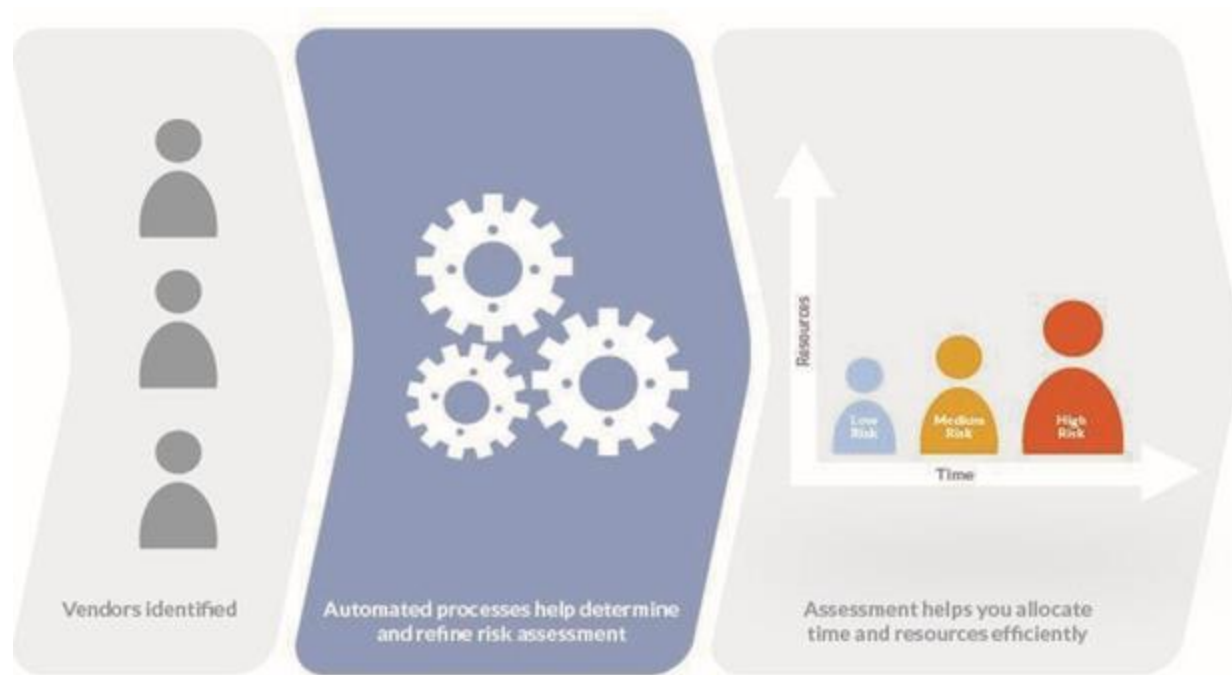
Privacy Risk Assessment



Privacy by Design

Privacy Risk Assessment

- GDPR requires that data controllers maintain an updated report on the **privacy risk assessment** on personal data collected



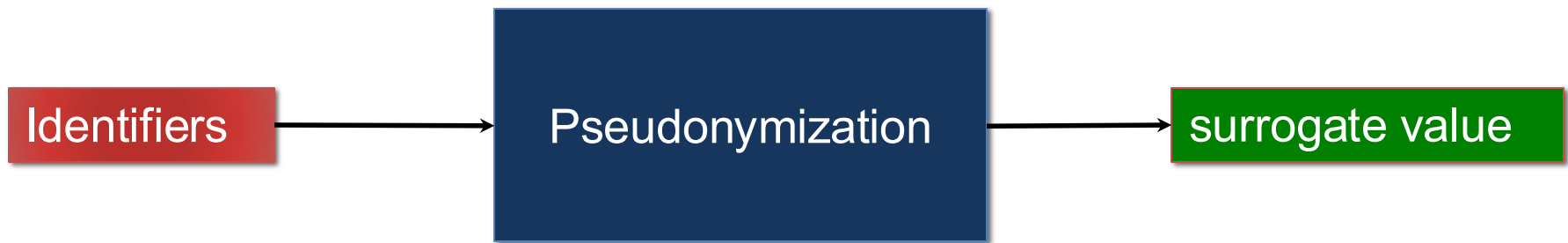
PSEUDONYMIZATION & ANONYMIZATION

Anonymization vs Pseudonymization

- Pseudonymization and Anonymization are two distinct terms often confused
- Anonymized data and pseudonymized data fall under very different categories in the regulation
- **Anonymization guarantees data protection** against the (direct and indirect) data subject re-identification
- **Pseudonymization substitutes the identity** of the data subject in such a way that additional information is required to re-identify the data subject

Pseudonymization

Substitute an **identifier** with a surrogate value called **token**



Substitute **unique names**, **fiscal code** or any attribute that identifies uniquely individuals in the data

Example of Pseudonymization

Name	Gender	DoB	ZIP Code	Diagnosis
Anna Verdi	F	1962	300122	Cancer
Luisa Rossi	F	1960	300133	Gastritis
Giorgio Giallo	M	1950	300111	Heart Attack
Luca Nero	M	1955	300112	Headache
Elisa Bianchi	F	1965	300200	Dislocation
Enrico Rosa	M	1953	300115	Fracture



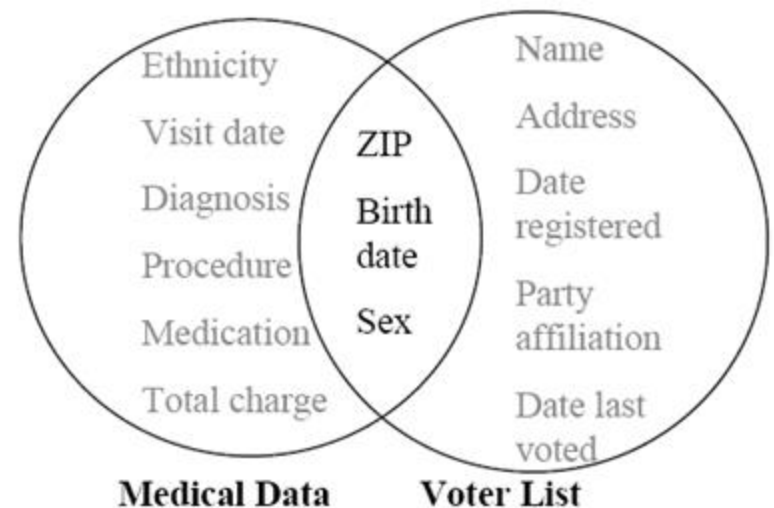
ID	Gender	DoB	ZIP CODE	DIAGNOSIS
11779	F	1962	300122	Cancer
12121	F	1960	300133	Gastritis
21177	M	1950	300111	Heart Attack
41898	M	1955	300112	Headache
56789	F	1965	300200	Dislocation
65656	M	1953	300115	Fracture

**Is Pseudonymization enough for
data protection?**

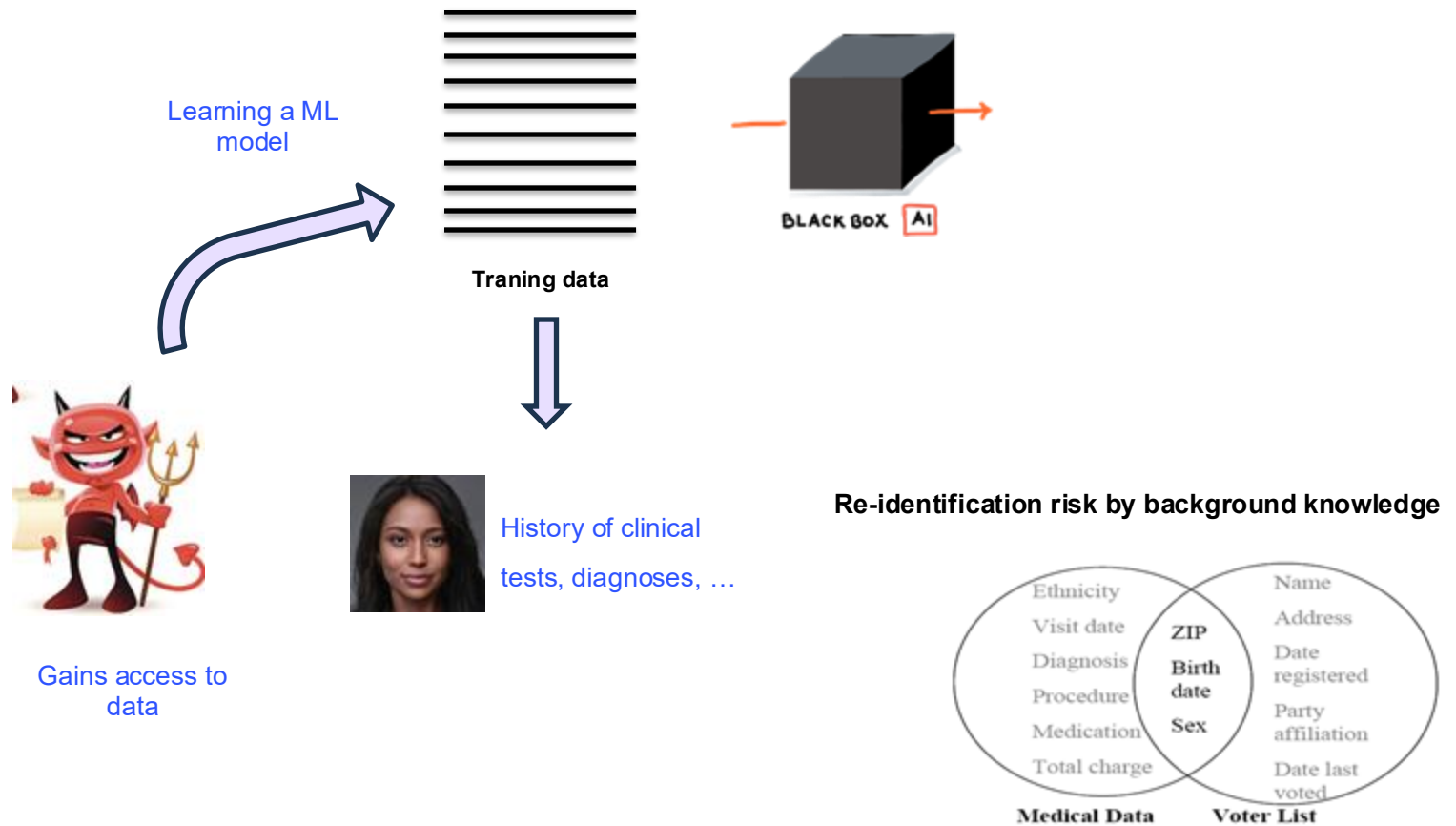
**Pseudonymized data are still
Personal Data!!**

Massachusetts' Governor

- Sweeney managed to re-identify the medical record of the governor of Massachusetts
 - MA collects and publishes sanitized medical data for state employees (microdata) **left circle**
 - voter registration list of MA (publicly available data) **right circle**
- looking for governor's record
- join the tables:
 - **6 people had his birth date**
 - **3 were men**
 - **1 in his zipcode**



Privacy Risk Assessment



Linking Attack

Governor: Birth Date = **1950**, ZIP = **300111**

ID	Gender	YoB	ZIP	DIAGNOSIS
1	F	1962	300122	Cancer
2	F	1960	300133	Gastritis
3	M	1950	300111	Heart Attack
4	M	1955	300112	Headache
5	F	1965	300200	Dislocation
6	M	1953	300115	Fracture

Which is the disease of the Governor?

Making data anonymous

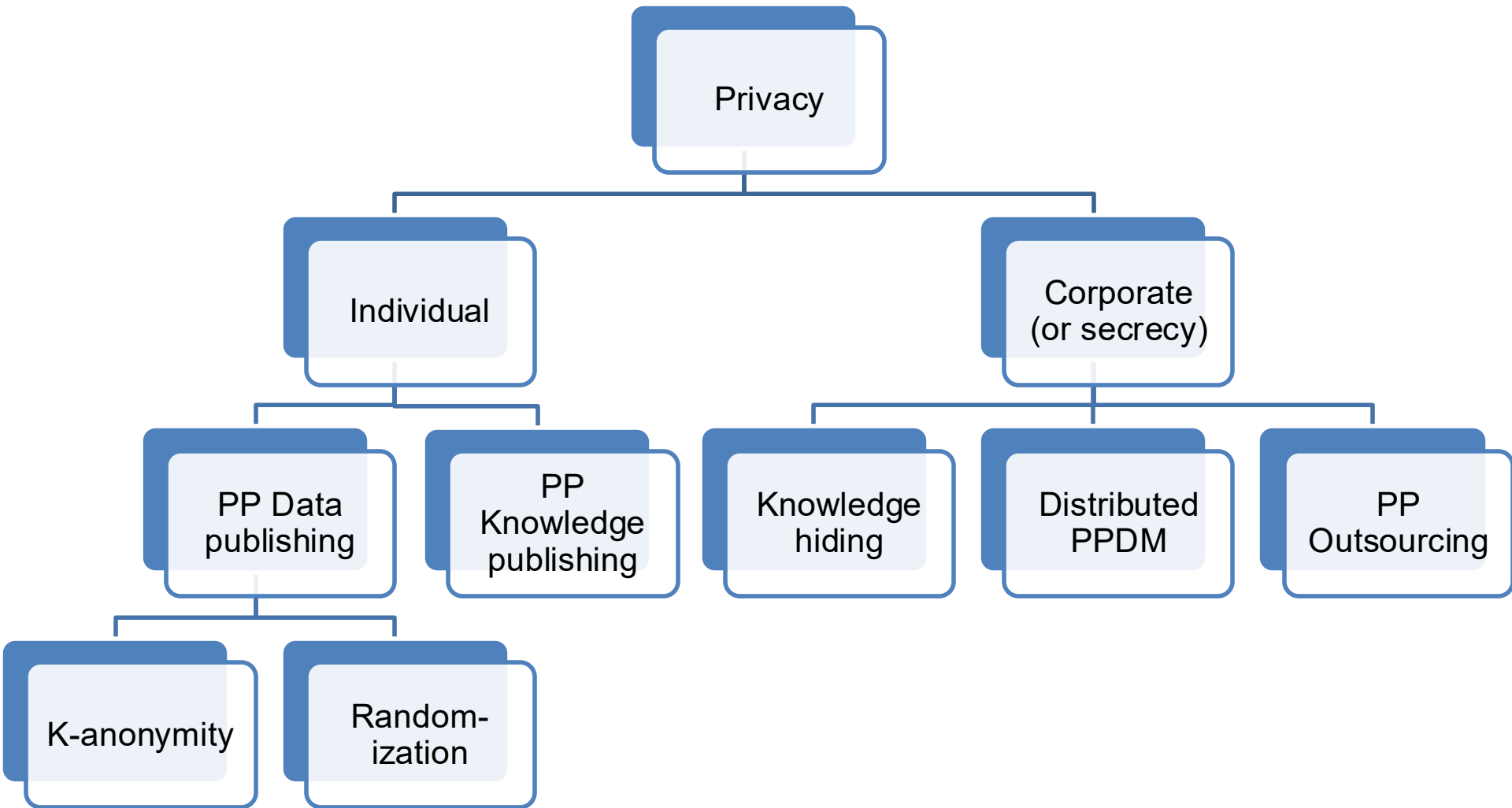
K-anonymity

Governor: Birth Date = 1950, ZIP = 300111

ID	Gender	YoB	ZIP	DIAGNOSIS
1	F	[1960-1956]	300***	Cancer
2	F	[1960-1956]	300***	Gastritis
3	M	[1950-1955]	30011*	Heart Attack
4	M	[1950-1955]	30011*	Headache
5	F	[1960-1956]	300***	Dislocation
6	M	[1950-1955]	30011*	Fracture

Which is the disease of the Governor?

Ontology of Privacy in Data Mining & AI



Attribute classification

Identifiers

Quasi-identifiers

Sensitive

ID	Gender	YoB	ZIP	DIAGNOSIS
1	F	1962	300122	Cancer
2	F	1960	300133	Gastritis
3	M	1950	300111	Heart Attack
4	M	1955	300112	Headache
5	F	1965	300200	Dislocation
6	M	1953	300115	Fracture

K-Anonymity

- **k-anonymity** hides each individual among $k-1$ others
 - each QI set should appear at least k times in the released data
 - linking cannot be performed with confidence $> 1/k$
- How to achieve this?
 - **Generalization**: publish more general values, i.e., given a domain hierarchy, roll-up
 - **Suppression**: remove tuples, i.e., do not publish outliers. Often the number of suppressed tuples is bounded
- Privacy vs utility tradeoff
 - do not anonymize more than necessary
 - Minimize the distortion

Vulnerability of K-anonymity

ID	Gender	DoB	ZIP	DIAGNOSIS
1	F	1962	300122	Cancer
2	F	1960	300133	Gastritis
3	M	1950	300111	Heart Attack
4	M	1950	300111	Heart Attack
5	M	1950	300111	Heart Attack
6	M	1953	300115	Fracture

/-Diversity

- Principle
 - Each equivalence class has at least k well-represented sensitive values
- Distinct k -diversity
 - Each equivalence class has at least k distinct sensitive values

ID	Gender	DoB	ZIP	DIAGNOSIS
1	F	1962	300122	Heart Attack
2	F	1960	300133	Headache
3	M	1950	300111	Dislocation
4	M	1950	300111	Fracture
5	M	1950	300111	Heart Attack
6	M	1953	300115	Headache

K-Anonymity

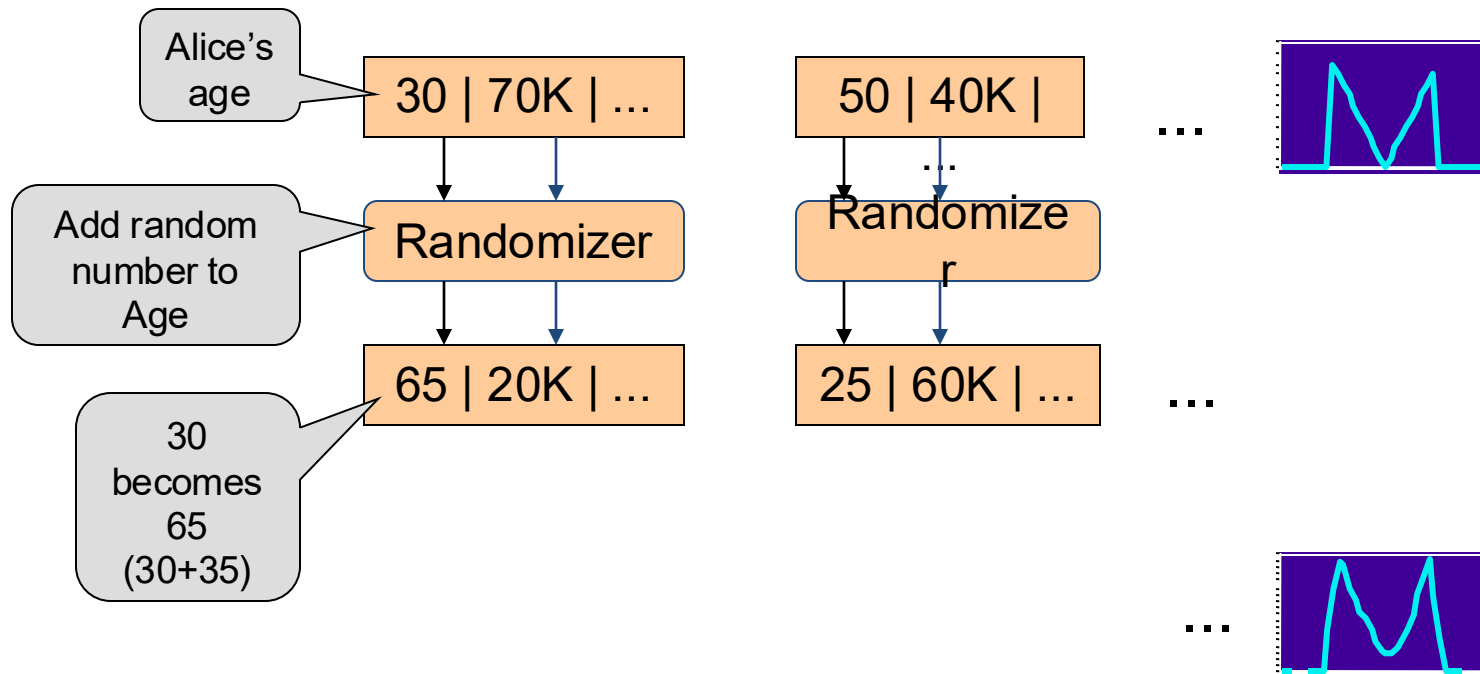
- Samarati, Pierangela, and Latanya Sweeney. “Generalizing data to provide anonymity when disclosing information (abstract).”
In PODS '98.
- Latanya Sweeney: *k-Anonymity: A Model for Protecting Privacy*. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 10(5): 557-570 (2002)
- Machanavajjhala, Ashwin, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkatasubramanian. “*l-diversity: Privacy beyond k-anonymity*.” *ACM Trans. Knowl. Discov. Data* 1, no. 1 (March 2007): 24.
- Li, Ninghui, Tiancheng Li, and S. Venkatasubramanian. “*t-Closeness: Privacy Beyond k-Anonymity and l-Diversity*.” *ICDE 2007*.

Randomization

- **Original values x_1, x_2, \dots, x_n**
 - from probability distribution X (unknown)
- **To hide these values, we use y_1, y_2, \dots, y_n**
 - from probability distribution Y
 - Uniform distribution between $[-\alpha, \alpha]$
 - Gaussian, normal distribution with $\mu = 0, \sigma$
- **Given**
 - $x_1 + y_1, x_2 + y_2, \dots, x_n + y_n$
 - the probability distribution of Y

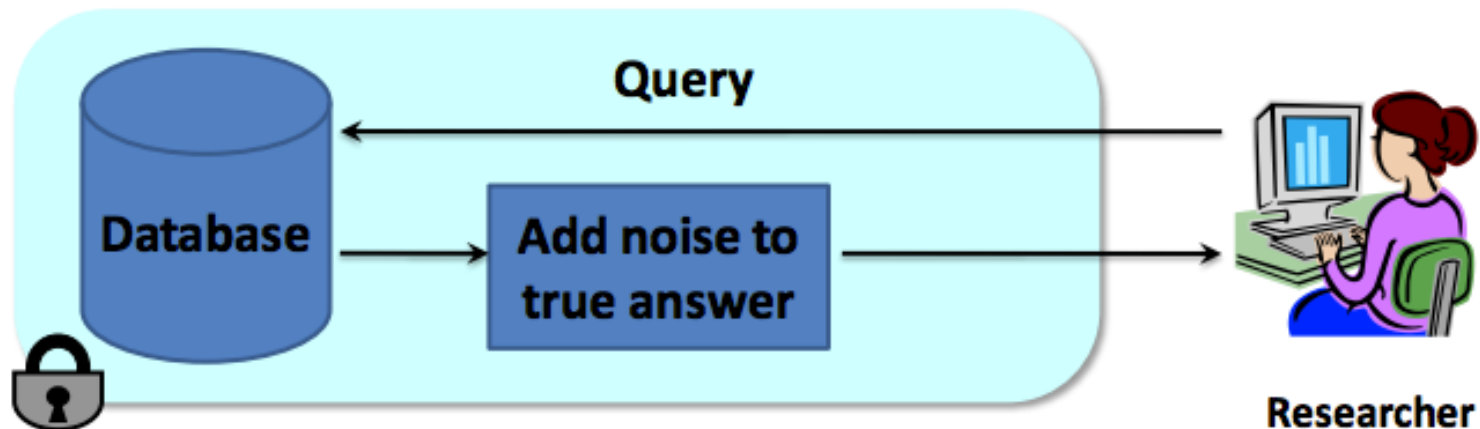
Estimate the probability distribution of X .

Randomization Approach Overview



Differential Privacy

- The risk to my privacy should not increase as a result of participating in a statistical database



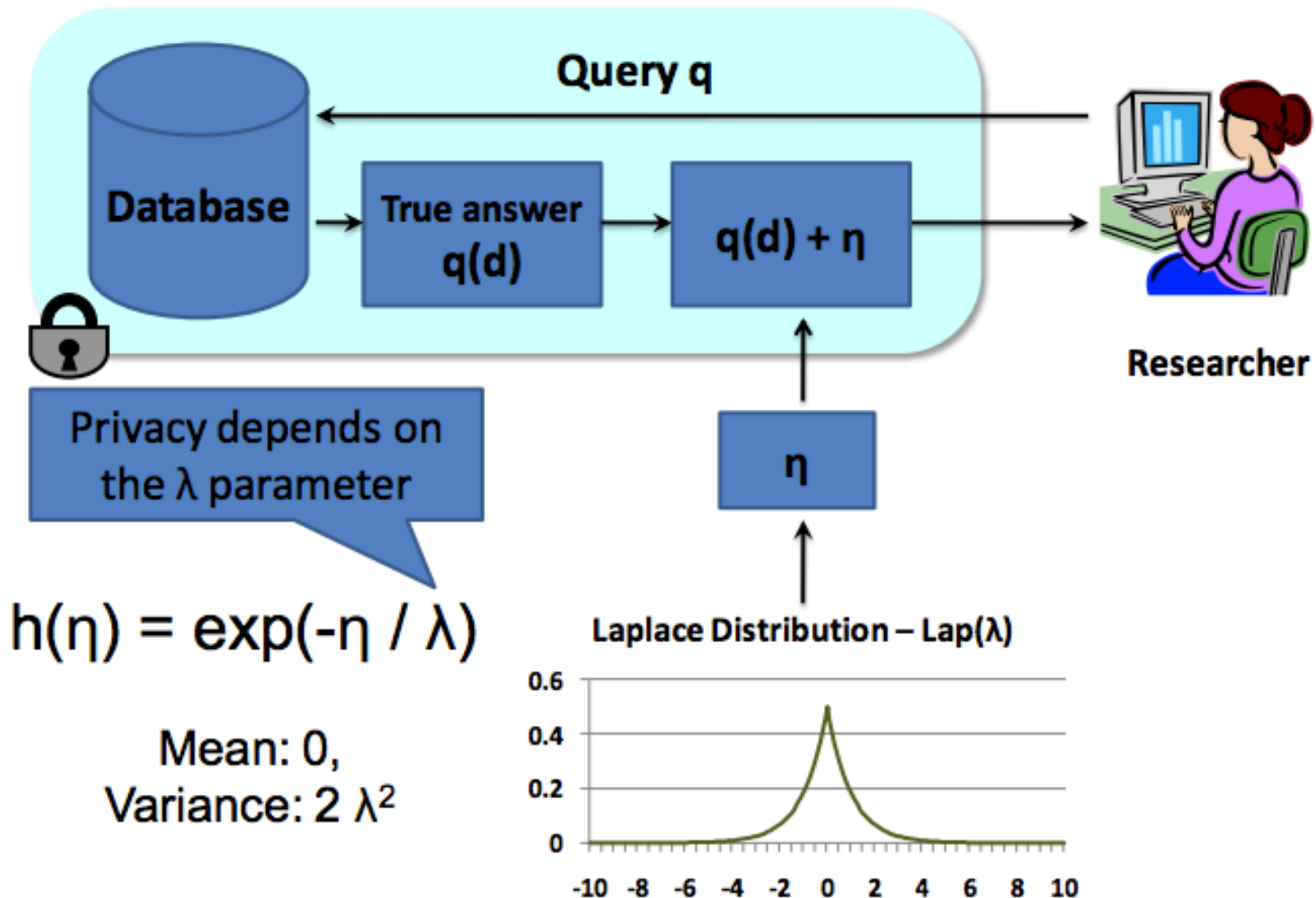
- Add noise to answers such that:
 - Each answer does not leak too much information about the database
 - Noisy answers are close to the original answers

Attack

Name	Has Diabetes
Alice	yes
Bob	no
Mark	yes
John	yes
Sally	no
Jack	yes

- 1) how many persons have Diabetes? 4
 - 2) how many persons, excluding Alice, have Diabetes? 3
- **So the attacker can infer that Alice has Diabetes.**
 - **Solution:** make the two answers similar
 - 1) the answer of the first query could be $4+1 = 5$
 - 2) the answer of the second query could be $3+2.5=5.5$

Differential Privacy



Randomization

- R. Agrawal and R. Srikant. [Privacy-preserving data mining](#). In Proceedings of SIGMOD 2000.
- D. Agrawal and C. C. Aggarwal. [On the design and quantification of privacy preserving data mining algorithms](#). In Proceedings of PODS, 2001.
- W. Du and Z. Zhan. [Using randomized response techniques for privacy-preserving data mining](#). In Proceedings of SIGKDD 2003.
- A. Evfimievski, J. Gehrke, and R. Srikant. [Limiting privacy breaches in privacy preserving data mining](#). In Proceedings of PODS 2003.
- A. Evfimievski, R. Srikant, R. Agrawal, and J. Gehrke. [Privacy preserving mining of association rules](#). In Proceedings of SIGKDD 2002.
- K. Liu, H. Kargupta, and J. Ryan. [Random Projection-based Multiplicative Perturbation for Privacy Preserving Distributed Data Mining](#). IEEE Transactions on Knowledge and Data Engineering (TKDE), VOL. 18, NO. 1.
- K. Liu, C. Giannella and H. Kargupta. [An Attacker's View of Distance Preserving Maps for Privacy Preserving Data Mining](#). In Proceedings of PKDD'06

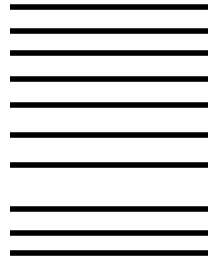
Differential Privacy

- Cynthia Dwork: [Differential Privacy](#). ICALP (2) 2006: 1-12
- Cynthia Dwork: [The Promise of Differential Privacy: A Tutorial on Algorithmic Techniques](#). FOCS 2011: 1-2
- Cynthia Dwork: [Differential Privacy in New Settings](#). SODA 2010: 174-183

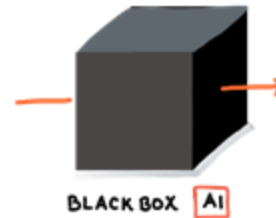
**Can we jeopardize individual
privacy without accessing data?**

Privacy risk of ML models

LEARNING A
ML MODEL



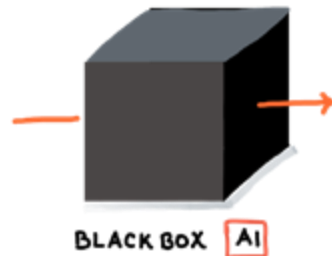
Training data



Infer she belongs to
confidential training
data



Query the BB
model



Get an answer



APPLY A ML
MODEL

The privacy attack: MIA

