

8-12th November 2021 – Politecnico Di Milano, Italy

Privacy-Preserving Data Processing

Luca Vassio, Martino Trevisan Performance 2021 November, 8

entsbroadband characterizationlongitudinal traffic analysisservice Complex NetworksDynamic Assistant Professor at Politecnico di Torino software puter defined Network measurements Mobile networksNetwork emulationData drivenRealistic achine User ModelingG Passive Real Time CommunicationsRTPClassification HTTP/3 Open Online Social NetworksInstagram Monitoring T CP networking Mobile networks Measurement Dataset Internet proxy encryption enerative Models ment NFVCloud **Online social networks** Access Line Capacity ea Privacy ന Iser Internet measurements Measureme data Classification Behavior Web Tracking Politics Performance Protocols applications Adult Content **QoE Metrics** (\mathbf{O}) **Traffic Management** Isage ornography Traffic monitoring Web Services online tracking Real Time Control Applications SIL Passive Measurements: Web Pornography

The Speakers

Martino Trevisan







The Speakers





Luca Vassio Assistant Professor at Politecnico di Torino

Characterizing Two way One way Free floating Recommendation systems Autoencoders Semi supervised I Car sharing Electric cars Smart city Shared economy Data driven analysis Clustering Passive Complex NetworksDynamic User ModelingGenerative Models Web pages SmartData Anomalies Social networkBig DataCOVID 19FacebookInstagram Uniform resource locators Data driven optimisation Measurements learning smart cities Network measurements Online social networks **User Behavior** Big user behaviour data Web po Modified Car Sharing Charging floating Politics passive traces sustenability Network Politics transportation Artificial human station Vehicle Trackers Network dynamics Data models Bee Colony behaviour Monitorin sage .ea placement B History Urban mobility passing bui Domain names alarm logs Influence data driven system 9 **Reinforcement learning** anomaly detection and prediction

3



- Understand why Privacy is an issue when processing data regarding people
- Have an overview of the most popular techniques for data anonymization
- Hands on: anonymize a dataset using the state-of-the-art libraries and techniques

Outline



- 1. The rise of Data-Driven approaches and problematics
- 2. Why we need Anonymization and why it is difficult
- 3. Privacy-Preserving Techniques
 - K-anonymity (and variants)
 - Differential Privacy
- 4. Open-Source tools for anonymization
- 5. Hands-on on Data



The goal of this session is to offer to the PhD students and early career Post-docs attending the Performance 2021 conference (even without presenting a paper) the opportunity to meet each other and get feedback on their work.

In this session PhD students and Post-docs can meet experienced researchers and get comments on their work.

The session is scheduled on Friday November 12 and will last about 1-1.5 hour starting at 11.00 NYC time. https://www.performance2021.deib.polimi.it/meet-the-star-event/



The rise of Data-Driven approaches and problematics

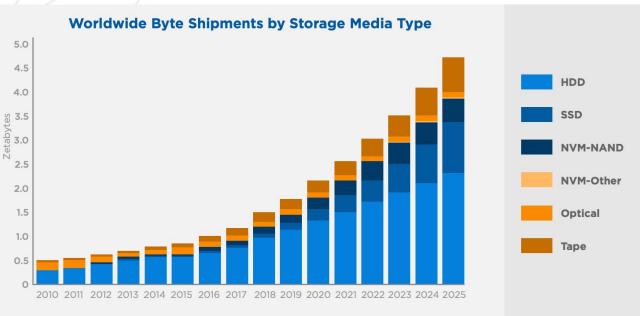
An Unprecedented Amount of Data

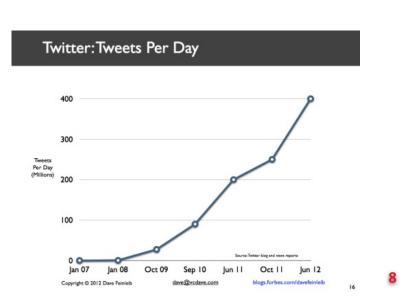




Continuous Growth

- Data volume increases exponentially over time •
- 44x increase from 2009 to 2020
 - Digital data 35 ZB in 2020





@StatistaCharts Source: Instagram

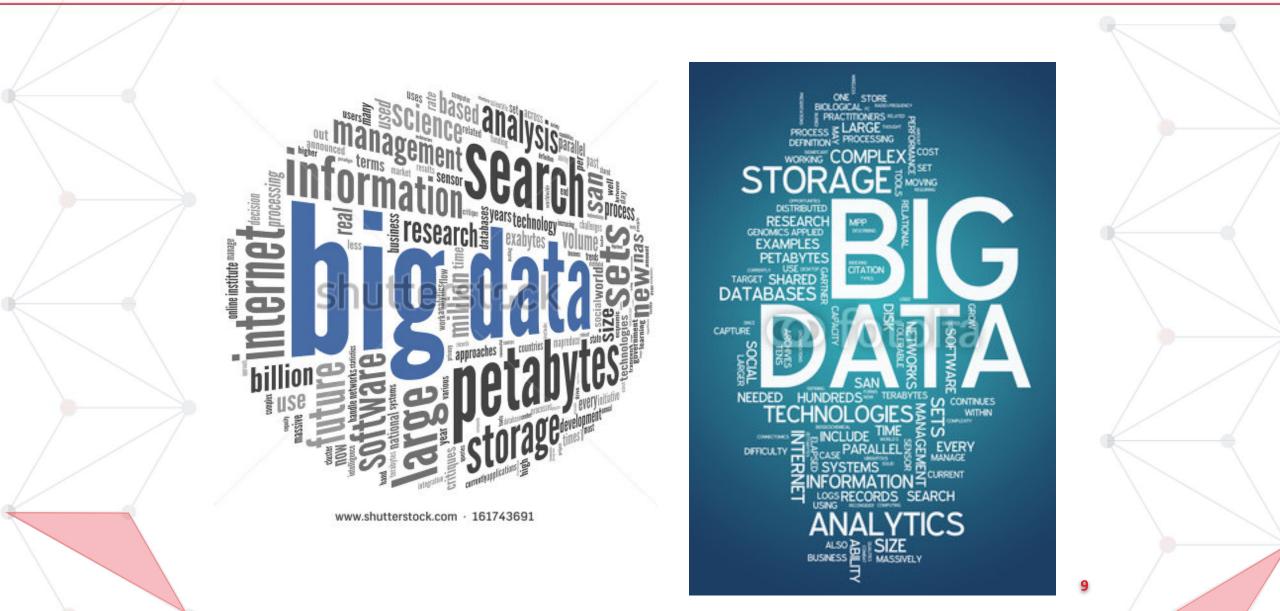
statista 🗹

2017 2018

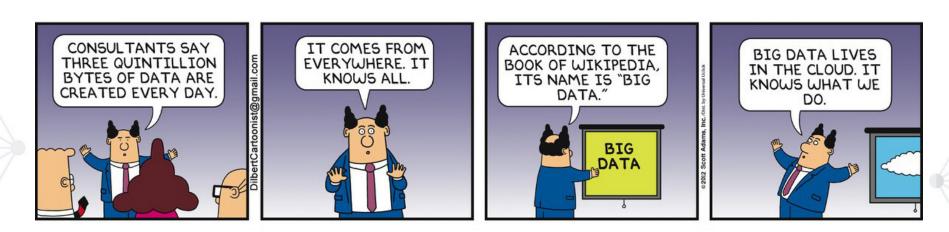
2016

What are Big data





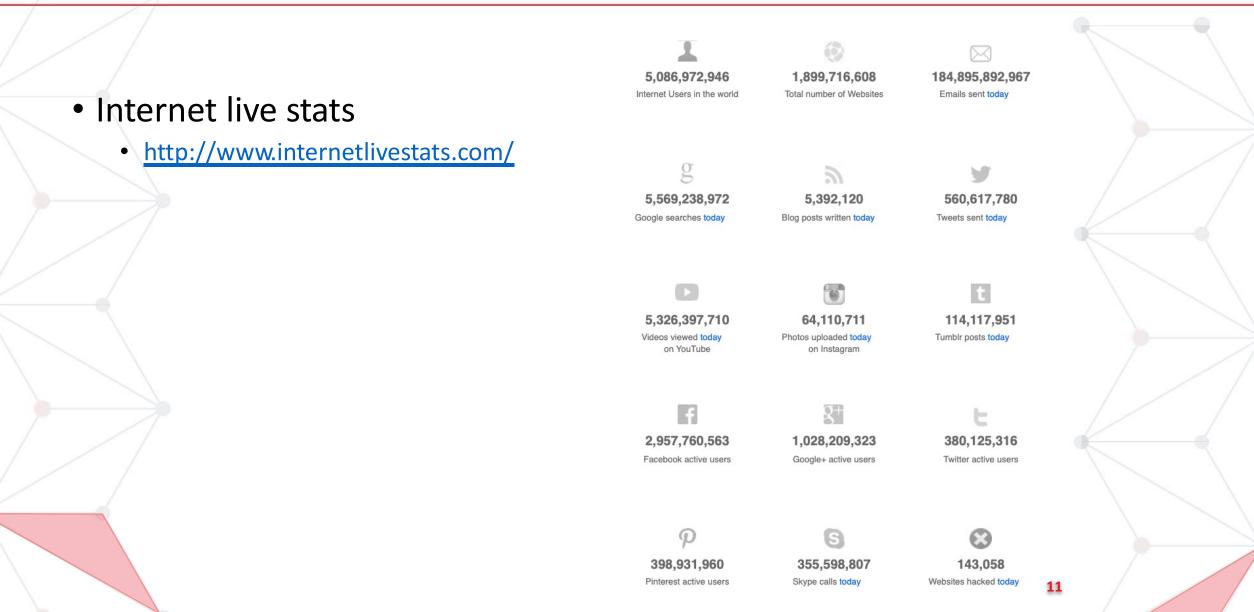




- Many different definitions
 - "Data whose scale, diversity and complexity require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"

Data on the Internet...





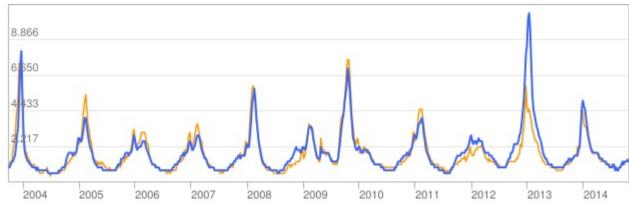
Google Flu trends





• February 2010

- Google detected flu outbreak two weeks ahead of CDC data (Centers for **Disease Control and** Prevention – U.S.A)
- Based on the analysis of Google search queries



12

The potential of data-driven economy



- According to McKinsey: «The age of analytics: Competing in a datadriven world»
- Even under the most conservative estimations, the figures for the potential of data-driven decision making are staggering
- Mobility alone could see benefits of close to **2.5 trillion** dollars by 2025, whereas multiple other sectors have annual benefits in the hundreds of billions (e.g., banking with **260 billion** per year).

The potential of data driven economy



	Potential impact: 2011 research	Value captured %	Major barriers	
Location-bas data	 \$100 billion+ revenues for service providers Up to \$700 billion value to end users 	50- 60	 Penetration of GPS-enabled smartphones globally 	
US retail ¹	 60%+ increase in net margin 0.5–1.0% annual productivity growth 	30- 40	 Lack of analytical talent Siloed data within companies 	
Manufacturin	 g² Up to 50% lower product development cost Up to 25% lower operating cost Up to 30% gross margin increase 	20-30	 Siloed data in legacy IT systems Leadership skeptical of impact 	
EU public sector ³	 ~€250 billion value per year ~0.5% annual productivity growth 	10-20	 Lack of analytical talent Siloed data within different agencies 	
US health car	 *e \$300 billion value per year ~0.7% annual productivity growth 	10-20	 Need to demonstrate clinical utility to gain acceptance Interoperability and data sharing 	

1 Similar observations hold true for the EU retail sector.

2 Manufacturing levers divided by functional application.

3 Similar observations hold true for other high-income country governments.

SOURCE: Expert interviews; McKinsey Global Institute analysis

Who generates data?



- User Generated Content (Web & Mobile)
 - E.g., Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube

Health and scientific computing







Who generates data?

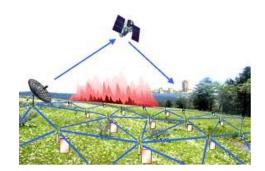


Log files

• Web server log files, machine system log files

	File Edit Format View Hulp			And a second
creat detector are an initial betting a lease to over a event that the second	vestimulum voi factitisis nella, pellentesque aper ri coren inqua dori tri succe locater rice de lacting mesenesa dapises arcs un risi elettene elementar. recein arts pertition nella succesar in aspara risis a liquam suas porter succe locater risis el de la vestimulum voi factitisis nella. Pellentesque eger ri altegam suas porter succe locater risis el del a recein agis porter succes aspara el del art mesenesa dapises arcs un risi elettene elementar. recein agis porter succes assecas in aspara risis mesenesa dapises arcs un risi elettene elementar.	lareet riss. arcs consected a brown to brown si dieferd el consected a brown to brown consected a brown to brown si dieferd el mascenas fon sit mascenas fon sit race lareet brown to brown sector brown to brown to brown to brown sector brown to brown to brown to brown to brown sector brown to brown to brown to brown to brown to brown to brown sector brown to brow	ot noll elefferd elementum. Nor nolla. Necetomis in magna risso. I ro da. Norde anno risso. International anno norde lancer risso el da losta amet, consecteur adjoiccing el ut nolla elefferd elementum. Nor nulla. Necetomis in magna risso. I risso lancera international risso. Internala. Pellettesque eget nog	nulla ente ufna, sagittis pe non augue condisentue i les et friegilla lorem pe it, Aenean at urna at ero trisus et dui sodales et it, Aenean at urna at ero turpis magna, placerat nulla ente urna, sagittis ue non augue condisentue i les et friegilla lorem pe
per tortor tincidunt eu. tibulum vel facilisis nulla. Peller [: zo - Noteped	the state of the s	the surname of the		i turpis magna, placerat milla ente urna, sagittis
		A second	The set of	and an effective constraints of the second s
Corper fortor tircle vestibulae vel facil	hrt es. isis nulla. Pellertesque eget neque non augue condisent	us posuere a ac augue. I	rraesent eu dolor orci. Mulli	laoreet risus et dui sod.

- Internet Of Things (IoT)
 - Sensor networks, RFID, smart meters



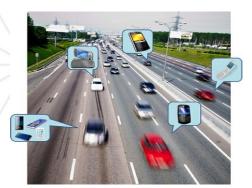




An example of Big data at work



Crowdsourcing

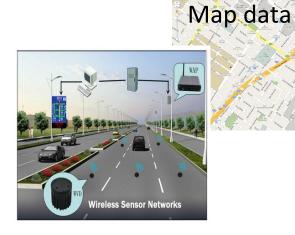




Computing



Real time traffic info



Sensing

Most of these data are generated by People This may threaten their Privacy!!!

17



Who collects data

- Data come from various sources lacksquareebay amazon **Online shops** Google • Search engines bing Navigation systems and ISPs **Google** Maps
 - Websites and Ads





Collection of information about users' online activity



- Fundamental pillar of Web industry
 - Essential to advertisement (Re-targeting, Online Behavioral Advertisement, etc.)

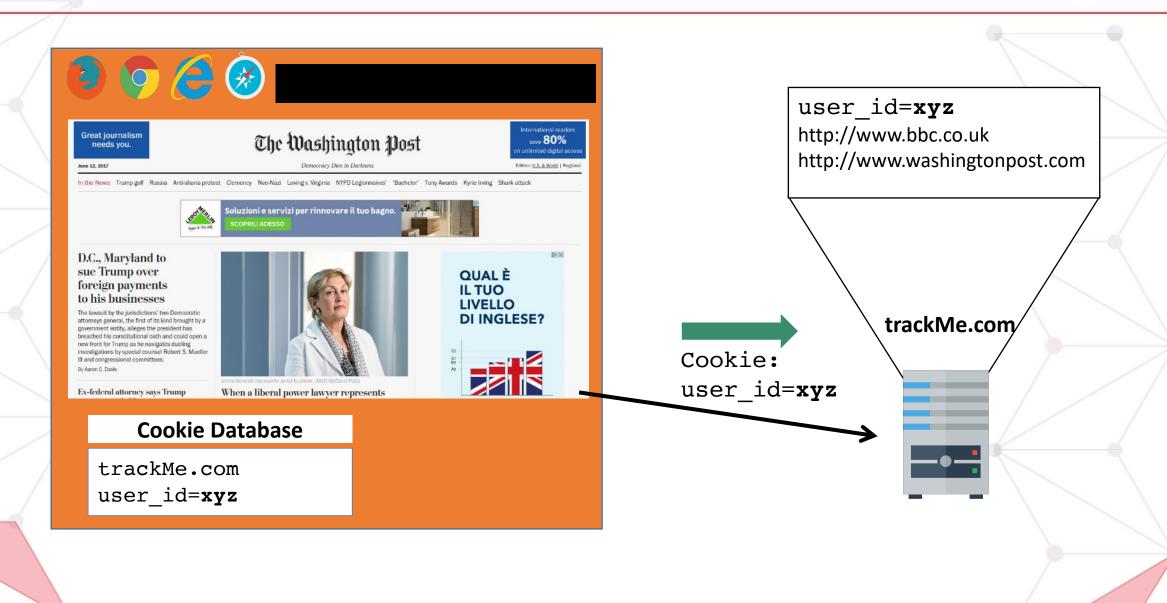
How it works





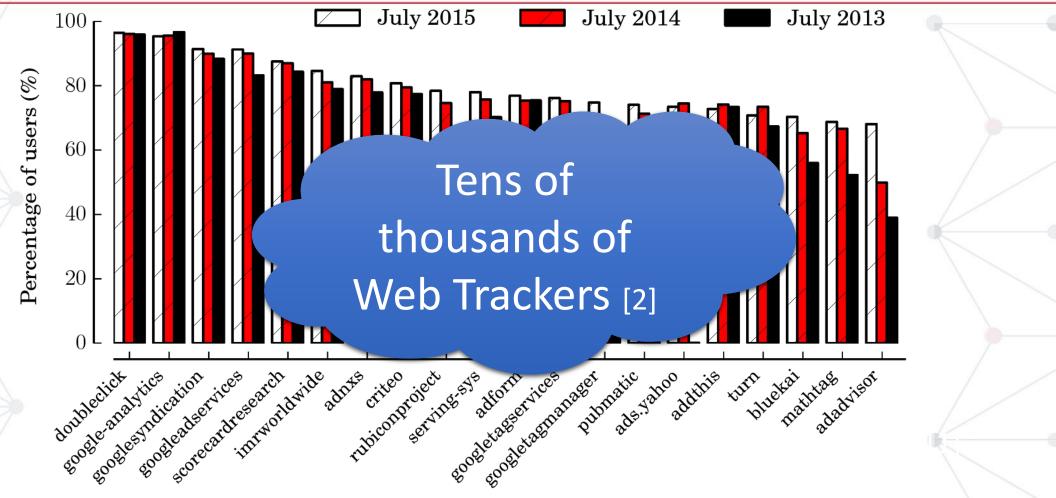
How it works





How many users are tracked?

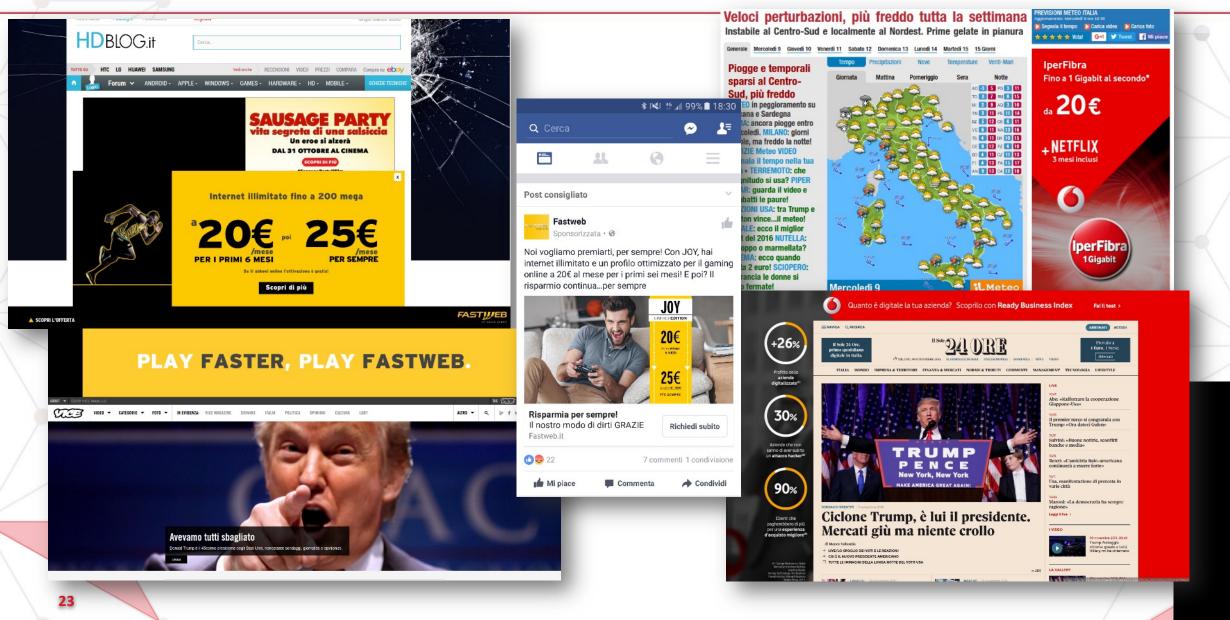




[1] Hassan Metwalley, Stefano Traverso, Marco Mellia, "Using Passive Measurements to Demystify Online Trackers", IEEE Computer "Communications and Privacy under Surveillance Issue", 2016
[2] Englehardt, Steven, and Arvind Narayanan. "Online tracking: A 1-million-site measurement and analysis", Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security

Tracking is annoying!





Tracking can be dangerous!





Kashmir Hill, Forbes Staff Welcome to The Not-So Private Parts where technology & privacy collide

Follow on Forbes Subscribe 36k

re technology & privacy collide

TECH | 2/16/2012 @ 11:02AM | 564,678 views

How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did



BUSINESS Insider

ADVERTISING

The NSA Is Using Google's Advertising Cookies To Track Its Targets



The National Security Agency is using the tracking data intended for Google's advertisers to locate its targets.

According to the Washington Post's new analysis of an internal presentation Edward Snowden leaked earlier in the year, the NSA has been using the numeric identifiers in Google's "PREF" cookies for hacking.

A cookie is a small piece of data sent from a website that is stored in a



NSA Director General Keith Alexander maintains that his agency's tactics are lawful and used against foreign adversaries.

Money International + Markets Economy Companies Tech Autos India Video Data brokers selling lists of rape victims, **AIDS** patients by Melanie Hicken @melhicken Recommend 1.6K **F S** (L) December 19, 2013: 12:38 PM ET Ikea gives employee paid pare RAPE VICTIM. champion protest POLICE OFFICERS Newly cor Yellen is a the U.S. e

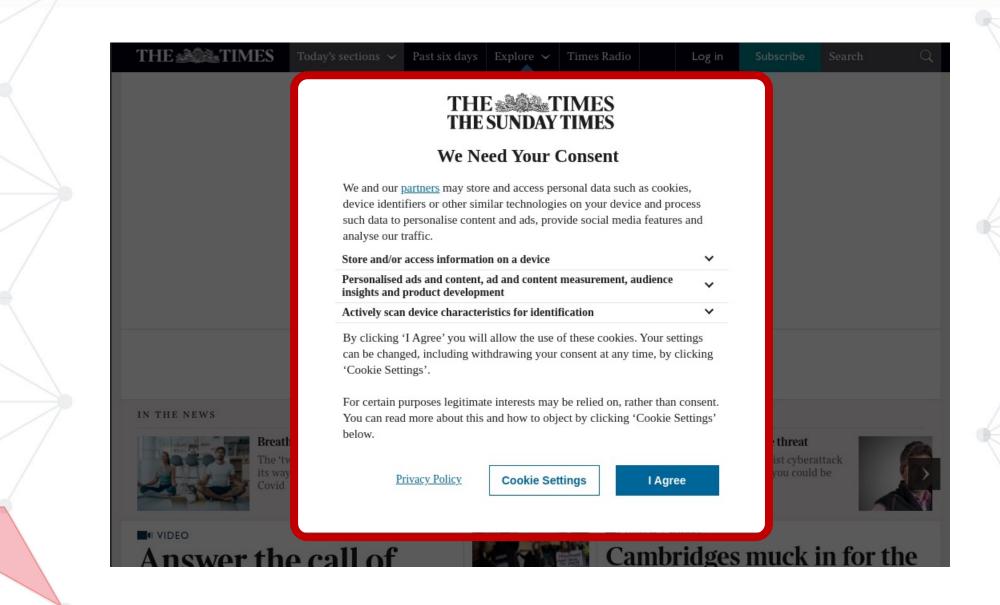
Laws on Privacy

- Legislators promulgated laws to defend users' privacy
- Cookie Law or Directive 2009/136/EC requires:
- Websites to ask "prior informed consent for storage or for access to information stored on a user's terminal equipment"
- In other words, a website must *explicitly* ask the visitor to authorize the storage and retrieval of data sent through *cookies* and similar tracking mechanisms *before* delivering and installing them.
- Its implementation has become evident to end-users because of the presence of a "Cookie Bar" on most of websites



According to EU, website should ask consent!







Regulation 2016/6791 – Entered into Force in May 2018

- Goal: protect users' privacy, and punish violations
- It applies to data of EU citizen, worldwide

What are Personal Data?

Personal data are any information which are related to an identified or identifiable natural person

- Can be identifiers (e.g., name, phone number)
- Or special characteristics (e.g., genetic, mental, cultural info)

General Data Protection Regulation



Principles:

- 1. Accountability: be responsible and adopt a safe behavior and document it
- 2. Privacy by design: assess the privacy impact before data processing
- 3. Privacy By Default: Cannot accumulate data without an objective
 - 1. Can accumulate only the data strictly needed for the goal
 - 2. Companies must always obtain consent of users
- 4. Transparency: simple, clear and complete information
 - 1. Users must be informed give the **consent freely**
 - Users should have access to data and be able to "download" it and delete it
- 5. Must inform users and authorities of **data breaches** within **72 hours**

Roles in the GDPR





A data subject is any natural person whose personal data is collected, retained or processed



Responsible for the data in their possession, such as personal data of employees, prospects / leads, customers or suppliers, among others.

DATA PROCESSOR



We refer to that service provider who must access personal data that are the responsibility of the Data Controller.

DPO

Data Protection Officer

- PIA and DPIA
 Privacy Impact Assessment (PIA)
 - Privacy Impact Assessment (PIA) is all about analyzing how an entity collects, uses, shares, and maintains personally identifiable information, related to existing risks.
 - used to protect privacy by design when an organization starts or acquires a new business, implements a new process, or launches a new product
 - Data Protection Impact Assessment (DPIA) is all about identifying and minimizing risks associated with the processing of personal data.
 - regularly applied to personal data processing, identifying, and mitigating risks.



POLITECNICO DI TORINO



Pitfalls in anonymization

Data Publishing and Lessons Learned

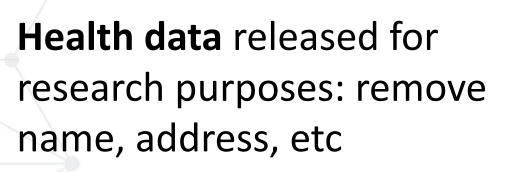


- Companies and public bodies have published data
 - For research, societal benefits challenges
- They involuntarily disclosed personal information
 - They tried to remove sensitive information
 - But this is harder than it appears

 Health data released for research purposes: remove name, address, etc

Anonymization

- Hide identity, remove identifying info
- **Proxy server / Network data:** remove source IP address





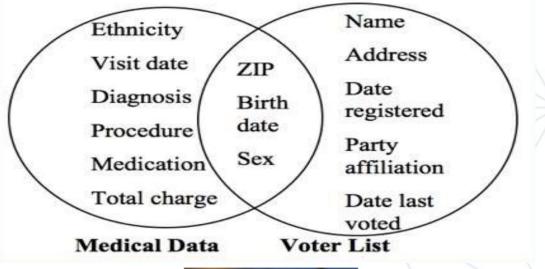




- In 1997 the Massachusetts Group Insurance
 Commission decided to release "anonymized" data on state employees that showed every single hospital visit.
- The goal was to help researchers, and the state spent time removing all obvious identifiers such as name, address, and Social Security number.

Popular cases of Data De-Anonymization

- Latanya Sweeney, a computer science researcher purchased the complete voter rolls from the city of Cambridge, a database containing, among other things, the name, address, ZIP code, birth date, and sex of every voter.
- By combining data, Sweeney found Governor Weld with ease.
 - Only six people in Cambridge shared his birth date, only three of them men, and of them, only he lived in his ZIP code.





born July 31, 1945 resident of 02138



Popular cases of Data De-Anonymization



- This demonstrates that removing IDs is not enough for privacy
- Latanya Sweeney in 2000 in "Uniqueness of Simple Demographics in the U.S. Population" stated that
 - 87% of the US population can be uniquely identified by gender, ZIP code and full date of birth.
- Philippe Golle in 2006 in *Revisiting the Uniqueness of Simple Demographics in the US Population*
 - Recompute this number to 63% of individuals
- Paul Ohm in 2009 in Broken Promises of Privacy: Responding to the Surprising Failure of Anonymization
 - "Easy reidentification represents a sea change not only in technology but in our understanding of privacy."



- In 2006, AOL researchers released a massive dataset of search queries
- They first "anonymized" the data by scrubbing user IDs and IP addresses.
- User IDs were scrubbed but were replaced with a number that uniquely identified each user.
 - Good for research 😳
 - An individual could be identified and matched to their account and search history by such information ⁽²⁾

98280 how does a male's cocaine use affect a fetus
98280 how does a male's cocaine use affect a fetus
98280 birth defects caused by father's cocaine use
98280 birth defects caused by father's cocaine use
98280 are chainletter scams ever successful

2006-04-10 1 2006-04-10 5 2006-04-10 1 2006-04-10 4 2006-04-10 0

Although the searchers were only identified by a numeric ID, some people's search results have become notable for various reasons.



Thelma Arnold

- "numb fingers"
- "60 single men"
- "dog that urinates on everything."

POLITECNICO DI TORINO

POLITECNICO DI TORINO

- Netflix Prize dataset, released 2006
 - 100,000,000 (private) ratings from 500,000 users
 - 10% of their users
 - average 200 ratings per user
- Competition to improve recommendations
 - i.e., if user X likes movies A,B,C, will also like D
- Anonymized: usernames replaced by a number



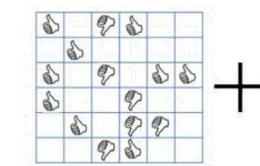
- Problem: can combine "private" ratings from Netflix with public reviews from IMDB to identify users in dataset
- May expose private info about members...



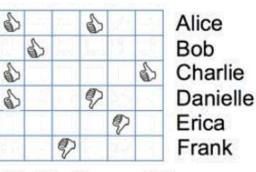




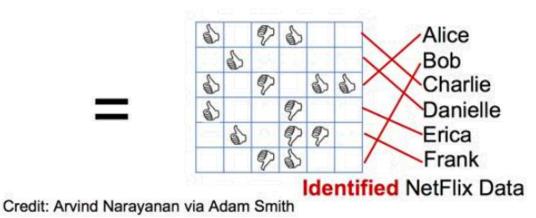
Use Public Reviews from IMDB.com



Anonymized NetFlix data



Public, incomplete IMDB data



Narayanan, Shmatikov, <u>Robust De-</u> anonymization of Large Datasets (How to <u>Break Anonymity of the Netflix Prize</u> <u>Dataset</u>), 2008





- Lesson: cannot always anonymize data simply by removing identifiers
- Vulnerable to aggregating data from multiple sources/networks
- Paul Ohm says: "For almost every person on earth, there is at least one fact about them stored in a computer database that an adversary could use to blackmail, discriminate against, harass, or steal the identity of him or her."
 - In "Broken promises of privacy: Responding to the surprising failure of anonymization"

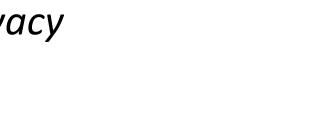


Other attacks:

- Su et al, De-anonymizing Web Browsing Data with Social Networks, 2017
 - Using links appearing in one's feed is unique
- Korolova, "Privacy Violations Using Microtargeted Ads: A Case Study", PADM
 - Attackers can instrument ad campaigns to identify individual users.
- Calandrino, Kilzer, Narayanan, Felten, Shmatikov, "You Might Also Like: Privacy Risks of Collaborative Filtering"
 - Attackers can infer customers' transaction with a limited amount of auxiliary data

What do we mean by privacy?

- Samuel Warren and Louis Brandeis (1890)
 - <u>"Right to be left alone</u>"
 - -In The Right to Privacy







- Alan Westin (1967)
 - "Right to control, edit, manage, and delete information about themselves and decide when, how, and to what extent information is communicated to others"



Credits



- Nikhil Jha Politecnico di Torino
- Eduardo Cuervo Oculos
- Amre Shakimov Duke University
- Frank McSherry Materialize, Inc.
- Krishnaram Kenthapadi AI @ LinkedIn
- Ilya Mironov Google Al
- Abhradeep Thakurta UC Santa Cruz



Perguntas Fragen DomandeGaldera Otázky Otazky Ouestions Spørgsmål Pertanyaan kysymykset Frågor Spørsmål Cwestiynau вопросыPreguntes Sorular Въпроси Vragen Pytania



8-12th November 2021 – Politecnico Di Milano, Italy

Privacy-Preserving Data Processing

Luca Vassio, Martino Trevisan November 8^{th,} 2021



- 1. The rise of Data-Driven approaches and problematics
- 2. Why we need Anonymization and why it is difficult
- 3. Privacy-Preserving Techniques
 - k-Anonymity
 - Differential Privacy
- 4. Open-Source tools for anonymization
- 5. Hands-on on Data

k-Anonymity

And other cluster-based methods

Classification of attributes



Key attributes

- Anything that identifies the person directly personally identifying information (PII)
- Name, address, phone number, email
- Always removed before release

Quasi-identifiers

- Not unique, but sufficiently well correlated with a person such that they can be combined with other quasi-identifiers to create a unique identifier
- (5-digit ZIP code, birth date, gender) uniquely identify 87% of the population in the U.S.
- Can be used for linking anonymized dataset with other datasets

4

Classification of attributes



Sensitive attributes

- Medical records, salaries, etc.
- These attributes is what the researchers need, so they are always released directly

Key Attribute	Qu	asi-identifier	Sensitive attribute	
Name	DOB	Gender	Zipcode	Disease
Andre	1/21/76	Male	53715	Heart Disease
Beth	4/13/86	Female	53715	Hepatitis
Carol	2/28/76	Male	53703	Brochitis
Dan	1/21/76	Male	53703	Broken Arm
Ellen	4/13/86	Female	53706	Flu
Eric	2/28/76	Female	53706	Hang Nail
5				

Re-identification by linking



[Latanya Sweeney, 1997]

Massachusetts hospital discharge dataset

Medical Data Released as Anonymous

SSN	Name	Ethnicity	Date Of Birth	Sex	ZIP	Marital Status	Problem
		asian	09/27/64	female	02139	divorced	hypertension
	2 2	asian	09/30/64	female	02139	divorced	obesity
		asian	04/18/64	male	02139	married	chest pain
	5 5	asian	04/15/64	male	02139	married	obesity
	3 3	black	03/13/63	male	02138	married	hypertension
		black	03/18/63	male	02138	married	shortness of breath
	21 B	black	09/13/64	female	02141	married	shortness of breath
		black	09/07/64	female	02141	married	obesity
	S 5	white	05/14/61	male	02138	single	chest pain
	10 I	white	05/08/61	male	02138	single	obesity
		white	09/15/61	female	02142	widow	shortness of breath

Public voter dataset



Voter List							
Name	Address	City	ZIP	DOB	Sex	Party	
Sue J. Carlson	1459 Main St.	Cambridge	02142	9/15/61	female	democrat	





- The information for each person contained in the released table cannot be distinguished from at least k-1 individuals whose information also appears in the release
 - Example: you try to identify a person in the released table, and the information you have is his birth date, ZIP and gender. However there are k people in the table with the same combination of birth, ZIP and gender
- Any combination of quasi-identifier present in the released table must appear in at least k records

First defined:

[Samarati and Sweeney. "Protecting privacy when disclosing information: *k*-anonymity and its enforcement through generalization and suppression", 1998]





[Samarati et al, 1998]

 Generalize, modify, or distort quasi-identifier values so that no individual is uniquely identifiable from a group of k
 In SQL, released table T is k-anonymous if each

```
SELECT COUNT(*)
FROM T
GROUP BY quasi-identifier combination
```

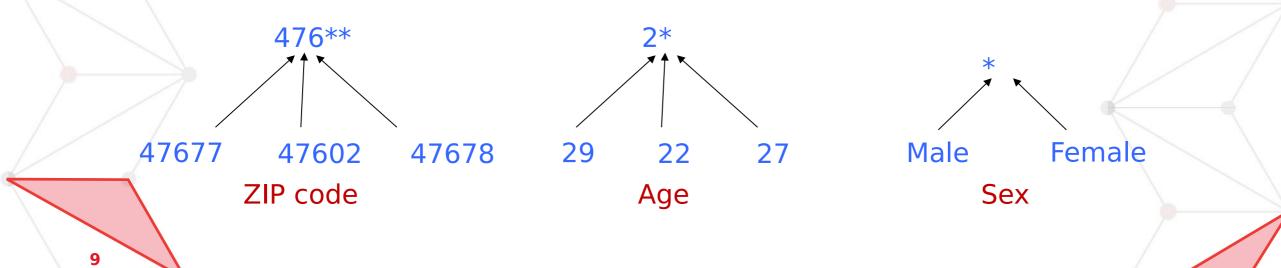
is **≥ k**

• Parameter k indicates the "degree" of anonymity

Achieving k-Anonymity: Generalization

POLITECNICO DI TORINO

- Goal of k-Anonymity
 - Each record is indistinguishable from at least k-1 other records
 - These k records form an equivalence class
- Generalization: replace quasi-identifiers with less specific, but semantically consistent values





	Race	Birth	Gender	ZIP	Problem
t1	Black	1965	m	0214*	short breath
t2	Black	1965	m	0214*	chest pain
t3	Black	1965	f	0213*	hypertension
t4	Black	1965	f	0213*	hypertension
t5	Black	1964	f	0213*	obesity
tб	Black	1964	f	0213*	chest pain
t7	White	1964	m	0213*	chest pain
t8	White	1964	m	0213*	obesity
t9	White	1964	m	0213*	short breath
t10	White	1967	m	0213*	chest pain
t11	White	1967	m	0213*	chest pain

Figure 2 Example of k-anonymity, where k=2 and QI={Race, Birth, Gender, ZIP}



Released table

						_
	Race	Birth	Gender	ZIP	Problem	
t1	Black	1965	m	0214*	short breath	
t2	Black	1965	m	0214*	chest pain	
t3	Black	1965	f	0213*	hypertension	
t4	Black	1965	f	0213*	hypertension	
t5	Black	1964	f	0213*	obesity	
tб	Black	1964	f	0213*	chest pain	
t7	White	1964	m	0213*	chest pain	
t8	White	1964	m	0213*	obesity	
t9	White	1964	m	0213*	short breath	
t10	White	1967	m	0213*	chest pain	
t11	White	1967	m	0213*	chest pain	
	W IIIC	1707	ш	0215	circa pair	

ZIP Name Birth Gender Race 1964 02135 White Andre m Beth 1964 f 55410 Black 90210 White Carol 1964 f Dan 1967 02174 White m 02237 Ellen 1968 f White

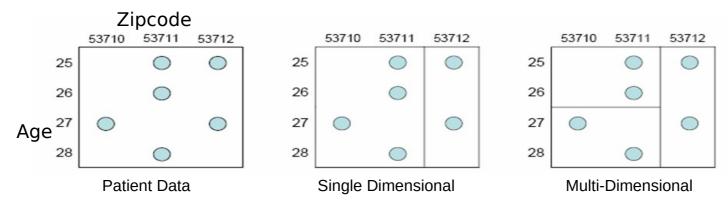
External data Source

By linking these 2 tables, you still don't learn Andre's problem

Generalization algorithms



- There are tens of k-anonymization algorithms
- Example: Greedy Partitioning Algorithm, Bucketization, ...
- **Problem:** find multi-dimensional partitions of the data, where each partition has two or more data points (i.e. k=2)



Optimal k-anonymous strict multi-dimensional partitioning is NP-hard
 Optimal = minimum information loss = maximum utility

Achieving k-Anonymity



Generalization

Replace specific quasi-identifiers with less specific values

Suppression

- When generalization causes too much information loss
 - This is common with "outliers"

Determining the value of k

The value of k depends on

- Number of records in the table
- Number of quasi-identifiers
- The distribution of each quasi-identifier
- The relationship between quasi-identifier

Rule of thumb:

- k increases \rightarrow privacy increases
- k increases \rightarrow utility decreases

[27-28]

2-Anonymized table

ID	Age	Sex	Zip	Disease
1	[25-26]	Male	53711	Flu
3	[25-26]	Male	53711	Brochities
2	[25-27]	Female	53712	Hepatitis
5	[25-27]	Female	53712	HIV
4	[27-28]	Male	[53710- 53711]	Broken Arm
6	[27-28]	Male	[53710- 53711]	Hang Nail

	3-Anonymized table						
ID	Age	Sex	Zip	Disease			
1	[25-26]	*	5371*	Flu			
3	[25-26]	*	5371*	Brochities			
2	[25-26]	*	5371*	Hepatitis			
5	[27-28]	*	5371*	HIV			
4	[27-28]	*	5371*	Broken Arm			

5371*

Hang Nai

4-Anonymized table

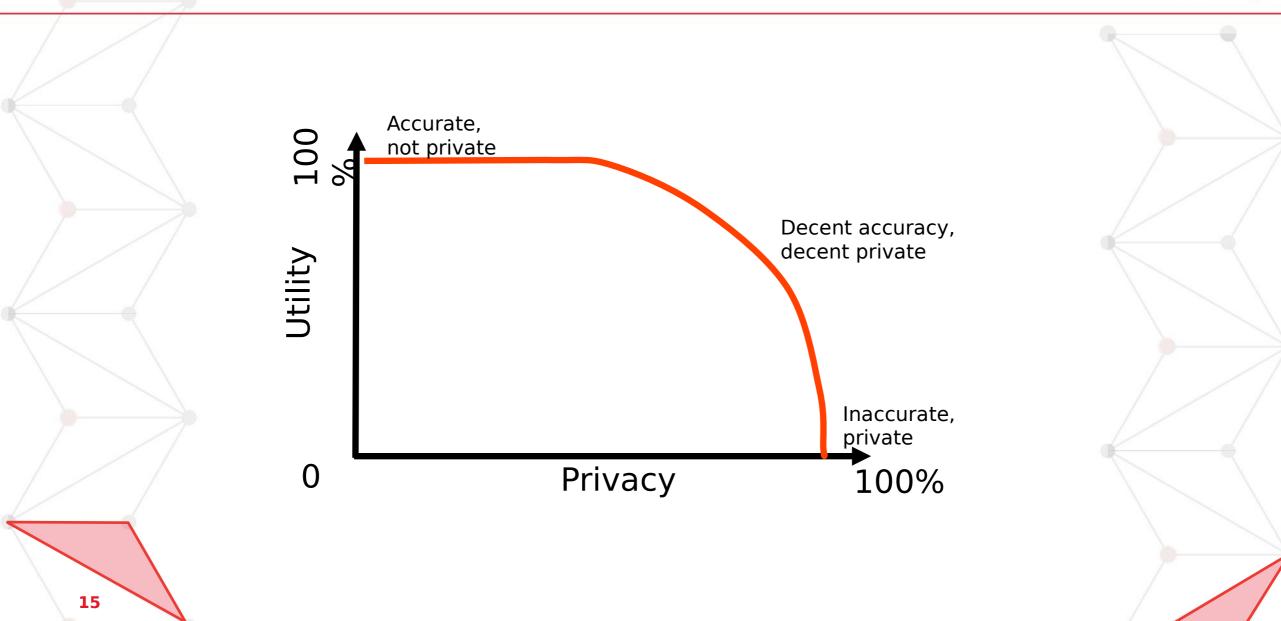
ID	Age	Sex	Zip	Disease
1	[25-28]	*	5371*	Flu
3	[25-28]	*	5371*	Brochities
2	[25-28]	*	5371*	Hepatitis
5	[25-28]	*	5371*	HIV
4	[25-28]	*	5371*	Broken Arm
6	[25-28]	*	5371*	Hang Nail

What is the best value of k?



Utility vs. Privacy





Curse of dimensionality



Generalization fundamentally relies on spatial locality

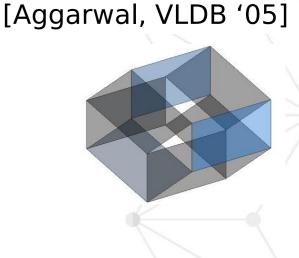
• Each record must have k close neighbors

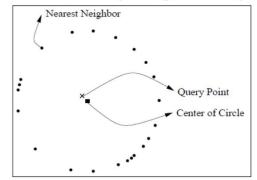
Real-world datasets are very sparse

- Many attributes (dimensions)
 - Netflix Prize dataset: 17000 dimensions
 - Amazon customer records: several million dimensions
- "Nearest neighbor" is very far

Projection to low dimensions loses all info

• k-anonymized datasets results useless





Attacks on k-Anonymity



k-Anonymity does not provide privacy if:

- Sensitive values in an equivalence class lack diversity
- The attacker has background knowledge

Homogeneity attack

Bob		
Zipcode	Age	
47678	27	

Background knowledge attack

Carl		
Zipcode	Age	
47673	36	

Carl does not have heart disease

A 3-anonymous patient table

	Zipcode	Age	Disease
(476**	2*	Heart Disease
	476**	2*	Heart Disease
	476**	2*	Heart Disease
	4790*	≥40	Flu
	4790*	≥40	Heart Disease
	4790*	≥40	Cancer
(476**	3*	Heart Disease
	476**	3*	Cancer
	476**	3*	Cancer

A solution: I-Diversity



[Machanavajjhala et al., ICDE '06]

Principle

18

- Each equi-class contains at least I well-represented sensitive values

Zipcode	Age	Disease
476**	2*	Acne
476**	2*	Heart disease
476**	2*	Flu
476**	2*	Heart disease
476**	2*	Flu
476**	2*	Flu
476**	3*	Flu
476**	3*	Acne
476**	3*	Cancer
476**	3*	Acne
476**	3*	Flu
476**	3*	Cancer

Sensitive attributes must be "diverse" within each quasi-identifier equivalence class

I-Diversity



Distinct I-Diversity

• Each equivalence class contains I distinct sensitive values

Probabilistic I-Diversity

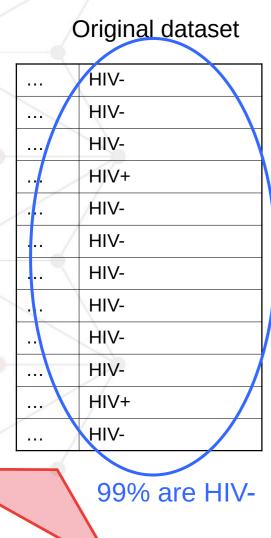
 The frequency of the most frequent value in an equivalence class is bounded by 1/l

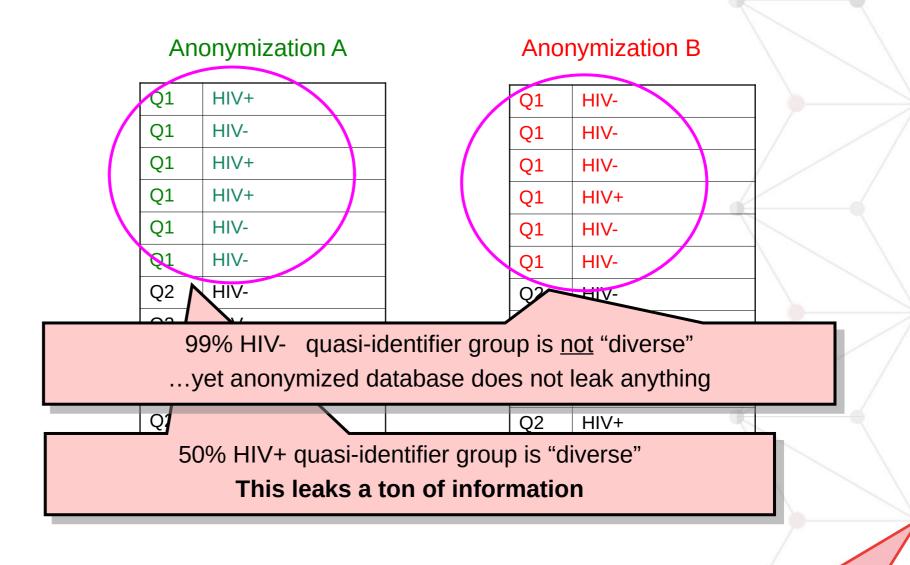
Entropy I-Diversity

 The entropy of the distribution of sensitive values in each equivalence class is at least log(I)

Neither sufficient, nor necessary









- Example: sensitive attribute is HIV+ (1%) or HIV- (99%)
 - Very different degrees of sensitivity!
- I-diversity is unnecessary
 - 2-diversity is unnecessary for an equivalence class that contains only HIV- records
- I-diversity is difficult to achieve
 - Suppose there are 10000 records in total
 - To have distinct 2-diversity, there can be at most 10000*1%=100 equivalence classes





- \prec probabilistic inference attacks
- Example: sensitive attribute is HIV+ (1%) or HIV- (99%)
- Consider an equivalence class that contains an equal number of HIV+ and HIV- records
 - Diverse, but potentially violates privacy!
- I-diversity does not differentiate:
 - Equivalence class 1: 49 HIV+ and 1 HIV-
 - Equivalence class 2: 1 HIV+ and 49 HIV-

I-diversity does not consider overall distribution of sensitive values!

Sensitive attribute disclosure



	Similarity attack Bob			Zincodo Ago Solary Di			Disease
				Zipcode	Age	Salary	Disease
	Zip	Age		476**	2*	20K	Gastric Ulcer
	-			476**	2*	30K	Gastritis
	47678	27		476**	2*	40K	Stomach Cancer
				4790*	≥40	50K	Gastritis
Conclu	ision		4790*	≥40	100K	Flu	
1 Bol	b's salary is	which is	4790*	≥40	70K	Bronchitis	
	atively low		476**	3*	60K	Bronchitis	
				476**	3*	80K	Pneumonia

476**

A 3-diverse patient table

3*

90K

Stomach Cancer

Bob has some stomach-related disease 2.

I-diversity does not consider semantics of sensitive values!



- k-anonymity prevents identity disclosure but not attribute disclosure
- To solve that problem **I-Diversity** requires that each eq. class has at least I values for each sensitive attribute
- t-Closeness requires that the distribution of a sensitive attribute in any eq. class is close to the distribution of a sensitive attribute in the overall table

t-Closeness protects against attribute disclosure but not identity disclosure

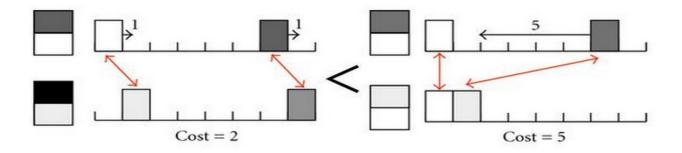
t-Closeness



Distribution of sensitive attribute value in each equi-class should be close to that of the overall dataset (within a distance t)

- Earth Mover Distance to capture semantic relationship among sensitive attribute values
- Privacy is measured by the information gain of an observer regarding distribution of sensitive attributes

Information Gain = Posterior Belief – Prior Belief



[Li, Li, Venkatasubramanian, ICDE'2007]



	ZIP Code	Age	Salary	Disease	
1	4767*	≤ 40	3K	gastric ulcer	
3	4767*	≤ 40	5K	stomach cancer	
8	4767*	≤ 40	9K	pneumonia	
4	4790*	≥ 40	6K	gastritis	
5	4790*	≥ 40	11K	flu	
6	4790*	≥ 40	8K	bronchitis	
2	4760*	≤ 40	4K	gastritis	
7	4760*	≤ 40	7K	DIOIICIIIIN	lestion: Why publish lentifiers at all??
9	4760*	≤ 40	10K	stomach cancer	

Table 5. Table that has 0.167-closeness w.r.t. Salary and 0.278-closeness w.r.t. Disease

k-Anonymous, I-Diverse t-Close dataset



		\bigcap	\bigcap
Caucas	787XX	HIV+	Flu
Asian/AfrAm	787XX	HIV-	Flu
Asian/AfrAm	787XX	HIV+	Shingles
Caucas	787XX	HIV-	Acne
Caucas	787XX	HIV-	Shingles
Caucas	787XX	HIV-	Acne
27		\bigcirc	\bigcirc

This is k-anonymous, I-diverse and t-close...

It is secure, right?

What does an attacker know?



Bob is Caucasian and I				M		
heard he was admitted to hospital with flu	Cauc	cas	787XX	HIV+	Flu	
	This is against t	/AfrAm he rules!	787XX	HIV-	Flu	\langle
	"flu" is not a qua	si-identifier	787XX	HIV+	Shingles	\backslash
	Yes and this is problem with k-a		787XX	HIV-	Acne	
	Cauc		787XX	HIV-	Shingles	
	Cauc	cas	787XX	HIV-	Acne	

k-Anonymity can be harmful



- Focuses on data transformation, not on what can be learned from the anonymized dataset
 - "k-anonymous" dataset can leak sensitive information
- "Quasi-identifier" fallacy
 - Assumes a priori that attacker will not know certain information about the target
- Relies on locality
 - Destroys utility of many real-world datasets

Differential Privacy

A rigorous framework for privacy-preserving analysis of datasets



It's hard to guess what capabilities attackers will have, especially in the future

- Future datasets
- Future techniques
- Future computational power

Analogy with crypto: Cryptosystems today are designed based on what quantum computers might be able to do in 30 years



Strong, quantifiable, composable mathematical privacy guarantee

It is (by design) resilient to known and unknown attack modes!

DP enables many computations with personal data while preserving personal privacy

First defined:

[Dwork, McSherry, Nissim, and Smith, Calibrating Noise to Sensitivity in Private Data Analysis, in Third Theory of Cryptography Conference, TCC 2006.] Earlier roots:

[Warner, Randomized Response 1965]

Learn from datasets

A research dataset shows that smoking causes cancer

- Smoker S's insurance premiums rise
- However, this is true even if S not in database!
- Learning that smoking causes cancer is the whole point
- Smoker S enrolls in a smoking cessation program...

Differential privacy: no harm in participation

 Outcome of any analysis is essentially equally likely, independent of whether any individual joins, or refrains from joining, the dataset





Differential Privacy

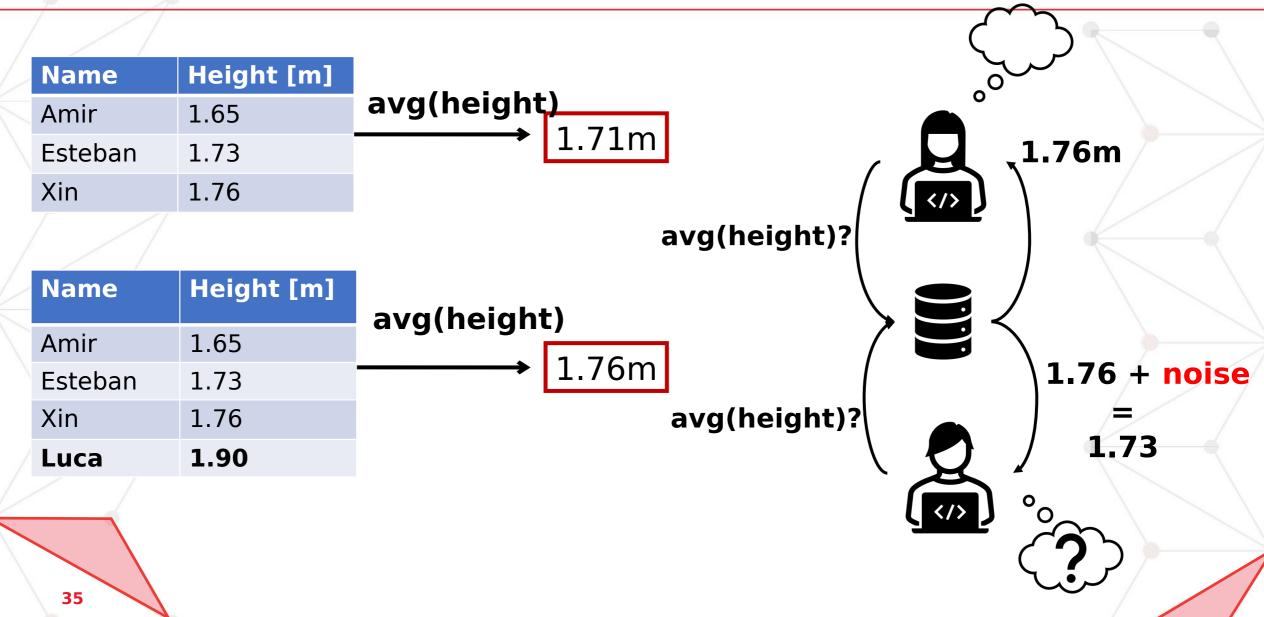


- Worst-case scenario: an attacker that knows everything, but not if the user is in the dataset
 - The attacker knows the value of every record, even the target user's one
- The attacker can perform **queries** to the dataset
 - The answer of the dataset should not hint the presence of the user

"Blatant non-privacy" [Dinur and Nissim, 2003] Overly accurate answers to too many questions destroys privacy

Example: is Luca in the dataset?





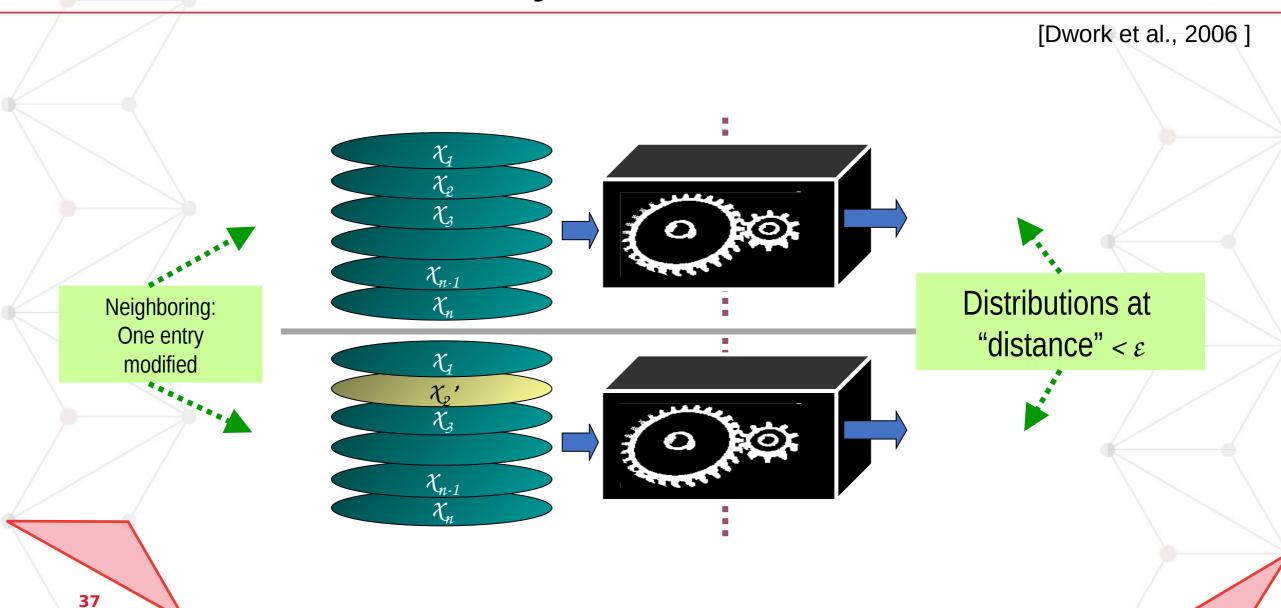
Differential privacy basic setup



- There is a database D which potentially contains sensitive information about individuals
- The **database curator** has access to the full database. We assume the curator is trusted
- The data analyst wants to analyze the data. She asks a series of queries to the curator, and the curator provides a response to each query
- The way in which the curator responds to queries is called the **mechanism**
- Two databases D and D' are neighbouring if they agree except for a single entry

Differential Privacy





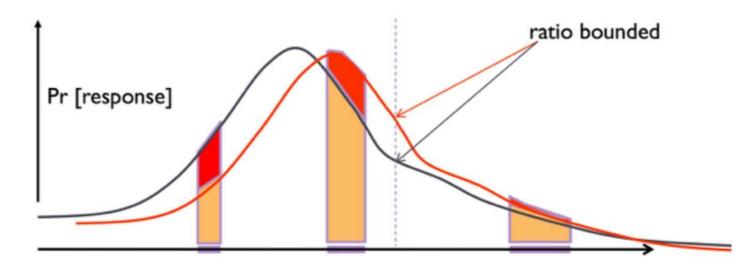
Formally...



[Dwork et al. 2006]

A query mechanism M is ϵ -differentially private if, for any two adjacent databases D and D' (differing in just one entry) and $C \subseteq range(M)$

 $\Pr(M(D) \in C) \le e^{\epsilon} \cdot \Pr(M(D') \in C)$



- Output does not overly depend on any single tuple
- Participation in the dataset poses **no additional risk**

Privacy parameter ε



A query mechanism M is ϵ -differentially private if, for any two adjacent databases D and D' (differing in just one entry) and $C \subseteq range(M)$

 $\Pr(M(D) \in C) \le e^{\epsilon} \cdot \Pr(M(D') \in C)$

 $\boldsymbol{\epsilon}$ is an arbitrary parameter and controls the privacy of the system

- Low ϵ , the two quantities are forced to be similar
- High ε, the two quantities are allowed to diverge
- However, all queries are not the same
 - Some of them may be more intrusive
 - What differentiate them is *sensitivity*



The sensitivity

- A query can be more or less intrusive on the privacy of the users
- Each query f has a sensitivity
 - Depends on the **difference** in output between f(D) and f(D')
 - Global sensitivity $GS_f = \max_{neighbors D,D'} ||f(D) f(D')||_1$
- The added noise depends on the query sensitivity
 - Sensitive query \rightarrow more noise
 - Insensitive query \rightarrow less noise

Implications



Global sensitivity: $GS_f = \max_{\text{neighbors } D, D'} ||f(D) - f(D')||_1$

• Example: Query f is the average for sets of numbers between 0 and 1

 $GS_{average} = 1/n$

Some queries, such as counting queries, can be answered relatively accurately

Since one tuple affects the result by at most 1 (GS=1)

A small amount of noise can be added to achieve DP

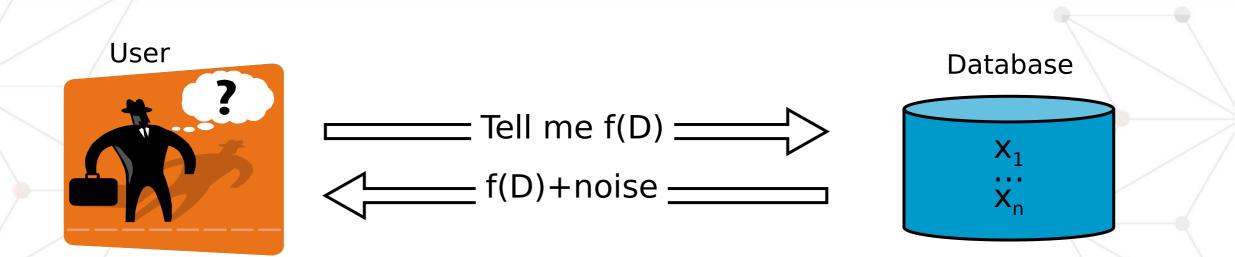
Some queries are hard to answer

E.g., max, since it can be greatly affected by a single tuple (GS unbounded) Challenge in using it

Find suitable queries to ask so that noisy answers provide most utility

A DP mechanism: output perturbation





Intuition: f(D) can be released accurately when f is insensitive to individual entries $x_1, \dots x_n$

We want f(D) + noise to be ε -indistinguishable

How this noise should be generated depending on ε and GS_f ?

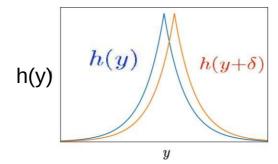
Laplace noise

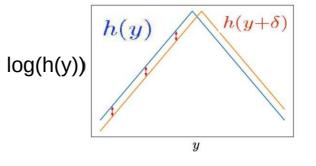


Theorem

If $A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$ then A is ε -indistinguishable.

Laplace distribution $Lap(\lambda)$ has density $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$





Laplace noise



Theorem

If $A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$ then A is ε -indistinguishable.

Laplace distribution $Lap(\lambda)$ has density $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$

It quantifies:

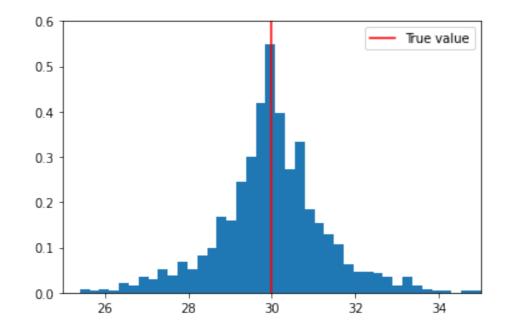
Noise to add when function is more sensitive (higher GS_f)

Noise to add if we want more privacy (lower ε)

Multiple queries



- What if we allow people to perform the same query over and over again?
- Eventually, the noise will cancel out and the true value will emerge



Properties of Differential Privacy



Group privacy

• Group privacy: ɛk privacy for a group of size k

Composability

- Applying the sanitization several time yields a graceful degradation
- If A₁ satisfies ε₁-DP, and A₂ satisfies ε₂-DP, then outputting both A₁ and A₂ satisfies (ε₁+ε₂)-DP
- even just proportional to square root of number of applications

Robustness to side information

- No need to specify **exactly** what the adversary knows
- Any post-processing cannot improve the attacker's knowledge

Privacy budget



- To avoid noise cancelling, ε becomes a *privacy budget* to query the dataset
- The privacy consumption is <u>additive</u>
 - With budget 10, I can choose
 - 10x queries with $\varepsilon = 1$
 - 5x queries with $\epsilon=2$
 - 1x query with $\epsilon = 10$

How to increase number of queries?

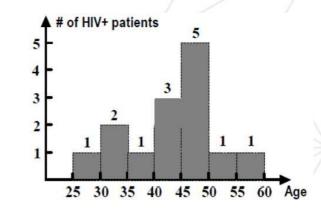


[Blum-Ligget-Roth, 2008]

Use Coordinated Noise

If noise is added in with careful coordination, rather than independently, DP can save on the budget Example: histogram queries with d bins:

- You can treat query independently: use d times the mechanism, hence noise Lap(d/ε)
- But actually only need Lap(1/ε), since sensitivity generalizes as max l₁ distance



A glimpse inside the rabbit hole



- What if we don't trust the data curator?
 - Local differential privacy, where data are anonymized before being uploaded to a platform
- What if we want to publish the dataset at once, without having to continuously answer to the queries?
 - Non-interactive differential privacy
- What if we have categorical values?
 - Other mechanisms exist: the exponential mechanism allows to choose the most-suited value inside a categorical set
- What if we assume the attacker does not know everything?
 - Relaxed differential privacy: only statistical knowledge of D and D', keeping precise knowledge on the element changing

DP: Pros & Cons





- **Rigorous** mathematical definition of privacy
- Flexible: several mechanisms are available
- Robust to postprocessing
- The level of privacy ϵ can be chosen by the system administrator



- Precision of the queries is affected
- Hard to explain
- Many non-trivial DP algorithms require really large datasets to be practically useful
- What ε and what privacy budget is **reasonable** for a dataset?

Do I need DP if I don't care about privacy?



- Almost any usage of the data that is not carefully crafted will leak something about it
- Statistics to generalize well should not be dependent on single instances
- As for machine learning: it should not overfit on training
- We can use DP to ensure statistical validity of exploratory data

Better Privacy = Better Data

Credits

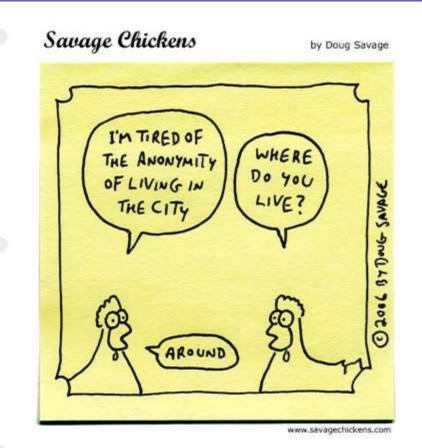


Nikhil Jha – Politecnico di Torino Ashwin Machanavajjhala - Duke University Vitaly Shmatikov - Cornell Tech Chris Clifton – Purdue Muhamad Felemban – KFUPM Moni Naor- Weizmann Institute of Science Adam Smith - Penn State Roger Grosse - University of Toronto Katrina Ligett - California Institute of Technology Cynthia Dwork - Harvard University Matteo Maffei – TU Wien

52



8-12th November 2021 – Politecnico Di Milano, Italy



53

Perguntas Fragen DomandeGaldera Otázky Ouestions Spørgsmål Pertanyaan, kysymykset Frågor Spørsmål Cwestiynau ВопросыPreguntes Sorular Въпроси Vragen Pytania



8-12th November 2021 – Politecnico Di Milano, Italy

Privacy-Preserving Data Processing

Luca Vassio, Martino Trevisan Performance 2021 November, 8

Outline



1. The rise of Data-Driven approaches and problematics 2. Why we need Anonymization and why it is difficult 3. Privacy-Preserving Techniques **K-anonymity** • Differential Privacy 4. Open-Source tools for anonymization 5. Hands-on on Data

Useful Tools



- 1. Arx
- 2. DiffPrivLib
- 3. Google Differential Privacy
- 4. P-PPA
- 5. [CNIL]



ARX is a comprehensive open-source software for anonymizing sensitive personal data. It supports a wide variety of

- privacy (and risk) models
- methods for transforming data
- methods for analyzing the usefulness of output data

Has a graphical interface and APIs Can read data in:

- CSV
- Excel
- SQL DBs (MySQL, PostgreSQL)

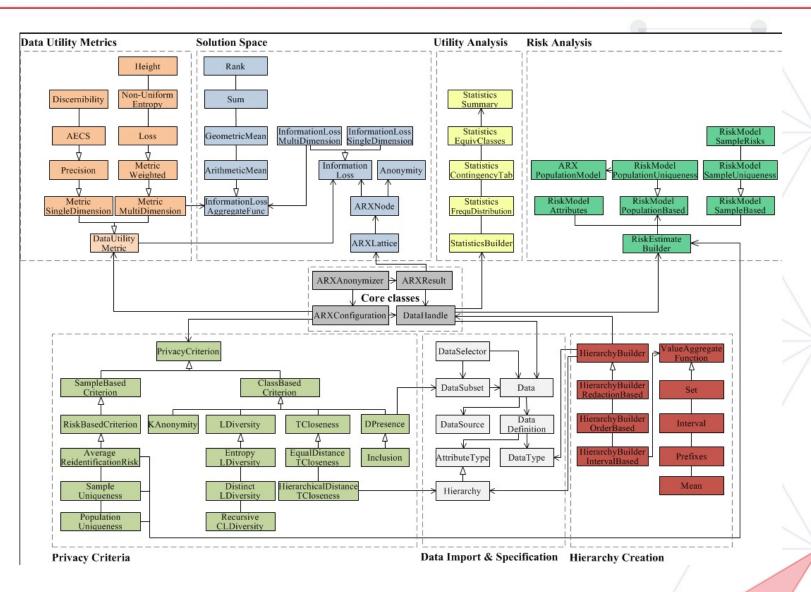
	A DAY BA	Test-Darryle		-	-		_	_			
6.14.	N 10 10 1	X V B. A.B.	ALERIA								-
Cont	igure that of part	nation fast	4 2 4 4	日至 9							
laput ch	100	nates & Espicera	into per Analyze and	Ry 🛛 # Analysis	44				Interior	NAMES AND DESCRIPTION OF TAXABLE PARTY AND ADDRESS OF TAXABLE PARTY.	
	-	-			1.10	5 E o	-				
1	3	Fack	Warbi day	- shartes	Alle Gard		- dpa	-		1118	
2	54	Reck	Descoil	Some-callege Rectedors	United Dates	Phani.	1	3.9	Bat I	a materiality electric procession	
	35	Back	Derod	Anno andro	Venui Genu	Pines	2	3-9	Det .	amont of presed tight objection facts from the facts of times of presed tight objection facts from the facts	8
	52	Reck	Descost	Motors.	United Dates United Dates	Pilet at	181	3.9	But	sprove or preased rapper advention. North Amazina . Manufac	
5	53	Back	Note maried	Som-callege	United-Barter	Pluty	* 5		Red Red	spread out pread higher education light investor. Manufer,	
7	12	Back	Separated Reported	Sona-callege	United States	Peak	. 6	3.0	Ball	server of preset type advantary facts investory. New Service server of preset types advantary. New Service	
	52	Black	Separated	Associations Some callege	United Matte	Perats	181	3.9	Bet	tence of passed light encoder light leases. No-lie	
	56	Back	Midnand	Sere-cellege	Verse Gen	Prim		3.9	Red Red	server of proof light elector. Neth Analia - Net-In- tence of proof light elector. Neth Analia - Net-In-	
30	52	illinge	Diverced	Sens-callege	Dated Brief	Federa	10	8.0	ilia.	geore or yeard fight stocator. Soft locator . Georee	
11	54	Minite Minite	Dvared Dvarest	Bathelon Moters	United States United States	Federal Local-	11	3.9	The state	genus of pears light abuilton light inversion Gaussing genus of pears light education light inversion. Gaussing	
12	52	and the second	Dearcel	Serve-callege	United Miles	Local-	12	3.9	-	grass styneed light studies light lensing forem	
34	38	HINK	Dearond	Baltelos	United Skills	Loui-	18	34.00	-	space or present higher advantage listed downey. Second	
15	56	INNE	Descol	Some-cellege	United Street	Sant-		3.9	1014	gener et prant light elucitor line inner	
4	Laura			_			*	-		ngang (Salagang Tang Augusta)	
ofe	e Indan Essan	a kaselig Marier	Dark M2	Olive	e Min e B sink	-	100 000 000 000	scars Parit	unter and		
-	_			Ofw Ofw				ions Arch	_		\ -
-	_	 Austrijschlander Austrijschlander 	92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M			_	_		
-	_		92 23	_	12.M	Í		_	_		
-	_		92 23	_	12.M	Ĩ		_	_		
-	_		92 23	_	12.M		10.0 0.00				
-	_		92 23	_	12.M		10.0 0.00				
-	_		92 23	_	12.M		10.0 0.00				
-	_		92 23	_	12.M		10.0 0.00				
-	_		92 23	_	12.M		10.0 0.00				
-	_		92 23	_	12.M		10.0 0.00		_		



Privacy models:

- K-anonymity
- I-diversity
- t-closeness
- Differential Privacy (initial)
- Some others

Supports various metrics for measuring information loss





Can be used as a library

- Written in Java 😕
- All functionalities (privacy models and metrics) available to developers
- Documentation: https://arx.deidentifier.org/wp-content/uploads/javadoc/current/api/index.html

There is a project called Arx as a Service that offers Arx API as Web API:

<u>https://navikt.github.io/arxaas/</u>

Limitations:

- Only Java API (but there are web api)
- Arx as a Service allows using Python, but the package is still a bit preliminary
 - We will use it in the lab



curl 'http://localhost:8080/api/anonymize' -i -X POST \ -H 'Content-Type: application/json' \ -d '{ "data" : [["age", "gender", "zipcode"], ["34", "male", "81667"], ["35", "female", "81668"], ["36", "male", "81669"], ["37", "female", "81670"], ["38", "male", "81671"], ["39", "female", "81672"], ["40", "male", "81673"], ["41", "female", "81674"], ["42", "male", "81675"], ["43", "female", "81676"], ["44", "male", "81677"]], "attributes" : [{ "field" : "age", "attributeTypeModel" : "IDENTIFYING", "hierarchy" : null }, { "field" : "gender", "attributeTypeModel" : "SENSITIVE", "hierarchy" : null }, { "field" : "zipcode", "attributeTypeModel" : "QUASIIDENTIFYING", "hierarchy" : [["81667", "8166*", "816**", "81***", "8****", "*****"], ["81668", "8166*", "816**", "81***", "8****", "*****"], ["81669", "8166*", "816**", "81***", "8****", "*****"], ["81670", "8167*", "816**" "81***", "8****", "*****"], ["81671", "8167*", "816**", "81***", "8****", "*****"], ["81672", "8167*", "816**", "81***", "8****", "*****"], ["81673", "8167*", "816**", "81***", "8****", "*****"], ["81674", "8167*", "816**", "81***", "8****", "*****"], ["81675", "8167*", "816**", "81***", "8***", "*****"], ["81676", "8167*", "816**", "81***", "8****", "*****"], ["81677", "8167*", "816**", "81***", "8****", "*****"] }], "privacyModels" : [{ "privacyModel" : "KANONYMITY", "params" : { "k" : "5" "privacyModel" : "LDIVERSITY DISTINCT", "params" : { "column_name" : "gender", "|" : "2"

HTTP/1.1 200 OK Vary: Origin Vary: Access-Control-Request-Method Vary: Access-Control-Request-Headers Content-Type: application/json Content-Length: 7813

{

"anonymizeResult" : {

"data": [["age", "gender", "zipcode"], ["*", "male", "816**"], ["*", "female", "816**"], "*", "male", "816**"], ["*", "female", "816**"], ["*", "male", "816**"]],

"riskProfile" : {
 "reldentificationRisk" : {
 "measures" : {
 "estimated_journalist_risk" : 0.090909090909090909090,
 "records_affected_by_highest_prosecutor_risk" : 1.0,
 "sample_uniques" : 0.0,
 "lowest_risk" : 0.09090909090909091,

POLITECNICO DI TORINO

How to use it in Python:

1.

2.

3.

Launch an ARX Server, using the Dockerized version

docker run -p 8080:8080 navikt/arxaas

Install pyarxaas

pip install pyarxaas

Use it in Python

```
arxaas = ARXaaS(url)
```

anon_result = arxaas.anonymize(dataset, [kanon])



Developed by IBM

"This is a library dedicated to **differential privacy** and machine learning. Its purpose is to allow experimentation, simulation, and implementation of differentially private models using a common codebase and building blocks" Written in Python Documentation at: https://diffprivlib.readthedocs.io/en/latest/

IBM DiffPrivLib

Supports:

- Basic Mechanisms for Differential Privacy
 - Laplacian mechanism
 - Geometric Mechanism
 - Exponential Mechanism
- Tools for Differential Privacy functions
 - Mean, Sum, Histogram

Differentially-private machine learning models

- Classification models
 - Gaussian Naive Bayes
 - Logistic Regression
 - Random Forest
- Regression models
 - Linear Regression
- Clustering models
 - K-Means
- Dimensionality reduction models
 - PCA
- Preprocessing
 - Standard Scaler



IBM DiffPrivLib

Differentially Private Histograms

Plotting the distribution of ages in Adult

In [1]: import numpy as np
from diffprivlib import tools as dp
import matplotlib.pyplot as plt

We first read in the list of ages in the Adult UCI dataset (the first column).

From:

https://github.com/IBM/differential-privacylibrary/blob/main/notebooks/histograms.ipynb

Differentially private histograms

Using diffprivlib, we can calculate a differentially private version of the histogram. For this example, we use the default settings:

- epsilon is 1.0
- range is not specified, so is calculated by the function on-the-fly. This throws a warning, as it leaks privacy about the data (from dp_bins, we know that there are people in the dataset aged 17 and 90).

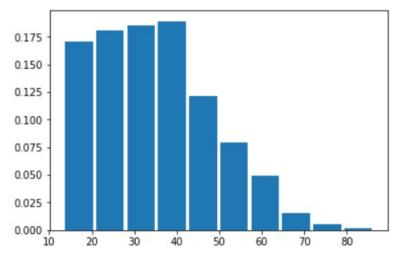
```
dp_hist, dp_bins = dp.histogram(ages_adult)
dp_hist = dp_hist / dp_hist.sum()
```

plt.bar(dp_bins[:-1], dp_hist, width=(dp_bins[1] - dp_bins[0]) * 0.9)
plt.show()

.../site-packages/diffprivlib/tools/histograms.py:131: PrivacyLeakWarnin g: Range parameter has not been specified. Falling back to taking range from the data.

To ensure differential privacy, and no additional privacy leakage, the r ange must be specified independently of the data (i.e., using domain kno wledge).

"specified independently of the data (i.e., using domain knowledge).", PrivacyLeakWarning)



IBM DiffPrivLib

Differentially Private Histograms

Plotting the distribution of ages in Adult

In [1]: import numpy as np
from diffprivlib import tools as dp
import matplotlib.pyplot as plt

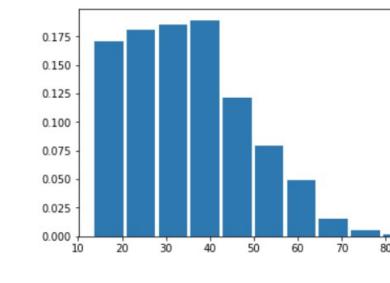
We first read in the list of ages in the Adult UCI dataset (the first column).

hist, bins = np.histogram(ages_adult) hist = hist / hist.sum()

Using matplotlib.pyplot, we can plot a barchart of the histogram distribution.

POLITECNICO DI TORINO

plt.bar(bins[:-1], hist, width=(bins[1]-bins[0]) * 0.9)
plt.show()

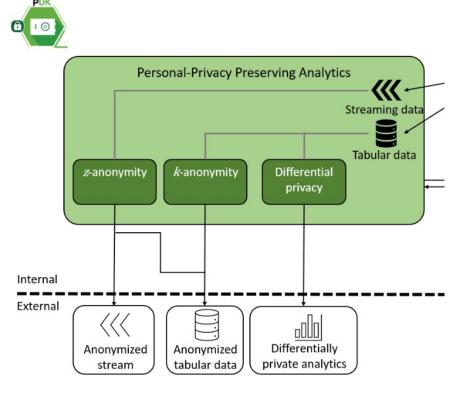




Library by Google to generate ε - and (ε , δ)-differentially private statistics over datasets Similar to IBM DiffPrivLib: <u>https://github.com/google/differential-privacy</u>

- Written in C, less mature Implements various mechanisms
- Laplace, Gaussian, etc. Various Statistics:
- Count, Sum, Mean, Variance, Quantiles... There exist Python bindings:
- https://github.com/OpenMined/PyDP
- Still preliminary

Developed by Politecnico di Torino In the context of the EU Project PIMCity (<u>https://www.pimcity-h2020.eu/</u>) It is a simple Python module that offers easy-to-use privacy models POLITECNICO DI TORINO





My

- The P-PPA can aggregate data from different sources and formats
 - Traditional RDBMS
 - MongoDB data lake
 - Data in CSV files



1	Α	B		с
1	Transaction Description	Expense Type	Amount	
2	Chester Diner	Restaurant	\$	24.22
3	Income Tax Payment	Taxes	\$	535.00
4	Ole Tymes Cafe	Restaurant	\$	12.58
5	Plane ticket to Melbourne	Travel	\$	654.32
6	Odessa's	Restaurant	\$	13.36
7	Car Rental in Australia	Travel	\$	185.55
8	K Crew	Clothing	\$	86.99
9	Ruby's Famous Bbq Joint	Restaurant	\$	5.66
10	Street Corner Market	Restaurant	\$	9.85
11	Airport Parking	Travel	\$	22.55
12	The Friendly Chef	Restaurant	\$	67.85
13	Floorgreen's	Personal Items	\$	24.55
14	Orange Democracy	Clothing	s	86.99
15	Car Care	Auto Expense	\$	24.22
16	The Narrow Lantern	Restaurant	s	101.33



- Written in Python
- Simple installation, only a few Python requirements
- Usage a Python module

from algorithms.kanonymity.mondrian.mondrian import Mondrian
mondrian = Mondrian(3, user_choice_index=[2,3,4,5])
k_anonymized_dataframe = mondrian.perform(input_dataframe)



• The P-PPA offers Web APIs to allow remote use

Personal-Pivacy Preserving Analytics (100 (100 (100)

The PIMCity P-PPA, it's a tool to allow data analysts and stakeholders to retrieve useful information from the data, while preserving the privacy of the users whose data are in the studied datasets. It follows some query example.

k-anonymity Provides k-anonymity privacy guarantees.

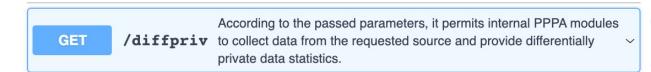


 \wedge

GET

/kanon According to the passed parameters, it permits internal PPPA modules to collect data from the requested source and provide k-anonymized data.

differential privacy Provides data statistics in a differentially private flavour.



CNIL-PIA software



The PIA software aims to help data controllers build and demonstrate compliance to the GDPR.

- The tools is available in French and in English.
- It facilitates carrying out a data protection impact assessment
- This tool also intends to ease the use of the PIA guides published by the CNIL.
- Portable and Web Version
- Website: <u>https://www.cnil.fr/en/open-source-pia-software-helps-carry-out-data-protection-impact-assesment</u>

Pla analyse d'impact sur la protection des données privacy impact assessment

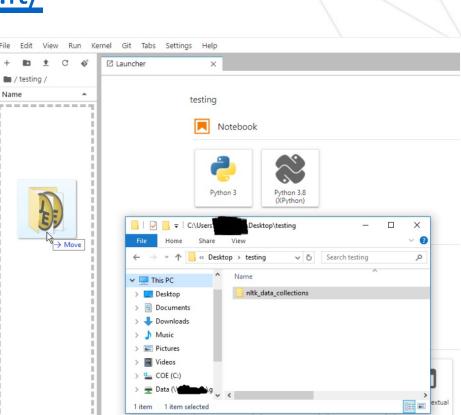




TOOL	Language	K-anon & similar	Diff Priv	Web API	Input	>
ARX	Java	Y	Partially	Y	CSV, Excel, SQL, JSON, Java API	/
DiffPrivLib	Python		Y		Numpy	
Google Differential Privacy	C/C++		Y		C vars	
P-PPA	Python	Υ	Y	Y	Pandas, CSV, SQL	

Access to the Lab

- Go with your browser on: <u>https://jupyter-dev.polito.it/</u>
- Login with:
 - Username: your mail
 - Password: your mail
- Download the notebook from: <u>https://www.dropbox.com/s/xqq66yjwejmuo2o/practice.ipynb?dl=1</u>
 - Can copy the link from chat
- Upload it on Jupyter (drag and drop works \odot)
- Complete the exercises!







Perguntas Fragen DomandeGaldera Otázky Otazky Ouestions Spørgsmål Pertanyaan kysymykset Frågor Spørsmål Cwestiynau вопросыPreguntes Sorular Въпроси Vragen Pytania