PRVACY ENHANCING TECHNOLOGIES

Privacy-Preserving Data Computations and Database Privacy

2022-02-03

JENNI REUBEN

Acknowledgement: Some of the slides in this set are adaptations of lecture slides of Dr. Olaf Hartig (Linköping University).





















- Hard Privacy
 - avoid or reduce as much as possible placing any trust in the parties
- Soft Privacy
 - but rather place some trust on the data collector. Latter the trust can be challenged.

takes the assumption that we lose control of the data therefore there is no choice





- Smart metering systems
 - the goal is to match the production and consumption
 - to remotely read fine-grained measurements from each smart meters, which enable the grid operators to balance load efficiently and offer adapted time-dependent tariffs.
- fine granular measurements are privacy invasive,
 - simulated attacks have shown to detect from the smart meter data, the presence/ absence of residents in a household.
- The utility provider wants to perform analysis for grid management and billing,
 - This is achieved by secure-multi party computations and homomorphic encryption





- Homomorphic encryption
 - It is a type of encryption that allows the receiver of the cipher text (the
- Secure multi-party computations
 - involved in the protocol.

encrypted smart meter data measurements) to compute an operation on these encrypted values like adding the daily fees without having to decrypt them.

It is a protocol that allows several parties to perform a common computation on their individual values without disclosing their respective values to the others







SMART GRID SETUP



2022-02-3

JENNI REUBEN

Centralized Smart Grid Architecture

FN Swedish Defence Research Agency





- Centralized Setup
 - communicates with each smart meter.
 - users may access to the stored data to get information about her consumption.
- **De-centralized Setup**
 - and billing on the metered data are distributed among consumers.
 - billing period and communicate to the energy suppliers.
 - Grid management and load balancing are performed collaboratively by the users.

• the smart meters sends measurements of short slot intervals to a central data storage that acts as a hub and

The aggregator database is then used for consumption calculation, load balancing calculation and billing. The

the smart meters play the role of an aggregator, the calculations such as total consumption, load balancing

• The meters perform a partial data aggregation themselves, calculating the total energy consumption for each



















Semi honest adversary

- Malicious

- Someone who deviates from the protocol, forging the interchanged messages to gain more information or to alter the output of the protocol.

Two – way trust relationship requirement

The customer disclose their private data to the energy supplier/grid operator and trust that the operator use it only for legitimate purpose. Conversely the trust relationship of the aggregator/operator focuses on the correctness of the data that the meters provide, so the customers provide a non-forged consumption values.

- Someone who follows the protocol but may try to infer information from the interchanged values.









Privacy preserving computation of total computation of a cell with three users Alice, Bob, Charlie along with a utility company UC.

1. Alice, Bob, Charlie splitting their measurements into random shares (3 shares in specific one share for each person)

2. Alice keeps her share and sends the other two to the utility company encrypting with Bob's and Charlie's public key. Bob and Charlie does the same

3. The UC using additive homomorphic encryption adds the shares that are encrypted using the same public key

4. The UC sends the result of the addition to the members, which is also encrypted

5. Upon receiving the summation of the shares from other parties, Alice decrypts it and adds her own share. The result of the addition is sent to the UC

6. The UC receiving all the sums of the shares from all the parties, compute the final addition to get the total consumption

2022-02-3 JENNI REUBEN







Respondent Privacy

corresponds to

<u>Owner Privacy</u>

query

End-user Privacy

Protecting the information of the individuals to which the records in a database

Protecting the information of each entities that are coming together for computing a

Protecting end-user's queries to an interactive databases such as search engines.









- the database represents
- exploited for variety of reasons such as disease control, market research, medical research
- we should be interested in the public availability of such data: results from such data can contribute to expanding our knowledge about diseases
- However, those datasets contain confidential information about the respondents who have given their information to the database
- Can the users (researchers, analysts) of such databases be trusted?

• enable its users to retrieve statistical knowledge from a subset of the population that







WHAT ARE THE RISKS? RECAP FROM LAST WEEK

- Anonymity in terms of unlinkability:
 - subject and this attribute [Pfitzmann17]

Two types of linkage from an adversary's perspective;

- the published data (that is presumably free of explicit identifiers)
- would have been possible without the access to the data.

The anonymity of a subject w.r.t an attribute may be defined as unlinkability of this

Record linkage: re-identify the individual that the records in the published database corresponds to, by linking the publicly available information to the information in

Attribute linkage: accurately infer the confidential attribute values of an individual or a set of individuals represented in the underlying database, such as inference







RECORD LINKAGE EXAMPLE

- In Massachusetts, USA, the Group Insurance Commission (GIC) is responsible for purchasing health insurance for state employees
- Sweeney paid \$20 to buy the voter registration list for Cambridge, MA
 - Former governor (William Weld) of MA lives in Cambridge, MA hence his record is in the Voters DB
 - 6 people in Voters DB shares his DOB
 - Of which only 3 of them were men
 - Of which only 1 record matches the Weld's ZIP code.
 - Mr. Weld's medical information learned!



[Fung10]









- Explicit Identifiers:
 - address, etc.
- Quasi Identifiers:
 - respondent. E.g., gender, age, telephone number, zip code etc.
- Sensitive attributes:
 - Attributes that contain sensitive information of the respondents. E.g., disease, salary. etc.
- Non-sensitive attributes:
 - > All other attributes that captures the respondents' non-sensitive information

Attributes that unambiguously identify the respondent. E.g., name, social security number, IP

> A set of non-sensitive attributes that when combined may lead to unambiguously identify the





THE CHALLENGE

- Statistical databases such as the databases of the U.S census Bureau contain confidential information such as age, sex, income, credit ratings, types of disease, etc.
- bow to publish statistics about the underlying population, which is based on their confidential attributes while not revealing anything about those individual. The privacy, utility trade-off
- We need a non-trivial way to limit the disclosure of confidential information
- Fact: 87% of the US population can be identified by the combination of ZIP, DOB and Sex.
- Statistical Disclosure Control (SDC) or Statistical Disclosure Limitation (SDL)
 - limits the disclosure of confidential information from the published statistics

2022-02-3 JENNI REUBEN







- - > X_{ii} is the value of the attribute *j* for respondent *i*.
- Non-perturbative
 - the true values of the respondents information.

 \triangleright Let X be a microdata, that is a $s \times t$ matrix, with s respondents and t attributes, then

 \triangleright Non-perturbative version of X is a modified version X', where X' is obtained from X by partial suppression or reduction of detail. The values represented in X' are





Perturbative

- affected.
- noise that is drawn from a distribution.
- Synthetic data generation

2022-02-3

 \triangleright Data perturbation: The perturbed version X' of X such that the X' preserves the statistical information of X, such that statistics computed on X' is not significantly

Query result perturbation: Queries are executed on the original datatable X, the results of the queries are perturbed by adding a calculated amount of random







K-ANONYMITY

- \triangleright A dataset or datable T is said to satisfy k-anonymity if each combination of values of the quasi-identifier attributes in T is shared by at least k-1 records.
- Let T be a table and X be a subset of the attributes of T. For every tuple t in T we write t[X] to denote the sequence of values that t has for the attributes in X.
- Example:
 - If $X = \{ZIP, Age, Sex\}$ and say t is the first tuple in T
 - then, t[X] is (12211, 18, M)
 - ▶ If $X = \{Z | P, Sex\}$, then t[X] is (12211, M)

2022-02-3 **JENNI REUBEN**

12211	18	М	Arthritis
12244	19	М	Cold
12245	27	М	Heart problem
12377	27	М	Flu
12377	27	F	Arthritis
12391	34	F	Diabetes
12391	45	F	Flu









we have $t[QI_T] = t1[QI_T] = t2[QI_T] = tk-1[QI_T]$.

12211	18	М	Arthritis	122**	18-19	М	Arthritis
12244	19	Μ	Cold	122**	18-19	М	Cold
12245	27	Μ	Heart problem	*	27	*	Heart problem
12377	27	Μ	Flu	*	27	*	Flu
12377	27	F	Arthritis	*	27	*	Arthritis
12391	34	F	Diabetes	12391	≥ 30	F	Diabetes
12391	45	F	Flu	12391	≥ 30	F	Flu

2022-02-3

JENNI REUBEN

Let T be a table and QI_T be the quasi-identifier of T. T satisfies k-anonymity if for every tuple t in T there exist (at least) k-1 other tuples t_1, t_2, \dots, t_{k-1} in T such that

2-anonymous table T*





Public Data

Chris	12211	18	М
Jack	19221	20	М

Chris	12211	18	М	Arthritis
Chris	12211	18	М	Cold

2022-02-3

JENNI REUBEN

Ol aroun / equivalence class	122**	18-19	М	Arthritis
Gi gioup / equivalence class	122**	18-19	М	Cold
	*	27	*	Heart problem
	*	27	*	Flu
	*	27	*	Arthritis
	12391	≥ 30	F	Diabetes
	12391	≥ 30	F	Flu
		2	-anony	mous table T

Chris is anonymous within his anonymity set







DATABASE RECONSTRUCTION ATTACK (DRA)

- released
- queries, one can recover the confidential data of the individuals in the underlying population.
- Take for example:
 - U.S census bureau database which contains answers given by the citizens of America

 - possible combinations that best fit the published statistics [Dinur03].

It turns out k-anonymity is not sufficient against inference attacks, so what if only aggregate data is

•But by simply observing the (perturbed or non-perturbed) the query answers/results of some random

• publishes statistics such as how many people belonging to a race, live in a particular block

The attack then is to guess using brute force computation, all the possible combinations of answers that people could have given to questions concerning race and block, and find out the









DATABASE RECONSTRUCTION ATTACK (DRA) EXAMPLE





JENNI REUBEN

Released Statistics

Mean Age	Median Age
30	38
30	33.5
51	48.5
51	53
35	35.6

Example taken from "Protecting privacy with math"









DATABASE RECONSTRUCTION ATTACK (DRA) CONT'D

Possible Ages for Mean 35 and Median 35.6

Female_ prof1	Female_prof2	Female_prof3
1	36	73
2	36	72
3	36	71
	•••	
6	36	68
	•••	
35	36	39
36	36	38

2022-02-3

JENNI REUBEN





DATABASE RECONSTRUCTION ATTACK (DRA) CONT'D

Possible Ages for Mean 35 and Median 35.6

Female_prof1	Female_prof 2	Female_prof 3	Female_prof 1	Female_prof 2	Female_prof 3
34	36	40	6	36	68
35	36	39	7	36	67
36	36	38	8	36	66



2022-02-3

JENNI REUBEN











	Count	Mean Age	Median Age
Total Population	7	30	38
Female	4	30	33.5
Ice cream lovers	4	51	48.5
Married Adults	4	51	53
Female Ice cream lovers	3	35	35.6

• Publishing less statistics, then there are little more plausible combinations of data that accurately fits the data

• Even lesser statistics are published, which increases the amount of data combinations that plausibly fit the released statistics.

JENNI REUBEN 2022-02-3







Observations from the above example,

- measure of loss of respondent privacy is the level of certainty in an attacker's ability in determining the plausibility of some possible combinations of data.
- Idea! to protect respondent privacy make all possible combinations of data from the respondents to be equally plausible.
- There is an inevitable trade-off between accuracy of the published results and not revealing information of the record owners in the underlying database.





A few possible data combinations are plausible



All possible data combinations are plausible







- How then to publish data for data analyses?
- query results' accuracy
- noisy results, which cancels out the noise.
- the cost of small loss in the accuracy of the results.

because increasing the uncertainty level of the adversaries, decreases the

Further, if random noise is added a bunch of times to a statistical query result, it is possible to get back the true results by taking the average of the

Differential privacy model that provides a strong privacy guarantee, yet at





The differential privacy model provides a way to **quantifies** the plausibility peak (i.e the loss of privacy) and **bounds** (that is to say the maximum) the loss of privacy for the individuals in the underlying dataset, as a consequence of publishing results computed on their data.



The plausibility/possibility plot with a few peaks that stands out

2022-02-3

JENNI REUBEN





•

DIFFERENTIAL PRIVACY EXAMPLE

Chris	Arthritis	
David	Cold	Q
Ethan	Heart problem	

be the same whether or not David is in the underlying database.

Observation:

- records are called database neighbors.
- getting answer 1 from D'.

2022-02-3 JENNI REUBEN



• Statistical Query: How many persons with a cold?, the answers from a differentially private computation will "nearly"

• The two databases where one contains David's data and the other do not contain his data - database neighbors. Generally speaking, any two databases D and D', which differ by at most one record but otherwise contain the same

• The results of the query over D and D' doesn't look the same, what it means here is that the probability distributions of the query result are the same. So, the likelihood of getting answer 1 when database is D is the same likelihood for





- Differential Privacy [Dwork06]:
 - A randomized query mechanism M_O for query Q provides ε -differential privacy if if for all databases D and D', where D and D' are database neighbors and • every subset O of the set of all possible outputs of M_O ,
- - We have that:

 $Pr[M_{Q}(D) \text{ in } O] \leq e^{\varepsilon} \cdot Pr[M_{Q}(D') \text{ in } O]$





Observation:

- Epsilon is the measure of peak that stand out in the plausibility plot (is the measure of information gain in adversaries ability to confidently choose one combination of data over the other), and the above definition bounds the loss of privacy from releasing the query results.
- **Composition** The future releases also guarantee ε -differential privacy

• if we publish the count of persons with cold with ε = 3 and publish the average age of persons with ε = 3, then the total privacy loss caused from the release of the two statistics is at most 6.





- Assume a query Q whose result Q(D) over any possible database instance D is a real number
- Randomized query mechanism M(Q) for Q, adds randomly selected noise η

• $M(Q) = Q(D) + \eta$

- Observation : the amount of noise depends both on ε and the sensitivity of the query being asked.
- The sensitivity of the query is a constant that captures the amount of maximum change any one individual may cause to the result of the query. Take our "how many persons with cold example, adding or removing a record will change the query result by at most a factor of 1.
- Less the epsilon, stronger the privacy







Definition: The sensitivity of a query Q is

$\Delta q = max \left| Q(D) - Q(D') \right|$ for any two neighboring databases D and D'

Examples:

• Δq for "count all patients diagnosed with cold" is: 1

2022-02-3

JENNI REUBEN





LAPLACE MECHANISM TO DIFFERENTIAL PRIVACY

- Idea: The noise to be added is drawn from the Laplace distribution Lap(λ), λ determines how flat the curve of the distribution is, from where the noise is drawn.
- Theorem [Dwork 2006]: Let M_Q be a mechanism for Q that returns $Q(D) + \eta$ where η is drawn randomly from Lap(λ) with $\lambda = \Delta q / \varepsilon$. M_O provides ε -differential privacy





Laplace distributions of varying scales from 1 to 4 the scale of the distribution depends on epsilon and Δq

Picture sources: https://commons.wikimedia.org/wiki/File:Laplace-verteilung.svg





Observations

- accuracy.
- However, for $\Delta q = 1$ and $\epsilon = 0.1$, we have $\lambda = 10$ (and $\lambda = 100$ if $\epsilon = 0.01$)
- Hence, for queries with higher sensitivity Δq , we have a higher value of λ thus, the noise n will typically be higher
- \blacktriangleright Likewise, for a smaller value of ε , the noise will be typically higher



The narrow the curve (Laplace distribution), the value drawn as noise is small, which implies the result of the query is changed by a small amount, narrow curve is good for





LAPLACE MECHANISM TO DIFFERENTIAL PRIVACY

- Given a sequence $Q_1, \dots, Q_m, \varepsilon$ -differential privacy can be achieved by drawing the noise for Q_m from $Lap(\lambda_m)$ where λ_m is the sum of all $\lambda_i = \Delta q_i / \varepsilon$ ($i = 1, \dots, m$)
 - Observation: The magnitude of the amount of noise added increases with every query.

• Theorem [Dwork 2006]: Let M_Q be a mechanism for Q that returns $Q(D) + \eta^k$ where η^k is a vector of size k whose elements are independently drawn randomly from Lap(λ) with $\lambda = \Delta q / \varepsilon$. M_Q provides ε –differential privacy









REFERENCES

[Pfitzmann17] - A Terminology for Talking about Privacy by Data Minimization, 2017 [Westin70] – A. Westin, Privacy and Freedom. Atheneum, New York, 170 [Dwork06] - C. Dwork, F. McSherry, K. Nissim, A. Smith, Calibrating Noise to Sensitivity in Private Data Analysis, 2006

[Dwork13] - C. Dwork, A. Roth, The Algorithmic Foundations of Differential Privacy, 2013

[Warren1890] - S. Warren, L. Brandeis, The right to privacy, Harvard Law Review, 1890]

handbook (3rd ed), 2017.

[Agre98] – P. Agre, M. Rotenberg, Technology and Privacy, 1998

SIGMOD-SIGACT-SIGART symposium on Principles of database systems, 2003

[Mcsherry07] - F. McSherry, K. Talware Mechanism Design via Differential Privacy

[Fung10] Privacy-Preserving Data Publishing: A Survey of Recent Developments

JENNI REUBEN 2022-02-3

- [Sim17] S.Fischer-Hübner, S. Berthold, Privacy-Enhancing Technologies, in Computer and Information Security
- [Dinur03] I. Dinur, K. Nissim, Revealing information while preserving privacy, in Proceedings of the 22nd ACM



