



MOTIVATION

- Privacy concerns are a limitation when sharing sensitive data, e.g., medical and census
- Differential privacy (DP)⁽¹⁾ safeguards against privacy attacks. An algorithm \mathcal{A} is (ϵ, δ) -DP if for any set S and neighboring dataset D, D' :

$$\mathbb{P}(\mathcal{A}(D) \in S) \leq \exp(\epsilon)\mathbb{P}(\mathcal{A}(D') \in S) + \delta \quad (1)$$

- This paper develops a novel methodology for high-quality DP data sanitation of large-scale tabular datasets

SOTA: MARGINAL-BASED APPROACHES

Algorithm: SOTA marginal-based DP data synthesis

Require Dataset \mathcal{D} , privacy parameters ϵ and δ

- select** set \mathcal{S} of subsets of $\{1, \dots, d\}$
- privatize** (discretized) marginals $\nu_S[\mathcal{D}]$: obtain (ϵ, δ) -DP copies $\hat{\nu}_S$ using e.g., the Gaussian mechanism
- generate** data from privatized marginals $\hat{\nu}_S$

return the DP dataset \mathcal{D}_{DP}

Step 3: graphical models (PGM)⁽²⁾ are the backbone of SOTA methods. Finds prob. dist. \hat{p} by approximately minimizing

$$\min_{\hat{p}} \sum_{S \in \mathcal{S}, x \in \mathcal{X}_S} (\hat{p}_S(\{x\}) - \hat{\nu}_S(\{x\}))^2 \quad \text{then } \mathcal{D}_{\text{DP}} \sim (\hat{p})^n$$

- ✓ robust and sample efficient, suitable for small ϵ and sample sizes n
- ✗ run-time increases exponentially in dimension d when selecting “too many” marginals!
- ✗ squared loss does not capture the “geometry” of the data, e.g., ordering
- ✗ limited abilities to incorporate additional domain-specific constraints

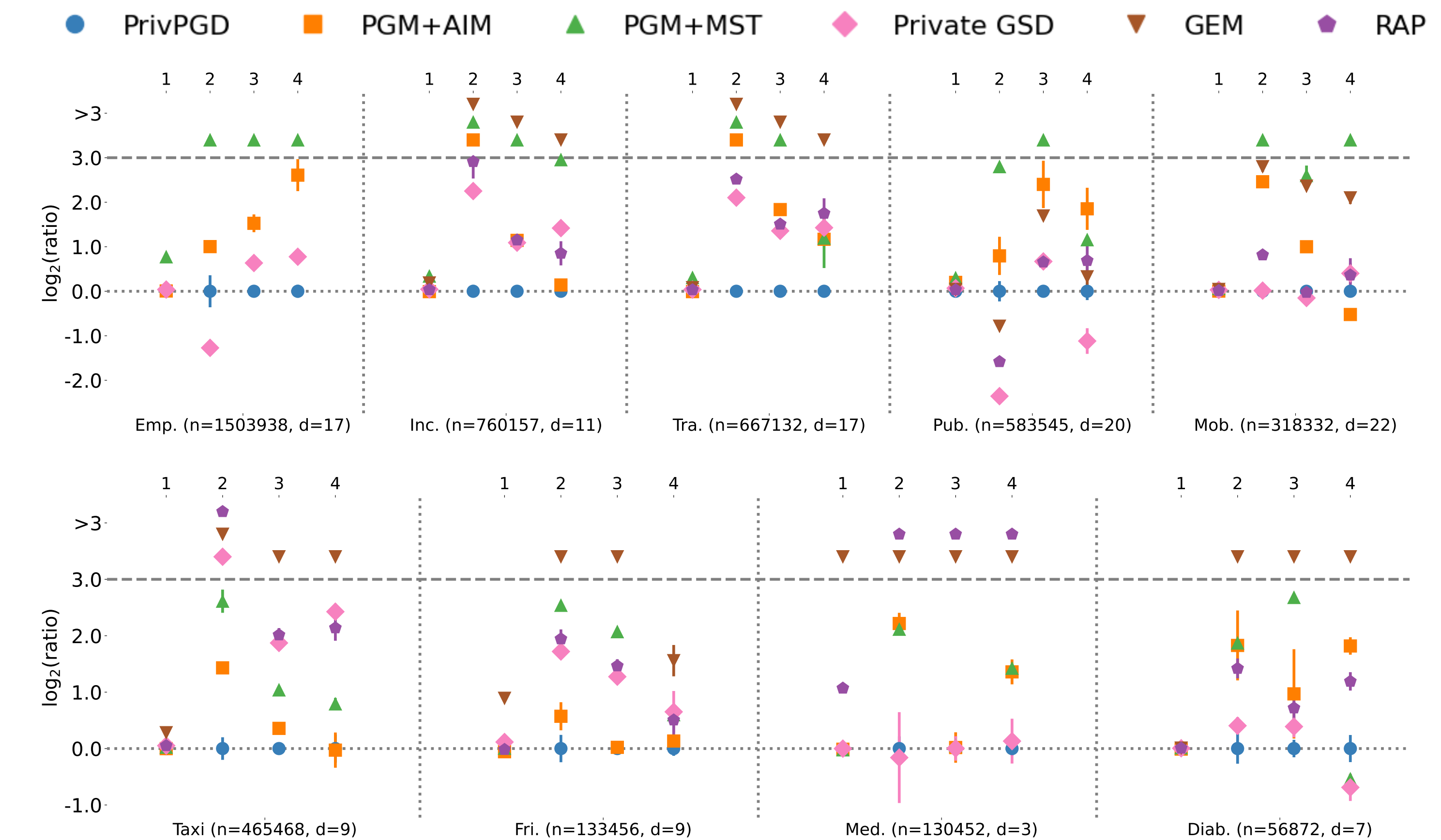
PRIVPGD OUTPERFORMS BASELINES IN A WIDE RANGE OF SETTINGS

Benchmark against SOTA methods on 9 real-world datasets ($\epsilon = 2.5$ and $\delta = 1e-5$)

- PGM+AIM/MST (marginal-based), Private GSD, RAP (query-based), and GEM (generator-based)

Diverse set of **metrics**:

- Downstream classification error
- Frobenius norm of differences of covariance matrix of data embedded in hypercube
- Error rate on 3-sparse counting queries
- Error rate on 3-sparse linear thresholding queries



PRIVATE PARTICLE GRADIENT DESCENT

Algorithm: Private Particle Gradient Descent (PrivPGD)

Require: DP marginals $\{\hat{\nu}_S\}_{S \in \mathcal{S}}$, additional (differentiable) loss $\hat{\mathcal{R}}$ capturing domain-specific constraints

- project:** construct probability measures $\hat{\mu}_S$ from noisy marginals $\hat{\nu}_S$
- optimize:** run gradient descent for particles $Z^{(0)} \in \Omega^m$ on squared sliced Wasserstein distance SW_2^2 :

$$\sum_{S \in \mathcal{S}_{\text{batch}}} \text{SW}_2^2(\mu_S[Z], \hat{\mu}_S) + \lambda \hat{\mathcal{R}}(Z) \quad (2)$$

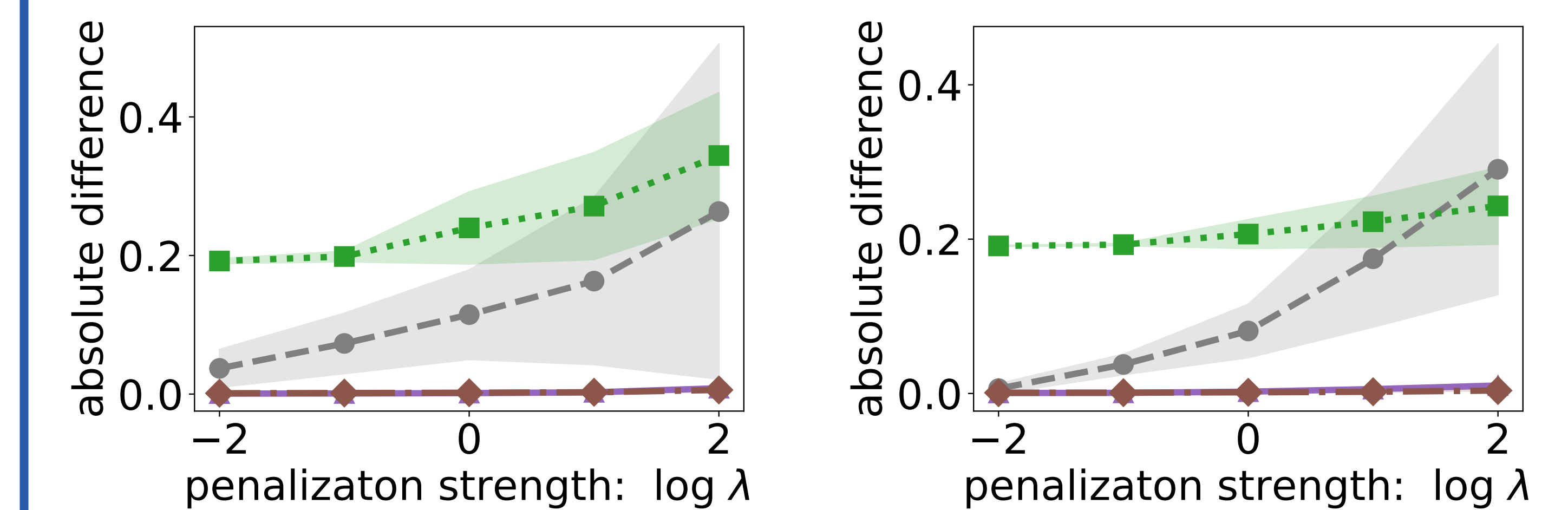
construct DP dataset \mathcal{D}_{DP} from $Z^{(T)}$

- Only linear run-time complexity in # marginals and # particles (up to \log) \Rightarrow allows to synthesize large scale datasets
- Efficient implementation of SW_2^2 using few random proj. \Rightarrow generate datasets of 15+ dim within minutes
- Captures the geometry of the original space and respects orderings of the data

INCORPORATING DOMAIN-SPEC. CONSTRAINTS

- Gradient descent-based generation offers great flexibility
- Example: achieve poor probing accuracy in spec. direction

—●— domain-specific query —▲— counting queries
 - - - test error —◆— thresholding queries



REFERENCES

- Dwork, C. (2006, July). Differential privacy. International colloquium on automata, languages, and programming. Berlin, Heidelberg: Springer Berlin Heidelberg.
- McKenna, R., Sheldon, D., Miklau, G. (2019, May). Graphical-model based estimation and inference for differential privacy. ICML.
- Liu, T., Tang, J., Vietri, G., Wu, S. (2023, July). Generating private synthetic data with genetic algorithms. ICML.
- Liu, T., Vietri, G., Wu, S. Z. (2021). Iterative methods for private synthetic data: Unifying framework and new methods. NeurIPS.