
Differentially Private Distributed Bayesian Linear Regression with MCMC

Bariş Alparslan¹ Sinan Yıldırım¹ Ş. İlker Birbil²

Abstract

We propose a novel Bayesian inference framework for distributed differentially private linear regression. We consider a distributed setting where multiple parties hold parts of the data and share certain summary statistics of their portions in privacy-preserving noise. We develop a novel generative statistical model for privately shared statistics, which exploits a useful distributional relation between the summary statistics of linear regression. We propose Bayesian estimation of the regression coefficients, mainly using Markov chain Monte Carlo algorithms, while we also provide a fast version that performs approximate Bayesian estimation in one iteration. The proposed methods have computational advantages over their competitors. We provide numerical results on both real and simulated data, which demonstrate that the proposed algorithms provide well-rounded estimation and prediction.

1. Introduction

Linear regression is a mathematical method that lies at the core of statistical research. Many researchers have been working on linear regression since the 19th century, and hence, many well-known solution methods exist. On a separate note, privacy-preserving statistical learning has gained popularity and importance in recent years, with *differential privacy* prevailing as the most commonly used definition for privacy (Dwork, 2006; Dwork et al., 2014a; Dankar & El Emam, 2013). As a result, there is a recent but growing interest in differentially private linear regression.

Many works in the data privacy literature do not mainly focus on regression but are motivated by or can be applied to regression. As an example, differentially private empiri-

cal risk minimisation (Chaudhuri et al., 2011; Bassily et al., 2014; Abadi et al., 2016; Kuru et al., 2022) can be applied to regression once it is cast as a data-driven optimisation problem. Many general-purpose Bayesian differentially private estimation methods can also be used in regression problems. Williams & Mcsherry (2010) is one of the first works that considered a hierarchical model for the privatised data and Bayesian estimation for the model parameters. Zhang et al. (2016) analyse several differential privacy mechanisms for posterior sampling and suggest using these mechanisms also for linear regression. Dimitrakakis et al. (2017) developed a posterior sampling query algorithm to combine differential privacy and Bayesian inference. Contrary to those one-sample approaches, general-purpose differentially private Markov chain Monte Carlo (MCMC) algorithms, which aim to identify the whole posterior distribution via iterative sampling, can also be applied to regression (Wang et al., 2015; Foulds et al., 2016; Wang et al., 2015; Yıldırım & Ermiş, 2019; Heikkilä et al., 2019; Gong, 2022; Alparslan & Yıldırım, 2022; Ju et al., 2022).

Several works in the literature are somewhat more directly related to differentially private regression. Zhang et al. (2012) have suggested a functional mechanism method, which is based on perturbing polynomial objective functions with privacy-preserving noise. As an alternative, Dwork et al. (2014b); Wang (2018) considered perturbation of summary statistics. Alabi et al. (2022) provide a technical discussion on different point estimation methods for differentially private simple linear regression, (that is when we have a single feature). Ferrando et al. (2022) present a method to compute confidence intervals for the coefficients of linear regression. Cai et al. (2021) study the rates of convergence for parameter estimation with differential privacy via output perturbation, where a non-private estimator is perturbed. All those works consider point estimation of the linear regression parameters.

In this paper, we focus on differentially private distributed Bayesian inference for the parameters of linear regression. We use a novel hierarchical model that relies on a distributional relationship (Proposition 3.1) between the summary statistics of linear regression, which, to the best of our knowledge, has not been exploited so far. We propose Bayesian inference algorithms that take perturbations of summary statistics as observations. The general inferential

¹Faculty of Engineering and Sciences, Sabancı University, Turkey ²Amsterdam Business School, University of Amsterdam, The Netherlands. Correspondence to: Bariş Alparslan <baris.alparslan@sabanciuniv.edu>.

tool we pick in this paper is MCMC, a well-known framework for iterative sampling from posterior distributions. As we shall see, the proposed MCMC algorithms in this paper already have lower computational complexities per iteration than their closest competitors in [Bernstein & Sheldon \(2019\)](#). Additionally, we also propose much faster Bayesian estimation methods that perform estimation in one iteration. Finally, for the sake of generality, we assume a distributed setting where the total dataset is shared among multiple parties (data nodes), who want to collaborate for the inference of a common parameter, see *e.g.*, [Heikkilä et al. \(2017\)](#) for such a setting. The non-distributed setting is just a special case (single data holder) for our methodology.

This paper has connections with several works in the literature, yet it has significant differences from each of those, as we shall explain below.

For the privacy-preserving mechanism, we consider adding noise to summary statistics of linear regression, similarly to [Wang \(2018\)](#); [Bernstein & Sheldon \(2019\)](#). The adaSSP framework of [Wang \(2018\)](#) also motivates the fast Bayesian estimation methods developed in this paper. However, adaSSP is a point estimation method while we aim for a posterior distribution. The latter work, [Bernstein & Sheldon \(2019\)](#), is particularly related to this paper as they also study Bayesian linear regression with differential privacy using perturbed statistics of data. However, there are some important differences between our work and that of [Bernstein & Sheldon \(2019\)](#). These differences stem from the choice of summary statistics and the consequent hierarchical structure used for modelling linear regression. Those modelling differences lead to significant differences in the inference methods as well as significant computational advantages for our methods. Specifically, the computational complexity of our methods is $\mathcal{O}(d^3)$, where d is the number of features. This order is much less than $\mathcal{O}(d^6)$ of [Bernstein & Sheldon \(2019\)](#). Finally, neither [Wang \(2018\)](#) nor [Bernstein & Sheldon \(2019\)](#) has considered a distributed learning setting as we do in this paper, although both works can be modified for the distributed setting after moderate modifications.

[Foulds et al. \(2016\)](#) and [Heikkilä et al. \(2017\)](#) are other differentially Bayesian inference methods that target posterior distributions of perturbed summary statistics of sensitive data. [Heikkilä et al. \(2017\)](#) is particularly interesting because they consider a distributed setting and present linear regression as their showcase example. However, we differ from those works in the way we model the perturbed statistics and in the choice of inference methods. Specifically, [Foulds et al. \(2016\)](#); [Heikkilä et al. \(2017\)](#) treat the perturbed statistics as if not perturbed, while we correctly incorporate the effect of perturbation in our model.

Recently, [Alparslan & Yıldırım \(2022\)](#) and [Ju et al. \(2022\)](#) employ data augmentation for modelling sensitive and pri-

vised data and propose MCMC for Bayesian inference, the latter having linear regression as a major application. Their methods have $\mathcal{O}(n)$ complexity per iteration in general where n is the number of instances in the data set, which can be slow when n is large. In contrast, our methods are scalable in data size since their computational complexities do not depend on n . We note that [Alparslan & Yıldırım \(2022, Section 4.2\)](#) also present an MCMC method scalable with n that exploits the approximate normality of additive summary statistics. However, a direct application of that would lead to an algorithm with $\mathcal{O}(d^6)$ computational complexity (per iteration), like in [Bernstein & Sheldon \(2019\)](#).

The paper is organised as follows: In [Section 2](#), we review differential privacy. In [Section 3](#), we lay out the hierarchical model for differentially private distributed linear regression with perturbed summary statistics. In [Section 4](#), we present and discuss the aspects of the proposed inference algorithms. In [Section 5](#), we provide numerical experiments. We conclude in [Section 6](#).

Notation: Matrices and vectors are shown in bold-face notation. For a matrix \mathbf{A} , its transpose, trace, and determinant (if they exist) are \mathbf{A}^T , $\text{tr}(\mathbf{A})$, and $|\mathbf{A}|$, respectively. \mathbf{I}_d is the $d \times d$ identity matrix. For any sequence $\{a_i\}_{i \geq 0}$, we write $a_{i:j}$ for (a_i, \dots, a_j) . We write $x \sim P$ to mean the random variable x has distribution P . $\mathcal{N}(\mathbf{m}, \Sigma)$ stands for the multivariate normal distribution with mean \mathbf{m} and covariance Σ . The Wishart and inverse-Wishart distributions, each with scale matrix Λ and κ degrees of freedom, are shown as $\mathcal{W}(\Lambda, \kappa)$ and $\mathcal{IW}(\Lambda, \kappa)$, respectively. $\mathcal{IG}(a, b)$ stands for the inverse-gamma distribution with shape and scale parameters a and b . We augment those notations, *e.g.*, with \mathbf{x} , to denote the respective probability density functions (pdf), *e.g.*, $\mathcal{N}(\mathbf{x}; \mathbf{m}, \Sigma)$.

2. Differential Privacy

Differential privacy ([Dwork, 2006; 2008](#)) concerns randomised algorithms that run on sensitive, or usually private, data. A randomised algorithm takes an input data set $D \in \mathcal{D}$ and returns a random output in \mathcal{O} where the randomness is intrinsic to the algorithm. A differentially private algorithm constrains the difference between the probability distributions of the output values obtained from neighbouring data sets. We say two data sets are *neighbours* if they have the same size and differ by a single element, corresponding to a single individual’s piece of data.

Definition 2.1 (Differential privacy). A randomised algorithm $M : \mathcal{D} \mapsto \mathcal{O}$ is (ϵ, δ) -differentially private (DP) if for any pair of neighbouring data sets $D, D' \in \mathcal{D}$ and for any subset $O \subseteq \mathcal{O}$ of the of support domain, it satisfies

$$\mathbb{P}[M(D) \in O] \leq e^\epsilon \mathbb{P}[M(D') \in O] + \delta.$$

The definition implies that we have more privacy with smaller (ϵ, δ) pair. Privacy-preserving algorithms often use noise-adding mechanisms. A popular noise-adding mechanism is the *Gaussian mechanism* (Dwork et al., 2006), which perturbs a function $f : \mathcal{D} \mapsto \mathbb{R}^k$ of the sensitive data, for some $k \geq 1$, with a random noise drawn from the Gaussian distribution. The amount of the added noise depends on the L_2 -sensitivity of the function, given by

$$\Delta_f = \max_{\text{neighbour } D_1, D_2 \in \mathcal{D}} \|f(D_1) - f(D_2)\|_2.$$

An (ϵ, δ) -DP Gaussian mechanism returns

$$f(D) + \Delta_f \sigma(\epsilon, \delta) \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k) \quad (1)$$

upon taking D as the input, where the quantity $\sigma(\epsilon, \delta)$ ensures (ϵ, δ) -DP. In this work, we take $\sigma(\epsilon, \delta)$ as the analytical solution given in Balle & Wang (2018, Algorithm 1) due to its tightness. The Gaussian mechanism is also central to other forms of privacy, such as zero-concentrated DP (Bun & Steinke, 2016) and Gaussian DP (Dong et al., 2022).

This paper considers (ϵ, δ) -DP as the type of privacy and the Gaussian mechanism to generate noisy observations. Moreover, the proposed methods in this paper never use the sensitive data given the noisy observations generated using the Gaussian mechanism, hence exploiting the *post-processing* property of differential privacy (Dwork & Roth, 2014).

Theorem 2.2 (Post-processing). *If $M : \mathcal{D} \mapsto \mathcal{O}$ be (ϵ, δ) -DP and let $f : \mathcal{O} \rightarrow \mathcal{O}'$ be another mapping independent of D given $M(D)$. Then $f_M : \mathcal{D} \mapsto \mathcal{O}'$ with $f_M(D) = f(M(D))$ is (ϵ, δ) -DP.*

3. Differentially Private Distributed Linear Regression

This section presents a new hierarchical model for differentially private distributed linear regression. For ease of exposition, we first present a model with a single data holder, then generalise the model for the distributed setting.

3.1. Basic Model and Privacy Setup

Suppose we have a sequence of random variables $\{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$, where $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^{d \times 1}$ are the feature vectors and $y_i \in \mathcal{Y} \subseteq \mathbb{R}$ is the i 'th response variable. We consider the normal linear regression to model the dependency between \mathbf{x}_i and y_i . Specifically,

$$y_i = \mathbf{x}_i^T \boldsymbol{\theta} + e_i, \quad e_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_y^2), \quad i = 1, \dots, n,$$

where $\boldsymbol{\theta} \in \mathbb{R}^d$ is the vector of the linear regression coefficients. We assume that the feature vectors \mathbf{x}_i 's are i.i.d. with distribution $P_{\mathbf{x}}$. A particular case of interest will be

one where $P_{\mathbf{x}}$ can be assumed to be a normal distribution. However, we will also present algorithms for general $P_{\mathbf{x}}$ that is not (even approximately) normal.

In matrix notation, the above can shortly be expressed as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}_n),$$

where $\mathbf{X} = [\mathbf{x}_1^T \dots \mathbf{x}_n^T]^T$ is the so-called design matrix, $\mathbf{y} = [y_1 \dots y_n]^T$. Additionally, we also define the summary statistics of \mathbf{X} and \mathbf{y} given by

$$\mathbf{S} := \mathbf{X}^T \mathbf{X} \text{ and } \mathbf{z} := \mathbf{X}^T \mathbf{y},$$

respectively. In this paper, we assume a data-sharing scenario where \mathbf{S} and \mathbf{z} are privately released as the noisy summary statistics $\hat{\mathbf{S}}$ and $\hat{\mathbf{z}}$, constructed as

$$\hat{\mathbf{S}} = \mathbf{S} + \sigma_s \mathbf{M}, \quad (2)$$

$$\hat{\mathbf{z}} = \mathbf{z} + \sigma_z \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad (3)$$

where \mathbf{M} is a $d \times d$ symmetric matrix with its upper triangular elements drawn from $\mathcal{N}(0, 1)$. Dwork et al. (2014b) arrange σ_s and σ_z so that both (2) and (3) are $(\epsilon/2, \delta/2)$ differentially private, leading to (ϵ, δ) -DP overall. Differently than Dwork et al. (2014b), we set

$$\sigma_s = \sigma_z = \Delta_{sz} \sigma(\epsilon, \delta),$$

where $\sigma(\epsilon, \delta)$ is given in Balle & Wang (2018, Algorithm 1), and Δ_{sz} is the overall L_2 sensitivity of $[\mathbf{S}, \mathbf{z}]$, given by

$$\Delta_{sz} = \sqrt{\|\mathbf{X}\|^4 + \|\mathbf{X}\|^2 \|\mathbf{Y}\|^2}$$

with $\|\mathbf{X}\| = \max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_2$ and $\|\mathbf{Y}\| = \max_{y \in \mathcal{Y}} |y|$.

Finally, we assign prior distributions for $\boldsymbol{\theta}, \sigma_y^2$ as

$$\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{m}, \mathbf{C}), \quad \sigma_y^2 \sim \mathcal{IG}(a, b). \quad (4)$$

Based on the above relations, we shall represent a hierarchical model that enables Bayesian inference of $\boldsymbol{\theta}$ given $\hat{\mathbf{S}}$ and $\hat{\mathbf{z}}$. One important element of our modelling approach is the following result that establishes the conditional distribution of \mathbf{z} given $\mathbf{S}, \boldsymbol{\theta}$, and σ_y^2 .

Proposition 3.1. *For the normal linear regression model, we have*

$$\mathbf{z} | \mathbf{S}, \boldsymbol{\theta}, \sigma_y^2 \sim \mathcal{N}(\mathbf{S}\boldsymbol{\theta}, \mathbf{S}\sigma_y^2).$$

Proof. First, note that,

$$\mathbb{E}[\mathbf{z} | \mathbf{X}, \boldsymbol{\theta}, \sigma_y^2] = \mathbb{E}[\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} + \mathbf{X}^T \mathbf{e}]$$

$$\text{Cov}(\mathbf{z} | \mathbf{X}, \boldsymbol{\theta}, \sigma_y^2) = \mathbf{X}^T \mathbf{X} \sigma_y^2$$

Hence, the conditional density of \mathbf{z} given $\mathbf{X}, \boldsymbol{\theta}$, and σ_y^2 is

$$p(\mathbf{z} | \mathbf{X}, \boldsymbol{\theta}, \sigma_y^2) = \mathcal{N}(\mathbf{z}; \mathbf{X}^T \mathbf{X} \boldsymbol{\theta}, \mathbf{X}^T \mathbf{X} \sigma_y^2). \quad (5)$$

Let ν and ω denote the probability distributions of \mathbf{x} and \mathbf{S} , respectively. By change of variables $\mathbf{S} = \mathbf{X}^T \mathbf{X}$, we have

$$\int f(\mathbf{S})\omega(d\mathbf{S}) = \int f(\mathbf{X}^T \mathbf{X})\nu(d\mathbf{X}) \quad (6)$$

for any real-valued measurable function $f : \mathcal{S}_d \mapsto \mathbb{R}$, from the set \mathcal{S}_d of $d \times d$ positive definite matrices. Also, for any pair of sets $A \subseteq \mathcal{S}_d$ and $B \in \mathbb{R}^{d \times 1}$, we can write

$$\mathbb{P}(\mathbf{S} \in A, \mathbf{z} \in B) = \int \mathbb{I}_A(\mathbf{X}^T \mathbf{X})\mathbb{P}(\mathbf{z} \in B | \mathbf{X})\nu(d\mathbf{X}).$$

The integrand above depends on \mathbf{X} through $\mathbf{X}^T \mathbf{X}$ only, since $\mathbb{P}(\mathbf{z} \in B | \mathbf{X}) = \int_B \mathcal{N}(\mathbf{z}; \mathbf{X}^T \mathbf{X} \boldsymbol{\theta}, \mathbf{X}^T \mathbf{X} \sigma_y^2) d\mathbf{z}$. Therefore, applying the relation in (6) to the RHS with the choice $f(\mathbf{S}) = \mathbb{I}_A(\mathbf{S}) \int_B \mathcal{N}(\mathbf{z}; \mathbf{S} \boldsymbol{\theta}, \mathbf{S} \sigma_y^2) d\mathbf{z}$, we can rewrite the joint probability as

$$\begin{aligned} \mathbb{P}(\mathbf{S} \in A, \mathbf{z} \in B) &= \int \mathbb{I}_A(\mathbf{S}) \int_B \mathcal{N}(\mathbf{z}; \mathbf{S} \boldsymbol{\theta}, \mathbf{S} \sigma_y^2) d\mathbf{z} \omega(d\mathbf{S}) \\ &= \int_A \int_B \mathcal{N}(\mathbf{z}; \mathbf{S} \boldsymbol{\theta}, \mathbf{S} \sigma_y^2) d\mathbf{z} \omega(d\mathbf{S}). \end{aligned}$$

This shows that the variables \mathbf{z}, \mathbf{S} have the joint probability distribution $\omega(d\mathbf{S}) \mathcal{N}(\mathbf{z}; \mathbf{S} \boldsymbol{\theta}, \mathbf{S} \sigma_y^2) d\mathbf{z}$ and that the conditional probability density is $p(\mathbf{z} | \mathbf{S}, \boldsymbol{\theta}, \sigma_y^2) = \mathcal{N}(\mathbf{z}; \mathbf{S} \boldsymbol{\theta}, \mathbf{S} \sigma_y^2)$, as claimed. \square

At this point, some important modelling differences between our work and Bernstein & Sheldon (2019) are worth discussing. In Bernstein & Sheldon (2019), the central limit theorem (CLT) is applied to $[\mathbf{S}, \mathbf{z}, \mathbf{y}^T \mathbf{y}]$, leading to a normality assumption for the whole vector. In contrast, we use the *exact* conditional distribution $p(\mathbf{z} | \mathbf{S}, \boldsymbol{\theta}, \sigma_y^2)$ thanks to Proposition 3.1. Moreover, unlike Bernstein & Sheldon (2019), we do *not* require a noisy version $\mathbf{y}^T \mathbf{y}$, hence have a slight advantage of using less privacy-preserving noise. In summary, our model has a different hierarchical structure and requires less privacy-preserving noise.

3.2. Distributed Setting

Here we extend our model to the distributed setting, where the total data are shared among $J \geq 1$ data holders as

$$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{X}_j, \mathbf{y}_j); j = 1, \dots, J\}. \quad (7)$$

We let n_i be number of rows in each \mathbf{x}_i , so that $n = n_1 + \dots + n_J$. Each data holder j shares their own summary statistics $\mathbf{S}_j = \mathbf{X}_j^T \mathbf{X}_j$, $\mathbf{z}_j = \mathbf{X}_j^T \mathbf{y}_j$ with privacy-preserving noise

$$\begin{aligned} \hat{\mathbf{S}}_j &= \mathbf{S}_j + \sigma_s \mathbf{M}_j, \\ \hat{\mathbf{z}}_j &= \mathbf{z}_j + \sigma_z^2 \mathbf{v}_j, \quad \mathbf{v}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \end{aligned} \quad (8)$$

Note that, to preserve a given (ϵ, δ) -DP overall, each party must provide that level of privacy for their data, hence σ_s^2

and σ_z^2 are the same as before. The hierarchical structure of the overall model (specified for normally distributed \mathbf{x}_i 's) is shown in Figure 1.

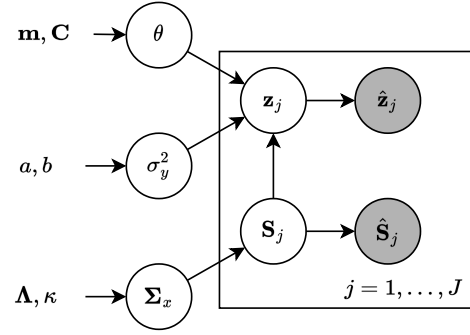


Figure 1. Differentially private distributed linear regression model (specified for normally distributed \mathbf{x}_i 's.)

The distributed setting deserves separate consideration than the single data holder case for a couple of reasons: Firstly, the node-specific observations $(\hat{\mathbf{S}}_1, \hat{\mathbf{z}}_1), \dots, (\hat{\mathbf{S}}_J, \hat{\mathbf{z}}_J)$ are altogether statistically *more informative* on θ than their aggregates $\sum_{j=1}^J \hat{\mathbf{S}}_j$ and $\sum_{j=1}^J \hat{\mathbf{z}}_j$. This is because the aggregate versions are *not* sufficient statistics of the node-specific observations $(\hat{\mathbf{S}}_1, \hat{\mathbf{z}}_1), \dots, (\hat{\mathbf{S}}_J, \hat{\mathbf{z}}_J)$ with respect to θ (even when σ_y^2 is known.) Therefore, when the node-specific observations are available, one should not, in principle, trivially aggregate them and apply an inference method designed for $J = 1$ using those aggregates.

Secondly, the partitioning of data as in (7) can be relevant to data privacy applications even *outside* the distributed learning framework, rendering the methodology in Section 4 useful in a broader sense. For example, batches of (\mathbf{x}, \mathbf{y}) -type of data may be donated to a common data collector as in (8). At this point, a particular and interesting relation exists with pan-privacy applications (Dwork et al., 2010). Imagine that sensitive data from individuals are collected sequentially in time and the data holder is concerned about possible intrusions into the memory where the sensitive data are stored. Then, one possible way to ensure the privacy of the data against such possible intrusions, which is the promise of pan-privacy, is to store the noisy statistics of every new batch of data and erase the original sensitive data. Then, at any time, the data collector has data of the form $(\hat{\mathbf{S}}_1, \hat{\mathbf{z}}_1), \dots, (\hat{\mathbf{S}}_J, \hat{\mathbf{z}}_J)$, each pair corresponding to a batch, the J th pair being the current batch. As a result, inference algorithms in Section 4 can be applied.

4. Algorithms for Bayesian Inference

Bayesian inference targets the posterior distribution of the latent variables of the model, in particular θ , given the observations $\hat{\mathbf{S}}_{1:J}$ and $\hat{\mathbf{z}}_{1:J}$. We present several Bayesian

inference algorithms for the hierarchical model described in the previous section. In addition to other concerns like computational budget, the choice among those approaches mainly depends on the specification of $P_{\mathbf{x}}$ as the distribution of \mathbf{S} directly depends on it. In this paper, we consider the following two cases and devise algorithms for each of them:

1. In some cases it may be adequate to specify $P_{\mathbf{x}} = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{x}})$. This leads to $\mathbf{S}|\Sigma_{\mathbf{x}} \sim \mathcal{W}(\Sigma_{\mathbf{x}}, n)$. Further, to account for the uncertainty about the covariance $\Sigma_{\mathbf{x}}$, one can treat it as a random variable with $\Sigma_{\mathbf{x}} \sim \mathcal{IW}(\Lambda, \kappa)$. Figure 1 shows the hierarchical structure of the distributed setting with those specifications. We defer discussing the conflict between the normality and boundedness assumptions to Remark 4.1 towards the end of Section 4.1.
2. As the second case, we assume a general (non-normal) $P_{\mathbf{x}}$. A normal approximation, based on the CLT, could be considered for the distribution \mathbf{S} (Wilson & Ghahramani, 2011). However, this would require the knowledge (or accurate estimation) of up to the fourth moments of $P_{\mathbf{x}}$ as well as expensive computations for sampling \mathbf{S} . We circumvent those difficulties by plugging in a point estimate of \mathbf{S} given $\hat{\mathbf{S}}$ and use it during the sampling process as if it is the true \mathbf{S} itself. Then, we develop two different algorithms for inference of $\boldsymbol{\theta}$, one being an MCMC algorithm and the other providing a closed form-solution for the posterior of $\boldsymbol{\theta}$ following a rough point-wise estimation of σ_y^2 . Note that these algorithms with fixed \mathbf{S} do not require a distribution for \mathbf{x} .

Next, we provide the details of our approaches and the resulting algorithms.

4.1. MCMC for Normally Distributed Features

In this section, we present an MCMC algorithm for Bayesian inference for the differentially private distributed linear regression model when $P_{\mathbf{x}} = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{x}})$ and $\Sigma_{\mathbf{x}} \sim \mathcal{IW}(\Lambda, \kappa)$. The latent variables involved in this variant are $\boldsymbol{\theta}, \Sigma_{\mathbf{x}}, \sigma_y^2, \mathbf{S}_{1:J}, \mathbf{z}_{1:J}$. Their posterior distribution given $\hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J}$ can be written as

$$p(\boldsymbol{\theta}, \sigma_y^2, \Sigma_{\mathbf{x}}, \mathbf{z}_{1:J}, \mathbf{S}_{1:J} | \hat{\mathbf{z}}_{1:J}, \hat{\mathbf{S}}_{1:J}) \propto p(\boldsymbol{\theta})p(\sigma_y^2)p(\Sigma_{\mathbf{x}}) \prod_{j=1}^J p(\mathbf{z}_j | \boldsymbol{\theta}, \sigma_y^2, \mathbf{S})p(\mathbf{S}_j | \Sigma_{\mathbf{x}})p(\hat{\mathbf{S}}_j | \mathbf{S}_j)p(\hat{\mathbf{z}}_j | \mathbf{z}_j). \quad (9)$$

One could design an MCMC algorithm for this posterior distribution that updates $\boldsymbol{\theta}, \sigma_y^2, \Sigma_{\mathbf{x}}, \mathbf{z}_{1:J}, \mathbf{S}_{1:J}$ in turn based on their full conditional distributions. However, such an algorithm suffers from poor convergence because of a high posterior correlation between $\boldsymbol{\theta}$ and $\mathbf{z}_{1:J}$ (as verified in our numerical studies). It is well known that highly correlated

variables result in poor convergence if they are updated one conditional on the other. To alleviate that problem, we work with the reduced model where $\mathbf{z}_{1:J}$ is integrated out. The reduced model has $\boldsymbol{\theta}, \Sigma_{\mathbf{x}}, \sigma_y^2$ as its latent variables, whose joint posterior distribution can be written as

$$p(\boldsymbol{\theta}, \sigma_y^2, \Sigma_{\mathbf{x}}, \mathbf{S} | \hat{\mathbf{z}}, \hat{\mathbf{S}}) \propto p(\boldsymbol{\theta})p(\sigma_y^2)p(\Sigma_{\mathbf{x}}) \prod_{j=1}^J p(\mathbf{S}_j | \Sigma_{\mathbf{x}})p(\hat{\mathbf{S}}_j | \mathbf{S}_j)p(\hat{\mathbf{z}}_j | \mathbf{S}_j, \boldsymbol{\theta}, \sigma_y^2), \quad (10)$$

where $p(\hat{\mathbf{z}} | \mathbf{S}, \boldsymbol{\theta}, \sigma_y^2) = \mathcal{N}(\hat{\mathbf{z}}; \mathbf{S}\boldsymbol{\theta}, \sigma_y^2\mathbf{S}\boldsymbol{\theta} + \sigma_z^2\mathbf{I}_d)$.

We would like to sample from the posterior distribution in (10) via MCMC that updates $\boldsymbol{\theta}, \sigma_y^2, \Sigma_{\mathbf{x}}, \mathbf{S}_{1:J}$ in turn based on their full conditional distributions. The variables $\boldsymbol{\theta}$ and $\Sigma_{\mathbf{x}}$ enjoy closed-form full conditional distributions (see Appendix A for the derivations):

$$\Sigma_{\mathbf{x}} | \mathbf{S}_{1:J}, \hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J} \sim \mathcal{IW}\left(\Lambda + \sum_{j=1}^J \mathbf{S}_j, \kappa + n\right), \quad (11)$$

$$\boldsymbol{\theta} | \sigma_y^2, \hat{\mathbf{z}}, \mathbf{S}_{1:J} \sim \mathcal{N}(\mathbf{m}_p, \Sigma_p), \quad (12)$$

where the posterior moments for $\boldsymbol{\theta}$ are

$$\Sigma_p^{-1} = \sum_{j=1}^J \mathbf{S}_j(\sigma_y^2\mathbf{S}_j + \sigma_z^2\mathbf{I}_d)^{-1}\mathbf{S}_j + \mathbf{C}^{-1},$$

$$\mathbf{m}_p = \Sigma_p \left(\sum_{j=1}^J \mathbf{S}_j(\sigma_y^2\mathbf{S}_j + \sigma_z^2\mathbf{I}_d)^{-1}\hat{\mathbf{z}}_j + \mathbf{C}^{-1}\mathbf{m} \right).$$

The full-conditional distributions of $\mathbf{S}_{1:J}$ and σ_y^2 have no closed form; hence, we design Metropolis-Hastings (MH) moves to update them. For σ_y^2 , one can simply use a random-walk MH move targeting $p(\sigma_y^2 | \boldsymbol{\theta}, \mathbf{S}_{1:J}, \hat{\mathbf{z}}_{1:J})$. For $\mathbf{S}_{1:J}$, their full conditional distribution can be factorised as

$$p(\mathbf{S}_{1:J} | \hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J}, \Sigma_{\mathbf{x}}, \sigma_y^2, \boldsymbol{\theta}) = \prod_{j=1}^J p(\mathbf{S}_j | \hat{\mathbf{S}}_j, \hat{\mathbf{z}}_j, \Sigma_{\mathbf{x}}, \sigma_y^2, \boldsymbol{\theta}),$$

where each factor is given by

$$p(\mathbf{S}_j | \hat{\mathbf{S}}_j, \hat{\mathbf{z}}_j, \Sigma_{\mathbf{x}}, \sigma_y^2, \boldsymbol{\theta}) \propto p(\hat{\mathbf{z}}_j | \mathbf{S}_j, \boldsymbol{\theta}, \sigma_y^2)p(\mathbf{S}_j | \Sigma_{\mathbf{x}})p(\hat{\mathbf{S}}_j | \mathbf{S}_j).$$

Thanks to that factorised form, each \mathbf{S}_j can be updated with an MH move independently and in parallel. For the MH algorithm to update one \mathbf{S}_j , we propose a new value from a Wishart distribution as $\mathbf{S}'_j \sim \mathcal{W}(\mathbf{S}_j/\alpha, \alpha)$, which has mean \mathbf{S}_j and variance determined by $\alpha > 0$. In our experiments, we adjust α using ideas from the adaptive

Algorithm 1 `MCMC-normalX` - one iteration

Input: Current values of $\mathbf{S}_{1:J}$, $\boldsymbol{\theta}$, σ_y^2 , $\boldsymbol{\Sigma}_x$; observations $\hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J}$; noise variances σ_s^2, σ_z^2 ; proposal parameters a, σ_q^2 ; hyperparameters $a, b, \kappa, \boldsymbol{\Lambda}, \mathbf{m}, \mathbf{C}$.

Output: New sample of $\boldsymbol{\Sigma}_x, \mathbf{S}, \sigma_y^2, \boldsymbol{\theta}$

Sample $\boldsymbol{\Sigma}_x$ using (11).

for $j = 1, 2, \dots, J$ **do**

 | Update \mathbf{S}_j via an MH move targeting $p(\mathbf{S}_j | \boldsymbol{\Sigma}_x, \boldsymbol{\theta}, \hat{\mathbf{z}}_j)$.

Sample $\boldsymbol{\theta}$ using (12).

Update σ_y^2 via an MH move targeting $p(\sigma_y^2 | \boldsymbol{\theta}, \mathbf{S}_{1:J}, \hat{\mathbf{z}}_{1:J})$.

MCMC framework (Andrieu & Thoms, 2008) to target an acceptance rate of around 0.2.

Algorithm 1 represents the overall MCMC algorithm for the hierarchical model for differentially Bayesian distributed linear regression when P_x is a normal distribution with a random covariance matrix having an inverse-Wishart distribution. We call this algorithm `MCMC-normalX`.

Remark 4.1. Admittedly, a potential concern is a conflict between the normality and boundedness assumptions (both for \mathbf{x} and y). However, we also note that the collected data often happen to have some natural boundaries (which can be exploited to determine the sensitivity of the shared statistics), and yet the normal distribution is still used for modelling and subsequent inference mainly for the sake of tractability. With the normality assumption, one can implement computationally efficient algorithms at the expense of minor modelling inaccuracies. While we acknowledge the methodologies in Alparslan & Yildirim (2022, Section 4.2) and Ju et al. (2022) that can correctly incorporate the effect of truncation into inference, we remark that those methods pay the price of exactness by having $\mathcal{O}(n)$ computational complexity per iteration.

4.2. MCMC for Non-normally Distributed Features

The normality assumption for \mathbf{x}_i 's in Section 4.1 may not be adequate for some data sets. Moreover, when d is large, updating \mathbf{S}_j 's can be the bottleneck of `MCMC-normalX` in Algorithm 1 in terms of computation time and convergence. We propose two algorithms to address both of those concerns. As it turns out, those algorithms provide accurate estimations even for the case of normally distributed features; see Section 5.1.

Our approach for non-normal \mathbf{x}_i 's is based on estimating \mathbf{S}_j 's from $\hat{\mathbf{S}}_j$'s at the beginning, using some principled estimation method, and fixing \mathbf{S}_j 's to those estimates during the whole course of the inference procedure. In that way, we obtain a faster version of `MCMC-normalX` that is also well-suited for non-normal \mathbf{x}_i 's. Indeed, we observed in our experiments that this method outperforms the other methods for most of the cases, especially when the total number of

Algorithm 2 `MCMC-fixedS` - one iteration

Input: Current values of $\boldsymbol{\theta}$, σ_y^2 ; estimates $\tilde{\mathbf{S}}_{1:J}$, observations $\hat{\mathbf{z}}_{1:J}$; noise variance σ_z^2 , and hyperparameters $a, b, \mathbf{m}, \mathbf{C}$.

Output: New sample of $\sigma_y^2, \boldsymbol{\theta}$.

Use $\mathbf{S}_{1:J} = \tilde{\mathbf{S}}_{1:J}$ throughout.

Sample $\boldsymbol{\theta}$ using (12).

Update σ_y^2 via an MH move targeting $p(\sigma_y^2 | \boldsymbol{\theta}, \mathbf{S}_{1:J}, \hat{\mathbf{z}}_{1:J})$.

nodes J increases. We call this variant `MCMC-fixedS` and present it in Algorithm 2.

As for estimating \mathbf{S}_j 's, one could simply take the privately shared $\hat{\mathbf{S}}_j$ as an estimator for \mathbf{S}_j , but $\hat{\mathbf{S}}_j$ is not necessarily a positive (semi-)definite matrix. Instead, we consider the nearest positive semi-definite matrix to $\hat{\mathbf{S}}_j$ in terms of the Frobenius norm as the estimator of \mathbf{S}_j . (The nearest *positive* definite matrix to $\hat{\mathbf{S}}_j$ does not exist.) To find the nearest positive semi-definite matrix, we follow Higham (1988) and apply the following procedure for each $j = 1, \dots, J$: (i) Calculate the eigendecomposition $\hat{\mathbf{S}}_j = \mathbf{E}\mathbf{D}\mathbf{E}^T$, where \mathbf{E} is a matrix of eigenvectors, and \mathbf{D} is a diagonal matrix consisting of the eigenvalues λ_i . (ii) The nearest symmetric positive semi-definite matrix is $\tilde{\mathbf{S}}_j = \mathbf{E}\mathbf{D}_+\mathbf{E}^T$, where \mathbf{D}_+ is a diagonal matrix with $\mathbf{D}_+(i, i) = \max\{\mathbf{D}(i, i), 0\}$.

Note that $\tilde{\mathbf{S}}_j$ found above is the maximum likelihood estimator of \mathbf{S}_j given $\hat{\mathbf{S}}_j$ (over the set of positive semi-definite matrices) since the conditional distribution of $\hat{\mathbf{S}}_j$ given \mathbf{S}_j is a normal distribution with mean \mathbf{S}_j .

Algorithm 2 is faster than Algorithm 1, since it avoids the step to update \mathbf{S}_j 's, which constitutes the main computational burden on Algorithm 1. However, Algorithm 2 can be made even faster by fixing σ_y^2 also. As a crude estimator, we used $\tilde{\sigma}_y^2 = \|Y\|/3$ throughout the experiments. We call the resulting algorithm `Bayes-fixedS-fast` and present it in Algorithm 3. Algorithm 3 does nothing but calculate the moments of the posterior distribution of $\boldsymbol{\theta}$ given $\sigma_y^2 = \tilde{\sigma}_y^2$, $\mathbf{S}_j = \tilde{\mathbf{S}}_j$ and $\hat{\mathbf{z}}_j$ for $j = 1, \dots, J$, and the prior parameters for $\boldsymbol{\theta}$.

4.3. Computational Cost

All our methods described in this section require $\mathcal{O}(d^3)$ computation (either per iteration for the iterative ones in Algorithms 1 and 2 or as a whole for the fast version in Algorithm 3) since they deal with $d \times d$ matrices. In contrast, since Bernstein & Sheldon (2019) apply CLT to the vector $[\mathbf{S}, \mathbf{z}, \mathbf{y}^T \mathbf{y}]$, their methods deal with covariance matrices of size $(d^2 + d + 1)$ *explicitly*, which leads to $\mathcal{O}(d^6)$ computation per MCMC iteration. For even moderate d , this computational difference between $\mathcal{O}(d^6)$ and $\mathcal{O}(d^3)$

Algorithm 3 `Bayes-fixedS-fast`

Input: Observations $\hat{S}_{1:J}, \hat{z}_{1:J}$; noise variance: σ_z^2 ; estimate $\tilde{\sigma}_y^2$ of σ_y^2 ; hyperparameters: \mathbf{m}, \mathbf{C} .

Output: Estimate $\hat{\theta}$.

for $j = 1, 2, \dots, J$ **do**

Calculate the estimate \tilde{S}_j for S_j using \hat{S}_j .

Calculate $\mathbf{U}_j = \tilde{S}_j(\tilde{\sigma}_y^2 \tilde{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \tilde{S}_j$.

Calculate $\mathbf{u}_j = \tilde{S}_j(\tilde{\sigma}_y^2 \tilde{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \hat{z}_j$.

return Posterior moments of θ : $\Sigma_{\text{post}}^{-1} = \sum_{j=1}^J \mathbf{U}_j + \mathbf{C}^{-1}$,

$\mathbf{m}_{\text{post}} = \Sigma_{\text{post}} \left(\mathbf{C}^{-1} \mathbf{m} + \sum_{j=1}^J \mathbf{u}_j \right)$.

becomes dramatic and the former may be prohibitive in practice.

We also note that the complexity of our methods does not depend on the data size n . This is in contrast to the $\mathcal{O}(n)$ complexity of general-purpose methods applicable to linear regression, such as Alparslan & Yildirim (2022, Section 4.3) and Ju et al. (2022).

4.4. Extensions

We mention two other variants of our methodology, deferring the details to Appendix B. One solution for dealing with non-normal P_x could be to average the feature vectors in \mathbf{X} and the corresponding response variables in \mathbf{y} , so that the averaged rows of \mathbf{X} can be modelled as approximately normal, due to CLT. This enables using the methods devised for normally distributed features. For the details of this approach, see Appendix B.1.

Secondly, if the features are normally distributed but not centred, we need to include the intercept parameter, which corresponds to appending \mathbf{x}_i with a one from the left, and MCMC-normalX does not directly apply. In that case, we can modify the hierarchical model that accommodates the non-centralised features and the intercept parameter, and we still benefit from the sampling techniques involved in MCMC-normalX in Algorithm 1. Appendix B.2 contains the details of the modified hierarchical model.

5. Numerical Experiments

We present several numerical evaluations of the proposed methods, MCMC-normalX, MCMC-fixedS, and Bayes-fixedS-fast, with simulated and real data. We compare our algorithms with two methods: adaSSP of Wang (2018) and the MCMC method of Bernstein & Sheldon (2019) for differentially private linear regression, which we will call MCMC-B&S. Note that adaSSP and MCMC-B&S were originally proposed for the non-distributed setting, that is, $J = 1$. For a comprehensive

comparison, we implemented their extensions for $J \geq 1$. The details of those extensions are provided in Appendix C. In particular, we carefully generalised the model in Bernstein & Sheldon (2019) for $J \geq 1$ (and for (ϵ, δ) -DP) in a similar fashion as we did to our model in Section 3.2. MCMC-B&S is the adaptation of Bernstein & Sheldon (2019, Algorithm 1) for this generalised model. Both for simulated and real data, we set $\|\mathbf{X}\|$ and $\|\mathbf{Y}\|$ to the maximum of the norms over the whole dataset (for real data, we perform centring and normalising first.) This procedure is equivalent to scaling the data to an interval and standard in existing work. The link for the code and data to replicate all of the experiments in this section is given at the end of Section 6.

5.1. Experiments with Simulated Data

We considered two different configurations, $(n = 10^5, d = 2)$ and $(n = 10^5, d = 5)$, for the problem size. For each (n, d) we simulated the data as follows. We generated $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma_x)$ where $\Sigma_x \sim \mathcal{IW}(\Lambda, \kappa)$, with $\kappa = d + 1$ and the scale matrix selected randomly as $\Lambda = \mathbf{V}^T \mathbf{V}$, where \mathbf{V} is a $d \times d$ matrix of i.i.d. variables from $\mathcal{N}(0, 1)$. The response variables \mathbf{y} are generated with $\sigma_y^2 = 1$. For inference, we used the same Λ, κ as above and $a = 20, b = 0.5, \mathbf{m} = \mathbf{0}_{d \times 1}, \mathbf{C} = b/(a - 1)\mathbf{I}_d$.

We evaluated the methods at all combinations of $J \in \{1, 5, 10\}$ and $\epsilon \in \{0.1, 0.2, 0.5, 1, 2, 5, 10\}$. All the MCMC algorithms were run for 10^4 iterations. For each (J, ϵ) pair, we ran each method for 50 times (each with different noisy observations) to obtain average performances.

For performance metrics, we looked at the mean squared errors (MSE) of (i) the estimates $\hat{\theta}$ and (ii) the predictions $\hat{y}(\mathbf{x}_{\text{test}})$ generated by the methods. For the Bayesian methods, $\hat{\theta}$ is taken as the posterior mean, which can be numerically estimated for the MCMC algorithms. For prediction performance, we calculated $\mathbb{E}[\hat{y}(\mathbf{x}_{\text{test}}) - y_{\text{test}}]^2$. For the Bayesian methods, $\hat{y}(\mathbf{x}_{\text{test}})$ is the posterior predictive expectation of y_{test} at \mathbf{x}_{test} . For adaSSP, we simply take $\hat{y}(\mathbf{x}_{\text{test}}) = \mathbf{x}_{\text{test}}^T \hat{\theta}$.

The results are summarised in Figure 2. We observe that MCMC-fixedS and Bayes-fixedS-fast outperform adaSSP and MCMC-B&S in almost all cases both in terms of estimation and prediction. Comparing the full-scale algorithms MCMC-normalX and MCMC-B&S (that involve updates of \mathbf{S}), we observe a clear advantage of MCMC-normalX at $d = 2$ but MCMC-B&S becomes more competitive at $d = 5$. This can be attributed to the fact that MCMC-B&S requires the extra statistic $\mathbf{y}^T \mathbf{y}$, unlike MCMC-normalX, which causes MCMC-B&S to use more noisy statistics. This difference becomes more significant at small d , where the relative effect of the presence of $\mathbf{y}^T \mathbf{y}$ on the sensitivity is more significant. Finally, all methods improve as ϵ grows, as expected.

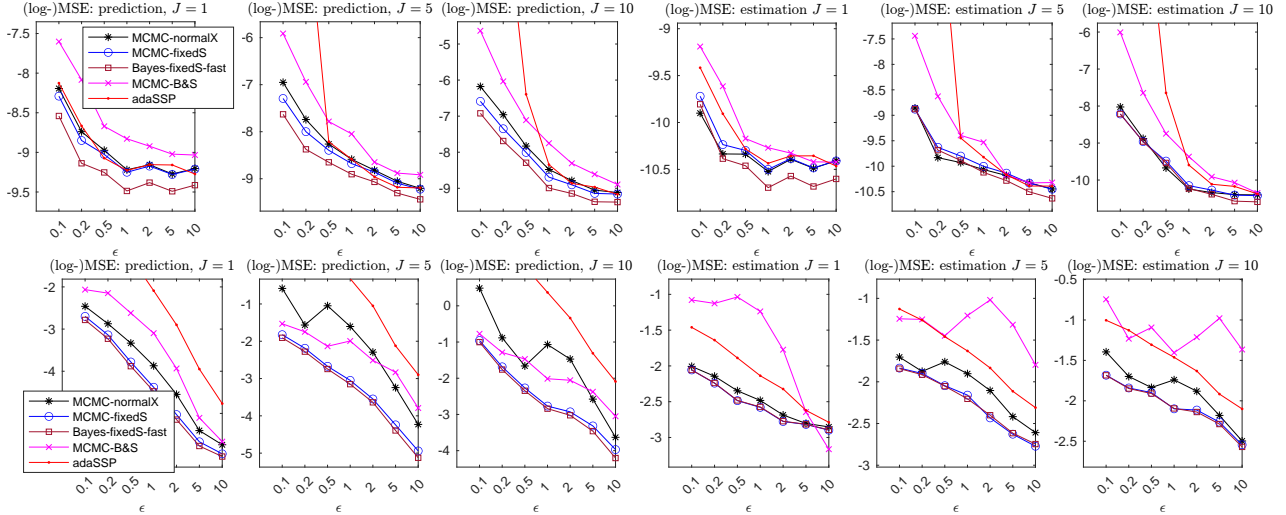


Figure 2. Averaged prediction and estimation performances (over 50 runs). Top row: $n = 10^5$, $d = 2$, Bottom row: $n = 10^5$, $d = 5$. For the non-private posterior, for $d = 2$ ($d = 5$, resp.), the obtained MSE values for estimation are 2.92×10^{-5} (1.54×10^{-5} , resp.) and for prediction are 9.60×10^{-5} (4.19×10^{-5} , resp.)

Table 1. Averaged prediction performances (over 50 runs) for the real datasets - $\epsilon = 1$

J	data sets	MCMC-normalX	MCMC-fixedS	Bayes-fixedS-fast	MCMC-B&S	adaSSP
$J = 1$	PowerPlant	0.0129	0.0129	0.0129	0.0128	0.0139
	BikeSharing	0.0024	0.0021	0.0021	0.0020	0.0107
	AirQuality	0.0060	0.0057	0.0057	0.0062	0.0066
	3droad	0.0229	0.0229	0.0229	0.0229	0.0229
$J = 5$	PowerPlant	0.0133	0.0134	0.0134	0.0136	0.0235
	BikeSharing	0.0174	0.0045	0.0045	0.0086	0.0382
	AirQuality	0.0142	0.0100	0.0099	0.0130	0.0227
	3droad	0.0229	0.0229	0.0229	0.0229	0.0229
$J = 10$	PowerPlant	0.0142	0.0143	0.0143	0.0143	0.0351
	BikeSharing	0.0812	0.0082	0.0082	0.0137	0.0526
	AirQuality	0.0985	0.0117	0.0117	0.0216	0.0314
	3droad	0.0229	0.0229	0.0229	0.0229	0.0229

We also compare the computation times of the MCMC algorithms MCMC-normalX, MCMC-fixedS, and MCMC-B&S¹. Figure 3 shows the run times of the algorithms vs. d . The drastic difference in computational loads explained in Section 4.3 is also visible in the figure. While MCMC-B&S may be improved in terms of accuracy as d increases, the $\mathcal{O}(d^6)$ dramatically slows it down.

In addition to the MSE, we also consider maximum mean discrepancy (MMD) for evaluating the calibration of the learned posteriors following the method in (Gretton et al., 2012). The estimated MMD values in Figure 4 in Appendix D show that all the methods converge to the non-private posterior as ϵ increases.

¹The algorithms were run in MATLAB 2021b on an Apple M1 chip with 8 cores and 16 GB LPDDR4 memory.

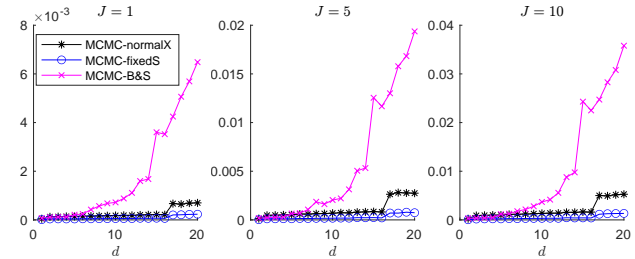


Figure 3. Run times per iteration for MCMC algorithms

5.2. Experiments with Real Data

For the real data case, we use four different data sets from the UCI Machine Learning Repository. We disregard the columns including string data or key values (ID, name, date, etc.), and we consider the most right-hand col-

umn as y . The finalised data sets are summarised below.

data set	n	d	hyperlinks
power plant energy	7655	4	view link
bike sharing	13904	14	view link
air quality	7486	12	view link
3d road	347900	3	view link

For prediction, we took 80% of the data for training and the rest for testing. We present the average prediction performances (out of 50 runs) in Table 1 for each dataset and J with $\epsilon = 1$. We observe that the prediction performances of the compared methods are close, while `MCMC-fixed-S` and `Bayes-fixed-S` are arguably the most stable ones. When $J > 1$, (the distributed data setting) those two methods beat `adaSSP` and `MCMC-B&S` more satisfactorily. To show the robustness of these conclusions to variation, we provide 90% confidence intervals for the prediction performances in Table 2 in Appendix D.

6. Conclusion

We propose a novel Bayesian inference framework, with MCMC being its main workhorse, for a differentially private distributed linear regression setting where the data is partitioned among the data holders and shared using summary statistics. We provide several Bayesian inference algorithms suited to the developed hierarchical model for linear regression. Those algorithms can be preferred one over the other depending on the computational budget, model specifics, or how much we know about the underlying statistical facts of the data. We exploit the conditional structure between the summary statistics of linear regression (Proposition 3.1), which leads to feasible algorithms with computational advantages over their competitors. The numerical experiments show that the proposed methods are competitive with their state-of-the-art alternatives in terms of accuracy. The extensions mentioned in Section 4.4 indicate potential future directions.

The experiments demonstrate that `MCMCM-fixedS` and `Bayes-fixedS-fast` are competitive even under normality. Hence, we can suggest users try those fast and effective versions on their first attempts. However, `MCMC-normalX` can potentially provide more insight since it infers P_x as well. There is also room for improvement of `MCMC-normalX`. We chose the most common MH moves to update σ_y^2 and S_j 's, without paying much attention to their efficiencies. Especially for large d , more advanced techniques such as Hamiltonian Monte Carlo or pseudo-marginal MCMC may be employed to facilitate the mixing of the algorithm.

Code for the experiments: Link to the code and the data for the experiments: https://github.com/sinanyildirim/Bayesian_DP_dist_LR.git.

Acknowledgements

The study was funded by the Scientific and Technological Research Council of Turkey (TÜBİTAK) ARDEB Grant No 120E534.

References

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, pp. 308–318, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4139-4. doi: 10.1145/2976749.2978318. URL <http://doi.acm.org/10.1145/2976749.2978318>.
- Alabi, D., McMillan, A., Sarathy, J., Smith, A., and Vadhan, S. Differentially private simple linear regression. *Proceedings on Privacy Enhancing Technologies*, 2022: 184–204, 04 2022. doi: 10.2478/popets-2022-0041.
- Alparslan, B. and Yıldırım, S. Statistic selection and MCMC for differentially private Bayesian estimation. *Statistics and Computing*, 32(5):66, 2022. doi: 10.1007/s11222-022-10129-8. URL <https://doi.org/10.1007/s11222-022-10129-8>.
- Andrieu, C. and Thoms, J. A tutorial on adaptive mcmc. *Statistics and Computing*, 18(4):343–373, 2008. doi: 10.1007/s11222-008-9110-y. URL <https://doi.org/10.1007/s11222-008-9110-y>.
- Balle, B. and Wang, Y.-X. Improving the Gaussian mechanism for differential privacy: Analytical calibration and optimal denoising. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 394–403. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/balle18a.html>.
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 464–473, Oct 2014. doi: 10.1109/FOCS.2014.56.
- Bernstein, G. and Sheldon, D. R. Differentially private bayesian linear regression. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/f90f2aca5c640289d0a29417bcb63a37-Paper.pdf>.

- Bun, M. and Steinke, T. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Proceedings, Part I, of the 14th International Conference on Theory of Cryptography - Volume 9985*, pp. 635–658, New York, NY, USA, 2016. Springer-Verlag New York, Inc. ISBN 978-3-662-53640-7. doi: 10.1007/978-3-662-53641-4_24. URL https://doi.org/10.1007/978-3-662-53641-4_24.
- Cai, T. T., Wang, Y., and Zhang, L. The cost of privacy: Optimal rates of convergence for parameter estimation with differential privacy. *The Annals of Statistics*, 49(5): 2825 – 2850, 2021. doi: 10.1214/21-AOS2058. URL <https://doi.org/10.1214/21-AOS2058>.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. Differentially private empirical risk minimization. *J. Mach. Learn. Res.*, 12:1069–1109, jul 2011. ISSN 1532-4435.
- Dankar, F. K. and El Emam, K. Practicing differential privacy in health care: A review. *Trans. Data Priv.*, 6(1): 35–67, 2013.
- Dimitrakakis, C., Nelson, B., Zhang, Z., Mitrokotsa, A., and Rubinstein, B. I. Differential privacy for bayesian inference through posterior sampling. *Journal of machine learning research*, 18(11):1–39, 2017.
- Dong, J., Roth, A., and Su, W. J. Gaussian differential privacy. *Journal of the Royal Statistical Society Series B*, 84(1):3–37, February 2022. doi: 10.1111/rssb.12454. URL <https://ideas.repec.org/a/bla/jorssb/v84y2022ilp3-37.html>.
- Dwork, C. Differential privacy. In Bugliesi, M., Preneel, B., Sassone, V., and Wegener, I. (eds.), *Automata, Languages and Programming*, pp. 1–12, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. ISBN 978-3-540-35908-1.
- Dwork, C. Differential privacy: A survey of results. In Agrawal, M., Du, D., Duan, Z., and Li, A. (eds.), *Theory and Applications of Models of Computation*, pp. 1–19, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg. ISBN 978-3-540-79228-4.
- Dwork, C. and Roth, A. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3–4):211–407, aug 2014. ISSN 1551-305X. doi: 10.1561/04000000042. URL <https://doi.org/10.1561/04000000042>.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography*, pp. 265–284. Springer, 2006.
- Dwork, C., Naor, M., Pitassi, T., Rothblum, G. N., and Yekhanin, S. Pan-private streaming algorithms. In *ICS*, pp. 66–80, 2010.
- Dwork, C., Roth, A., et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014a.
- Dwork, C., Talwar, K., Thakurta, A., and Zhang, L. Analyze gauss: Optimal bounds for privacy-preserving principal component analysis. In *Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing*, STOC ’14, pp. 11–20, New York, NY, USA, 2014b. Association for Computing Machinery. ISBN 9781450327107. doi: 10.1145/2591796.2591883. URL <https://doi.org/10.1145/2591796.2591883>.
- Ferrando, C., Wang, S., and Sheldon, D. Parametric bootstrap for differentially private confidence intervals. In Camps-Valls, G., Ruiz, F. J. R., and Valera, I. (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 1598–1618. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/ferrando22a.html>.
- Foulds, J., Geumlek, J., and an Kamalika Chaudhuri, M. W. On the theory and practice of privacy-preserving Bayesian data analysis. Technical report, arxiv:1603.07294, mar 2016.
- Gong, R. Exact inference with approximate computation for differentially private data via perturbations. *Journal of Privacy and Confidentiality*, 12(2), 2022. doi: 10.29012/jpc.797. URL <https://par.nsf.gov/biblio/10387759>.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773, 2012. URL <http://jmlr.org/papers/v13/gretton12a.html>.
- Heikkilä, M., Lagerspetz, E., Kaski, S., Shimizu, K., Tarkoma, S., and Honkela, A. Differentially private Bayesian learning on distributed data. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/dfce06801e1a85d6d06f1fdd4475dacd-Paper.pdf>.
- Heikkilä, M. A., Jälkö, J., Dikmen, O., and Honkela, A. Differentially private Markov chain Monte Carlo. In *NeurIPS*, 2019.
- Higham, N. J. Computing a nearest symmetric positive semidefinite matrix. *Linear Algebra and its*

- Applications*, 103:103–118, 1988. ISSN 0024-3795. doi: [https://doi.org/10.1016/0024-3795\(88\)90223-6](https://doi.org/10.1016/0024-3795(88)90223-6). URL <https://www.sciencedirect.com/science/article/pii/0024379588902236>.
- Ju, N., Awan, J., Gong, R., and Rao, V. Data augmentation MCMC for bayesian inference from privatized data. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=tTWCQrgjuM>.
- Kuru, N., Birbil, S. I., Gürbüzbalaban, M., and Yıldırım, S. Differentially private accelerated optimization algorithms. *SIAM Journal on Optimization*, 32(2):795–821, 2022. doi: 10.1137/20M1355847. URL <https://doi.org/10.1137/20M1355847>.
- Wang, Y.-X. Revisiting differentially private linear regression: optimal and adaptive prediction & estimation in unbounded domain. In *UAI*, 2018.
- Wang, Y.-X., Fienberg, S., and Smola, A. Privacy for free: Posterior sampling and stochastic gradient Monte Carlo. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2493–2502, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/wangg15.html>.
- Williams, O. and Mcsherry, F. Probabilistic inference and differential privacy. In Lafferty, J., Williams, C., Shawe-Taylor, J., Zemel, R., and Culotta, A. (eds.), *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010. URL <https://proceedings.neurips.cc/paper/2010/file/fb60d411a5c5b72b2e7d3527cfc84fd0-Paper.pdf>.
- Wilson, A. G. and Ghahramani, Z. Generalised Wishart processes. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, UAI’11, pp. 736–744, Arlington, Virginia, USA, 2011. AUAI Press. ISBN 9780974903972.
- Yıldırım, S. and Ermiş, B. Exact MCMC with differentially private moves. *Statistics and Computing*, 29(5):947–963, 2019. doi: 10.1007/s11222-018-9847-x. URL <https://doi.org/10.1007/s11222-018-9847-x>.
- Zhang, J., Zhang, Z., Xiao, X., Yang, Y., and Winslett, M. Functional mechanism: Regression analysis under differential privacy. *Proc. VLDB Endow.*, 5(11):1364–1375, jul 2012. ISSN 2150-8097. doi: 10.14778/2350229.2350253. URL <https://doi.org/10.14778/2350229.2350253>.
- Zhang, Z., Rubinstein, B., and Dimitrakakis, C. On the differential privacy of Bayesian inference. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1), Mar. 2016. doi: 10.1609/aaai.v30i1.10254. URL <https://ojs.aaai.org/index.php/AAAI/article/view/10254>.

A. Derivations for MCMC-normalX

We reserve this section for the derivations for our algorithm MCMC-normalX.

A.1. Full Conditional Distribution of Σ_x

We can write the full conditional distribution for Σ_x as

$$\begin{aligned}
 p(\Sigma_x | \mathbf{S}_{1:J}, \hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J}) &\propto p(\Sigma_x) \prod_{j=1}^J p(\mathbf{S}_j | \Sigma_x) \\
 &= \frac{|\Lambda|^{dk/2}}{2^{dk/2} \Gamma_d(\frac{k}{2})} |\Sigma_x|^{-(d+\kappa+1)/2} e^{-\frac{1}{2} \text{tr}(\Lambda \Sigma_x^{-1})} \prod_{j=1}^J \frac{|\mathbf{S}_j|^{(n_j-d-1)/2} e^{-\frac{1}{2} \text{tr}(\Sigma_x^{-1} \mathbf{S}_j)}}{2^{n_j d/2} |\Sigma_x|^{n_j/2} \Gamma_d(n_j/2)} \\
 &\propto |\Sigma_x|^{-\frac{n}{2} - \frac{(d+\kappa+1)}{2}} e^{-\frac{1}{2} (\sum \text{tr}(\Sigma_x^{-1} \mathbf{S}_j) + \text{tr}(\Lambda \Sigma_x^{-1}))} \\
 &\propto |\Sigma_x|^{-\frac{(d+\kappa+n+1)}{2}} e^{-\frac{1}{2} \text{tr}((\sum \mathbf{S}_j + \Lambda) \Sigma_x^{-1})}.
 \end{aligned}$$

Therefore, we have

$$\Sigma_x | \mathbf{S}_{1:J}, \hat{\mathbf{S}}_{1:J}, \hat{\mathbf{z}}_{1:J} \sim \mathcal{IW} \left(\Lambda + \sum_{j=1}^J \mathbf{S}_j, \kappa + n \right).$$

A.2. Full Conditional Distribution of θ

The full conditional distribution for θ can be written as

$$p(\theta | \mathbf{S}_{1:J}, \sigma_y^2, \hat{\mathbf{z}}_{1:J}) \propto \mathcal{N}(\theta; \mathbf{m}, \mathbf{C}) p(\hat{\mathbf{z}}_{1:J} | \mathbf{S}_{1:J}, \theta, \Sigma_y^2).$$

For the second factor, we have

$$\begin{aligned}
 p(\hat{\mathbf{z}}_{1:J} | \mathbf{S}_{1:J}, \theta, \sigma_y^2) &\propto \prod_{i=1}^J p(\hat{z}_j | \mathbf{S}_j, \theta, \sigma_y^2) = \prod_{i=1}^J p(\hat{z}_j; \mathbf{S}_j \theta, \sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d) \\
 &\propto \prod_{i=1}^J \exp \left\{ -\frac{1}{2} (\hat{z}_j - \mathbf{S}_j \theta)^T (\sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} (\hat{z}_j - \mathbf{S}_j \theta) \right\} \\
 &\propto \exp \left\{ -\frac{1}{2} \left[\theta^T \left(\sum_j \mathbf{S}_j (\sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \mathbf{S}_j \right) \theta - 2\theta^T \left(\sum_j \mathbf{S}_j (\sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \right) \hat{\mathbf{z}}_j \right] \right\}.
 \end{aligned}$$

Reorganising the terms, we end up with

$$p(\theta | \mathbf{S}_{1:J}, \sigma_y^2, \hat{\mathbf{z}}_{1:J}) \propto \exp \left\{ -\frac{1}{2} [\theta^T \Sigma_p^{-1} \theta - 2\theta^T \Sigma_p^{-1} \mathbf{m}_p] \right\},$$

where $\Sigma_p = [\sum_j \mathbf{S}_j (\sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \mathbf{S}_j + \mathbf{C}^{-1}]^{-1}$ and $\mathbf{m}_p = \Sigma_p [\sum_j \mathbf{S}_j (\sigma_y^2 \mathbf{S}_j + \sigma_z^2 \mathbf{I}_d)^{-1} \hat{\mathbf{z}}_j + \mathbf{C}^{-1} \mathbf{m}]$. Therefore, $\theta | \mathbf{S}_{1:J}, \sigma_y^2, \hat{\mathbf{z}}_{1:J} \sim \mathcal{N}(\theta; \mathbf{m}_p, \Sigma_p)$.

A.3. Acceptance Ratio for the MH Update of \mathbf{S}_j

We drop the index j from \mathbf{S}_j for simplicity. When $\mathbf{S}' \sim \mathcal{W}(\mathbf{S}/\alpha, \alpha)$, the proposal density is

$$q(\mathbf{S}' | \mathbf{S}) = \frac{|\mathbf{S}'|^{(\alpha-d-1)/2} e^{-\text{tr}[\alpha \mathbf{S}^{-1} \mathbf{S}']/2}}{|\mathbf{S}/\alpha|^{\alpha/2} 2^{\alpha d/2} \Gamma_d(\frac{\alpha}{2})} = \frac{|\mathbf{S}'|^{(\alpha-d-1)/2} e^{-\text{tr}[\alpha \mathbf{S}^{-1} \mathbf{S}']/2}}{|\mathbf{S}|^{\alpha/2} 2^{\alpha d/2} \Gamma_d(\frac{\alpha}{2})} \alpha^{\alpha/2}.$$

Therefore, the acceptance ratio corresponding to this proposal is

$$\min \left\{ 1, \frac{q(\mathbf{S}|\mathbf{S}')}{q(\mathbf{S}'|\mathbf{S})} \frac{\mathcal{W}(\mathbf{S}'; n_j \Sigma_x, \kappa) p(\hat{\mathbf{S}}|\hat{\mathbf{S}}') \mathcal{N}(\hat{\mathbf{z}}; \mathbf{S}'\boldsymbol{\theta}, \sigma_y^2 \mathbf{S}'\boldsymbol{\theta} + \sigma_z^2 \mathbf{I}_d)}{\mathcal{W}(\mathbf{S}; n_j \Sigma_x, \kappa) p(\hat{\mathbf{S}}|\hat{\mathbf{S}}) \mathcal{N}(\hat{\mathbf{z}}; \mathbf{S}\boldsymbol{\theta}, \sigma_y^2 \mathbf{S}\boldsymbol{\theta} + \sigma_z^2 \mathbf{I}_d)} \right\},$$

where the ratio of proposals becomes

$$\frac{q(\mathbf{S}|\mathbf{S}')}{q(\mathbf{S}'|\mathbf{S})} = \frac{|\mathbf{S}|^{(\alpha-d-1)/2} |\mathbf{S}'|^{\alpha/2} e^{-\text{tr}[\alpha \mathbf{S}'^{-1} \mathbf{S}]/2}}{|\mathbf{S}'|^{(\alpha-d-1)/2} |\mathbf{S}|^{\alpha/2} e^{-\text{tr}[\alpha \mathbf{S}^{-1} \mathbf{S}']/2}} = \left(\frac{|\mathbf{S}|}{|\mathbf{S}'|} \right)^{\alpha-(d+1)/2} e^{\alpha(\text{tr}[\mathbf{S}^{-1} \mathbf{S}'] - \text{tr}[\mathbf{S}'^{-1} \mathbf{S}])/2}.$$

A.4. Acceptance Ratio for the MH Update of σ_y^2

To update σ_y^2 , we use a random walk proposal $\sigma_y^{2'} \sim \mathcal{N}(\sigma_y^2, \sigma_q^2)$. The resulting acceptance ratio is

$$\min \left\{ 1, \frac{\text{IG}(\sigma_y^{2'}; a, b) \prod_{j=1}^J \mathcal{N}(\hat{\mathbf{z}}_j; \mathbf{S}_j \boldsymbol{\theta}, \sigma_y^{2'} \mathbf{S}_j \boldsymbol{\theta} + \sigma_z^2 \mathbf{I}_d)}{\text{IG}(\sigma_y^2; a, b) \prod_{j=1}^J \mathcal{N}(\hat{\mathbf{z}}_j; \mathbf{S}_j \boldsymbol{\theta}, \sigma_y^2 \mathbf{S}_j \boldsymbol{\theta} + \sigma_z^2 \mathbf{I}_d)} \right\}.$$

B. Other Variants

This appendix is reserved for the details of the other variants mentioned in Section 4.4.

B.1. Approximating Normality by Averaging

Assume $J = 1$ for simplicity. When the feature vectors $\mathbf{x}_i, i = 1, \dots, n$ are not normal, another approach that we consider is based on modifying the data to make the rows of the feature matrix averages of multiple original features, and therefore, approximately normal, by the CLT. Specifically, let $k > 1$ be an integer that divides n so that $m = n/k$ is also an integer. Consider the $m \times n$ matrix

$$\mathbf{A} = \frac{1}{\sqrt{k}} \begin{bmatrix} \mathbf{1}_{1 \times k} & \mathbf{0}_{1 \times k} & \dots & \mathbf{0}_{1 \times k} \\ \mathbf{0}_{1 \times k} & \mathbf{1}_{1 \times k} & \dots & \mathbf{0}_{1 \times k} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{1 \times k} & \mathbf{0}_{1 \times k} & \dots & \mathbf{1}_{1 \times k} \end{bmatrix}_{m \times n},$$

where k is some number that divides n so that $m = n/k$ is an integer. Then, the matrix $\mathbf{X}_{\text{av}} = \mathbf{A}\mathbf{X}$ corresponds to constructing a shorter $m \times d$ matrix whose i 'th column is the average of the rows $(i-1)k + 1, \dots, ik$ of \mathbf{X} (scaled by $1/\sqrt{k}$). When k is large enough, we can make normality assumptions for the rows of \mathbf{X}_{av} . Further, consider

$$\mathbf{y}_{\text{av}} := \mathbf{A}\mathbf{y} = \mathbf{X}_{\text{av}}\boldsymbol{\theta} + \mathbf{A}\mathbf{e},$$

whose mean is $\mathbf{X}_{\text{av}}\boldsymbol{\theta}$ and covariance $\mathbf{A}\mathbf{A}^T \sigma_y^2$. But we have $\mathbf{A}\mathbf{A}^T = \mathbf{I}_m$, so the covariance is $\sigma_y^2 \mathbf{I}_m$. Therefore, the same hierarchical model in Figure 1 can be used for $\mathbf{X}_{\text{av}}, \mathbf{y}_{\text{av}}$ with their respective summary statistics

$$\mathbf{z}_{\text{av}} = (\mathbf{X}_{\text{av}})^T \mathbf{y}_{\text{av}}, \quad \mathbf{S}_{\text{av}} = (\mathbf{X}_{\text{av}})^T \mathbf{X}_{\text{av}},$$

as well as the noisy versions of those summary statistics to provide a given level of privacy. Note that \mathbf{S}_{av} and \mathbf{z}_{av} have the same sensitivities as \mathbf{S} and \mathbf{z} , hence the same noise variances are needed for privacy. However, \mathbf{S}_{av} and \mathbf{z}_{av} bear less information about $\boldsymbol{\theta}$ than \mathbf{S} and \mathbf{z} due to averaging.

B.2. Including the Intercept

Again, assume $J = 1$ for simplicity. If we include the intercept parameter, which corresponds to appending \mathbf{x}_i with a 1 from the left, the design matrix will be changed from \mathbf{S} to $\mathbf{S}_0 = \begin{bmatrix} n & n\bar{\mathbf{x}}^T \\ n\bar{\mathbf{x}} & \mathbf{S} \end{bmatrix}$, where $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. Also, note that $\mathbf{S} = (n-1)\widehat{\Sigma}_x + n\bar{\mathbf{x}}\bar{\mathbf{x}}^T$ where $\widehat{\Sigma}_x$ is the sample covariance of $\mathbf{x}_1, \dots, \mathbf{x}_n$. Under the normality assumption for \mathbf{x}_i 's, we have $\bar{\mathbf{x}} \sim \mathcal{N}(\mathbf{m}, \Sigma_x/n)$ and $\widehat{\Sigma}_x \sim \mathcal{W}(n-1, \Sigma_x)$ are independent and have known distributions. Therefore, we can

write a model that includes the additional variables ($\mathbf{b} = \bar{\mathbf{x}}, \widehat{\Sigma}_x, \mathbf{S}_0$) such that \mathbf{b} and $\widehat{\Sigma}_x$ are independent and have known distributions and

$$\mathbf{S}_0 = \begin{bmatrix} n & n\mathbf{b}^T \\ n\mathbf{b} & (n-1)\widehat{\Sigma} + n\mathbf{b}\mathbf{b}^T \end{bmatrix}$$

replaces \mathbf{S} in the standard model. More specifically, we have the following hierarchical model:

$$\begin{aligned} \boldsymbol{\theta} &\sim \mathcal{N}(\mathbf{m}, \mathbf{C}), \quad \Sigma_x \sim \mathcal{IW}(\boldsymbol{\Lambda}, \kappa), \quad \widehat{\Sigma}_x | \Sigma_x \sim \mathcal{W}(n-1, \Sigma_x), \quad \mathbf{b} | \Sigma_x \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma_x/n), \\ \mathbf{z} | \boldsymbol{\theta}, \sigma_y^2, \widehat{\Sigma}_x, \mathbf{b} &\sim \mathcal{N}(\mathbf{S}_0 \boldsymbol{\theta}, \mathbf{S}_0 \sigma_y^2), \quad \widehat{\mathbf{S}} | \widehat{\Sigma}_x, \mathbf{b} \sim \mathcal{N}(\mathbf{S}_0, \sigma_{s,\epsilon}^2 \mathbf{I}_{d+1}), \quad \hat{\mathbf{z}} | \mathbf{z} \sim \mathcal{N}(\mathbf{z}, \sigma_z^2 \mathbf{I}_{d+1}). \end{aligned}$$

C. Compared Methods

Here, we provide the details of the methods which we compare with the proposed methods in this paper. Those methods are originally proposed for $J = 1$. However, for comparison, we implemented their natural extensions to the general (distributed) case $J \geq 1$. The implementations of those methods can be found in the code package provided for this paper.

C.1. MCMC of Bernstein & Sheldon (2019) Adapted to the Distributed Setting

In Bernstein & Sheldon (2019), $J = 1$ is considered only and the vector $\mathbf{s}\mathbf{s} = [\text{vec}(\mathbf{S}), \mathbf{z} = \mathbf{X}^T \mathbf{y}, u = \mathbf{y}^T \mathbf{y}]$ is perturbed with privacy-preserving noise to generate the observations of the model. For $J \geq 1$, we consider the following natural extension for generating perturbed observations $\widehat{\mathbf{s}}\mathbf{s} = [\text{vec}(\widehat{\mathbf{S}}_j), \hat{\mathbf{z}}_j, \hat{u}_j]$, where

$$\widehat{\mathbf{S}}_j = \mathbf{S}_j + \sigma_{dp} \mathbf{M}_j, \quad \hat{\mathbf{z}}_j = \mathbf{z}_j + \mathbf{v}_j, \quad \mathbf{v}_j \sim \mathcal{N}(\mathbf{0}, \sigma_{dp}^2 \mathbf{I}_d), \quad \hat{u}_j = u_j + w_j, \quad w_j \sim \mathcal{N}(0, \sigma_{dp}^2), \quad (13)$$

where $\sigma_{dp} = \sigma(\epsilon, \delta) \Delta_{\mathbf{s}\mathbf{s}}$ with $\Delta_{\mathbf{s}\mathbf{s}} = \sqrt{\|\mathbf{X}\|^4 + \|\mathbf{X}\|^2 \|\mathbf{Y}\|^2 + \|\mathbf{Y}\|^4}$.

For completeness, we provide the further specifics of the model: We take $(\boldsymbol{\theta}, \sigma_y^2) \sim \mathcal{NIG}(a_0, b_0, \mathbf{m}, \mathbf{C})$ and $P_{\mathbf{x}} = \mathcal{N}(\mathbf{0}, \Sigma_x)$, where $\Sigma_x \sim \mathcal{IW}(\boldsymbol{\Lambda}, \kappa)$.

During the comparisons, we set the $a_0, b_0, \mathbf{m}, \mathbf{C}, \boldsymbol{\Lambda}, \kappa$ to the same values for both this model and our proposed model that assumes normally distributed features, *i.e.*, $P_{\mathbf{x}} = \mathcal{N}(\mathbf{0}, \Sigma_x)$. Then, we apply an extension of the method of Bernstein & Sheldon (2019) suited to those observations. One iteration of that algorithm includes the following steps in order:

- Calculate the $D \times 1$ mean vector and $D \times D$ covariance matrix

$$\boldsymbol{\mu}_{\text{prior}}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}}[\mathbf{s}\mathbf{s}], \quad \Sigma_{\text{prior}}(\boldsymbol{\theta}) = \text{Cov}_{\boldsymbol{\theta}}[\mathbf{s}\mathbf{s}]$$

This step requires the fourth moments $\mathcal{N}(\mathbf{0}, \Sigma_x)$.

- For $j = 1, \dots, J$, sample $\mathbf{s}\mathbf{s}_j \sim \mathcal{N}(\boldsymbol{\mu}_{\text{post}}^{(j)}, \Sigma_{\text{post}}^{(j)})$ with

$$\Sigma_{\text{post}}^{(j)} = ([n_j \Sigma_{\text{prior}}(\boldsymbol{\theta})]^{-1} + (1/\sigma_{dp}^2) \mathbf{I}_D)^{-1}, \quad \text{and} \quad \boldsymbol{\mu}_{\text{post}}^{(j)} = \Sigma_{\text{post}}^{(j)} (\Sigma_{\text{prior}}(\boldsymbol{\theta})^{-1} \boldsymbol{\mu}_{\text{prior}}(\boldsymbol{\theta}) + \widehat{\mathbf{s}}\mathbf{s}_j / \sigma_{dp}^2).$$

- Sample $\Sigma_x \sim \mathcal{IW}(\sum_{j=1}^J \mathbf{S}_j + \boldsymbol{\Lambda}, n + \kappa)$.
- Sample $(\boldsymbol{\theta}, \sigma^2) \sim \mathcal{NIG}(a_n, b_n, \boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n)$ by sampling $\sigma^2 \sim \mathcal{IG}(a_n, b_n)$, followed by sampling $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_n, \sigma_y^2 \boldsymbol{\Lambda}_n^{-1})$ with

$$\boldsymbol{\Lambda}_n = \sum_{j=1}^J \mathbf{S}_j + \boldsymbol{\Lambda}_0, \quad \boldsymbol{\mu}_n = \boldsymbol{\Lambda}_n \left(\sum_{j=1}^J \mathbf{z}_j + \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0 \right), \quad a_n = a_0 + n/2, \quad b_n = 0.5u + \boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0 \boldsymbol{\mu}_0 - \boldsymbol{\mu}_n^T \boldsymbol{\Lambda}_n \boldsymbol{\mu}_n.$$

C.2. A Variant of adaSSP for the Distributed Setting

The adaSSP algorithm of (Wang, 2018) is originally designed for a single data holder. In adaSSP, a differentially private estimate of $\boldsymbol{\theta}$ is released as

$$\hat{\boldsymbol{\theta}} = (\widehat{\mathbf{S}} + \lambda \mathbf{I}_d)^{-1} \hat{\mathbf{z}}. \quad (14)$$

Here $\hat{\mathbf{S}}$ and $\hat{\mathbf{z}}$ are the privatised versions of \mathbf{S} and \mathbf{z} as in (2) and (3), except that ϵ and δ must be changed to $2\epsilon/3$ and $2\delta/3$ in those equations to provide ϵ, δ -DP. This is because `adaSSP` uses another parameter, λ , which is also calculated from the sensitive data and a third of the privacy budget is spent for privatising that calculation. With $v \sim \mathcal{N}(0, 1)$, λ is specifically calculated as

$$\lambda = \max \left\{ 0, \sigma \sqrt{d \ln(6/\delta) \ln(2d^2/\rho)} - \tilde{\lambda}_{\min} \right\}$$

with $\sigma = \|\mathbf{X}\|^2/(\epsilon/3)$, $\lambda_{\min} = \min(\text{eig}(\mathbf{S}))$, and $\tilde{\lambda}_{\min} = \max \left\{ \lambda_{\min} + \sqrt{\ln(6/\delta)}\sigma v - \ln(6/\delta)\sigma v, 0 \right\}$. We consider an extension of (Wang, 2018) for $J \geq 1$. To perform the extension, we reflect on its tendency to approximate a (regularised) least square solution and consider the following estimate

$$\hat{\theta} = \left(\sum_{j=1}^J \hat{\mathbf{S}}_j + \mathbf{I}_d \sum_{j=1}^J \lambda_j \right)^{-1} \left(\sum_{j=1}^J \hat{\mathbf{z}}_j \right). \quad (15)$$

Here $\hat{\mathbf{S}}_j$, $\hat{\mathbf{z}}_j$ and λ_j are calculated in data node j separately from the other nodes. The estimation procedure in (15) does not properly account for the Bayesian paradigm but aggregates the shared $\hat{\mathbf{S}}_j$'s and $\hat{\mathbf{z}}_j$'s to approximate the (regulated) least squares solution. Note that each node has separate λ_j because it depends on δ_j , and so the number of rows for each node.

D. Additional results

Figure 4 shows the MMD estimates for $d = 2$. The (squared) MMD between two distributions can be estimated unbiasedly using i.i.d. samples from those distributions. Non-private posterior and private posteriors of `Bayes-fixedS-fast` are in closed form and can be sampled easily. For the MCMC models, we use every 50th sample of the chain to avoid autocorrelation and thus obtain nearly independent samples.

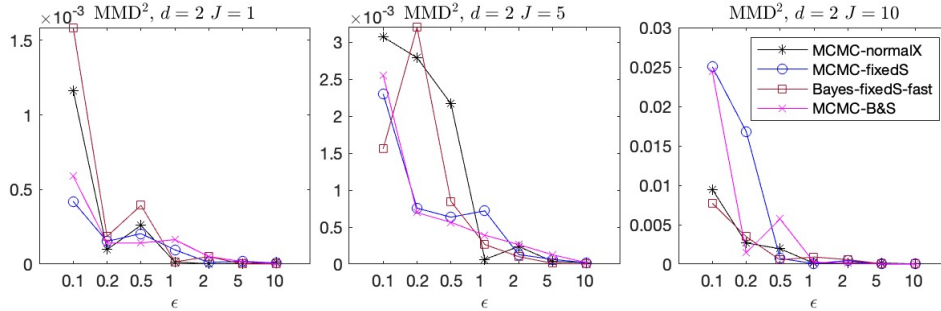


Figure 4. MMD² estimates for each J and $d = 2$.

Table 2 shows the confidence intervals for the prediction MSE for the real-data experiments.

Table 2. 90% confidence interval for mean prediction MSE (over 50 runs) for the real-data experiments - $\epsilon = 1$

J	data sets	MCMC-normalX	MCMC-fixedS	Bayes-fixedS-fast	MCMC-B&S	adaSSP
$J = 1$	PowerPlant	[0.0128, 0.0129]	[0.0128, 0.0129]	[0.0128, 0.0129]	[0.0128, 0.0129]	[0.0137, 0.0140]
	BikeSharing	[0.0021, 0.0027]	[0.0018, 0.0024]	[0.0018, 0.0024]	[0.0017, 0.0022]	[0.0106, 0.0108]
	AirQuality	[0.0051, 0.0069]	[0.0048, 0.0066]	[0.0048, 0.0066]	[0.0053, 0.0071]	[0.0065, 0.0067]
	3droad	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]
$J = 5$	PowerPlant	[0.0132, 0.0135]	[0.0132, 0.0136]	[0.0132, 0.0136]	[0.0135, 0.0138]	[0.0234, 0.0236]
	BikeSharing	[0.0137, 0.0210]	[0.0041, 0.0049]	[0.0040, 0.0049]	[0.0076, 0.0095]	[0.0380, 0.0383]
	AirQuality	[0.0109, 0.0175]	[0.0089, 0.010]	[0.0089, 0.0109]	[0.0109, 0.0151]	[0.0226, 0.0229]
	3droad	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]
$J = 10$	PowerPlant	[0.0139, 0.0145]	[0.0140, 0.0146]	[0.0140, 0.0146]	[0.0141, 0.0146]	[0.0349, 0.0353]
	BikeSharing	[0.0671, 0.0954]	[0.0072, 0.0092]	[0.0072, 0.0092]	[0.0116, 0.0158]	[0.0524, 0.0527]
	AirQuality	[0.0733, 0.1236]	[0.0099, 0.0135]	[0.0099, 0.0135]	[0.0175, 0.0257]	[0.0313, 0.0315]
	3droad	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]	[0.0229, 0.0229]