Advertising systems in social media (3)

Mainack Mondal

CS 60017 Autumn 2021



The story so far ...

- Social advertising systems
 - Why bother about them?
 - The curious case of Facebook ads
 - How can we leverage these systems for doing good
- Abuse of the advertising systems
 - Why is targeted advertising bad?
 - Privacy risks with PII based targeting

The story so far ...

- Social advertising systems
 - Why bother about them?
 - The curious case of Facebook ads
 - How can we leverage these systems for doing good
- Abuse of the advertising systems
 - Why is targeted advertising bad?
 - Privacy risks with PII based targeting
- Now, how to prevent abuse of advertising systems and provide data privacy?

Preserving privacy of social data

- Two broad dimensions
 - Preserving privacy from the background actors, e.g., advertisers or even the social media platform
 - Preserving privacy of data from other users, e.g., your ex

Preserving privacy from background actors

What are we going to talk about?

- Mechanisms for hiding privacy sensitive attributes in databases
 - K-anonymity
 - Differential privacy

- Slides heavily borrowed from
 - Vitaly Smatikov from Cornell
 - Li Xiong from Emory

Public Data Conundrum

- Health-care datasets
 - Clinical studies, hospital discharge databases ...
- Genetic datasets
 - \$1000 genome, HapMap, deCode ...
- Demographic datasets
 - U.S. Census Bureau, sociology studies ...
- Search logs, recommender systems, social networks, blogs ...
 - AOL search data, social networks of blogging sites, Netflix movie ratings, Amazon ...

What About Privacy?

- First thought: anonymize the data
- How?
- Remove "personally identifying information" (PII)
 - Name, Social Security number, phone number, email, address... what else?
 - Anything that identifies the person directly
- Is this enough?

Re-identification by Linking

网络古马科学校 化合金 医外部的 化过程的 化过程的 医布莱尔氏 化合金 医内部的 化过程的 化过程的 医布莱斯氏试验 化合金 医外部的 化过程的 医外外的 化化合金 医子科学校 化合金

Microdata

ID	QID			SA
Name	Zipcode	Age	Sex	Disease
Alice 🔇	47677	29		Ovarian Cancer
Betty	47602	22	F	Ovarian Cancer
Charles	47678	27	М	Prostate Cancer
David	47905	43	М	Flu
Emily	47909	52	F	Heart Disease
Fred	47906	47	М	Heart Disease

Voter registration data

2012/11/14

Name	Zipcode	Age	Sex
Alice 🤇	47677	29	F
Bob	47983	65	М
Carol	47677	22	F
Dan	47532	23	М
Ellen 46789		43	F

Latanya Sweeney's Attack (1997)

Massachusetts hospital discharge dataset

SSN	Name	vnicity	Date Of Birth	Sex	ZIP	Marital Status	Problem
			09/27/64	female	02139	divorced	hypertension
	8		09/30/64	female	02139	divorced	obesity
		asian	04/18/64	male	02139	married	chest pain
	ŝ	asian	04/15/64	male	02139	married	obesity
	8	black	03/13/63	male	02138	married	hypertension
		black	03/18/63	male	02138	married	shortness of breath
	£.	black	09/13/64	female	02141	married	shortness of breath
		black	09/07/64	female	02141	married	obesity
	\$1	white	05/14/61	male	02138	single	chest pain
	Q	white	05/08/61	male	02138	single	obesity
		white	09/15/61	female	02142	widow	shortness of breath

Voter List

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

Figure *k* -dentifying anonymous data by linking to external data

Public voter dataset

Quasi-Identifiers

Key attributes

- Name, address, phone number uniquely identifying!
- Always removed before release
- Quasi-identifiers
 - (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

Classification of Attributes

Sensitive attributes

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

Key Attribute	Qı	uasi-identif	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

K-Anonymity: Intuition

- 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 man in the released table, but the only information you have is his birth date and gender. There are k men in the table with the same birth date and gender.

 Any quasi-identifier present in the released table must appear in at least k records

Generalization

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



Achieving k-Anonymity

Generalization

- Replace specific quasi-identifiers with less specific values until get k identical values
- Partition ordered-value domains into intervals

Example of a k-Anonymous Table

	Race	Birth	Gender	ZIP	Problem
t1	Black	1965	m	0214*	short breath
t2	Black	1965	m	0214*	chest pain
t3	Black	1965	Í	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 Ql={Race, Birth, Gender, ZIP}

At least two people With same attributes

At least two people QI = quasi identifier tuple

Curse of Dimensionality

[Aggarwal VLDB '05]

- Generalization fundamentally relies on spatial locality
 - Each record must have k close neighbors
- Real-world datasets are very sparse
 - Many attributes (dimensions)
 - Amazon customer records: several million dimensions
 - Not possible to create k close neighbors
- ◆ Projection to low dimensions loses all info ⇒ k-anonymized datasets are useless



Two (and a Half) Interpretations

 Membership disclosure: Attacker cannot tell that a given person in the dataset

- Sensitive attribute disclosure: Attacker cannot tell that a given person has a certain sensitive attribute
- Identity disclosure: Attacker cannot tell which record corresponds to a given person

This interpretation is correct, assuming the attacker does not know anything other than quasi-identifiers <u>But this does not imply any privacy!</u> Example: k clinical records, all HIV+

Attacks on k-Anonymity

k-Anonymity does not provide privacy if

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



k-Anonymity Considered Harmful

Syntactic

- 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 his target
- Relies on locality
 - Destroys utility of many real-world datasets

What are we going to talk about?

- Mechanisms for hiding privacy sensitive attributes in databases
 - K-anonymity
 - Differential privacy

- Slides heavily borrowed from
 - Vitaly Smatikov from Cornell
 - Li Xiong from Emory

Statistical Databases



Statistical Data Privacy

- Non-interactive vs interactive
- Privacy goal: individual is protected
- Utility goal: statistical information useful for analysis



- Promise: an individual will not be affected, adversely or otherwise, by allowing his/her data to be used in any study or analysis, no matter what other studies, datasets, or information sources, are available"
- Paradox: learning nothing about an individual while learning useful statistical information about a population

• Statistical outcome is indistinguishable regardless whether a particular user (record) is included in the data



• Statistical outcome is indistinguishable regardless whether a particular user (record) is included in the data



Differential privacy: an example



Original records

Original histogram

Perturbed histogram with differential privacy



Why all pairs of datasets ...?



Guarantee holds no matter what the other records are.

Why all outputs?

Should not be able to distinguish whether input was D_1 or D_2 no matter what the output



Privacy Parametere



Controls the degree to which D_1 and D_2 can be distinguished. Smaller the ε more the privacy (and better the utility)

Can deterministic algorithms satisfy differential privacy?

Output Randomization



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

[DMNS 06]

Laplace Mechanism



Laplace Distribution

• PDF:
$$f(x \mid \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$$

- Denoted as Lap(b) when u=0
- Mean u
- Variance 2b²



How much noise for privacy?

[Dwork et al., TCC 2006]

Sensitivity: Consider a query q: I → R. S(q) is the smallest number s.t. for any neighboring tables D, D',

$$|q(D) - q(D')| \leq S(q)$$

Theorem: If **sensitivity** of the query is **S**, then the algorithm $A(D) = q(D) + Lap(S(q)/\epsilon)$ guarantees ϵ -differential privacy

Sensitivity

- Semantically Sensitivity is
 - Given a query, what the maximum amount that the output will change by adding a row?

Example 1

- Let's consider a simple count query
 - Number of people clicking on an ad / having a disease?
 - What is the sensitivity?

Example: COUNT query

- Number of people having disease
- Sensitivity = 1

- Solution: 3 + η, where η is drawn from Lap(1/ε)
 - Mean = 0
 - Variance = $2/\epsilon^2$



Example 2

- Let's consider another count query
 - Number of people clicking on an ad / having a disease rounded to nearest multiple of 10?
 - What is the sensitivity?

Privacy of Laplace Mechanism

- Consider neighboring databases D and D'
- Consider some output O

$$\frac{\Pr\left[A(D)=O\right]}{\Pr\left[A(D')=O\right]} = \frac{\Pr\left[q(D)+\eta=O\right]}{\Pr\left[q(D')+\eta=O\right]} \qquad \lambda = \text{variance} = S(q)/\epsilon$$
$$= \frac{e^{-|O-q(D)|/\lambda}}{e^{-|O-q(D')|/\lambda}}$$

$$\leq e^{|q(D)-q(D')|/\lambda} \leq e^{S(q)/\lambda} = e^{\varepsilon}$$

Utility of Laplace Mechanism

- Laplace mechanism works for any function that returns a real number
- Error: E(true answer noisy answer)²
 = Var(Lap(S(q)/ε))
 = 2*S(q)² / ε²

- Where is there room for improvement?
 - The Laplace mechanism adds independent noise to every coordinate...
 - What happens if the user asks (essentially) the same question in every coordinate?
 - Read [Dinur,Nissim03]: a computationally efficient attack that gives blatant non-privacy for a mechanism that adds noise bounded by $o(\sqrt{n})$