

¹ Modular High Performance Computing and Artificial Intelligence, Systems Medicine, German Center for Neurodegenerative Diseases (DZNE), Bonn, North Rhine-Westphalia, Germany ² CISPA Helmholtz Center for Information Security, Saarbrücken, Saarland, Germany



Category: Deep generative model **Attribute Type**: Continuous **DP Sanitization**: Iterative ³ Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio 2014. Generative adversarial nets. Advances in neural information processing systems 27 (2014). **Private-PGM⁴**

Category: Graphical model **Attribute Type**: Discrete only **DP Sanitization**: One-shot ⁴ Ryan McKenna, Daniel Sheldon, and Gerome Miklau. 2019. Graphical-model based estimation and inference for differential privacy. Ir nternational Conference on Machine Learning. PMLR, 4435–4444 **PrivSyn⁵**

Category: Marginal Attribute Type: Discrete only **DP Sanitization**: One-shot ⁵Zhikun Zhang, Tianhao Wang, Ninghui Li, Jean Honorio, Michael Backes, Shibo He, Jiming Chen, and Yang Zhang. 2021. {PrivSyn}: Differentially Private Data Synthesis. In 30th USENIX Security Symposium (USENIX Security 21). 929–946.

METRICS 3

represents the biological characteristics of real-world data.







enes across *n* biological samples

Classification

Downstream Utility

Histogram Intersection Distance to Closest Record

Statistical Properties

Scan Me!

Paper accepted @

PoPETS 2024



Biological Utility

Differential Gene Expression



Gene Co-Expression

GitHub



CONCLUSION 6

We provide the first systematic analysis of non-private and differentially private generation of gene expression data that covers five diverse modeling approaches. Our analysis encompasses a diverse set of metrics that shed light on the quality of the generated data in terms of statistical and biological properties as well as down-stream utility. Key Messages:

- A broad evaluation is necessary in order to understand the limitations of current generators.
- Simple estimators fall behind in performance but equally very complex models like GAN are suffering from the low sample regime typically encountered in bio-medical applications.
- While downstream utility can be strong, the synthetic data itself might not retain statistical nor biological properties.
- Adding privacy guarantees amplifies these problems.

ACKNOWLEDGEMENTS 7

Dingfan Chen was partially supported by Qualcomm Innovation Fellowship Europe.



