

【ICDE2015勉強会】

Session 22:

## Data Privacy and Security 2

担当：曹洋 D2 (京都大学 吉川・馬研究室)

# Outline

---

## 1. A Hybrid Private Record Linkage Scheme: Separating Differentially Private Synopses From Matching Records

- Jianneng Cao (Institute for Infocomm Research), Fang-Yu Rao (Purdue University), Elisa Bertino (Purdue University), Murat Kantarcioglu (University of Texas at Dallas)

## 2. Conservative or Liberal? Personalized Differential Privacy

- Zach Jorgensen (North Carolina State University), Ting Yu (North Carolina State University), Graham Cormode (University of Warwick)

## 3. Differentially Private Frequent Sequence Mining via Sampling-based Candidate Pruning

- Shengzhi Xu (Beijing University of Posts and Telecommunications), Sen Su (Beijing University of Posts and Telecommunications), Xiang Cheng (Beijing University of Posts and Telecommunications), Zhengyi Li (Beijing University of Posts and Telecommunications), Li Xiong (Emory University)




# Differential Privacy

## ▶ $\epsilon$ -差分プライバシー ( $\epsilon$ -Differential Privacy, $\epsilon$ -DP)

### ▶ Privacy = “Right to be let alone”

- ▶ database  $D$  : each record is an individual's data
- ▶ user : “Please to analyze  $D$  **except my data**”

**D**

User	Disease
u1 	心臓病
u2 	HIV
u3 	HIV

# Differential Privacy




## ▶ $\epsilon$ -差分プライバシー ( $\epsilon$ -Differential Privacy, $\epsilon$ -DP)

### ▶ Privacy = “Right to be let alone”

- ▶ database  $D$  : each record is an individual's data
- ▶ user : “Please to analyze  $D$  **except my data**”

### ▶ DP is a formal definition of this idea

- ▶ For Database  $D$ , Query  $Q$ ; given  $\epsilon > 0$ ;
- ▶  $D'$  (**except one individual's data**) is denoted as “Neighboring Database” of  $D$ .

$D$	
User	Disease
u1 	心臓病
u2 	HIV
u3 	HIV



Neighboring Database  $D'$




Database  $D$

# Differential Privacy

## ▶ ε-差分プライバシー (ε-Differential Privacy, ε-DP)

### ▶ Privacy = “Right to be let alone”

- ▶ database **D** : each record is an individual’s data
- ▶ user : “Please to analyze **D** **except my data**”

<b>D</b>	
User	Disease
u1 	心臓病
u2 	HIV
u3 	HIV

### ▶ DP is a formal definition of this idea

- ▶ For Database **D**, Query **Q**; given  $\epsilon > 0$ ;

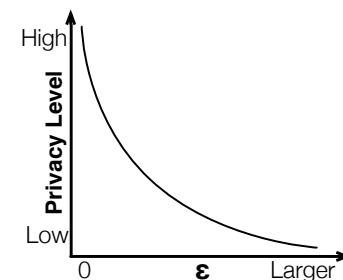
▶ **D'** (except one individual’s data) is denoted as “Neighboring Database” of **D**.

▶ If Mechanism **M** guarantees result **r** of  $Q(D)$  and  $Q(D')$  are “similar”:

$$\frac{\Pr(M(Q(D)) = r)}{\Pr(M(Q(D')) = r)} \leq e^\epsilon$$

e.g.,  $\frac{\Pr(M(Q(\text{u1, u2, u3})) = r)}{\Pr(M(Q(\text{u2, u3})) = r)} \leq e^\epsilon$

- ▶ then **M** satisfies ε-DP. ε is a specified privacy level:

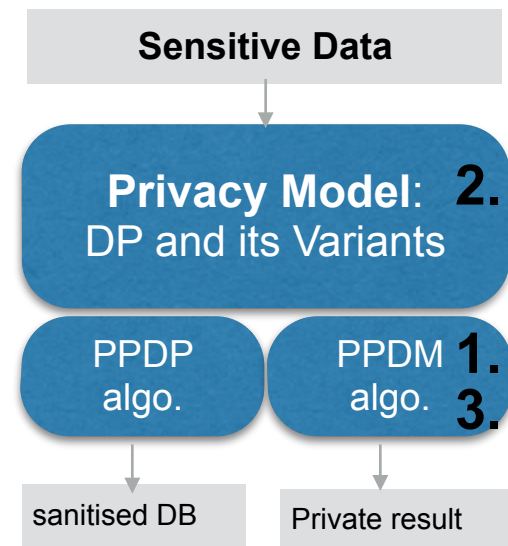


# Differential Privacy

## ▶ $\epsilon$ -差分プライバシー ( $\epsilon$ -Differential Privacy, $\epsilon$ -DP)

### ▶ DP-based data analysis

- ▶ Privacy Preserving Data Publishing (PPDP)
  - sensitive data → private (“safer”) data
- ▶ Privacy Preserving Data Mining (PPDM)
  - DM algo.  $\Rightarrow$  DP-algo
  - sensitive data → private result



(research scope)

1. A Hybrid Private Record Linkage Scheme:  
Separating Differentially Private Synopses From Matching Records
2. Conservative or Liberal? Personalized Differential Privacy
3. Differentially Private Frequent Sequence Mining via Sampling-based Candidate Pruning

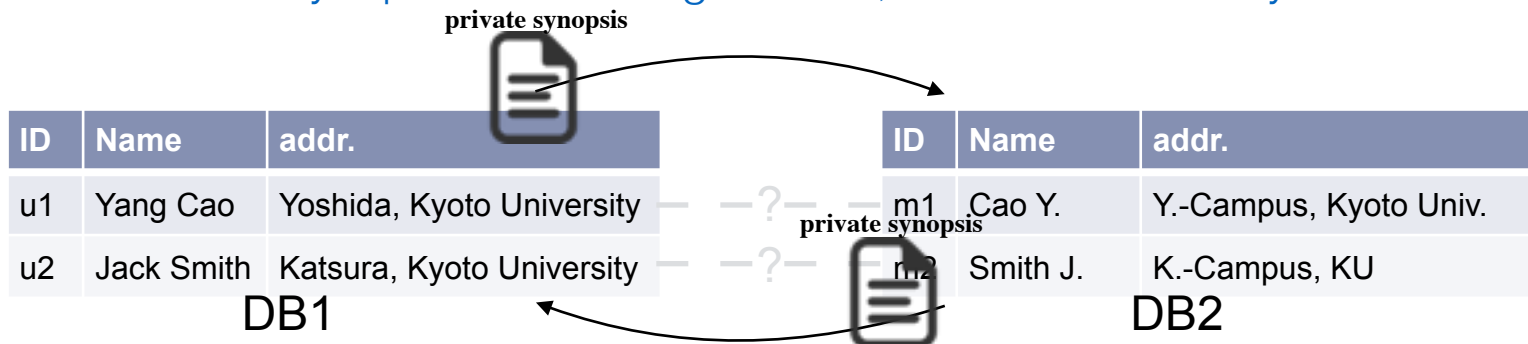


# 1. A Hybrid Private Record Linkage Scheme: Separating Differentially Private Synopses From Matching Records

## ▶ Motivation

### ▶ drawback of SMC-based PPRL

- ▶  **$O(m*n)$**  computational time ( $m, n$  are size of two owners' DB)
- ▶  $\Rightarrow$  Hybrid approach [1,2](SMC+private synopsis) to make it more scalable:
  - 1) release **private synopsis** ([1] by k-anonymity [2] by DP) to each other
    - **synopses**: Spatial indexing techniques (BSP-Tree, KD-Tree, or R-Tree) are used to form sub-sets (hyper-rectangles)
  - 2) compute similar records by private synopsis
- ▶ But [2] is **not totally differential private!** [this paper proved] (省略)
  - when DP-synopsis + matching records, DP do not hold anymore



[1] A. Inan, M. Kantarcioglu, E. Bertino, and M. Scannapieco, "A Hybrid Approach to Private Record Linkage," ICDE'08

[2] A. Inan, M. Kantarcioglu, G. Ghinita, and E. Bertino, "Private record matching using differential privacy," in EDBT'10



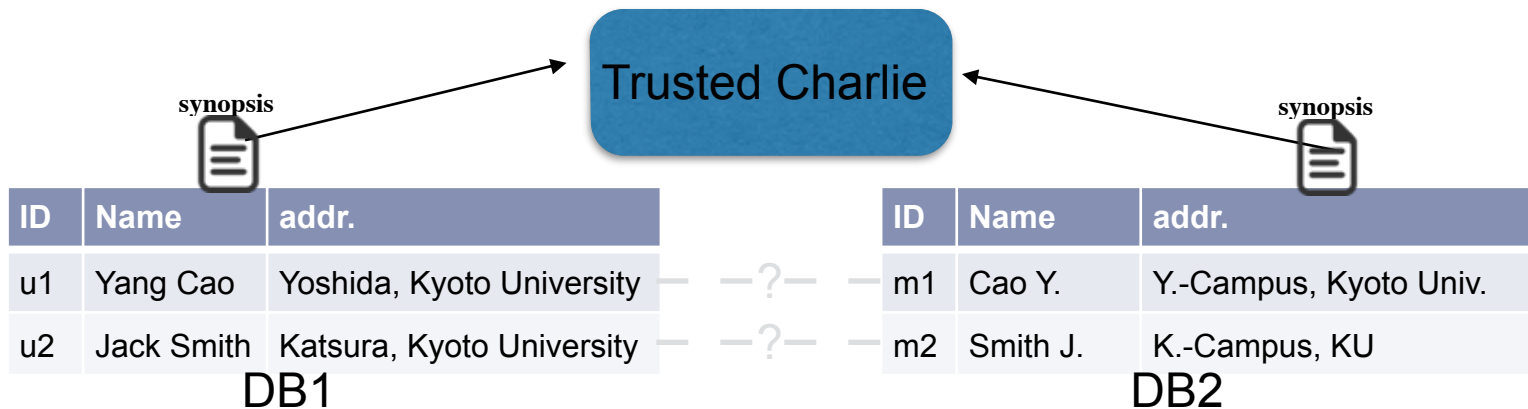
# 1. A Hybrid Private Record Linkage Scheme: Separating Differentially Private Synopses From Matching Records

## ▶ Proposed solution

- ▶ Separating Differentially Private Synopses From Matching Records
- ▶ two parties: Alice, Bob (previous study) ⇒ three parties: + Trusted Charlie (proposed)
  - ▶ Charlie prunes the record matching based on the received synopses.
- ▶ Building blocks
  - ▶ Secure Multi-party Computation (based on *Paillier Cryptosystem*)
  - ▶ Differential Private synopsis

## ▶ Contribution

- ▶ approximate PPRL framework
- ▶ prove the safety of this framework (satisfy differential privacy & SMC)



## 2. Conservative or Liberal? Personalized Differential Privacy

Recall:  $\epsilon_i \sim$  privacy level  
smaller  $\epsilon \sim$  higher privacy  
larger  $\epsilon \sim$  lower privacy

### ► Motivation

- In differential privacy, **all users** in dataset receive **the same privacy guarantee**

► if  $\frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 1, 3}))=r)} \leq e^{\epsilon_1}$     $\frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 2, 3}))=r)} \leq e^{\epsilon_2}$     $\frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 1, 4}))=r)} \leq e^{\epsilon_3}$

DP  
only consider  
"worst case"

then M need to satisfies  $\epsilon_{\min}$ -DP ( $\epsilon_{\min} = \text{Min}\{\epsilon_1, \epsilon_2, \epsilon_3\}$ ):

$$\frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 1, 3}))=r)} \leq e^{\epsilon_{\min}} \quad \frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 2, 3}))=r)} \leq e^{\epsilon_{\min}} \quad \frac{\Pr(M(Q(\text{User 1, 2}))=r)}{\Pr(M(Q(\text{User 1, 4}))=r)} \leq e^{\epsilon_{\min}}$$

- In reality, **different individuals** often have **different privacy preferences**
  - some users are under-protected, while others are over-protected.
- How to make "Personalized DP" possible? (exactly satisfy  $\epsilon_i$ )

## 2. Conservative or Liberal? Personalized Differential Privacy

### ▶ Proposed Mechanisms for PDP

$$\frac{\Pr(\mathbf{M}(Q(\text{User 1, 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 1}))=r)} \leq e^{\epsilon_1} \quad \frac{\Pr(\mathbf{M}(Q(\text{User 1, 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 2}))=r)} \leq e^{\epsilon_2} \quad \frac{\Pr(\mathbf{M}(Q(\text{User 1, 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 1, 2, 3}))=r)} \leq e^{\epsilon_3}$$

#### ▶ *Baseline*

- ▶ achieve **Min**{ $\epsilon_1, \epsilon_2, \epsilon_3$ }-DP, because of DP  $\Rightarrow$  PDP

#### ▶ *Threshold*

- ▶ give a threshold  $\epsilon_{\min} \leq t \leq \epsilon_{\max}$ , delete all data of  $u$  whose  $\epsilon_u < t$
- ▶ then achieve  $t$ -DP by traditional DP mechanisms

例

in the above example:

given  $\epsilon_1=1, \epsilon_2=2, \epsilon_3=3, t=2.5$ ,

By *threshold mechanism*, first deleting data of  $u_1, u_2$

then only need to achieve :

$$\frac{\Pr(\mathbf{M}(Q(\text{User 3}))=r)}{\Pr(\mathbf{M}(Q(\text{User}))=r)} \leq e^{\epsilon_3}$$

		D
User		Disease
<del><math>u_1</math></del>	<del></del>	<del>心臟病</del>
<del><math>u_2</math></del>	<del></del>	<del>HIV</del>
$u_3$		HIV

## 2. Conservative or Liberal? Personalized Differential Privacy

### ▶ Proposed Mechanisms for PDP

$$\frac{\Pr(\mathbf{M}(Q(\text{User 1, User 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 1}))=r)} \leq e^{\epsilon_1} \quad \frac{\Pr(\mathbf{M}(Q(\text{User 1, User 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 2}))=r)} \leq e^{\epsilon_2} \quad \frac{\Pr(\mathbf{M}(Q(\text{User 1, User 2}))=r)}{\Pr(\mathbf{M}(Q(\text{User 1, User 2, User 3}))=r)} \leq e^{\epsilon_3}$$

#### ▶ *Baseline*

- ▶ achieve **Min**{ $\epsilon_1, \epsilon_2, \epsilon_3$ }-DP, because of DP  $\Rightarrow$  PDP

#### ▶ *Threshold*

- ▶ give a threshold  $\epsilon_{\min} \leq \mathbf{t} \leq \epsilon_{\max}$ , delete all data of  $u$  whose  $\epsilon_u < \mathbf{t}$
- ▶ then achieve  $\mathbf{t}$ -DP by traditional DP mechanisms

#### ▶ *Sample*

- ▶ give a threshold  $\epsilon_{\min} \leq \mathbf{t} \leq \epsilon_{\max}$ , “carefully” sample a subset of  $D$
- ▶ then achieve  $\mathbf{t}$ -DP by traditional DP mechanisms

#### ▶ *Exponential-like*

- ▶ essentially, it is a **modified Exponential Mechanism**

## 2. Conservative or Liberal? Personalized Differential Privacy

---

- ▶ **Contributions:**
- ▶ Formally define a **new privacy model** based on DP: **Personalized Differential Privacy (PDP)**.
  - ▶ each  $u_i$  has a privacy preference  $\epsilon_i$
  - ▶ PDP is a **generalization** of DP: algo A satisfy DP  $\Rightarrow$  A satisfy PDP.
- ▶ Proposed several **Mechanisms** to achieve PDP
  - ▶ any existing DP-algorithm can be transformed to PDP-algorithm by these mechanisms
- ▶ **Experiments** on real/artificial dataset to transform several simple DP-algorithms to PDP-algorithms
  - ▶ PDP-Median (synthetic data) *Sample.M is best*
  - ▶ PDP-Count (synthetic data) *Exponential-like.M is best*
  - ▶ PDP-Linear regression (2012 U.S. census data) *Sample.M is best*

### 3. Differentially Private Frequent Sequence Mining via Sampling-based Candidate Pruning

#### Motivation

- Frequent Sequence Mining (FSM) - a fundamental DM algo.

ID	sequence
seq1	abc
seq2	abcd
seq3	abcde
seq4	abe

A sequence DB



	frequent sequences
3-seq	abc(3)
2-seq	ab(4), ac(3), bc(3)
1-seq	a(4), b(4), c(3)

FSM result

support

C3={abc,bac,...}

C2={ab,ac,bc,ba..}

- When sequences is sensitive : DP version of FSM (DP-FSM)

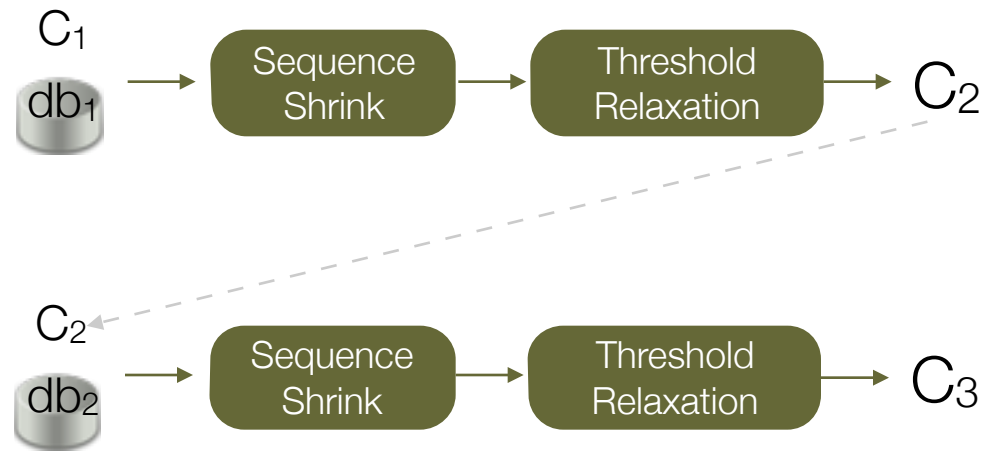
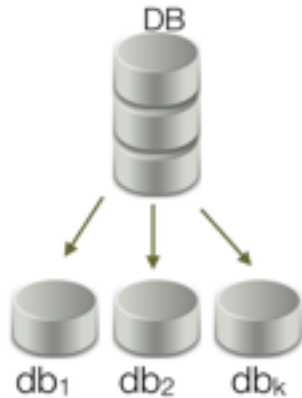
#### How to design **DP-FSM**?

$$\frac{\Pr(\text{DP-FSM}(D, r) = r)}{\Pr(\text{DP-FSM}(D', r) = r)} \leq e^\epsilon$$

- problem: e.g., removing seq3 will significantly effect 1-seq,2-seq...
- idea: reduce candidates of i-seq → (i+1)-seq , author use **heuristic methods** to detect which i-seq is potentially impossible (i+1)-seq.

### 3. Differentially Private Frequent Sequence Mining via Sampling-based Candidate Pruning

#### ▶ Proposed Solution



#### ▶ **Sequence Shrink**

- ▶ Irrelevant item detection
- ▶ Consecutive Pattern Compression
- ▶ Sequence Reconstruction

#### ▶ **Threshold Relaxation**

- ▶ decrease the probability of misestimating in db<sub>k</sub>

### 3. Differentially Private Frequent Sequence Mining via Sampling-based Candidate Pruning

---

#### ▶ **Related work:**

#### ▶ **Frequent Consecutive Sequence Mining[\*]**

▶ difference:

□ 1) [\*] consider consecutive

□ 2) [\*] is interactive publish, this paper is non-interactive publish

publish APIs

publish sanitised DB

#### ▶ Line of work in DP-based Frequent Pattern Mining

▶ DP-Frequent Graph Mining (FGM) [3]

▶ DP-Frequent Item Mining (FIM) [4,5]

▶ DP-Frequent Sequence Mining (FSM) [6]

[\*]L. Bonomi and L. Xiong, "A Two-phase Algorithm for Mining Sequential Patterns with Differential Privacy," CIKM,2013.

[3] E. Shen and T. Yu, "Mining Frequent Graph Patterns with Differential Privacy," KDD 2013.

[4] R. Bhaskar, S. Laxman, A. Smith, and A. Thakurta, "Discovering Frequent Patterns in Sensitive Data," KDD, 2010.

[5] R. Chen, et al., "Publishing set-valued data via differential privacy," VLDB, 2011.

[6] R. Chen, et al., "Differentially Private Sequential Data Publication via Variable-Length N-Grams," ACM CCS, 2012.