# A  Reviewer 2 Experimental Results

## A.1  Additional Dataset and Metrics

**Texas.**  The Texas Hospital Discharge dataset[1] is a large public use data file provided by the Texas Department of State Health Services. We used the preprocessed which consists of 100,000 records uniformly selected from a pre-processed file containing patient data from 2013[2] version from Stadler et al. (2022). We retain 18 attributes and assume a binary classification task by predicting only minor and major mortality risk following the setup of Afonja et al. (2023). Duplicates where also removed. The final size of the dataset was therefore reduced to 75,105 which was split to non-overlapping train/test/validation. Validation size is fixed to 1000 for all dataset.

|         | # Train | # Test | # N | # C | # Total | Class Ratio |
|---------|---------|--------|-----|-----|---------|-------------|
| Adult   | 30932   | 16858  | 6   | 9   | 15      | 76:24       |
| Airline | 103904  | 24976  | 19  | 5   | 24      | 57:43       |
| Texas   | 60127   | 13978  | 7   | 11  | 18      | 81:19       |

Table 1: **Tabular Dataset statistics.** # N and # C are the numbers of numerical and categorical columns, respectively.

**Pairwise Correlation Similarity Accuracy (CorAcc).**  We evaluate the correlation between data columns using the approach described by Tao et al. (2021) and Afonja et al. (2023). Specifically, we use Cramer's V with bias correction for categorical columns, the Correlation Ratio for numerical-categorical columns, and the Pearson Correlation Coefficient (absolute values) for numerical columns. The ranges for these measures are as follows: Cramer's V and Correlation Ratio are bounded between 0 and 1, while the Pearson Correlation Coefficient spans -1 to 1. Following Tao et al. (2021), correlation values are discretized into four levels: low $[0, 0.1)$, weak $[0.1, 0.3)$, medium $[0.3, 0.5)$, and strong $[0.5, 1)$. The *CorAcc* metric quantifies the similarity between synthetic and original data by measuring the fraction of column pairs where the assigned correlation levels match.

**Pairwise Attribute Distribution Similarity (Pair).**  This metric extends the Normalized Histogram Intersection (HIST) by calculating the histogram intersection for all two-way marginals and averaging the results across all attribute pairs. For numerical columns, we discretize the values into bins of size 20 and 50 before computing the intersections.

|                  | Method | **Adult** | | | | | | **Airline** | | | | | | **Texas** | | | | | |
|------------------|--------|------|------|------|--------|------|------|------|------|------|--------|------|------|------|------|------|--------|------|------|
|                  |        | **F1** | **AUC** | **ACC** | **CorAcc** | **Pair** | **HIST** | **F1** | **AUC** | **ACC** | **CorAcc** | **Pair** | **HIST** | **F1** | **AUC** | **ACC** | **CorAcc** | **Pair** | **HIST** |
|                  | Real data | $69.9_{\pm2}$ | $91.7_{\pm1}$ | $84.0_{\pm4}$ | $97.3$ | $97.5_{\pm0}$ | $99.1_{\pm0}$ | $90.6_{\pm8}$ | $96.2_{\pm5}$ | $91.8_{\pm7}$ | $97.2$ | $98.4_{\pm0}$ | $99.4_{\pm0}$ | $86.6_{\pm2}$ | $98.8_{\pm0}$ | $95.0_{\pm1}$ | $96.3_{\pm0}$ | $98.9_{\pm0}$ | $99.5_{\pm0}$ |
| $\varepsilon=\infty$ | CTGAN | $59.5_{\pm6}$ | $\underline{88.5_{\pm0}}$ | $80.2_{\pm3}$ | $74.9_{\pm4}$ | $\mathbf{85.0_{\pm1}}$ | $\underline{91.2_{\pm1}}$ | $\underline{87.2_{\pm3}}$ | $\underline{94.7_{\pm2}}$ | $\underline{88.9_{\pm3}}$ | $\underline{84.4_{\pm1}}$ | $\mathbf{89.3_{\pm1}}$ | $\mathbf{94.4_{\pm1}}$ | $\underline{80.8_{\pm4}}$ | $\underline{96.8_{\pm2}}$ | $\underline{92.8_{\pm2}}$ | $\mathbf{72.8_{\pm4}}$ | $\mathbf{88.3_{\pm1}}$ | $93.1_{\pm0}$ |
|                  | TVAE   | $\underline{63.2_{\pm2}}$ | $87.5_{\pm1}$ | $77.7_{\pm4}$ | $\underline{76.6_{\pm2}}$ | $\underline{84.5_{\pm1}}$ | $\mathbf{91.5_{\pm1}}$ | $85.8_{\pm5}$ | $93.0_{\pm6}$ | $87.2_{\pm6}$ | $78.0_{\pm3}$ | $82.8_{\pm3}$ | $90.3_{\pm2}$ | $76.7_{\pm7}$ | $95.3_{\pm3}$ | $90.9_{\pm3}$ | $\underline{61.2_{\pm4}}$ | $63.4_{\pm7}$ | $\underline{93.5_{\pm1}}$ |
|                  | VAE    | $53.8_{\pm10}$ | $86.6_{\pm1}$ | $\underline{80.5_{\pm1}}$ | $65.3_{\pm2}$ | $60.2_{\pm2}$ | $73.3_{\pm3}$ | $79.8_{\pm1}$ | $91.1_{\pm1}$ | $80.0_{\pm1}$ | $67.8_{\pm5}$ | $60.4_{\pm1}$ | $76.8_{\pm0}$ | $74.3_{\pm3}$ | $95.6_{\pm1}$ | $90.6_{\pm2}$ | $51.5_{\pm7}$ | $57.5_{\pm1}$ | $89.3_{\pm1}$ |
|                  | GPT-2  | $\mathbf{68.9_{\pm0}}$ | $\mathbf{90.7_{\pm0}}$ | $\mathbf{83.7_{\pm2}}$ | $\mathbf{79.9_{\pm1}}$ | $83.7_{\pm0}$ | $90.7_{\pm1}$ | $\mathbf{89.6_{\pm5}}$ | $\mathbf{95.9_{\pm3}}$ | $\mathbf{91.4_{\pm0}}$ | $\mathbf{86.5_{\pm4}}$ | $\underline{85.5_{\pm2}}$ | $\underline{90.8_{\pm1}}$ | $\mathbf{83.9_{\pm2}}$ | $\mathbf{98.6_{\pm0}}$ | $\mathbf{93.6_{\pm1}}$ | $\underline{72.8_{\pm1}}$ | $\underline{82.8_{\pm3}}$ | $\mathbf{94.9_{\pm0}}$ |

Table 2: **Non-DP Benchmark.** The Real data baseline represents the optimal achievable performance, determined by evaluating metrics using real training data compared to real test data. **F1**, **AUC** and **ACC** are reported as averages of two ML Models (XGBoost and Logistic Regression). **Pair**, and **Hist** are reported as averages of two bin sizes (Bins 20 and 50). Each model run five times and four synthetic datasets generated per run with standard deviation reported after $\pm$. The best value per column for each $\varepsilon$ is shown in **bold** while the second best value is underlined.

## A.2  Marginal-based DP Baselines

**AIM.**  Proposed by (McKenna et al., 2022), AIM is a marginal-based model for generating differentially private synthetic data. It is a workload-adaptive algorithm that follows a three-step process: selecting a set

---

[1] https://www.dshs.texas.gov/thcic/
[2] https://github.com/spring-epfl/synthetic_data_release/blob/master/data/texas.csv

of queries, privately measuring those queries, and generating synthetic data from the noisy measurements. AIM employs innovative techniques to iteratively prioritize the most useful measurements, considering both their relevance to the workload and their importance in approximating the input data.

**MST.** Proposed by (McKenna et al., 2021), MST was the winning mechanism of the 2018 NIST Differential Privacy Synthetic Data Competition. It is a general approach for differentially private synthetic data generation that follows three main steps: (1) selecting a collection of low-dimensional marginals, (2) measuring these marginals using a noise addition mechanism, and (3) generating synthetic data that accurately preserves the measured marginals.

| | Method | Adult | | | | | | Airline | | | | | | Texas | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | F1 | AUC | ACC | CorAcc | Pair | HIST | F1 | AUC | ACC | CorAcc | Pair | HIST | F1 | AUC | ACC | CorAcc | Pair | HIST |
| $\varepsilon=1$, | AIM | $59.6_{\pm6}$ | $86.8_{\pm1}$ | $80.3_{\pm2}$ | $86.4_{\pm1}$ | $77.1_{\pm9}$ | $88.4_{\pm5}$ | $77.3_{\pm5}$ | $88.9_{\pm4}$ | $78.2_{\pm5}$ | $91.8_{\pm1}$ | $46.7_{\pm3}$ | $68.1_{\pm2}$ | $84.5_{\pm1}$ | $98.3_{\pm0}$ | $94.2_{\pm1}$ | $81.0_{\pm2}$ | $93.0_{\pm5}$ | $98.9_{\pm0}$ |
| $\delta=10^{-5}$ | MST | $39.6_{\pm19}$ | $76.8_{\pm1}$ | $72.8_{\pm2}$ | $70.0_{\pm1}$ | $74.6_{\pm10}$ | $87.0_{\pm5}$ | $72.2_{\pm6}$ | $83.4_{\pm5}$ | $75.2_{\pm4}$ | $72.7_{\pm0}$ | $46.3_{\pm3}$ | $68.2_{\pm2}$ | $81.7_{\pm0}$ | $94.8_{\pm0}$ | $93.2_{\pm0}$ | $77.0_{\pm0}$ | $97.5_{\pm0}$ | $99.0_{\pm0}$ |
| | DP-GAN | $\underline{33.5}_{\pm20}$ | $\underline{67.7}_{\pm9}$ | $64.2_{\pm10}$ | $39.9_{\pm3}$ | $41.2_{\pm4}$ | $63.7_{\pm3}$ | $40.2_{\pm24}$ | $63.9_{\pm13}$ | $59.8_{\pm6}$ | $37.4_{\pm9}$ | $22.2_{\pm13}$ | $44.7_{\pm12}$ | $13.7_{\pm16}$ | $58.3_{\pm12}$ | $78.3_{\pm8}$ | $36.1_{\pm7}$ | $34.6_{\pm5}$ | $68.3_{\pm7}$ |
| | DP-CTGAN | $\mathbf{42.2}_{\pm20}$ | $\mathbf{78.0}_{\pm7}$ | $\mathbf{75.7}_{\pm3}$ | $51.3_{\pm3}$ | $59.2_{\pm2}$ | $75.7_{\pm2}$ | $\underline{67.1}_{\pm8}$ | $\underline{76.7}_{\pm8}$ | $\underline{68.0}_{\pm6}$ | $31.7_{\pm2}$ | $62.2_{\pm2}$ | $78.7_{\pm2}$ | $63.9_{\pm11}$ | $91.4_{\pm3}$ | $82.6_{\pm19}$ | $43.9_{\pm6}$ | $\underline{66.9}_{\pm6}$ | $84.7_{\pm3}$ |
| | DP-VAE | $0.0_{\pm0}$ | $50.0_{\pm0}$ | $\underline{75.6}_{\pm0}$ | $48.8_{\pm1}$ | $40.3_{\pm1}$ | $61.8_{\pm2}$ | $26.5_{\pm28}$ | $57.9_{\pm13}$ | $57.3_{\pm6}$ | $46.6_{\pm1}$ | $20.6_{\pm0}$ | $41.8_{\pm1}$ | $0.0_{\pm0}$ | $50.0_{\pm0}$ | $82.5_{\pm0}$ | $62.1_{\pm1}$ | $43.9_{\pm1}$ | $77.9_{\pm1}$ |
| | *GPT-2* | | | | | | | | | | | | | | | | | | |
| | DP-Standard | $27.8_{\pm15}$ | $58.5_{\pm7}$ | $65.2_{\pm9}$ | $55.0_{\pm1}$ | $68.4_{\pm1}$ | $85.7_{\pm2}$ | $60.5_{\pm7}$ | $65.3_{\pm9}$ | $62.4_{\pm7}$ | $64.0_{\pm2}$ | $77.0_{\pm2}$ | $90.3_{\pm3}$ | $55.4_{\pm10}$ | $90.2_{\pm5}$ | $77.1_{\pm11}$ | $\underline{70.3}_{\pm1}$ | $60.6_{\pm1}$ | $92.3_{\pm1}$ |
| | DP-2Stage-U | $21.2_{\pm12}$ | $48.9_{\pm6}$ | $61.9_{\pm13}$ | $55.0_{\pm1}$ | $\mathbf{76.1}_{\pm1}$ | $86.7_{\pm1}$ | $\mathbf{68.5}_{\pm9}$ | $\mathbf{77.8}_{\pm10}$ | $\mathbf{72.1}_{\pm7}$ | $65.3_{\pm1}$ | $\mathbf{80.8}_{\pm1}$ | $\underline{90.7}_{\pm1}$ | $23.5_{\pm14}$ | $59.8_{\pm14}$ | $67.3_{\pm17}$ | $68.2_{\pm0}$ | $\mathbf{80.7}_{\pm6}$ | $\mathbf{93.4}_{\pm0}$ |
| | DP-2Stage-O | | | | | | | | | | | | | | | | | | |
| | +adult | - | - | - | - | - | - | $55.2_{\pm18}$ | $62.5_{\pm19}$ | $60.0_{\pm16}$ | $\mathbf{66.8}_{\pm1}$ | $\underline{80.1}_{\pm1}$ | $\mathbf{92.5}_{\pm1}$ | $\mathbf{74.8}_{\pm4}$ | $\mathbf{96.7}_{\pm1}$ | $\mathbf{89.1}_{\pm3}$ | $69.9_{\pm2}$ | $60.5_{\pm2}$ | $91.7_{\pm1}$ |
| | +airline | $30.4_{\pm17}$ | $61.6_{\pm8}$ | $66.7_{\pm8}$ | $\underline{55.4}_{\pm1}$ | $\underline{72.3}_{\pm1}$ | $\mathbf{88.5}_{\pm1}$ | - | - | - | - | - | - | $\underline{74.3}_{\pm5}$ | $\underline{96.4}_{\pm0}$ | $\underline{88.8}_{\pm3}$ | $\mathbf{70.9}_{\pm1}$ | $62.0_{\pm2}$ | $\underline{93.2}_{\pm0}$ |
| | +texas | $31.6_{\pm13}$ | $60.5_{\pm7}$ | $66.4_{\pm8}$ | $\mathbf{55.6}_{\pm1}$ | $71.3_{\pm1}$ | $\underline{86.9}_{\pm1}$ | $52.5_{\pm13}$ | $61.0_{\pm13}$ | $58.4_{\pm10}$ | $\underline{66.1}_{\pm2}$ | $78.7_{\pm2}$ | $90.4_{\pm1}$ | - | - | - | - | - | - |

Table 3: **DP Benchmark.** Utility metrics (F1, AUC, and ACC) are presented as the averages of two ML Models (XGBoost and Logistic Regression). **Pair**, and **Hist** are reported as averages of two bin sizes (Bins 20 and 50). Results are averaged across five model runs and four synthetic datasets per run with standard deviation reported after ±. The best value per column (excluding the marginal-based metrics, AIM and MST) for $\varepsilon=1$ is shown in **bold** while second best value is underlined.