

《differential privacy and its application》

Sunday, October 20, 2019 11:46 PM

<<Differential privacy and Its Application>>

* Abstract:

* Problems:

- ① high sensitivity queries → large amount of noise
→ decrease the utility of dataset
- ② sparsity of the dataset → induce redundant noise (using randomized mechanism)
- ③ coupled records will release more info
→ decrease the privacy guarantee.

Schemes:

- ① application-aware sensitivity
(determined by the change of each record)
- ② shrink the randomized domain.
take records into groups and limit the randomized domain within each group.
- ③ re-measure the sensitivity to capture the relation among records.
→ coupled sensitivity < global sensitivity
- ④ iteration based coupled releasing mechanism
save privacy budgets

Chapter 1 Introduction

1. Differential Privacy:

Figure 1.1

Data collection

↓ trust boundary

Data analysis

- | data release | interactive setting
- | min-interactive setting
- | data mining, Interface Based data mining
- | fully access data mining

1.1 Data release

- | 1.1a interactive setting
- | 1.1b non-interactive setting

+2 Data mining

- | 1.1c Xtrust Miner
- | 1.1d vTrust Miner.

Chapter 2 Differential Privacy theory

2.1 notation

dataset size: $|D|$

query: $f: D \rightarrow R$

F : group query, $D \rightarrow R$.

M : randomization mechanism

maximal difference: Δf - sensitivity

randomized answer of querying: $f(D)$

noise = λ

outcome: R - output range

Definition 1: (ϵ, δ) -Differential privacy
 $\Pr[M(D) \in S] \leq e^\epsilon \cdot \Pr[M(D') \in S] + \delta$
 ϵ : privacy budget
 δ : violation probability
 Definition 2: Sequential Composition
 $M = \{M_1, \dots, M_m\}$ sequential performed on the dataset,
 the M will provide $(m \cdot \epsilon)$ -differential privacy
 each M_i provides ϵ privacy guarantee.
 Definition 3: Parallel Composition
 $\max\{\epsilon_1, \dots, \epsilon_m\}$ -differential privacy
 2.2 The sensitivity, global sensitivity
 local sensitivity
 Definition 1: Global sensitivity
 $\Delta f_{\text{G}} = \max_{D, D' : d(D, D')=1} \|f(D) - f(D')\|_1$.
 only works well with low sensitivity
 Definition 2: Local sensitivity
 $\Delta f_{\text{L}} = \max_D \|f(D) - f(D')\|_1$.
 $\therefore \Delta f_{\text{G}} = \max_S f_{\text{L}}$.
 generates less noise for queries with high sensitivity
 2.3 Mechanism: Laplace: numeric queries
 Exponential: non-numeric
 • Laplace probability: $p(x) = \frac{1}{2b} e^{-\frac{|x|}{b}}$
 definition 1: Laplace mechanism.
 $M(D) = f(D) + \text{Lap}(\frac{4\epsilon}{b})$.
 Definition 2: Exponential mechanism.
 score function: $g(D, \psi)$, represent how good an output ψ is for D .
 $M(D) = \{ \text{return } r \text{ with the probability } \propto \exp(\frac{g_r(D)}{2\sigma}) \}$
 2.4 Differential Private Data Release
 1 interactive setting
 non-interactive setting: improve the accuracy of query results
 research lines
 ① query release, batch queries release
 group-based release
 ② group-based tech in data release
 ③ contingency table release
 ④ release a synthetic dataset
 2.4.1 Interactive Data Release
 query f_i will only be answered after the answer of previous query f_{i-1} released.
 $\Pr_{f \in F} [f(D) - M(f(D)) \leq \alpha] \geq 1 - \delta$
 perturbation on query result.

2.4.2 Mechanism Design for Interactive Data Release
 Naive Laplace
 Query Separation
 Iteration
 ① Naive Laplace: most popular
 pros: maximal number(query) is limited in sub-linear n
 cons: noise magnitude will be large when dealing with continuous attribute
 ② Query separation:
 median mechanism, separating queries into hard and easy ones
 pros: answer exponentially more queries than Naive Laplace
 cons: time complexity
 ③ Iteration:
 1) Private Multiple Weights (PMW)
 pros: saved the privacy budget & decrease noise.
 2) Iterative Database Construct (IDC)
 when a significant difference occurred, the mechanism updates the current dataset for next iteration.
 pros: more general than PMW.

2.4.3 Histogram Release.
 Definition 1: Histogram Representation.
 Laplace: Pros - large noise
 ① partition the attributes into several mutual groups.
 ② kd-tree based partition method, to generate nearly uniform partitions
 ↓
 ③ Optimal partition method by minimizing the sum of squared Error (SSE)
 1) Noise first
 2) Structure first

④ maintain the consistency:
 constrained inference, sorted constraints
 hierarchical constraints.

2.4.4 Non-Interactive Data Release.
 pros: more queries, higher flexibility
 cons: too much noise, utility
 k batch queries release
 research lines
 ① contingency table release.
 ② group-based release.
 Sanitized dataset release
 ③ batch queries release
 ④ group-base release.
 ⑤ contingency table release
 ⑥ Synthetic Dataset Release.
 few many records → very accurate result
 algorithm → obtain the synthetic dataset in polynomial time.
 Definition 1: $(\alpha - \delta)$ -usefulness.
 $\Pr \{ |f(M(D)) - f(D)| \leq \alpha \} \geq 1 - \delta$.

2.5 Differentially Private Data Mining
 data mining with interface -
 data mining with full access -
 2.5.1. data mining with interface.
 1) sub-linear Queries (SULQ)
 2) interface privacy Integrated Queries platform (PINQ)
 ① SULQ: $\sum_i f_i(r_i) + N(0, R)$
 ② PINQ: Laplace, Exponential

2.5.2 classification

- ① SULQ framework - ID3 algorithm:
 - cons: ① privacy budget arrangement
 - ② dealing with continuous attributes
- ② PINQ improves:
 - ① exponential mechanism in the attributes choosing step
 - ② divided into ranges by every possible splitting value

2.5.3 clustering

- ① SULQ + k-means clustering

$$\bar{s}_j = \begin{cases} \text{SULQ}(f(r_i)) & := 1 \text{ if } j = \arg \min_j \\ 0 & \text{otherwise} \end{cases}$$

$$\bar{c}_j = \begin{cases} \text{SULQ}(f(r_i)) & := r_i, \text{ if } j = \arg \min_j \\ 0 & \text{otherwise} \end{cases}$$
- ② Local sensitivity on k-means clustering

2.5.4 Data Mixing with Full access.

- ① ID3
- ② logistic regression
 - ① inject Laplace noise in classifier w
 - ② add noise to the loss function
 - ③ LID.

2.5.5 Frequent Itemset Mining

- to discover itemsets that frequently appear in dataset
- ① top-k frequent itemsets mining algorithm DiffFIM
 - ① private selection = truncated frequency
 - ② frequency perturbation

- 2.6 Applications of Differential privacy
- applications: recommender system,
 - tagging recommender system
 - coupled differential privacy
- 2.6.1 Privacy Preserving recommender system
- CF: collaborative Filtering
 - neighborhood-based methods
 - model-based methods
 - ? privacy leak risk, kNN attack
- 2.6.2 Privacy Preserving Tagging recommender system.
- ① tag suppression approach.
 - model the user's profile using a tagging histogram
 - and eliminate sensitive tags from this profile.
 - cons: only release an incomplete dataset.
- 2.6.3 Differential privacy for a coupled dataset.

Chapter 3 Differentially Private Neighborhood-based Collaborative Filtering 时间轴

3.1 Introduction.

- application-aware sensitivity in the context of neighborhood-based collaborative filtering

① collaborative filtering (CF), neighborhood-based methods

3.2 fundamentals of a recommender system.

- D can be represented by column vectors:

$$D = [t_1^T, t_2^T, \dots, t_m^T]^T, \text{ and } t_i = [r_{i1}, r_{i2}, \dots, r_{in}]^T$$

$s_{(i,j)}$ denotes the similarity between t_i with t_j

① collaborative Filtering (CF)

- neighborhood-based methods
- neighbor selection
- rating prediction

3.3 research lines

- ① How to preserve neighborhood privacy - private Neighbor selection
- ② How to define sensitivity for recommendation purpose
 - recommendation-aware sensitivity
- ③ How to design the exponential mechanism for CF?

Similarity: Pearson correlation coefficient (PCC)
 | cosine-based similarity (cos)

a. Pearson correlation coefficient:

$$\text{sim}(i,j) = \frac{\sum_{k \in U} (r_{ki} - \bar{r}_i)(r_{kj} - \bar{r}_j)}{\sqrt{\sum_{k \in U} (r_{ki} - \bar{r}_i)^2} \sqrt{\sum_{k \in U} (r_{kj} - \bar{r}_j)^2}}$$

 \bar{r} is the average rating given by relative users.

b. cosine-based similarity

$$\text{sim}(i,j) = \frac{r_i \cdot r_j}{\|r_i\|_2 \|r_j\|_2}$$

Rating prediction stage:
 item-based: $r_{ui} = \frac{\sum_{j \in U} s_{ui,j} \cdot r_{uj}}{\sum_{j \in U} |s_{ui,j}|}$
 user-based: $r_{vu} = \bar{r}_v + \frac{\sum_{u \in U} s_{vu} (r_u - \bar{r}_u)}{\sum_u |s_{vu}|}$
 \bar{r}_v, \bar{r}_u is the average rating given by user v and u .
 s_{vu} denotes the similarity between v to u .

- ② known attack to CF., $O(\log n)$.
- ③ related work
 - Privacy Preserving Recommender System.
 - Differential Private Recommender System

3.3 private Neighbor Collaborative Filtering
 ① overview of algorithm. (PNCF)

private neighbor selection: prevent inferring neighbor
 PNCF | perturbation = prevent inferring rating

② the private neighbor selection
 means privately select k neighbors from a list of candidates for the privacy preserving purpose.
 A. recommendation-Aware sensitivity, $R_{sc(i,j)}$.
 B. private neighbor selection + accuracy

truncated similarity:

$$S_{sc(i,j)} = \max(s_{ci,j}, s_{ck(i,:)}) - w$$

 ensures no item in $N_k(t_i)$ has a similarity less than $s_{ci,j} - w$. and is greater than $(s_{ci,j}) + w$

3.4 Privacy and utility analysis.

① Privacy analysis.

$$S_{noise}(i,j) = S_{ci,j} + \text{Lap}\left(\frac{2 \cdot R_{sc(i,j)}}{\epsilon}\right)$$

② Utility analysis.

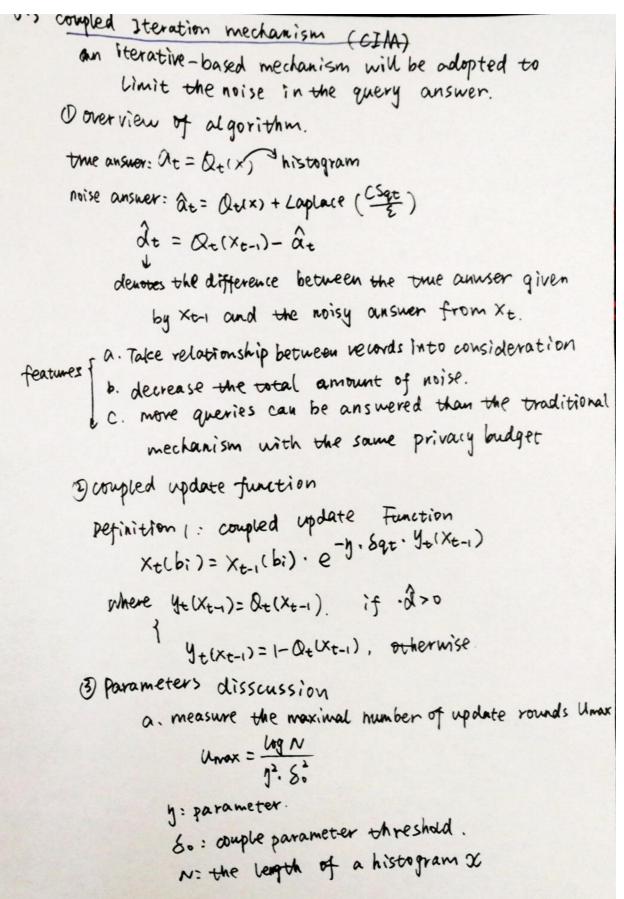
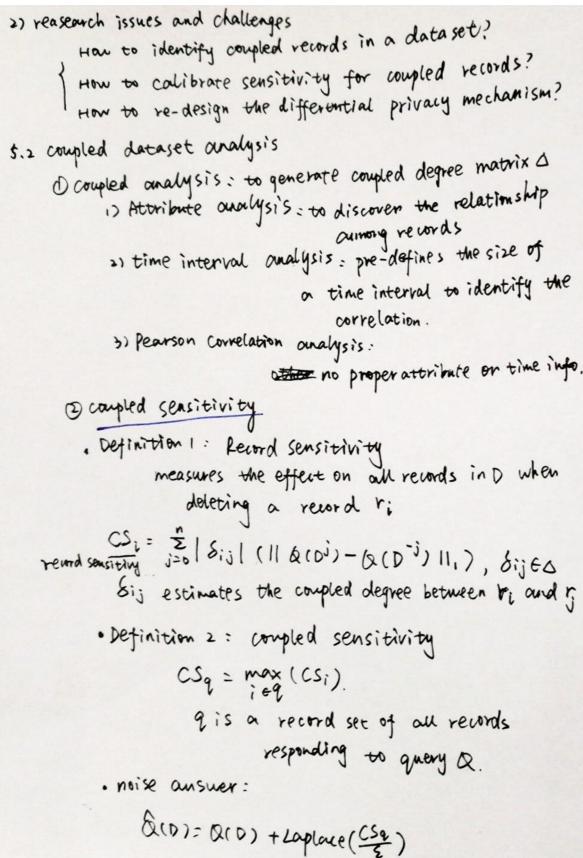
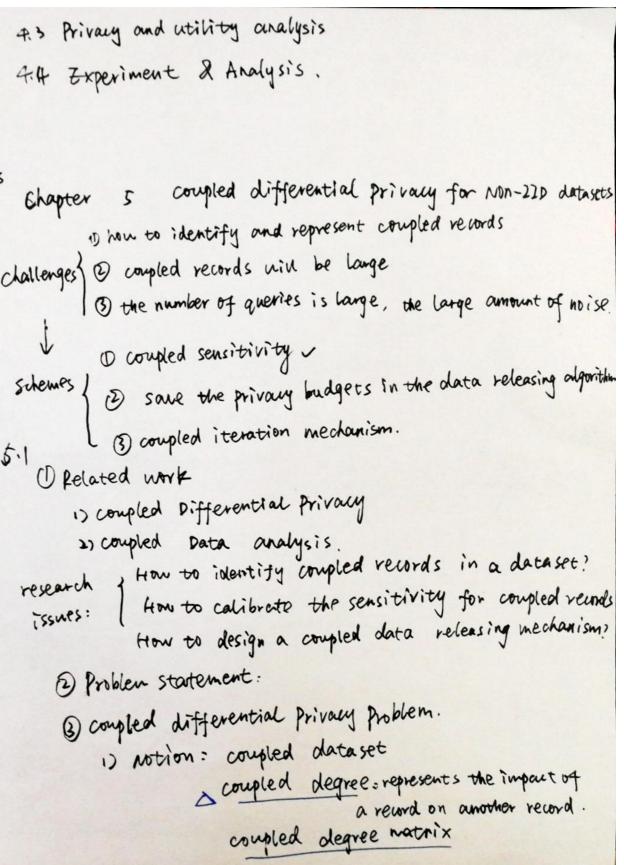
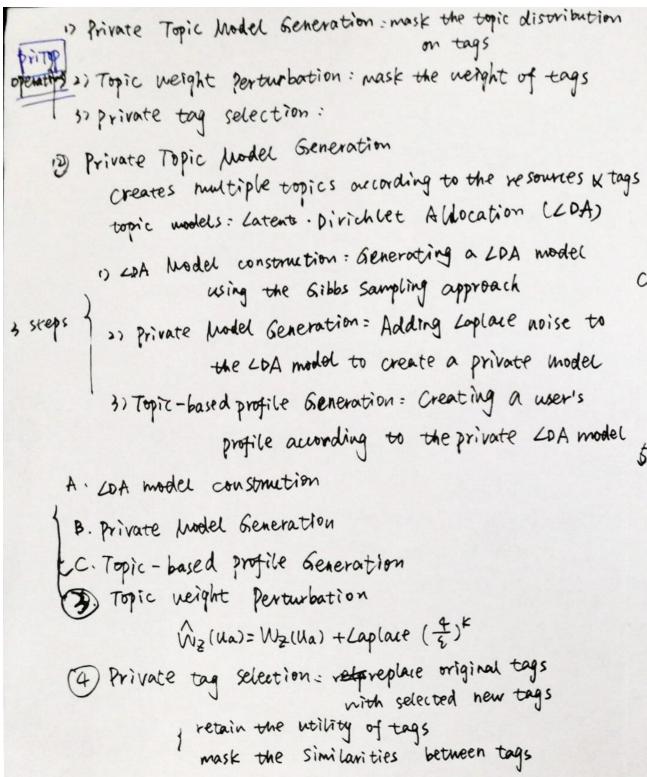
3.5 Experiment and analysis

	① How does PNCF perform on the privacy preserving issue from the perspective of neighbors?
	② How does the proposed recommendation-Aware sensitivity satisfy the recommendation purpose?

- ① How does PNCF perform compared to other-related methods?
- ② How does the privacy budget affect the performance of PNCF?
- ① Datasets and measurements
ex: Netflix, MovieLens
- Mean Absolute Error (MAE):
$$\text{MAE} = \frac{1}{|T|} \sum_{(u,i) \in T} |r_{ui} - \hat{r}_{ui}|$$
- ② Performance of PNCF.
- ③ performance of recommendation-Aware sensitivity
- ④ performance of the private neighbor selection mechanism
- ⑤ Effect of privacy budget
to determine privacy preserving level

- chapter 4 Differentially Private Tagging Recommender System.
- | | |
|--|--|
| | ① ignore the relationship between items/sources |
| | ② randomized mechanism will result in a large magnitude of noise with million tags |
| | ③ generate a synthetic dataset which retains the relationship |
- ↓
 schemes
- | | |
|--|---|
| | ② shrink the randomized domain |
| | ① utility v. PriTop. |
| | ② model-based method, & shrink the randomized domain. |
- contribution:
- | | |
|--|--------------------|
| | ③ privacy utility. |
|--|--------------------|

- 4.1 Fundamental of Tagging Recommender System.
- ① Notation:
 $T(u_a, r_b)$: represent all tags flagged by the user u_a on resource r_b .
 $P(u_a)$: a user u_a 's profile.
 $W(u_a)$: tagging's weights
- ② Tagging recommender system
- 4.2 Private Tagging Release.
 Private Topic-based Tagging Release (PriTop)
- ① overview of PriTop algorithm



a estimate the possible number of update rounds U_2 for a query set Q
 $P_1 = e^{-\frac{2\alpha|T-\alpha|}{CSq}}$ the accuracy of C2M
 probability of update
 $P_2 = 1 - e^{-\frac{2\alpha|T-\alpha|}{CSq}}$
 probability of non-update
The accuracy of C2M is related to T.

5.4 mechanism analysis.

① privacy analysis.

② utility analysis

definition 1: (α, β) -accuracy for C2M.

$$\max_t |C2M_t(x_i) - Q_t(x_i)| \leq \alpha.$$

↓
is the output of C2M in round t.

$$\text{when } \alpha \geq \frac{CSq}{2\delta_0} \left(\log \frac{P_1 P_2 L}{\beta} \right) + \frac{\epsilon}{2},$$

C2M is satisfied with (α, β) -accuracy.

5.5 Experiments.

Mean Square Error:

$$MSE = \frac{1}{|Q|} \sum_{i=1}^{|Q|} (\hat{Q}_i(x) - Q_i(x))^2.$$

lower MSE implies a better utility

③ the privacy budget ϵ serves as a key parameter to determine privacy

Chapter 6: Conclusion.

Future work:

① Differential privacy for continuous releasing

② Distributed Differential privacy

③ Synthetic dataset releasing for arbitrary Mining purposes

④ privacy & Game theory

Summary of differential privacy

Summary & Logic

Problems

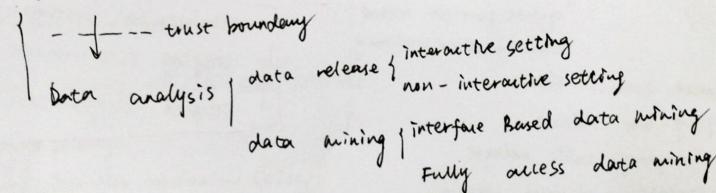
- ① high sensitivity queries → large amount of noise
→ decrease the utility of dataset
- ② sparsity of the dataset → induce redundant noise.
- ③ coupled records will release more info (privacy data).

schemes

- ① application-aware sensitivity
- ② using randomized mechanism & shrink the randomized domain.
- ③ remeasure the sensitivity to capture the relation among records
 - a. coupled sensitivity
 - b. iteration based coupled releasing mechanism

Differential privacy scenarios:

Data collection



Theory: notion; (ϵ, δ) -Differential privacy

privacy budget violation probability

sensitivity: global sensitivity

local sensitivity

Mechanism: Laplace - numeric

Exponential - non-numeric

Differential private Data mining:

data mining with interface | sub-linear queries (SubQ),
Privacy Integrated queries platform (PIQ)

data mining with full access

ID3.

logistic regression.
frequent Itemset Mining (DiffFIM).

Differentially Private data release:

interactive setting: mechanism of naive Laplace

query separation

Iteration

Histogram release: partition

kd-tree based partition

optimal partition method

maintained the consistency

Applications of differential privacy

recommender system

tagging recommender system

coupled differential privacy

non-interactive setting

batch queries release

group-base release

contingency table release

Synthetic dataset release

1. Differential private neighborhood-based collaborative Filtering (PnCF).

Similarity = { Pearson correlation coefficient (PCC)
cosine-based similarity (cos)

rating selection { item-based
user-based.

Δ | PnCF | { Private neighbor Selection
perturbation: exponential mechanism.

2. Differential Private Tagging system

| PriTop | { Private Topic Model generation { LDA model construction
Topic weight perturbation | private model Generation
Topic-based profile Generation.
Private tag selection.

3. Coupled differential privacy for non-IID data sets

Coupled differential privacy:
↳ Coupled Iteration mechanism (CIM)