

A Reviewer 1 Experimental Results

Requested Changes. “I also would like to make the following request. If the experimental results are convincing, I will recommend the paper for publication. Could the authors also evaluate the data synthetic methods using pairwise metrics, such as ‘Pair’ and ‘Corr’ in this paper. Could the authors also evaluate previous DP-focused tabular data synthetic methods such as DPSyn, PrivMRF and Private-PGM?”

A.1 Additional Metrics

		Adult				Airline			
	Method	F1	CorAcc	Pair	Hist	F1	CorAcc	Pair	Hist
$\varepsilon = \infty$	Real data	69.9 \pm 2	95.6	97.7 \pm 0	99.1 \pm 0	90.6 \pm 8	94.4	98.4 \pm 0	99.4 \pm 0
	CTGAN	59.5 \pm 6	72.7 \pm 3	85.0\pm1	<u>91.2\pm1</u>	87.2 \pm 3	74.2 \pm 1	89.3\pm1	94.4\pm1
	TVAE	<u>63.2\pm2</u>	75.1 \pm 3	<u>84.5\pm1</u>	91.5\pm1	85.8 \pm 5	70.3 \pm 3	82.8 \pm 3	90.3 \pm 2
	VAE	53.8 \pm 10	64.0 \pm 1	60.2 \pm 2	73.3 \pm 3	79.8 \pm 1	62.8 \pm 6	60.4 \pm 1	76.8 \pm 0
	GPT-2	68.9\pm0	79.4\pm1	83.7 \pm 0	90.7 \pm 1	89.6\pm5	79.8\pm5	<u>85.5\pm2</u>	<u>90.8\pm1</u>

Table 1: **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** is 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 synthetic data generator is run five times and four synthetic datasets generated per run. Standard deviation reported after \pm . The best value per column is shown in **bold** while the second best value is underlined.

		Adult				Airline			
	Method	F1	CorAcc	Pair	Hist	F1	CorAcc	Pair	Hist
$\varepsilon = 1,$ $\delta = 10^{-5}$	DP-GAN	33.5 \pm 20	39.1 \pm 3	41.2 \pm 4	63.7 \pm 3	40.2 \pm 24	37.5 \pm 5	22.2 \pm 13	44.7 \pm 12
	DP-CTGAN	42.2\pm20	49.3 \pm 3	59.2 \pm 2	75.7 \pm 2	<u>67.1\pm8</u>	31.6 \pm 3	62.2 \pm 2	78.7 \pm 2
	DP-VAE	0.0 \pm 0	45.2 \pm 1	40.3 \pm 1	61.8 \pm 2	26.5 \pm 28	42.4 \pm 1	20.6 \pm 0	41.8 \pm 1
	<i>GPT-2</i>								
	DP-Standard	27.8 \pm 15	49.6\pm2	68.4 \pm 1	85.7 \pm 2	60.5 \pm 7	49.9 \pm 2	77.0 \pm 2	90.3 \pm 3
	DP-2Stage-U	21.2 \pm 12	49.6\pm2	76.1\pm1	<u>86.7\pm1</u>	68.5\pm9	<u>51.7\pm2</u>	80.8\pm1	<u>90.7\pm1</u>
	DP-2Stage-O	30.4 \pm 17	<u>49.5\pm2</u>	<u>72.3\pm1</u>	88.5\pm1	55.2 \pm 18	52.1\pm2	<u>80.1\pm1</u>	92.5\pm1
$\varepsilon = 8,$ $\delta = 10^{-5}$	DP-GAN	19.6 \pm 20	42.9 \pm 3	11.1 \pm 7	33.3 \pm 9	-	-	-	-
	DP-CTGAN	46.5\pm18	49.0 \pm 2	65.6 \pm 4	80.0 \pm 2	<u>67.7\pm4</u>	31.2 \pm 4	54.7 \pm 2	76.8 \pm 1
	DP-VAE	0.0 \pm 0	45.2 \pm 1	40.9 \pm 1	62.1 \pm 1	51.9 \pm 25	44.3 \pm 2	18.3 \pm 1	40.0 \pm 1
	<i>GPT-2</i>								
	DP-Standard	31.3 \pm 15	52.0\pm1	70.6 \pm 1	84.5 \pm 1	64.9 \pm 6	51.8 \pm 3	78.2 \pm 3	89.8 \pm 3
	DP-2Stage-U	22.4 \pm 15	51.4 \pm 2	77.1\pm1	<u>86.9\pm1</u>	71.9\pm9	56.1\pm2	<u>80.9\pm1</u>	<u>90.4\pm1</u>
	DP-2Stage-O	<u>33.4\pm16</u>	<u>51.7\pm2</u>	<u>74.5\pm1</u>	87.9\pm1	64.2 \pm 11	<u>55.1\pm3</u>	81.1\pm1	92.3\pm1

Table 2: **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** is 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 synthetic data generator is run five times and four synthetic datasets generated per run. Standard deviation reported after \pm . The best value per column for each ε is shown in **bold** while second best value is underlined.

A.2 Additional Baselines

		Adult				Airline			
	Method	F1	CorAcc	Pair	Hist	F1	CorAcc	Pair	Hist
$\varepsilon = 1,$ $\delta = 10^{-5}$	<i>Marginal</i>								
	AIM	59.6 ± 6	82.7 ± 1	77.1 ± 9	<u>88.4</u> ± 5	77.3 ± 5	83.7 ± 2	46.7 ± 3	68.1 ± 2
	MST	39.6 ± 19	<u>65.1</u> ± 1	74.6 ± 10	87.0 ± 5	<u>72.2</u> ± 6	<u>59.9</u> ± 2	46.3 ± 3	68.2 ± 2
	<i>DNN</i>								
	DP-GAN	33.5 ± 20	39.1 ± 3	41.2 ± 4	63.7 ± 3	40.2 ± 24	37.5 ± 5	22.2 ± 13	44.7 ± 12
	DP-CTGAN	<u>42.2</u> ± 20	49.3 ± 3	59.2 ± 2	75.7 ± 2	67.1 ± 8	31.6 ± 3	62.2 ± 2	78.7 ± 2
	DP-VAE	0.0 ± 0	45.2 ± 1	40.3 ± 1	61.8 ± 2	26.5 ± 28	42.4 ± 1	20.6 ± 0	41.8 ± 1
	<i>GPT-2</i>								
	DP-Standard	27.8 ± 15	49.6 ± 2	68.4 ± 1	85.7 ± 2	60.5 ± 7	49.9 ± 2	77.0 ± 2	90.3 ± 3
	DP-2Stage-U	21.2 ± 12	49.6 ± 2	<u>76.1</u> ± 1	86.7 ± 1	68.5 ± 9	51.7 ± 2	80.8 ± 1	<u>90.7</u> ± 1
	DP-2Stage-O	30.4 ± 17	49.5 ± 2	72.3 ± 1	88.5 ± 1	55.2 ± 18	52.1 ± 2	<u>80.1</u> ± 1	92.5 ± 1

Table 3: **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** is 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 synthetic data generator is run five times and four synthetic datasets generated per run. Standard deviation reported after \pm . The best value per column for each ε is shown in **bold** while second best value is underlined.

Dataset			XGB (F1)	XGB (AUC)	XGB (ACC)	LR (F1)	LR (AUC)	LR (ACC)	HIST (bin=50)	HIST (bin=20)
$\varepsilon = 1$ $\delta = 10^{-5}$	Adult	AIM	54.0 ± 4	86.3 ± 1	81.8 ± 0	65.2 ± 0	87.3 ± 0	78.8 ± 1	83.5 ± 0	93.3 ± 1
		MST	20.5 ± 3	76.7 ± 1	74.5 ± 0	<u>58.7</u> ± 0	77.0 ± 2	71.0 ± 0	81.8 ± 2	<u>92.2</u> ± 0
		DP-GAN	27.4 ± 24	67.6 ± 10	69.1 ± 9	39.6 ± 14	67.8 ± 9	59.4 ± 8	61.9 ± 2	65.5 ± 3
		DP-CTGAN	38.6 ± 21	<u>77.2</u> ± 7	<u>76.0</u> ± 3	45.8 ± 19	<u>78.8</u> ± 7	75.3 ± 3	75.0 ± 2	76.4 ± 2
		DP-VAE	0.0 ± 0	50.0 ± 0	75.6 ± 0	0.0 ± 0	50.0 ± 0	<u>75.6</u> ± 0	60.1 ± 0	63.5 ± 0
		DP-Standard	13.9 ± 7	55.5 ± 6	73.0 ± 1	41.6 ± 5	61.4 ± 7	57.4 ± 6	85.1 ± 2	86.2 ± 2
		DP-2Stage-U	10.3 ± 5	48.4 ± 4	74.5 ± 1	32.1 ± 6	49.3 ± 8	49.4 ± 4	86.3 ± 1	87.1 ± 1
		DP-2Stage-O	15.0 ± 7	55.9 ± 6	73.5 ± 1	45.9 ± 5	67.3 ± 6	59.8 ± 5	88.2 ± 1	88.8 ± 1
	Airline	AIM	73.0 ± 4	85.2 ± 3	74.3 ± 4	81.7 ± 0	92.7 ± 0	82.1 ± 0	65.9 ± 0	70.2 ± 0
		MST	<u>69.2</u> ± 8	<u>80.4</u> ± 6	<u>74.0</u> ± 5	<u>75.2</u> ± 0	<u>86.3</u> ± 0	<u>76.4</u> ± 0	66.1 ± 0	70.3 ± 0
		DP-GAN	38.3 ± 25	65.1 ± 14	60.3 ± 6	42.1 ± 24	62.7 ± 12	59.2 ± 6	43.4 ± 12	46.0 ± 12
		DP-CTGAN	64.1 ± 10	74.0 ± 8	65.9 ± 6	70.1 ± 5	79.4 ± 7	70.1 ± 5	78.5 ± 2	79.0 ± 1
		DP-VAE	1.0 ± 3	54.4 ± 11	56.7 ± 1	52.0 ± 14	61.4 ± 14	58.0 ± 8	41.5 ± 0	42.2 ± 0
		DP-Standard	55.8 ± 3	62.1 ± 7	60.0 ± 6	65.1 ± 6	68.4 ± 10	64.9 ± 8	89.6 ± 4	91.1 ± 3
		DP-2Stage-U	64.5 ± 10	73.8 ± 9	70.0 ± 7	72.5 ± 6	81.8 ± 8	74.1 ± 6	<u>90.1</u> ± 0	<u>91.3</u> ± 1
		DP-2Stage-O	53.5 ± 16	61.8 ± 16	59.2 ± 13	56.9 ± 20	63.3 ± 22	60.8 ± 18	92.0 ± 1	93.0 ± 1

Table 4: **DP Benchmarks.** For each dataset, the best value is highlighted in **bold** while second best is underlined. Results are averaged across five model runs with varying random seeds, with four synthetic datasets generated per run. LR refers to the Logistic Regression model, and XGB represents the XGBoost model.