Co-clustering for differentially private synthetic data generation

Tarek Benkhelif, Françoise Fessant, Fabrice Clérot and Guillaume Raschia January 23, 2018

Orange Labs & LS2N Journée thématique EGC & IA : Données personnelles, vie privée et éthique





Context

- Releasing data, either in their original or aggregated form
- Protecting individuals represented in the data
- Providing sufficient utility













- Same format as the original data
- Multidimensional data
- Independent of the data mining task



- Same format as the original data
- Multidimensional data
- Independent of the data mining task

Differential Privacy: Intuition



Differential Privacy

- It should not harm you or help you as an individual to enter or to leave the dataset.
- To ensure this property, we need a mechanism whose output is nearly unchanged by the presence or absence of a single respondent in the database.
- In constructing a formal approach, we concentrate on pairs of databases (D_1, D_2) differing on only one row, with one a subset of the other and the larger database containing a single additional row.

 ε -Differential Privacy [Dwo06] A data release mechanism \mathcal{A} satisfies ε -differential privacy if for all neighboring database D_1 and D_2 , and released output O, $Pr[\mathcal{A}(D_1) = O] \leq e^{\varepsilon} \times Pr[\mathcal{A}(D_2) = O]$.

Achieving ε -DP : Laplace mechanism

Adds random noise to the true answer of a query Q, $\mathcal{A}_Q(D) = Q(D) + \tilde{N}$, where \tilde{N} is the Laplace noise. The magnitude of the noise depends on the privacy levels and the query's sensitivity

Existing approaches



1. Discretize attribute domain into cells



- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)



- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...



- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...
 - 3.1 Answer queries directly (assume distribution is uniform within cell)



- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...
 - 3.1 Answer queries directly (assume distribution is uniform within cell)
 - 3.2 Generate synthetic data (derive distribution from counts and sample)



Limitations

• Granularity of discretization

- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...
 - 3.1 Answer queries directly (assume distribution is uniform within cell)
 - 3.2 Generate synthetic data (derive distribution from counts and sample)



- Granularity of discretization
 - Coarse: detail lost

- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...
 - 3.1 Answer queries directly (assume distribution is uniform within cell)
 - 3.2 Generate synthetic data (derive distribution from counts and sample)



- Granularity of discretization
 - Coarse: detail lost
 - Fine: noise overwhelms signal

- 1. Discretize attribute domain into cells
- 2. Add noise to cell counts (Laplace mechanism)
- 3. Use noisy counts to either...
 - 3.1 Answer queries directly (assume distribution is uniform within cell)
 - 3.2 Generate synthetic data (derive distribution from counts and sample)

DP multidimensional data release approaches

Approach	Dimension	Mixed data type	Parameter(s)
DPCube [XXFG12]	Multi-D	×	Variance threshold
DP-MHMD [RKS16]	Multi-D	×	Attribute grouping
DiffGen [MCFY11]	Multi-D	\checkmark	 Attributes taxonomy n^{br} of specializations
PrivBayes [ZCP ⁺ 14]	Multi-D	\checkmark	Bayesian network degree



Proposition: DPCocGen

Bi-clustering

Simultaneously partition the rows and columns of a data matrix.

D-clustering

Simultaneously partition the *d*-dimensions of a data hyper cube.





Capture the interaction (underlying structure) between the *d* entities.

Grouping

Discover the best reordering and grouping of the data cube ¹ that:

• maximize the mutual information between the *d*-clusterings

Aggregation

Aggregation ability which allows to decrease the number of clusters in a greedy optimal way



¹Boullé , M.: Functional data clustering via piecewise constant nonparametric density estimation.



Original data





Original data









Evaluation of DPCocGen

Criteria

1. Joint distribution preservation

Criteria

- 1. Joint distribution preservation
- 2. Relative error for random range queries

Criteria

- 1. Joint distribution preservation
- 2. Relative error for random range queries
- 3. Performance in classification with a classifier that learns from synthetic data

Criteria

- 1. Joint distribution preservation
- 2. Relative error for random range queries
- 3. Performance in classification with a classifier that learns from synthetic data

To observe

1. Impact of the privacy budget ε

Criteria

- 1. Joint distribution preservation
- 2. Relative error for random range queries
- 3. Performance in classification with a classifier that learns from synthetic data

- 1. Impact of the privacy budget ε
- 2. Impact of the aggregation level (number of cells)

Criteria

- 1. Joint distribution preservation
- 2. Relative error for random range queries
- 3. Performance in classification with a classifier that learns from synthetic data

- 1. Impact of the privacy budget ε
- 2. Impact of the aggregation level (number of cells)
- 3. Comparison with the base line algorithm and PrivBayes

Adult

- The dataset ² contains 48,842 instances and has 14 different attributes. The characteristics of the attributes are both numeric and nominal
- The attributes {age, workclass, education, relationship, sex} are retained
- We discretize continuous attributes into data-independent equi-width partitions

²UC Irvine Machine Learning Repository

The Hellinger distance between two discrete probability

distributions
$$P = (p_1, ..., p_k)$$
 and $Q = (q_1, ..., q_k)$ is given by :
 $D_{Hellinger}(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^k (\sqrt{p_i} - \sqrt{q_i})^2}$

Experiment

- Compute the multivariate distribution vector *P* of the original dataset

The Hellinger distance between two discrete probability distributions $P = (p_1, ..., p_k)$ and $Q = (q_1, ..., q_k)$ is given by : $D_{Hellinger}(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{k} (\sqrt{p_i} - \sqrt{q_i})^2}$

- Compute the multivariate distribution vector *P* of the original dataset
- Compute the multivariate distribution vector *Q* of the synthetic data generated using DPCocGen

The Hellinger distance between two discrete probability distributions $P = (p_1, ..., p_k)$ and $Q = (q_1, ..., q_k)$ is given by : $D_{Hellinger}(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{k} (\sqrt{p_i} - \sqrt{q_i})^2}$

- Compute the multivariate distribution vector *P* of the original dataset
- Compute the multivariate distribution vector *Q* of the synthetic data generated using DPCocGen
- Compute the multivariate distribution vector Q' of the synthetic data generated using Base line

The Hellinger distance between two discrete probability distributions $P = (p_1, ..., p_k)$ and $Q = (q_1, ..., q_k)$ is given by : $D_{Hellinger}(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^{k} (\sqrt{p_i} - \sqrt{q_i})^2}$

- Compute the multivariate distribution vector *P* of the original dataset
- Compute the multivariate distribution vector *Q* of the synthetic data generated using DPCocGen
- Compute the multivariate distribution vector Q' of the synthetic data generated using Base line
- Compute $D_{Hellinger}(P, Q)$ and $D_{Hellinger}(P, Q')$

Results: Multivariate distribution preservation



50 datasets are generated for each configuration

- Generate 100 random queries
- Compute all the queries and report their average error
- Iterate over 15 runs

Results: Random range queries



- Randomly divide the original dataset into 2 sets :
 - Training set: contains 80% of the data
 - Test set: contains 20% of the data

- Randomly divide the original dataset into 2 sets :
 - Training set: contains 80% of the data
 - Test set: contains 20% of the data
- Generate synthetic data using DPCocGen, Base line and PrivBayes on the Training set

- Randomly divide the original dataset into 2 sets :
 - Training set: contains 80% of the data
 - Test set: contains 20% of the data
- Generate synthetic data using DPCocGen, Base line and PrivBayes on the Training set
- Learn a naive Bayes classifier from the synthetic data to predict the value of the attribute *Sex*

- Randomly divide the original dataset into 2 sets :
 - Training set: contains 80% of the data
 - Test set: contains 20% of the data
- Generate synthetic data using DPCocGen, Base line and PrivBayes on the Training set
- Learn a naive Bayes classifier from the synthetic data to predict the value of the attribute *Sex*
- Measure classification performances of the trained models on the Test set

Classification : predict Sex



Figure 1: Average AUC, across 15 runs

Conclusion

Advantages

- 1. Parameter-free
- 2. Preserves utility

Limits

- 1. Limited dimension
- 2. Requires a discretization step

Perspectives

1. Using differentially private dimension reduction strategies, to tackle the dimension limitation

Thank you!



Differential privacy.

In Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener, editors, *Automata, Languages and Programming*, volume 4052 of *Lecture Notes in Computer Science*, pages 1–12. Springer Berlin Heidelberg, 2006.

📔 Noman Mohammed, Rui Chen, Benjamin Fung, and Philip S Yu.

Differentially private data release for data mining.

In Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 493–501. ACM, 2011.

Harichandan Roy, Murat Kantarcioglu, and Latanya Sweeney.

Practical differentially private modeling of human movement data.

In IFIP Annual Conference on Data and Applications Security and Privacy, pages 170–178. Springer, 2016.

Yonghui Xiao, Li Xiong, Liyue Fan, and Slawomir Goryczka. Dpcube: differentially private histogram release through multidimensional partitioning. arXiv preprint arXiv:1202.5358, 2012.

Jun Zhang, Graham Cormode, Cecilia M Procopiuc, Divesh Srivastava, and Xiaokui Xiao. Privbayes: Private data release via bayesian networks.

In Proceedings of the 2014 ACM SIGMOD international conference on Management of data, pages 1423–1434. ACM, 2014.