

Data Privacy

Garantia de la Informació i Seguretat

Departament d'Enginyeria de la Informació i les Comunicacions
Universitat Autònoma de Barcelona

1

Privacy (and anonymity)

- A very wide concept:
 - Privacy in communication networks
 - Data privacy (document confidentiality, private data publishing...)
 - Lots of other more
- **Privacy Enhancing Technologies (PET)**
 - used to widely refer to technologies used to provide privacy.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

2

2

Some examples of PETs

- The TOR network,
- Mix networks (Mixnets),
- Encryption,
- Use of pseudonyms,
- Anonymous credentials,
- Location privacy,
- ...

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

3

3

Private information publishing

- We will deal with a specific privacy problem:
data publishing
- Includes several disciplines, but mainly:
 - Statistical Disclosure Control
 - Privacy Preserving Data mining

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

4

4

Motivating example



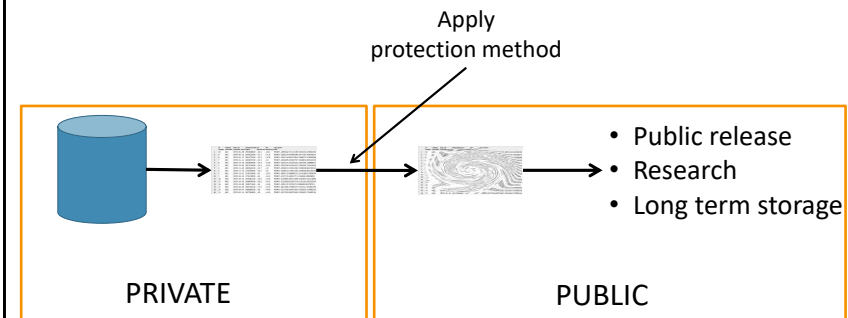
2021-05-02

[DEIC - UAB] G. Navarro-Arribas

5

5

Need to ensure privacy



2021-05-02

[DEIC - UAB] G. Navarro-Arribas

6

6

Definitions: anonymity

Anonymity of a subject means that the subject is not identifiable with a set of subjects, the **anonymity set**.

Unlinkability of two or more items of interest (e.g. subjects, messages, actions, ...) from an attacker's perspective means that within the system, the attacker cannot sufficiently distinguish whether these items are related or not

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

7

7

Definitions: disclosure

Disclosure takes place when attackers take advantages of the observation of available data to improve their knowledge on some confidential information about an item of interest.

• Identity disclosure

- The attacker can correctly identify a particular entity in the system, or can link some protected data to an specific entity.

• Attribute disclosure

- The attacker can learn something new about an attribute of an entity.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

8

8

Example of disclosure

- This table has been naively anonymized by removing identifier attributes (name, national ID number, etc.)

City	Age	Profession	Income
Campeche	30	Teacher	200
Campeche	30	Teacher	250
Campeche	30	Teacher	300
Palenque	27	Physician	350
Palenque	45	Police	250

- But, what happens if an attacker knows:
 - Alice is 30 years old and lives in Campeche
 - Bob is 45 years old and lives in Palenque

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

9

9

Example 2

Name	Passport	Profession	City	Age	Salary
Alfredo Pareja	AS3682	Teacher	Guayaquil	23	600
Eugenio Espejo	RX3453	Teacher	Quito	48	500
Raul Perez	TS5645	Policeman	Loja	20	734
Francisco Proaño	XB3456	Writer	Guayaquil	32	543
Cesar Davila	AE4324	Teacher	Cuenca	38	890
Ignacio Laso	FF3455	Writer	Quito	41	678
Raul Serrano	LO0903	Policeman	Portoviejo	21	399
Luz Argentina	KO8264	Policeman	Quevedo	44	943
Abdon Ubidia	RC3766	Teacher	Machala	31	632

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

10

10

Example 2: naïve anonymization

	Name	Passport	Profession	City	Age	Salary
Are you from Quevedo?	Identity disclosure	XXXXXX	Teacher	Guayaquil	32	600
		XXXXXX	Teacher	Quito	48	500
		XXXXXX	Teacher	Loja	20	734
		XXXXXX	Writer	Guayaquil	32	543
Are you from Guayaquil?	Attribute disclosure	XXXXXX	Teacher	Cuenca	38	890
		XXXXXX	Writer	Quito	41	678
		XXXXXX	Policeman	Portoviejo	21	399
		XXXXXX	Policeman	Quevedo	44	943
		XXXXXX	Teacher	Machala	31	632

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

11

11

Example: AOL Case

- AOL publicly releases in 2006: query logs by 650,000 users over 3 months
- Logs where apparently anonymized

```
24969 orioles tickets 2006-05-31 12:31:57 2 http://www.greatseats.com
24969 jennifer craford my space.com 2006-05-31 19:15:02
24969 jennifer crawford my space.com 2006-05-31 19:16:05
14423 boston redsoxweb page.com 2006-03-28 17:51:55
14423 www.bostonredsox 2006-03-28 18:12:26 1 http://boston.redsox.mlb.com
...
```

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

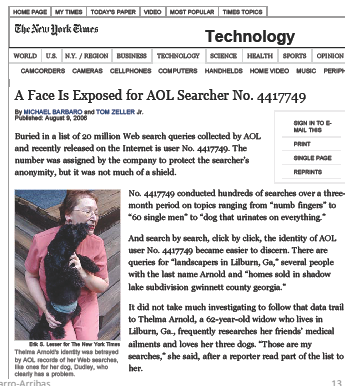
12

12

Example: AOL case

A user (Thelma Arnold) was easily re-identified

- <https://www.nytimes.com/2006/08/09/technology/09aol.html>
- https://en.wikipedia.org/wiki/AOL_search_data_leak



2021-05-02

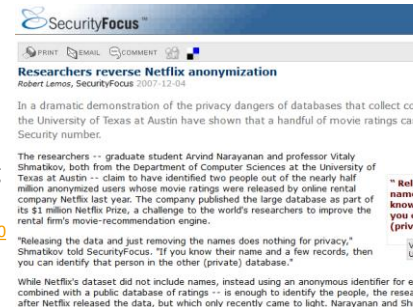
[DEIC - UAB] G. Navarro-Arribas

13

13

Example: Netflix

- Netflix:
 - \$1M prize for 10% improvement of its recommendation system.
 - Released an *anonymized* data to be used as training set.
- Researchers could re-identify users by linking the anonymized training data to IMDB reviews!
- https://www.cs.utexas.edu/~shmat/shmat_oak08netflix.pdf



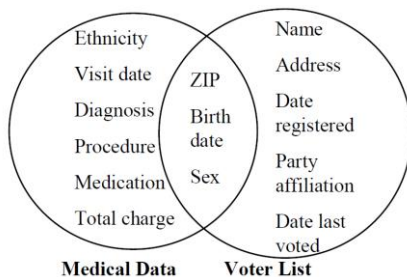
2021-05-02

[DEIC - UAB] G. Navarro-Arribas

14

14

Example: US Census and Medical data



Researcher found that:

- date of birth
- gender
- 5-digit ZIP

Uniquely identifies 87.1% of USA population!!!

<https://dataprivacylab.org/projects/identifiability/paper1.pdf>

Link two public anonymized datasets

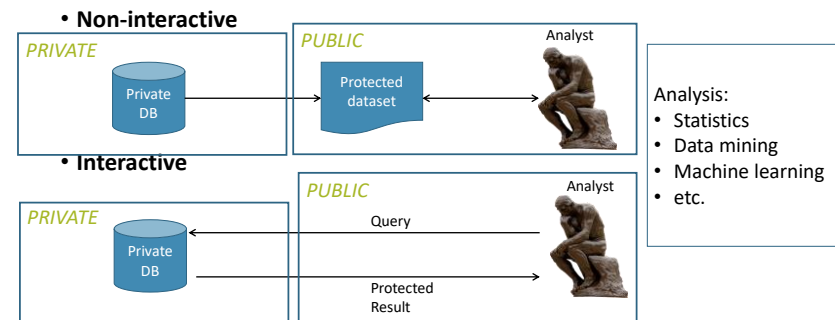
2021-05-02

[DEIC - UAB] G. Navarro-Arribas

15

15

Privacy Protection Scenarios



2021-05-02

[DEIC - UAB] G. Navarro-Arribas

16

16

Privacy Models

- k-anonymity
- Differential Privacy

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

17

17

Type of attributes related to an entity

- **Identifiers**: unambiguously identify the entity (passport number, full name, etc.)
 - Usually removed or encrypted for public release.
- **Quasi-identifiers**: in combination, can be linked with external information to reidentify an entity (age, zip code, city, ...)
 - Usually a protection mechanism is applied to these attributes.
- **Confidential**: contain sensitive information about the entity (salary, religion, political affiliation, health condition, etc.)
 - Usually not modified (object of study).

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

18

18

Example

Identifiers		Quasi-identifiers		Confidential	
Name	Passport	Profession	City	Age	Salary
Alfredo Pareja	AS3682	Teacher	Guayaquil	23	600
Eugenio Espejo	RX3453	Teacher	Quito	48	500
Raul Perez	TS5645	Policeman	Loja	20	734
Francisco Proaño	XB3456	Writer	Guayaquil	32	543
Cesar Davila	AE4324	Teacher	Cuenca	38	890
Ignacio Laso	FF3455	Writer	Quito	41	678
Raul Serrano	LO0903	Policeman	Portoviejo	21	399
Luz Argentina	KO8264	Policeman	Quevedo	44	943
Abdon Ubidia	RC3766	Teacher	Machala	31	632

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

19

19

Definition of k-anonymity

A dataset X satisfies **k-anonymity** with respect to a set of quasi-identifiers when the projection of X in this set results into a partition of X in sets of at least k indistinguishable records

- There are **k** indistinguishable items
- An item (individual) cannot be distinguished from other **k-1** items from the dataset.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

20

20

Example of k-anonymity

City	Age	Salary
Guayaquil	23	600
Quito	48	500
Loja	20	734
Guayaquil	32	543
Cuenca	38	890
Quito	41	678
Portoviejo	21	399
Quevedo	44	943
Machala	31	632

3-anonymous table

City	Age	Salary
(Ecuador)	[20, 29]	600
(Ecuador)	[40, 49]	500
(Ecuador)	[20, 29]	734
(Ecuador)	[30, 39]	543
(Ecuador)	[30, 39]	890
(Ecuador)	[40, 49]	678
(Ecuador)	[20, 29]	399
(Ecuador)	[40, 49]	943
(Ecuador)	[30, 39]	632

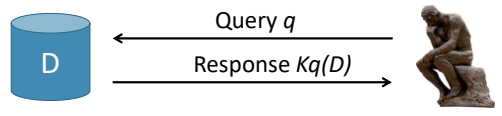
Quasi-identifiers: City, Age
Confidential: Salary

2021-05-02 [DEIC - UAB] G. Navarro-Arribas 21

21

Differential Privacy

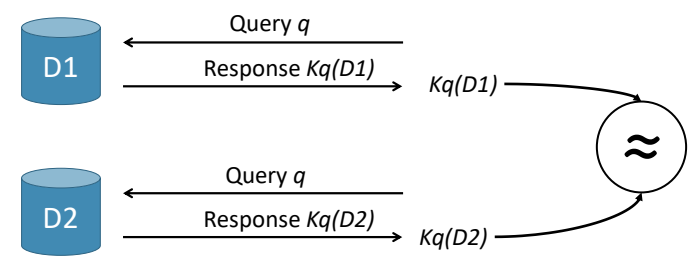
A function K_q for a query q gives ϵ -differential privacy if for all data sets D_1 and D_2 differing in at most one element, and all $S \subseteq \text{Range}(K_q)$,

$$\frac{\Pr[K_q(D_1) \in S]}{\Pr[K_q(D_2) \in S]} \leq e^\epsilon$$


2021-05-02 [DEIC - UAB] G. Navarro-Arribas 22

22

(Very) broadly speaking...



2021-05-02 [DEIC - UAB] G. Navarro-Arribas 23

23

Differential privacy

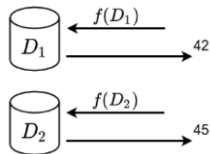
- For any D_1, D_2 : two versions of the database D that differ in at most one item (element, record, ...)
- If we perform the same query to both, the result should be approximately the same.
- ϵ is the level of privacy (smaller $\epsilon \Rightarrow$ greater privacy). E.g. for values:
 - 0: total privacy (both queries return the same answer)
 - 0.01: both probabilities differ by <1%
 - 1: both probabilities differ by e

2021-05-02 [DEIC - UAB] G. Navarro-Arribas 24

24

Differential privacy example

Table 1: D_1			Table 2: D_2		
nombre	edad	salario	nombre	edad	salario
Ataúlfo	43	100	Ataúlfo	43	100
Sigerico	15	300	Walia	50	200
Walia	50	200	Teodoro	51	100
Teodoro	51	100	Turismundo	35	500
Turismundo	35	500	Tedrico II	40	300
Tedrico II	40	300	Eurico	44	700
Eurico	44	700	Alarico II	49	300
Alarico II	49	300	Gesaleico	69	500
Gesaleico	69	500	Amalarico	24	200
Amalarico	24	200			



2021-05-02

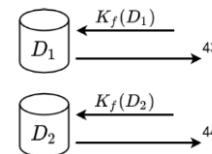
[DEIC - UAB] G. Navarro-Arribas

25

25

Cont.

- Imagine we have $Kf = f + \text{noise}$. Noise is a number randomly chosen from $\{-2, -1, 0, 1, 2\}$.



2021-05-02

[DEIC - UAB] G. Navarro-Arribas

26

26

Privacy budget in differential privacy

- In differential privacy the loss of privacy is accumulative
 - E.g. the same query to the same DB:
 - One query: 1/5 probability of guessing the real value.

Consulta	$K_f(D_1)$	$f(D_1)$?
1	43	41, 42, 43, 44, 45
2	40	38, 39, 40, 41, 42

- Now, we know the real values is 41 or 42, prob: 1/2
- We have a privacy budget which limits the number of queries!!

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

27

27

Protection
methods

Original Data

X

Protected Data

X'

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

28

28

Protection methods classification

- **Perturbative:** the original data X is distorted. Protected data X' contain some erroneous information.
- **Non-perturbative:** replace original values X by others which are less specific. Protected data X' do not contain erroneous information.
- **Synthetic data generators:** new artificial data is generated to substitute the original data.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

29

29

Some protection methods

We will see some examples of common protection methods

- Perturbative
 - Rank Swapping
 - Microaggregation
 - Additive and multiplicative Noise
- Non-perturbative
 - Generalization

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

30

30

Rank Swapping

- For each single variable V :
 1. Sort values $V = (a_1, \dots, a_n)$, such that $a_i \leq a_j$ for all $1 \leq i \leq j \leq n$
 2. Swap each value a_i with a_l , where l is randomly chosen from $[i + 1, \min(n, i + p * |X|/100)]$
 3. Undo the sorting step (1)

- Each value is swapped with another one randomly chosen within a restricted range, determined by p .
- p is the percentage from the total number of values.

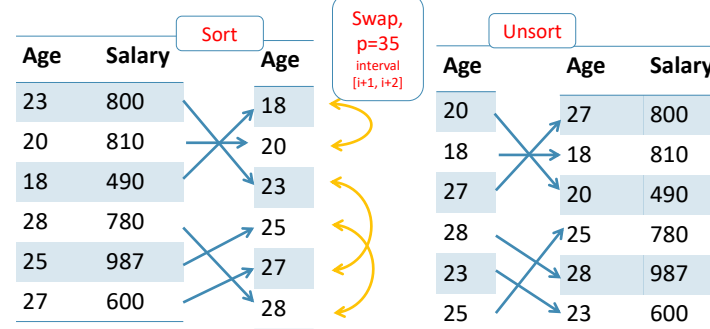
2021-05-02

[DEIC - UAB] G. Navarro-Arribas

31

31

Rank swapping: example



2021-05-02

[DEIC - UAB] G. Navarro-Arribas

32

32

Microaggregation

- Build small microclusters and replace each original data by its cluster representative.
- Usual method:
 - Partition: group similar records into cluster of at least k records
 - Aggregation: for each compute the cluster representative (centroid)
 - Replacement: each record is replaced by the cluster representative.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

33

33

Microaggregation: example

Age	Salary	Partition $k=2$		Aggregation, replacement	
Age	Salary	Age	Age	Age	Salary
23	800	23	24	24	800
20	810	20	19	19	810
18	490	18	19	19	490
28	780	28	27.5	27.5	780
25	987	25	24	24	987
27	600	27	27.5	27.5	600

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

34

34

Microaggregation: some notes

- Can be applied to multiple variables at the same time (multivariate).
- Can be used to achieve k -anonymity

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

35

35

Microaggregation: example (multivariate)

Weight	Age	Pathogen	Microaggregation $k=3$ applied to quasi- identifiers: Weight and Age			Weight	Age	Pathogen
70	33	Virus				74.3	36	Virus
84	40	Bacteria				74.3	36	Bacteria
59	15	Virus				60.6	20	Virus
69	35	Protozoa				74.3	36	Protozoa
63	25	Virus				60.6	20	Virus
60	20	Bacteria				60.6	20	Bacteria

3-anonymous table

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

36

36

Additive and Multiplicative Noise

- **Additive Noise:** add noise to the original data:

$$X' = X + \varepsilon$$

- ε is noise, following a certain distribution

- **Multiplicative Noise:** multiply noise to the original data:

$$X' = X * \varepsilon$$

- ε is noise, following a certain distribution

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

37

37

Additive Noise: example

Weight	Age	Pathogen
70	33	Virus
84	40	Bacteria
69	35	Protozoa
59	15	Virus
63	25	Virus
60	20	Bacteria

$$X' = X + N(0, p \sigma^2(X))$$

Noise as a Normal Distribution with mean $\mu = 0$, standard deviation $\sigma^2 = p \sigma^2(X)$, with parameter $p=0.2$
Applied to Weight and Age

Weight	Age	Pathogen
71	36	Virus
81	40	Bacteria
68	30	Protozoa
58	16	Virus
62	29	Virus
59	23	Bacteria

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

38

38

Additive noise: example comments

- Bigger p , means more distortion, more privacy.
- This example uses **uncorrelated** noise:
 - Preserves means and covariances from the original data
- Multiplicative noise with a Laplace distribution is commonly used to provide differential-privacy.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

39

39

Generalization

- Generalization of values following a given scheme or hierarchy.

Age	ZIP	Profession
30	08193	Nurse
36	08176	Paramedic
32	08191	Veterinarian
25	08034	Mathematician
29	08022	Physics
27	08010	Engineer

Age	ZIP	Profession
[30, 40)	081**	Health prof.
[30, 40)	081**	Health prof.
[30, 40)	081**	Health prof.
[20, 30)	080**	Science prof.
[20, 30)	080**	Science prof.
[20, 30)	080**	Science prof.

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

40

40

Evaluation of protection methods

- **Privacy:** level of privacy or anonymity achieved by the protection.
 - Privacy model: does it guarantees some privacy model? (k-anonymity, differential-privacy)
 - Reidentification: simulate an attack. How many users can the attacker re-identify?

Usually greater privacy implies greater distortion or generalization of the original data.

- We need to measure how data is deteriorated

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

41

41

Utility and Information loss

Information loss (IL) determines the information lost due to the anonymization.

$$IL(X, X') = \text{divergence}(X, X')$$

divergence is a way to compare both datasets X, and X'

- If X is very similar to X': lower IL
- If X is very different from X': higher IL

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

42

42

Information loss: Example

- **Mean Square Error (MSE):** [for numerical data] consider the dataset as a matrix M and the protected version as M', and c(M) is the number of elements in the matrix M, then:

$$\text{divergence1}(M, M') = \text{MSE}(M, M') = \frac{\sum_{ij}(M_{ij} - M'_{ij})^2}{c(M)}$$

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

43

43

Information Loss: example MSE

X = M		X' = M'	
30.0	2.0	29.0	2.5
22.0	2.0	23.5	1.5
25.0	1.0	23.5	1.5
28.0	3.0	29.0	2.5

$$\text{MSE}(M, M') = \frac{\sum_{ij}(M_{ij} - M'_{ij})^2}{c(M)}$$

Microaggregation k=2

$$= ((30.0 - 29.0)^2 + (2.0 - 2.5)^2 + (22.0 - 23.5)^2 + (2.0 - 1.5)^2 + (25.0 - 23.5)^2 + (1.0 - 1.5)^2 + (28.0 - 29.0)^2 + (3.0 - 2.5)^2) / 8 = \mathbf{0.9375}$$

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

44

44

Information Loss: example MSE 2

X = M		X' = M'	
30.0	2.0	26.5	2.0
22.0	2.0	26.5	2.0
25.0	1.0	26.5	2.0
28.0	3.0	26.5	2.0

$$MSE(M, M') = \frac{\sum_{ij}(M_{ij} - M'_{ij})^2}{c(M)}$$

Microaggregation k=4

$$= ((30.0 - 29.0)^2 + (2.0 - 2.5)^2 + (22.0 - 23.5)^2 + (2.0 - 1.5)^2 + (25.0 - 23.5)^2 + (1.0 - 1.5)^2 + (28.0 - 29.0)^2 + (3.0 - 2.5)^2) / 8 = \mathbf{4.875}$$

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

45

45

Information loss: example

- Other generic divergence measures can be used:
 - mean absolute error, mean relative error,
 - different between several static characteristics: between correlation matrices, etc.
 - They can be combined (aggregated)
- Specific measures for machine learning
 - Us a specific machine learning model (decision tree, clustering, regression, ...) and compare the performance of the model using training data from the original data X and the protected data X'

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

46

46

Guillermo Navarro-Arribas

2021-05-02

[DEIC - UAB] G. Navarro-Arribas

47

47