

---

# Differentially Private Synthetic Control

---

Saeyoung Rho  
Columbia University

Rachel Cummings  
Columbia University

Vishal Misra  
Columbia University

## Abstract

Synthetic control is a causal inference tool used to estimate the treatment effects of an intervention by creating synthetic counterfactual data. This approach combines measurements from other similar observations (i.e., *donor pool*) to predict a counterfactual time series of interest (i.e., *target unit*) by analyzing the relationship between the target and the donor pool before the intervention. As synthetic control tools are increasingly applied to sensitive or proprietary data, formal privacy protections are often required. In this work, we suggest the first algorithms for differentially private synthetic control with explicit error bounds based on the analysis of the sensitivity of the synthetic control query. Our approach builds upon tools from non-private synthetic control and differentially private empirical risk minimization. We empirically evaluate the performance of our algorithms and show favorable results in a variety of parameter regimes.

## 1 INTRODUCTION

The fundamental problem of causal inference is that for an individual unit, we can only observe one of the relevant outcomes – with a particular treatment or without (Rubin, 1974). To estimate the (causal) effect of a treatment, one has to produce a counterfactual of the test arm, which is typically done at a population- and distributional-level via randomized control trials (RCTs) and A/B testing, yielding average treatment effects. However, controlled trials are often impossible to implement, and only observational data are available. Synthetic control is a powerful causal inference tool to estimate the treatment effect of interventions using only observational data. It has been used both at an aggregate population level (e.g., countries/cities/cohorts

of patients etc.), as well as at an individual unit level (Amjad et al., 2019; Agarwal et al., 2021a), and has been called “arguably the most important innovation in the policy evaluation literature in the last 15 years” (Athey and Imbens, 2017).

Recently, synthetic control has increasingly been used in clinical trials where running a randomized control trial presents logistical challenges (e.g., rare diseases), or ethical issues (e.g., oncology trials enrolling patients for placebos with life threatening diseases) (Thorlund et al., 2020). Synthetic control has been successfully used to achieve regulatory approval for new medical treatments for lung cancer (Petrone, 2018) and rare forms of leukemia (Gökbuget et al., 2016), where RCTs would otherwise have been impossible. Since these synthetic control analyses are deployed in real-world medical applications, preserving privacy of sensitive patient data is paramount.

Differential privacy by Dwork et al. (2006) has emerged as the de facto gold-standard in privacy-preserving data analysis. It is a mathematically rigorous parameterized privacy notion, which bounds the maximum amount that can be learned about any data donor based on analysis of her data. Differentially private algorithms have been designed for a wide variety of optimization, learning, and data-driven decision-making tasks, and have been deployed in practice by several major technology companies and government agencies. Despite the growing maturity of the differential privacy toolkit, the pressing need for a private synthetic control solution has thus far gone unaddressed.

In this paper, we propose the first algorithms for differentially private synthetic control (Algorithm 2 and 3). These two algorithms naturally extend existing non-private techniques for synthetic control by first privately estimating the regression coefficients  $\hat{f}$  that relate a (target) pre-intervention observation of interest  $y_{pre}$  to other similar (donor) observations  $X_{pre}$ . This is done using output perturbation and objective perturbation techniques for differentially private empirical risk minimization (DP-ERM) by Chaudhuri et al. (2011) and Kifer et al. (2012). The algorithm then combines the private regression coefficients  $\hat{f}$  with privatized post-intervention donor observations  $\tilde{X}_{post}$  (also via output perturbation) to predict the post-intervention target outcome  $\hat{y}_{post} = \tilde{X}_{post}^\top \hat{f}$ .

We provide privacy and accuracy guarantees for each algorithm. For privacy (Theorems 3.1 and 3.4), although our algorithmic techniques rely on existing approaches, prior results on privacy do not apply in our setting. DP methods add noise that scales with the *sensitivity* of the function being computed, which is defined as the maximum change in the function’s output that can be caused by changing a single donor’s data. However, synthetic control performs a regression in a vertical way, treating each time point, rather than one donor’s data point, as one sample – thus, the *transposed* setting changes the definition of neighboring databases, completely altering the impact of a single donor’s data. The majority of our privacy analysis is devoted to computing sensitivity of this new method.

For accuracy guarantees (Theorems 3.2 and 3.5), we bound the root mean squared error (RMSE) of the algorithm’s output compared to the post-intervention target signal. Our bounds are comparable to those for non-private synthetic control (e.g., Amjad et al. (2018)), and in Section C.1.1, we explicitly show that the RMSE of our algorithm relative to a non-private version is only greater by a factor of  $O(1/\epsilon)$  for output perturbation, which is unavoidable in most analysis tasks. To better interpret our bounds in terms of natural problem parameters such as number of samples and length of observations, we also provide Corollaries C.5 and E.2, which give explicit closed-form upper bounds on the RMSE under mild assumptions on the underlying data distribution.

## 1.1 Related Work

**Synthetic Control.** Synthetic control (SC) was originally proposed by Abadie and Gardeazabal (2003) to evaluate the effects of intervention by creating synthetic counterfactual data. Its first application was measuring the economic impact of the 1960s terrorist conflict in Basque Country, Spain by combining GDP data from other Spanish regions prior to the conflict to construct a *synthetic* GDP dataset for Basque Country in the counterfactual world without the conflict. Synthetic control has since been applied to a wide array of topics such as estimating the effect of California’s tobacco control program (Abadie et al., 2010), estimating the effect of the 1990 German reunification on per capita GDP in West Germany (Abadie et al., 2015), evaluating health policies (Kreif et al., 2016), forecasting weekly sales at Walmart stores (Amjad et al., 2019), and predicting cricket score trajectories (Amjad et al., 2019).

The core algorithm of synthetic control lies on finding a relationship between the target time series (e.g., GDP of Basque Country) and the donor pool (e.g., GDP of other Spanish regions). The original method by Abadie and Gardeazabal (2003) used linear regression with a simplex constraint on the weights: the regression coefficients should be non-negative and sum to one. Since its first intro-

duction, the synthetic control literature has evolved to include a richer set of techniques, including tools to deal with multiple treated units (Abadie and L’Hour, 2021; Dube and Zipperer, 2015), to correct bias (Abadie and L’Hour, 2021; Ben-Michael et al., 2021), to use Lasso and Ridge regression instead of linear regression with simplex constraints (Doudchenko and Imbens, 2016; Amjad et al., 2018), and to incorporate matrix completion techniques (Athey et al., 2021; Amjad et al., 2018, 2019). See a review paper by Abadie (2021) for a detailed survey of these techniques.

The most relevant extension for our work is *robust synthetic control* (RSC) by Amjad et al. (2018), which comprises of two steps: first de-noising the data via hard singular value thresholding (HSVT), and then learning and projecting via regression. It assumes a latent variable model and applies HSVT before running the regression, which reduces the rank of the data. RSC also relaxes the simplex constraints on the regression coefficients and applies unconstrained Ridge regression. Because of the de-noising step, RSC can be viewed as an instantiation of principal component regression (PCR) and the possibility of differentially private PCR has been briefly discussed by Agarwal et al. (2021b). However, no formal algorithm or analysis has been put forth until this paper.

## Differentially Private Empirical Risk Minimization.

Chaudhuri et al. (2011) first proposed methods for differentially private empirical risk minimization (ERM) for supervised regression and classification. Our first algorithm uses the *output perturbation* method by Chaudhuri et al. (2011), which first computes coefficients to minimize the loss function between data features and labels, and then perturbs the coefficients using a high-dimensional variant of the Laplace Mechanism by Dwork et al. (2006). Our second algorithm uses the *objective perturbation* method by Chaudhuri et al. (2011) and Kifer et al. (2012), which adds noise directly to the loss function and then exactly optimizes the noisy loss. This method tends to provide better theoretical accuracy guarantees but requires the loss function to satisfy additional structural properties. These methods were later extended by Bassily et al. (2014) to include *gradient perturbation* in stochastic gradient descent, which uses a noisy version of randomly sampled points’ contribution to the gradient at each update. This technique provides tighter error bounds, assuming Lipschitz convex loss and bounded optimization domain. Wang et al. (2017) followed up with a faster gradient perturbation algorithm that provided a tighter upper bound on error and lower gradient complexity.

Although the framework by Chaudhuri et al. (2011) is more general, the analysis and applications focused only on methods for binary classification. The analysis was later extended to include ridge regression by Cummings et al. (2015), which we use in our algorithms. Our algorithms

for differentially private synthetic control apply DP-ERM methods to a ridge regression loss function. However, synthetic control applies regression in the transposed dimension of the data (i.e., along columns rather than rows of the database), while privacy protections are still required along the rows, which requires novel analysis to ensure differential privacy and accuracy.

## 2 MODEL AND PRELIMINARIES

In this section, we first present our model (Section 2.1) and then provide relevant background on synthetic control (Section 2.2) and differential privacy (Section 2.3).

### 2.1 Our Model

Our model follows the synthetic control framework illustrated in Figure 1. We consider a database  $X \in \mathbb{R}^{n \times T}$ , also called the *donor pool*. The donor pool  $X$  consists of  $n$  time series, each observed at times  $t = 1, \dots, T$ . We denote the column vectors of  $X$  as  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^n$ , where each  $\mathbf{x}_t$  contains observations from all donor time series at time  $t$ . We assume an intervention occurred at a known time  $T_0 + 1 < T$ . The first  $T_0$  columns of  $X$  are collectively referred to as  $X_{pre}$ , and the remaining  $T - T_0$  columns from data after the intervention are collectively denoted  $X_{post}$ , respectively corresponding to the pre- and post-intervention donor data. We are also given a *target unit*  $\mathbf{y} \in \mathbb{R}^T$ , which can be divided as  $\mathbf{y}_{pre} = (y_1, \dots, y_{T_0})$  and  $\mathbf{y}_{post} = (y_{T_0+1}, \dots, y_T)$ .

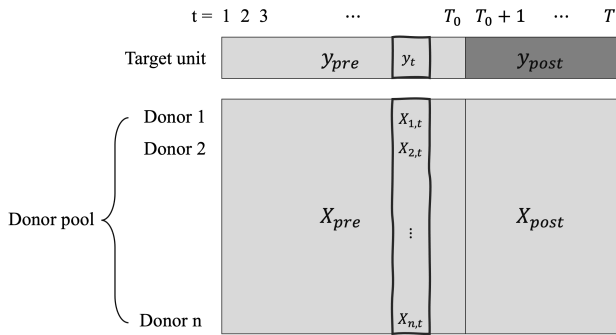


Figure 1: General data structure for synthetic control. The donor pool ( $X$ ) and the target unit ( $\mathbf{y}$ ) are divided into pre- and post-intervention periods. Synthetic control first performs a *vertical* regression using the pre-intervention column vectors  $\mathbf{x}_t$  as features for the label  $y_t$  for  $t \in [T_0]$  to estimate regression coefficients  $\hat{\mathbf{f}}$ , and then uses this to project the post-intervention column vectors and predict  $\hat{\mathbf{y}}_{post}$ .

The underlying assumption is that time series in the donor pool that behave similarly to  $\mathbf{y}$  before the intervention will remain similar after the intervention. In this paper, we use the latent variable model, the same as Amjad et al. (2018),

for the underlying distribution of the data. Our donor data and target data are noisy versions of the true signal (denoted  $M$  and  $\mathbf{m}$  respectively), and can be written as follows:

$$X = M + Z, \quad \mathbf{y} = \mathbf{m} + \mathbf{z}, \quad (1)$$

where  $Z \in \mathbb{R}^{n \times T}$  is a noise matrix where each element is sampled i.i.d. from some distribution with zero-mean,  $\sigma^2$ -variance, and support  $[-s, s]$ , and  $\mathbf{z} \in \mathbb{R}^T$  is a noise vector with elements sampled from the same distribution.

The signals  $M$  and  $\mathbf{m}$  can be expressed in terms of a latent function  $g$ :

$$M_{i,t} = g(\theta_i, \rho_t) \quad m_t = g(\theta_0, \rho_t), \quad \theta_i \in [n], t \in [T],$$

where  $\theta_i$  and  $\rho_t$  are latent feature vectors capturing unit  $i$ 's and time  $t$ 's intrinsic characteristics, respectively. We note that if the intervention is effective, one would expect to see a change in  $\rho_t$  before and after  $T_0$ . We make no assumptions on the latent function  $g$ , except in Sections C.2 and E.2, where we assume  $M$  is low rank, i.e.,  $\text{rank}(M) = k$  for some  $k \leq \min\{n, T\}$ .

Finally, we assume a linear relationship between the features of  $M_{i,t}$  and the label  $m_t$  at all times  $t \in [T_0]$ ; that is, there exists an  $\mathbf{f} \in \mathbb{R}^n$  such that,

$$m_t = \sum_{i=1}^n M_{i,t} f_i, \quad \text{for all } t \in [T_0]. \quad (2)$$

We assume that all entries of  $X$ ,  $M$ ,  $\mathbf{y}$ , and  $\mathbf{m}$  lie in a bounded range, which we rescale to  $[0, 1]$  WLOG, and that  $\mathbf{f}$  has  $\ell_1$ -norm bounded by 1, as is standard in the synthetic control literature (Abadie and Gardeazabal, 2003). Formally, we assume:

$$\begin{aligned} |x_{i,t}| &\leq 1, |M_{i,t}| \leq 1 \quad \forall t \in [T], i \in [n], \\ |y_t| &\leq 1, |m_t| \leq 1 \quad \forall t \in [T], \text{ and } \|\mathbf{f}\|_1 = \sum_{k=1}^n |f_k| \leq 1. \end{aligned} \quad (3)$$

### 2.2 Synthetic Control

The goal of synthetic control is to predict  $\mathbf{y}_{post}$  given  $X$  and  $\mathbf{y}_{pre}$ . The general approach, outlined in algorithm 1, is to first use the pre-intervention data  $D_1 := (X_{pre}, \mathbf{y}_{pre})$  to learn an estimate  $\hat{\mathbf{f}}$  of the true coefficient vector  $\mathbf{f}$ . For each  $t \in [T_0]$ , the column vector  $\mathbf{x}_t = (X_{1,t}, \dots, X_{n,t})^\top$  is treated as a feature vector for label  $y_t$ . This setup distinguishes synthetic control from the classic regression setting, as the regression is performed vertically rather than horizontally. The estimate  $\hat{\mathbf{f}}$  is then used along with the post-intervention donor data to predict the counterfactual outcome of the target:  $\hat{\mathbf{y}}_{post} = X_{post}^\top \hat{\mathbf{f}}$ , where  $\hat{y}_t = \mathbf{x}_t^\top \hat{\mathbf{f}}$   $\forall t \in [T_0 + 1, \dots, T]$ .



and  $DPSC_{obj}$  (Algorithm 3). Similar to non-private synthetic control algorithms (e.g., Algorithm 1), both algorithms are divided into two high-level steps: first the algorithm learns an estimate of the regression coefficients  $\mathbf{f}$ , and then it uses these coefficients to predict the post-intervention target unit  $\mathbf{y}_{post}$ . To ensure differential privacy of the overall algorithm, both of these steps must be performed privately. The second step remains the same for both, and only the first part differs:  $DPSC_{out}$  adds privacy noise directly to the output of the algorithm (output perturbation), whereas  $DPSC_{obj}$  perturbs the objective function and minimize the noisy objective (objective perturbation).

### 3.1 DPSC via Output Perturbation $DPSC_{out}$

Our first algorithm is  $DPSC_{out}$  (Algorithm 2), which utilizes *output perturbation* to achieve differential privacy. The learning step of this algorithm formalizes synthetic control as an instance of empirical risk minimization with the Ridge regression loss function given in Equation (4). This enables us to apply the Output Perturbation method by Chaudhuri et al. (2011) for DP-ERM. The algorithm first learns non-private regression coefficients  $\mathbf{f}^{reg}$  as in Algorithm 1. It then samples a noise vector  $\mathbf{v}$  from a high-dimensional Laplace distribution with parameter  $\Delta \mathbf{f}^{reg}/\epsilon_1$ , as described in Section 2.3. Finally, the privatized regression coefficient vector is  $\mathbf{f}^{out} = \mathbf{f}^{reg} + \mathbf{v}$ .

The prediction step uses this coefficient vector to predict  $\mathbf{y}_{post}$ . A simple approach would be to directly predict  $\hat{\mathbf{y}}_{post} = X_{post}^\top \mathbf{f}^{out}$ ; however, this approach would not provide privacy for the post-intervention donor data  $X_{post}$ . Instead, we again apply the high-dimensional Laplace Mechanism to privatize  $X_{post}$  by adding a noise matrix  $W$  sampled from a high-dimensional Laplace distribution with parameter  $\Delta X_{post}/\epsilon_2$ . The privatized version of donor data is  $\tilde{X}_{post} = X_{post} + W$ , which is then used along with  $\mathbf{f}^{out}$  to produce the private prediction of the post-intervention target unit:  $\mathbf{y}^{out} = \tilde{X}_{post}^\top \mathbf{f}^{out}$ .

The entire algorithm is then  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private by composition of these two steps. We remark that the algorithm does not output  $\mathbf{f}^{out}$ , simply because this vector is typically not of interest in most cases, and is instead considered only an intermediate analysis step. However, this vector could be output if desired with no additional privacy loss because Step 1 of the  $DPSC_{out}$  algorithm is  $\epsilon_1$ -differentially private (Theorem B.1), and this privacy loss is already accounted for in the composition step.

We provide two main results on the privacy and accuracy of  $DPSC_{out}$ . First, Theorem 3.1 shows that this algorithm is differentially private. Although our algorithm relies on algorithmic techniques by Chaudhuri et al. (2011) for DP-ERM, the vertical regression setup in synthetic control requires novel sensitivity analysis for  $\mathbf{f}^{reg}$ , which constitutes the bulk of the work required to prove Theorem 3.1. The-

---

**Algorithm 2** DPSC via Output Perturbation  
 $DPSC_{out}(X_{pre}, X_{post}, \mathbf{y}_{pre}, n, T, T_0, \lambda, \epsilon_1, \epsilon_2)$

---

**Step 1: Learn regression coefficients**

Learn the regression coefficient  $\mathbf{f}^{reg}$  using Ridge regression with parameter  $\lambda = 0$ :

$$\mathbf{f}^{reg} = \arg \min_{\mathbf{f} \in \mathbb{R}^n} \frac{1}{T_0} \|\mathbf{y}_{pre} - X_{pre}^\top \mathbf{f}\|_2^2 + \frac{\lambda}{2T_0} \|\mathbf{f}\|_2^2.$$

$$\text{Let } a = \frac{\mathbf{f}^{reg}}{\epsilon_1} = \frac{4T_0 \sqrt{8+n}}{\lambda \epsilon_1}$$

Sample  $\mathbf{v}$  according to pdf  $p(\mathbf{v}; a) \propto \exp\left(-\frac{\|\mathbf{v}\|_2}{a}\right)$

$$\text{Let } \mathbf{f}^{out} = \mathbf{f}^{reg} + \mathbf{v}$$

**Step 2: Predict  $\mathbf{y}_{post}$  via projection**

$$\text{Let } b = \frac{2\sqrt{T-T_0}}{\epsilon_2}$$

Sample each entry of  $W \in \mathbb{R}^{n \times (T-T_0)}$  i.i.d. according to pdf  $p(W; b) \propto \exp\left(-\frac{\|W\|_F}{b}\right)$

$$\text{Let } \tilde{X}_{post} = X_{post} + W$$

$$\text{Output } \mathbf{y}^{out} = \tilde{X}_{post}^\top \mathbf{f}^{out}.$$


---

orem 3.2 shows that our  $DPSC_{out}$  algorithm produces an accurate prediction of the post-intervention target unit, as measured by the standard metric of root mean squared error (RMSE) with respect to the true signal vector  $\mathbf{m}$ . In Section C.2, we also extend Theorem 3.2 to to remove the dependence on distributional parameters and provide an expression of RMSE that depends only on the input parameters, under some mild additional assumptions on the distribution of data. Full proofs for Theorems 3.1 and 3.2 are respectively presented in Sections B and C.

**Theorem 3.1**  $DPSC_{out}$  of Algorithm 2 is  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private.

**Theorem 3.2** The estimator  $\mathbf{y}^{out}$  output by Algorithm 2 satisfies:

$$\begin{aligned} RMSE(\mathbf{y}^{out}) &= \frac{\|M_{post}\|_2}{T - T_0} \mathbb{E}[\|\mathbf{f}^{reg} - \mathbf{f}\|_2] + \frac{4T_0 \rho_{\frac{8+n}{\lambda \epsilon_1}}}{\lambda \epsilon_1} \\ &+ \frac{\rho_{\frac{8+n}{\lambda \epsilon_1}}}{n\sigma^2} + \frac{\rho_{\frac{8+n}{\lambda \epsilon_1}}}{\epsilon_2} \rho_{\frac{8+n}{\lambda \epsilon_1}} + \frac{4T_0 \rho_{\frac{8+n}{\lambda \epsilon_1}}}{\lambda \epsilon_1}, \end{aligned}$$

where  $\|\mathbf{f}^{reg}\|_\infty \leq \psi$  for some  $\psi > 0$ , and  $RMSE$  is the root mean squared error of the estimator, defined as  $RMSE(\mathbf{y}^{out}) = \frac{1}{\sqrt{T-T_0}} \mathbb{E}[\|\mathbf{y}^{out} - \mathbf{m}_{post}\|_2]$ .

**Remark 3.3** The accuracy bound grows as  $O(n)$ , which is shown to be necessary in Section B.1.1. While this might be undesirable in most other learning domains,  $n$  does not grow with the problem size in synthetic control settings for several reasons. Typically,  $M$  is assumed to be a low-rank

matrix and hence  $X$  is approximately low rank (Amjad et al., 2018, 2019). This is not only an assumption, but true in most cases (Udell and Townsend, 2019). Therefore, there exists a saturation point where adding additional donors does not meaningfully improve accuracy (see Section 3.3 for more details). The remaining dependence on  $T_0$  can be handled by setting  $\lambda = O(T_0)$  (see Section C.1.1).

### 3.2 DPSC via Objective Perturbation $DPSC_{obj}$

We next present our second algorithm for differentially private synthetic control,  $DPSC_{obj}$  (Algorithm 3), based on objective perturbation. While Step 2 remains unchanged relative to Algorithm 2, Step 1 is modified to perturb the objective function itself and then exactly optimize the perturbed objective, instead of first computing the optimal non-private coefficients and then adding noise. Objective perturbation has been shown to outperform output perturbation in the standard private ERM setting when the loss function is strongly convex (Chaudhuri et al., 2011).

The algorithm augments the objective function with two terms. The first is an additional regularization term to ensure  $\frac{\lambda}{T_0}$ -strong convexity (compared to  $\frac{\lambda}{T_0}$  with the regularization term of Algorithm 2). The  $\Delta$  parameter is tuned by the algorithm to ensure that it can still satisfy  $(\epsilon_1, \delta)$ -DP in Step 1, even when  $\epsilon_1$  is small. The second is the noise term  $\mathbf{b}^\top \mathbf{f}$  to ensure privacy, where  $\mathbf{b}$  is sampled from a high-dimensional Laplace distribution if  $(\epsilon, 0)$ -DP is desired (i.e., if  $\delta = 0$ ), and from a multi-variate Gaussian distribution if  $(\epsilon, \delta)$ -DP is desired (i.e., if  $\delta > 0$ ).

The algorithm then exactly optimizes this new objective function, where the noise term  $\mathbf{b}$  ensures that this minimization satisfies differential privacy. Although the algorithmic procedure in Step 1 is similar to that of Objective Perturbation algorithms for DP-ERM by Chaudhuri et al. (2011) and Kifer et al. (2012), the sensitivity and privacy analysis again requires substantial novelty because the definition of neighboring databases change and previous work cannot be immediately applicable to the transposed regression setting. Finally, Algorithm 3 maintains the same Step 2 process as Algorithm 2 to predict  $\mathbf{y}_{post}$ , based on  $\mathbf{f}^{obj}$  computed from Step 1. Algorithm 3 is  $(\epsilon_1 + \epsilon_2, \delta)$ -DP by composition of privacy guarantees from these two steps.

$DPSC_{obj}$  requires an additional parameter  $c$  that is used in the analysis to bound the maximum absolute eigenvalue of  $2(X'_{pre} X'_{pre}{}^\top - X_{pre} X_{pre}{}^\top)$ , which is closely related to  $\| \mathbf{f} \|^2$ . Because  $X_{pre}$  and  $X'_{pre}$  are neighboring databases, the matrix of interest will only have one column and one row that are non-zero. In our setting, we use the fact that all entries of  $X$  are in  $[ -1, 1 ]$  to derive an upper bound on this matrix and its eigenvalues. In general, an analyst can use domain expertise or prior knowledge of the data distribution to choose an appropriate value of  $c$ .

---

#### Algorithm 3 DPSC via Objective Perturbation $DPSC_{obj}(X_{pre}, X_{post}, y_{pre}, n, T, T_0, \lambda, \epsilon_1, \epsilon_2, \delta, c)$

---

##### Step 1: Learn regression coefficients

if  $\epsilon_1 > \log(1 + \frac{2c}{\lambda} + \frac{c^2}{\lambda^2})$  then

Let  $\epsilon_0 = \epsilon_1 - \log(1 + \frac{2c}{\lambda} + \frac{c^2}{\lambda^2})$  and  $\Delta = 0$

else

$\epsilon_0 = \frac{\epsilon_1}{2}$  and  $\Delta = \frac{c}{e^{(\frac{1}{\lambda^2})} - 1} \lambda$

end if

if  $\delta > 0$  then

Sample  $\mathbf{b} \geq \mathbb{R}^n$  from  $\mathcal{N}(0, \beta^2 I_n)$ ,

where  $\beta = \frac{4T_0\sqrt{8+n}}{\epsilon_0} \frac{2 \log^2 + 2\epsilon_0}{\epsilon_0}$

else

Sample  $\mathbf{b} \geq \mathbb{R}^n$  from  $p(\mathbf{b}; \beta) \propto \exp \left( -\frac{\|\mathbf{b}\|_2}{\beta} \right)$ ,

where  $\beta = \min \left\{ \frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n} + 4T_0}{\epsilon_0} g \right\}$

end if

Learn  $\mathbf{f}^{obj}$  by minimizing

$$J = \frac{1}{T_0} \mathbf{y}_{pre}^\top \mathbf{X}_{pre}^\top \mathbf{f} \mathbf{f}^\top \mathbf{f} + \frac{\lambda + \Delta}{2T_0} \mathbf{f} \mathbf{f}^\top \mathbf{f} + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f}.$$

##### Step 2: Predict $\mathbf{y}_{post}$ via projection

Let  $b = \frac{2\sqrt{T-T_0}}{\epsilon_2}$

Sample each entry of  $W \geq \mathbb{R}^{n \times (T-T_0)}$  i.i.d. according to pdf  $p(W; b) \propto \exp \left( -\frac{\|W\|_F}{b} \right)$

Let  $\tilde{X}_{post} = X_{post} + W$

Output  $\mathbf{y}^{obj} = \tilde{X}_{post}^\top \mathbf{f}^{obj}$

---

We provide two main results on the privacy and accuracy of  $DPSC_{obj}$ . First, Theorem 3.4 shows that our algorithm is differentially private. To prove privacy in Step 1, we must consider two cases based on the value of  $\Delta$ , which adds additional strong convexity to the loss function if it is needed. The privacy budget must be allocated differently within the analysis in the two cases of  $\Delta = 0$  and  $\Delta > 0$ .

Theorem 3.5 shows that  $DPSC_{obj}$  produces an accurate prediction of the post-intervention target unit, as measured as RMSE between its output  $\mathbf{y}^{obj}$  and the target unit's post-intervention signal vector  $\mathbf{m}_{post}$ . As with  $DPSC_{out}$ , we also extend Theorem 3.5 in Section E.2 to provide an explicit closed-form bound on RSME that does not depend on the distributional parameters. Full proofs for Theorems 3.4 and 3.5, along with their extensions, are respectively presented in Sections D and E.

**Theorem 3.4**  $DPSC_{obj}$  of Algorithm 3 is  $(\epsilon_1 + \epsilon_2, \delta)$ -differentially private.

**Theorem 3.5** The estimator  $\mathbf{y}^{obj}$  output by Algorithm 3

satisfies:

$$\begin{aligned}
 RMSE(\mathbf{y}^{obj}) & \frac{jjM_{post}jj_2}{T} E[|j(\mathbf{f}^{reg} - \mathbf{f})j|_2] \\
 & + \frac{2}{\lambda + \Delta} E[|j\mathbf{b}j|_2] + 1 \neq 0 \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} 2T_0^2 \rho_{\bar{n}} \\
 & + \frac{\rho_{\bar{n}}}{n\sigma^2} + \frac{\rho_{\bar{n}}}{\epsilon_2} \rho_{\bar{n}} \psi + \frac{2}{\lambda + \Delta} E[|j\mathbf{b}j|_2] \\
 & + 1 \neq 0 \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} 2T_0^2 \rho_{\bar{n}} ,
 \end{aligned}$$

where  $jj\mathbf{f}^{reg}jj_{\infty} = \frac{\psi}{\epsilon_0}$  for some  $\psi > 0$ , and  $E[|j\mathbf{b}j|_2] = \frac{\rho_{\bar{n}}}{nT_0 4\sqrt{8+n} 2 \log^2 + \epsilon_0}$  for Gaussian noise ( $\delta > 0$  case)

and  $E[|j\mathbf{b}j|_2] = \min\{f^{\frac{4T_0\sqrt{8+n}}{\epsilon_0}}, \frac{c\sqrt{n}+4T_0}{\epsilon_0} g\}$  for Laplace noise ( $\delta = 0$  case), and  $\epsilon_0$ , and  $\Delta$  are computed internally by the algorithm.

As in Section 3.1, we remark that while the accuracy bound of Theorem 3.5 grows as  $O(n)$ , in our setting  $n$  does not typically grow substantially with the problem size, both in theory (Amjad et al., 2018, 2019) and in practice (Udell and Townsend, 2019).

### 3.3 Comparing DP-ERM and DPSC

In this section, we compare the results of DP-ERM by Chaudhuri et al. (2011) with our approach. Consider a Ridge regression task in  $p$ -dimensional space with  $q$  samples (i.e., covariates  $x_k \in \mathbb{R}^p$  and labels  $y_k \in \mathbb{R}$ ,  $8k \in [q]$ ). The regression coefficient  $\theta \in \mathbb{R}^p$  is learned by a standard empirical risk minimization process with a regularizer  $\lambda |j\thetaj|_2^2$ . In a typical regression setup where the privacy goal is to protect one sample  $x_k$ , corresponding to one individual's data, the sensitivity of the coefficient is  $\Delta\theta = \frac{2}{q\lambda}$  (Chaudhuri et al., 2011). It does not depend on the dimension  $p$ , and the sensitivity decreases as the number of samples  $q$  increases. Intuitively, adding or removing one person's data should exhibit diminishing marginal effect on the final model  $\theta$  as the training sample size grows.

On the other hand, in our transposed setting of synthetic control, the privacy goal is to protect the  $i$ -th entry of each  $x_k$  (i.e., an individual's data are spread across all samples), the sensitivity is  $\Delta\theta = \frac{4q\sqrt{8+p}}{\lambda}$  (Lemma 3.7). In this setting, each dimension of the coefficient  $\theta$  captures how important the corresponding donor is for explaining the target; hence the impact of changing  $i$ -th person's data will have a significant on the  $i$ -th dimension of  $\theta$ , regardless of the number of individuals in the donor pool. This is at the crux of why it is more difficult to guarantee privacy in the transposed setting of synthetic control, relative to the standard regression setting.

### 3.4 Proof Sketch for Privacy Guarantees

In this section, we outline the proof for privacy guarantees for both algorithms. The full versions with all omitted proofs are presented in Appendices B and D.

The proof of Theorem 3.1 relies on the privacy of  $\mathbf{f}^{out}$  in the learning phase, and then  $\tilde{X}_{post}$  in the prediction phase. At a high level,  $\mathbf{f}^{out}$  is  $\epsilon_1$ -DP through a (non-trivial) application of the Output Perturbation algorithm by Chaudhuri et al. (2011). In the prediction phase, we must show that sufficient noise is added to ensure  $\tilde{X}_{post}$  is an  $\epsilon_2$ -DP version of  $X_{post}$ . Then privacy of  $\mathbf{y}^{out}$  comes from the composition of these two private estimates.

We start by proving that  $\mathbf{f}^{out}$  is  $\epsilon_1$ -DP.

**Theorem 3.6** *Step 1 of Algorithm 2 that computes  $\mathbf{f}^{out}$  is  $(\epsilon_1, 0)$ -differentially private.*

Step 1 of Algorithm 2 instantiates the Laplace mechanism, and the crux of the proof lies in obtaining the sensitivity bound of the synthetic control query, which is fundamentally different from the setting of a traditional DP-ERM (Chaudhuri et al., 2011).

**Lemma 3.7** *The  $\ell_2$  sensitivity of  $\mathbf{f}^{reg}$  is  $\Delta\mathbf{f}^{reg} = \frac{4T_0\sqrt{8+n}}{\lambda}$ .*

The asymptotic dependence on  $n$  and  $T_0$  may seem undesirable, but we show that it is unavoidable.

**Lemma 3.8** *The  $\ell_2$  sensitivity of  $\mathbf{f}^{reg}$  is  $\Delta\mathbf{f}^{reg} = \Omega(\rho_{\bar{n}})$ .*

Next we move to privacy of  $\tilde{X}_{post}$  and its role in ensuring privacy of  $\mathbf{y}^{out}$ .

**Lemma 3.9** *The computation of  $\tilde{X}_{post}$  in Step 2 of Algorithm 2 is  $(\epsilon_2, 0)$ -differentially private.*

$\tilde{X}_{post}$  is privatized through another instantiation of the Laplace Mechanism by Dwork et al. (2006). Thus to prove Lemma 3.9, we only need to bound the sensitivity of  $X_{post}$ . We first note that the Frobenius norm of a matrix  $X \in \mathbb{R}^{n \times (T-T_0)}$  is equal to the  $\ell_2$  norm of the equivalent flattened vector  $X \in \mathbb{R}^{n(T-T_0)}$  (Horn and Johnson, 2012). Thus implementing the matrix-valued Laplace Mechanism with noise parameter calibrated to the  $\ell_2$  sensitivity of the flattened matrix-valued query over  $\epsilon$  will ensure  $(\epsilon, 0)$ -differential privacy. Since all entries in  $X_{post}$  are bounded in  $[-1, 1]$ , each entry can change by at most 2 between two neighboring databases, which can differ in at most  $T - T_0$  entries, hence the  $\ell_2$  sensitivity of flattened  $X_{post}$  is  $2\sqrt{T - T_0}$ .

Finally, we combine Theorem B.1 and Lemma 3.9 to complete the proof of Theorem 3.1. The estimates  $\mathbf{f}^{out}$  and  $\tilde{X}_{post}$  are together  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private by DP

composition, and then  $\mathbf{y}^{out}$  is  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private by post-processing.

Since the prediction step of Algorithm 3 is identical to that of Algorithm 2, we only need to show that  $\mathbf{f}^{obj}$  is computed in an  $(\epsilon_1, \delta)$ -DP manner (Theorem D.1) to complete the proof of Theorem 3.4.

**Theorem 3.10** *Step 1 of Algorithm 3 that computes  $\mathbf{f}^{obj}$  is  $(\epsilon_1, \delta)$ -differentially private.*

We show a proof sketch here and defer the full proof to Appendix D. At a high-level, the privacy of  $\mathbf{f}^{obj}$  comes from a carefully modified instantiation of the Objective Perturbation algorithms by Chaudhuri et al. (2011) and Kifer et al. (2012). The Objective Perturbation method modifies the standard Ridge Regression objective function  $J(\mathbf{f})$  by adding an additional regularization term and a noise term to ensure privacy:

$$\begin{aligned} J^{obj}(\mathbf{f}) &= J(\mathbf{f}) + \frac{1}{2T_0} \|\mathbf{f}\|_2^2 + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f} \\ &= L(\mathbf{f}) + \frac{\lambda^+}{2T_0} \|\mathbf{f}\|_2^2 + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f}, \end{aligned}$$

where  $\mathbf{b}$  is a random vector drawn from a high-dimensional Laplace distribution ( $\delta = 0$ ) or a multivariate Gaussian distribution ( $\delta > 0$ ).

Notice that  $J^{obj}(\mathbf{f})$  is strongly convex (for any  $\Delta > 0$ ) and differentiable. Hence, for any given input dataset  $D = \{X_{pre}, y_{pre}\}$  and any fixed parameters  $(\lambda, \epsilon_1, \epsilon_2, \delta)$ , there exists a bijection between a realized value of the noise term  $\mathbf{b}$  and  $\mathbf{f}^{obj} := \arg \min_{\mathbf{f}} J^{obj}(\mathbf{f})$  given that realized  $\mathbf{b}$ .<sup>1</sup> Let  $\mathbf{b}(\alpha; D)$  be noise value that must have been realized when database  $D$  was input and  $\mathbf{f}^{obj} = \alpha$  was the output. We can then use this bijection to analyze the distribution over outputs on neighboring databases via the (explicitly given) noise distribution. Then, we can express the ratio between the two distributions over  $\mathbf{f}^{obj}$  trained from neighboring databases, and tune the amount of noise depending on the privacy budget.

The ratio is a product of two terms,  $\Gamma(\alpha)$  and  $\Phi(\alpha; \Delta)$ . The parameter  $\Delta$  serves a role to divide the  $\epsilon_1$  budget between these two terms, by distinguishing between two cases. In the first case,  $\epsilon_1$  is large enough that we can choose  $\Delta = 0$  and still have some privacy budget ( $\epsilon_0$ ) remaining to bound  $\Gamma(\alpha)$ . In the other case, if  $\epsilon_1$  is too small to bound  $\Phi(\alpha; \Delta)$  with  $\Delta = 0$ , then we divide the privacy budget equally between bounding  $\Gamma(\alpha)$  and  $\Phi(\alpha; \Delta)$ , and find an appropriate value for  $\Delta > 0$ . For both cases, we prove that  $\Phi(\alpha; \Delta)$  is upper bounded by  $e^{\epsilon_1 - \epsilon_0}$  (Lemma D.2) and  $\Gamma(\alpha)$  is upper bounded by  $e^{\epsilon_0}$  (Lemma D.5). Product of the two upper bounds provides the upper bound of the ratio, completing the proof of Theorem D.1.

<sup>1</sup>For a simple analogy, consider the one-dimensional Laplace Mechanism on query  $f$  and database  $x$ , which outputs  $y = f(x) + Lap(\cdot; f)$ . Given  $f$  and  $x$ , there is a bijection between noise terms and outputs since the noise term must equal  $y - f(x)$ .

### 3.5 Proof Sketch for Accuracy Guarantees

Theorem 3.2 and 3.5 presents accuracy guarantee of the two algorithms in terms of root mean squared error (RMSE), defined as follows:  $RMSE(\mathbf{y}^{out}) = \frac{1}{\sqrt{T-T_0}} \mathbb{E}[\|\mathbf{y}^{out} - \mathbf{m}_{post}\|_2]$ . We note that while it may seem most natural to bound the difference between  $\mathbf{y}^{out}$  and  $\mathbf{y}_{post}$ , we instead use  $\mathbf{m}_{post}$  for two reasons. Firstly,  $\mathbf{y}_{post}$  may not match  $\mathbf{y}_{pre}$  due to the intervention. Secondly,  $\mathbf{m}_{post}$  captures the true signal that we are trying to estimate.

The proof aims to bound the expectation using the submultiplicative norm property (i.e., for any matrix  $A$  and vector  $\mathbf{x}$ ,  $\|A\mathbf{x}\|_2 \leq \|A\|_2 \|\mathbf{x}\|_2$ ) and the known distributions of noise terms  $\mathbf{v}$  (Step 1 of Algorithm 2),  $\mathbf{b}$  (Step 1 of Algorithm 3), and  $W$  (Step 2 of both algorithms). For example, we can decompose  $\mathbf{y}^{out} = \tilde{X}_{post}^\top \mathbf{f}^{out} = (M_{post}^\top + Z^\top + W^\top)(\mathbf{f}^{reg} + \mathbf{v})$ , and compare against  $\mathbf{m}_{post} = M_{post}^\top \mathbf{f}$ .

The full proof of Theorem 3.2 and 3.5 are presented in Appendices C and E.

## 4 EMPIRICAL PERFORMANCE

This section presents the empirical performance of both  $DPSC_{out}$  and  $DPSC_{obj}$  on synthetic datasets. We create four synthetic datasets following the modeling assumptions of Section 2.1, with the number of pre-intervention observations  $T_0 \in \{10, 100\}$  and the number of donors  $n \in \{10, 100\}$ . We set  $T = T_0 + 3$ , meaning that the performance will be measured by RMSE of the next three data points. More details about the data generation process can be found in Appendix A.1

For a fair comparison, we use  $\delta = 0$  for the objective perturbation, so the only privacy parameter we consider is  $\epsilon = \epsilon_1 + \epsilon_2$ . We show two experiments, one with fixed  $\epsilon$  and varying  $\lambda$ , the other with fixed  $\lambda$  and varying  $\epsilon$ . Each experiment was repeated 500 times on each dataset, and the error bands in all figures show 95% confidence intervals taken over the randomness in the algorithms. We present four graphs highlighting the main findings and defer all detailed plots and additional analyses to Appendix A.

### 4.1 Results

For the first set of experiments, we fix  $\epsilon_1 = \epsilon_2 = 50$  and vary the regularization parameter  $\lambda \in \{5, 10, 20, 50, 100, 200, 500, 1000, 2000, 5000\}$ . The top two graphs in Figure 3 show that the optimal choice of  $\lambda$  roughly remains  $\lambda = T_0$  for all sizes of database considered. The orange and blue curves in the figure corresponding to the two databases with  $T_0 = 10$  are approximately minimized at  $\lambda = 10$ , while the green and red curves with



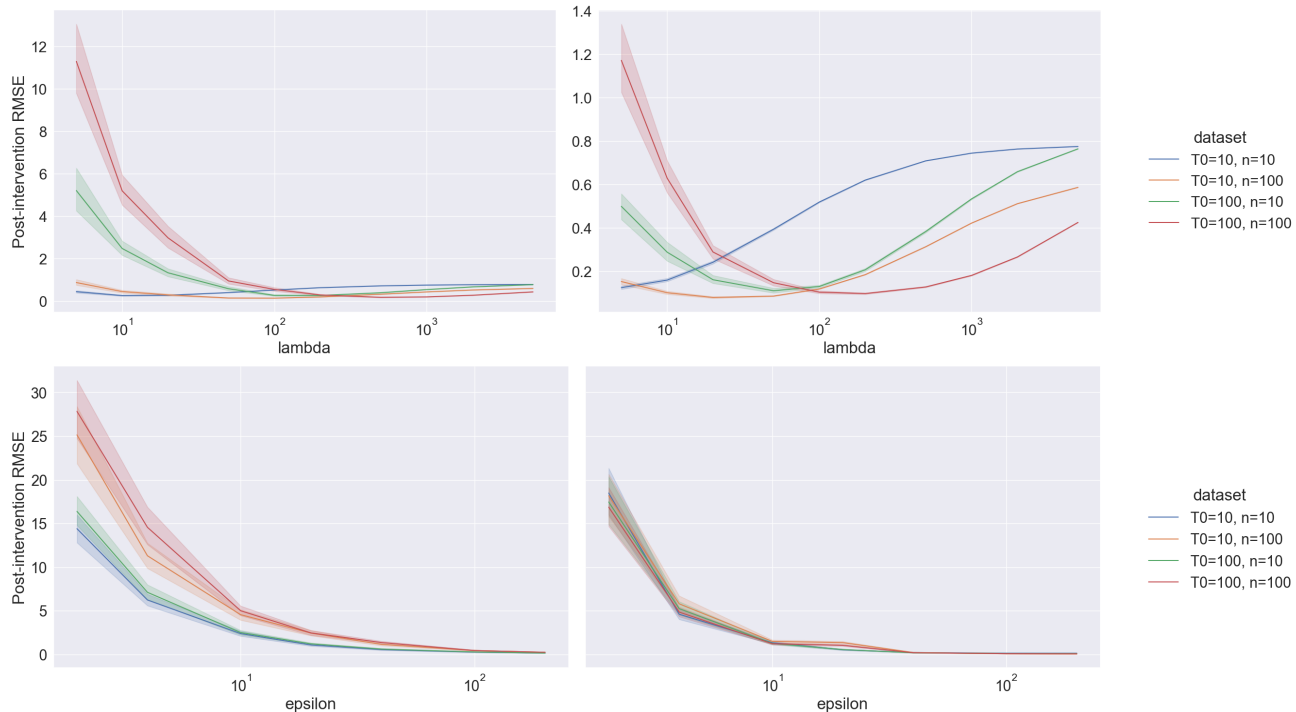


Figure 2: Post-intervention RMSE of  $DPSC_{out}$  (left) and  $DPSC_{obj}$  (right) as a function of  $\lambda$  (top) and as a function of  $\epsilon$  (bottom) on four synthetic datasets of varying size.

$T_0 = 100$  are approximately minimized at  $\lambda = 100$ . While this is clearer in the figure for Objective Perturbation, it is also true for Output Perturbation, although the U-shape is less visible due to the larger scale of the  $y$ -axis.

Additionally, we observe that the RMSE of the private methods is larger than that of the non-private method for smaller  $\lambda = 20$ , with objective perturbation substantially outperforming output perturbation. Also, the empirical performance of both algorithms is dramatically better than the theoretical bounds may suggest, which implies potential room for improvements in the bound of Theorem 3.5. A more detailed discussion can be found in Appendix A.2.

For the next set of experiments, we fix  $\lambda = T_0$  and run more experiments with varying  $\epsilon \in \{2, 4, 10, 20, 40, 100, 200\}g$ , where the privacy budget is split equally between Step 1 and Step 2 (i.e.,  $\epsilon_1 = \epsilon_2 = \epsilon/2$ .) The bottom two plots in Figure 2 show that the RMSE diminishes as  $\epsilon$  grows, as expected. We continue to observe  $DPSC_{obj}$  outperforming  $DPSC_{out}$  for most  $\epsilon$  values. At  $\epsilon = 10$  for  $DPSC_{out}$ , we see that RMSE on datasets with  $n = 100$  (red and orange lines) is higher than that of databases with  $n = 10$  (blue and green lines). This is consistent with our theoretical analysis in Section C.1.1 that the RMSE of  $DPSC_{out}$  is  $O(\frac{n}{\epsilon})$ . For  $DPSC_{obj}$ , the accuracy bound has an additional dependency on  $T_0$  (Corollary E.2), so the ordering of performance by database size is less clear.

## 4.2 Guidance for Hyperparameter Tuning

Both the  $DPSC_{out}$  and  $DPSC_{obj}$  algorithms require tuning the hyperparameter  $\lambda$ . This parameter plays an important role in determining the amount of noise, since it appears in the sensitivity of  $f$ , i.e.,  $\Delta f = \frac{4T_0\sqrt{8+n}}{\lambda}$ . In theory, the optimal choice of this parameter is recommended to be  $\lambda = O(T_0)$  because the regression coefficient in the objective function is  $\frac{\lambda}{2T_0}$ , and the importance of the regularizer should not diminish as  $T_0$  increases (i.e., as we have more training data points). We confirm this empirically by plotting the post-intervention RMSE as a function of  $\lambda$  (Figure 2, top), which is near-optimal around  $\lambda = T_0$ .

## 5 Conclusion

This paper is the first to propose differentially private versions of the synthetic control algorithm. We provide algorithms based on output perturbation and objective perturbation, and provide formal privacy and accuracy guarantees for each. Our main technical contribution is a novel analysis of the sensitivity of regression in the *transposed* setting, which also impacted our accuracy analysis. To enable practical use of the new tools, we provide a closed-form accuracy bound for both algorithms under distributional assumptions and guidance to practitioners for tuning the parameters of each algorithm.

## Acknowledgements

The first and third authors were supported in part by Novartis AG. The first and second author was supported in part by NSF grant CNS-1942772 (CAREER), a Mozilla Research Grant, a JPMorgan Chase Faculty Research Award, and an Apple Privacy-Preserving Machine Learning Award.

## References

- Abadie, A. (2021). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*, 59(2):391–425.
- Abadie, A., Diamond, A., and Hainmueller, J. (2010). Synthetic control methods for comparative case studies: Estimating the effect of california’s tobacco control program. *Journal of the American statistical Association*, 105(490):493–505.
- Abadie, A., Diamond, A., and Hainmueller, J. (2015). Comparative politics and the synthetic control method. *American Journal of Political Science*, 59(2):495–510.
- Abadie, A. and Gardeazabal, J. (2003). The economic costs of conflict: A case study of the basque country. *American economic review*, 93(1):113–132.
- Abadie, A. and L’Hour, J. (2021). A penalized synthetic control estimator for disaggregated data. *Journal of the American Statistical Association*, pages 1–18.
- Agarwal, A., Misra, V., Shah, D., Shen, D., 1, A., and 2, A. (2021a). Identifying personalized treatment effect from population-level clinical trials data for alzheimer’s patient. *Working Paper*.
- Agarwal, A., Shah, D., Shen, D., and Song, D. (2021b). On robustness of principal component regression. *Journal of the American Statistical Association*, (just-accepted):1–34.
- Amjad, M., Misra, V., Shah, D., and Shen, D. (2019). mrsc: Multi-dimensional robust synthetic control. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(2):1–27.
- Amjad, M., Shah, D., and Shen, D. (2018). Robust synthetic control. *Journal of Machine Learning Research*, 19(22):1–51.
- Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2021). Matrix completion methods for causal panel data models. *Journal of the American Statistical Association*, pages 1–15.
- Athey, S. and Imbens, G. W. (2017). The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives*, 31(2):3–32.
- Bassily, R., Smith, A., and Thakurta, A. (2014). Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pages 464–473. IEEE.
- Ben-Michael, E., Feller, A., and Rothstein, J. (2021). The augmented synthetic control method. *Journal of the American Statistical Association*, (just-accepted):1–34.
- Billingsley, P. (1995). *Measure and probability*.
- Boyd, S. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. (2011). Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3).
- Cummings, R., Ioannidis, S., and Ligett, K. (2015). Truthful linear regression. In *Conference on Learning Theory*, pages 448–483. PMLR.
- Doudchenko, N. and Imbens, G. W. (2016). Balancing, regression, difference-in-differences and synthetic control methods: A synthesis. Technical report, National Bureau of Economic Research.
- Dube, A. and Zipperer, B. (2015). Pooling multiple case studies using synthetic controls: An application to minimum wage policies.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. In *Proceedings of the 3rd Conference on Theory of Cryptography*, TCC ’06, pages 265–284.
- Gökbuget, N., Kelsh, M., Chia, V., Advani, A., Bassan, R., Dombret, H., Doubek, M., Fielding, A. K., Giebel, S., Haddad, V., et al. (2016). Blinatumomab vs historical standard therapy of adult relapsed/refractory acute lymphoblastic leukemia. *Blood cancer journal*, 6(9):e473–e473.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer series in statistics. Springer.
- Horn, R. A. and Johnson, C. R. (2012). *Matrix Analysis*. Cambridge University Press, USA, 2nd edition.
- Kearns, M., Pai, M., Roth, A., and Ullman, J. (2014). Mechanism design in large games: Incentives and privacy. In *Proceedings of the 5th Conference on Innovations in Theoretical Computer Science*, ITCS ’14, pages 403–410. ACM.
- Kifer, D., Smith, A., and Thakurta, A. (2012). Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pages 25–1. JMLR Workshop and Conference Proceedings.
- Kreif, N., Grieve, R., Hangartner, D., Turner, A. J., Nikolova, S., and Sutton, M. (2016). Examination of the synthetic control method for evaluating health policies with multiple treated units. *Health economics*, 25(12):1514–1528.

- Petrone, J. (2018). Roche pays \$1.9 billion for flatiron's army of electronic health record curators. *Nature Biotechnology*, 36(4):289–291.
- Rubin, D. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688–701.
- Thorlund, K., Dron, L., Park, J., and Mills, E. (2020). Synthetic and external controls in clinical trials – a primer for researchers. *Clinical Epidemiology*, Volume 12:457–467.
- Udell, M. and Townsend, A. (2019). Why are big data matrices approximately low rank? *SIAM Journal on Mathematics of Data Science*, 1(1):144–160.
- Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press.
- Wang, D., Ye, M., and Xu, J. (2017). Differentially private empirical risk minimization revisited: Faster and more general. *Advances in Neural Information Processing Systems*, 30.

## A Experimental Settings and Results

In this section, we provide more details about the experiments that were omitted from the main text due to the page limit.

### A.1 Synthetic Data Generation

We use synthetic datasets in our experiments, which enables us to observe the impact of varying the relevant parameters in the data and to match the modeling assumptions of Section 2.1. We create use  $T_0 \in \{10, 100\}$  and  $n \in \{10, 100\}$ , corresponding to both smaller and larger number of donors and observations, and we always use  $T = T_0 + 3$ , meaning that the synthetic control algorithm must predict the next three data points.

The true signals  $M$  and  $\mathbf{m}$  are generated according to a linear model with random slope, formalized as:

$$M_{i,t} = \theta_i t \quad \text{and} \quad m_t = \theta_0 t, \quad \theta_i \in [n], t \in [T],$$

where the  $\theta_i$  are sampled i.i.d. from a truncated Gaussian with mean 4, variance 1, and support  $[3, 5]$ . Elements of the noise terms  $Z$  and  $\mathbf{z}$  are sampled i.i.d. from a truncated Gaussian with mean zero, variance 0.1 and support  $[-1, 1]$ . Following Equation 1, the donor and target data were respectively  $X = M + Z$  and  $\mathbf{y} = \mathbf{m} + \mathbf{z}$ . Figure 3 shows an example synthetic dataset generated in this way, with the donor data in grey and the target in red.

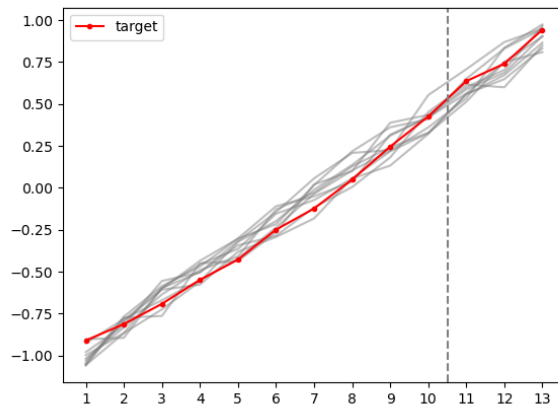


Figure 3: Illustration of example synthetic dataset generated with  $T_0 = 10$  and  $n = 10$ . The target time series is in red, and the donor time series are all in grey.

In each experiment with a fixed  $T_0$  and  $n$ , a single database was generated, and then algorithms were run 500 times on each dataset. We evaluate post-intervention RMSE as the accuracy measure of interest, as in our theoretical results. Error bands in all figures show 95% confidence intervals, taken over the randomness in the algorithms.

### A.2 Optimizing regularization parameter $\lambda$

The first question we aim to address in our experiments is the impact of the parameter  $\lambda$  on performance, and guidance for analysts in their choice of optimal  $\lambda$ . In our first set of experiments, we fixed  $\epsilon_1 = \epsilon_2 = 50$ ,  $T_0 = 10$ , and  $n = 10$ —other values of  $\epsilon$  and  $(T_0, n)$  are considered respectively in Sections A.3 and ??—and empirically measured pre- and post-intervention RMSE as a function of  $\lambda$ .

Figure 4 shows the post-intervention RMSE of  $DPSC_{out}$ ,  $DPSC_{obj}$ , and non-private synthetic control as a function of  $\lambda$ , for values  $\lambda \in \{1, 2, 5, 10, 20, 50, 100, 200, 500, 1000, 2000, 5000\}$ .

We observe that the performance of the three methods converges as  $\lambda$  grows large, but that the RMSE of the private methods is larger than that of the non-private method for smaller  $\lambda \leq 20$ , with Objective Perturbation substantially outperforming Output Perturbation.

To aid the analyst in choosing an optimal  $\lambda$ , we observe that the RMSE of DPSC is minimized around  $\lambda = T_0$  for all four datasets (see Figure 2 in Section 4.1). This is consistent with our theoretical recommendations that  $\lambda$  should be  $O(T_0)$ .

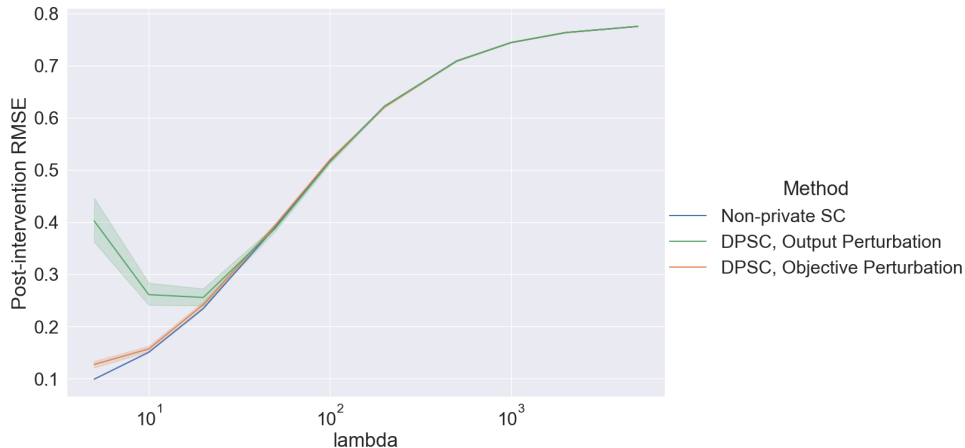


Figure 4: Behavior of post-intervention RMSE over  $\lambda$ , tested on a synthetic dataset with  $T_0 = 10, n = 10$  for the synthetic control methods of non-private SC (blue),  $DPSC_{out}$  (green), and  $DPSC_{obj}$  (orange).

The U-shape has a natural theoretical explanation: smaller  $\lambda$  increases sensitivity and thus privacy noise and RMSE, while larger  $\lambda$  increases the weight of the regularization term in the loss function, which will cause all three regularized methods to converge to each other.

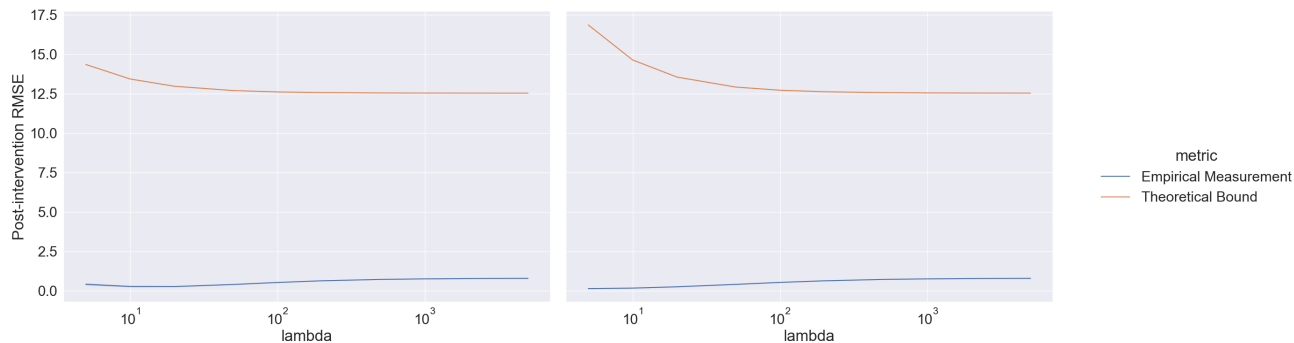


Figure 5: Comparison of post-intervention RMSE in theory versus in practice, using  $DPSC_{out}$  (left) and  $DPSC_{obj}$  (right) on a dataset of size  $n = 10, T_0 = 10$ .

Figure 5 compares the empirical post-intervention RMSE of  $DPSC_{out}$  and  $DPSC_{obj}$  with the theoretical guarantees of Theorem 3.2 and 3.5 instantiated with parameters of our experiments. We observe that the empirical performance of both algorithms is dramatically better than the theoretical bounds may suggest. We also observe that  $DPSC_{obj}$  (right) has lower empirical error than  $DPSC_{out}$  (left), which diverges from our theoretical predictions. This suggests potential room for theoretical improvements in the bound of Theorem 3.5.

### A.3 Effect of privacy parameter $\epsilon$

Next, we address the effect of  $\epsilon$  in the performance of both  $DPSC_{out}$  and  $DPSC_{obj}$ . In these experiments, we use  $\lambda = T_0$  based on the findings in Section A.2 and consider overall privacy budget  $\epsilon = \epsilon_1 + \epsilon_2$  with  $\epsilon_1 = \epsilon_2 = \epsilon/2$ . That is, the privacy budget is split evenly between the regression and projection steps in both algorithms. Results are presented for  $\epsilon \in \{2, 4, 10, 20, 40, 100, 200\}$ ; stronger privacy guarantees (i.e.,  $\epsilon = 2$ ) were tested but excluded from the plots due to substantially higher RMSE values.

Figure 6 shows the post-intervention RMSE of  $DPSC_{out}$  and  $DPSC_{obj}$ . As is to be expected, error diminishes with larger  $\epsilon$ . We also continue to observe  $DPSC_{obj}$  outperforming  $DPSC_{out}$  for most  $\epsilon$  values, as in Section A.2.  $DPSC_{out}$  performs slightly better than  $DPSC_{obj}$  at  $\epsilon = 2$  in this dataset ( $T_0 = 10$  and  $n = 10$ ); however, it is not the case for all datasets (See Figure 2 in Section 4.1).

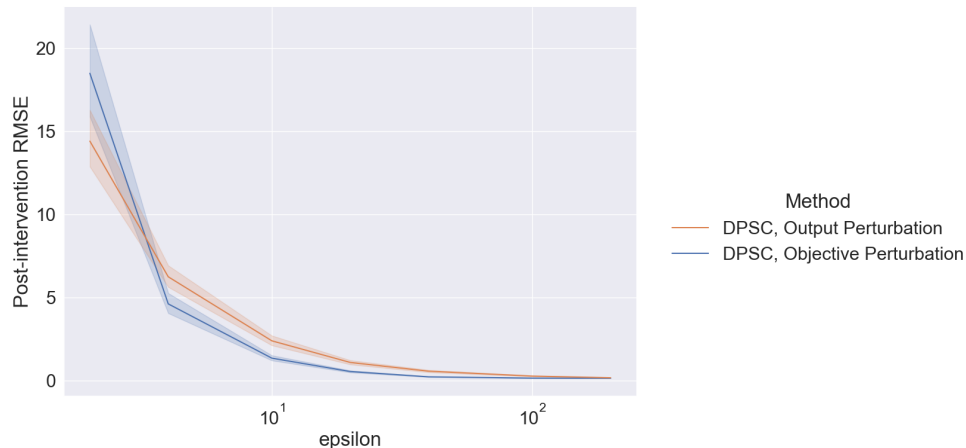


Figure 6: Post-intervention RMSE of  $DPSC_{out}$  (blue), and  $DPSC_{obj}(red)$  for varying  $\epsilon$ , tested on a dataset with  $T_0 = 10$  and  $n = 10$ .

For epsilon-regimes that are closer to the values chosen in practice (i.e.,  $\epsilon \approx 4$ ), the empirical RMSE was too high for practical use. We suggest a few methods to remedy this in future work. First, a rejection sampling step can be introduced between the learning and projection steps of each algorithm that compares the noisy  $\mathbf{f}^{out}$  and the original  $\mathbf{f}^{reg}$ . This step must also be done differentially privately to maintain the overall privacy guarantee. Additionally, our experiments only considered pure-DP with  $\delta = 0$ ; relaxing to approximate-DP with  $\delta > 0$  would likely yield lower RMSE.

## B Privacy Guarantees of $DPSC_{out}$

In this section, we will prove Theorem 3.1, that DPSC is  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private. The proof relies on the privacy of  $\mathbf{f}^{out}$  in the learning phase, and then  $\tilde{X}_{post}$  in the prediction phase. At a high level,  $\mathbf{f}^{out}$  is  $\epsilon_1$ -DP through a (non-trivial) application of the Output Perturbation algorithm of Chaudhuri et al. (2011). The non-triviality comes from the vertical regression used in synthetic control, rather than the horizontal regression classically used in empirical risk minimization (as illustrated in Figure 1), which requires novel sensitivity analysis of the function  $\mathbf{f}^{reg}$ . In the prediction phase, we must show that sufficient noise is added to ensure  $\tilde{X}_{post}$  is an  $\epsilon_2$ -DP version of  $X_{post}$ . Then privacy of  $\mathbf{y}^{out}$  comes from post-processing and composition of these two private estimates.

### B.1 Privacy of $\mathbf{f}^{out}$

Let us begin by proving that  $\mathbf{f}^{out}$  is  $\epsilon_1$ -DP.

**Theorem B.1** *Step 1 of Algorithm 2 that computes  $\mathbf{f}^{out}$  is  $(\epsilon_1, 0)$ -differentially private.*

It might seem that Theorem 3.1 should follow immediately from the privacy guarantees of Output Perturbation in Chaudhuri et al. (2011). Indeed, Theorem 6 of Chaudhuri et al. (2011) states that a similar algorithm is  $(\epsilon, 0)$ -DP under certain technical conditions. However, the proof of this result relies on sensitivity analysis of classical empirical risk minimization (see Corollary 8 of Chaudhuri et al. (2011)) which does not hold in the synthetic control setting. The crux of the difference comes from the vertical regression (i.e., along the columns) of synthetic control as illustrated in Figure 1, while privacy must still be maintained along the rows. Thus the sensitivity of  $\mathbf{f}^{reg}$  to a change in a single *donor row* is fundamentally different from the sensitivity in a standard empirical risk minimization setting. See Remark B.5 for a more technical exploration of this difference. Additionally, while the ERM framework of Chaudhuri et al. (2011) is fully general, their results (including Theorem 6 and Corollary 8) apply only to the problem setting of binary classification via logistic regression, by assuming a specific loss function  $L$  in the analysis.

Instead, we prove Theorem B.1 primarily using first-principles (i.e., direct sensitivity analysis and the Laplace Mechanism of Dwork et al. (2006), which also underpins the results of Chaudhuri et al. (2011)) starting with Lemma 3.7. The proof of Lemma 3.7 and Theorem B.1 will be augmented with one intermediate result for output perturbation from Chaudhuri et al. (2011) that does apply to our setting, and one fact from Cummings et al. (2015), which extended the binary classification result of Chaudhuri et al. (2011) to the Ridge regression loss function that we use.

**Lemma 3.7** The  $\ell_2$  sensitivity of  $\mathbf{f}^{reg}$  is  $\Delta \mathbf{f}^{reg} \leq \frac{4T_0\sqrt{8+n}}{\lambda}$ .

To prove Lemma 3.7, we will first use the following lemma from Chaudhuri et al. (2011), which bounds the sensitivity of  $\mathbf{f}^{reg}$  as a function of the strong convexity parameter of the loss function  $L$ .

**Lemma B.2 (Chaudhuri et al. (2011), Lemma 7)** Let  $G(\mathbf{f})$  and  $g(\mathbf{f})$  be two vector-valued functions, which are continuous, and differentiable at all points. Moreover, let  $G(\mathbf{f})$  and  $G(\mathbf{f}) + g(\mathbf{f})$  be  $\lambda$ -strongly convex. If  $\mathbf{f}_1 = \arg \min_{\mathbf{f}} G(\mathbf{f})$  and  $\mathbf{f}_2 = \arg \min_{\mathbf{f}} G(\mathbf{f}) + g(\mathbf{f})$ , then

$$\|\mathbf{f}_1 - \mathbf{f}_2\|_2 \leq \frac{1}{\lambda} \max_{\mathbf{f}} \|\nabla g(\mathbf{f})\|_2.$$

We instantiate this lemma by defining

$$G(\mathbf{f}) = L(\mathbf{f}, D) \quad \text{and} \quad g(\mathbf{f}) = L(\mathbf{f}, D') - L(\mathbf{f}, D), \quad (5)$$

for two arbitrarily neighboring databases  $D, D'$  and defining the following two maximizers:

$$\mathbf{f}_1 = \arg \min L(\mathbf{f}, D) = \arg \min G(\mathbf{f}) \quad \text{and} \quad \mathbf{f}_2 = \arg \min L(\mathbf{f}, D') = \arg \min G(\mathbf{f}) + g(\mathbf{f}).$$

Then,

$$\Delta \mathbf{f}^{reg} = \max_{D, D' \text{ neighbors}} \|\mathbf{f}_1 - \mathbf{f}_2\|_2.$$

To apply Lemma B.2, we must show that  $G(\mathbf{f})$  and  $g(\mathbf{f})$  are continuous and differentiable.  $G(\mathbf{f})$  is simply the Ridge regression loss function, which is known to be continuous and differentiable Hastie et al. (2009). Since  $g(\mathbf{f})$  is the difference between two continuous and differentiable functions, then it is also continuous and differentiable Boyd and Vandenberghe (2004). We must also show strong convexity of  $G(\mathbf{f})$  and  $G(\mathbf{f}) + g(\mathbf{f})$ . The following lemma from Cummings et al. (2015) immediately gives that these two functions are both  $\lambda$ -strongly convex.

**Lemma B.3 (Cummings et al. (2015), Lemma 32)** The Ridge regression loss function with regularizer  $\frac{\lambda}{2T_0}$  is  $\frac{\lambda}{T_0}$ -strongly convex.

Thus by Lemma B.2, the sensitivity  $\Delta \mathbf{f}^{reg} = \max_{D, D' \text{ neighbors}} \|\mathbf{f}_1 - \mathbf{f}_2\|_2 \leq \frac{T_0}{\lambda} \max_{\mathbf{f}} \|\nabla g(\mathbf{f})\|_2$ . All that remains is to bound  $\|\nabla g(\mathbf{f})\|_2$ . A proof of the following lemma is deferred to Appendix F.

**Lemma B.4** Let  $g(\mathbf{f}) = L(\mathbf{f}, D') - L(\mathbf{f}, D)$  for two arbitrarily neighboring databases  $D, D'$ . Then,

$$\max_{\mathbf{f}} \|\nabla g(\mathbf{f})\|_2 \leq \frac{2}{\lambda} \sqrt{8+n}.$$

**Remark B.5** If we were instead considering simple linear regression in the classical setting (i.e., as in Chaudhuri et al. (2011)) using  $T_0$  data points with  $n$  dimensional features,  $g(\mathbf{f})$  would only contain one term in the error, namely, the one data point  $(\mathbf{x}_i, y_i)$  that differed across two neighboring databases. This yields

$$g(\mathbf{f}) = \frac{1}{T_0} ((\mathbf{x}'_i \quad \mathbf{x}_i)^\top \mathbf{f} - (y'_i \quad y_i))^2$$

with gradient

$$\nabla g(\mathbf{f}) = \frac{2}{T_0} ((\mathbf{x}'_i \quad \mathbf{x}_i)(\mathbf{x}'_i \quad \mathbf{x}_i)^\top \mathbf{f} - (y'_i \quad y_i)(\mathbf{x}'_i \quad \mathbf{x}_i)),$$

which can be bounded by  $O(\frac{1}{T_0})$ . This result does not depend on the dimension of the features ( $n$ ) and only depends on the number of data points ( $T_0$ ).

However, in synthetic control, terms do not cancel as neatly across neighboring databases, and instead,

$$g(\mathbf{f}) = \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} \mathbf{x}_t^\top \mathbf{f} - y_t + \mathbf{x}_t^\top \mathbf{f} - x_{i,t} f_i + x'_{i,t} f_i - y_t - (x'_{i,t} - x_{i,t}) f_i.$$

Through a more involved analysis of this expression, we get the bound of Lemma B.4, which depends on  $n$ , rather than  $T_0$ .

Using these lemmas, we can now bound the sensitivity of our query, to complete the proof of Lemma 3.7.

$$\Delta \mathbf{f}^{reg} = \max_{\mathcal{D}, \mathcal{D}^o \text{ neighbors}} \|\mathbf{f}(\mathcal{D}) - \mathbf{f}(\mathcal{D}^o)\|_2 = \max_{\mathcal{D}, \mathcal{D}^o} \|\mathbf{f}_1 - \mathbf{f}_2\|_2 = \frac{4T_0 \rho}{\lambda} \frac{1}{8+n}. \quad (6)$$

Theorem B.1 then follows from the privacy guarantee of the high-dimensional Laplace Mechanism instantiated with the appropriate sensitivity value.

### B.1.1 Dependence on $n$

One might wonder whether the asymptotic dependence on  $n$  and  $T_0$  in the sensitivity is necessary. In practice, one should set  $\lambda = O(T_0)$  (as discussed in greater detail in Section C.1.1), so the dependence on  $T_0$  will not affect the accuracy of the algorithm. However, as we show next in Lemma 3.8, the dependence on  $n$  is asymptotically tight.

**Lemma 3.8** *The  $\ell_2$  sensitivity of  $\mathbf{f}^{reg}$  is  $\Delta \mathbf{f}^{reg} = \Omega(\frac{\rho}{\sqrt{n}})$ .*

Consider two neighboring databases  $(X, \mathbf{y})$  and  $(X', \mathbf{y})$ , where  $\mathbf{y} = \mathbf{1} \in \mathbb{R}^{T_0}$ ,  $X \in \mathbb{R}^{n \times T_0}$  has all entries  $1/n$ , except the first row, which is all 1s. Neighboring database  $X'$  differs from  $X$  only in the first row, which is instead all 0s, and all other entries and  $1/n$ . The dimensions in this example are chosen to be  $T_0 = n$ , and we choose  $\lambda = 2T_0$ , so that the regularization coefficient is 1.

Computing the minimizers of the loss functions under each neighboring database using the closed-form expression yields  $\mathbf{f}^{reg} = (X X^\top + I)^{-1} X \mathbf{y}$  with the first coordinate equal to  $\frac{n^2}{n^2 + 2n - 1}$ , and all other coordinates are  $\frac{n}{n^2 + 2n - 1}$ , and  $\mathbf{f}^{reg'} = (X' X'^\top + I)^{-1} X' \mathbf{y}$  with first coordinate 0 and all other coordinates  $\frac{n}{1 - 2n}$ . This yields  $\ell_2$  difference of,

$$\|\mathbf{f}^{reg} - \mathbf{f}^{reg'}\|_2 = \sqrt{\frac{n^2}{n^2 + 2n - 1}^2 + (n - 1) \frac{n^3}{(n^2 + 2n - 1)(1 - 2n)}^2} = \Theta(\frac{\rho}{\sqrt{n}}).$$

Since we have a pair of neighboring databases with  $\ell_2$  distance in their output of  $\Theta(\frac{\rho}{\sqrt{n}})$ , then the sensitivity of  $\mathbf{f}^{reg}$  cannot be  $o(\frac{\rho}{\sqrt{n}})$ .

**Remark B.6** *We note that while the example in Lemma 3.8 is mathematically valid, such a degenerate case where all the donors are identical except for one person and the (exact) rank of the donor matrix is 1 is unlikely to happen in practical settings. This suggests that with additional domain knowledge on the selection criteria for donors, practitioners may be able to reduce the sensitivity and thus add less noise for privacy in special restricted cases of interest.*

### B.2 Privacy of $\tilde{X}_{post}$ and $\mathbf{y}^{out}$

Next we move to privacy of  $\tilde{X}_{post}$  and its role in ensuring privacy of  $\mathbf{y}^{out}$ .

**Lemma 3.9** *The computation of  $\tilde{X}_{post}$  in Step 2 of Algorithm 2 is  $(\epsilon_2, 0)$ -differentially private.*

$\tilde{X}_{post}$  is privatized through a simple application of the Laplace Mechanism of Dwork et al. (2006). Thus to prove Lemma 3.9, we need only to bound the sensitivity of  $X_{post}$  to show that the algorithm adds sufficient noise. We first note that the Frobenius norm of a matrix  $X \in \mathbb{R}^{n \times (T - T_0)}$  is equal to the  $\ell_2$  norm of the equivalent flattened vector  $X \in \mathbb{R}^{n(T - T_0)}$  Horn and Johnson (2012). Thus implementing the matrix-valued Laplace Mechanism with noise parameter calibrated to the  $\ell_2$  sensitivity of the flattened matrix-valued query over  $\epsilon$  will ensure  $(\epsilon, 0)$ -differential privacy.

**Lemma B.7** *The  $\ell_2$  sensitivity of flattened  $X_{post}$  is  $2 \sqrt{\frac{\rho}{(T - T_0)}}$ .*

Changing one donor unit in  $X_{post}$  can change at most  $T - T_0$  entries in the matrix. Since all entries in  $X_{post}$  are bounded in  $[-1, 1]$ , each data point can change by at most 2 between two neighboring databases. Thus viewing  $X_{post}$  as a flattened matrix, this will change the  $\ell_2$ -norm of  $X_{post}$  by at most  $2 \sqrt{\frac{\rho}{(T - T_0)}}$ .

Finally, we can combine Theorem B.1 and Lemma 3.9 to complete the proof of Theorem 3.1. The estimates  $\mathbf{f}^{out}$  and  $\tilde{X}_{post}$  are together  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private by DP composition, and then  $\mathbf{y}^{out}$  is  $(\epsilon_1 + \epsilon_2, 0)$ -differentially private by post-processing. We note that if one wanted to publish  $\mathbf{f}^{out}$ , this would not incur any additional privacy loss.



## C Accuracy Guarantees of $DPSC_{out}$

In this section we will analyze the accuracy of  $DPSC_{out}$ . We first prove Theorem 3.2, restated below for convenience.

**Theorem 3.2** *The estimator  $\mathbf{y}^{out}$  output by Algorithm 2 satisfies:*

$$RMSE(\mathbf{y}^{out}) \leq \frac{\|M_{post}\|_2}{\sqrt{T-T_0}} \mathbb{E}[\|\mathbf{f}^{reg} - \mathbf{f}\|_2] + \frac{4T_0 \frac{\rho}{8+n}}{\lambda\epsilon_1} + \frac{\rho}{n\sigma^2} + \frac{\rho \frac{1}{2}}{\epsilon_2} + \frac{\rho}{n\psi} + \frac{4T_0 \frac{\rho}{8+n}}{\lambda\epsilon_1},$$

where  $\|\mathbf{f}^{reg}\|_\infty \leq \psi$  for some  $\psi > 0$ , and  $RMSE$  is the root mean squared error of the estimator, defined as  $RMSE(\mathbf{y}^{out}) = \frac{1}{\sqrt{T-T_0}} \mathbb{E}[\|\mathbf{y}^{out} - \mathbf{m}_{post}\|_2]$ .

This theorem gives bounds on the predicted post-intervention target vector  $\mathbf{y}^{out}$ , as measured by RMSE. This result is stated in full generality with respect to the distribution of data and the latent variables, and thus the bound depends on terms such as  $\|M_{post}\|_{2,2}$  and  $\mathbb{E}[\|\mathbf{f}^{reg} - \mathbf{f}\|_2]$ . This is consistent with comparable bounds on the RMSE of robust synthetic control Amjad et al. (2018) which also depended on these terms (although the stated bounds of Amjad et al. (2018) suppress dependence on  $n$ ). Section C.1 provides a proof of this main result.

Analysts may still wonder about the full asymptotic performance of  $DPSC_{out}$  algorithm. To this end, in Section C.2, we additionally derive closed-form bounds for these distribution-dependent terms (under some mild assumptions). We present Corollary C.5, which gives a bound on RMSE of  $\mathbf{y}^{out}$  that depends only on input parameters of the algorithm and the model.

### C.1 Accuracy of post-intervention prediction $\mathbf{y}^{out}$

We will prove Theorem 3.2 by showing that the prediction vector  $\mathbf{y}^{out}$  output by  $DPSC_{out}$  in Algorithm 2 is close to the true values, as measured by Root Mean Squared Error (RMSE), defined as follows:

$$RMSE(\mathbf{y}^{out}) = \frac{1}{\sqrt{T-T_0}} \mathbb{E}[\|\mathbf{y}^{out} - \mathbf{m}_{post}\|_2]. \quad (7)$$

We note that while it may seem most natural to bound the difference between  $\mathbf{y}^{out}$  and  $\mathbf{y}_{post}$ , we instead use  $\mathbf{m}_{post}$  for two reasons. Firstly,  $\mathbf{y}_{post}$  may not even match  $\mathbf{y}_{pre}$  due to the intervention. Secondly,  $\mathbf{m}_{post}$  captures the true signal that we are trying to estimate, which is the counterfactual outcome without the intervention.

We begin by bounding the expected  $\ell_2$  difference between  $\mathbf{y}^{out}$  and  $\mathbf{m}_{post}$ . Using the fact that

$$\mathbf{y}^{out} = \tilde{X}_{post}^\top \mathbf{f}^{out} = (X_{post}^\top + W^\top)(\mathbf{f}^{reg} + \mathbf{v}),$$

and that  $X_{post} = M_{post} + Z$  and  $\mathbf{m} = M_{post}^\top \mathbf{f}$  (by Equation (2)), we can expand the expectation as follows:

$$\begin{aligned} \mathbb{E}[\|\mathbf{y}^{out} - \mathbf{m}_{post}\|_2] &= \mathbb{E}[\|(X_{post}^\top + W^\top)(\mathbf{f}^{reg} + \mathbf{v}) - M_{post}^\top \mathbf{f}\|_2] \\ &= \mathbb{E}[\|(M_{post}^\top + Z^\top + W^\top)(\mathbf{f}^{reg} + \mathbf{v}) - M_{post}^\top \mathbf{f}\|_2] \\ &= \mathbb{E}[\|M_{post}^\top (\mathbf{f}^{reg} - \mathbf{f})\|_2 + \|Z^\top + W^\top\|_2 \|\mathbf{f}^{reg}\|_2 + \|M_{post}^\top + Z^\top + W^\top\|_2 \|\mathbf{v}\|_2] \\ &= \mathbb{E}[\|M_{post}^\top (\mathbf{f}^{reg} - \mathbf{f})\|_2] + \mathbb{E}[\|Z^\top + W^\top\|_2 \|\mathbf{f}^{reg}\|_2] + \mathbb{E}[\|M_{post}^\top + Z^\top + W^\top\|_2 \|\mathbf{v}\|_2] \end{aligned} \quad (8)$$

We next proceed to bound each of the terms in Equation (8) separately, making use of the following submultiplicative norm property, which holds for any matrix  $A$  and vector  $\mathbf{x}$ :

$$\|A\mathbf{x}\|_2 \leq \|A\|_2 \|\mathbf{x}\|_2 \leq \|A\|_F \|\mathbf{x}\|_2, \quad (9)$$

where  $\|A\|_2 = \|A\|_{2,2}$  is the spectral norm of  $A$ ,  $\|A\|_F$  is the Frobenius norm of  $A$ , and  $\|\mathbf{x}\|_2$  is the  $\ell_2$  norm of  $\mathbf{x}$ .

We also know the distribution of the norms of noise terms  $\mathbf{v}$  and  $W$  that were added to preserve privacy, because they were constructed explicitly within Algorithm 2:

$$\mathbb{E}[\|\mathbf{v}\|_2] = \frac{4T_0 \frac{\rho}{8+n}}{\lambda\epsilon_1} \quad \text{and} \quad \mathbb{E}[\|W\|_F] = b = \frac{2 \frac{\rho}{T-T_0}}{\epsilon_2}. \quad (10)$$

Using these facts, we can obtain bounds for the three terms in (8). A complete proof of Lemma C.1 can be found in Appendix F.2.

**Lemma C.1** *The three terms in Equation (8) can be bounded as follows:*

$$\begin{aligned} & \mathbb{E}[jjM_{post}^\top(\mathbf{f}^{reg} - \mathbf{f})jj_2] \quad jjM_{post}jj_{2,2} \quad \mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2], \\ & \mathbb{E}[jj(Z^\top + W^\top)\mathbf{f}^{reg}jj_2] \quad \rho_{n\psi} \quad \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho}{T - T_0} \quad \epsilon_2, \text{ and} \\ & \mathbb{E}[jj(M_{post}^\top + Z^\top + W^\top)\mathbf{v}jj_2] \quad jjM_{post}jj_{2,2} + \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho}{T - T_0} \quad \frac{4T_0}{\epsilon_2} \frac{\rho}{8+n} \frac{1}{\lambda\epsilon_1}. \end{aligned}$$

Applying the bounds of Lemma C.1 to Equation (8) yields,

$$\begin{aligned} \mathbb{E}[jj\mathbf{y}^{out} - \mathbf{m}jj_2] & \quad jjM_{post}jj_{2,2} \quad \mathbb{E}[\mathbf{f}^{reg} - \mathbf{f}jj_2] + \rho_{n\psi} \quad \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho}{T - T_0} \quad \epsilon_2 \\ & \quad + \quad jjM_{post}jj_{2,2} + \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho}{T - T_0} \quad \frac{4T_0}{\epsilon_2} \frac{\rho}{8+n} \frac{1}{\lambda\epsilon_1} \\ & \quad jjM_{post}jj_{2,2} \quad \mathbb{E}[\mathbf{f}^{reg} - \mathbf{f}jj_2] + \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{8+n} \\ & \quad + \quad \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho}{T - T_0} \quad \rho_{n\psi} + \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{8+n} \end{aligned}$$

Combining this with Equation (7) gives the desired bound for Theorem 3.2:

$$RMSE(\mathbf{y}^{out}) \quad \frac{jjM_{post}jj_{2,2}}{\rho \frac{\rho}{T - T_0}} \quad \mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2] + \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{8+n} \quad + \quad \rho_{n\sigma^2} + \frac{\rho_{\frac{\rho}{2}}}{\epsilon_2} \quad \rho_{n\psi} + \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{8+n}.$$

### C.1.1 Cost of privacy in synthetic control

To understand the additional error incurred due to privacy, compare the bound of Theorem 3.1 to the RMSE of the equivalent non-private prediction,  $\mathbf{y}^{reg} = X_{post}^\top \mathbf{f}^{reg}$ .

$$\begin{aligned} RMSE(\mathbf{y}^{reg}) &= \rho \frac{1}{T - T_0} \mathbb{E}[jjX_{post}^\top \mathbf{f}^{reg} - \mathbf{m}jj_2] \\ &= \rho \frac{1}{T - T_0} \mathbb{E}[jj(M_{post}^\top + Z_{post}^\top)\mathbf{f}^{reg} - M_{post}^\top \mathbf{f}jj_2] \\ &= \rho \frac{1}{T - T_0} \mathbb{E}[jjM_{post}^\top(\mathbf{f}^{reg} - \mathbf{f})jj_2 + jjZ_{post}^\top \mathbf{f}^{reg}jj_2] \\ &= \frac{jjM_{post}jj_{2,2}}{\rho \frac{\rho}{T - T_0}} (\mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2]) + \rho_{n\psi} \quad \rho_{n\sigma^2} \end{aligned} \quad (11)$$

Lemma C.6 in the next section shows that  $\mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2] = O(\rho_{\bar{n}})$ . Then the first term of Equation (11) can be easily bounded using the following fact,

$$jjM_{post}jj_{2,2} \quad jjM_{post}jj_F \quad \rho \frac{\rho}{n(T - T_0)},$$

so  $\frac{jjM_{post}jj_{2,2}}{\sqrt{T - T_0}} \rho_{\bar{n}}$ . Thus we see that  $RMSE(\mathbf{y}^{reg}) = O(n)$ .

Comparing Equation (11) with the bound on  $RMSE(\mathbf{y}^{out})$  in Theorem 3.1, we observe that the additional terms induced by privacy are:

$$\frac{jjM_{post}jj_{2,2}}{\rho \frac{\rho}{T - T_0}} \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{8+n} + \frac{4T_0}{\lambda\epsilon_1} \frac{\rho}{(8+n)n\sigma^2} + \frac{\rho_{2n\psi}}{\epsilon_2} + \frac{4T_0}{\lambda\epsilon_1\epsilon_2} \frac{\rho}{2(8+n)}. \quad (12)$$

Then, using the fact that  $\frac{jjM_{post}jj_{2,2}}{\sqrt{T - T_0}} \rho_{\bar{n}}$  and setting  $\epsilon := \epsilon_1 = \epsilon_2$  and  $\lambda = O(T_0)$ , Equation (12) can be bounded by,

$$\frac{4T_0}{\lambda\epsilon} \frac{\rho}{(8+n)n} + \frac{4T_0}{\lambda\epsilon} \frac{\rho}{(8+n)n\sigma^2} + \frac{\rho_{2n\psi}}{\epsilon} + \frac{4T_0}{\lambda\epsilon^2} \frac{\rho}{2(8+n)} = O \left( \frac{n}{\epsilon} + \frac{\rho_{\bar{n}}}{\epsilon^2} \right) = O \left( \frac{n}{\epsilon} \right) \quad \text{for } \epsilon \geq 1/\rho_{\bar{n}}.$$

Thus we conclude that the cost of privacy in the  $DPSC_{out}$  algorithm is at most a factor of  $O(\frac{1}{\epsilon})$ . The restriction to  $\epsilon \leq \frac{1}{\sqrt{n}}$  is consistent with standard practice in both theoretical and practical deployments of differential privacy, and thus is effectively without loss.

## C.2 Closed-form bound on RMSE of Output Perturbation

In this section, we impose assumptions on the underlying data distribution to extend Theorem 3.2 to provide an explicit closed-form bound on the RMSE. Throughout this section, we make the following three mild assumptions of the distribution of  $X$ , which are required to achieve this closed-form expression:

**Assumption C.2**  $X_{pre}$  takes values in a  $k$ -dimensional subspace  $E$  for some small  $k \leq \min\{fn, T_0\}$ .

**Assumption C.3** The distribution of  $X_{pre}$  over  $E$  is isotropic, hence the covariance matrix  $Cov(X_{pre}) = \Sigma = P_E$  where  $P_E$  is an orthogonal projection matrix onto  $E$ .

**Assumption C.4** The distribution of  $\mathbf{x}_t \in \mathbb{R}^n$  is supported in some centered Euclidean ball with radius  $O(\frac{\rho}{\sqrt{k}})$ .

These assumptions are only slightly stronger than those commonly made in theory Amjad et al. (2018, 2019) and that typically hold in practice Udell and Townsend (2019). The first assumption means that  $X_{pre}$  is low rank. Assuming  $X_{pre}$  to be *approximately* low rank is a common practice in synthetic control literature Amjad et al. (2018, 2019). Indeed, most large matrices in practice are approximately low-rank Udell and Townsend (2019). Hence, we only further assume that it is *exactly* rank  $k$  for some small  $k$ . The second assumption allows us to apply useful mathematical properties:  $\|P_E\|_2 = 1$  and  $\text{trace}(P_E) = k$ . Then,  $E[X_{pre}X_{pre}^\top] = \text{trace}(P_E) = k$  and, using Markov's inequality, we know that most of the distribution mass should be within a ball of radius  $\frac{\rho}{m}$  for  $m = O(k)$ . Hence, the third assumption asserts that not *most* but *all* the probability mass should lie within that ball, i.e.,  $\|X\|_{2,2} = O(\frac{\rho}{\sqrt{k}})$  almost surely.

Corollary C.5 provides a closed-form bound on the RMSE of  $\mathbf{y}^{out}$  under these assumptions.

**Corollary C.5** If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \geq (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t^2}$ , if  $T_0 \geq C(t/\xi)^2 k \log n$ , we have

$$RMSE(\mathbf{y}^{out}) \leq \frac{\rho}{n} \frac{(\frac{\rho}{2n\sigma^2} + \frac{\rho}{2n\sigma^2 s^2})T_0 + \frac{\lambda}{2T_0}}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}} + \frac{4T_0 \frac{\rho}{8+n}}{\lambda \epsilon_1} + \frac{\rho}{n\sigma^2} + \frac{\rho}{\epsilon_2} \leq \frac{\rho}{n} \psi + \frac{4T_0 \frac{\rho}{8+n}}{\lambda \epsilon_1}.$$

To derive Corollary C.5 from Theorem 3.2, we only need to derive and apply bounds on  $\|M_{post}\|_{2,2}$  and  $E[\|f^{reg} - f\|_2]$ . As we did before, we bound the first term using

$$\|M_{post}\|_{2,2} \leq \|M_{post}\|_F \leq \frac{\rho}{n(T - T_0)},$$

and thus  $\frac{\|M_{post}\|_{2,2}}{\sqrt{T - T_0}} \leq \frac{\rho}{n}$ . Therefore, the key step is to bound  $E[\|f^{reg} - f\|_2]$ . The following lemma provides the required bound on this term to prove Corollary C.5. The remainder of this section will be devoted to providing a proof sketch for Lemma C.6. A full proof is presented in Appendix F.3.

**Lemma C.6** Let  $\mathbf{f}^{reg} = (X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1}X_{pre}\mathbf{y}_{pre}$  be the Ridge regression coefficients and let  $\mathbf{f}$  be the true coefficients. If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \geq (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t^2}$ , if  $T_0 \geq C(t/\xi)^2 k \log n$ , we have,

$$E[\|f^{reg} - f\|_2] \leq \frac{(\frac{\rho}{2n\sigma^2} + \frac{\rho}{2n\sigma^2 s^2})T_0 + \frac{\lambda}{2T_0}}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}.$$

[Proof sketch of Lemma C.6.] First we can expand  $E[\|f^{reg} - f\|_2]$ :

$$\begin{aligned} E[\|f^{reg} - f\|_2] &= E[\|f^{reg} - E[f^{reg}] + E[f^{reg}] - f\|_2] \\ &= E[\|f^{reg} - E[f^{reg}]\|_2] + E[\|E[f^{reg}] - f\|_2] \\ &= E[\|f^{reg} - E[f^{reg}]\|_2] + E[\|Bias(\mathbf{f}^{reg})\|_2]. \end{aligned} \tag{13}$$

We can bound these two terms separately as:

$$\begin{aligned} \text{Bias}(\mathbf{f}^{reg}) &= \frac{\lambda}{2T_0} \mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2} \quad \text{and} \\ \mathbb{E}[\mathbb{J}\mathbf{f}^{reg} \quad \mathbf{f}j_{2,2}] &= \mathbb{E}[\mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2} \quad \mathbb{J}X_{pre}\mathbf{z} \quad X_{pre}Z^\top \mathbf{f}j_{2,2}]. \end{aligned}$$

This can be combined back with Equation (13) to yield,

$$\mathbb{E}[\mathbb{J}\mathbf{f}^{reg} \quad \mathbf{f}j_{2,2}] = \mathbb{E}[\mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2} \quad (\mathbb{J}X_{pre}\mathbf{z} \quad X_{pre}Z^\top \mathbf{f}j_{2,2} + \frac{\lambda}{2T_0})]. \quad (14)$$

Next, we use our assumptions on the data distribution to prove the following lemma about  $\mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2}$ .

**Lemma C.7** *If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \geq (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t^2}$  and  $T_0 \geq C(t/\xi)^2 k \log n$ , we have*

$$\mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2} \leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}.$$

To prove Lemma C.7, we use the following lemma about concentration of random matrices.

**Lemma C.8 (Corollary 5.52 of Vershynin (2010))** *Consider a distribution in  $\mathbb{R}^n$  with covariance matrix  $\Sigma$ , and supported in some centered Euclidean ball whose radius we denote  $\sqrt{m}$ . Let  $T_0$  be the number of samples and define the sample covariance matrix  $\Sigma_{T_0} = \frac{1}{T_0}XX^\top$ . Let  $\xi \geq (0, 1)$  and  $t \geq 1$ . Then with probability at least  $1 - n^{-t^2}$ , one has,*

$$\text{If } T_0 \geq C(t/\xi)^2 \mathbb{J}\Sigma\mathbb{J}j_{2,2}^{-1} m \log n \text{ then } \mathbb{J}\Sigma_{T_0} \quad \Sigma\mathbb{J}j_{2,2} \quad \xi\mathbb{J}\Sigma\mathbb{J}j_{2,2},$$

where  $C$  is an absolute constant.

We instantiate Lemma C.8 using our assumptions that  $\mathbb{J}\Sigma\mathbb{J}j_{2,2} = \mathbb{J}P_E\mathbb{J}j_{2,2} = 1$  and the distribution is supported within some centered Euclidean ball with radius  $\sqrt{O(k)}$  to get that with probability at least  $1 - n^{-t^2}$  and  $T_0 \geq C(t/\xi)^2 k \log n$ ,

$$\mathbb{J}\frac{1}{T_0}X_{pre}X_{pre}^\top \quad \Sigma\mathbb{J}j_{2,2} \quad \xi.$$

We then use this to show that

$$\mathbb{J}X_{pre}X_{pre}^\top \quad T_0\mathbb{J}j_{2,2} \quad \xi T_0,$$

and thus all eigenvalues of  $(X_{pre}X_{pre}^\top \quad T_0I)$  must be at most  $\xi T_0$ , so

$$(1 - \xi)T_0 \leq \lambda_{\min}(X_{pre}X_{pre}^\top) \leq (1 + \xi)T_0.$$

Finally, we can complete the proof of Lemma C.7, by observing that,

$$\mathbb{J}(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} \mathbb{J}j_{2,2} = \frac{1}{\mathbb{J}\lambda_{\min}(X_{pre}X_{pre}^\top) + \frac{\lambda}{2T_0}} \leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}.$$

Returning to Equation (14), we can use this bound — along with the model properties specified in Equation (1) that each element of  $\mathbf{z}$  and  $Z$  has mean 0, variance  $\sigma^2$ , and support  $[-s, s]$  — to obtain the desired bound:

$$\begin{aligned} \mathbb{E}[\mathbb{J}\mathbf{f}^{reg} \quad \mathbf{f}j_{2,2}] &\leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}} \mathbb{E}[\mathbb{J}X_{pre}\mathbf{z} \quad X_{pre}Z^\top \mathbf{f}j_{2,2} + \frac{\lambda}{2T_0}] \\ &\leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}} \left( \mathbb{P} \frac{1}{nT_0} + \mathbb{P} \frac{1}{nT_0\sigma^2} \right) \mathbb{E}[\mathbb{J}\mathbf{z} \quad Z^\top \mathbf{f}j_{2,2}] + \frac{\lambda}{2T_0} \\ &\leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}} \left( \mathbb{P} \frac{1}{nT_0} + \mathbb{P} \frac{1}{nT_0\sigma^2} \right) \mathbb{P} \frac{1}{2T_0\sigma^2} + \frac{\lambda}{2T_0} \\ &\leq \frac{\left( \mathbb{P} \frac{1}{2n\sigma^2} + \mathbb{P} \frac{1}{2n\sigma^2 s^2} \right) T_0 + \frac{\lambda}{2T_0}}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}. \end{aligned}$$

## D Privacy Guarantees of $DPSC_{obj}$

In this section, we prove Theorem 3.4, that  $DPSC_{obj}$  is  $(\epsilon_1 + \epsilon_2, \delta)$ -differentially private. This proof relies on composition of the  $(\epsilon_1, \delta)$ -DP learning step and the  $(\epsilon_2, 0)$ -DP prediction step. The prediction step is identical to that of Algorithm 2, so the privacy of this step follows immediately from Lemma 3.9 (that  $\tilde{X}_{post}$  is computed in an  $(\epsilon_2, 0)$ -DP manner) and post-processing on the DP output of Step 1. All that remains to be shown is that  $\mathbf{f}^{obj}$  is computed in an  $(\epsilon_1, \delta)$ -DP manner (Theorem D.1), and then Theorem 3.4 will follow by basic composition.

**Theorem D.1** *Step 1 of Algorithm 3 that computes  $\mathbf{f}^{obj}$  is  $(\epsilon_1, \delta)$ -differentially private.*

At a high-level, the privacy of  $\mathbf{f}^{obj}$  comes from a carefully modified instantiation of the Objective Perturbation algorithms of Chaudhuri et al. (2011); Kifer et al. (2012), with novel sensitivity analysis, again due to the transposed regression setting of synthetic control (i.e., along columns not rows), where privacy is still required along the rows.

More formally, we start with the standard Ridge Regression objective function  $J(\mathbf{f})$ , that can be separated into the MSE loss function  $L(\mathbf{f})$  and the regularization term  $r(\mathbf{f}) = \frac{\lambda}{2T_0} \|\mathbf{f}\|_2^2$  as follows:

$$J(\mathbf{f}) = L(\mathbf{f}) + r(\mathbf{f}) = \frac{1}{T_0} \mathbf{y}_{pre}^\top \mathbf{X}_{pre}^\top \mathbf{f} \mathbf{f}^\top \mathbf{f} + \frac{\lambda}{2T_0} \|\mathbf{f}\|_2^2.$$

The Objective Perturbation method modifies  $J(\mathbf{f})$  by adding two terms: an additional regularization term and a noise term to ensure privacy:

$$J^{obj}(\mathbf{f}) = J(\mathbf{f}) + \frac{\Delta}{2T_0} \|\mathbf{f}\|_2^2 + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f} = L(\mathbf{f}) + \frac{\lambda + \Delta}{2T_0} \|\mathbf{f}\|_2^2 + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f},$$

where  $\mathbf{b}$  is a random vector drawn from a high-dimensional Laplace distribution if  $\delta = 0$ , and from a multivariate Gaussian distribution if  $\delta > 0$ .

Notice that  $J^{obj}(\mathbf{f})$  is strongly convex (for any  $\Delta \geq 0$ ) and differentiable. Hence, for any given input dataset  $D = (X_{pre}, y_{pre})$  and any fixed parameters  $(\lambda, \epsilon_1, \epsilon_2, \delta)$ , there exists a bijection between a realized value of the noise term  $\mathbf{b}$  and  $\mathbf{f}^{obj} := \arg \min_{\mathbf{f}} J^{obj}(\mathbf{f})$  given that realized  $\mathbf{b}$ .<sup>2</sup> We can then use this bijection to analyze the distribution over outputs on neighboring databases via the (explicitly given) noise distribution.

To observe this bijection concretely, let  $\mathbf{b}(\alpha; D)$  be noise value that must have been realized when database  $D$  was input and  $\alpha = \arg \min_{\mathbf{f}} J^{obj}(\mathbf{f})$  was output. We can derive a closed-form expression for  $\mathbf{b}(\alpha; D)$  by computing the gradient of  $J^{obj}(\mathbf{f})$ , which should be zero when evaluated at  $\mathbf{f} = \alpha$  since  $\alpha$  is defined to be the minimizer of  $J^{obj}(\mathbf{f})$ :

$$\nabla_{\mathbf{f}} J^{obj}(\mathbf{f}) \Big|_{\mathbf{f}=\alpha} = \nabla_{\mathbf{f}} L(\alpha) + \nabla_{\mathbf{f}} r(\alpha) + \frac{\Delta}{T_0} \alpha + \frac{\mathbf{b}(\alpha; D)}{T_0} \stackrel{!}{=} \mathbf{0}.$$

Rearranging the equation yields

$$\mathbf{b}(\alpha; D) = (T_0 \nabla_{\mathbf{f}} L(\alpha) + T_0 \nabla_{\mathbf{f}} r(\alpha) + \Delta \alpha).$$

Now, consider two arbitrary neighboring databases  $D$  and  $D'$  and an arbitrary output value  $\alpha$ . Similar to Chaudhuri et al. (2011), we can use, e.g., Billingsley (1995) to express the ratio of the probabilities of outputting  $\alpha$  on neighboring  $D$  and  $D'$  as:<sup>3</sup>

$$\frac{\Pr(\mathbf{f}^{obj} = \alpha \mid D)}{\Pr(\mathbf{f}^{obj} = \alpha \mid D')} = \frac{\Pr(\mathbf{b}(\alpha; D)) \cdot |\det(\nabla_{\mathbf{b}} \mathbf{b}(\alpha; D))|}{\Pr(\mathbf{b}(\alpha; D')) \cdot |\det(\nabla_{\mathbf{b}} \mathbf{b}(\alpha; D'))|} := \Gamma(\alpha) \cdot \Phi(\alpha; \Delta),$$

where we define  $\Gamma(\alpha) := \frac{\Pr(\mathbf{b}(\alpha; D))}{\Pr(\mathbf{b}(\alpha; D'))}$  and  $\Phi(\alpha; \Delta) := \frac{|\det(\nabla_{\mathbf{b}} \mathbf{b}(\alpha; D))|}{|\det(\nabla_{\mathbf{b}} \mathbf{b}(\alpha; D'))|}$ . In the remainder of the proof, we will bound  $\Gamma(\alpha) \leq e^{\epsilon_0}$  and  $\Phi(\alpha; \Delta) \leq e^{\epsilon_1 - \epsilon_0}$  so that the product is bounded by  $e^{\epsilon_1}$ .

The parameter  $\Delta$  serves a role to divide the  $\epsilon_1$  budget between these two terms, by distinguishing between two cases. In the first case,  $\epsilon_1$  is large enough that we can choose  $\Delta = 0$  and still have some privacy budget ( $\epsilon_0$ ) remaining to bound  $\Gamma(\alpha)$ . In the other case, if  $\epsilon_1$  is too small to bound  $\Phi(\alpha; \Delta)$  with  $\Delta = 0$ , then we divide the privacy budget equally between bounding  $\Gamma(\alpha)$  and  $\Phi(\alpha; \Delta)$ , and find an appropriate value for  $\Delta > 0$ .

First, we will show  $\Phi(\alpha; \Delta)$  is upper bounded by  $e^{\epsilon_1 - \epsilon_0}$ .

<sup>2</sup>For a simple analogy, consider the one-dimensional Laplace Mechanism on query  $f$  and database  $x$ , which outputs  $y = f(x) + \text{Lap}(f; \cdot)$ . Given  $f$  and  $x$ , there is a bijection between noise terms and outputs since the noise term must equal  $y - f(x)$ .

<sup>3</sup>with abuse of notation to let  $\Pr$  denote pdf for simplicity of presentation.

**Lemma D.2** If  $\Delta = 0$  and  $\epsilon_0 = \epsilon_1 - \log(1 + \frac{2c}{\lambda} + \frac{c^2}{\lambda^2})$ , or if  $\Delta = \frac{c}{e^{\epsilon_1 - 4} - 1} - \lambda$  and  $\epsilon_0 = \frac{\epsilon_1}{2}$ , then  $\Phi(\alpha; \Delta) \leq e^{\epsilon_1 - \epsilon_0}$ .

We start with Lemma D.3 (proved in Appendix G.1), which bounds  $\Phi(\alpha; \Delta)$  as a function of  $\lambda$ ,  $c$ , and  $\Delta$ .

**Lemma D.3** For any  $\Delta \geq 0$ ,  $\Phi(\alpha; \Delta) = \frac{|\det(\nabla \mathbf{b}(\alpha; \mathcal{D}^\delta))|}{|\det(\nabla \mathbf{b}(\alpha; \mathcal{D}))|} (1 + \frac{c}{\lambda})^2$ .

Next, we use this result to prove our desired bound that  $\Phi(\alpha; \Delta) \leq e^{\epsilon_1 - \epsilon_0}$ . We do this by considering two cases. First, when  $\Delta = 0$ , then  $\Phi(\alpha; \Delta = 0) \leq 1 + \frac{2c}{\lambda} + \frac{c^2}{\lambda^2} \leq e^{\epsilon_1 - \epsilon_0}$  by design, where the first inequality comes from Lemma D.3 and the second inequality is a rearrangement of our choice of  $\epsilon_0 = \epsilon_1 - \log(1 + \frac{2c}{\lambda} + \frac{c^2}{\lambda^2})$  in this case. In the second case,  $\Delta = \frac{c}{e^{\epsilon_1 - 4} - 1} - \lambda$ . Plugging this  $\Delta$  value into the bound of Lemma D.3 gives  $\Phi(\alpha; \Delta) \leq e^{\epsilon_1/2} = e^{\epsilon_1 - \epsilon_0}$ , where the second inequality comes from our choice of  $\epsilon_0 = \epsilon_1/2$ . Hence, in both cases,  $\Phi(\alpha; \Delta) \leq e^{\epsilon_1 - \epsilon_0}$ .

Then, we bound  $\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; \mathcal{D}))}{\Pr(\mathbf{b}(\alpha; \mathcal{D}^\delta))}$ . Note that this term depends only on the noise distribution, and not on the value of  $\Delta$ . Algorithm 3 offers two options of noise distributions: Laplace noise when  $\delta = 0$ , and Gaussian noise when  $\delta > 0$ .

In the case of Laplace noise, the bound that  $\Gamma(\alpha) \leq e^{\epsilon_0}$  follows immediately from the Laplace mechanism instantiated with privacy parameter  $\epsilon_0$  and Lemma B.4 to bound the sensitivity. The following lemma is proved in Appendix G.2.

**Lemma D.4** When  $\mathbf{b}$  is sampled according to pdf  $p(\mathbf{b}; \beta) \propto \exp(-\frac{\|\mathbf{b}\|_2}{\beta})$ , where  $\beta = \min\{\frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n+4T_0}}{\epsilon_0}\}$ , then  $\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; \mathcal{D}))}{\Pr(\mathbf{b}(\alpha; \mathcal{D}^\delta))} \leq e^{\epsilon_0}$ .

The two different  $\beta$  values come from two different upper bounds on the sensitivity, and the minimum value will give a tighter bound.

In the case where  $\delta > 0$  and the Gaussian Mechanism is used, we cannot simply bound  $\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; \mathcal{D}))}{\Pr(\mathbf{b}(\alpha; \mathcal{D}^\delta))}$  with probability 1. Instead, the bound must incorporate the  $\delta$  term to bound  $\Gamma(\alpha)$  with probability  $1 - \delta$  over the internal randomness of the algorithm, as in Lemma D.5, formally proven in Appendix G.3.

**Lemma D.5** When  $\mathbf{b} \sim N(0, \beta^2 I_n)$ , where  $\beta = \frac{4T_0\sqrt{8+n}}{\epsilon_0} \frac{\rho}{2 \log 2 + \epsilon_0}$ , then  $\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; \mathcal{D}))}{\Pr(\mathbf{b}(\alpha; \mathcal{D}^\delta))} \leq e^{\epsilon_0}$  with probability at least  $1 - \delta$ .

Finally, we combine the bounds on  $\Phi(\alpha; \Delta)$  and  $\Gamma(\alpha)$  to complete the proof. When  $\delta = 0$  with Laplace noise, Lemmas D.2 and D.4 combine immediately to give that  $\Phi(\alpha; \Delta)\Gamma(\alpha) \leq e^{\epsilon_1 - \epsilon_0 + \epsilon_0} = e^{\epsilon_1}$ . When  $\delta > 0$  and Gaussian noise is used, we define  $G$  to be the good event that  $\Gamma(\alpha) \leq e^{\epsilon_0}$ , which we know from Lemma D.5 will happen with at least probability  $1 - \delta$ . Then conditioned on  $G$  we have,

$$\frac{\Pr(\mathbf{f}^{obj} = \alpha_j D, G)}{\Pr(\mathbf{f}^{obj} = \alpha_j D', G)} = \Gamma(\alpha) \Phi(\alpha; \Delta) \leq e^{\epsilon_0} e^{\epsilon_1 - \epsilon_0} = e^{\epsilon_1}.$$

We can then use this fact to derive our desired (unconditioned) privacy bound:

$$\begin{aligned} \Pr(\mathbf{f}^{obj} = \alpha_j D) &= \Pr(G) \Pr(\mathbf{f}^{obj} = \alpha_j D, G) + \Pr(\bar{G}) \Pr(\mathbf{f}^{obj} = \alpha_j D, \bar{G}) \\ &\leq e^{\epsilon_1} \Pr(G) \Pr(\mathbf{f}^{obj} = \alpha_j D', G) + \delta \\ &\leq e^{\epsilon_1} \Pr(\mathbf{f}^{obj} = \alpha_j D') + \delta. \end{aligned}$$

Hence,  $\mathbf{f}^{obj}$  in Algorithm 3 is  $(\epsilon_1, \delta)$ -DP and the final output  $\mathbf{y}^{obj}$  is  $(\epsilon_1 + \epsilon_2, \delta)$ -DP by composition.

## E Accuracy Guarantees of $DPSC_{obj}$

In this section we analyze the accuracy of  $DPSC_{obj}$ . We first prove Theorem 3.5, restated below for convenience.

**Theorem 3.5** The estimator  $\mathbf{y}^{obj}$  output by Algorithm 3 satisfies:

$$\begin{aligned} RMSE(\mathbf{y}^{obj}) &\leq \frac{\|j\|_2 M_{post} \|j\|_2}{T - T_0} \mathbb{E}[\|j(\mathbf{f}^{reg} - \mathbf{f})\|_2] + \frac{2}{\lambda + \Delta} \mathbb{E}[\|j\mathbf{b}j\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} + 2T_0^2 \frac{\rho}{n} \\ &+ \frac{\rho}{n\sigma^2} + \frac{\rho}{2} \frac{1}{\epsilon_2} \frac{\rho}{n\psi} + \frac{2}{\lambda + \Delta} \mathbb{E}[\|j\mathbf{b}j\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} + 2T_0^2 \frac{\rho}{n}, \end{aligned}$$

where  $\|f^{reg}\|_{\infty} \leq \psi$  for some  $\psi > 0$ , and  $E[\|b\|_2] = \frac{\rho}{\epsilon_0} \frac{1}{nT_0 4\sqrt{8+n} \sqrt{2 \log \frac{2}{\epsilon_0}}}$  for Gaussian noise ( $\delta > 0$  case) and  $E[\|b\|_2] = \min\left\{f^{\frac{4T_0\sqrt{8+n}}{\epsilon_0}}, \frac{c\sqrt{n+4T_0}}{\epsilon_0}g\right\}$  for Laplace noise ( $\delta = 0$  case), and  $\epsilon_0$ , and  $\Delta$  are computed internally by the algorithm.

This theorem gives bounds on the predicted post-intervention target vector  $\mathbf{y}^{obj}$ , as measured by RMSE. Similar to Theorem 3.2, this result is stated in full generality with respect to the distribution of data and the latent variables, and thus the bound depends on terms such as  $\|M_{post}\|_{2,2}$  and  $E[\|f^{reg}\|_2]$ . Section E.1 provides proof of this result, with omitted detailed deferred to Appendix G.

Comparing the bound of Theorem 3.5 to that of Theorem 3.2 for output perturbation, we see that the difference comes only from the respective terms  $E[\|(\mathbf{f}^{out} - \mathbf{f}^{reg})\|_2]$ . For output perturbation, the error  $\mathbf{f}^{out} - \mathbf{f}^{reg}$  is simply the noise directly added to the output, so the expected norm of the error is simply the expected norm of the noise,  $a = \frac{4T_0\sqrt{8+n}}{\lambda\epsilon_1}$ . For objective perturbation, the interpretation of these error terms is less straightforward and is instead bounded using Lemma E.1. As a simple case for comparison, when  $\Delta = 0$  and  $\delta = 0$  (i.e., using Laplace noise), the expected difference becomes  $E[\|(\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2] = \min\left\{f^{\frac{8T_0\sqrt{8+n}}{\lambda\epsilon_0}}, \frac{2c\sqrt{n+8T_0}}{\lambda\epsilon_0}g\right\}$ . If the first term is the smaller of the two, then  $E[\|(\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2]$  is bigger than  $E[\|(\mathbf{f}^{out} - \mathbf{f}^{reg})\|_2]$  since the denominator is smaller ( $\epsilon_0 < \epsilon_1$ , assuming the same  $\epsilon_1$  values for comparison) and the numerator is bigger due to the factor of 2. If the second term is the minimum, then the upper bound on error is hard to compare as both the denominator and the numerator are (asymptotically) bigger for output perturbation. In case of  $\epsilon_0 = \frac{\epsilon_{ps1}}{2\rho}$ , the expected difference of output perturbation becomes  $O(T_0 \frac{\rho}{n})$  and that of objective perturbation becomes  $O(T_0 + \frac{\rho}{n})$ . In this case, we may expect the objective perturbation to yield a better RMSE for a reasonably big  $T_0$  and  $n$ .

### E.1 Accuracy of post-intervention prediction via objective perturbation $\mathbf{y}^{obj}$

We will prove Theorem 3.5, which upper bounds the Root Mean Squared Error (RMSE) of  $\mathbf{y}^{obj}$ , defined as:

$$RMSE(\mathbf{y}^{obj}) = \frac{1}{T - T_0} E[\|\mathbf{y}^{obj} - \mathbf{m}_{post}\|_2].$$

Using the facts that  $\mathbf{y}^{obj} = \tilde{X}_{post}^\top \mathbf{f}^{obj}$ ,  $\tilde{X}_{post} = X_{post} + M_{post} + Z_{post}$ , and  $\mathbf{m}_{post} = M_{post}^\top \mathbf{f}$  (by Equation (2)), we can bound the expectation as follows:

$$\begin{aligned} E[\|\mathbf{y}^{obj} - \mathbf{m}_{post}\|_2] &= E[\|\tilde{X}_{post}^\top \mathbf{f}^{obj} - M_{post}^\top \mathbf{f}\|_2] \\ &= E[\|\tilde{X}_{post}^\top \mathbf{f}^{obj} - \tilde{X}_{post}^\top \mathbf{f}^{reg} + \tilde{X}_{post}^\top \mathbf{f}^{reg} - M_{post}^\top \mathbf{f}\|_2] \\ &= E[\|(M_{post} + Z_{post} + W_{post})^\top (\mathbf{f}^{obj} - \mathbf{f}^{reg}) + (M_{post} + Z_{post} + W_{post})^\top \mathbf{f}^{reg} - M_{post}^\top \mathbf{f}\|_2] \\ &= E[\|(M_{post} + Z_{post} + W_{post})^\top (\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2] \\ &\quad + E[\|M_{post}^\top (\mathbf{f}^{reg} - \mathbf{f})\|_2] + E[\|(Z_{post} + W_{post})^\top \mathbf{f}^{reg}\|_2], \end{aligned} \quad (15)$$

where the first equality is due to the definition of  $\mathbf{y}^{obj}$ , the second equality adds and subtracts the same term, the third equality collects terms and plugs in the expression for  $\tilde{X}_{post}$ , and the final step is due to triangle inequality.

Lemma C.1 already bounds the last two terms because they do not involve  $\mathbf{f}^{obj}$  and Step 2 of Algorithms 2 and 3 are the same. Specifically, we know that,

$$E[\|M_{post}^\top (\mathbf{f}^{reg} - \mathbf{f})\|_2] \leq \|M_{post}\|_{2,2} E[\|\mathbf{f}^{reg} - \mathbf{f}\|_2] \text{ and } E[\|(Z_{post} + W_{post})^\top \mathbf{f}^{reg}\|_2] \leq \frac{\rho}{n\psi} \left( \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\sqrt{T - T_0}}{\epsilon_2} \right).$$

Thus we only need to bound the first term:

$$\begin{aligned} &E[\|(M_{post} + Z_{post} + W_{post})^\top (\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2] \\ &= E[\|M_{post}^\top (\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2] + E[\|(Z_{post} + W_{post})^\top (\mathbf{f}^{obj} - \mathbf{f}^{reg})\|_2] \\ &= \|M_{post}\|_{2,2} E[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] + E[\|Z_{post} + W_{post}\|_2] E[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] \\ &= \|M_{post}\|_{2,2} E[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] + \left( \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2\rho\sqrt{T - T_0}}{\epsilon_2} \right) E[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2], \end{aligned} \quad (17)$$

where the first step is simply triangle inequality, the second step is due to the independence of  $Z, W$  and  $\mathbf{f}^{obj}$ , and the third step comes from the proof of Lemma C.1 (see Appendix F.2), where  $\mathbb{E}[\|Z_{post} + W_{post}\|_2]$  was bounded as an intermediate step.

Thus we only need to derive a bound on  $\mathbb{E}[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2]$ , which we do in Lemma E.1 (formally proven in Appendix G.4) to complete the proof.

**Lemma E.1** *The  $\ell_2$  distance between  $\mathbf{f}^{obj}$  and  $\mathbf{f}^{reg}$  satisfies:*

$$\mathbb{E}[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] \leq \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^-,$$

where  $\mathbf{b}$  and  $\Delta$  are computed internally by Algorithm 3.

Combining Equations (15), (16), and (17) with Lemma E.1 completes the proof of Theorem 3.5:

$$\begin{aligned} RMSE(\mathbf{y}^{obj}) & \leq \frac{\|M_{post}\|_2}{T - T_0} \mathbb{E}[\|\mathbf{f}^{reg} - \mathbf{f}\|_2] + \mathbb{E}[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] \\ & + \frac{\rho_n^-}{n\sigma^2} + \frac{\rho_n^-}{\epsilon_2} \rho_n^- \psi + \mathbb{E}[\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2] \\ & \leq \frac{\|M_{post}\|_2}{T - T_0} \mathbb{E}[\|\mathbf{f}^{reg} - \mathbf{f}\|_2] + \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^- \\ & + \frac{\rho_n^-}{n\sigma^2} + \frac{\rho_n^-}{\epsilon_2} \rho_n^- \psi + \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^-, \end{aligned}$$

where  $\mathbb{E}[\|\mathbf{b}\|_2] = \frac{\rho_n^-}{n\beta} = \frac{\rho_n^-}{nT_0\zeta} \frac{\rho_n^-}{\epsilon_0^{2\log 2 + \epsilon_0}}$  for Gaussian noise and  $\mathbb{E}[\|\mathbf{b}\|_2] = \min\left\{\frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n+4T_0}}{\epsilon_0}g\right\}$  for Laplace noise.

## E.2 Closed-form bound on RMSE of Objective Perturbation

Using similar analysis as in Section C.2, we can extend Theorem 3.5 to obtain the following closed-form accuracy bound that depends only on explicit input parameters, under the same distributional assumptions.

**Corollary E.2** *If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \geq (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t^2}$ , if  $T_0 \geq C(t/\xi)^2 k \log n$ , we have*

$$\begin{aligned} RMSE(\mathbf{y}^{out}) & \leq \frac{\rho_n^-}{n} \frac{(\frac{\rho_n^-}{2n\sigma^2} + \frac{\rho_n^-}{2n\sigma^2 s^2})T_0 + \frac{\lambda}{2T_0}}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}} + \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^- \\ & + \frac{\rho_n^-}{n\sigma^2} + \frac{\rho_n^-}{\epsilon_2} \rho_n^- \psi + \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^-, \end{aligned}$$

where  $\mathbb{E}[\|\mathbf{f}^{reg}\|_2] \leq \psi$  for some  $\psi > 0$ , and  $\mathbb{E}[\|\mathbf{b}\|_2] = \frac{\rho_n^-}{nT_0\zeta} \frac{\rho_n^-}{\epsilon_0^{2\log 2 + \epsilon_0}}$  for Gaussian noise ( $\delta > 0$  case) and  $\mathbb{E}[\|\mathbf{b}\|_2] = \min\left\{\frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n+4T_0}}{\epsilon_0}g\right\}$  for Laplace noise ( $\delta = 0$  case), and  $\epsilon_0$  and  $\Delta$  are computed internally by the algorithm.

The additional terms that arise due to the noise required to guarantee differential privacy in this setting, relative to the bound on  $RMSE(\mathbf{y}^{reg})$  in Equation (11), are:

$$\left(\frac{\rho_n^-}{n} + \frac{\rho_n^-}{n\sigma^2} + \frac{\rho_n^-}{\epsilon_2}\right) \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \cdot 2T_0^2 \rho_n^- + \frac{\rho_n^-}{\epsilon_2} \psi. \quad (18)$$

To analyze this expression in a simplified way, assume the regularization parameter is  $\lambda = O(T_0)$  so  $\frac{T_0}{\lambda} = O(1)$ , and that Laplace noise was used (i.e.,  $\delta = 0$ ), so that  $\mathbb{E}[\|\mathbf{b}\|_2] = O\left(\frac{T_0\sqrt{n}}{\epsilon_0}\right)$ . Then the first parenthesis of Equation (18) becomes



$O(\frac{D}{n} + \frac{1}{\epsilon_2})$ , the second parenthesis becomes  $O(\frac{T_0\sqrt{n}}{\epsilon_0} + T_0 \frac{D}{n})$ , and the additive term becomes  $O(\frac{\sqrt{n}}{\epsilon_2})$ . Since  $\epsilon_0 < \epsilon_1$ , we replace  $\epsilon_0$  by  $\epsilon_1$  in the bounds. Then Equation (18) is  $O(\frac{T_0 n}{\epsilon_1} + \frac{T_0\sqrt{n}}{\epsilon_1\epsilon_2})$  from the product of two parentheses, and omitting the additive term, which is asymptotically dominated by the others.

Comparing to the cost of privacy in Output Perturbation in Corollary C.5, we see that the bound in Corollary C.5 does not depend on  $T_0$ . This additional dependence on  $T_0$  arises for Objective Perturbation from the second parenthesis containing  $\mathbb{E}[\|j\|_2]$  and the indicator function, which is absent in the output perturbation case.

## F Omitted Proofs for $DPSC_{out}$

### F.1 Proof of Lemma B.4

**Lemma B.4** *Let  $g(\mathbf{f}) = L(\mathbf{f}, D') - L(\mathbf{f}, D)$  for two arbitrarily neighboring databases  $D, D'$ . Then,*

$$\max_{\mathbf{f}} \|g(\mathbf{f})\| \leq \frac{D}{4\sqrt{8+n}}.$$

We first re-arrange  $g(\mathbf{f})$  in a way that makes it easier to compute the gradient. Let  $i$  be the index of the record that differs between  $D$  and  $D'$ .

$$\begin{aligned} g(\mathbf{f}) &= L(\mathbf{f}, D') - L(\mathbf{f}, D) \\ &= \frac{1}{T_0} \sum_{t=1}^T \sum_{k=1}^n (x'_{k,t} f_k - y_t) - \frac{1}{T_0} \sum_{t=1}^T \sum_{k=1}^n (x_{k,t} f_k - y_t) \\ &= \frac{1}{T_0} \sum_{t=1}^T \sum_{j \neq i} (x'_{j,t} f_j - y_t + x'_{i,t} f_i) - \frac{1}{T_0} \sum_{t=1}^T \sum_{j \neq i} (x_{j,t} f_j - y_t + x_{i,t} f_i) \\ &= \frac{1}{T_0} \sum_{t=1}^T \sum_{j \neq i} (x_{j,t} f_j - y_t) + (x'_{i,t} - x_{i,t}) f_i + (x'^2_{i,t} - x^2_{i,t}) f_i^2 \\ &= \frac{1}{T_0} \sum_{t=1}^T \sum_{j \neq i} (x_{j,t} f_j - y_t) + (x'_{i,t} + x_{i,t}) f_i + (x'_{i,t} - x_{i,t}) f_i \\ &= \frac{1}{T_0} \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{f} - y_t + \mathbf{x}_t^\top \mathbf{f} \cdot (x'_{i,t} f_i + x_{i,t} f_i) - y_t + (x'_{i,t} - x_{i,t}) f_i \end{aligned} \quad (19)$$

The second equality comes from the definition of the Ridge regression loss function  $L$ ; in the third step we pull out the record  $i$  that differs between  $D$  and  $D'$ ; the fourth step combines the sums and cancels terms, including the observation that  $\sum_{j \neq i} x_{j,t} f_j = \sum_{j \neq i} x'_{j,t} f_j$ . The final two steps also involve rearranging terms.

For notational ease, we define two additional terms,

$$D_t := \mathbf{x}_t^\top \mathbf{f} - y_t \quad \text{and} \quad E_t := (x'_{i,t} - x_{i,t}) f_i.$$

Then, Equation (19) becomes

$$g(\mathbf{f}) = \frac{1}{T_0} \sum_{t=1}^T (2D_t + E_t) E_t.$$

We will take the partial derivatives of  $D_t$  and  $E_t$  with respect to both  $f_i$  (the index of the data entry that differs between  $D$  and  $D'$ ) and  $f_j$  for  $j \neq i$ , and then combine these to arrive at the gradient of  $g(\mathbf{f})$ :

$$\frac{\partial D_t}{\partial f_i} = x_{i,t}; \quad \frac{\partial D_t}{\partial f_j} = x_{j,t}; \quad \frac{\partial E_t}{\partial f_i} = x'_{i,t} - x_{i,t}; \quad \frac{\partial E_t}{\partial f_j} = 0.$$

Now, we compute the derivative of  $g(\mathbf{f})$  with respect to  $f_i$ .

$$\begin{aligned}
 \frac{\partial g(\mathbf{f})}{\partial f_i} &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} 2 \frac{\partial D_t}{\partial f_i} + \frac{\partial E_t}{\partial f_i} E_t + (2D_t + E_t) \frac{\partial E_t}{\partial f_i} \\
 &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} (x'_{i,t} + x_{i,t}) E_t + (2D_t + E_t) (x'_{i,t} - x_{i,t}) \\
 &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} 2x'_{i,t} (x'_{i,t} - x_{i,t}) f_i + 2(\mathbf{x}_t^\top \mathbf{f} - y_t) (x'_{i,t} - x_{i,t}) \\
 &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} 2 (x'_{i,t} f_i + \mathbf{x}_t^\top \mathbf{f} - y_t) (x'_{i,t} - x_{i,t})
 \end{aligned} \tag{20}$$

Next, we compute the derivative of  $g(\mathbf{f})$  with respect to  $f_j$  where  $j$  is the index of unchanged donors ( $j \neq i$ ). There are fewer term in this derivative because  $\frac{\partial E_t}{\partial f_j} = 0$ .

$$\begin{aligned}
 \frac{\partial g(\mathbf{f})}{\partial f_j} &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} f(2x_{j,t}) E_t g \\
 &= \frac{1}{T_0} \sum_{t=1}^{\mathcal{X}_0} 2x_{j,t} (x'_{i,t} - x_{i,t}) f_i
 \end{aligned} \tag{21}$$

Finally, we can use (20) and (21) to derive an upper bound for  $\|j\|_2^2 g(\mathbf{f})\|_2$ .

$$\begin{aligned}
 \|j\|_2^2 g(\mathbf{f})\|_2^2 &= \sum_{j \neq i} \left( \frac{\partial g(\mathbf{f})}{\partial f_i} \right)^2 + \sum_{j \neq i} \left( \frac{\partial g(\mathbf{f})}{\partial f_j} \right)^2 \\
 &= \sum_{t=1}^{\mathcal{X}_0} \frac{1}{T_0} \left( 2 (x'_{i,t} f_i + \mathbf{x}_t^\top \mathbf{f} - y_t) (x'_{i,t} - x_{i,t}) \right)^2 + \sum_{j \neq i} \sum_{t=1}^{\mathcal{X}_0} \frac{1}{T_0} \left( 2x_{j,t} (x'_{i,t} - x_{i,t}) f_i \right)^2 \\
 &= \sum_{t=1}^{\mathcal{X}_0} \frac{1}{T_0} \left( 2 (x'_{i,t} f_i + \mathbf{x}_t^\top \mathbf{f} - y_t) (x'_{i,t} - x_{i,t}) \right)^2 + \sum_{j \neq i} \sum_{t=1}^{\mathcal{X}_0} \frac{1}{T_0} [2x_{j,t} (x'_{i,t} - x_{i,t}) f_i]^2 \\
 &= \sum_{t=1}^{\mathcal{X}_0} \frac{1}{T_0} \left( 4 (x'_{i,t} f_i + \mathbf{x}_t^\top \mathbf{f} - y_t)^2 (x'_{i,t} - x_{i,t})^2 + 4x_{j,t}^2 (x'_{i,t} - x_{i,t})^2 f_i^2 \right) \\
 &= \sum_{t=1}^{\mathcal{X}_0} \frac{4}{T_0} (x'_{i,t} - x_{i,t})^2 \left( (x'_{i,t} f_i + \mathbf{x}_t^\top \mathbf{f} - y_t)^2 + x_{j,t}^2 f_i^2 \right) \\
 &= \sum_{t=1}^{\mathcal{X}_0} \frac{4}{T_0} (x'_{i,t} - x_{i,t})^2 \left( (\mathbf{x}_t^\top \mathbf{f} - y_t)^2 + 2x'_{i,t} f_i (\mathbf{x}_t^\top \mathbf{f} - y_t) + \sum_{k=1}^{\mathcal{X}_0} x_{k,t}^2 f_i^2 \right)
 \end{aligned} \tag{22}$$

The second equality comes from plugging in the partial derivatives computed in (20) and (21), the following inequality comes from applying Jensen's inequality, and the final three steps come from rearranging, expanding, and simplifying terms.

We can proceed by bounding the individual terms in (22) using the our modeling assumptions of Equation (3), which give us that:

$$(x'_{i,t} - x_{i,t})^2 \leq 4, \quad \text{and} \quad (\mathbf{x}_t^\top \mathbf{f} - y_t)^2 \leq 4, \quad \text{and} \quad 2x'_{i,t} f_i (\mathbf{x}_t^\top \mathbf{f} - y_t) \leq 4, \quad \text{and} \quad \sum_{k=1}^{\mathcal{X}_0} x_{k,t}^2 f_i^2 \leq n.$$

Then  $\|j\|_2^2 g(\mathbf{f})\|_2^2 \leq 128 + 16n$  and  $\|j\|_2 g(\mathbf{f})\|_2 \leq 4 \sqrt{8 + n}$ .

## F.2 Proof of Lemma C.1

**Lemma C.1** *The three terms in Equation (8) can be bounded as follows:*

$$\begin{aligned} & \mathbb{E}[jjM_{post}^\top(\mathbf{f}^{reg} - \mathbf{f})jj_2] \leq jjM_{post}jj_{2,2} \mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2], \\ & \mathbb{E}[jj(Z^\top + W^\top)\mathbf{f}^{reg}jj_2] \leq \frac{\rho_-}{n\psi} \left( \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2^\rho \overline{T - T_0}}{\epsilon_2} \right), \text{ and} \\ & \mathbb{E}[jj(M_{post}^\top + Z^\top + W^\top)\mathbf{v}jj_2] \leq jjM_{post}jj_{2,2} + \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2^\rho \overline{T - T_0}}{\epsilon_2} \frac{4T_0 \rho}{\lambda\epsilon_1 \overline{8+n}}. \end{aligned}$$

We prove these three bounds separately. Most steps follow from the sub-multiplicative norm property of Equation (9) and the bounds on the noise terms of Equation (10).

First,

$$\begin{aligned} \mathbb{E}[jjM_{post}^\top(\mathbf{f}^{reg} - \mathbf{f})jj_2] & \leq \mathbb{E}[jjM_{post}jj_{2,2} jj\mathbf{f}^{reg} - \mathbf{f}jj_2] \\ & \leq jjM_{post}jj_{2,2} \mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2]. \end{aligned}$$

Next,

$$\begin{aligned} \mathbb{E}[jj(Z_{post}^\top + W^\top)\mathbf{f}^{reg}jj_2] & \leq \mathbb{E}[jjZ_{post} + Wjj_2 jj\mathbf{f}^{reg}jj_2] \\ & \leq \mathbb{E}[jjZ_{post} + Wjj_2 \frac{\rho}{n} jj\mathbf{f}^{reg}jj_\infty] \\ & \leq \mathbb{E}[jjZ_{post} + Wjj_2 \frac{\rho_-}{n\psi}] \\ & \leq \frac{\rho_-}{n\psi} \mathbb{E}[jjZ_{post}jj_2 + jjWjj_2] \\ & \leq \frac{\rho_-}{n\psi} \mathbb{E}[jjZ_{post}jj_F + jjWjj_F] \\ & \leq \frac{\rho_-}{n\psi} \left( \rho \frac{\rho}{n(T - T_0)\sigma^2} + b \right) \\ & \leq \frac{\rho_-}{n\psi} \left( \rho \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2^\rho \overline{T - T_0}}{\epsilon_2} \right), \end{aligned}$$

where the second step comes from the relationship between the  $\ell_2$  norm and the  $\ell_\infty$  norm, and the third step comes from our definition that  $jj\mathbf{f}^{reg}jj_\infty \leq \psi$  for some  $\psi > 0$ .

Finally,

$$\begin{aligned} \mathbb{E}[jj(M_{post}^\top + Z_{post}^\top + W^\top)\mathbf{v}jj_2] & \leq \mathbb{E}[jjM_{post} + Z_{post} + Wjj_{2,2} jj\mathbf{v}jj_2] \\ & = \mathbb{E}[jjM_{post} + Z_{post} + Wjj_{2,2}] \mathbb{E}[jj\mathbf{v}jj_2] \\ & \leq \mathbb{E}[jjM_{post}jj_{2,2} + jjZ_{post}jj_{2,2} + jjWjj_{2,2}] \frac{4 \rho}{\lambda\epsilon_1 \overline{8+n}} \\ & \leq (jjM_{post}jj_{2,2} + \mathbb{E}[jjZ_{post}jj_F + jjWjj_F]) \frac{4 \rho}{\lambda\epsilon_1 \overline{8+n}} \\ & \leq jjM_{post}jj_{2,2} + \frac{\rho}{n(T - T_0)\sigma^2} + \frac{2^\rho \overline{T - T_0}}{\epsilon_2} \frac{4 \rho}{\lambda\epsilon_1 \overline{8+n}}, \end{aligned}$$

where the second step holds because  $Z_{post}$ ,  $W$  and  $\mathbf{v}$  are all independent of each other.

## F.3 Proof of Lemma C.6

**Lemma C.6** *Let  $\mathbf{f}^{reg} = (X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1}X_{pre}\mathbf{y}_{pre}$  be the Ridge regression coefficients and let  $\mathbf{f}$  be the true coefficients. If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \in (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t}$ , if  $T_0 \geq C(t/\xi)^2 k \log n$ , we have,*

$$\mathbb{E}[jj\mathbf{f}^{reg} - \mathbf{f}jj_2] \leq \frac{(\frac{\rho}{2n\sigma^2} + \frac{\rho}{2n\sigma^2 s^2})T_0 + \frac{\lambda}{2T_0}}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}.$$

First we can expand  $\mathbb{E}[\|j\mathbf{f}^{reg} - \mathbf{f}\|_2]$ :

$$\begin{aligned} \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbf{f}\|_2] &= \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}] + \mathbb{E}[\mathbf{f}^{reg}] - \mathbf{f}\|_2] \\ &= \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}]\|_2] + \mathbb{E}[\|\mathbb{E}[\mathbf{f}^{reg}] - \mathbf{f}\|_2] \\ &= \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}]\|_2] + \mathbb{E}[\|Bias(\mathbf{f}^{reg})\|_2], \end{aligned} \quad (23)$$

where

$$Bias(\mathbf{f}^{reg}) = \mathbb{E}[\mathbf{f}^{reg}] - \mathbf{f} = \lambda(X_{pre}X_{pre}^\top + \lambda I)^{-1} \mathbf{f}.$$

Hence, we only need to bound the two terms:  $\|jBias(\mathbf{f}^{reg})\|_2$  and  $\mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}]\|_2]$ , which we do next. First,

$$\begin{aligned} \|jBias(\mathbf{f}^{reg})\|_2 &= \|j\lambda(X_{pre}X_{pre}^\top + \lambda I)^{-1} \mathbf{f}\|_2 \\ &= \lambda \|j\mathbf{f}\|_2 \|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2} \\ &= \lambda \|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2}, \end{aligned}$$

where the last inequality uses the fact that the  $\ell_1$  norm of  $\mathbf{f}$  is 1, which also upper bound the  $\ell_2$  norm. Next,

$$\begin{aligned} \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}]\|_2] &= \mathbb{E}[\|j\mathbf{f}^{reg} - (\mathbf{f} + Bias(\mathbf{f}^{reg}))\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}X_{pre}\mathbf{y}_{pre} - \mathbf{f} + \lambda(X_{pre}X_{pre}^\top + \lambda I)^{-1} \mathbf{f}\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}X_{pre}(M_{pre}^\top \mathbf{f} + \mathbf{z}_{pre}) - \mathbf{f} + \lambda(X_{pre}X_{pre}^\top + \lambda I)^{-1} \mathbf{f}\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}X_{pre}(X_{pre}^\top \mathbf{f} - Z_{pre}^\top \mathbf{f} + \mathbf{z}_{pre}) - \mathbf{f} + \lambda(X_{pre}X_{pre}^\top + \lambda I)^{-1} \mathbf{f}\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}(X_{pre}X_{pre}^\top + \lambda I) \mathbf{f} - \mathbf{f} \\ &\quad + (X_{pre}X_{pre}^\top + \lambda I)^{-1}(X_{pre}\mathbf{z}_{pre} - X_{pre}Z_{pre}^\top \mathbf{f})\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}(X_{pre}\mathbf{z}_{pre} - X_{pre}Z_{pre}^\top \mathbf{f})\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2} \|jX_{pre}\mathbf{z}_{pre} - X_{pre}Z_{pre}^\top \mathbf{f}\|_2], \end{aligned}$$

where the first four steps come respectively from plugging in expressions of  $\mathbb{E}[\mathbf{f}^{reg}]$ ,  $(\mathbf{f}^{reg}$  and  $Bias(\mathbf{f}^{reg}))$ ,  $\mathbf{y}_{pre}$ , and  $M_{pre}$ . The fifth and sixth steps come from rearranging and canceling terms, and the final inequality comes from the submultiplicative norm property of Equation (9).

Plugging everything back to (23) yields,

$$\begin{aligned} \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbf{f}\|_2] &= \mathbb{E}[\|j\mathbf{f}^{reg} - \mathbb{E}[\mathbf{f}^{reg}]\|_2] + \mathbb{E}[\|jBias(\mathbf{f}^{reg})\|_2] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2} \|jX_{pre}\mathbf{z}_{pre} - X_{pre}Z_{pre}^\top \mathbf{f}\|_2 + \lambda \|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2}] \\ &= \mathbb{E}[\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2} (\|jX_{pre}\mathbf{z}_{pre} - X_{pre}Z_{pre}^\top \mathbf{f}\|_2 + \lambda)] \end{aligned} \quad (24)$$

Next, we use our assumptions on the data distribution to prove the following lemma about  $\|j(X_{pre}X_{pre}^\top + \lambda I)^{-1}\|_{2,2}$ .

**Lemma C.7** *If Assumptions C.2, C.3, and C.4 hold, then for all  $\xi \geq (0, 1)$  and  $t \geq 1$ , with probability at least  $1 - n^{-t^2}$  and  $T_0 \geq C(t/\xi)^2 k \log n$ , we have*

$$\|j(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1}\|_{2,2} \leq \frac{1}{(1 - \xi)T_0 + \frac{\lambda}{2T_0}}.$$

[Proof of Lemma C.7] A key component of the proof of Lemma C.7 is the following lemma about concentration of random matrices.

**Lemma C.8 (Corollary 5.52 of Vershynin (2010))** *Consider a distribution in  $\mathbb{R}^n$  with covariance matrix  $\Sigma$ , and supported in some centered Euclidean ball whose radius we denote  $\sqrt{m}$ . Let  $T_0$  be the number of samples and define the sample covariance matrix  $\Sigma_{T_0} = \frac{1}{T_0}XX^\top$ . Let  $\xi \geq (0, 1)$  and  $t \geq 1$ . Then with probability at least  $1 - n^{-t^2}$ , one has,*

$$\text{If } T_0 \geq C(t/\xi)^2 j\Sigma j_{2,2}^{-1} m \log n \text{ then } \|j\Sigma_{T_0} - \Sigma\|_{2,2} \leq \xi j\Sigma j_{2,2},$$

where  $C$  is an absolute constant.

To instantiate Lemma C.8, we view the data  $X_{pre}$  as  $T_0$  samples corresponding to the columns  $\mathbf{x}_t \in \mathbb{R}^n$ ,  $t \in \{1, 2, \dots, T_0\}$ . We use our assumptions that  $X$  takes values in a  $k$ -dimensional subspace  $E$ , and  $\Sigma = P_E$  where  $P_E$  is the orthogonal projection from  $\mathbb{R}^n$  onto  $E$ . Then, the *effective rank* of  $\Sigma$  is  $r(\Sigma) = \frac{\text{trace}(\Sigma)}{\|\Sigma\|_2} = k$  by definition, because  $\text{tr}(\Sigma) = \text{tr}(P_E) = k$  and  $\|\Sigma\|_2 = 1$ , since eigenvalues of an orthogonal projection matrix are either 0 or 1 as shown in Lemma 19 of Amjad et al. (2018). Then,  $\mathbb{E}[\text{tr}(X_{pre} X_{pre}^\top)] = \text{trace}(\Sigma) = k \text{tr}(\Sigma) = k \text{tr}(P_E) = k$ . Using Markov's inequality, most of the distribution should be within a ball of radius  $\frac{k}{m}$  where  $m = O(k)$ . Finally, let us assume that all the probability mass is within that ball, i.e.,  $\text{tr}(X_{pre} X_{pre}^\top) = O(k)$  almost surely. Then, Lemma C.8 holds with  $T_0 = C(t/\epsilon)^2 k \log n$  samples. This is also noted in Remark 5.53 of Vershynin (2010).

To translate this to our setting, we see that with probability at least  $1 - n^{-t^2}$ , if  $T_0 = C(t/\xi)^2 k \log n$ , then

$$\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - \Sigma\right) \leq \xi. \quad (25)$$

Since  $\Sigma = P_E$  is an orthogonal projection matrix,  $\text{tr}(P_E) = k$ . We apply triangle inequality to obtain,

$$\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - P_E\right) \leq \text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - \Sigma\right) + \text{tr}(\Sigma - P_E) = \text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - \Sigma\right) + 0 \leq \xi.$$

Combining this with Equation (25), we can bound

$$\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - I\right) \leq \xi, \quad \text{or equivalently,} \quad \text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top - T_0 I\right) \leq \xi T_0. \quad (26)$$

We will use this latter expression to obtain a lower bound on the minimum singular value of  $X_{pre} X_{pre}^\top$ , and then use it to bound  $\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top + \lambda I\right)^{-1}$  from above.

Note that since  $\text{tr}(A)$  is the sum of singular values of matrix  $A$ , the upper bound of  $\xi T_0$  of Equation (26) should hold for all singular values of  $A$ . For symmetric matrices such as  $X_{pre} X_{pre}^\top + T_0 I$ , the singular values are also the absolute values of its eigenvalues. This means that all eigenvalues  $\lambda_*$  of  $X_{pre} X_{pre}^\top + T_0 I$  must satisfy  $|\lambda_*| \leq \xi T_0$ . Therefore, this bound must also hold for the smallest eigenvalue  $\lambda_{\min}(\cdot)$ :

$$\begin{aligned} \lambda_{\min}(X_{pre} X_{pre}^\top + T_0 I) &\geq T_0 - \xi T_0 \\ \lambda_{\min}(X_{pre} X_{pre}^\top) &\geq T_0(1 - \xi) \end{aligned}$$

By plugging in the lower bound on the minimum singular value of  $X_{pre} X_{pre}^\top$ , we arrive at the desired bound to complete the proof of Lemma C.7.

$$\begin{aligned} \text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top + \lambda I\right)^{-1} &= \frac{1}{\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top + \lambda I\right)} \\ &= \frac{1}{\lambda_{\min}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top + \lambda I\right)} \\ &= \frac{1}{\lambda_{\min}(X_{pre} X_{pre}^\top) + \lambda} \\ &\leq \frac{1}{(1 - \xi)T_0 + \lambda}. \end{aligned}$$

Returning to Equation (24), we can use this bound to obtain,

$$\begin{aligned} \mathbb{E}[\text{tr}(\mathbf{f} \mathbf{f}^\top)] &\leq \mathbb{E}[\text{tr}\left(\frac{1}{T_0} X_{pre} X_{pre}^\top + \lambda I\right)^{-1} \text{tr}(X_{pre} Z^\top \mathbf{f} \mathbf{f}^\top X_{pre} Z + \lambda)] \\ &\leq \frac{1}{(1 - \xi)T_0 + \lambda} \mathbb{E}[\text{tr}(X_{pre} Z^\top \mathbf{f} \mathbf{f}^\top X_{pre} Z + \lambda)] \end{aligned} \quad (27)$$

The expectation term in Equation (27) becomes,

$$\begin{aligned}
 \mathbb{E}[jjX_{pre}z_{pre} \quad X_{pre}Z_{pre}^\top fjj_2 + \lambda] &= \mathbb{E}[jj(M_{pre} + Z_{pre})z_{pre} \quad (M + Z_{pre})Z_{pre}^\top fjj_2 + \lambda] \\
 &= \mathbb{E}[jjM_{pre}(z_{pre} \quad Z_{pre}^\top f)jj_2 + jjZ_{pre}(z_{pre} \quad Z_{pre}^\top f)jj_2 + \lambda] \\
 &= jjM_{pre}jj_{2,2} \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] + \mathbb{E}[jjZ_{pre}jj_{2,2} \quad jjz_{pre} \quad Z_{pre}^\top fjj_2] + \lambda \\
 &= jjM_{pre}jj_F \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] + \mathbb{E}[jjZ_{pre}jj_F \quad jjz_{pre} \quad Z_{pre}^\top fjj_2] + \lambda \\
 &= \rho \frac{1}{nT_0} \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] + \rho \frac{1}{nT_0s^2} \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] + \lambda \\
 &= \left( \rho \frac{1}{nT_0} + \rho \frac{1}{nT_0s^2} \right) \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] + \lambda,
 \end{aligned}$$

where the first step is plugging in for  $X_{pre}$ , the second step is triangle inequality, the third and fourth steps are due to the submultiplicative norm property, the fifth step comes from the definition of the Frobenius norm, the fact that  $M_{pre}$  and  $Z_{pre}$  are both of dimension  $n \times T_0$ , and bounds on data entries. The final step collects terms.

Finally, we need only to obtain a bound on  $\mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2]$ .

$$\begin{aligned}
 \mathbb{E}[jjz_{pre} \quad Z_{pre}^\top fjj_2] &= \mathbb{E} \left[ \sum_{t=1}^{T_0} \frac{1}{\mathcal{X}_0} (z_t \quad Z_t^\top f)^2 \right] \\
 &\stackrel{(a)}{\leq} \sum_{t=1}^{T_0} \frac{1}{\mathcal{X}_0} \mathbb{E}[(z_t \quad Z_t^\top f)^2] \\
 &= \sum_{t=1}^{T_0} \frac{1}{\mathcal{X}_0} \mathbb{E}[z_t^2 \quad 2z_t Z_t^\top f + (Z_t^\top f)^2] \\
 &\stackrel{(b)}{=} \sum_{t=1}^{T_0} \frac{1}{\mathcal{X}_0} (\sigma^2 + \mathbb{E}[(Z_t^\top f)^2]) \\
 &= \sum_{t=1}^{T_0} \frac{1}{T_0\sigma^2 + \sum_{i=1}^{T_0} \mathcal{X}_0^t \mathbb{E}[(z_i f_i)^2]} \\
 &\stackrel{(c)}{=} \sum_{t=1}^{T_0} \frac{1}{T_0\sigma^2 + \sum_{i=1}^{T_0} \mathcal{X}_0^t \mathbb{E}[z_i^2 f_i^2]} \\
 &= \sum_{t=1}^{T_0} \frac{1}{T_0\sigma^2 + \sum_{i=1}^{T_0} \mathcal{X}_0^t \sigma^2 f_i^2} \\
 &= \sum_{t=1}^{T_0} \frac{1}{T_0\sigma^2 + \sum_{i=1}^{T_0} \mathcal{X}_0^t \sigma^2 jjfjj_2^2} \\
 &\stackrel{(d)}{\leq} \sum_{t=1}^{T_0} \frac{1}{T_0\sigma^2 + \sum_{i=1}^{T_0} \mathcal{X}_0^t \sigma^2} \\
 &= \rho \frac{1}{2T_0\sigma^2}
 \end{aligned}$$

Inequality (a) is due to Jensen's inequality. The step in (b) is because  $\mathbb{E}[z_t Z_t^\top f] = \mathbb{E}[z_t] \mathbb{E}[Z_t^\top f] = 0$  by independence of noise terms. The step in (c) is by the same logic as in (b), since all cross-terms  $f_i f_j$  for  $i \neq j$  are zero in expectation. Lastly, we bound the  $\ell_2$  norm of  $f$  by  $\ell_1$  norm instead in (d) (i.e.,  $jjfjj_2 \leq jjfjj_1 \leq 1$ ).

Hence,

$$\begin{aligned}
 \mathbb{E}[jjX_{pre}z_{pre} \quad X_{pre}Z_{pre}^\top fjj_2 + \lambda] &= \left( \rho \frac{1}{nT_0} + \rho \frac{1}{nT_0s^2} \right) \rho \frac{1}{2T_0\sigma^2} + \lambda \\
 &= T_0 \frac{\rho}{2n\sigma^2} + T_0 \frac{\rho}{2n\sigma^2 s^2} + \lambda
 \end{aligned}$$

Finally, combining this with Equation (27) gives the desired bound to complete the proof of Lemma C.6.

$$\mathbb{E}[\|J\mathbf{f}^{reg} - J\mathbf{f}\|_2] \leq \frac{(\rho_{2n\sigma^2} + \rho_{2n\sigma^2s^2})T_0 + \lambda}{(1 - \xi)T_0 + \lambda}.$$

## G Omitted Proofs for $DPS_{Obj}$

### G.1 Proof of Lemma D.3

**Lemma D.3** For any  $\Delta \geq 0$ ,  $\Phi(\boldsymbol{\alpha}; \Delta) = \frac{|\det(\nabla \mathbf{b}(\boldsymbol{\alpha}; \mathcal{D}^\delta))|}{|\det(\nabla \mathbf{b}(\boldsymbol{\alpha}; \mathcal{D}))|} \left(1 + \frac{\Delta}{\lambda}\right)^2$ .

Recall that  $\mathbf{b}(\boldsymbol{\alpha}; D)$  is the noise value that must have been realized when database  $D$  was input and  $\boldsymbol{\alpha} = \arg \min_{\mathbf{f}} J^{obj}(\mathbf{f})$  was output. Since  $J^{obj}(\mathbf{f})$  is strongly convex for any  $\Delta$  and is differentiable, the closed-form expression for  $\mathbf{b}(\boldsymbol{\alpha}; D)$  is derived by computing the gradient of  $J^{obj}(\mathbf{f})$ , which should be zero when evaluated at its minimizer  $\mathbf{f} = \boldsymbol{\alpha}$ :

$$\nabla_{\mathbf{f}} J^{obj}(\mathbf{f})|_{\mathbf{f}=\boldsymbol{\alpha}} = r^2 L(\boldsymbol{\alpha}) + r^2 r(\boldsymbol{\alpha}) + \frac{\Delta}{T_0} \boldsymbol{\alpha} + \frac{\mathbf{b}(\boldsymbol{\alpha}; D)}{T_0} \stackrel{!}{=} \mathbf{0}.$$

Rearranging the equation yields

$$\mathbf{b}(\boldsymbol{\alpha}; D) = (T_0 r^2 L(\boldsymbol{\alpha}; D) + T_0 r^2 r(\boldsymbol{\alpha}) + \Delta \boldsymbol{\alpha}).$$

For ease of notation, let  $A = r^2 \mathbf{b}(\boldsymbol{\alpha}; D)$  and  $E = r^2 \mathbf{b}(\boldsymbol{\alpha}; D) - r^2 \mathbf{b}(\boldsymbol{\alpha}; D')$ . Then,

$$\Phi(\boldsymbol{\alpha}; \Delta) = \frac{j \det(r^2 \mathbf{b}(\boldsymbol{\alpha}; D'))}{j \det(r^2 \mathbf{b}(\boldsymbol{\alpha}; D))} = \frac{j \det(-r^2 \mathbf{b}(\boldsymbol{\alpha}; D'))}{j \det(r^2 \mathbf{b}(\boldsymbol{\alpha}; D))} = \frac{j \det(A + E)}{j \det(A)}.$$

By definition,  $A = r^2 \mathbf{b}(\boldsymbol{\alpha}; D) = T_0(r^2 L(\boldsymbol{\alpha}; D) + r^2 r(\boldsymbol{\alpha})) + \Delta I_n$ . Using the Hessians  $r^2 L(\boldsymbol{\alpha}; D) = \frac{2}{T_0} X_{pre} X_{pre}^\top$  and  $r^2 r(\boldsymbol{\alpha}) = \frac{\lambda}{T_0} I_n$ ,  $A$  can be expressed as

$$A = 2X_{pre} X_{pre}^\top + (\lambda + \Delta) I_n.$$

To express  $E$  succinctly, let neighboring databases  $D = (X, y)$  and  $D' = (X', y)$  differ in the  $j$ -th row. Then,

$$E = 2(X'_{pre} X_{pre}^\top - X_{pre} X_{pre}^\top) = \begin{cases} \geq 2(jj\mathbf{x}'_j j^2 - jj\mathbf{x}_j j^2) & (j, j) \\ \geq 2(\mathbf{x}'_j - \mathbf{x}_j)^\top \mathbf{x}_i & (j, i) \text{ or } (i, j), \delta i \geq 2 [n], i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

where  $\mathbf{x}_i$  (resp.  $\mathbf{x}'_i$ ) denotes the  $i$ -th person's data, which is the  $i$ -th row of  $X_{pre}$  (resp.  $X'_{pre}$ ).

Note that all eigenvalues of  $A$  are at least  $\lambda + \Delta > 0$  (i.e.,  $\lambda_{\min}(A) \geq \lambda + \Delta$ ) because  $X_{pre} X_{pre}^\top$  is positive-semi-definite, and thus  $A$  is full rank. Also,  $\text{rank}(E) = 2$ . This allows us to apply the following lemma.

**Lemma G.1 (Lemma 2 of Chaudhuri et al. (2011))** If  $A$  is full rank and  $E$  has rank at most 2,

$$\frac{\det(A + E)}{\det(A)} = \lambda_1(A^{-1}E) + \lambda_2(A^{-1}E) + \lambda_1(A^{-1}E)\lambda_2(A^{-1}E),$$

where  $\lambda_i(Z)$  is  $i$ -th eigenvalue of matrix  $Z$ .

Let  $\lambda_{|max|}(Z) = \max_i |\lambda_i(Z)|$ , the maximum absolute of eigenvalue of matrix  $Z$ . Instantiating Lemma G.1 yields:

$$\begin{aligned} \Phi(\boldsymbol{\alpha}; \Delta) &= \frac{j \det(A + E)}{j \det(A)} \\ &= \frac{\det(A + E)}{\det(A)} = 1 \\ &= j[1 + \lambda_1(A^{-1}E) + \lambda_2(A^{-1}E) + \lambda_1(A^{-1}E)\lambda_2(A^{-1}E)] \\ &= j[1 + j\lambda_1(A^{-1}E) + j\lambda_2(A^{-1}E) + j\lambda_1(A^{-1}E)\lambda_2(A^{-1}E)] \\ &= 1 + 2\lambda_{|max|}(A^{-1}E) + \lambda_{|max|}(A^{-1}E)^2, \end{aligned}$$

where the first inequality is simply triangle inequality, and the second inequality bounds all absolute eigenvalues by the maximum one  $\lambda_{|max|}$ .

Assume that  $\lambda_{|max|}(E) \leq c$  for some constant  $c$ . Since  $E$  is a real-valued matrix, such a finite  $c$  exist. In Algorithm 3,  $c$  is explicitly taken as an input parameter. Then,

$$\lambda_{|max|}(A^{-1}E) \leq \frac{\lambda_{|max|}(E)}{\lambda_{min}(A)} \leq \frac{c}{\lambda + \Delta}.$$

Finally,

$$\Phi(\alpha; \Delta) \leq 1 + 2\lambda_{|max|}(A^{-1}E) + \lambda_{|max|}(A^{-1}E)^2 \leq 1 + \frac{2c}{\lambda + \Delta} + \frac{c^2}{(\lambda + \Delta)^2} \leq 1 + \frac{c}{\lambda + \Delta}.$$

## G.2 Proof of Lemma D.4

**Lemma D.4** When  $\mathbf{b}$  is sampled according to pdf  $p(\mathbf{b}; \beta) \propto \exp\left(-\frac{\|\mathbf{b}\|_2}{\beta}\right)$ , where  $\beta = \min\left\{\frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n}+4T_0}{\epsilon_0}\right\}$ , then

$$\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; D))}{\Pr(\mathbf{b}(\alpha; D'))} \leq e^{\epsilon_0}.$$

We can start by re-writing  $\Gamma(\alpha)$  as follows, where the first line directly comes from the pdf  $\Pr(\mathbf{b}; \beta)$ , the second line is due to reverse triangle inequality, and the third line is from the definition of  $\mathbf{b}(\alpha; D)$  and canceling terms that occur in both  $\mathbf{b}(\alpha; D)$  and  $\mathbf{b}(\alpha; D')$ :

$$\begin{aligned} \Gamma(\alpha) &= \exp\left(-\frac{1}{\beta} \|\mathbf{b}(\alpha; D)\|_2 + \frac{1}{\beta} \|\mathbf{b}(\alpha; D')\|_2\right) \\ &= \exp\left(\frac{1}{\beta} \|\mathbf{b}(\alpha; D) - \mathbf{b}(\alpha; D')\|_2\right) \\ &= \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2\right). \end{aligned} \quad (29)$$

Next, we can continue to bound Equation (29) in two different ways, corresponding to the two possible values of  $\beta$ . The two values come from two different upper bounds on the sensitivity, and the minimum value will give a tighter bound.

The first upper bound uses Lemma B.4, and its notation of  $g(\mathbf{f}) = L(\mathbf{f}, D') - L(\mathbf{f}, D)$  for neighboring databases  $D, D'$ . Then we can bound:

$$\begin{aligned} (29) &\leq \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2\right) \\ &\leq \exp\left(\frac{1}{\beta} 4T_0 \sqrt{8+n}\right). \end{aligned}$$

Hence, setting  $\beta = \frac{4T_0\sqrt{8+n}}{\epsilon_0}$  makes  $\Gamma(\alpha) \leq e^{\epsilon_0}$ .

The second upper bound is based on  $c$ , and will yield a tighter bound when  $c$  is small. Recall that matrix  $E$  is defined in Equation (28), and that  $c$  is the upper bound  $\lambda_{|max|}(E) \leq c$ . By plugging in  $\mathbf{r}(\alpha) = \frac{1}{T_0} (2X_{pre}X_{pre}^\top \alpha - 2X_{pre}\mathbf{y}_{pre})$ , we can alternatively bound:

$$\begin{aligned} (29) &= \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2\right) \\ &= \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2\right) + \frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2 \\ &= \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2 + \frac{4T_0}{\beta}\right) \\ &= \exp\left(\frac{1}{\beta} \|\mathbf{r}(\alpha; D') - \mathbf{r}(\alpha; D)\|_2 + \frac{4T_0}{\beta}\right) \\ &= \exp\left(\frac{c\sqrt{n} + 4T_0}{\beta}\right), \end{aligned}$$



where the second step is due to triangle inequality, the third step is plugging in the definition of  $E$  and bounding the second term based on the worst-case  $X'_{pre} = X_{pre}$ , which is all zeros with just one row with all 2's, and worst-case  $y_{pre}$ , which is all 1's). The fourth step is the submultiplicative property of operator norms, and the final step is due to the fact that  $\|E\|_{2,2} = \lambda_{|max|}(E) = c$  and that all elements of  $\alpha \in [0, 1]^n$  are bounded by 1. Then setting  $\beta = \frac{c\sqrt{n}+4T_0}{\epsilon_0}$  ensures  $\Gamma(\alpha) \leq e^{\epsilon_0}$ .

If either of the above conditions on  $\beta$  holds, then  $\Gamma(\alpha) \leq e^{\epsilon_0}$  as desired. Thus we can choose  $\beta = \min\left\{\frac{4T_0\sqrt{8+n}}{\epsilon_0}, \frac{c\sqrt{n}+4T_0}{\epsilon_0}\right\}$  that at least one will be satisfied. Taking the minimum rather than just one allows for a lower  $\beta$  and hence lower noise magnitude, while still satisfying the privacy requirement.

### G.3 Proof of Lemma D.5

**Lemma D.5** When  $\mathbf{b} \sim N(0, \beta^2 I_n)$ , where  $\beta = \frac{4T_0\sqrt{8+n}}{\epsilon_0}$ , then  $\Gamma(\alpha) = \frac{\Pr(\mathbf{b}(\alpha; \mathcal{D}))}{\Pr(\mathbf{b}(\alpha; \mathcal{D}'))} \leq e^{\epsilon_0}$  with probability at least  $1 - \delta$ .

The proof of Lemma D.5 follows a similar structure to Lemma 14 of Kifer et al. (2012). We include the full proof for completeness. Let the noise term  $\mathbf{b}$  be sampled from a multivariate Gaussian distribution  $N(0, \beta^2 I_n)$ , and let  $D$  and  $D'$  be two arbitrary neighboring databases. Let  $h(\alpha) = \mathbf{b}(\alpha; D') - \mathbf{b}(\alpha; D)$ . Then, we can express  $\Gamma(\alpha)$  as,

$$\begin{aligned} \Gamma(\alpha) &= \frac{\exp\left(-\frac{\|\mathbf{b}(\alpha; \mathcal{D})\|_2^2}{2\beta^2}\right)}{\exp\left(-\frac{\|\mathbf{b}(\alpha; \mathcal{D}')\|_2^2}{2\beta^2}\right)} \\ &= \exp\left(\frac{1}{2\beta^2}(\|\mathbf{b}(\alpha; D')\|_2^2 - \|\mathbf{b}(\alpha; D)\|_2^2)\right) \\ &= \exp\left(\frac{1}{2\beta^2}(\|\mathbf{b}(\alpha; D) + h(\alpha)\|_2^2 - \|\mathbf{b}(\alpha; D)\|_2^2)\right) \\ &= \exp\left(\frac{1}{2\beta^2}(2\mathbf{b}(\alpha; D), h(\alpha)) + \|h(\alpha)\|_2^2\right), \end{aligned} \quad (30)$$

where the first step is from the distribution of noise  $\mathbf{b}$ , the final step is a binomial expansion applied to norms.

Note that,

$$\begin{aligned} h(\alpha) &= \mathbf{b}(\alpha; D') - \mathbf{b}(\alpha; D) \\ &= T_0(r L(\alpha; D') - r L(\alpha; D)) \\ &= T_0 r g(\alpha), \end{aligned}$$

where  $g(\alpha) = L(\alpha; D') - L(\alpha; D)$ , as defined in Equation (5). By Lemma B.4, we know that  $\|r g(\alpha)\|_2 \leq \frac{4\sqrt{8+n}}{\epsilon_0}$ , so also

$$\|r h(\alpha)\|_2 \leq 4T_0 \frac{4\sqrt{8+n}}{\epsilon_0}. \quad (31)$$

Similarly, because  $\mathbf{b}$  is sampled from a multivariate Gaussian distribution  $N(0, \beta^2 I_n)$  and sum of Gaussian variables is also Gaussian, then,

$$\mathbf{b}(\alpha; D), h(\alpha) \sim N(0, \beta^2 \|r h(\alpha)\|_2^2).$$

Since the exact distribution is known, we use a Gaussian tail bound to find a *well-behaving* set of  $\mathbf{b}$ .

**Lemma G.2 (Chernoff bound for Gaussian Wainwright (2019))** Let  $Z \sim N(0, \sigma^2)$ . Then, for all  $t > \sigma$ ,

$$P[Z \geq t] \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

We instantiate Lemma G.2 with  $Z = \mathbf{b}(\alpha; D), h(\alpha)$  and  $t = \beta \|r h(\alpha)\|_2$ . Note that  $t > \sigma$  for any  $\delta > 1/2$ .

Then,

$$\Pr\left[\mathbf{b}(\alpha; D), h(\alpha) \geq \beta \|r h(\alpha)\|_2\right] \leq \frac{\delta}{2},$$

which, by Equation (31) implies that,

$$\Pr[\|\mathbf{h}(\boldsymbol{\alpha}; D) - \mathbf{h}(\boldsymbol{\alpha})\|_2 \leq \beta(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}}] \geq 1 - \delta. \quad (32)$$

Define a set of values of  $\mathbf{b}$ , corresponding to the good event described by Equation (32):  $\mathbf{GOOD} = \{\mathbf{b} \mid \|\mathbf{h}(\boldsymbol{\alpha}; D) - \mathbf{h}(\boldsymbol{\alpha})\|_2 \leq \beta(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}}\}$ . By definition,  $\Pr[\mathbf{b} \in \mathbf{GOOD}] \geq 1 - \delta$ . That is, with probability at least  $1 - \delta$ , the noise vector  $\mathbf{b}$  is in the well-behaving set  $\mathbf{GOOD}$ .

When  $\mathbf{b} \in \mathbf{GOOD}$ , then we can complete the bound on  $\Gamma(\boldsymbol{\alpha})$  from Equation (30), combining the bound on  $\|\mathbf{h}(\boldsymbol{\alpha}) - \mathbf{h}(\boldsymbol{\alpha}; D)\|_2$  from Equation (31):

$$\Gamma(\boldsymbol{\alpha}) = \exp\left(\frac{1}{2\beta^2} [2\|\mathbf{h}(\boldsymbol{\alpha}; D) - \mathbf{h}(\boldsymbol{\alpha})\|_2 + \|\mathbf{h}(\boldsymbol{\alpha})\|_2]^2\right) \leq \exp\left[\frac{1}{2\beta^2} [2\beta(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}} + (4T_0)^{\frac{1}{8+n}}]^2\right].$$

Finally, the goal is to bound  $\Gamma(\boldsymbol{\alpha}) \leq e^{\epsilon_0}$ , in the case where  $\mathbf{b} \in \mathbf{GOOD}$ . Solving the expression above for  $\beta$  yields

$$\begin{aligned} \beta &\geq \frac{1}{2} \sqrt{\frac{(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}}}{\epsilon_0} + \frac{(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}}}{\epsilon_0^2} + \frac{(4T_0)^{\frac{1}{8+n}}}{\epsilon_0}} \\ &= \frac{1}{2} \sqrt{\frac{4T_0^{\frac{1}{8+n}}}{\epsilon_0} \sqrt{2 \log \frac{2}{\delta}} + \frac{4T_0^{\frac{1}{8+n}}}{\epsilon_0} \sqrt{2 \log \frac{2}{\delta}} + \epsilon_0} \end{aligned} \quad (33)$$

Note that choosing

$$\beta = \frac{(4T_0)^{\frac{1}{8+n}} \sqrt{2 \log \frac{2}{\delta}}}{\epsilon_0}$$

satisfies the bound of Equation (33).

Thus  $\Gamma(\boldsymbol{\alpha}) \leq e^{\epsilon_0}$ , conditioned on  $\mathbf{b} \in \mathbf{GOOD}$ , which occurs with probability at least  $1 - \delta$ .

#### G.4 Proof of Lemma E.1

**Lemma E.1** *The  $\ell_2$  distance between  $\mathbf{f}^{obj}$  and  $\mathbf{f}^{reg}$  satisfies:*

$$\|\mathbf{f}^{obj} - \mathbf{f}^{reg}\|_2 \leq \frac{2}{\lambda + \Delta} \mathbb{E}[\|\mathbf{b}\|_2] + 1 \leq \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \sqrt{2T_0^{\frac{1}{8+n}}},$$

where  $\mathbf{b}$  and  $\Delta$  are computed internally by Algorithm 3.

Recall the objective functions  $J^{obj}$  and  $J^{reg}$ :

$$J^{obj}(\mathbf{f}) = L(\mathbf{f}) + \frac{\lambda + \Delta}{2T_0} \|\mathbf{f}\|_2^2 + \frac{1}{T_0} \mathbf{b}^\top \mathbf{f} \quad \text{and} \quad J^{reg}(\mathbf{f}) = L(\mathbf{f}) + \frac{\lambda}{2T_0} \|\mathbf{f}\|_2^2,$$

with their respective minimizers  $\mathbf{f}^{obj}$  and  $\mathbf{f}^{reg}$ . Define another objective function  $J^\#$  and its minimizer  $\mathbf{f}^\#$ ,

$$J^\#(\mathbf{f}) = L(\mathbf{f}) + \frac{\lambda + \Delta}{2T_0} \|\mathbf{f}\|_2^2$$

which is a noise-free variant of  $J^{obj}$ .

We will express the difference between  $\mathbf{f}^{reg}$  and  $\mathbf{f}^{obj}$  using  $\mathbf{f}^\#$  as an intermediate value:

$$\|\mathbf{f}^{reg} - \mathbf{f}^{obj}\|_2 = \|\mathbf{f}^{reg} - \mathbf{f}^\# + \mathbf{f}^\# - \mathbf{f}^{obj}\|_2 \leq \|\mathbf{f}^{reg} - \mathbf{f}^\#\|_2 + \|\mathbf{f}^\# - \mathbf{f}^{obj}\|_2. \quad (34)$$

We will bound these two terms separately, starting with  $k\mathbf{f}^\# - \mathbf{f}^{obj}k_2$ . It is known that  $J^{obj}$  is  $(\frac{\lambda+\Delta}{T_0})$ -strongly convex, and that the gradient of  $J^{obj}$  evaluated at its minimizer  $\mathbf{f}^{obj}$  is zero. Then by the definition of strong convexity,

$$k\mathbf{f}^\# - \mathbf{f}^{obj}k_2 \leq J^{obj}(\mathbf{f}^\#) - J^{obj}(\mathbf{f}^{obj}) \leq \frac{2T_0}{\lambda + \Delta}. \quad (35)$$

We can proceed to bound the difference in the objective function  $J^{obj}$  at these two points:

$$\begin{aligned} J^{obj}(\mathbf{f}^\#) - J^{obj}(\mathbf{f}^{obj}) &= J^\#(\mathbf{f}^\#) + \frac{1}{T_0}b^\top \mathbf{f}^\# - J^\#(\mathbf{f}^{obj}) + \frac{1}{T_0}b^\top \mathbf{f}^{obj} \\ &= J^\#(\mathbf{f}^\#) - J^\#(\mathbf{f}^{obj}) + \frac{1}{T_0}b^\top \mathbf{f}^\# - \frac{1}{T_0}b^\top \mathbf{f}^{obj} \\ &= 0 + \frac{1}{T_0}k\mathbf{b}k_2k\mathbf{f}^\# - \mathbf{f}^{obj}k_2 \end{aligned}$$

where the inequality is due to the fact that  $J^\#(\mathbf{f}^\#) \leq J^\#(\mathbf{f}^{obj})$ , since  $\mathbf{f}^\#$  is the minimizer of  $J^\#$ .

Plugging this into Equation (35) gives

$$k\mathbf{f}^\# - \mathbf{f}^{obj}k_2 \leq \frac{1}{T_0}k\mathbf{b}k_2k\mathbf{f}^\# - \mathbf{f}^{obj}k_2 \leq \frac{2T_0}{\lambda + \Delta},$$

or equivalently,

$$k\mathbf{f}^\# - \mathbf{f}^{obj}k_2 \leq k\mathbf{b}k_2 \frac{2}{\lambda + \Delta}.$$

To bound the first term of Equation (34), we observe that if  $\Delta = 0$ , then  $J^\# = J^{reg}$  and thus  $\mathbf{f}^\# = \mathbf{f}^{reg}$ , so  $k\mathbf{f}^{reg} - \mathbf{f}^\#k_2 = 0$ . Thus we only need to bound the distance when  $\Delta \neq 0$ .

We can write  $\mathbf{f}^{reg}$  and  $\mathbf{f}^\#$  using their closed-form expressions,

$$\mathbf{f}^{reg} = (X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1}X_{pre}\mathbf{y}_{pre} \quad \text{and} \quad \mathbf{f}^\# = (X_{pre}X_{pre}^\top + \frac{\lambda + \Delta}{2T_0}I)^{-1}X_{pre}\mathbf{y}_{pre},$$

and use these to bound the difference:

$$\begin{aligned} k\mathbf{f}^{reg} - \mathbf{f}^\#k_2 &= k(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1} - (X_{pre}X_{pre}^\top + \frac{\lambda + \Delta}{2T_0}I)^{-1}X_{pre}\mathbf{y}_{pre}k_2 \\ &= (k(X_{pre}X_{pre}^\top + \frac{\lambda}{2T_0}I)^{-1}k_2 + k((X_{pre}X_{pre}^\top + \frac{\lambda + \Delta}{2T_0}I)^{-1})^{-1}k_2)kX_{pre}\mathbf{y}_{pre}k_2 \end{aligned} \quad (36)$$

The spectral norm of a general form  $k(X_{pre}X_{pre}^\top + \lambda I)^{-1}k_2$  can be bounded by the inverse of minimum singular value of the matrix  $X_{pre}X_{pre}^\top + \lambda I$ , which is positive semi-definite and has minimum singular value at least  $\lambda$ :

$$k(X_{pre}X_{pre}^\top + \lambda I)^{-1}k_2 \leq \frac{1}{\sigma_{min}(X_{pre}X_{pre}^\top + \lambda I)} \leq \frac{1}{\lambda}.$$

Using this fact, we can further bound Equation (36) as,

$$\begin{aligned} k\mathbf{f}^{reg} - \mathbf{f}^\#k_2 &\leq \frac{2T_0}{\lambda} + \frac{2T_0}{\lambda + \Delta}kX_{pre}\mathbf{y}_{pre}k_2 \\ &\leq 2T_0 \left( \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \right) \|X_{pre}\|_F \|\mathbf{y}_{pre}\|_2 \\ &\leq 2T_0 \left( \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \right) \sqrt{\frac{1}{nT_0}} \sqrt{\frac{1}{T_0}} \\ &= \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \sqrt{2T_0^2 \frac{1}{n}}. \end{aligned}$$

Finally, we combine Equation (34) with bounds on both terms to yield:

$$\begin{aligned} \mathbb{E}[k\mathbf{f}^{reg} - \mathbf{f}^{obj}k_2] &\leq \mathbb{E}[k\mathbf{f}^\# - \mathbf{f}^{obj}k_2] + \mathbb{E}[k\mathbf{f}^{reg} - \mathbf{f}^\#k_2] \\ &\leq \frac{2}{\lambda + \Delta} \mathbb{E}[k\mathbf{b}k_2] + \mathbb{E}[k\mathbf{f}^\# - \mathbf{f}^{obj}k_2] + \frac{1}{\lambda} + \frac{1}{\lambda + \Delta} \sqrt{2T_0^2 \frac{1}{n}}. \end{aligned}$$