

IS API ACCESS TO LLMs USEFUL FOR GENERATING PRIVATE SYNTHETIC TABULAR DATA?

Anonymous authors

Paper under double-blind review

ABSTRACT

Differentially private (DP) synthetic data is a versatile tool for enabling the analysis of private data. Recent advancements in large language models (LLMs) have inspired a number of algorithm techniques for improving DP synthetic data generation. One family of approaches uses DP finetuning on the foundation model weights; however, the model weights for state-of-the-art models may not be public. In this work we propose two DP synthetic tabular data algorithms that only require API access to the foundation model. We adapt the Private Evolution algorithm (Lin et al., 2023; Xie et al., 2024)—which was designed for image and text data—to the tabular data domain. In our extension of Private Evolution, we define a query workload-based distance measure, which may be of independent interest. We propose a family of algorithms that use one-shot API access to LLMs, rather than adaptive queries to the LLM. Our findings reveal that API-access to powerful LLMs does not always improve the quality of DP synthetic data compared to established baselines that operate without such access. We provide insights into the underlying reasons and propose improvements to LLMs that could make them more effective for this application.

1 INTRODUCTION

Synthetic data has long been a “holy grail” for performing computations on sensitive data, with the allure of protecting privacy while supporting typical data queries and regular data workflows out-of-the-box. Unfortunately, without a rigorous treatment of privacy, the synthetic dataset may inadvertently reveal information about the sensitive data from which it is derived.

Differential privacy (DP) (Dwork et al., 2006; 2016) has emerged as the gold standard for quantifying privacy leakage by algorithms that process sensitive data records from users. At a high level, a (randomized) algorithm satisfies differential privacy if the algorithm’s output distribution is not affected very much by a single person’s data, regardless of what the other data records are. This ensures that the mechanism’s output reveals little about any individual person’s data as a result of their participation in the data analysis, even after arbitrary post-processing of the mechanism output.

Many algorithms have been developed for DP synthetic data, particularly for tabular data (McKenna et al., 2022; Tao et al., 2021; Liu et al., 2021b; Aydore et al., 2021; Liu et al., 2021a; Cai et al., 2021; Zhang et al., 2021). With the advancement of large language models (LLMs), a number of recent works propose improved DP synthetic data algorithms that use LLMs trained on public data¹. Among these are two broad categories of methods: those which privately finetune a foundation model, and those which only use API access to the foundation model. Sablayrolles et al. (2023), Tran & Xiong (2024), and Afonja et al. (2024) use private finetuning on generative language models to generate private tabular synthetic data, and Kurakin et al. (2023) similarly do private LoRA finetuning on an LLM to generate synthetic text data. Similarly, Ghalebikesabi et al. (2023) employ DP finetuning of diffusion models for generating DP synthetic images.

Despite their power, these DP finetuning methods have significant hurdles. First, finetuning algorithms require white-box access to the model, as the weights need to be directly adjusted. This is a

¹The extent to which data used to train LLMs is considered *public* and compatible with privacy goals is hotly contested (Tramèr et al., 2022). We sidestep this question and assume public models are fair to treat as non-private, but we acknowledge it remains an important question.

054 problem because many state-of-the-art models are proprietary, with weights that remain confidential.
055 Only a limited set of researchers are able to even experiment with DP finetuning on such models.
056 Secondly, the resources needed for DP finetuning scales with model dimensionality; time and en-
057 ergy costs quickly become prohibitive. These hurdles motivate alternative ways of using foundation
058 models. In particular, even many proprietary models have a publicly accessible API.

059 A series of works in the synthetic image (Lin et al., 2023) and text (Xie et al., 2024) domains use
060 only API access to foundation models. The algorithm, Private Evolution, combines adaptive queries
061 to the foundation model with a genetic algorithm to privately generate synthetic image and text data.
062 These methods were further extended to the federated setting by Hou et al. (2024). Yet a different
063 approach (Amin et al., 2024) uses private prediction combined with other privacy budget saving
064 tricks on the foundation model to generate DP synthetic text.

065 In light of these recent successes for image and text data, we ask: *Can API access to an LLM improve*
066 *algorithms for generating DP synthetic tabular data?*
067

068 *A priori* it isn't obvious why an LLM would be useful at all for generating synthetic tabular data that
069 it wasn't trained on; however, in initial experiments we found that with descriptive column names,
070 the LLM we used has a reasonable prior over realistic-looking data records. This prior is a powerful
071 source of information we harness in our algorithms.

072 We design and evaluate two types of DP synthetic tabular generation algorithms that leverage LLM
073 API access. In Section 2, we adapt Private Evolution (Lin et al., 2023; Xie et al., 2024) to the tabular
074 domain. A key part of our solution uses a workload-aware family of distance functions, which may
075 be of independent interest, to align the genetic algorithm with the final workload error. In Section 3
076 we introduce a new class of private synthetic data algorithms that use one-shot API access to the
077 foundation model. Unlike prior methods, which require adaptive queries to the foundation model
078 or finetuning the model's weights, our method consumes only one (offline) round of queries to the
079 foundation model. Along the way, we evaluate our two approaches against a number of accuracy
080 baselines to determine whether they advance the state-of-the-art for DP synthetic tabular data.

081 We evaluated our algorithms with Gemini 1.0 Pro (Gemini Team Google, 2023), which allowed us
082 to constrain the outputs to structured tabular records. In our evaluations, we find that the proposed
083 methods fail to consistently beat our baselines. Despite this, we think our attempts are instructive to
084 the research community and could inform the development of state-of-the-art methods, especially
085 as foundation models improve. In light of our findings, we share our key take-aways:

086 **The role of data domain.** The state-of-the-art for DP synthetic data generation is highly domain
087 specific. In particular, DP tabular synthetic data has been very well-studied compared to image and
088 text, so the state of the art for tabular data is much harder to improve on. Additionally, prior work on
089 Private Evolution relies on public image and text embeddings to measure the fidelity of the synthetic
090 data, but similar embeddings do not exist for tabular data. Our workload-aware distance function in
091 Section 2 is one substitute, but surely other solutions exist as well.
092

093 **The importance of appropriate baselines.** In the tabular data domain, there is no single algorithm
094 that dominates on all datasets, query workloads, and privacy budgets. Any new algorithm in this area
095 requires extensive comparison to the handful of algorithms that dominate the state-of-the-art, as well
096 as naive baselines. In Section 3, we show that combining Gemini-generated data with JAM (Fuentes
097 et al., 2024) outperforms all other methods; however, in testing other baselines we find that this holds
098 regardless of the public data we give JAM. Without this naive baseline, we would have reached a
099 false conclusion that the Gemini-generated data was improving the state-of-the-art.
100

101 2 ADAPTING PRIVATE EVOLUTION TO TABULAR DATA

102

103 We adapt Private Evolution (PE) (Lin et al., 2023; Xie et al., 2024) to the tabular data domain. Private
104 Evolution works in rounds, by maintaining a set of *candidates* S_t generated by the foundation model
105 and using a distance function together with a differentially private histogram to have each private
106 record individually vote for candidates. The best performing synthetic candidates become part of an
107 *elite set* for that round S_t^e ; at the end of each round, the foundation model is prompted to generate
more examples similar to the elite set, which then become the new candidates S_{t+1} .

The first set of candidates are populated by a `Random_API`, which prompts the model to generate some prespecified number of initial candidates adhering to the column names and datatypes of the private dataset. Each subsequent set of candidates are generated via the `Variation_API` which takes the current elite set of candidates and prompts the model to generate some number of additional candidates that are similar.

Algorithm 1 Private Evolution (Lin et al., 2023; Xie et al., 2024)

Input: Private samples S_{priv} , Number of iterations T , Number of generated samples N_{synth} , Distance function $\text{dist}_\varepsilon(\cdot, \cdot)$, Noise multiplier σ
Output: Synthetic data S_{synth}

- 1: $S_1 \leftarrow \text{Random_API}(2 \cdot N_{\text{synth}})$
- 2: **for** $t = 1$ to T **do**
- 3: $H = []$ ▷ Initialize histogram over S_t
- 4: **for** $x_{\text{priv}} \in S_{\text{priv}}$ **do**
- 5: $i = \arg \min_{j \in [n]} \text{dist}_\varepsilon(x_{\text{priv}}, S_t)$ ▷ Compute closest synthetic candidate
- 6: $H[i] = H[i] + 1$
- 7: $H \leftarrow H + \mathcal{N}(0, \sigma I_{2 \cdot N_{\text{synth}}})$ ▷ Add noise to ensure DP
- 8: $H \leftarrow \max(0, H)$ ▷ Post-process element-wise
- 9: $\mathcal{P}_t \leftarrow H / \text{sum}(H)$ ▷ Compute empirical distribution on S_t
- 10: $S'_t \leftarrow \text{draw } N_{\text{synth}} \text{ samples with replacement from } \mathcal{P}_t$
- 11: $S_{t+1} \leftarrow \text{Variation_API}(S'_t)$
- 12: **return** S_T

2.1 WORKLOAD-AWARE DISTANCE FUNCTION

Prior methods that applied Private Evolution to image (Lin et al., 2023) and text (Xie et al., 2024) data used public text and image embeddings, respectively, to measure the distance between candidate synthetic examples and the private examples. Choosing a sensible distance function for tabular records is less straightforward: public tabular embeddings (if they exist) likely wouldn’t capture the features of unseen data, simple ℓ_p distance fails to account for differences in scale among columns.

Instead, we derive a workload-aware distance function. A private synthetic dataset is typically optimized for and evaluated on a particular *workload* of (linear) queries $W = \{q_1, \dots, q_k\}$. The workload error is typically some ℓ_p variation on:

$$\text{WError}(S_{\text{priv}}, S_{\text{synth}}) = \sum_{i \in [k]} |q_i(S_{\text{priv}}) - q_i(S_{\text{synth}})|.$$

Note that workload error is a function of pairs of datasets; however, the distance function we require is a function of pairs of individual records. We unpack the workload error further: assuming the queries are *linear*, then they correspond to a sum over a predicate on data records $q_i(\mathbf{x}) = \sum_{j \in [n]} \psi_i(x_j)$. Thus, for the given predicates $\psi = (\psi_1, \dots, \psi_k)$ corresponding to the queries in W , we will define the workload-aware distance function between a private record and synthetic candidate:

$$\text{Wdist}_\psi(x, c) = \sum_{i \in k} |\psi_i(x) - \psi_i(c)|.$$

A dataset of synthetic candidates with low workload-aware distance will have low workload error compared to the private data.

2.2 EXPERIMENTAL RESULTS

We evaluate our adapted private evolution algorithm on a modified version of NYC Taxi and Limousine Commission data using Gemini 1.0 Pro. We use data from January 2024, to avoid data contamination between the public and private evaluation data (Google Cloud, 2023).

For our workload, we use a scaled ℓ_1 -distance on numerical variables for each combination of categorical variables. We rescale the numerical variables to account for different value ranges—for

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

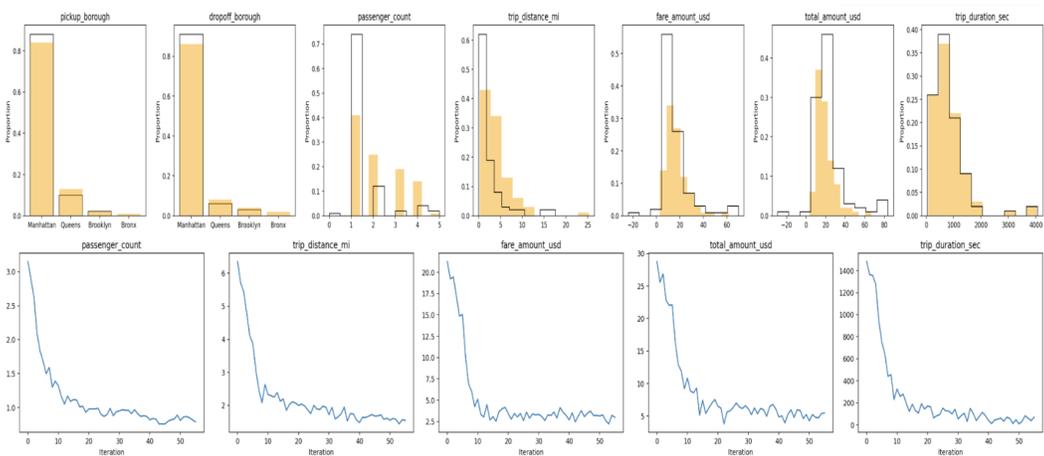


Figure 1: **Top** 1-way marginals on private (outlined) and synthetic (yellow) data. **Bottom** workload error of synthetic data over time for Private Evolution with $\epsilon = \infty$.

example, trip distance (in miles) versus trip duration (in seconds). As an initial experiment, we ran the algorithm without any privacy constraints, and we found that the workload error converges and the 1-way marginals of the synthetic data converges to the 1-way marginals on the private data as well (see Figure 1). It’s worth noting that, depending on the expressivity of the foundation model, it’s not a given that the PE algorithm would converge even without privacy.

We then ran the same experiments but with a DP histogram instead of a nonprivate one. We experimented with various hyperparameters: how many initial random examples to use, how many iterations to use, how to split the budget across iterations, etc. We found that using increasing budget across runs worked better than an even or decreasing budget; additionally, having more candidates relative to N_{synth} worked best, and finally, using fewer iterations worked best—in fact, using *only a single iteration* was the optimal setting we could find.

With differential privacy, the private evolution algorithm failed to beat two simple baselines: *independent* which privately computes all 1-way marginals and samples data from the product over the private marginals, and *DP workload* which directly computes the workload queries with DP, without generating any synthetic data. These two baselines are not the only two which we’d hope to beat, rather, they’re the bare minimum. Moreover, querying any large foundation model hundreds of times is relatively slow.

What we learned The observation that the workload-aware private evolution algorithm performs best with *one shot* data generation implies that: *whatever marginal gains we get from iterating multiple times, they are outweighed by the privacy cost of composing over iterations*. Additionally, while PE was developed for image and text domains where finetuning a foundation model is the main alternative for DP synthetic data, there is a vast literature on algorithms for private synthetic *tabular* data that do not require access to generative models. These lessons paved the way for our second attempt, which proved to be more successful, though still did not beat the current state-of-the-art.

3 USING GEMINI-GENERATED RECORDS AS PUBLIC DATA

There is a substantial body of work on DP synthetic tabular data. Some state-of-the-art algorithms within this space make use of *public data* to improve the accuracy or efficiency of the algorithm on private data; however, for many applications such public data may not be available in the format required², as is discussed in-depth in Liu et al. (2021b)[Section 6.1]. Our second approach uses Gemini generated data in lieu of this public data.

²This is especially true for algorithms that assume the public dataset has the same (or substantially overlapping) columns, or even is distributed similarly to the private data.

3.1 APPROACH OVERVIEW

Using Gemini’s structured output functionality, we prompt Gemini to generate data records with a response schema matching the column names and datatypes of our private dataset. Importantly, none of the private records influence the prompts to Gemini—only the column names and datatypes do. This data generation occurs “offline” and in one shot with no loss of privacy budget. We call this dataset Gem_Synth. Later, we plug this synthetic public dataset into various DP synthetic tabular algorithms that use public data as well as the private data.

One major benefit of the one-shot nature of this method (rather than querying Gemini interactively in a loop) is we can generate *many* synthetic public data records and reuse the same generated records when trying different approaches. This is not possible when the records are generated adaptively as in Private Evolution. Thus, this method takes advantage of our observations about PE. We begin with a high-level overview of how public data is incorporated into two DP synthetic data algorithms. For both, we consider what happens when we use Gemini-generated tabular data as the public data source for these algorithms.

PMW^{pub} Liu et al. (2021a) The PMW^{pub} algorithm is an improvement of MWEM (Hardt et al., 2012), which we will not discuss in detail. The basic idea is to use public data to initialize the generating distribution over synthetic records and iteratively refine this distribution to reduce the workload error. The public records reduce the number of iterations required, by providing a “warm start” for the synthetic data distribution, along with reducing the data domain over which the distribution is estimated.

A key sub-routine of PMW^{pub} is to estimate a distribution that approximately matches some noisy statistics. Specifically, let Q denote a collection of linear queries and let $\tilde{y} = Q(D) + \xi = \sum_{x \in D} Q(x) + \xi$ be the noisy answers to those queries on the sensitive data. PMW^{pub} finds a distribution supported on the “public” data Gem_Synth, and finds the weights to assign to each public record to minimize the ℓ_2 squared error to the noisy observations.

$$w^* = \arg \min_{w \in \mathbb{R}_+} \left\| \sum_{x \in \text{Gem_Synth}} w_x Q(x) - \tilde{y} \right\|_2^2$$

When the public records are sufficiently representative, this method can work quite well. However, with small or unrepresentative public datasets, this method may not find a good distribution even in the absence of noise.

Gemini inference We use the Gemini-generated records as the public records for the subroutine of PMW^{pub}, calling this “Gemini inference”, setting Q to be the query workload.

MST modified to take public data The standard MST algorithm (McKenna et al., 2021) has three phases: selecting marginal queries, measuring the marginals with DP, and lastly using Private-PGM to post-process the noisy marginals and generate a synthetic dataset. We modify the final step (generation), replacing Private-PGM with the subroutine from PMW^{pub} that utilizes Gem_Synth. This method differs from the Gemini inference approach primarily in how the queries Q are selected.

JAM The JAM-PGM mechanism (Fuentes et al., 2024) was developed for marginal queries, and utilizes public data in a different manner. It privately decides whether to measure each marginal query on the public data or the private data in order to minimize the overall workload error. This mechanism has the benefit that it can utilize public data that is accurate on some, but not necessarily all, marginals. We run this mechanism as-is, using Gem_Synth as the public data.

3.2 BASELINES FOR COMPARISON

Because there are a wealth of methods for generating private synthetic data with and without public data, we have a fair number of baselines that we need to compare any new methods to. A number of works that privately finetune foundation models for tabular data omit comparisons to state-of-the-art methods for generating DP tabular data, so it is unclear if they outperform existing approaches.

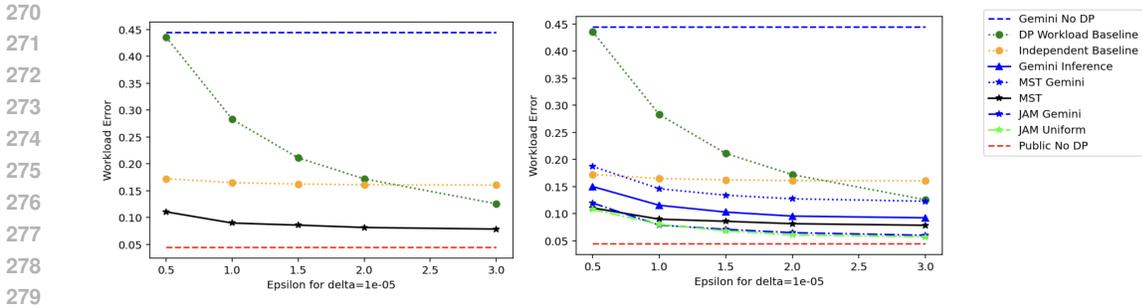


Figure 2: (Left) Workload error for baseline methods for generating tabular synthetic data without use of Gemini. (Right) Workload error for baseline methods and our one-shot methods that use API access to Gemini.

We study two baselines that require no privacy. First, we consider an **in-distribution public dataset**: publically data drawn from the same distribution as the private data. This is essentially the lowest error we could hope for, up to sampling error; however, in-distribution public data is usually not available. Second, we consider using the **Gemini data with no DP** to answer workload queries.

Next, we consider a number of baselines that do not require public data. **DP workload**: we compute the queries directly using DP. **Independent baseline**: privately compute the 1-way marginals and sample records from the corresponding product distribution over marginals. **MST algorithm**: a tabular synthetic data algorithm that does not use public data. **PMW^{pub} with uniform data** and **JAM with uniform data**: using data that is drawn uniformly from the domain as public data for these algorithms.

3.3 EXPERIMENTAL SETUP

Our private dataset is UCI Adult (Becker & Kohavi, 1996). Using this structured output constraint, we sample Gemini with top- $k=1$ and temperature=1 to generate 131,000 records in Gem_Synth. We use 2-way marginals as our query workload to evaluate the fidelity of the DP synthetic data.

3.4 RESULTS

Figure 2 (Left) shows the workload error versus epsilon for the baseline methods discussed. Among these methods, MST achieves the lowest workload error (except the in-distribution public dataset which is our unachievable “best-case baseline”). Figure 2 (Right) shows all of the results for the baseline methods in addition to the methods that use Gem_Synth. Notice that JAM with Gem_Synth performs best overall; however, JAM performs equally well with uniform data. This is because JAM is simply using the private data to compute answers to the queries rather than utilizing the public data. Thus, JAM with Gem_Synth is not better than the state-of-the-art methods on this dataset and query workload. In general, Gem_Synth may capture 1-way marginals on the data better than uniform, however it is generally inaccurate on k -way marginals.

4 CONCLUSION AND FUTURE WORK

We evaluated two methods for incorporating API access to Gemini for generating DP synthetic tabular data. While our methods did not beat state-of-the-art methods, this work motivates a number of future directions. First, as foundation models continue to improve, combining our methods with better models (e.g. models trained on more tabular data) could potentially improve the final accuracy, especially if the models are trained specifically for the tabular setting. Additionally, because Gemini uses word embeddings, perhaps doing some finetuning on publicly available tabular data could improve the quality of the Gemini-generated tabular records fed into our one-shot method. Lastly, perhaps there are ways to achieve better accuracy by combining Private Evolution and our one-shot approach. Using foundation models for DP synthetic data generation is still a very new area of research, with many avenues for improvements and breakthroughs.

REFERENCES

- 324
325
326 Tejumade Afonja, Hui-Po Wang, Raouf Kerkouche, and Mario Fritz. Dp-2stage: Adapting language
327 models as differentially private tabular data generators, 2024. URL <https://arxiv.org/abs/2412.02467>.
328
- 329 Kareem Amin, Alex Bie, Weiwei Kong, Alexey Kurakin, Natalia Ponomareva, Umar Syed, Andreas
330 Terzis, and Sergei Vassilvitskii. Private prediction for large-scale synthetic text generation. *arXiv preprint arXiv:2407.12108*, 2024.
331
332
- 333 Sergul Aydore, William Brown, Michael Kearns, Krishnaram Kenthapadi, Luca Melis, Aaron Roth,
334 and Ankit A Siva. Differentially private query release through adaptive projection. In *International Conference on Machine Learning*, pp. 457–467. PMLR, 2021.
335
- 336 Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI:
337 <https://doi.org/10.24432/C5XW20>.
338
- 339 Kuntai Cai, Xiaoyu Lei, Jianxin Wei, and Xiaokui Xiao. Data synthesis via differentially private
340 markov random fields. *Proceedings of the VLDB Endowment*, 14(11):2190–2202, 2021.
- 341 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity
342 in private data analysis. In *Theory of cryptography conference*, pp. 265–284. Springer, 2006.
343
- 344 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity
345 in private data analysis. *Journal of Privacy and Confidentiality*, 7(3):17–51, 2016.
- 346 Miguel Fuentes, Brett C Mullins, Ryan McKenna, Gerome Miklau, and Daniel Sheldon. Joint
347 selection: Adaptively incorporating public information for private synthetic data. In *International Conference on Artificial Intelligence and Statistics*, pp. 2404–2412. PMLR, 2024.
348
349
- 350 Gemini Team Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
351
- 352 Sahra Ghalebikesabi, Leonard Berrada, Sven Gowal, Ira Ktena, Robert Stanforth, Jamie Hayes,
353 Soham De, Samuel L Smith, Olivia Wiles, and Borja Balle. Differentially private diffusion models
354 generate useful synthetic images. *arXiv preprint arXiv:2302.13861*, 2023.
355
- 356 Google Cloud. Generative AI on Vertex AI Documentation: Google Models, 2023.
357 URL <https://cloud.google.com/vertex-ai/generative-ai/docs/learn/models#gemini-1.0-pro>.
358
- 359 Moritz Hardt, Katrina Ligett, and Frank McSherry. A simple and practical algorithm for differentially
360 private data release. *Advances in neural information processing systems*, 25, 2012.
361
- 362 Charlie Hou, Akshat Shrivastava, Hongyuan Zhan, Rylan Conway, Trang Le, Adithya Sagar, Giulia
363 Fanti, and Daniel Lazar. Pre-text: Training language models on private federated data in the age
364 of llms. *arXiv preprint arXiv:2406.02958*, 2024.
- 365 Alexey Kurakin, Natalia Ponomareva, Umar Syed, Liam MacDermed, and Andreas Terzis. Harnessing
366 large-language models to generate private synthetic text. *arXiv preprint arXiv:2306.01684*,
367 2023.
368
- 369 Zinan Lin, Sivakanth Gopi, Janardhan Kulkarni, Harsha Nori, and Sergey Yekhanin. Differentially
370 private synthetic data via foundation model apis 1: Images. *arXiv preprint arXiv:2305.15560*,
371 2023.
- 372 Terrance Liu, Giuseppe Vietri, Thomas Steinke, Jonathan Ullman, and Steven Wu. Leveraging
373 public data for practical private query release. In *International Conference on Machine Learning*,
374 pp. 6968–6977. PMLR, 2021a.
375
- 376 Terrance Liu, Giuseppe Vietri, and Steven Z Wu. Iterative methods for private synthetic data: Uni-
377 fying framework and new methods. *Advances in Neural Information Processing Systems*, 34:
690–702, 2021b.

378 Ryan McKenna, Gerome Miklau, and Daniel Sheldon. Winning the nist contest: A scalable and
379 general approach to differentially private synthetic data. *arXiv preprint arXiv:2108.04978*, 2021.
380

381 Ryan McKenna, Brett Mullins, Daniel Sheldon, and Gerome Miklau. Aim: An adaptive and iterative
382 mechanism for differentially private synthetic data. *arXiv preprint arXiv:2201.12677*, 2022.
383

384 Alexandre Sablayrolles, Yue Wang, and Brian Karrer. Privately generating tabular data using lan-
385 guage models. *arXiv preprint arXiv:2306.04803*, 2023.

386 Yuchao Tao, Ryan McKenna, Michael Hay, Ashwin Machanavajjhala, and Gerome Miklau.
387 Benchmarking differentially private synthetic data generation algorithms. *arXiv preprint*
388 *arXiv:2112.09238*, 2021.

389 Florian Tramèr, Gautam Kamath, and Nicholas Carlini. Position: Considerations for differentially
390 private learning with large-scale public pretraining. In *Forty-first International Conference on*
391 *Machine Learning*, 2022.
392

393 Toan V Tran and Li Xiong. Differentially private tabular data synthesis using large language models.
394 *arXiv preprint arXiv:2406.01457*, 2024.

395 Giuseppe Vietri, Cedric Archambeau, Sergul Aydore, William Brown, Michael Kearns, Aaron Roth,
396 Ankit Siva, Shuai Tang, and Steven Z Wu. Private synthetic data for multitask learning and
397 marginal queries. *Advances in Neural Information Processing Systems*, 35:18282–18295, 2022.
398

399 Chulin Xie, Zinan Lin, Arturs Backurs, Sivakanth Gopi, Da Yu, Huseyin A Inan, Harsha Nori, Hao-
400 tian Jiang, Huishuai Zhang, Yin Tat Lee, et al. Differentially private synthetic data via foundation
401 model apis 2: Text. *arXiv preprint arXiv:2403.01749*, 2024.

402 Mengmeng Yang, Chi-Hung Chi, Kwok-Yan Lam, Jie Feng, Taolin Guo, and Wei Ni. Tabular data
403 synthesis with differential privacy: A survey. *arXiv preprint arXiv:2411.03351*, 2024.
404

405 Zhikun Zhang, Tianhao Wang, Ninghui Li, Jean Honorio, Michael Backes, Shibo He, Jiming Chen,
406 and Yang Zhang. {PrivSyn}: Differentially private data synthesis. In *30th USENIX Security*
407 *Symposium (USENIX Security 21)*, pp. 929–946, 2021.
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

432 A PRELIMINARIES

433
434 We begin by presenting the definition of differential privacy, which is a constraint on an algorithm
435 \mathcal{A} that processes a dataset $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ of user records, one per user. Two datasets are called
436 *neighbors* if they differ on one person’s record. At a high level, differential privacy requires that for
437 any pair of neighboring datasets, the algorithm’s output distributions are similar when run on each
438 dataset.

439 **Definition 1 (Differential Privacy (Dwork et al., 2006; 2016))** *A randomized algorithm $\mathcal{A} : \mathcal{U}^n \rightarrow \mathcal{Y}$ is ε -differentially private if for every pair of neighboring datasets $\mathbf{x}, \mathbf{x}' \in \mathcal{U}^n$, and for all outputs $y \in \mathcal{Y}$,*

$$441 \Pr[\mathcal{A}(\mathbf{x}) = y] \leq e^\varepsilon \cdot \Pr[\mathcal{A}(\mathbf{x}') = y] + \delta,$$

442
443 where the probability is taken over the internal coins of \mathcal{A} .

444
445 The differential privacy guarantee is parameterized by $\varepsilon > 0$, where algorithms with lower values
446 have less privacy leakage and higher values of epsilon denote more privacy leakage from the algo-
447 rithm’s output. DP gives a worst-case guarantee (over the algorithm’s inputs and outputs) on how
448 much information an algorithm leaks about its input.

449 A.1 PRIOR WORK

450
451 **GAN-based methods for DP synthetic data** Many prior works have proposed synthetic data
452 mechanisms based on generative adversarial networks. See Yang et al. (2024) for a nice survey of
453 these and other approaches. These mechanisms generally work by fitting the parameters of the model
454 via DP-SGD, and then using the model to generate synthetic data after training. These techniques
455 are typically best suited for unstructured data like images or text.

456
457 **Marginal-based methods for DP synthetic data** Many mechanisms for DP synthetic data gen-
458 eration work by adding noise to low-dimensional marginals of the data distribution McKenna et al.
459 (2021; 2022); Cai et al. (2021); Aydore et al. (2021); Fuentes et al. (2024); Vietri et al. (2022); Liu
460 et al. (2021b;a); Zhang et al. (2021). Some mechanisms in this space are also designed to lever-
461 age public data when it’s available Fuentes et al. (2024); Liu et al. (2021b;a). Benchmarks have
462 confirmed these approaches work very well in tabular data settings Tao et al. (2021).