

# The Name of the Title Is Hope

Chih-Cheng Rex Yuan

hello@rexyuan.com

Institute of Information Science, Academia Sinica

Taipei, Taiwan

Bow-Yaw Wang

bywang@iis.sinica.edu.tw

Institute of Information Science, Academia Sinica

Taipei, Taiwan

## Abstract

abstract

### ACM Reference Format:

Chih-Cheng Rex Yuan and Bow-Yaw Wang. 2024. The Name of the Title Is Hope. In . ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 Introduction

## 2 Related Work

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference'17, July 2017, Washington, DC, USA*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

### 3 Preliminaries

A row  $r_i$  is a lookup table or dictionary. A database  $D = \{r_1, r_2, \dots\}$  is a collection of rows. The set of all databases is denoted  $\mathcal{D}$ . The attributes of a  $D$  is  $\mathcal{A} = \{A_1, A_2, \dots\}$ . The domain of  $A_i$  is  $\Omega_i$ .

#### 3.1 Fairness Measures

For fairness measures[10, 13], let  $Y$  to denote the ground truth of an outcome, let  $\hat{Y}$  to denote the predicated result of an outcome, let  $S$  denote the protected attribute, and let  $\epsilon$  denote some threshold. For non-binary prediction, such as a score, we use  $\hat{V}$ .

Fairness measures can be broadly categorized into independence, separation, and sufficiency, which are defined by conditional independence in Table 1.

**Table 1: Fairness categories.**

Category	Definition
Independence	$S \perp \hat{Y}$
Separation	$S \perp \hat{Y}   Y$
Sufficiency	$S \perp Y   \hat{Y}$

These categories can be expanded into forms of probability. For example, the definition of separation is expanded to

$$P[\hat{Y} = 1 | S = 1, Y = 1] = P[\hat{Y} = 1 | S \neq 1, Y = 1]$$

$$P[\hat{Y} = 1 | S = 1, Y = 0] = P[\hat{Y} = 1 | S \neq 1, Y = 0]$$

The definition can be relaxed. Its relaxation, for some parameter  $\epsilon$ , is

$$|P[\hat{Y} = 1 | S = 1, Y = 1] - P[\hat{Y} = 1 | S \neq 1, Y = 1]| \leq \epsilon$$

$$|P[\hat{Y} = 1 | S = 1, Y = 0] - P[\hat{Y} = 1 | S \neq 1, Y = 0]| \leq \epsilon$$

which is also the definition of a fairness measure called equalized odds.

We consider in this work various fairness measures listed in Table 2.

#### 3.2 Differential Privacy

A *randomized mechanism* is a randomized algorithm  $M : \mathcal{D} \rightarrow \mathcal{R}$  that takes a database and, after introducing noise, outputs some results.

**Definition 3.1 (Gaussian Mechanism[2]).** Let  $f : \mathcal{D} \rightarrow \mathbb{R}^p$  be a function that takes a database and outputs a vector. The Gaussian Mechanism  $M$  adds i.i.d. Gaussian noise with scale  $\sigma$  to each of the  $p$  outputs:

$$M(D) = f(D) + \mathcal{N}(0, \sigma^2 \mathbb{I})$$

**Definition 3.2 (Rényi Differential Privacy (RDP)).** A randomized mechanism  $M$  satisfies  $(\alpha, \gamma)$ -RDP for  $\alpha \geq 1$  and  $\gamma \geq 1$  if, for all databases  $D_1, D_2$  that differ in exactly one row, we have

$$D_\alpha(M(D_1) || M(D_2)) \leq \gamma$$

where  $D_\alpha$  is the Rényi divergence[12] of order  $\alpha$ .

**THEOREM 3.3 (RDP OF THE GAUSSIAN MECHANISM[3, 8]).** The Gaussian Mechanism satisfies  $(\alpha, \alpha \frac{\Delta_f^2}{2\sigma^2})$ -RDP.

#### 3.3 Differentially Private Synthetic Data

Let  $C \subseteq \mathcal{A}$  be a subset of attributes. Let  $\Omega_C = \prod_{i \in C} \Omega_i$ . A *marginal*[1, 6] of  $C$  is a vector  $\mu \in \mathbb{R}^{|\Omega_C|}$ , indexed by domain element  $t \in \Omega_C$ , such that each entry is a count  $\mu_t = \sum_{x \in D} \mathbb{1}[x_C = t]$  where  $\mathbb{1}$  is the indicator function; that is, it is the vector of the count of each possible element.

Let  $M_C(D)$  be the function that computes the marginal of  $C$  on  $D$ , i.e.,  $\mu = M_C(D)$ . We call marginals of  $|C| = n$  attributes  $n$ -way marginals.

The task of differentially private synthetic data is, given a database  $D$ , adding some noise such that it satisfies differential privacy guarantee and outputting another database  $D'$ , such that the  $L_1$  errors between some selected marginals  $C_1, C_2, \dots$  of  $D$  and  $D'$  is small; that is, their marginals  $(M_{C_1}(D), M_{C_1}(D')), (M_{C_2}(D), M_{C_2}(D')), \dots$  are similar.

For example, suppose we have a database with attributes sex and race. The 2-way marginals of the original database and the synthetic database are shown in Table 3. The marginals of the synthetic data is supposed to be similar to that of the original database.

### 4 Auditing Framework

Our auditing framework is tripartite. It consists of three parties: the data provider, the model maker, and the third-party auditor.

The data provider is responsible for supplying the raw datasets which should originate from trustworthy sources, such as government agencies like a census bureau.

The model maker develops AI models. These are AI companies or research labs specialized in training and optimizing AI models.

The third-party auditor acts as an evaluator, using our framework to audit the AI models for fairness issues by combining both the datasets and the models. These may be investigative journalists or regulatory bodies.

In the framework of our previous work[13], after obtaining real data from the data provider, the 3rd party auditor holds onto the real data for performing fairness audits, and it supposedly retains it indefinitely for the possibility of any future audits.

However, this practice introduces security concerns. It creates a vulnerability to data security threats. A breach at the auditor's end could result in compromises of individuals' privacy.

Moreover, the storage of the datasets also raises privacy concerns. Holding large amounts of sensitive data for an extended period opens the door to the risk of misuse. The auditor may misuse the data for unauthorized purposes.

Thus, we introduce a new framework where the auditor generates synthetic data based on real data upon retrieval of the real data, and then holds onto the synthetic data and discards the real data, preventing further privacy breaches.

### 5 Methodology

We employed the tools of the winner of the 2018 NIST Differential Privacy Synthetic Data Challenge competition[9] by Ryan McKenna[4–7, 11] and the fairness checker tool from our previous research[13].

This research is implemented in Python Jupyter notebooks and is publicly available.

**Table 2: Fairness measures.**

Category	Fairness Measure	Definition
Independence	Disparate Impact	$\frac{P[\hat{Y}=1 S=1]}{P[\hat{Y}=1 S=0]} \geq 1 - \epsilon$
	Demographic Parity	$ P[\hat{Y} = 1 S = 1] - P[\hat{Y} = 1 S \neq 1]  \leq \epsilon$
	Conditional Statistical Parity	$ P[\hat{Y} = 1 S = 1, L = l] - P[\hat{Y} = 1 S \neq 1, L = l]  \leq \epsilon$
	Mean Difference	$ E[\hat{Y} S = 1] - E[\hat{Y} S \neq 1]  \leq \epsilon$
Separation	Equalized Odds	$ P[\hat{Y} = 1 S = 1, Y = 0] - P[\hat{Y} = 1 S \neq 1, Y = 0]  \leq \epsilon$
	Equal Opportunity	$ P[\hat{Y} = 1 S = 1, Y = 1] - P[\hat{Y} = 1 S \neq 1, Y = 1]  \leq \epsilon$
	Predictive Equality	$ P[\hat{Y} = 1 S = 1, Y = 1] - P[\hat{Y} = 1 S \neq 1, Y = 1]  \leq \epsilon$
	Predictive Equality	$ P[\hat{Y} = 1 S = 1, Y = 0] - P[\hat{Y} = 1 S \neq 1, Y = 0]  \leq \epsilon$
Sufficiency	Conditional Use Accuracy Equality	$ P[Y = 1 S = 1, \hat{Y} = 1] - P[Y = 1 S \neq 1, \hat{Y} = 1]  \leq \epsilon$
	Predictive Parity	$ P[Y = 0 S = 1, \hat{Y} = 0] - P[Y = 0 S \neq 1, \hat{Y} = 0]  \leq \epsilon$
	Equal Calibration	$ P[Y = 1 S = 1, \hat{Y} = 1] - P[Y = 1 S \neq 1, \hat{Y} = 1]  \leq \epsilon$
	Equal Calibration	$ P[Y = 1 S = 1, \hat{Y} = 0] - P[Y = 1 S \neq 1, \hat{Y} = 0]  \leq \epsilon$
N/A	Overall Accuracy Equality	$ P[Y = \hat{Y} S = 1] - P[Y = \hat{Y} S \neq 1]  \leq \epsilon$
	Positive Balance	$ E[\hat{Y} Y = 1, S = 1] - E[\hat{Y} Y = 1, S \neq 1]  \leq \epsilon$
	Negative Balance	$ E[\hat{Y} Y = 0, S = 1] - E[\hat{Y} Y = 0, S \neq 1]  \leq \epsilon$

**Table 3: Example marginals.**

(a) Marginal of original data.		(b) Marginal of synthetic data.	
Attributes	Count	Attributes	Count
Male,White	24	Male,White	22
Female,White	33	Female,White	35
Male,Black	13	Male,Black	10
Female,Black	47	Female,Black	46

## 5.1 Data Synthesis

The synthesis framework is three-fold, namely, select-measure-generate[5]. We first select the important marginals to preserve, measure them by adding differentially private noise, and then generate synthetic data.

Underneath the hood, the tool employs a Markov random field (MRF). The select step corresponds to marking cliques in an MRF, and the generate step corresponds to sampling from the fitted MRF.

By default, all 1-way marginals are selected to preserve the quantity of each attribute element. We can further preserve correlations by adding  $n$ -way marginals. For example, if we want to preserve the relationship between sex and race, we may add the clique (sex,race).

In a perfect world where all correlation information is to be preserved, we may wish to make a completely connected graph. However, this was found to be intractable as the complexity of the problem would skyrocket.

To circumvent the complexity explosion, instead, Ryan McKenna devised a technique where the mutual information(MI) of all the database attribute pairs is calculated, and then a maximum spanning tree(MST) algorithm was run with edge weights being the MIs to obtain a skeleton MST MRF.

For the competition, Ryan McKenna further manually added certain cliques based on his investigation of the competition dataset. For example, he manually added the clique (sex,city,income). In

addition, he would add some edges based on some heuristics tailored to that particular dataset.

We developed an alternative heuristic for a general-purpose workflow. Since MIs are bounded by the minimum of the pair's respective Shannon entropies, we add additional edges according to a fraction of these upper bounds.

For the measure step, we followed the examples provided in the tool's repository. Gaussian noises are added to the selected marginals. Half of the privacy budget is spent on all 1-way marginals and the other half on the selected cliques. These marginals are then fed to the tool to fit the MRF. By [6], this procedure satisfies  $(\alpha, \frac{\alpha}{2\sigma^2})$ -RDP for all  $\alpha \geq 1$ .

## 5.2 Fairness Checking

## 5.3 Test Models

## 6 Results

### 6.1 Adult Income Dataset

### 6.2 COMPAS Dataset

### 6.3 One More Dataset

## 7 Discussion

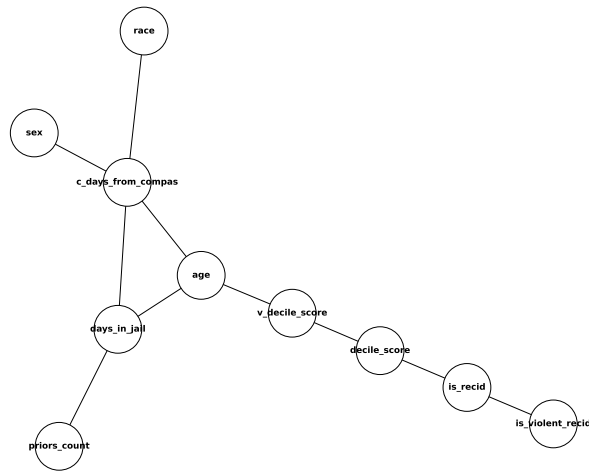
### 7.1 Accuracy

### 7.2 Impossibility

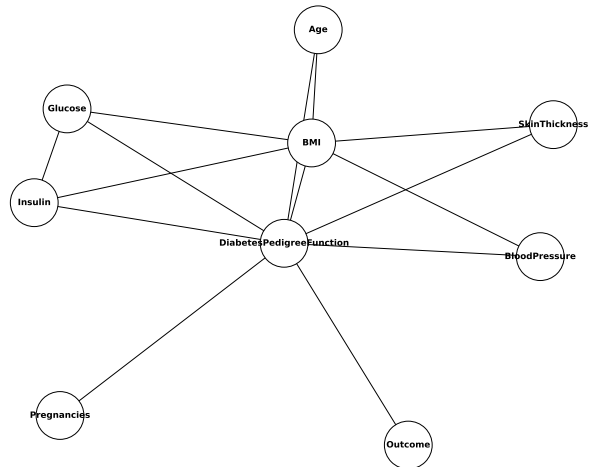
## 8 Conclusion

## References

- [1] Boaz Barak, Kamalika Chaudhuri, Cynthia Dwork, Satyen Kale, Frank McSherry, and Kunal Talwar. 2007. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. 273–282.
- [2] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.



**Figure 1: 1907 Franklin Model D roadster. Photograph by Harris & Ewing, Inc. [Public domain], via Wikimedia Commons. (<https://goo.gl/VLCRBB>).**



**Figure 2: 1907 Franklin Model D roadster. Photograph by Harris & Ewing, Inc. [Public domain], via Wikimedia Commons. (<https://goo.gl/VLCRBB>).**

- [8] Ilya Mironov. 2017. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium (CSF)*. IEEE, 263–275.
- [9] National Institute of Standards and Technology. 2018. 2018 Differential Privacy Synthetic Data Challenge. <https://www.nist.gov/ctl/pscr/open-innovation-prize-challenges/past-prize-challenges/2018-differential-privacy-synthetic> Accessed: 2024-11-10.
- [10] Dana Pessach and Erez Shmueli. 2022. A review on fairness in machine learning. *ACM Computing Surveys (CSUR)* 55, 3 (2022), 1–44.
- [11] Jonathan Ullman. 2022. What is Synthetic Data? <https://differentialprivacy.org/synth-data-0/> Accessed: 2024-11-10.
- [12] Tim Van Erven and Peter Harremoës. 2014. Rényi divergence and Kullback-Leibler divergence. *IEEE Transactions on Information Theory* 60, 7 (2014), 3797–3820.
- [13] Chih-Cheng Rex Yuan and Bow-Yaw Wang. 2024. Ensuring Fairness with Transparent Auditing of Quantitative Bias in AI Systems. In *2024 Pacific Neighborhood Consortium Annual Conference and Joint Meetings (PNC)*. IEEE, 25–32.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

- [3] Vitaly Feldman, Ilya Mironov, Kunal Talwar, and Abhradeep Thakurta. 2018. Privacy amplification by iteration. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE, 521–532.
- [4] Ryan McKenna. 2023. private-pgm: A library for private probabilistic graphical models. <https://github.com/ryan112358/private-pgm> Accessed: 2024-11-10.
- [5] Ryan McKenna and Terrance Liu. 2022. A Simple Recipe for Private Synthetic Data Generation. <https://differentialprivacy.org/synth-data-1/> Accessed: 2024-11-10.
- [6] Ryan McKenna, Gerome Miklau, and Daniel Sheldon. 2021. Winning the NIST Contest: A scalable and general approach to differentially private synthetic data. *arXiv preprint arXiv:2108.04978* (2021).
- [7] Ryan McKenna, Daniel Sheldon, and Gerome Miklau. 2019. Graphical-model based estimation and inference for differential privacy. In *International Conference on Machine Learning*. PMLR, 4435–4444.