

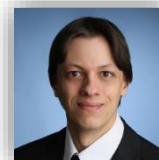
Ethics in Natural Language Processing – SS 2022



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Lecture 9 Privacy & Security II

Dr. Thomas Arnold
Aniket Pramanik



Ubiquitous Knowledge Processing Lab
Technische Universität Darmstadt

Slides and material from Yulia Tsvetkov



Carnegie Mellon University
Language Technologies Institute

Syllabus (tentative)

<u>Nr.</u>	<u>Lecture</u>
01	Introduction, Foundations I
02	Foundations II
03	Bias I
04	Bias II
05	Incivility and Hate Speech I
06	NO LECTURE – Christi Himmelfahrt
07	Incivility and Hate Speech II
08	Low-Resource NLP
09	NO LECTURE - Fronleichnam
10	Privacy and Security I
11	Privacy and Security II
12	Language of Manipulation I
13	Language of Manipulation II

Outline

Recap

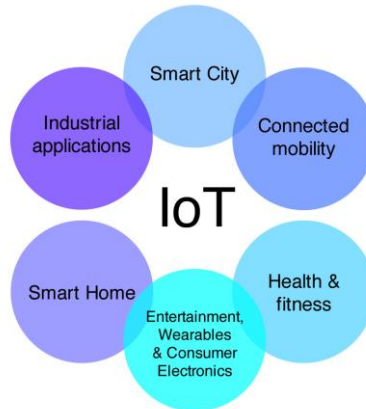
Anonymizing Data

Differential Privacy

Profiling



TSA Pre✓



What is Privacy

<https://en.wikipedia.org/wiki/Privacy> is the ability of an individual or group to seclude themselves or information about themselves, and thereby express themselves selectively

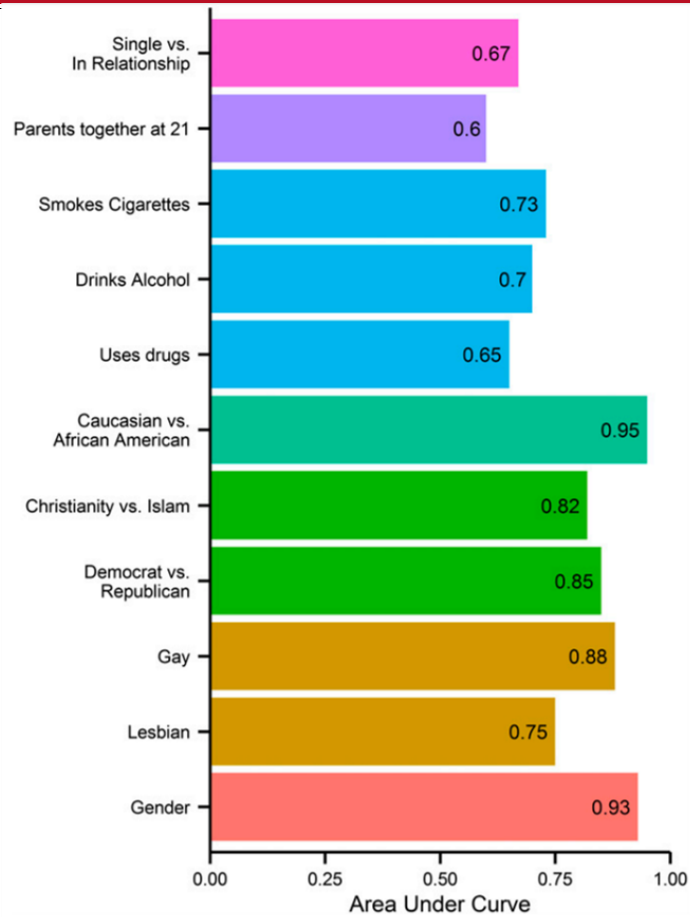


Fig. 2. Prediction accuracy of classification for dichotomous/dichotomized attributes expressed by the AUC.

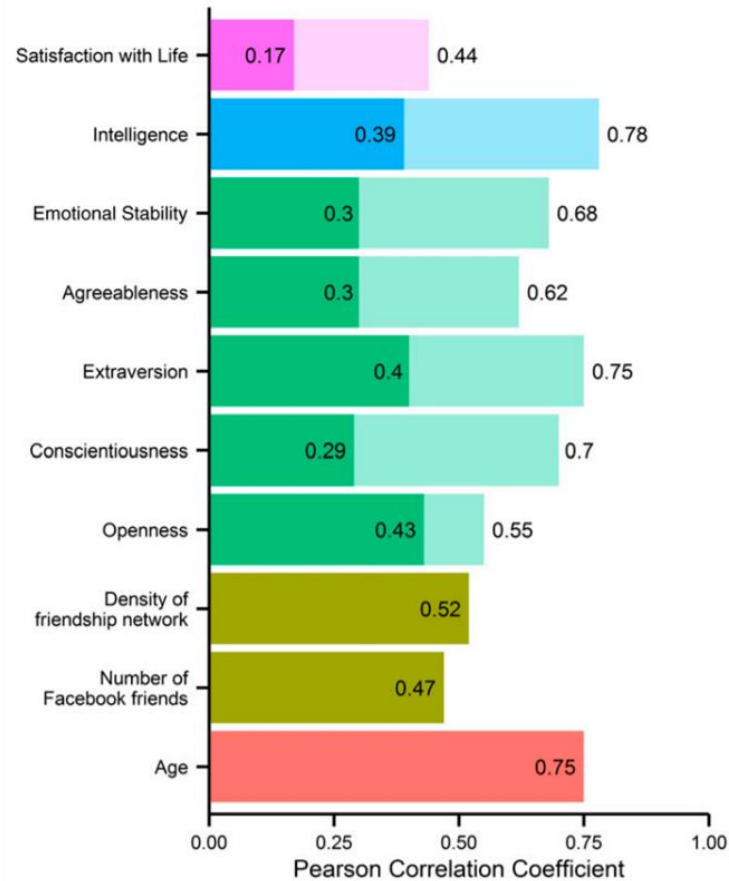


Fig. 3. Prediction accuracy of regression for numeric attributes and traits expressed by the Pearson correlation coefficient between predicted and actual attribute values; all correlations are significant at the $P < 0.001$ level. The

Privacy vs. Utility



high utility,
no privacy



high privacy,
no utility

Learning Goals

After hearing this lecture, you should be able to...

- determine if a dataset has a certain level of k-anonymity, l-diversity or t-closeness
- describe possible attacks on tables with k-anonymity
- explain a simple form of differential privacy, and calculate true distributions from distributions with random noise

Outline

Recap

Anonymizing Data

Differential Privacy

Privacy and Anonymity

- Being on-line without giving up everything about you
- Ensuring collected data doesn't reveal its users data
- Privacy in
 - Structured Data: k-anonymity, differential privacy
 - Text: obfuscating authorship
 - Speech: speaker id and de-identification

Companies Getting Your Data

- They actually don't want your data, they want to upsell
 - They want to be able to do tasks (recommendations)
 - They actually don't care about the individual you
- Can they process data to never have identifiable content?
 - Cumulated statistics
 - Averages, counts, for classes
- How many examples before it is anonymous?

k-anonymity

- Latanya Sweeney and Pierangela Samarati 1998
- Given some table for data with features and values
- Release data that guarantees individuals can't be identified
 - **Suppression:** Delete entries that are too “unique”
 - **Generalization:** relax specificity of fields,
e.g. age to age-range or city to region

k-anonymity

Name	Age	Gender	State of domicile	Religion	Disease
Ramsha	29	Female	Tamil Nadu	Hindu	Cancer
Yadu	24	Female	Kerala	Hindu	Viral infection
Salima	28	Female	Tamil Nadu	Muslim	TB
Sunny	27	Male	Karnataka	Parsi	No illness
Joan	24	Female	Kerala	Christian	Heart-related
Bahuksana	23	Male	Karnataka	Buddhist	TB
Rambha	19	Male	Kerala	Hindu	Cancer
Kishor	29	Male	Karnataka	Hindu	Heart-related
Johnson	17	Male	Kerala	Christian	Heart-related
John	19	Male	Kerala	Christian	Viral infection

- From wikipedia: K-anonymity

k-anonymity

Name	Age	Gender	State of domicile	Religion	Disease
*	$20 < \text{Age} \leq 30$	Female	Tamil Nadu	*	Cancer
*	$20 < \text{Age} \leq 30$	Female	Kerala	*	Viral infection
*	$20 < \text{Age} \leq 30$	Female	Tamil Nadu	*	TB
*	$20 < \text{Age} \leq 30$	Male	Karnataka	*	No illness
*	$20 < \text{Age} \leq 30$	Female	Kerala	*	Heart-related
*	$20 < \text{Age} \leq 30$	Male	Karnataka	*	TB
*	$\text{Age} \leq 20$	Male	Kerala	*	Cancer
*	$20 < \text{Age} \leq 30$	Male	Karnataka	*	Heart-related
*	$\text{Age} \leq 20$	Male	Kerala	*	Heart-related
*	$\text{Age} \leq 20$	Male	Kerala	*	Viral infection

- From wikipedia: K-anonymity

k-anonymity

A dataset has k -anonymity if the information for each person cannot be distinguished from at least $k - 1$ other individuals

- Optimal k -anonymity is an NP-hard problem

k-anonymity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

k-anonymity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...



Personal attributes



Sensitive attribute(s)

k-anonymity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Diagram illustrating k-anonymity with k=3. The table shows personal attributes (Name, Age, Gender, City, Religion, Crime). The first three rows (Parking Violation, Murder, Speeding) form an equivalence class of size 3, indicated by a bracket labeled '3'. The next two rows (Speeding, Robbery) form an equivalence class of size 2, indicated by a bracket labeled '2'.

Equivalence class: Entries that have the same personal attributes

k-anonymity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Diagram illustrating k-anonymity with k=2. The table is divided into two groups of rows, each indicated by a bracket and a number (3 and 2). The first group (3 rows) has a common value of 31 – 40 for Age, Male for Gender, and Griesheim for City. The second group (2 rows) has a common value of 31 – 40 for Age, Male for Gender, and Darmstadt for City. The Crime column shows different values for each row in both groups, demonstrating that the table is 2-anonymous.

If all **equivalence classes** are at least size 2, this table has 2-anonymity

k-anonymity

A dataset has k -anonymity if the information for each person cannot be distinguished from at least $k - 1$ other individuals

- Optimal k -anonymity is an NP-hard problem
- Homogeneity Attack: All sensitive values within a set can be identical
- Background Knowledge Attack: Association between one or more quasi-identifier attributes with the sensitive attribute

Homogeneity Attack

I am male, 39 and live in Griesheim. What was my crime?

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Homogeneity Attack

I am male, 39 and live in Griesheim. What was my crime?

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Background Knowledge Attack

I am male, 39 and live in Griesheim. I always go by bike. What was my crime?

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Background Knowledge Attack

I am male, 39 and live in Griesheim. I always go by bike. What was my crime?

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

very
unlikely

I-diversity

An equivalence class has l -diversity if there are at least l "well-represented" values for the sensitive attribute. A dataset has l -diversity if every equivalence class of the dataset has l -diversity.

An equivalence class has I-diversity if there are at least I "well-represented" values for the sensitive attribute. A dataset has I-diversity if every equivalence class of the dataset has I-diversity.

What means "well-represented" values?

- **Distinct I-diversity: At least I distinct values (simplest definition)**
- Entropy I-diversity: Calculates entropy of sensitive values (most complex)
- Recursive I-diversity: Compromise definition

(Distinct) I-diversity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

(Distinct) I-diversity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Diagram illustrating the concept of (Distinct) I-diversity. The table shows data rows with asterisks (*) indicating missing or unknown values. To the left of the table, two vertical brackets with arrows point to the first and fourth rows, each accompanied by a question mark (?), indicating the task of identifying distinct sensible values for each class.

How many distinct **sensible values** are in each eq. Class?

(Distinct) I-diversity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Diagram illustrating the grouping of rows for I-diversity calculation:

- Group 1 (indicated by a bracket and arrow labeled '1') contains the first three rows (all 'Murder' crimes).
- Group 2 (indicated by a bracket and arrow labeled '2') contains the next two rows ('Speeding' and 'Robbery' crimes).

This dataset has 1-diversity (lowest value)

(Distinct) I-diversity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

(Distinct) I-diversity

Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Diagram illustrating the concept of (Distinct) I-diversity. The table shows data rows grouped by the 'Age' column. The first three rows (Age 31 – 40, Male, Griesheim) are grouped together, and the next three rows (Age 31 – 40, Male, Darmstadt) are grouped together. The 'Crime' column shows distinct values for each group: Parking Violation, Murder, Speeding, Speeding, and Robbery. The diagram indicates that the first group has 3 distinct values and the second group has 2 distinct values, both of which are at least 2, satisfying the condition for (Distinct) I-diversity.

If all **equivalence classes** have at least 2 distinct sensible values,
this table has 2-diversity

t-closeness

An equivalence class has t -closeness if the distance between the distribution of the sensitive attribute in the class and the distribution of the attribute in the whole data set is no more than threshold t .

A dataset has t -closeness if every eq. class of the dataset has t -closeness.

t-closeness

An equivalence class has t -closeness if the distance between the distribution of the sensitive attribute in the class and the distribution of the attribute in the whole data set is no more than threshold t .

A dataset has t -closeness if every eq. class of the dataset has t -closeness.

t is a tradeoff between security and utility!

0.0-closeness: most secure, no utility

1.0-closeness: lowest security, highest utility

There are several ways to measure the distance between distributions. The easiest is the variational distance:

For two distributions $P = (p_1, p_2, \dots, p_m)$, $Q = (q_1, q_2, \dots, q_m)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

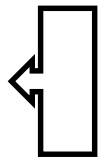
t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

t-closeness



Name	Age	Gender	City	Religion	Crime
...
*	31 – 40	Male	Griesheim	*	Parking Violation
*	31 – 40	Male	Griesheim	*	Murder
*	31 – 40	Male	Griesheim	*	Speeding
*	31 – 40	Male	Darmstadt	*	Speeding
*	31 – 40	Male	Darmstadt	*	Robbery
...

Distribution $Q = (0, 0, 0.5, 0.5)$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

Distribution $Q = (0, 0, 0.5, 0.5)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

Distribution $Q = (0, 0, 0.5, 0.5)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

$$= 0.5 * (|0.5 - 0| + |0.1 - 0| + |0.2 - 0.5| + |0.2 - 0.5|)$$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

Distribution $Q = (0, 0, 0.5, 0.5)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

$$= 0.5 * (|0.5 - 0| + |0.1 - 0| + |0.2 - 0.5| + |0.2 - 0.5|)$$

$$= 0.5 * (0.5 + 0.1 + 0.3 + 0.3)$$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

Distribution $Q = (0, 0, 0.5, 0.5)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

$$= 0.5 * (|0.5 - 0| + |0.1 - 0| + |0.2 - 0.5| + |0.2 - 0.5|)$$

$$= 0.5 * (0.5 + 0.1 + 0.3 + 0.3)$$

$$= 0.5 * 1.2 = \underline{\underline{0.6}}$$

t-closeness

Example: The whole data set contains 4 sensitive attribute classes:

500 Parking Violation, 100 Murder, 200 Speeding, 200 Robbery

Distribution $P = (0.5, 0.1, 0.2, 0.2)$

Distribution $Q = (0, 0, 0.5, 0.5)$

$$D(P, Q) = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|$$

$$= 0.5 * (|0.5 - 0| + |0.1 - 0| + |0.2 - 0.5| + |0.2 - 0.5|)$$

$$= 0.5 * (0.5 + 0.1 + 0.3 + 0.3)$$

$$= 0.5 * 1.2 = \underline{\underline{0.6}}$$

This eq. class has 0.6 closeness (and higher)

Take-Home Message

- k-anonymity provides some anonymity, but can be vulnerable to certain weaknesses (Homogeneity, Background Knowledge)
- l-diversity improves k-anonymity by adding constraints to the diversity of the sensitive values
- t-closeness compares the distribution of sensitive values to the overall distribution (no explicit statements about eq. class size)

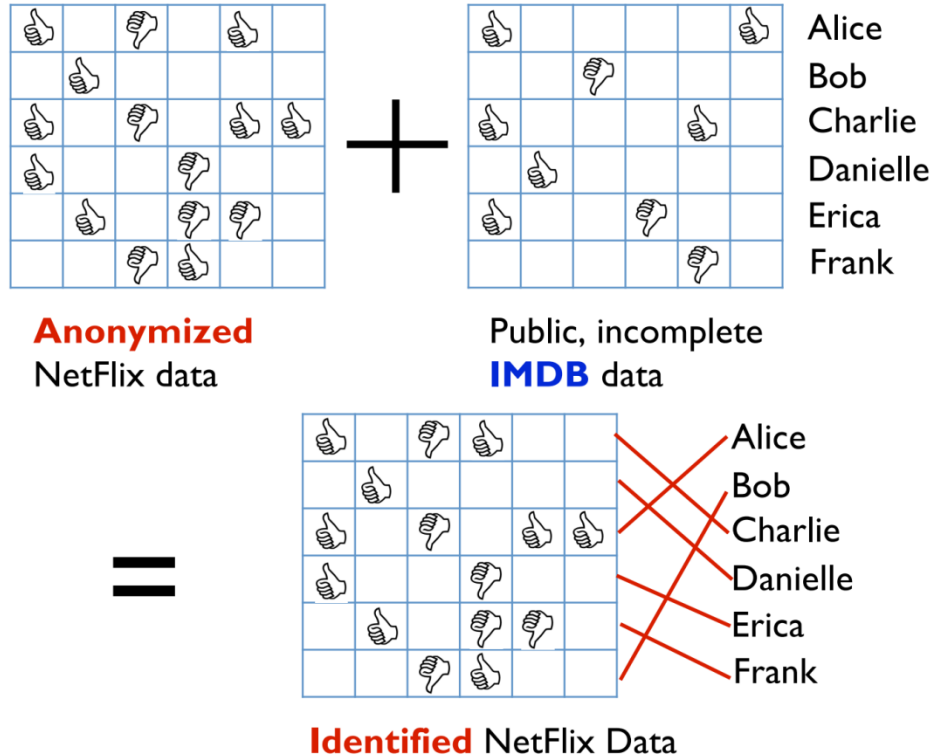
Outline

Recap

Anonymizing Data

Differential Privacy

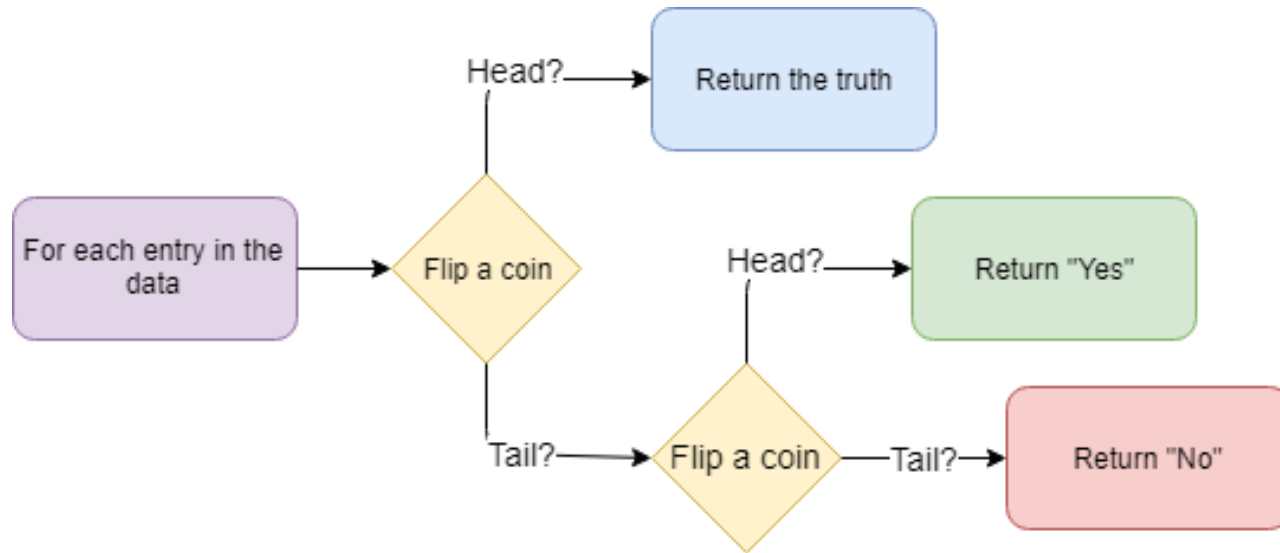
Linkage Attacks Are Still Possible

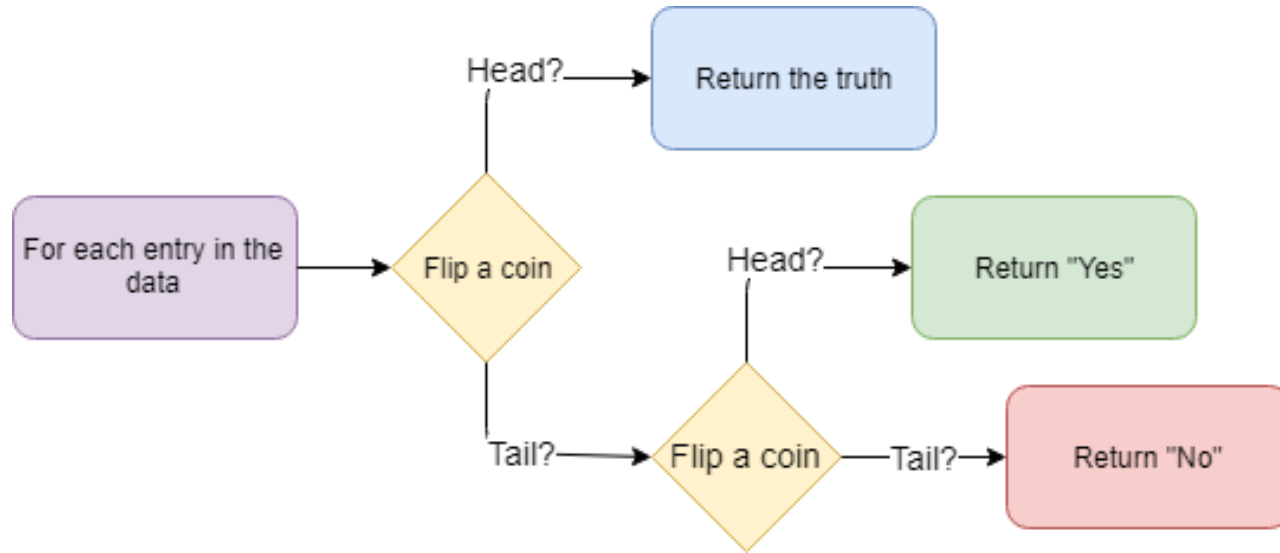


Is there a better way to hide identification?

Basic idea: Introduce randomness

- Coin Toss Example: When asked about feature x for record y
 - Toss a coin: if heads give right answer
 - If tails: throw coin again, answer yes if heads, no if tails
- Still has accuracy at some level of confidence
- Still has privacy at some level of confidence (plausible deniability)





First coin toss: Privacy parameter

Always heads: No privacy, perfect accuracy

Always tails: Perfect privacy, no accuracy

Differential Privacy

The distribution of attributes can still be estimated

If person X has attribute A , then

$$P(A|X) = 0.75$$

$$P(\sim A|X) = 0.25$$

If p is the true proportion of people with attribute A , then we expect

$$(1/4) + p/2 \quad \text{positive responses}$$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution** \mathbf{p} of people that like ice-cream?

$P(\text{Likes ice-cream}) =$

$P(\text{Coin} = \text{tails}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{tails})$

+

$P(\text{Coin} = \text{heads}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{heads})$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution** p of people that like ice-cream?

$P(\text{Likes ice-cream}) =$

$P(\text{Coin} = \text{tails}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{tails})$

+

$P(\text{Coin} = \text{heads}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{heads})$

$p = \text{true ice-cream lovers}$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution p** of people that like ice-cream?

$P(\text{Likes ice-cream}) =$

$P(\text{Coin} = \text{tails}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{tails})$

+

$P(\text{Coin} = \text{heads}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{heads})$

0.5 (fair coin toss)

p = true ice-cream lovers

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution p** of people that like ice-cream?

$P(\text{Likes ice-cream}) =$

$$\begin{aligned} &P(\text{Coin} = \text{tails}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{tails}) \\ &+ \\ &P(\text{Coin} = \text{heads}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{heads}) \end{aligned}$$

0.5 (fair coin toss) 0.5 (decided by another coin toss) $p = \text{true ice-cream lovers}$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution p** of people that like ice-cream?

0.7 (700 out of 1000)

$P(\text{Likes ice-cream}) =$

$$P(\text{Coin} = \text{tails}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{tails}) + P(\text{Coin} = \text{heads}) * P(\text{"I like ice-cream"} \mid \text{Coin} = \text{heads})$$

0.5 (fair coin toss)

0.5 (decided by another coin toss)

$p = \text{true ice-cream lovers}$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution** p of people that like ice-cream?

$$0.7 = 0.5 * 0.5 + 0.5 * p$$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution p** of people that like ice-cream?

$$0.7 = 0.5 * 0.5 + 0.5 * p$$

$$0.7 = 0.25 + 0.5 * p$$

$$0.45 = 0.5 * p$$

$$\mathbf{P = 0.9}$$

Differential Privacy

Example: 700 people say they like ice-cream, 300 say they do not like ice-cream

What is the **estimated true distribution p** of people that like ice-cream?

$$0.7 = 0.5 * 0.5 + 0.5 * p$$

$$0.7 = 0.25 + 0.5 * p$$

$$0.45 = 0.5 * p$$

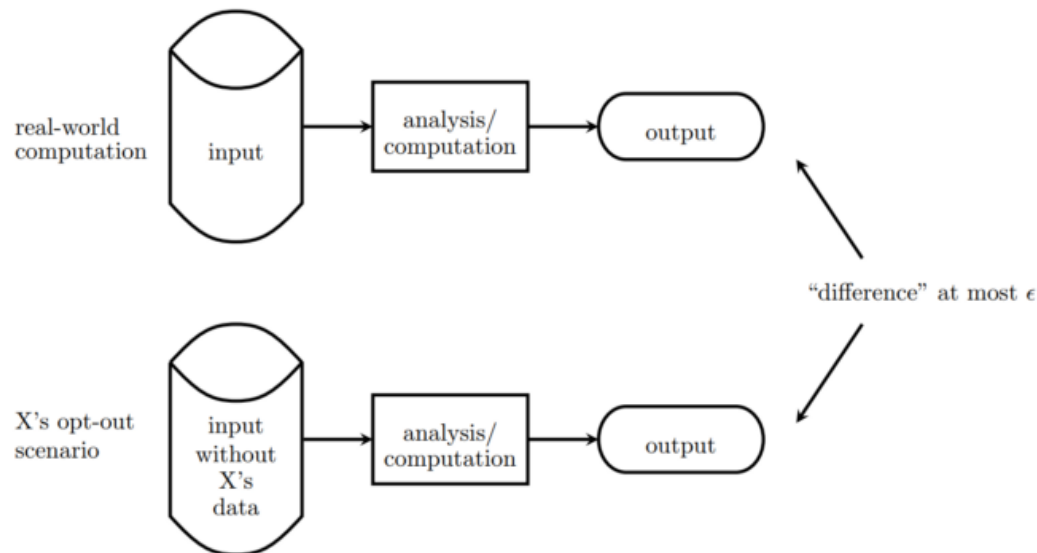
$$\mathbf{P = 0.9}$$

Approximately 90 % of the people liked ice-cream

Differential Privacy

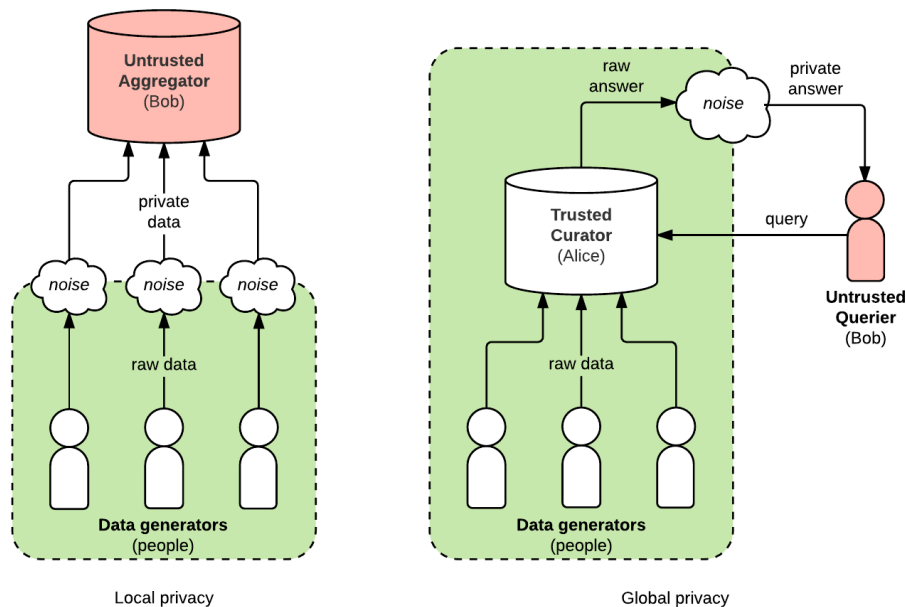
Coin toss is a very simplified version of Differential Privacy

Main idea:



Differential Privacy

Either submit data with noise (like the coin toss) or add global noise



Next Lecture

Language of Manipulation I

Now

Lecture Evaluation
The link is in Moodle
Thank you for your feedback!