

PRIVACY-PRESERVING DATA PUBLISHING

Masooma Iftikhar masooma.iftikhar@anu.edu.au

School of Computing College of Engineering and Computer Science The Australian National University, Canberra, Australia

JANUARY 27, 2022



- Introduction
- Related Work
- Research Goals and Challenges
- Contributions and Publications
- Privacy-Preserving Data Publishing (PPDP) Mechanisms
 - Relational Data Publishing
 - Graph Data Publishing

INTRODUCTION



Publishing data about individuals poses a privacy threat.



Related Work

PPDP Paradigms







DP bounds a shift in the output distribution of a randomized mechanism \mathcal{K} that can be caused by a small change in its input.





One way to satisfy **DP** is to add controlled **noise** to the output of a query *f* using a random distribution (e.g., Laplace).



Noise is calibrated according to sensitivity (Δ) of f and ε .

- Δ is maximum variation in f between X and Y.
- **\varepsilon** > **o** is the privacy parameter.

RESEARCH GOALS AND CHALLENGES



To develop privacy-preserving data publishing (PPDP) mechanisms.

- **Privacy:** Provide privacy guarantee of differential privacy.
- Utility: Enhance the accuracy of published data while providing differential privacy.



- How to determine the right amount of noise that guarantees both privacy and accuracy under DP?
- How to reduce the amount of noise needed to achieve DP by controlling sensitivity?
- How to enhance data utility under different data structures while providing DP guarantee?

CONTRIBUTIONS AND PUBLICATIONS

CONTRIBUTIONS



■ Part I: Relational Data Publishing

- Reduce sensitivity by incorporating microaggregation into DP.
- Enhance overall utility by reducing both error produced during microaggregation and DP.
- **Part II:** Graph Data Publishing
 - Preserve topological structure of a graph through adding controlled perturbation to its edges.
 - Publish higher-order network statistics while providing DP guarantee to nodes in a network.
 - Consider personalization to achieve personal data protection under DP while enhancing utility.

PUBLICATIONS



- Part I: Relational Data Publishing
 - M.Iftikhar, Q.Wang, Y.Lin. Publishing Differentially Private Datasets via Stable Microaggregation. EDBT 2019.
 - M.Iftikhar, Q.Wang, Y.Li, Y.Lin, J.Shao. Differentially Private Data Release via α-Stable Microaggregation. Under journal preparation.
- Part II: Graph Data Publishing
 - M.Iftikhar, Q.Wang, Y.Lin. dK-Microaggregation: Anonymizing Graphs with Differential Privacy Guarantees. PAKDD 2020.
 - M.Iftikhar, Q.Wang. dK-Projection: Publishing Graph Joint Degree Distribution with Node Differential Privacy. PAKDD 2021.
 - M.Iftikhar, Q.Wang, Y.Li. dK-Personalization: Publishing Network Statistics with Personalized Differential Privacy. Accepted by PAKDD 2022.

Relational Data Publishing



Name	Zip	Age	Nationality	Disease	Name	Zip	Age	Nationality	Disease
Eve	13053	28	Russian	Heart	Eve	13053	28	Russian	Heart
Alice	13068	29	American	Heart	Alice	13068	29	American	Heart
Bob	13068	21	Japanese	Flu	Bob	13068	21	Japanese	Flu
Amy	14853	50	Indian	Cancer	Farhan	26001	29	Pakistan	Covid-19

Х

Y

X and Y differ in one record



- **Problem:** Given a dataset *X*, we want to generate X_{ε} that can provide DP guarantee.
- Key Observation: Microaggregation can help to reduce sensitivity for improving data utility.



• A microaggregated dataset \overline{X} is added between X and X_{ε} .



MICROAGGREGATION



A microaggregation algorithm \mathcal{M} consists of two phases:





An arbitrary \mathcal{M} could not reduce sensitivity when incorporated into DP.



- *M* may generate considerably different clusters.
- Leading to a much larger $\Delta(f \circ \mathcal{M})$ than $\Delta(f)$.



Stable microaggregation characterizes a certain correspondence of clusters in microaggregated datasets to control sensitivity.





There is a bijection between C_X and C_Y such that at most two pairs of corresponding clusters differ in a single record.



 $\blacksquare \Delta(f \circ \mathcal{M}) \text{ is reduced to } (2 \times \Delta(f)/k) \text{ from } (n/k \times \Delta(f)).$

Needs $k \ge 2$ to reduce noise regardless of dataset size n.



We proposed a unified framework called α -stable microaggregation which generalized stable microaggregation:

- At most α pairs of corresponding clusters in C_X and C_Y differ in a single record to enhance within cluster homogeneity.
- α indicates the trade-off between error introduced during microaggregation and achieving DP.
- $\Delta(f \circ \mathcal{M})$ is $(\alpha \times \Delta(f)/k)$.



- 1. Sequential α -stable microaggregation algorithm: Performs record-level search such that swapping records leads to enhance within cluster homogeneity.
- Decisional α-stable microaggregation algorithm: Performs cluster-level search followed by record-level within the selected cluster to enhance within cluster homogeneity.



Three datasets:

- (1) CENSUS contains 1,080 records.
- (2) EIA contains 4,092 records.
- (2) Tarragona contains 834 records.

Two measures:

- IL1s measures information loss between the original and differentially private datasets [8].
- RL measures the percentage of record linkage between the original and differentially private datasets [7].



Does the proposed framework yield less *IL*1s and **RL** in microaggregated datasets?



- Significantly reduced *IL*1s during microaggregation by providing better within cluster homogeneity.
- Attain **RL** below 5% for $k \ge 20$.



Does the proposed framework yield less *IL*1s and *RL* in differentially private datasets?



- Reduced *IL*1s as $\Delta(f \circ \mathcal{M})$ is $\alpha \times \Delta(f)/k$.
- When $k \ge \alpha$, noise is reduced regardless of n.
- Attain **RL** below 5%



What kind of trade-off exists between utility and privacy while generating differentially private datasets?

- Error caused by DP that depend on $\Delta(f \circ \mathcal{M})$ dominates the impact on data utility as compared to microaggregation error.
- Reducing sensitivity can increase the data utility but it is not straightforward.
- Adding more noise provides better privacy but less utility and vise versa.

GRAPH DATA PUBLISHING

NEIGHBORING GRAPHS







- Graph data is highly sensitive to structural changes.
- Directly perturbing graph data for achieving DP often leads to inject a large amount of noise.
- Preserving topological structures of an original graph while achieving DP is not straightforward.



- **Problem:** Given a graph G, we want to publish graph statistics under the guarantee of DP.
- **Key Observation:** *dK*-distributions can serve as a good basis for representing graph statistics. The <u>dK-graph model</u> [5] provides a systematic way of extracting *dK*-distributions from *G*.



dK-distributions are a set of reproducible graph properties, which capture degree correlations within *d*-sized subgraphs of a graph.



When d = |V|, a dK-distribution specifies the entire graph.

 γ^{dK}(G) queries the dK-distribution of G.



DP has two variants when applying to graph data:

- **Edge-DP:** Hide the presence and absence of a single edge in a graph.
- Node-DP: Hide the presence and absence of a single node and the set of edges incident to that node.

Node-DP can provide stronger privacy protection than edge-DP.

Achieving node-DP is more challenging than for edge-DP as graph data is highly sensitive under node-DP.



dK-Microaggregation Framework: Microaggregation helps to reduce the overall noise needed to achieve edge-DP.



• $\gamma^{dK} \circ \mathcal{M}$ is $(4 \times g + 1) \times n$, where $g = max(\{deg(G), deg(G')\})$, and n is the number of clusters generated by \mathcal{M} .



A microaggregation algorithm partitions *dK*-distribution into clusters and then aggregates frequency values of tuples in each cluster.



- 1. **MDAV-dK algorithm:** Partition *dK*-distribution with a fixed-size constraint such that each cluster has at least *k* tuples.
- 2. **MPDC-dK algorithm:** Partition *dK*-distribution such that every pair of tuples in a cluster satisfies a distance constraint.



Three network datasets:

- (1) *polbooks* contains 105 nodes and 441 edges.
- (2) *ca-GrQc* contains 5,242 nodes and 14,496 edges.
- (3) *ca-HepTh* contains 9,877 nodes and 25,998 edges.

Two measures:

- Euclidean distance measures network structural error between original and perturbed dK-distributions [6].
- Sum of absolute error measures within-cluster homogeneity of clustering algorithms [3].



Does the proposed framework reduce the amount of noise added into dK-distributions while still providing edge-DP guarantee?



Lead to less structural error.

■ Introduce overall less noise to achieve edge-DP.



How do our microaggregation algorithms **perform** in providing better within cluster homogeneity for dK-distributions?

Datasets	Measures	k = 1	k=3	$_{k}$	k = 5		k = 7	
polbooks	SAE	0	144.6	184.67		224.84		
	# Clusters	161	53	32		23		
ca- $GrQc$	SAE	0	1073.3	1476		1810.5		
	# Clusters	1233	411	24	246		176	
ca-HepTh	SAE	0	968.72		304		1599.8	
	# Clusters		1295 431		259		185	
Datasets	Measures	$\tau = 1$	$\tau = 3$		$\tau = 5$		$\tau = 7$	
polbooks	SAE	90.72	2 192.1		5 328.9		424.2	
	# Clusters	68	25	25			8	
ca- $GrQc$	SAE	725.38	3 1732	.1	2630.6		3470.6	
	# Clusters	483	178		98		61	
ca- $HepTh$	SAE	841.87	1761	.8	2773.3		3721.4	
	# Clusters	412	140		73		37	

Produce clusters with less sum of absolute error.

Reduce error due to microaggregation.



dK-Projection Framework: Projection helps to reduce sensitivity by bounding maximum degree in *G*.



■ $\gamma^{2K} \circ \mathcal{P}$ is reduced to $(2 \times \theta + 1) \times \theta$ from $(2 \times deg(G) + 1) \times |E^+|$



The **algorithm** projects a graph using a two-level ordering: (i) global node ordering, and (ii) local neighborhood ordering.

- Assume a sequence of edges ordered by two-level ordering, and let $\theta = 1$.



34



Four network datasets:

- (1) *Facebook* contains 4,039 nodes and 88,234 edges.
- (2) Wiki-Vote contains 7,115 nodes and 103,689 edges.
- (3) *Ca-HepPh* contains 12,008 nodes and 118,521 edges.
- (4) Email-Enron contains 36,692 nodes and 183,831 edges.
- Three utility metrics [1]:
 - Preserved edge ratio measures the ratio of edges being preserved by graph projection.
 - L1 distance measures the network structural error between an original dK-distribution and its perturbed dK-distribution.
 - KS distance quantifies the closeness between an original dK-distribution and its perturbed dK-distribution.



Does the proposed graph projection algorithm yield more **utility** in projected graphs?



Dataset	$\theta = 1$	6	$\theta = 32$	2	$\theta = 64$		
	EAD	SER	EAD	SER	EAD	SER	
Facebook	0.27	0.61	0.44	0.71	0.66	0.84	
Wiki-Vote	0.19	0.59	0.32	0.66	0.50	0.76	
Ca- $HepPh$	0.16	0.61	0.24	0.68	0.31	0.77	
Email-Enron	0.17	0.52	0.22	0.61	0.29	0.71	

- Preserves more edges.
- Leads to less network structural error.
- Generates dK-distributions that are more similar to their original dK-distributions.



Does the proposed graph projection algorithm yield more **utility** in differentially private datasets?



- Yields less network structural error.
- For smaller values of θ, differentially private dK-distributions are more similar to their original dK-distributions.



Limitations of DP:

- Uniform privacy level (i.e., ε) is assigned to each individual while performing perturbation.
- DP may lead to provide insufficient protection for some individuals, while over-protecting others.

Personalized differential privacy provides freedom to individuals to set their own privacy parameter ε .



■ **Problem:** Given a graph G, we want to publish graph statistics under the guarantee of DP while considering personalization.

Key Challenges:

- Graph is a structure of connections between nodes, thus publishing data about one node may leads to violate privacy of others under personalization.
- Each individual (node) has its own privacy preference whereas each entry in data distribution reflects information about more than one node.



We analyze the sensitivity (Δ) of a single *dK*-distribution entry, i.e., degree query γ_q rather than the entire dK-distribution γ^{dK} .

- $\Delta(\gamma_q)$ of is $|E^+| + 1$ over 1K(G) under node-DP.
- $\Delta(\gamma_q)$ of $(deg(G) + 1) \times |E^+|$ over 2K(G) under node-DP.
- $\Delta(\gamma_q)$ of is 2 over 1K(G) under Edge-DP.
- $\Delta(\gamma_q)$ of is $2 \times deg(G) + 1$ over 2K(G) under Edge-DP.

We observe that the sensitivity of γ_q is half as compared to γ^{dK} .



Local Least Based Personalized Perturbation: LL-dK perturbs entries with the strongest local ε .



The frequency value 2 in 1*K*(*G*), and the frequency value 3 in 2*K*(*G*) are perturbed with $\varepsilon = min(\Phi^B, \Phi^F)$, and $\varepsilon = min(\Phi^A, \Phi^C, \Phi^D, \Phi^E)$, respectively.



Threshold Projection Based Personalized Perturbation: TP-dK transforms a graph into a θ -bounded graph then removes all nodes with $\varepsilon < \tau$.



Since $deg(G) \le \theta$, the sensitivity of γ_q is reduced.

With threshold τ all nodes with high privacy are removed.



Sampling Based Personalized Perturbation: ST-dK first splits entries, and then samples them with non-uniform probabilities.



Inclusion probability for each entry depends on corresponding ε and global threshold τ .

Sampled dK-distribution is perturbed with τ .



Aggregation Based Personalized Perturbation: AG-dK computes corresponding ε values to performs aggregation over dK-distribution.



Entries are perturbed with the strongest local ε corresponding to each partition.

 γ_q is approximated to $\gamma_q \circ \mathcal{M}$.



Four network datasets:

- (1) Facebook contains 4,039 nodes and 88,234 edges.
- (2) Wiki-Vote contains 7,115 nodes and 103,689 edges.
- (3) *Ca-HepPh* contains 12,008 nodes and 118,521 edges.
- (4) Email-Enron contains 36,692 nodes and 183,831 edges.

Two utility metrics [1]:

- L1 distance measures the network structural error between an original dK-distribution and its perturbed dK-distribution.
- KS distance quantifies the closeness between an original dK-distribution and its perturbed dK-distribution.

45



Does the proposed personalized approaches yield more **utility** in 1K-distribution under **edge-PDP** and node-PDP?



- Our methods yield less network structural error.
- AG-dK outperforms under edge-PDP and LL-dK outperforms under node-PDP by generating more similar 1K-distributions.



Does the proposed personalized approaches yield yield more **util**ity in 2K-distribution under edge-PDP and node-PDP?



- Our methods yield less network structural error.
- AG-dK outperforms under edge-PDP and LL-dK outperforms under node-PDP by generating more similar 2K-distributions.



What kind of trade-off exists between utility and privacy while generating personalized differentially private dK-distributions?

- The error caused by sensitivity (Δ) and the privacy preference ε dominates the impact on output utility.
- Increasing ε and decreasing Δ can help to reduce error.
- Reducing sensitivity is more challenging under node-PDP than for edge-PDP as graph data is highly sensitive under node-DP.

50



PPDP mechanisms for publishing differentially private data:

- Relational Data: We present novel framework that outperforms the state-of-the-art methods in terms of preserving output utility while guarantee DP.
- Graph Data: We present novel PPDP mechanisms to publish higher-order graph statistics under edge, node and personalized DP while enhancing output utility.



- To publish graph statistics under local differential privacy while considering personalization.
- To develop differentially private mechanisms for continual release of graph statistics in dynamic graphs.
- To release statistics about social groups in a network while protecting privacy of individuals under zero knowledge privacy (ZKP).

REFERENCES



WEI-YEN DAY, NINGHUI LI, AND MIN LYU. PUBLISHING GRAPH DEGREE DISTRIBUTION WITH NODE DIFFERENTIAL PRIVACY. In SIGMOD, pages 123-138, 2016. CYNTHIA DWORK, FRANK MCSHERRY, KOBBI NISSIM, AND ADAM SMITH. CALIBRATING NOISE TO SENSITIVITY IN PRIVATE DATA ANALYSIS. In TCC, pages 265-284, 2006. VLADIMIR ESTIVILL-CASTRO AND JIANHUA YANG. FAST AND ROBUST GENERAL PURPOSE CLUSTERING ALGORITHMS. In PRICAI, pages 208-218, 2000. 1 PRIYA MAHADEVAN, CALVIN HUBBLE, DMITRI KRIOUKOV, BRADLEY HUFFAKER, AND AMIN VAHDAT. ORBIS: RESCALING DEGREE CORRELATIONS TO GENERATE ANNOTATED INTERNET TOPOLOGIES. In SIGCOMM, pages 325-336, 2007. PRIYA MAHADEVAN, DMITRI KRIOUKOV, KEVIN FALL, AND AMIN VAHDAT. SYSTEMATIC TOPOLOGY ANALYSIS AND GENERATION USING DEGREE CORRELATIONS. In SIGCOMM, pages 135-146, 2006. ALESSANDRA SALA, XIAOHAN ZHAO, CHRISTO WILSON, HAITAO ZHENG, AND BEN Y ZHAO. SHARING GRAPHS USING DIFFERENTIALLY PRIVATE GRAPH MODELS. In SIGCOMM, pages 81-98, 2011. IORDI SORIA-COMAS, IOSEP DOMINGO-FERRER, DAVID SÁNCHEZ, AND SERGIO MARTÍNEZ. ENHANCING DATA UTILITY IN DIFFERENTIAL PRIVACY VIA MICROAGGREGATION-BASED K-ANONYMITY. The VLDB Journal, 23(5):771-794, 2014. WILLIAM E YANCEY, WILLIAM E WINKLER, AND ROBERT H CREECY. DISCLOSURE RISK ASSESSMENT IN PERTURBATIVE MICRODATA PROTECTION. In Inference control in statistical databases, pages 135–152. 2002.



THANKS FOR YOUR ATTENTION!

ANY QUESTIONS