

Quantitative Auditing of Fairness Measures with Differentially Private Synthetic Data

ANONYMOUS AUTHOR(S)

abstract

ACM Reference Format:

Anonymous Author(s). 2024. Quantitative Auditing of Fairness Measures with Differentially Private Synthetic Data. 1, 1 (November 2024), 11 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

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.

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

Manuscript submitted to ACM

Manuscript submitted to ACM

3 Preliminaries

Let $\prod_i x_i$ denotes the Cartesian product of x_i s.

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 Ω_i .

Let $X \in \mathcal{X}$ be a discrete random variable and $p(x) := \Pr[X = x]$. Its *marginal Shannon entropy* is $H(X) := -\sum_{x \in \mathcal{X}} p(x) \log p(x)$. It quantifies the level of uncertainty of X . Let $Y \in \mathcal{Y}$ be another discrete random variable. Their *joint Shannon entropy* is $H(X, Y) := -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y)$.

Let X, Y be discrete random variables. Their *mutual information* is $I(X; Y) := H(X) + H(Y) - H(X, Y)$. It quantifies the level of dependence between X and Y . By definition[22], the joint entropy is greater than the marginal entropies $H(X, Y) \geq \max(H(X), H(Y))$. Hence, we have an upper bound of mutual information $I(X; Y) \leq \min(H(X), H(Y))$.

TODO: mrf definition

3.1 Fairness Measures

For fairness measures[19, 26], 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 ϵ 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. $X \perp Y | Z$ denotes the conditional independence between X and Y conditioning on Z .

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 ϵ , 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 Rényi 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.

Manuscript submitted to ACM

Table 2. Fairness measures.

Category	Fairness Measure	Definition
Independence	Disparate Impact	$\frac{P[\hat{Y}=1 S \neq 1]}{P[\hat{Y}=1 S=1]} \geq 1 - \epsilon$
	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$
		$ P[\hat{Y} = 1 S = 1, Y = 1] - P[\hat{Y} = 1 S \neq 1, Y = 1] \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 = 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$
		$ P[Y = 0 S = 1, \hat{Y} = 0] - P[Y = 0 S \neq 1, \hat{Y} = 0] \leq \epsilon$
	Predictive Parity	$ 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{V} = v] - P[Y = 1 S \neq 1, \hat{V} = v] \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{V} Y = 1, S = 1] - E[\hat{V} Y = 1, S \neq 1] \leq \epsilon$
	Negative Balance	$ E[\hat{V} Y = 0, S = 1] - E[\hat{V} Y = 0, S \neq 1] \leq \epsilon$

Definition 3.1 (Gaussian Mechanism[3]). 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 σ 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 (α, γ) -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 $D_\alpha(P \| Q)$ is the Rényi divergence[8, 24, 25] of order α between discrete probability distributions P and Q , defined on \mathcal{X} :

$$D_\alpha(P \| Q) := \frac{1}{\alpha - 1} \log \sum_{x \in \mathcal{X}} P(x)^\alpha Q(x)^{1-\alpha}$$

THEOREM 3.3 (RDP OF THE GAUSSIAN MECHANISM[4, 14]). The Gaussian Mechanism satisfies $(\alpha, \alpha \frac{\Delta_f^2}{2\sigma^2})$ -RDP, where Δ_f denotes the sensitivity[3] of f .

TODO: define sensitivity?

3.3 Differentially Private Synthetic Data

TODO: indicator thing

Let $C \subseteq \mathcal{A}$ be a subset of attributes. Let $\Omega_C = \prod_{i \in C} \Omega_i$. A *marginal*[1, 12] 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[9, 11, 17, 23] is, given a database D , adding some noise to some marginals of D such that it satisfies some differential privacy guarantees and outputting another database D' , such that the L_1 errors between some marginals C_i s of D and D' is small; that is, their marginals $M_{C_i}(D), M_{C_i}(D')$ 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.

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

4 Motivation

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 tool 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[26], 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.

TODO: security => unauthorized access. bulwark database privacy => authorized access, no info leak

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. The third-party thus retains the ability to audit all incoming future models as needed.

5 Methodology

We employed the tools of the winner of the 2018 NIST Differential Privacy Synthetic Data Challenge competition[16] by [10, 12, 13] and the fairness checker tool from our previous research[26].

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

5.1 Data Synthesis

The synthesis framework is three-fold, namely, select-measure-generate[11]. 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. The select step corresponds to marking cliques in a Markov random field, and the generate step corresponds to sampling from the fitted Markov random field.

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, [10, 12] devised a technique where the mutual information of all the database attribute pairs is calculated, and then a maximum spanning tree algorithm was run with edge weights being the mutual information to obtain a skeleton spanning-tree-shaped Markov random field.

For the competition, [10, 12] further manually added certain cliques based on his investigation of the competition dataset. For example, they manually added the clique (sex,city,income). In addition, they would add some edges based on some sophisticated heuristics tailored to that particular dataset. Meanwhile, his other approach where the auditor does access the real dataset does not fit our framework.

We hence developed an alternative heuristic for the general-purpose workflow. As mentioned in Section 3, the mutual information of two random variables is bounded by the pair's respective Shannon entropy. Using this property, we add additional edges with weights exceeding a fraction of the minimum of these upper bounds. As a rule of thumb, we have found setting the fraction to be 0.1 to be effective.

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 Markov random field. By [12], this procedure satisfies $(\alpha, \frac{\alpha}{2\sigma^2})$ -RDP for all $\alpha \geq 1$.

5.2 Fairness Checking

After synthesizing datasets, we used the fairness checker from [26] to compute their fairness measures. To test the viability of our method, we compare the metrics computed from the synthetic dataset against those of the original dataset. We used various datasets with fairness concerns mentioned in [19].

The fairness checker evaluates datasets based on multiple fairness metrics, such as demographic parity and equalized odds. These metrics are computed on some sensitive attributes, predicted outcomes, and ground truths. Examples of sensitive attributes are race and sex. Examples of predicted outcomes and ground truths are loan approval and criminal recidivism.

By comparing these measures between the synthetic and original datasets, we aim to ensure that the synthetic data preserves the fairness properties of the original data. The comparison process is three-fold. It goes as follows.

The dataset is first processed so it can be fed into the synthetic data generator. Some marginals are selected as described in the previous section, and the synthetic data generator model is fitted to the original data according to the marginals. Then the generator is run multiple times to obtain multiple sets of synthetic data.

TODO: section experiment vs Methodology. move results to experiment

Next, several AI models are extracted from various real life authors from Kaggle. They are finetuned to perform well on the original dataset. For one, a random forest model is finetuned by searching hyperparameters settings[5]. For another, a logistic regression model is finetuned by performing principal component analysis[20].

Several AI models and both the original dataset and the rounds of synthetic datasets are fed to the fairness checker. Sensitive attributes are identified based on manual examination with common sense or by referring to [19]. Then, all applicable fairness measures are computed using the checker for both the original and the synthetic.

Finally, we analyze the discrepancies between the fairness properties of the original and the synthetic by calculating the difference and the ratio of their perspective fairness measure values. The sum of the difference and the average of the ratio serve as a summary of the analysis.

6 Experiments

TODO: dataset sizes

We looked at several publicly available datasets, such as adult[2, 21], COMPAS[7, 18], and diabetes[6, 15]. The adult dataset comes from the 1994 census in the United States. The COMPAS dataset comes from an investigative report by ProPublica of the COMPAS criminal recidivism assessment system. The diabetes dataset comes from the hospital readmission data published in the 1994 AI in Medicine journal.

6.1 Adult Income Dataset

For the adult income dataset, the shape of the maximum spanning tree is very shallow, almost resembling a star; it has one internal node and all but one of the leaves have a depth of one. After introducing edges according to our heuristic, we observed an increase in the pairwise edges of the leaves, forming many 3-cliques and two 4-cliques. The resulting graph is shown in Figure 1.

Marginals based on this graph are then passed to the synthetic data generator for model fitting.

After generating ten rounds of synthetic data and passing them to the checker, their fairness measure values are averaged. Then, we compare them against the values of the original data. The results are shown in Table 4.

We observed that across all examined fairness measures, their difference all fall below 0.1. The sum of their differences is 0.259 and the average of their ratios is 0.952, which we consider quite satisfactory.

Table 4. Fairness measures experiment results of the adult dataset. Sum of differences is 0.259. Average of ratios is 0.952.

Measure	Original	Synthetic	Diff	Ratio
Demographic Parity	0.172	0.104	0.067	1.651
Accuracy Equality	0.047	0.117	0.069	0.404
Equalized Odds 1	0.057	0.079	0.021	0.723
Equalized Odds 2	0.166	0.122	0.044	1.365
Accuracy Equality 1	0.100	0.132	0.031	0.759
Accuracy Equality 2	0.119	0.146	0.027	0.814

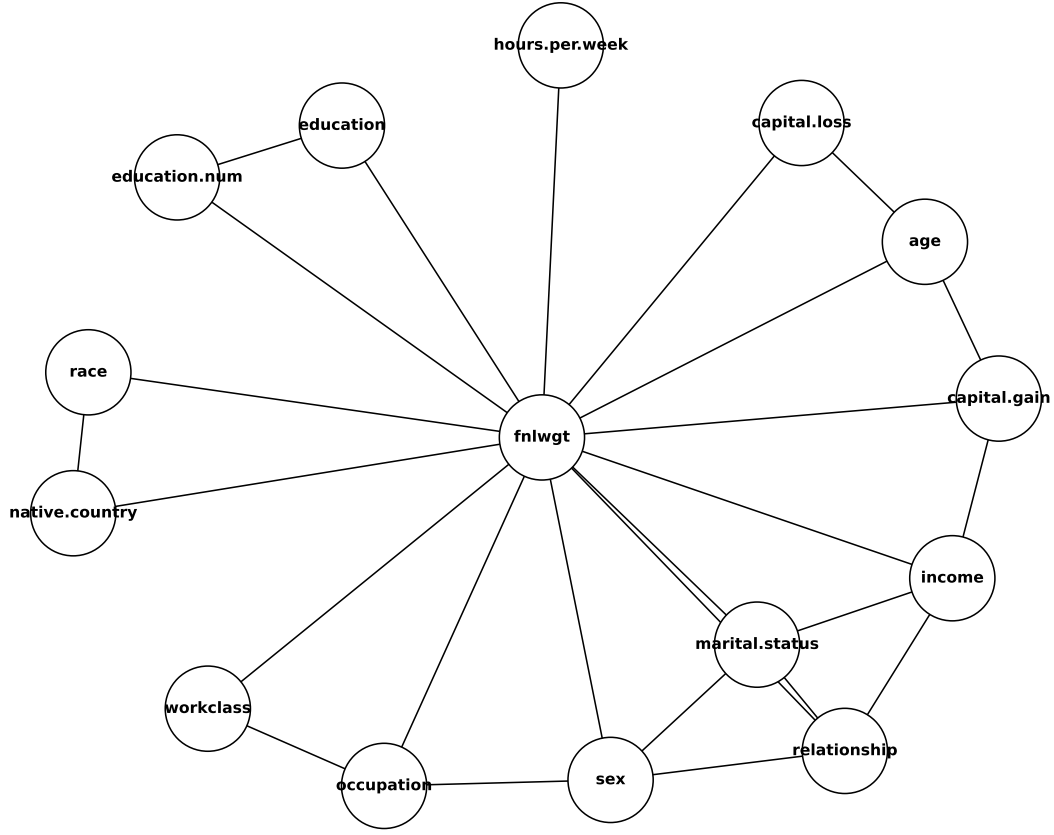


Fig. 1. Markov random field for marginals of the adult dataset.

6.2 COMPAS Dataset

For the COMPAS dataset, the initial spanning tree has a long tail, which is not surprising because, upon closer inspection, they all are related to the original COMPAS risk scores. The heuristic edge addition did not change the graph significantly. It only introduced one 3-clique triangle. The resulting graph is shown in Figure 2.

Marginals of this graph are then too passed to the synthetic data generator for fitting.

We ran the same workflow as adult dataset for the COMPAS dataset. The comparison results are shown in Table 5.

The results showed an increase in error on the sufficiency measure values. In particular, one of the measures has an error ratio as high as 6. The sum of their differences is 0.419, and the average of their ratios is 1.745, which is worse than the adult dataset.

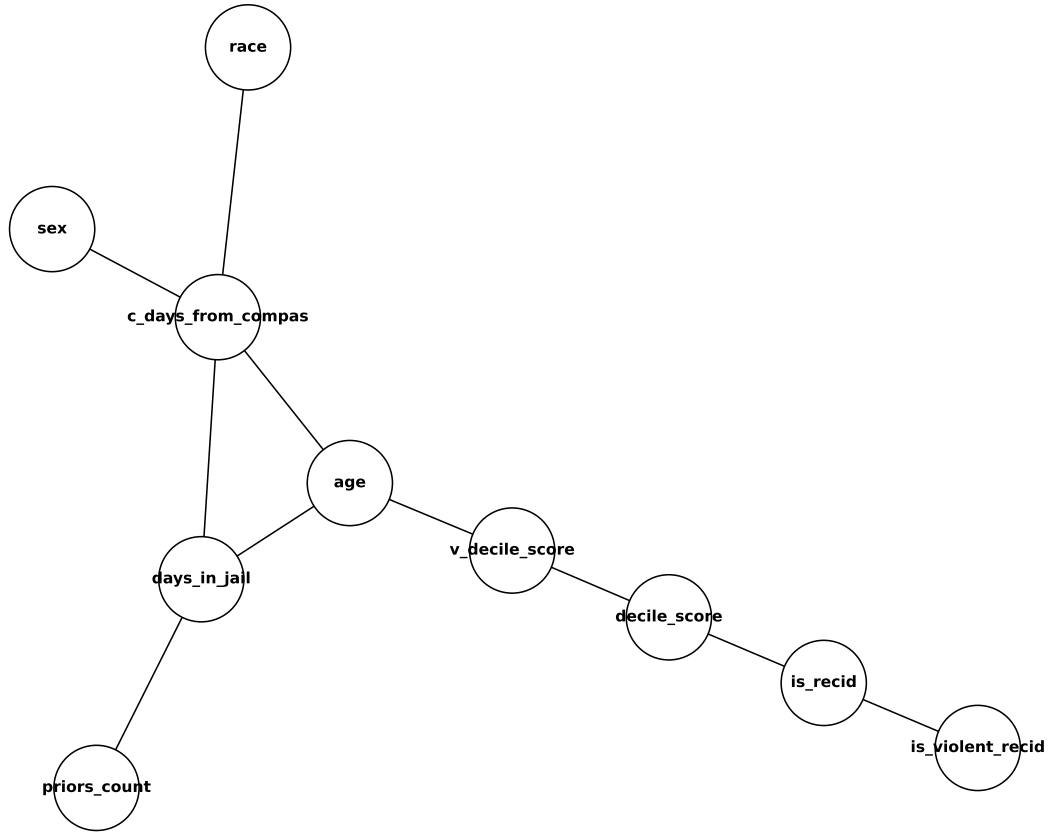


Fig. 2. Markov random field for marginals of the COMPAS dataset.

Table 5. Fairness measures experiment results of the COMPAS dataset. Sum of differences is 0.419. Average of ratios is 1.745.

Measure	Original	Synthetic	Diff	Ratio
Demographic Parity	0.131	0.098	0.032	1.329
Accuracy Equality	0.007	0.013	0.005	0.575
Equalized Odds 1	0.024	0.099	0.074	0.249
Equalized Odds 2	0.017	0.097	0.079	0.182
Accuracy Equality 1	0.170	0.082	0.088	2.082
Accuracy Equality 2	0.169	0.027	0.141	6.058

6.3 Diabetes Dataset

The tree grown from the diabetes dataset did not appear to have any particular characteristics. The root of the tree is placed in the BMI value, which reasonably captures most information. The heuristic edge addition process introduced some 3-clique triangles. The resulting graph is shown in Figure 3.

The same process is conducted to fit the synthetic data generator model.

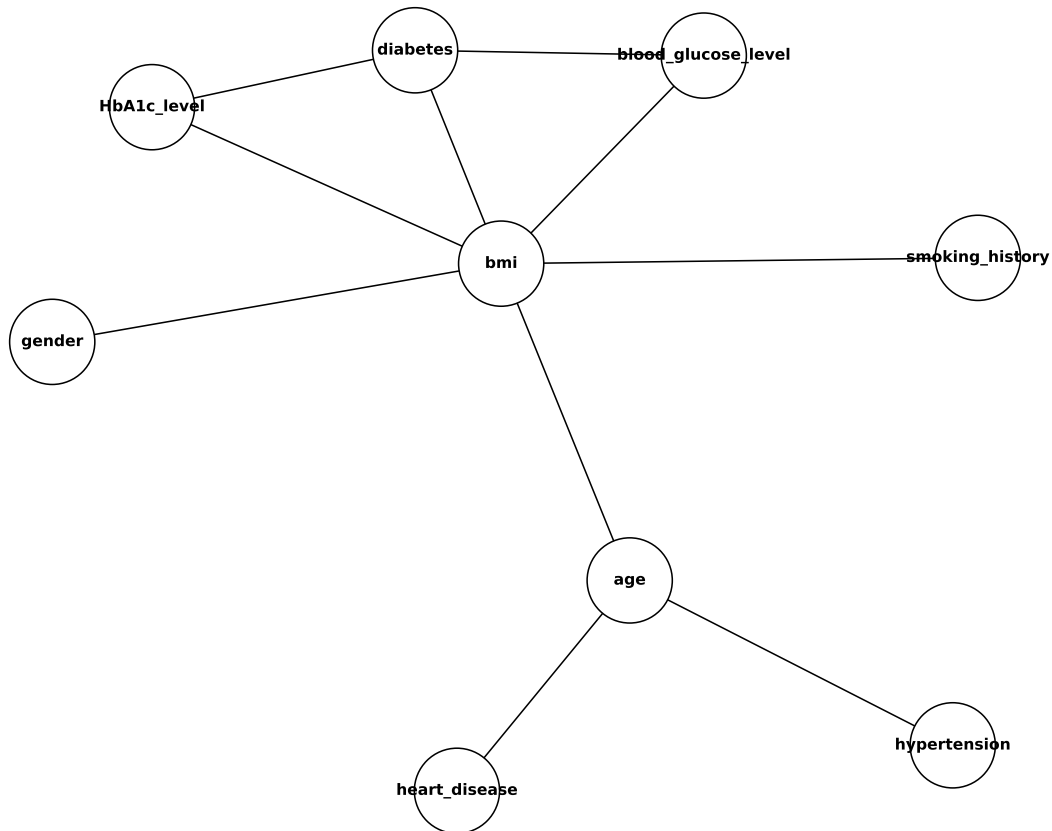


Fig. 3. Markov random field for marginals of the diabetes dataset.

The same workflow was done as on previous datasets. The comparison results are shown in Table 6.

There was also an increase of error on some sufficiency measure values. The sum of their differences is 0.189 and the average of their ratios is 0.489. However, this comparison may be skewed because one of the fairness measures in the original data was calculated to be zero, thus making the average ratio seem lower than otherwise.

TODO: only compare diff. remove ratios. use avg of diff.

Table 6. Fairness measures experiment results of the COMPAS dataset. Sum of differences is 0.189. Average of ratios is 0.489.

Measure	Original	Synthetic	Diff	Ratio
Demographic Parity	0.013	0.015	0.002	0.867
Accuracy Equality	0.007	0.009	0.002	0.782
Equalized Odds 1	0.000	0.003	0.003	0.000
Equalized Odds 2	0.008	0.106	0.097	0.081
Accuracy Equality 1	0.013	0.097	0.084	0.136
Accuracy Equality 2	0.021	0.020	0.001	1.068

7 Evaluation

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] Barry Becker and Ronny Kohavi. 1996. Adult. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5XW20>.
- [3] 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.
- [4] 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.
- [5] Ipbyrne. [n. d.]. Income Prediction (84.369% Accuracy). <https://www.kaggle.com/code/ipbyrne/income-prediction-84-369-accuracy> Accessed: 2024-11-10.
- [6] Michael Kahn. [n. d.]. Diabetes. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5T59G>.
- [7] Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. 2016. ProPublica Compas Analysis—Data and Analysis for ‘Machine Bias’. <https://github.com/propublica/compas-analysis>.
- [8] Yingzhen Li and Richard E Turner. 2016. Rényi divergence variational inference. *Advances in neural information processing systems* 29 (2016).
- [9] Ryan McKenna. 2021. Marginal-based Methods for Differentially Private Synthetic Data. <https://www.youtube.com/watch?v=UKzh9QgNRxA>. Accessed: 2024-11-26.
- [10] Ryan McKenna. 2023. private-pgm: A library for private probabilistic graphical models. <https://github.com/ryan112358/private-pgm> Accessed: 2024-11-10.
- [11] Ryan McKenna and Terrance Liu. 2022. A Simple Recipe for Private Synthetic Data Generation. <https://differentialprivacy.org/synth-data-1/> Accessed: 2024-11-10.
- [12] 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).
- [13] 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.
- [14] Ilya Mironov. 2017. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium (CSF)*. IEEE, 263–275.
- [15] Tz Mustafa. [n. d.]. Diabetes Prediction Dataset. <https://www.kaggle.com/datasets/iammustafatz/diabetes-prediction-dataset> Accessed: 2024-11-26.
- [16] 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.
- [17] Joseph Near and David Darais. 2021. Differentially Private Synthetic Data. <https://www.nist.gov/blogs/cybersecurity-insights/differentially-private-synthetic-data>. Accessed: 2024-11-26.
- [18] Dan Ofer. [n. d.]. COMPAS Dataset. <https://www.kaggle.com/datasets/danofer/compass> Accessed: 2024-11-26.
- [19] Dana Pessach and Erez Shmueli. 2022. A review on fairness in machine learning. *ACM Computing Surveys (CSUR)* 55, 3 (2022), 1–44.
- [20] Prashant111. [n. d.]. EDA, Logistic Regression, PCA. [https://www.kaggle.com/code/prashant111/eda-logistic-regression-pca#Introduction-to-Principal-Component-Analysis-\(PCA\)](https://www.kaggle.com/code/prashant111/eda-logistic-regression-pca#Introduction-to-Principal-Component-Analysis-(PCA)) Accessed: 2024-11-10.

- [21] UC Irvine Machine Learning Repository. [n. d.]. Adult Census Income Dataset. <https://www.kaggle.com/datasets/uciml/adult-census-income/data> Accessed: 2024-11-26.
- [22] Madhur Tulsiani and Shubhendu Trivedi. 2014. Information and Coding Theory Lecture 2. <https://home.ttic.edu/~madhurt/courses/infotheory2014/l2.pdf> Accessed: 2024-11-29.
- [23] Jonathan Ullman. 2022. What is Synthetic Data? <https://differentialprivacy.org/synth-data-0/> Accessed: 2024-11-10.
- [24] Tim van Erven and Peter Harremoës. 2010. Rényi divergence and majorization. In *2010 IEEE International Symposium on Information Theory*. IEEE, 1335–1339.
- [25] Tim Van Erven and Peter Harremos. 2014. Rényi divergence and Kullback-Leibler divergence. *IEEE Transactions on Information Theory* 60, 7 (2014), 3797–3820.
- [26] 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.