

$pMSE$ Mechanism: Differentially Private Synthetic Data with Maximal Distributional Similarity

Joshua Snoke and Aleksandra Slavković

Department of Statistics
 Pennsylvania State University, University Park, PA 16802, USA
 {snoke,sesa}@psu.edu

Abstract. We propose a method for the release of differentially private synthetic datasets. In many contexts, data contain sensitive values which cannot be released in their original form in order to protect individuals' privacy. Synthetic data is a protection method that releases alternative values in place of the original ones, and differential privacy (DP) is a formal guarantee for quantifying the privacy loss. We propose a method that maximizes the distributional similarity of the synthetic data relative to the original data using a measure known as the $pMSE$, while guaranteeing ϵ -differential privacy. Additionally, we relax common DP assumptions concerning the distribution and boundedness of the original data. We prove theoretical results for the privacy guarantee and provide simulations for the empirical failure rate of the theoretical results under typical computational limitations. We also give simulations for the accuracy of linear regression coefficients generated from the synthetic data compared with the accuracy of non-differentially private synthetic data and other differentially private methods. Additionally, our theoretical results extend a prior result for the sensitivity of the Gini Index to include continuous predictors.

Keywords: differential privacy, synthetic data, classification trees

1 Introduction

In many contexts, researchers wish to gain access to data which are restricted due to privacy concerns. While there are many proposed methods for allowing researchers to fit models or receive query responses from the data, there are other cases where either due to methodological familiarity or modeling flexibility, researchers desire to have an entire dataset rather than a set of specific queries. This paper proposes a method for releasing synthetic datasets under the framework of ϵ -differential privacy, which formally quantifies and guarantees the privacy loss from these releases.

Differential privacy (DP), originally proposed by [Dwork et al. \(2006\)](#), is a formal method of quantifying the privacy loss related to any release of information based on private data; for a more in-depth review see [Dwork et al. \(2014\)](#), and for a non-technical primer see [Nissim et al. \(2017\)](#). Since its inception, DP has spawned a large literature in computer science and some in statistics. It has been explored in numerous contexts such as machine learning algorithms (e.g., [Blum et al. \(2005\)](#); [Kasiviswanathan et al. \(2011\)](#); [Kifer et al. \(2012\)](#)), categorical data (e.g., [Barak et al. \(2007\)](#); [Li et al. \(2010\)](#)), dimension reduction (e.g., [Chaudhuri et al. \(2012\)](#); [Kapralov and Talwar \(2013\)](#)), performing statistical analysis (e.g., [Wasserman and Zhou \(2010\)](#); [Karwa et al. \(2016\)](#)), and streaming data (e.g., [Dwork et al. \(2010\)](#)), to name only a few applications.

While DP is a rigorous risk measure, it has lacked flexible methods for modeling and generating synthetic data. Non-differentially private synthetic data methods (e.g., see [Raghuathan et al. \(2003\)](#); [Reiter \(2005, 2002\)](#); [Drechsler \(2011\)](#); [Raab et al. \(2017\)](#)) while not offering provable privacy, provide good tools for approximating accurate generative models reflecting the original data. Our proposed method maintains the flexible modeling approach of synthetic data methodology, and in addition maximizes a metric of distributional similarity, the $pMSE$, between the released synthetic data and the original data, subject to satisfying ϵ -DP. We also do *not* require one of the most common DP assumptions concerning the input data, namely that it is bounded, and we do not limit ourselves to only categorical or discrete data. We find that our method produces good results in simulations, and it provides a new avenue for releasing DP datasets for potentially a wide-range of applications.

Our specific contributions are: (1) the combination of the flexible synthetic data modeling framework with the guarantee of ϵ -DP, (2) the relaxation of DP assumptions concerning boundedness or discreteness of the input data, (3) the embedding of a metric within our mechanism guaranteeing maximal distributional similarity between the synthetic and original data, and (4) a proof for the sensitivity bound of the Gini Index for CART models in the presence of continuous predictors.

The rest of the paper is organized as follows. Section 2 gives a review of important results from differential privacy that we rely on and provides a review of related methods to ours which we use for comparison in our simulation study. Section 3 details our proposed methodology for sampling differentially private data with maximal distributional similarity. Section 4 provides theoretical results for the privacy guarantees of our algorithm. Section 5 shows simulations that support our theoretical findings and provide an empirical estimate of the privacy loss under standard computational practices. Section 6 provides simulation results for the comparison of the accuracy of linear regression coefficients calculated from our method and other DP methods. Section 7 details conclusions and future considerations.

2 Differential Privacy Preliminaries

Differentially Privacy is a formal framework for quantifying the disclosure risk associated with the release of statistics or raw data derived from private input data. The general concept relies on defining a randomized algorithm which has similar output probabilities regardless of the presence or absence of any given record, as formalized in Definition 1. We replicate the definitions and theorems here using notation assuming $X \in \mathbb{R}^{n \times q}$ and $\theta \in \mathbb{R}^k$. X is an original data matrix, and we wish to release instead a private version, X^s , with the same dimension. θ is a vector of parameters corresponding to a chosen parametric model, which can be used to generate synthetic data that reflects the generative distribution of X . Further restrictions may be placed on θ depending on the parametric model.

Definition 1 (Dwork et al. (2006)). *A randomized algorithm, \mathcal{M} , is ϵ -Differentially Private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all X, X' such that $\delta(X, X') = 1$:*

$$\frac{P(\mathcal{M}(X') \in \mathcal{S})}{P(\mathcal{M}(X) \in \mathcal{S})} \leq \exp(\epsilon).$$

The privacy is controlled by the ϵ parameter, with values closer to zero offering stronger privacy. Relaxations of ϵ -DP have been proposed to reframe the privacy definition or to improve the statistical utility. A few examples are approximate differential privacy (also known as (ϵ, δ) -DP, see Dwork et al. (2006)), probabilistic differential privacy (Machanavajjhala et al. (2008)), on-average K-L privacy (Wang et al. (2016)), or concentrated privacy (Dwork and Rothblum (2016); Bun and Steinke (2016)). We do not cover these relaxations further, since our work relies on the stronger ϵ -DP.

A general example of an ϵ -DP mechanism, which produces private outputs, is the Exponential Mechanism defined by McSherry and Talwar (2007); see Definition 2. For a given θ that we wish to release, this mechanism defines a distribution from which private samples, $\tilde{\theta}_i$, can be made and released in place of the original vector.

Definition 2 (McSherry and Talwar (2007)). *The mechanism that releases values with probability proportional to*

$$\exp\left(\frac{-\epsilon u(X, \theta)}{2 \Delta u}\right),$$

where $u(X, \theta)$ is a quality function that assigns values to each possible output, θ , satisfies ϵ -DP.

Δu is the global sensitivity, and it is defined as the greatest possible change in the u function for any two inputs, differing in one row. Note that some definitions of DP use the addition or deletion of a row, but here we assume X and X' have the same dimension. More formally:

Definition 3 (Dwork et al. (2006)). *For all X, X' such that $\delta(X, X') = 1$, the Global Sensitivity (GS) of a function $u : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}_{\geq 0}$ is defined as:*

$$\Delta u = \sup_{\theta} \sup_{\delta(X, X')=1} |u(X, \theta) - u(X', \theta)|$$

We also rely on two theorems, known as post-processing and sequential composition. The first, stated in Proposition 1, says that any function applied to the output from a differentially private algorithm is also differentially private. We use this to generate synthetic data based on private parameters, rather than directly generating differentially private data.

Proposition 1 (Dwork et al. (2006); Nissim et al. (2007)). *Let \mathcal{M} be any randomized algorithm, such that $\mathcal{M}(X)$ satisfies ϵ -differential privacy, and let g be any function. $g(\mathcal{M}(X))$ also satisfies ϵ -differential privacy.*

Sequential composition, stated in Theorem 1, says that for multiple outputs from a differentially private algorithm, one must compose the ϵ values for each output to produce the overall privacy loss of the process. We need to compose the privacy if we make multiple draws of private parameters from which we produce multiple private synthetic datasets. We may want to release multiple synthetic datasets for better accuracy in estimates based on the data. Estimates based on multiple datasets are calculated using appropriate combining rules; see Raab et al. (2017) for details.

Theorem 1 (McSherry (2009)). *Suppose a randomized algorithm, \mathcal{M}_j , satisfies ϵ_j -differential privacy for $j = 1, \dots, q$. The sequence $\mathcal{M}_j(X)$ carried out on the same X provides $(\sum_j \epsilon_j)$ -differential privacy.*

These theorems and definitions lay the groundwork for our method. Next, we give a brief overview of some related methods to ours which we use for comparison in our simulations in Section 6.

2.1 Review of Related Methods

A number of different methods have been proposed for releasing differentially private synthetic datasets, although only a few are focused on real-valued, $n \times q$ matrices. Bowen and Liu (2016) proposed drawing data from a noisy Bayesian Posterior Predictive Distribution (BPPD), and Wasserman and Zhou (2010) generate data from smooth histograms. Bowen and Liu (2016) provides a fairly comprehensive list of DP synthetic data methods. Many of these methods are limited to specific data types, such as categorical data (e.g., Abowd and Vilhuber (2008); Charest (2011)), or network data (e.g., Karwa et al. (2016, 2017)), or they are computationally infeasible for data with a reasonable number of dimensions, such as the Empirical CDF (Wasserman and Zhou (2010)), NoiseFirst/StructureFirst (Xu et al. (2013)), or the PrivBayes (Zhang et al. (2017)), though some recent work has proposed ways to reduce the computation time and improve the accuracy (Li et al. (2018)).

The noisy BPPD method from Bowen and Liu (2016) is similar to ours in the sense that focuses on drawing generative model parameters from a noisy distribution and then using these private parameters generates private synthetic data according to post-processing. In this case private parameters are drawn from posterior distribution $f(\theta|s^*)$ where s is the Bayesian sufficient statistic and s^* is the statistic perturbed according to the Laplace mechanism (e.g., see Dwork et al. (2006)). Bowen and Liu (2016) recommend drawing multiple sets of private parameters and producing a synthetic dataset for each one, which requires composing ϵ .

The smooth histogram method from Wasserman and Zhou (2010) works non-parametrically by binning the data, using these bins to estimate a consistent density function, and applying smoothing to the function which guarantees DP before drawing new samples. The DP smooth histogram is defined as:

$$\hat{f}_K^*(\mathbf{x}) = (1 - \lambda)\hat{f}_K(\mathbf{x}) + \lambda\Omega \tag{1}$$

where K is the total number of bins, $\Omega = \left(\prod_{j=1}^p (c_{j1} - c_{j0}) \right)^{-1}$, $\lambda \geq \frac{K}{K+n(e^{\epsilon/n}-1)}$, and $\hat{f}_K(\mathbf{x})$ is a mean-squared consistent density histogram estimator. This method does not need to generate multiple datasets, since it is not redrawing model parameters, and accordingly does not need to split ϵ across multiple synthetic datasets.

However, both of these methods require bounded data. This can be seen explicitly in the smooth histogram formulation where we assume the j_{th} variable has bounds $[c_{j0}, c_{j1}]$. We also need to assume bounds in order to create (and sample from) a finite K bins. The boundedness assumption is less explicit

in the noisy BPPD method, but it comes into play when calculating the sensitivity of the statistics. If the data were unbounded, the sensitivity could be infinite, which would mean we have to sample the noise from a Laplace distribution with infinite variance.

We want to avoid this assumption in our method because assuming bounds is problematic when it comes to approximating the underlying generative distribution. Continuous data may be naturally unbounded, and at best in many real data scenarios we do not know what the bounds should be. If our bounds are too loose, we introduce more noise than necessary through the privacy process, limiting our accuracy. On the other hand, we introduce bias if we set the bounds too low because that truncates the generative distribution below its true range. We further explore the effect of these assumed bounds using a simulation study in Section 6.

Our method avoids the bounding problem by sampling from a distribution that shrinks in probability as we move to the tails. We do not bound the sample space, but we have low probability of sampling values which are far from the truth. This allows us to produce private data that accurately reflects the natural range of the data. We describe this further in Section 3.1.

Furthermore, the smooth histogram method suffers from computational limitations as the number of variables increases, since it divides the data matrix into bins, the number of which grows $O(p^p)$. Our method has computational limitations too, though of a different nature, which we discuss further in Section 3. One nice aspect of the noisy BPPD method is that it is computationally fast.

Finally, our method improves over these methods by incorporating a measure of distributional similarity on the resulting synthetic data. The noisy BPPD and smooth histogram add noise to the generative model for the synthetic data. These mechanisms concern only minimizing the noise added to the parameters, but they guarantee nothing concerning the data generated using these parameters. Our method on the other hand, finds the private parameters which can be used to generate synthetic data that will maximize the distributional similarity with respect to the original data. Regardless of the dataset, our method finds the best private parameters for an assumed synthesizing model. Sections 3 and 4 give a detailed explanation and theoretical results for our method.

3 Sampling Differentially Private Synthetic Data via the $pMSE$ Mechanism

We propose to release a DP synthetic data matrix, X^s , in place of the original data X and with the same dimension, $n \times q$. In practice it is infeasible to sample a matrix of this size using the Exponential Mechanism. More simply, we sample private model parameters and then generate synthetic data based on these noisy parameters. We know from the results on post-processing, see Proposition 1, that data generated based on these DP parameters are also DP. Based on the exponential mechanism we draw DP samples using

$$f(\theta) \propto \exp\left(\frac{-\epsilon u(X, \theta)}{2\Delta u}\right), \quad (2)$$

where ϵ is our privacy parameter, $u(X, \theta)$ is our quality (or utility) function, and Δu is the sensitivity of the quality function. In practice we use a Monte Carlo Markov Chain (MCMC) approach, using the Metropolis algorithm to generate samples from this unnormalized density, since we do not know the value of the u function a priori. Next we define our quality function and derive a bound on its sensitivity.

3.1 Defining the Quality Function using the $pMSE$

We base our quality function on the $pMSE$ statistic developed in Woo et al. (2009) and Snoke et al. (2018):

$$pMSE = \frac{1}{N} \sum_{i=1}^N (\hat{p}_i - 0.5)^2,$$

where \hat{p}_i are predicted probabilities (i.e., propensity scores) from a chosen classification algorithm. Algorithm 1 gives the steps for calculating the $pMSE$ statistic. The $pMSE$ is simply the mean-squared error of

the predicted probabilities from this classification task, and it is a metric to assess how well we are able to discern between datasets based on a classifier. If we are unable to discern, the two datasets have high distributional similarity. A $pMSE = 0$ means every $\hat{p}_i = 0.5$, and it implies the highest utility. There has been much work dedicated to tuning models for out-of-sample prediction, but for our purposes we only use the classifier to get estimates of the in-sample predicted probabilities.

Algorithm 1 General Method for Calculating the $pMSE$

- 1: stack the n rows of original data X and the n rows of masked data X^s to create X^{comb} with $N = 2n$ rows
 - 2: add an indicator variable, I , to X^{comb} s.t. $I = \{1 : x_i^{comb} \in X^s\}$
 - 3: fit a model to predict I using predictors $Z = f(X^{comb})$.
 - 4: find predicted probabilities of class 1, \hat{p}_i , for each row of X^{comb}
 - 5: obtain the $pMSE = \frac{1}{N} \sum_{i=1}^N (\hat{p}_i - 0.5)^2$
-

To make our quality function a function of θ , the vector of parameters we wish to sample, we use the expected value of the $pMSE$ given θ , i.e.,

$$u(X, \theta) = E[pMSE(X, X^s) | X, \theta] \quad (3)$$

where $X^s \sim g(\theta)$. In practice we approximate this by generating m datasets for a given set of parameters and calculate the average $pMSE$ across each data set. This approximation does not affect the privacy guarantee (as shown in the proof for Theorem 2), but for accuracy m should be large enough to give satisfactory results for estimating the expected value of $u(X, \theta)$.

As we mentioned before, this quality function makes no assumptions concerning whether the original data are categorical, discrete, or continuous. Secondly, because the $pMSE$ is a function of the predicted probabilities, \hat{p} , which are bounded $\in [0, 1]$, the $pMSE$ is bounded $\in [0, 0.25]$. This is true regardless of the range of the data, X , so we do not need to assume any kinds of bounds on the data.

We refer to our method as the $pMSE$ Mechanism, since we rely on the $pMSE$ for our quality function in the exponential mechanism. Algorithm 2 outlines the steps of the $pMSE$ mechanism. The main assumption we need is that a reasonable generative model for the data, $g(\theta)$, exists.

Algorithm 2 Sampling DP Synthetic Data via the $pMSE$ Mechanism

Input: Original dataset: X , chosen value: ϵ , synthesis model: $g(\theta)$, quality function: $u(X, \theta)$

- 1: Sample l vectors $\{\tilde{\theta}_1, \dots, \tilde{\theta}_l\}$ from a density proportional to equation 2
 - 2: For each $\tilde{\theta}_i$ generate synthetic data set $X_i^s \sim g(\tilde{\theta}_i)$, giving m total synthetic datasets $\{X_1^s, \dots, X_l^s\}$ each with the same dimension as X
 - 3: Releasing $\{X_1^s, \dots, X_l^s\}$ satisfies $(l\epsilon)$ -DP
-

3.2 Estimating the $pMSE$ using Classification and Regression Tree (CART) Models

A key component to defining our quality function is the classification model used to estimate the predicted probabilities, \hat{p} , used in computing the $pMSE$. We choose the classification trees (Breiman et al. (1984)) fit using the Gini Index, for two primary reasons. First, we need a tight bound on the sensitivity of $u(X, \theta)$. While other machine learning models have been shown to outperform CART in many applications, we would have a far weaker bound on the sensitivity and would need to add much more noise. Secondly, as was shown in Snoke et al. (2018), CART models exhibit at least satisfactory performance in determining the distributional similarity. Future work may prove desirable bounds on the sensitivity of the $pMSE$ when using stronger classifiers, in which case those models should certainly be adopted.

We use the impurity function known as Gini Index from [Breiman et al. \(1984\)](#), defined as:

$$GI = \operatorname{argmin} \sum_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{m_i}\right), \quad (4)$$

where m_i are the total number of observations in each node, a_i are the number of observations labeled 1 in each node, and D is the total number of nodes. In practice these models are fit in a greedy manner for computational purposes. The process is to make the first optimal split that minimizes the impurity for two nodes, and then to make proceeding splits and adding additional nodes if doing so continues to minimize the impurity according to a chosen cost function. If computation is not a concern, it would also be possible for any fixed D to do a full grid search to determine the optimal D splits that minimize the impurity over $D + 1$ nodes. The difference between globally optimal and greedy trees is important for our theoretical results. In our theoretical results in section 4 we prove the sensitivity bound when trees are fit based on the globally optimal Gini Index, and in our simulations in section 5 we perform an empirical examination of how frequently the greedy fitting violates our theoretical results.

4 Theoretical Results for the Sensitivity Bound

In order to sample from the exponential mechanism, we need to give a bound on the sensitivity of the quality function. The $pMSE$ function is naturally bounded, but in Theorem 2 we prove a much tighter bound.

Theorem 2. *Given $u(X, \theta) = E_{\theta}[pMSE(X, X_{\theta}^s)|X, \theta]$ where $pMSE = \sum_{i=1}^{2n} \frac{(\hat{p}_i - 0.5)^2}{2n}$ with \hat{p}_i estimated from a classification tree with optimal splits found using the Gini Index. Then*

$$\Delta u = \sup_{\theta} \sup_{\delta(X, X')=1} |u(X, \theta) - u(X', \theta)| \leq \frac{1}{n},$$

where $X, X_{\theta}^s \in \mathbb{R}^{n \times q}$.

The proof is given in the Appendix. Intuitively, the proof follows from the fact that we can relate the $pMSE$ to the Gini Index. We can then bound the change in Gini Index given a change in one row of the input data due to the fact that we are finding the global optimum. We will not suddenly do much better or much worse. In fact, we can quantify exactly how much better or worse we can do, which leads to the bound.

This bound is nice because it decreases with n , meaning the noise added decreases as the number of observations increase. This bound matches the results derived for the sensitivity of the Gini Index when assuming discrete predictors from [Friedman and Schuster \(2010\)](#). Our proof shows this the bound remains the same when using continuous predictors. The result in [Friedman and Schuster \(2010\)](#) was used for performing classification under differential privacy, rather than producing synthetic data, and we see our extension of the proof to include continuous predictors as a useful side result of this paper.

It is important to note that this proof is for the theoretical case when we can find the optimal partitioning for any number of nodes. The greedy method can violate the bound because we can no longer control how much the Gini Index can change after changing one row. While it would be possible to use our method with a full grid search, computationally this is a poor idea. On the other hand, it is necessary in order to satisfy ϵ -DP. An alternative method could be to use adaptive composition, i.e., fit the CART models greedily but in a way that satisfies DP. We could then compose the privacy between fitting the CART model and sampling from the exponential distribution, which we explore in future work.

5 Empirical Failure Rate of the Sensitivity Bound

These simulations show the empirical rate for which the greedy fitting violates the bound. We can also view the maximum simulated value as an empirical estimate of the sensitivity for this particular dataset, but we are more interested in the failure rate. We generated datasets, X , with $q = 2$ and $n = 5000$.

$X_1 \sim N(2, 10)$ and $X_2 \sim N(-2.5 + 0.5x_1, 3)$, and we produced X' by taking X and adding random Gaussian noise, $N(0, 25)$, to each variable for one observation. We then drew a synthetic dataset X^s with $X_1^s \sim N(\theta_1, \theta_2)$ and $X_2^s \sim N(\theta_3 + \theta_4 x_1^s, \theta_5)$ where $\theta_i \sim N(0, 10)$. We estimated the $pMSE$ with respect to X^s for both X and X' and calculated the difference. Recall the theoretical sensitivity bound is $1/n = 0.0002$, and any values larger than this violate the bound. We repeated this process 1,000,000 times each using CARTs of depths 1, 2, 5, and unlimited for the $pMSE$ model. For all trees we included a complexity parameter (cp) requiring a certain percentage improvement in order to make an additional split. This parameter is necessary in order to not produce trees that are fully saturated (one terminal node per observation) when there is no depth limitation.

Table 1: Empirical failure rates of 1,000,000 simulations for the sensitivity bound when using the greedy CART fitting algorithm for different tree sizes and different complexity parameters.

Tree Depth	cp	Percentage Violating Bound
Depth 1	0.01	0.0%
Depth 2	0.01	0.3%
Depth 5	0.01	0.6%
Depth Unlimited	0.01	0.7%
Depth 1	0.001	0.0%
Depth 2	0.001	0.5%
Depth 5	0.001	2.0%
Depth Unlimited	0.001	2.5%

Figure 1 shows a sample of the simulated empirical sensitivity results. There are four groupings, for the trees with different depths, and darker points denote those violating the theoretical bound due to the greedy fitting algorithm. Table 1 gives the full results. The percentage of simulations which violate the bound increases with tree depth size, and in the unlimited case for $cp = 0.01$ the empirical failure rate is $\leq 1\%$. As expected, there are no results which violate the bound when only one split is made. This confirms our theoretical results because with only one split, greedy is equivalent to optimal, so the bound is never violated in simulation. The empirical sensitivity also depends on the cp , so we ran simulations for two different values. A lower cp will lead to larger trees (subject to depth constraints), which means we are making more greedy splits and increasing the chance of violating the bound.

6 Empirical Evaluation of Differentially Private Linear Regression

In order to assess the practical statistical utility of our method, we ran simulations testing the accuracy of an estimated linear regression model. Our method guarantees maximal distributional similarity of the synthetic data based on the $pMSE$ metric, but many researchers may be interested in more specific comparisons such as regression outputs.

We simulate datasets, X , in the same way as in Sections 5 and 7. Using this data, we regress X_2 on X_1 and get ordinary least squares (OLS) estimates of the intercept ($\hat{\beta}_0$) and slope ($\hat{\beta}_1$) coefficients. We calculate the absolute difference between these estimates and the corresponding estimates we get by fitting the same model with the differentially private methods, i.e., $|\hat{\beta} - \hat{\beta}^{priv}|$. Our comparison methods are the noisy Bayesian method from Bowen and Liu (2016) and the smooth histogram from Wasserman and Zhou (2010), both described in Section 2.1. We also compare with methods that do not produce synthetic data but produce DP regression estimates. Specifically, we use the Functional Mechanism of Zhang et al. (2012) and Awan and Slavkovic (2018), considering l_1 and l_∞ mechanisms, respectively. Note that these methods require the data to be bounded in the same way as the noisy Bayes and smooth histogram methods. Finally we compare with estimates from non-DP synthetic data, sampled from the unperturbed BPDD.

For the $pMSE$ mechanism, we carry out the simulations using trees of depths 1, 2, 5, and unlimited with a $cp = 0.01$. We see that the utility significantly improves as we move from depth 1 to 2 and from

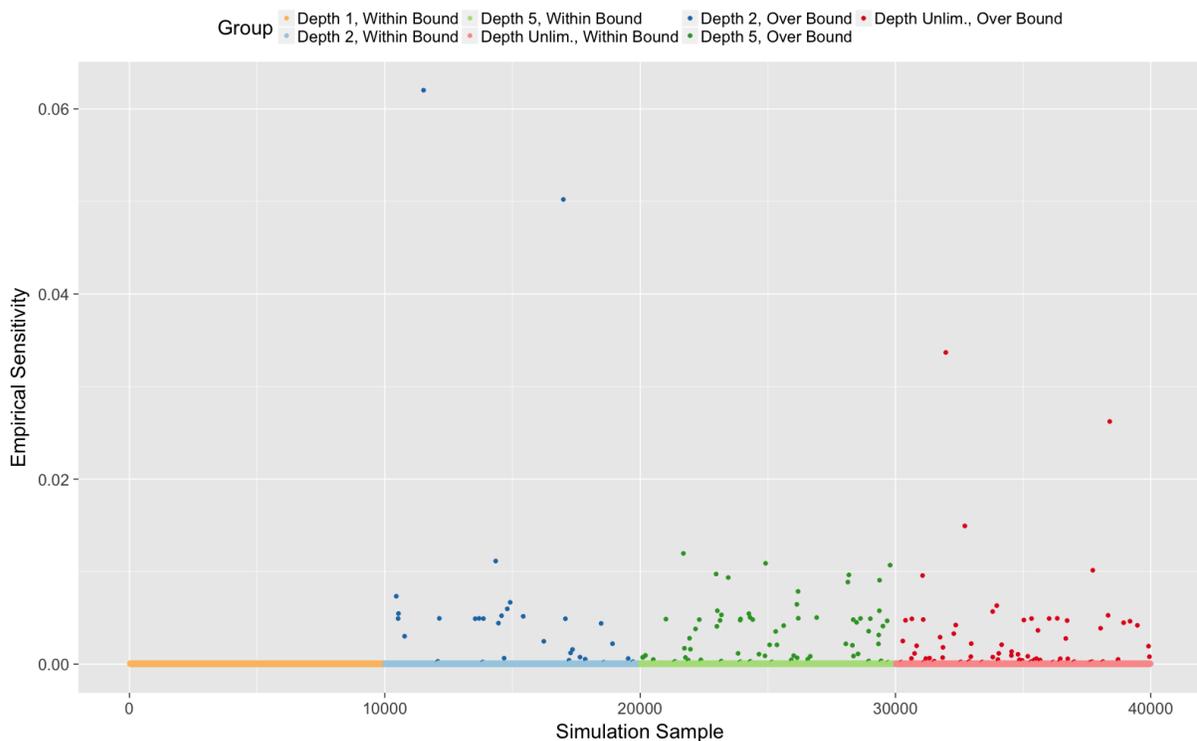


Fig. 1: Random sample of 10,000 simulations for each tree depth. Values shown are differences between $u(X, \theta)$ calculated with X and X' . Darker points violate the theoretical bound. $Cp = 1\%$.

2 to 5, but there is little change from 5 to unlimited. This is likely because trees of depth 5 are large enough to evaluate this dataset. The tree size is a potential tuning parameter for future work using this method. Astute readers may have noticed that the unnormalized distribution we propose for the $pMSE$ mechanism does not necessarily exist, since the probability in the tails remains very flat. To fix this, we add a very flat prior, $N(0, 100,000)$, to each of our parameters when sampling. The flatness ensures it does not affect the utility, but by adding it we also ensure the probability in the tails eventually goes to zero.

We run the mechanisms with values $\epsilon \in \{0.25, 0.5, 1\}$. For the $pMSE$ mechanism and the noisy BPDD we generate $l = 10$ private datasets each satisfying (ϵ/l) -DP, and for the smooth histogram and functional mechanisms we produce only one output satisfying ϵ -DP. This ensures all mechanisms satisfy the same level of privacy. The non-DP synthetic data method does not guarantee any privacy.

For the noisy Bayes, smooth histogram, and functional mechanism methods, we ran the simulations truncating the data at different assumed bounds. For both variables, we set these bounds at two, four, five, or ten times the standard deviation. Four or five can be thought of as roughly the appropriate bounds, since this is Gaussian data and most observations will fall into those ranges. Two was chosen to be a range that is too narrow and excludes part of the true distribution, and ten was chosen for a looser bound. We see from the results that the smaller bound achieves better results on the regression, even when it is more narrow even than the truth. This is an artefact of the model we chose, and if an even tighter bound was chosen it may have greater adverse effects. The loose bound (10 times) performs quite poorly, since we must add much more noise. Figure 2 visualizes the better performing results (limited only for readability). Full results for all tree sizes and bounds are shown in Figure 4 in the Appendix.

Overall our method outperforms the other two synthesis methods when using trees of depth 5 or unlimited, regardless of the bounds chosen. Even trees of depth 2 perform roughly the same or better than the other methods. For deeper trees, our method performs almost as well as the functional mechanism.

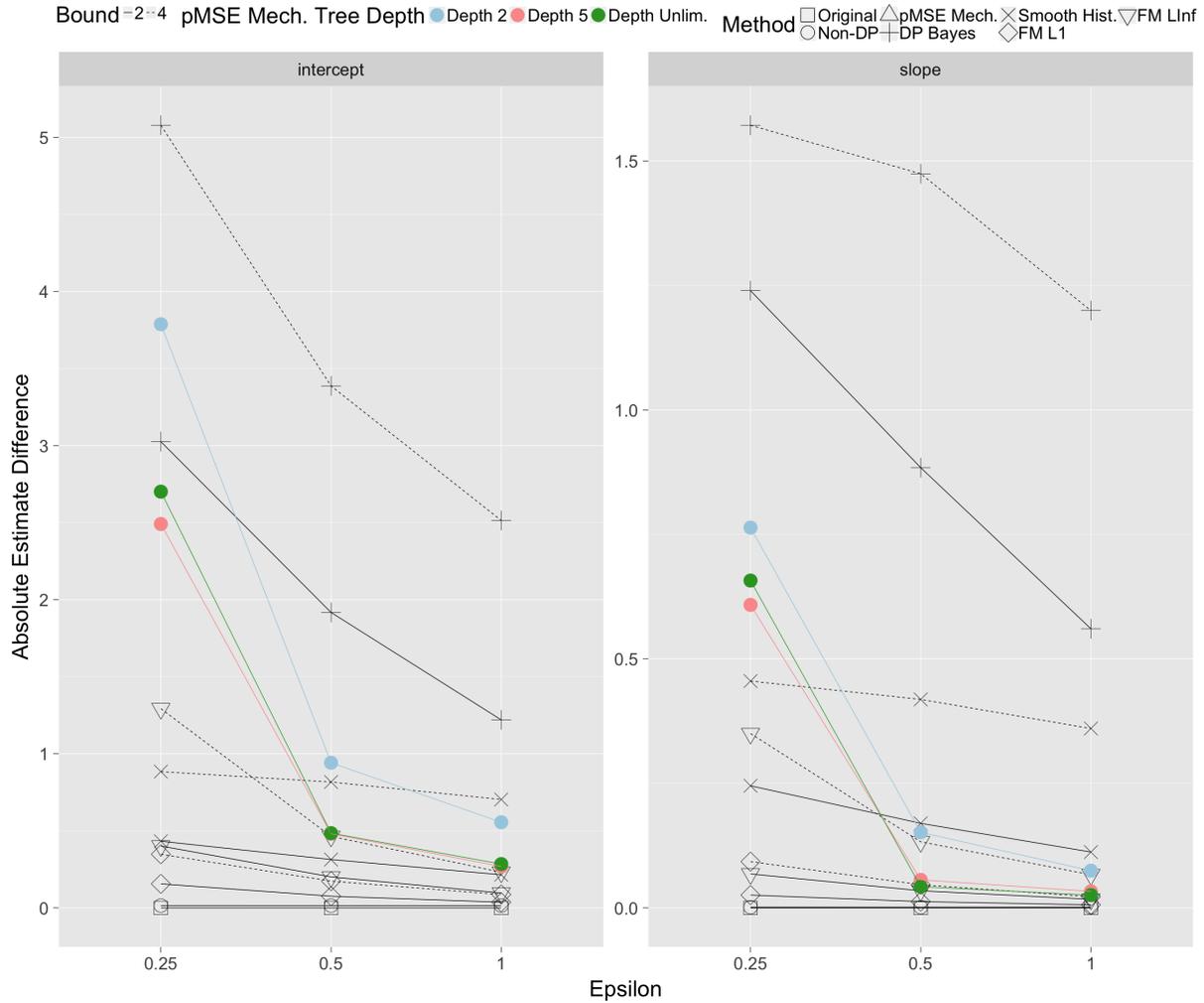


Fig. 2: Lineplots showing the mean simulation results. X-axis indicates different values of ϵ . Lines are also subdivided within methods by the tree depth and the bound.

This is quite encouraging, since that method focuses only on providing regression estimates rather than entire synthetic datasets. For a slight decrease in utility from the functional mechanism, our method provides an entire synthetic dataset, which can be used to fit any number of models using our synthetic data without changing the privacy guarantee. For the functional mechanism, it requires further splitting of the privacy parameter to estimate a different model.

These results show good performance, and further work should consider simulations with larger numbers of variables or a mixture of categorical and continuous variables. We expect our method will only improve against the other methods with more variables, since theoretically our method maximizes similarity on the entire distribution.

7 Empirical Evaluation of the $pMSE$

We guarantee theoretical maximization of the $pMSE$ for the differentially private synthetic data produced from the $pMSE$ mechanism, but as many practitioners know empirical tests often look slightly different from theory. To evaluate this, we ran simulations to estimate the $pMSE$ from datasets generated using our method, two other DP synthesis methods, and a standard non-DP synthesis method.

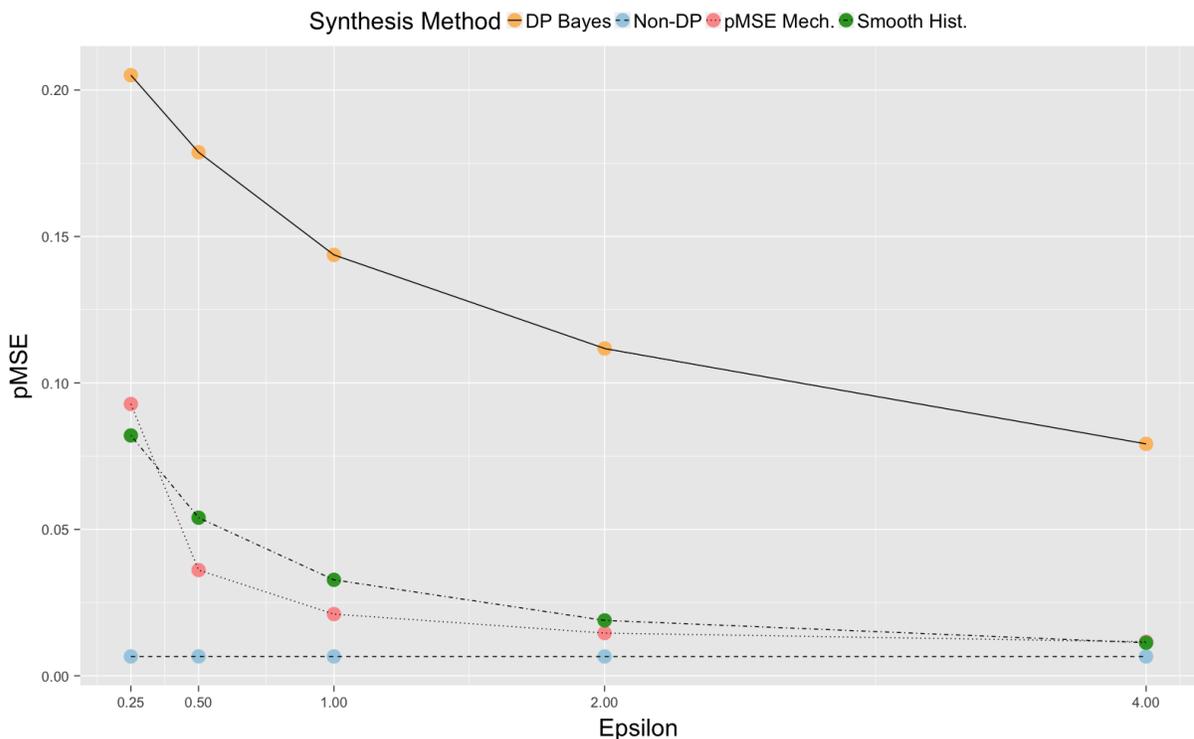


Fig. 3: Simulations results showing the mean $pMSE$ calculated using synthetic producing according to four different methods. $pMSE$ is calculated from comparison with original data, with values closer to 0 implying higher utility.

We again simulate datasets, X , in the same way as in Sections 5 and 6. For each X we generated synthetic datasets X^s and then calculated the $pMSE$ using X and X^s . Our four synthesis methods were the $pMSE$ mechanism, the noisy Bayesian method from [Bowen and Liu \(2016\)](#), the smooth histogram from [Wasserman and Zhou \(2010\)](#), and sampling from the non-differentially private BPPD using fully conditional sequential models. We ran 2,500 simulations each for five different values of $\epsilon \in \{0.25, 0.5, 1, 2, 4\}$. For our method we used CART trees with unlimited depth, and for the other two DP methods we assumed a bound on the data of four times the standard deviation, which is roughly the correct bound given that it is Gaussian data. Figure 3 shows the results for the mean simulations results and Table 2 shows the full mean and variance of the results.

We see that as expected, the $pMSE$ mechanism offers either the best or one of the best values of the $pMSE$ among the methods guaranteeing differential privacy. Interestingly the smooth histogram method is fairly good as well, even offering comparable values at $\epsilon = 0.25$ or $\epsilon = 4$. The noisy BPPD method on the other hand is quite bad, even at high levels of ϵ , so should be used with caution.

We also see that the variance in the estimated $pMSE$ values changes quite a bit depending on the method and level of ϵ . Both the $pMSE$ mechanism and the smooth histogram show higher variances for either low (0.25) or high (4) values of ϵ , while the noisy BPPD method increases in variance as ϵ increases. This variance is something to keep in mind both in choosing a protection method and in developing the practical implementation. We could likely improve our current implementation of the $pMSE$ mechanism in order to better sample noisy parameters and generate synthetic data with less variance in the resulting $pMSE$. On the other hand, it should also be expected that the variance decreases some as ϵ grows because we are adding less noise through the privacy mechanism.

Comparing the DP methods to the traditional synthetic data approach, we see that the best method at $\epsilon = 4$ produces an average $pMSE$ roughly two times that from the non-DP synthesis method, which is

Table 2: Simulation results giving the mean and variance of the $pMSE$ values calculated using four different synthesis methods and five different levels of ϵ .

ϵ	Simulated Values	Non-DP	$pMSE$ Mech.	DP Bayes	Smooth Hist.
0.25	$pMSE$ Mean	0.00660	0.09281	0.20509	0.08206
	$pMSE$ Var.	8.681e-07	1.347e-03	3.079e-03	2.826e-05
0.5	$pMSE$ Mean	0.00663	0.03610	0.17876	0.05398
	$pMSE$ Var.	8.929e-07	8.020e-05	5.914e-03	1.809e-05
1	$pMSE$ Mean	0.00661	0.02107	0.14372	0.03278
	$pMSE$ Var.	8.342e-07	1.648e-05	8.579e-03	1.069e-05
2	$pMSE$ Mean	0.00660	0.01459	0.11177	0.01892
	$pMSE$ Var.	8.342e-07	1.648e-05	8.579e-03	1.069e-05
4	$pMSE$ Mean	0.00660	0.01161	0.07919	0.01129
	$pMSE$ Var.	8.671e-07	3.329e-06	9.087e-03	5.964e-06

producing synthetic data from the correct generative model. This is actually quite good considering we are adding the strong guarantee of DP. Even for $\epsilon = 1$ our method produces $pMSE$ values only roughly three times that of the non-DP synthetic data.

8 Conclusions and Future Work

The $pMSE$ mechanism we propose provides a novel flexible method to produce high-quality synthetic datasets guaranteeing ϵ -DP. By sampling generative model parameters from the exponential mechanism and using the $pMSE$ as our quality function, we produce synthetic data with maximal distributional similarity to the original data. By using the $pMSE$, we ensure the sensitivity depends neither on the dimension nor the range of the data, and the bound decreases as we increase the sample size. This allows us to use this mechanism for continuous data, and the amount of noise we add will not grow with the dimension (apart from sampling from a more complex distribution).

Our simulations in Sections 6 and 7 confirm that the $pMSE$ mechanism generally performs as well or better than the other standard DP synthesis methods. In the case of linear regression the $pMSE$ mechanism even performs roughly as well as methods that produce estimates of regression coefficients only rather than entire synthetic datasets. In the case of the empirical $pMSE$, as expected our method performs worse than non-DP synthetic data, but the utility cost seems reasonable for the privacy gain.

The $pMSE$ mechanism relies on defining an appropriate form for the generative distribution from which to draw synthetic values. It is possible to misspecify this model, which would lead to poor utility. This is one drawback of the synthetic approach as opposed to simply adding noise. Fortunately, this aspect has been addressed in great detail in the synthetic data literature, so we feel that finding an appropriate model is possible without too much difficulty.

Our primary limitation is the computational feasibility to ensure the theoretical sensitivity bound. From the empirical simulations we saw that the bound does not always hold when using typical greedy fitting algorithms. Fitting the models using the global optimum would ensure the theoretical bound and guarantee ϵ -DP. Proposals have been made to carry out machine learning using global optimums, such as [Bertsimas and Dunn \(2017\)](#), so methods may exist to aid the computation.

An alternative implementation of our method would be to consider fitting the CART models in a way that satisfies ϵ -DP and then composing this ϵ with that from sampling from the $pMSE$ mechanism. This is similar to the approach of [Li et al. \(2018\)](#). This is desirable because we could use any standard CART software to implement the method. Other future work could consider using different impurity measures than the Gini Index, deriving measures of choosing the best tree size, or best practices for sampling from the unnormalized density we get through the exponential mechanism.

Bibliography

- Abowd, J. M. and L. Vilhuber (2008). How protective are synthetic data? In *International Conference on Privacy in Statistical Databases*, pp. 239–246. Springer.
- Awan, J. and A. Slavkovic (2018). Structure and sensitivity in differential privacy: Comparing k-norm mechanisms. *arXiv preprint arXiv:1801.09236*.
- Barak, B., K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. 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*, pp. 273–282. ACM.
- Bertsimas, D. and J. Dunn (2017). Optimal classification trees. *Machine Learning* 106(7), 1039–1082.
- Blum, A., C. Dwork, F. McSherry, and K. Nissim (2005). Practical privacy: the sulq framework. In *Proceedings of the twenty-fourth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 128–138. ACM.
- Bowen, C. M. and F. Liu (2016). Comparative study of differentially private data synthesis methods. *arXiv preprint arXiv:1602.01063*.
- Breiman, L., J. H. Friedman, R. A. Olshen, and C. J. Stone (1984). *Classification and Regression Trees*. Belmont, Wadsworth, CA.
- Bun, M. and T. Steinke (2016). Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pp. 635–658. Springer.
- Charest, A.-S. (2011). How can we analyze differentially-private synthetic datasets? *Journal of Privacy and Confidentiality* 2(2), 3.
- Chaudhuri, K., A. Sarwate, and K. Sinha (2012). Near-optimal differentially private principal components. In *Advances in Neural Information Processing Systems*, pp. 989–997.
- Drechsler, J. (2011). *Synthetic datasets for statistical disclosure control: theory and implementation*, Volume 201. Springer Science & Business Media.
- Dwork, C., K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor (2006). Our data, ourselves: Privacy via distributed noise generation. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pp. 486–503. Springer.
- Dwork, C., F. McSherry, K. Nissim, and A. Smith (2006). Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography*, pp. 265–284. Springer.
- Dwork, C., M. Naor, T. Pitassi, G. N. Rothblum, and S. Yekhanin (2010). Pan-private streaming algorithms. In *ICS*, pp. 66–80.
- Dwork, C., A. Roth, et al. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9(3–4), 211–407.
- Dwork, C. and G. N. Rothblum (2016). Concentrated differential privacy. *arXiv preprint arXiv:1603.01887*.
- Friedman, A. and A. Schuster (2010). Data mining with differential privacy. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 493–502. ACM.
- Kapralov, M. and K. Talwar (2013). On differentially private low rank approximation. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1395–1414. SIAM.
- Karwa, V., P. N. Krivitsky, and A. B. Slavković (2017). Sharing social network data: differentially private estimation of exponential family random-graph models. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 66(3), 481–500.
- Karwa, V., A. Slavković, et al. (2016). Inference using noisy degrees: Differentially private β -model and synthetic graphs. *The Annals of Statistics* 44(1), 87–112.
- Kasiviswanathan, S. P., H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith (2011). What can we learn privately? *SIAM Journal on Computing* 40(3), 793–826.
- Kifer, D., A. Smith, and A. Thakurta (2012). Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pp. 25–1.
- Li, Karwa, Slavković, and Steorts (2018). Release of differentially private high dimensional histograms. *Pre-print*.

- Li, C., M. Hay, V. Rastogi, G. Miklau, and A. McGregor (2010). Optimizing linear counting queries under differential privacy. In *Proceedings of the twenty-ninth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 123–134. ACM.
- Machanavajjhala, A., D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber (2008). Privacy: Theory meets practice on the map. In *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, pp. 277–286. IEEE.
- McSherry, F. and K. Talwar (2007). Mechanism design via differential privacy. In *Foundations of Computer Science, 2007. FOCS’07. 48th Annual IEEE Symposium on*, pp. 94–103. IEEE.
- McSherry, F. D. (2009). Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pp. 19–30. ACM.
- Nissim, K., S. Raskhodnikova, and A. Smith (2007). Smooth sensitivity and sampling in private data analysis. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pp. 75–84. ACM.
- Nissim, K., T. Steinke, A. Wood, M. Altman, A. Bembenek, M. Bun, M. Gaboardi, D. R. O’Brien, and S. Vadhan (2017). Differential privacy: A primer for a non-technical audience.
- Raab, G. M., B. Nowok, and C. Dibben (2017). Practical data synthesis for large samples. *Journal of Privacy and Confidentiality* 7(3), 4.
- Raghunathan, T. E., J. P. Reiter, and D. B. Rubin (2003). Multiple imputation for statistical disclosure limitation. *Journal of Official Statistics* 19(1), 1–17.
- Reiter, J. P. (2002). Satisfying disclosure restrictions with synthetic data sets. *Journal of Official Statistics* 18(4), 531–544.
- Reiter, J. P. (2005). Using cart to generate partially synthetic, public use microdata. *Journal of Official Statistics* 21(3), 441–462.
- Snoke, J., G. M. Raab, B. Nowok, C. Dibben, and A. Slavković (2018). General and specific utility for synthetic data. *Journal of the Royal Statistical Society Series A: Statistics in Society*.
- Wang, Y.-X., J. Lei, and S. E. Fienberg (2016). On-average kl-privacy and its equivalence to generalization for max-entropy mechanisms. In *International Conference on Privacy in Statistical Databases*, pp. 121–134. Springer.
- Wasserman, L. and S. Zhou (2010). A statistical framework for differential privacy. *Journal of the American Statistical Association* 105(489), 375–389.
- Woo, M.-J., J. P. Reiter, A. Oganian, and A. F. Karr (2009). Global measures of data utility for microdata masked for disclosure limitation. *Journal of Privacy and Confidentiality* 1, 111–124.
- Xu, J., Z. Zhang, X. Xiao, Y. Yang, G. Yu, and M. Winslett (2013). Differentially private histogram publication. *The VLDB Journal* 22(6), 797–822.
- Zhang, J., G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao (2017). Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)* 42(4), 25.
- Zhang, J., Z. Zhang, X. Xiao, Y. Yang, and M. Winslett (2012). Functional mechanism: regression analysis under differential privacy. *Proceedings of the VLDB Endowment* 5(11), 1364–1375.

9 Appendix: Proof of Theorem 4.1

Proof. We first show that using the expected value, and approximating it, can be bounded above by the supremum across all possible datasets X^s generated using θ .

$$\Delta u = \sup_{\theta} \sup_{\delta(X, X')=1} |u(X, \theta) - u(X', \theta)| \quad (5)$$

can be rewritten as

$$\Delta u = \sup_{\theta} \sup_{\delta(X, X')=1} |E_{\theta}[pMSE(X, X_{\theta}^s)|X, \theta] - E_{\theta}[pMSE(X', X_{\theta}^s)|X, \theta]| \quad (6)$$

where $u(X, \theta) = E_\theta[pMSE(X, X_\theta^s)|X, \theta]$. Since the absolute value is a convex function, we can apply Jensen's inequality and get

$$\leq \sup_{\theta} \sup_{\delta(X, X')=1} E_\theta[|pMSE(X, X_\theta^s) - pMSE(X', X_\theta^s)||X, \theta]. \quad (7)$$

Then by taking the supremum over any data set X_θ^s , we obtain

$$\leq \sup_{X_\theta^s} \sup_{\delta(X, X')=1} |pMSE(X, X_\theta^s) - pMSE(X', X_\theta^s)|. \quad (8)$$

This also bounds our approximation of the expected value that we propose to use in practice, since the supremum is also greater than or equal to the sample mean.

Now writing this explicitly in terms of the CART model, we get

$$\sup_{a_i, m_i, a'_i, m'_i} \frac{1}{2n} \left| \sum_{i=1}^{D+1} m_i \left(\frac{a_i}{m_i} - 0.5 \right)^2 - m'_i \left(\frac{a'_i}{m'_i} - 0.5 \right)^2 \right| \quad (9)$$

where a_i, m_i , and D are defined as before, and a'_i and m'_i are the corresponding values for the model fit using X' . Expanding this we get

$$\sup_{a_i, m_i, a'_i, m'_i} \frac{1}{2n} \left| \sum_{i=1}^{D+1} \left(\frac{a_i^2}{m_i} - a_i - 0.25m_i \right) - \left(\frac{a_i'^2}{m_i'} - a'_i - 0.25m_i' \right) \right| \quad (10)$$

and we can cancel the third terms because $\sum_{i=1}^{D+1} m_i = \sum_{i=1}^{D+1} m_i'$. When we multiple by $2n$, the remaining inside term is equivalent to the sensitivity of the impurity, i.e.,

$$\sup_{a_i, m_i, a'_i, m'_i} \left| GI(X, X^s, D) - GI(X', X^s, D) \right| = \Delta GI \quad (11)$$

By bounding the impurity, we bound the $pMSE$. We can rewrite the above as

$$\left| \min_D GI(X, X^s, D) - \min_D GI(X', X^s, D) \right| \quad (12)$$

since the optimal CART model finds the minimum impurity across any D . The greatest possible difference then is the difference between these two optimums. And we can bound this above by

$$\leq \sup_D \left| GI(X, X^s, D) - GI(X', X^s, D) \right|. \quad (13)$$

Let X^{comb} and X'^{comb} be the combined data matrices as described in Algorithm 1, including the 0, 1 outcome variable. Recall that only one record has changed between X^{comb} and X'^{comb} (total number of records staying fixed), and it is labeled 0. We know that for a given D optimal split points producing $D + 1$ nodes on X^{comb} , there are a_i records labeled 1 and \tilde{m}_i total records in each bin, such that $\exists j \neq k \neq l_1 \neq \dots \neq l_{D-1}$ s.t. $\tilde{m}_j - m_j = m_k - \tilde{m}_k = 1$, $\tilde{m}_{l_v} = m_{l_v}$ for $v = \{1, \dots, D - 1\}$. In the same way, for a given D optimal split points producing $D + 1$ nodes on X'^{comb} , there are a'_i records labeled 1 and \tilde{m}'_i total records in each bin, such that $\exists j' \neq k' \neq l'_1 \neq \dots \neq l'_{D-1}$ s.t. $\tilde{m}'_{j'} - m'_{j'} = m'_{k'} - \tilde{m}'_{k'} = 1$, $\tilde{m}'_{l'_v} = m'_{l'_v}$ for $v = \{1, \dots, D - 1\}$. What this simply means is that after changing one record, the discrete counts in the nodes change by at most one in two of the nodes and does not change in the other $D - 1$ nodes.

Due to the fact that the CART model produces the D splits that minimize the impurity, we know both that

$$\sum_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{m'_i} \right) \leq \sum_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{\tilde{m}_i} \right) \quad (14)$$

and

$$\Sigma_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{m_i}\right) \leq \Sigma_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{\tilde{m}'_i}\right). \quad (15)$$

The inequality 14 implies that after changing one record, if new split points are chosen, the impurity must be equivalent or better than simply keeping the previous splits and changing the counts. The inequality 15 implies that the first split points chosen must be equivalent or better than using the new splits with the changed counts. If this were not the case, the first split points would have never been made in the first place. These lead to the final step.

Because we have an absolute value, we consider two cases.

Case 1: $GI(X, X^s, D) \geq GI(X', X^s, D)$

$$\begin{aligned} \sup_D \left| \Sigma_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{m_i}\right) - \Sigma_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{m'_i}\right) \right| &\leq \\ \sup_D \left| \Sigma_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{\tilde{m}'_i}\right) - \Sigma_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{m'_i}\right) \right| &= \\ \left| a'_{j'} \left(1 - \frac{a'_{j'}}{\tilde{m}'_{j'}}\right) - a'_{j'} \left(1 - \frac{a'_{j'}}{m'_{j'}}\right) + a'_{k'} \left(1 - \frac{a'_{k'}}{\tilde{m}'_{k'}}\right) - a'_{k'} \left(1 - \frac{a'_{k'}}{m'_{k'}}\right) \right| &= \\ \left| \frac{a'^2_{j'} (\tilde{m}'_{j'} - m'_{j'})}{\tilde{m}'_{j'} m'_{j'}} + \frac{a'^2_{k'} (\tilde{m}'_{k'} - m'_{k'})}{\tilde{m}'_{k'} m'_{k'}} \right| = \left| \frac{a'^2_{j'}}{\tilde{m}'_{j'} m'_{j'}} - \frac{a'^2_{k'}}{\tilde{m}'_{k'} m'_{k'}} \right| &\leq 2 \quad (16) \end{aligned}$$

The last step we know because $a_i \leq m_i$, and $\frac{n^2}{n(n-1)} \leq 2$.

Case 2: $GI(X', X^s, D) \geq GI(X, X^s, D)$

$$\begin{aligned} \sup_D \left| \Sigma_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{m_i}\right) - \Sigma_{i=1}^{D+1} a'_i \left(1 - \frac{a'_i}{m'_i}\right) \right| &\leq \\ \sup_D \left| \Sigma_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{\tilde{m}_i}\right) - \Sigma_{i=1}^{D+1} a_i \left(1 - \frac{a_i}{m_i}\right) \right| &= \\ \left| a_j \left(1 - \frac{a_j}{\tilde{m}_j}\right) - a_j \left(1 - \frac{a_j}{m_j}\right) + a_k \left(1 - \frac{a_k}{\tilde{m}_k}\right) - a_k \left(1 - \frac{a_k}{m_k}\right) \right| &= \\ \left| \frac{a^2_j (\tilde{m}_j - m_j)}{\tilde{m}_j m_j} + \frac{a^2_k (\tilde{m}_k - m_k)}{\tilde{m}_k m_k} \right| = \left| \frac{a^2_j}{\tilde{m}_j m_j} - \frac{a^2_k}{\tilde{m}_k m_k} \right| &\leq 2 \quad (17) \end{aligned}$$

Finally, this gives us $\Delta GI \leq 2 \implies \frac{\Delta GI}{2n} = \Delta u \leq \frac{1}{n}$.

10 Appendix: Full Simulation Results



Fig. 4: Boxplots showing simulation results. The rows indicate the different coefficients, and the columns indicate different values of ϵ . Boxplots are also subdivided within methods by the tree depth (for the *pMSE* mechanism method) and the bound (for others).