

Big Data Privacy

Open DataBase

Deterministic database:

→ Database with concrete information

pseudo	job	sex	age	discrete
U1	engineer	m	35	high
U2	engineer	m	38	high
U3	lawyer	f	38	low
U4	writer	f	30	low
U5	writer	f	30	high
U6	dancer	f	30	high
U7	dancer	f	30	high

Issue: if $f_1 \rightarrow (f_2, f_3, f_4)$
we know she has HV

b-anonymity

$$k := \min_{o \in O} |A(o)|$$

number of users having the
same sensitive characteristics &

→ can be done by:

- Aggregating data
- Grouping users
- Noise addition

pseudo	job	sex	age	discrete
U1	professor	m	35-40	high
U2	professor	m	35-40	high
U3	professor	m	35-40	high
U4	artist	f	30-35	low
U5	artist	f	30-35	low
U6	artist	f	30-35	low
U7	artist	f	30-35	low

Issue: lack of diversity in sensible

l-diversity

- sensible attributes are "well distributed"
- Ideal: uniformly
- $l = \# \text{ of different sensible in each group}$



t-closeness

The distribution of sensitive attributes in an anonymity set is "close" to the distribution of the full table

Probabilistic database

→ can capture real world complexity
 $A(o) = 2 \neq 1$

pseudo	comedy	drama
U1	—	—
U2	68 (true 60%)	35 (true 40%)

→ If we have enough observations we
can guess who is the user
→ How much / observation?

Hypothetical Information:

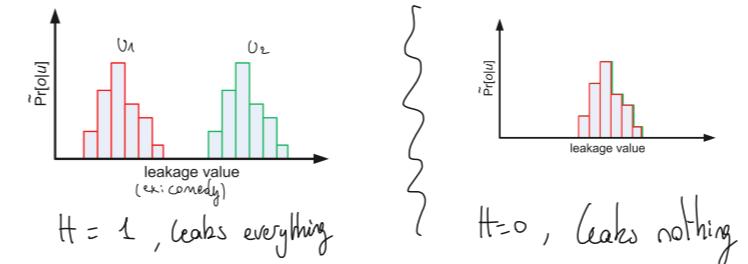
Idea: → How much information 1 observation looks
→ Based on internal structure
→ If obs follows this prob model, how much info is
there? Would it be a leakage carry?

$$HI(U; O) = H[U] + \sum_{u \in U} \Pr[u] \cdot \sum_{o \in DB} \Pr[u|o] \cdot \log_2 \Pr[u|o]$$

with $\Pr[X = x] := \Pr[x]$ & $H[U] = -\sum_u \Pr[u] \cdot \log_2 \Pr[u]$

→ Re-Identification Success Rate: $1 - (1 - H)^n$

Two examples:



Perceived Information

→ important bc one might act on
his beliefs

→ Given what the attacker believes how
much info do they think they gained

→ even if they actually did not gain much knowledge

PC might be high → simulates attacker's model

• Split the DB in two parts [DB₁ and DB₂] → simulates the leakage

• Where p stands for profiling and t for testing

• Build a model $\Pr[a|u] \leftarrow DB$, for every user

$$\Pr(U; O) = H[U] + \sum_{u \in U} \Pr[u] \cdot \sum_{a \in A} \Pr[a|u] \cdot \log_2 \Pr[a|u]$$

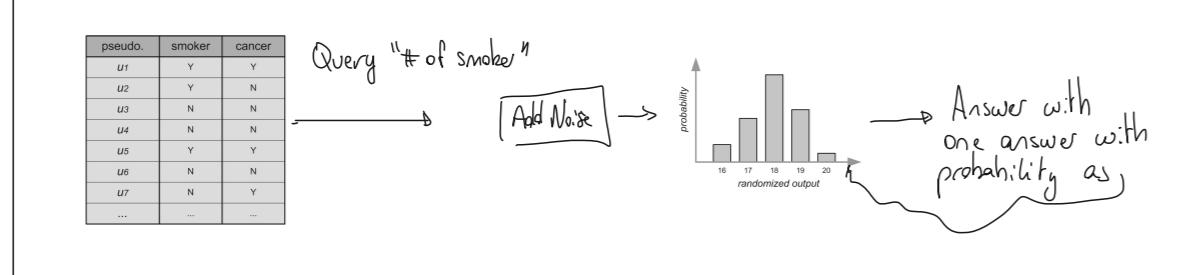
$$\Pr(U; O) = H[U] + \sum_{u \in U} \Pr[u] \cdot \sum_{a \in A} \frac{1}{|DB|} \cdot \log_2 \Pr[a|u]$$

↓ for every u
it's been profiling the user

Trusted DB with sanitized answer

Sanitization:

→ We will add some noise to the answers of
the queries:



Differential Privacy

→ Idea: "For two neighbor DB's (= 1 person has been added or removed) the results are indistinguishable"

→ Definition: ϵ -DP for Mechanism (= query) M if for
two neighbor DB's D_1, D_2 and $c \in \text{range}(M)$ then:

$$\Pr[M(O_1) = c] \leq (1+\epsilon) \Pr[M(O_2) = c] \quad (\text{are mult. close})$$

| (to deal with $\Pr = 0$)

$$(\epsilon - \delta) - DP \quad \Pr[M(O_1) = c] \leq (1+\epsilon) \Pr[M(O_2) = c] + \delta$$

→ Sensitivity: $\Delta f = \max_{D_1, D_2} |f(D_1) - f(D_2)|$

→ Laplacian noise: $b = \frac{\Delta f}{\epsilon}$ achieves ϵ -DP

→ Note: if $\epsilon_1 - \epsilon_2$ DP followed by $\epsilon_1 - \epsilon_2$ then
 $(\epsilon_1 + \epsilon_2)$ DP which grows linearly

• Often the noise needed for privacy will make
the data not useful

Untrusted DB

* PGC is better (as seen in previous chapter)

* Fully Homomorphic encryption:

• idea: Decipher Encrypted - for m

• We need additive homomorphism (easy)

• the following multiplicative homomorphism which has

Problem 1: (basic homomorphic) decryption has to
know the coefficients of a quadratic polynomial in the
variables of x \Rightarrow quadratic blowup of ciphertext size

Problem 2: If $\sigma^2(c) = \epsilon$, how much is $\sigma^2(c \cdot c')$?
 \Rightarrow ciphertexts are rapidly too noisy to be decrypted
(But are highly technical and so far extremely expensive)