
Optimal Differentially Private Model Training with Public Data

Andrew Lowy¹ Zeman Li² Tianjian Huang² Meisam Razaviyayn²

Abstract

Differential privacy (DP) ensures that training a machine learning model does not leak private data. In practice, we may have access to auxiliary public data that is free of privacy concerns. In this work, we assume access to a given amount of public data and settle the following fundamental open questions: 1. *What is the optimal (worst-case) error of a DP model trained over a private data set while having access to side public data?* 2. *How can we harness public data to improve DP model training in practice?* We consider these questions in both the local and central models of pure and approximate DP. To answer the first question, we prove tight (up to log factors) lower and upper bounds that characterize the optimal error rates of three fundamental problems: mean estimation, empirical risk minimization, and stochastic convex optimization. We show that the optimal error rates can be attained (up to log factors) by either discarding private data and training a public model, or treating public data like it is private and using an optimal DP algorithm. To address the second question, we develop novel algorithms that are “even more optimal” (i.e. better constants) than the asymptotically optimal approaches described above. For local DP mean estimation, our algorithm is optimal including constants. Empirically, our algorithms show benefits over the state-of-the-art.

1. Introduction

Training machine learning models on people’s data can leak sensitive information, violating their privacy (Fredrikson et al., 2015; Shokri et al., 2017; Carlini et al., 2021). *Differential Privacy (DP)* prevents such leaks by providing

¹University of Wisconsin-Madison, Wisconsin Institute of Discovery, Madison, WI, USA ²Department of Industrial & Systems Engineering, University of Southern California, Los Angeles, CA, USA. Correspondence to: Andrew Lowy <alowy@wisc.edu>.

a rigorous guarantee that no attacker can learn too much about any individual’s data (Dwork et al., 2006). DP has been successfully deployed by various companies (Apple, 2016; Thakurta et al., 2017; Úlfar Erlingsson et al., 2014; Ding et al., 2017), and by government agencies (U.S. Census Bureau, 2020). However, a major hindrance to more widespread adoption of DP is that DP-trained models are less accurate than their non-private counterparts.

Leveraging *public data*—that is free of privacy concerns—appears to be a promising and practically important avenue for closing the accuracy gap between DP and non-private models (Papernot et al., 2017; Avent et al., 2017; Feldman et al., 2018; Amid et al., 2022). For example, large language models (LLMs) are often pre-trained on public data and fine-tuned on private data (Kerrigan et al., 2020b; Li et al., 2021b; Yu et al., 2021a). Public data may be provided by people who volunteer (e.g. product developers or early testers) (Church, 2005; Feldman et al., 2018) or sell their data. Data that is generated synthetically (Torkzadehmahani et al., 2019; Vietri et al., 2020; Boediardjo et al., 2022; He et al., 2023) or released through a legal process (Klimt & Yang, 2004) may serve as additional sources of public data.

The power and limitations of public data vary depending on the particular learning problem and loss function/hypothesis class. To calibrate the effectiveness of a public-data-assisted DP algorithm, we can compare it against two *naïve baselines*: 1) “throw away” the private data and run an optimal non-private algorithm on the public data; 2) use an optimal DP algorithm on the full data set, treating the public data like it is private data. Some works have identified problems where significant improvements over the naïve approaches are possible. For example, (Alon et al., 2019) show that for agnostic PAC learning with a hypothesis class of finite VC-dimension, it is possible to achieve asymptotically smaller sample complexity than the naïve baselines. (Bassily et al., 2020) show a similar result for private query release with finite VC-dimension. On the other hand, for certain problems it is impossible to do better than the naïve approaches: e.g., releasing binary decision stumps (Bassily et al., 2020). Understanding what improvements, if any, over the naïve approaches are possible for other problems (e.g. optimization) and function classes is interesting.

This work considers DP *model training* with side access to

public data: given a loss function, find model parameters to approximately minimize the expected training or test loss. We consider *empirical risk minimization* (ERM) and *stochastic convex optimization* (SCO), which correspond to minimizing training loss and expected test loss, respectively. For population mean estimation and SCO, we assume access to in-distribution public data. For ERM, the public data may be out-of-distribution. We answer a fundamental question:

Question 1. What is the optimal (minimax) error of DP model training with side access to public data? Is it possible to achieve smaller error than the naïve baselines?

Contribution 1: Limitations of Public Data for DP Model Training To answer **Question 1**, we characterize the optimal minimax error (up to constants or logarithms) of *semi-DP* (Beimel et al., 2013; Alon et al., 2019) algorithms—algorithms that are DP w.r.t. private data, but not necessarily DP w.r.t. public data (Definition 3). We provide tight lower and upper bounds for three fundamental training problems: mean estimation of bounded random variables¹, ERM and SCO with Lipschitz convex functions.² We consider both the *local* (Kasiviswanathan et al., 2011; Duchi et al., 2013) and *central* (Dwork et al., 2006) models of *pure* ($\delta = 0$) and *approximate* ($\delta \neq 0$) semi-DP. We prove nine sets of lower and upper bounds: see Figs. 1, 2 and 7. Our lower bounds imply that *it is impossible to obtain asymptotic improvements over the naïve approaches for semi-DP model training in the worst case.*

In light of these negative results, it is natural to wonder whether/how one can harness public data for more effective DP model training. This leads us to **Question 2**:

Question 2. Can we provide improved performance (e.g. theoretically smaller error and/or superior empirical results) over the naïve baselines?

Some prior works have tackled **Question 2** by imposing additional assumptions and/or shrinking the problem class. For example, (Amid et al., 2022) shows that under certain distributional assumptions, public data permits benefits over DP-SGD in linear regression. Also, (Zhou et al., 2020; Kairouz et al., 2021) show that by imposing certain “low-rank subspace” assumptions on the gradients in DP-SGD, public data can help attain dimension-independent rates for DP ERM. In contrast, we consider **Question 2** *without imposing any additional assumptions and without shrinking the loss function/data distribution class.*

Contribution 2: Power of Public Data for DP Model Training To address **Question 2**, we develop novel (central and local) semi-DP algorithms that add less noise than

¹Our ϵ -semi-DP analysis extends to unbounded/heavy-tailed distributions with bounded k -th order moment.

²Our results for ERM also cover non-convex loss functions.

would be necessary to privatize the full data set (including public). By doing so, we can achieve *optimal (worst-case) error bounds with significantly improved constants* over the asymptotically optimal naïve algorithms. Our local semi-DP mean estimation algorithm is optimal including constants.

We complement our theoretical analyses with extensive numerical experiments. Our experiments show that *our algorithms outperform the naïve approaches, even when the optimal DP algorithm is pre-trained on the public data.* For example, our Algorithm 1 achieves a significant improvement in CIFAR-10 image classification tasks, reducing test error by 8–9% for logistic regression and by at most 18.9% for Wide-ResNet, compared with the naïve approaches. We also identify a linear regression problem in which *DP-SGD diverges, but our algorithm converges* with small error using $n_{\text{pub}} = 0.1n$ public samples: see Figure 6. Moreover, *our algorithm consistently outperforms the state-of-the-art public-data-assisted mirror-descent (PDA-MD) of (Amid et al., 2022).*

1.1. Techniques and Challenges

Lower bounds: We develop and utilize a variety of techniques to prove our lower bounds.

To prove our central ϵ, δ, q -semi-DP SCO lower bound in Theorem 12, we build upon the techniques of Dwork et al. (2015); Bassily et al. (2020). A key challenge is finding a distribution whose mean has small norm and proving that this distribution is still hard enough for any semi-DP algorithm to estimate in ℓ_2 -distance. To accomplish this, we make three main innovations: First, we modify the *tracing attack* of (Dwork et al., 2015) by incorporating more aggressive truncation; this is used to infer membership of many individuals in the data set, even under weak ℓ_2 -accuracy guarantees. Second, we construct a novel distribution in which the prior (mean) is drawn from a shorter interval than the posterior (data). This innovation is crucial for obtaining our tight bounds, but also complicates the analysis. Thus, we provide a novel generalization of the *fingerprinting lemma* (Bun et al., 2014; 2017; Kamath et al., 2022b) that permits analysis of our construction. We provide more details on these proofs in Appendix E.

For our central pure ϵ -semi-DP population lower bounds in Theorems 32 and 42, we develop a *Semi-DP Fano’s method* (Theorem 33). In combination with the reduction from estimation to testing, Theorem 33 generalizes DP Fano’s method (Acharya et al., 2021) and strengthens classical Fano’s method (Yu, 1997). We build on the tools of Barber & Duchi (2014) in our proof of Theorem 33.

Our ERM lower bound (Theorem 10) uses a novel semi-DP *packing argument*: we construct $2^{d/2}$ “hard” data sets with well-separated sample means, and use the group privacy

Optimal Differentially Private Model Training with Public Data

Learning problem	Semi-DP error	When is semi-DP error less than DP error?	Learning problem	Semi-LDP error	When is semi-LDP error less than LDP?
Mean Estimation (Pop. MSE)	$\min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n^2 \varepsilon^2} + \frac{1}{n} \right\}$ (Theorem 4)	$n_{\text{pub}} > \frac{n^2 \varepsilon^2}{d}$ or $n_{\text{pub}} = \Theta(n)$	Mean Estimation (Pop. MSE)	$\min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n \varepsilon^2} \right\}$ (Theorem 14 & Remark 15)	$n_{\text{pub}} > \frac{\varepsilon^2 n}{d}$ or $n_{\text{pub}} = \Theta(n)$
SCO (Excess pop. risk)	$\min \left\{ \frac{1}{\sqrt{n_{\text{pub}}}}, \frac{\sqrt{d}}{n \varepsilon} + \frac{1}{\sqrt{n}} \right\}$ (Theorem 12)	$n_{\text{pub}} > \frac{n^2 \varepsilon^2}{d}$ or $n_{\text{pub}} = \Theta(n)$	SCO (Excess pop. risk)	$\min \left\{ \frac{1}{\sqrt{n_{\text{pub}}}}, \sqrt{\frac{d}{n \varepsilon^2}} \right\}$ (Theorem 18 & Remark 15)	$n_{\text{pub}} > \frac{\varepsilon^2 n}{d}$ or $n_{\text{pub}} = \Theta(n)$

Figure 1. Minimax optimal error rates for central $\rho\varepsilon, \delta q$ -semi-DP (up to logs) and (local) $\rho\varepsilon, \delta q$ -semi-LDP. $n = n_{\text{priv}} + n_{\text{pub}}$, where n_{priv} (n_{pub}) denotes the number of private (public) samples. Dependence on δ , range and Lipschitz parameters, constraint set diameter omitted. See Appendix for strongly convex SCO results.

Learning problem	Semi-DP error	When is semi-DP error less than DP error?
ERM (Excess emp. Risk)	$\min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n \varepsilon} \right\}$ (Theorem 10)	$n_{\text{pub}} > n - d/\varepsilon$

Figure 2. Minimax optimal error rates for ε -semi-DP ERM. See Table 7 in Appendix for more ε -semi-(L)DP results (e.g. SCO).

property of DP to show that any semi-DP algorithm must make large error on at least one of these data sets. Such arguments have long been used to prove pure DP lower bounds (Hardt & Talwar, 2010; Bassily et al., 2014). However, to the best of our knowledge, our proof is the first to extend packing arguments to the semi-DP setting. The main challenge is in carefully constructing the private and public data sets so as to force semi-DP A to make large error, even when A has access to public data.

For our local semi-DP lower bounds (Theorems 14 and 18), we build on the sophisticated techniques of Duchi & Rogers (2019). The main idea of our proofs is to combine Assouad’s method (Duchi, 2021) with bounds on the mutual information between the input and output of semi-DP algorithms.

Algorithms: We develop novel algorithms to obtain smaller error (including constants) than the asymptotically optimal naïve algorithms. Our algorithms are simple. For example, we propose a central $\rho\varepsilon, \delta q$ -semi-DP estimator that puts different weights on the private and public samples and adds (Gaussian) noise that is calibrated to the private weight. We choose the weights depending on the privacy level, dimension, and number of public samples to trade off smaller sensitivity (less noise) with larger variance on the public data. An optimal choice of weights minimizes the variance of our unbiased estimator, leading to smaller error than the optimal DP estimator. Our local semi-DP algorithm simply applies the optimal local randomizer of Bhowmick et al. (2018) only to the private samples, and averages the

noisy private data with the raw public data. For SCO, we develop semi-DP stochastic gradient methods that use our mean estimation algorithms to estimate the gradient of the loss function in each iteration of training.

1.2. Preliminaries and Notation

Let $\|\cdot\|$ be the ℓ_2 norm and X denote a data universe. Function $g : W \rightarrow \mathbb{R}$ is μ -strongly convex if $g(w) \geq \alpha \|w - w^*\|^2$ for all $w, w^* \in W$. If $\mu = 0$, we say g is convex. Function $f : W \rightarrow \mathbb{R}$ is uniformly L -Lipschitz in w if $\sup_{x, x'} |f(w, x) - f(w', x)| \leq L \|w - w'\|$. $B_{\ell_2}(x, r)$ denotes the unit ℓ_2 -ball in \mathbb{R}^d .

Definition 1 (Differential Privacy (Dwork et al., 2006)). Let $\varepsilon \in (0, 1]$, $\delta \in [0, 1]$. Randomized algorithm $A : X^n \rightarrow W$ is $\rho\varepsilon, \delta q$ -differentially private (DP) if for all $X, X' \in X^n$ differing in one sample and all measurable subsets $S \subseteq W$, we have $\mathbb{P}(A(X) \in S) \leq e^{\rho\varepsilon} \mathbb{P}(A(X') \in S) + \delta$.

Definition 1 prevents attackers from learning much more about an individual’s data than if the data had not been used for training. If $\delta = 0$, we write ε -DP and say “pure” DP.

Definition 2 (Zero-Concentrated Differential Privacy (zCDP) (Bun & Steinke, 2016)). A randomized algorithm $A : X^n \rightarrow W$ satisfies ρ -zero-concentrated differential privacy (ρ -zCDP) if for all pairs of adjacent data sets $X, X' \in X^n$ and all $\alpha \in [1, \infty)$, we have $D_{\alpha}(A(X) \| A(X')) \leq \rho \alpha$, where $D_{\alpha}(A(X) \| A(X'))$ is the α -Rényi divergence³ between the distributions of $A(X)$ and $A(X')$.

zCDP is weaker than ε -DP, but stronger than $\rho\varepsilon, \delta q$ -DP: see Proposition 20 in Appendix B.

We now define semi-DP (Beimel et al., 2013; Alon et al.,

³For distributions P and Q with probability density/mass functions p and q , $D_{\alpha}(P \| Q) = \frac{1}{\alpha - 1} \ln \int p(x)^{\alpha} q(x)^{1-\alpha} dx$ (Rényi, 1961).

2019), a relaxation of DP that permits A to violate the privacy of the public data:

Definition 3 (Semi-Differential Privacy (Beimel et al., 2013; Alon et al., 2019)). Let $n = n_{\text{priv}} + n_{\text{pub}}$. Consider an algorithm $A : X^n \rightarrow \mathbb{R}^d$ that takes data $X = (X_{\text{priv}}, X_{\text{pub}}) \in \mathcal{X}^{n_{\text{priv}}} \times \mathcal{X}^{n_{\text{pub}}}$ as input. A is (centrally) (ϵ, δ) -semi-DP if $A_{\text{priv}}, X_{\text{pub}}$ is (ϵ, δ) -DP for all $X_{\text{pub}} \in \mathcal{X}^{n_{\text{pub}}}$.

We define ρ -semi-zCDP analogously. We define semi-LDP in Section 3, which is a similar relaxation of local DP (LDP) (Kasiviswanathan et al., 2011; Duchi et al., 2013). To distinguish Definition 3 from semi-LDP, we sometimes refer to algorithms that satisfy Definition 3 as *centrally* semi-DP.

1.3. Roadmap

We study the central model of semi-DP in Section 2. We give tight error bounds for mean estimation in Section 2.1 and Appendix E.1.2, and a novel algorithm with better constants than the asymptotically optimal algorithms in Section 2.2. In Sections 2.3 and 2.4, we characterize the optimal excess risk of semi-DP ERM and SCO respectively. We give an improved semi-DP algorithm for SCO in Section 2.5. In Section 3, we turn to the local model of semi-DP. We characterize the optimal error rates for semi-LDP mean estimation in Section 3.1 and SCO in Section 3.4. We give semi-LDP algorithms with improved error in Sections 3.2 and 3.5. We experimentally evaluate our algorithms in Section 4 and Appendix G. In Appendix A, we discuss related works in more detail. Due to the page limit, some results and all proofs are presented in the Appendix.

2. Optimal Centrally Private Model Training with Public Data

2.1. Optimal Semi-DP Mean Estimation

In this section, we determine the minimax optimal semi-DP error rates for estimating the mean of a bounded distribution.

Consider the following problem: given n_{priv} private samples $X_{\text{priv}} \in \mathbb{B}$ and n_{pub} public samples $X_{\text{pub}} \in \mathbb{B}$, drawn i.i.d. from an unknown distribution P on \mathbb{B} , estimate the population mean $\mathbb{E}_X [x]$ subject to the constraint that A satisfies semi-DP. Defining $n = n_{\text{priv}} + n_{\text{pub}}$, we will characterize the minimax squared error of population mean estimation under (ϵ, δ) -semi-DP:

$$\mathcal{M}_{\text{pop}}(\epsilon, \delta, n_{\text{priv}}, n, d) := \inf_{A \in \mathcal{A}_{\text{semi-DP}}} \sup_{P \in \mathcal{P}(\mathbb{B})} \mathbb{E}_{A; X} \{ \mathbb{E}_X [\|Ax - \mathbb{E}_X [x]\|^2] \}, \quad (1)$$

where $\mathcal{A}_{\text{semi-DP}}$ denotes the set of all (ϵ, δ) -semi-DP estimators $A : \mathbb{B}^n \rightarrow \mathbb{R}^d$, and $|\mathcal{X}_{\text{priv}}| = n_{\text{priv}}$.

Theorem 4. Let $\epsilon \in (0, 1]$, $\delta \in (0, 1/n_{\text{priv}}]$. Then, there

is a constant $C \geq 0$ such that

$$\mathcal{M}_{\text{pop}}(\epsilon, \delta, n_{\text{priv}}, n, d) \leq C \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n^2 \epsilon^2}, \frac{1}{n} \right\} \cdot \frac{1}{n} \cdot \frac{1}{\epsilon} \cdot \frac{1}{\delta} \cdot \frac{1}{n_{\text{priv}}} \cdot \frac{1}{n}, \quad (2)$$

where $\frac{1}{\epsilon} \cdot \frac{1}{\delta}$ is logarithmic in d and n .

Appendix E.1.2 has the proof of Theorem 4 and the pure ϵ -semi-DP result.

Remark 5. Technically, our lower bound proof requires us to assume that $A = (A^1, \dots, A^d)$ is symmetric, meaning $A^j = A^l$ for all $j, l \in [d]$. This is a very reasonable assumption: to our knowledge, every algorithm that has been proposed in the literature (for ℓ_2 geometry) is symmetric. Further, the concurrent work of Ullah et al. (2024) gives an alternative proof that eliminates this assumption.⁴

Naïve algorithms attain the optimal rates: the throw-away estimator $A_{\text{priv}} X_{\text{priv}} / n_{\text{priv}} + X_{\text{pub}} / n_{\text{pub}}$ has MSE $O(1/n_{\text{pub}})$, and the DP (hence semi-DP) Gaussian mechanism has MSE $O(d \ln(1/\delta) \ln(1/\epsilon))$. In the next subsection, we show that these two algorithms have suboptimal constants: we provide improved (smaller error) estimators.

2.2. An “Even More Optimal” Semi-DP Algorithm for Mean Estimation

Before presenting our improved semi-DP algorithms, we precisely describe the worst-case error of the optimal naïve algorithms discussed in the preceding subsection. We will consider ρ -semi-zCDP, which facilitates a sharp characterization of the privacy of the Gaussian mechanism. Note that the lower bound in (2) also holds for semi-zCDP, since $\epsilon^2 \ln(1/\delta)$ -zCDP implies (ϵ, δ) -DP, by Proposition 20.

Definition 6. Let $\mathcal{P}(\mathbb{B}, V)$ be the collection of all distributions P on \mathbb{R}^d such that for any $x \in \mathbb{B}$, we have $\text{Var}_P(x) \leq V$ and $\mathbb{E}_P[x] \in \mathbb{B}$, P -almost surely.

Lemma 7. The error of the ρ -semi-zCDP throw-away algorithm $A_{\text{priv}} X_{\text{priv}} / n_{\text{priv}} + X_{\text{pub}} / n_{\text{pub}}$ is

$$\sup_{P \in \mathcal{P}(\mathbb{B}, V)} \mathbb{E}_X \{ \mathbb{E}_X [\|Ax - \mathbb{E}_X [x]\|^2] \} \leq \frac{V^2}{n_{\text{pub}}}.$$

Further, let X be the average of the public and private samples. The minimax error of the ρ -zCDP Gaussian mech-

⁴The work of Ullah et al. (2024) appeared on arXiv March 6, 2024. The first version of our paper appeared on arXiv on June 26, 2023, while v2 added the d -dimensional (ϵ, δ) -semi-DP lower bounds (Theorem 4 and Theorem 12) and appeared on February 14, 2024.

anism $G_{\rho, X} : X \rightarrow \mathbb{R}^d$ is

$$\inf_{\text{-zCDP}} \sup_{G : \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{G; X} \{ \mathbb{E}_{\rho} \{ \text{MSE} \} \} = \frac{2dB^2}{\rho n^2} \frac{V^2}{n}.$$

Intuitively, it seems like the naïve estimators in Lemma 7 do not harness the public and private data in the most effective way possible, despite being optimal up to constants: Throw-away fails to utilize the private data at all, while the Gaussian mechanism gives equal weight to X_{priv} and X_{pub} (regardless of ρ, d, n_{priv}), and provides unnecessary privacy for X_{pub} . We now present a ρ -semi-zCDP estimator that is “even more optimal” than the naïve estimators, meaning our estimator has smaller worst-case error (accounting for constants). We define the family of *Weighted-Gaussian* estimators:

$$A_r : X \rightarrow \mathbb{R}^d, \quad A_r(x) = \frac{1}{n_{\text{pub}}} \sum_{x \in X_{\text{pub}}} x + r \sum_{x \in X_{\text{priv}}} x, \quad (3)$$

for $r \in [0, 1]$. This estimator can recover both the throw-away and standard Gaussian mechanisms by choosing $r = 0$ or $r = 1$. Intuitively, as ρ shrinks, the accuracy cost of adding privacy noise grows, so we should choose smaller r to reduce the sensitivity of A_r . On the other hand, smaller r increases the variance of A_r on X_{pub} . By choosing r optimally (depending on $\rho, d, n_{\text{priv}}, B, V$), A_r achieves smaller MSE than both throw-away and the Gaussian mechanism:⁵

Proposition 8. A_r is ρ -semi-zCDP, and $\exists r^* \in [0, 1]$ such that

$$\sup_{G : \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{G; X} \{ \text{MSE} \} = \min \left\{ \frac{V^2}{n_{\text{pub}}}, \frac{2dB^2}{\rho n^2} \frac{V^2}{n} \right\}. \quad (4)$$

Further, if $\frac{V^2}{n_{\text{pub}}} \leq \frac{2dB^2}{\rho n^2}$, then the advantage of A_r is

$$\sup_{G : \mathbb{R}^d \rightarrow \mathbb{R}^d} \mathbb{E}_{G; X} \{ \text{MSE} \} \leq \frac{q}{q - s^2} \min \left\{ \frac{V^2}{n_{\text{pub}}}, \frac{2dB^2}{\rho n^2} \frac{V^2}{n} \right\}, \quad (5)$$

where $q = 2 \frac{n_{\text{priv}} V^2}{dB^2}$ and $s = \frac{V}{B} \frac{n_{\text{priv}}}{n_{\text{pub}}}$.

When $\frac{V^2}{n_{\text{pub}}} \leq \frac{2dB^2}{\rho n^2}$, the throw-away estimator outperforms the DP Gaussian mechanism and our Weighted Gaussian

⁵We find the optimal choice of r explicitly in the proof of Proposition 8 in Appendix E.1.3.

estimator outperforms both of these estimators by a factor of at least $q/(q - s^2)$. Also, $q/(q - s^2) \geq 1$ for allowable s, q . For example, if $n = 10,000, d = 100, B = 25V, \rho = 0.1$, and $n_{\text{pub}} = 0.008n$, then the MSE of our Weighted Gaussian A_r is smaller than the MSE of throw-away and standard Gaussian by a multiplicative factor of ≈ 1.98 .

Figures 12-14 in Appendix G.1.2 show that our estimator outperforms both naïve baselines for d -dimensional Bernoulli data with $\rho = 0.5$ (regardless of whether or not throw-away outperforms the Gaussian mechanism).

For pure ε -semi-DP, using Laplace noise instead of Gaussian noise in (3) yields an estimator with smaller error than the ε -DP Laplace mechanism and throw-away.

2.3. Optimal Semi-DP Empirical Risk Minimization

For a given (fixed) $X = (X_{\text{priv}}, X_{\text{pub}}) \in \mathbb{R}^d \times \mathbb{R}^d$ and parameter domain W , consider the ERM problem:

$$\min_{w \in W} \mathbb{P}_X \rho(w) = \frac{1}{n} \sum_{j=1}^n f(w, x_j),$$

where $f : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a loss function and $n_{\text{priv}} = |X_{\text{priv}}|$ samples are private. We discuss practical applications of semi-DP ERM beyond ML in Appendix E.2. We measure the (in-sample) performance of a training algorithm $A : X^n \rightarrow W$ on the data set X by its *excess empirical risk*

$$\mathbb{E}_A \mathbb{P}_X \rho(A(X)) - \mathbb{P}_X \rho^* = \mathbb{E}_A \mathbb{P}_X \rho(A(X)) - \min_{w \in W} \mathbb{P}_X \rho(w).$$

Definition 9. Let $F_{L, D}$ be the set of all functions $f : W \times \mathbb{R}^d \rightarrow \mathbb{R}$ that are uniformly L -Lipschitz and μ -strongly convex ($\mu \neq 0$) in w for some convex compact $W \subseteq \mathbb{R}^d$ with ℓ_2 -diameter bounded by $D \geq 0$ and some set X .

Let $\mathcal{A}^{\varepsilon}$ contain all ε -semi-DP algorithms $A : X^n \rightarrow W$ for some X, W . Define the minimax excess empirical risk of ε -semi-DP (strongly) convex ERM as

$$R_{\text{ERM}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \quad (6)$$

$$:= \inf_{A \in \mathcal{A}^{\varepsilon}} \sup_{F \in F_{L, D}} \sup_{(X, W)} \mathbb{E}_A \mathbb{P}_X \rho(A(X)) - \mathbb{P}_X \rho^*.$$

Theorem 10. There are absolute constants $0 < c \leq C$ s.t.

$$cLD \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\} \leq R_{\text{ERM}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \leq C \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\}. \quad (7)$$

See Appendix E.2 for the $\mu \geq 0$ result and proofs.

Remark 11. The same minimax risk bound (7) holds up to a logarithmic factor if we replace $F_{0, L, D}$ by set of all Lipschitz non-convex loss functions in the definition (6). However, the optimal semi-DP algorithms are inefficient for non-convex loss functions. See Appendix E.2 for details.

2.4. Optimal Semi-DP Stochastic Convex Optimization

In stochastic convex optimization (SCO), we are given n i.i.d. samples from an unknown distribution $X \sim P^n$ (with n_{priv} of them being private), and aim to approximately minimize the expected population loss $F(w) = \mathbb{E}_X \int \langle w, x \rangle f(x, w)$. We measure the quality of a learner A by its excess population risk

$$R_{A;X} = \mathbb{E}_X \int \langle w, x \rangle f(x, w) - \min_{w \in \mathcal{W}} \mathbb{E}_X \int \langle w, x \rangle f(x, w)$$

Denote the minimax optimal semi-DP excess risk by

$$R_{\text{SCO}}(\epsilon, \delta, n_{\text{priv}}, n, d, L, D, \mu) = \inf_{A \in \mathcal{A}(\epsilon, \delta)} \sup_{P \in \mathcal{P}_F} \sup_{L, D} \mathbb{E}_X \int \langle w, x \rangle f(x, w) - \min_{w \in \mathcal{W}} \mathbb{E}_X \int \langle w, x \rangle f(x, w), \quad (8)$$

where $\mathcal{A}(\epsilon, \delta)$ contains all ϵ, δ -semi-DP algorithms $A : X^n \rightarrow \mathcal{W}$ for some X, \mathcal{W} , and $|X_{\text{priv}}| = n_{\text{priv}}$.

Theorem 12. *Let $\epsilon \in \{1\} \cup \{2^{-k} \mid k \in \mathbb{N}\}$ and $\delta \in \{1\} \cup \{2^{-k} \mid k \in \mathbb{N}\}$. Then, there is a constant $C \geq 0$ such that*

$$R_{\text{SCO}}(\epsilon, \delta, n_{\text{priv}}, n, d, L, D, \mu) \leq C \left(\frac{1}{n_{\text{pub}}} + \frac{d}{n\epsilon} + \frac{1}{n} \right) + \frac{1}{n} \left(\frac{d \ln(1/\delta)}{n\epsilon} + \frac{1}{n} \right),$$

where $\ln(1/\delta)$ is logarithmic in d and n .

We provide the $\delta = 0$ and μ -strongly convex results ($\mu \geq 0$), and proofs in Appendix E.3. Remark 5 also applies to Theorem 12

Let us compare the semi-DP bound for SCO in Theorem 12 with the ERM bound in Theorem 10 when $d = 1 \dots L \dots D$. Depending on the values of ϵ and n_{priv} , the minimax excess population risk (“test loss”) of SCO may either be larger or smaller than the excess empirical risk (“training loss”) of ERM. For example, if $\epsilon = 1$, then the semi-DP excess empirical risk $\frac{1}{n} \mathbb{E} \int \langle w, x \rangle f(x, w)$ is smaller than the excess population risk $\frac{1}{n} \mathbb{E} \int \langle w, x \rangle f(x, w)$. On the other hand, suppose $\epsilon = 2^{-k}$ and $n_{\text{priv}} = n^{2k}$: then the semi-DP excess empirical risk $\frac{1}{n} \mathbb{E} \int \langle w, x \rangle f(x, w)$ is larger than the excess population risk $\frac{1}{n} \mathbb{E} \int \langle w, x \rangle f(x, w)$. This is surprising: for both non-private learning and DP learning (with $n_{\text{pub}} = 0$), the optimal error of ERM is never larger than that of SCO. While it may seem counter-intuitive that minimizing the training loss can be harder than minimizing test loss, there is a natural explanation: For SCO, a small amount of public data gives us free information about the private data, since $X \sim P^n$ is i.i.d. by assumption. By contrast, for ERM, the public data does not give us any information about the private data, since X is not i.i.d.

2.5. Semi-DP SCO with an “Even More Optimal” Gradient Estimator

Our Algorithm 1 is a noisy stochastic gradient method that uses the “even more optimal” Weighted-Gaussian estimator (33) to estimate $\int \langle w, x \rangle f(x, w)$ in iteration t .⁶ In Algorithm 1, $\text{clip}_C \int \langle w, x \rangle f(x, w) = \arg \min_{y \in \mathbb{R}^d, \|y\|_2 \leq C} \int \langle w, x \rangle f(x, w) - y$ is the Euclidean projection onto the centered ℓ_2 -ball of radius C .

We give privacy and excess risk guarantees for Algorithm 1 and describe an accelerated variant of Algorithm 1 in Appendix E.3.1. The excess risk of our algorithm is smaller than the state-of-the-art excess risk for a linear-time DP algorithm whose privacy analysis does not require convexity (Lowy & Razaviyayn, 2023b). We empirically evaluate our algorithm in Section 4.

Algorithm 1 Semi-DP-SGD via Weighted-Gaussian Gradient Estimation

- 1: **Input:** $T \in \mathbb{N}$, clip threshold $C \geq 0$, stepsizes $\eta_t \in \mathbb{R}^+$, batch sizes $K_{\text{priv}} \in \mathbb{N}$, $K_{\text{pub}} \in \mathbb{N}$, weight parameter $\alpha \in \mathbb{R}^+$, noise parameter $\sigma^2 \geq 0$.
 - 2: Initialize $w_0 \in \mathcal{W}$.
 - 3: **for** $t \in \{0, 1, \dots, T-1\}$ **do**
 - 4: Draw random batch of K_{priv} private samples B_t^{priv} .
 - 5: Draw random batch of K_{pub} public samples B_t^{pub} .
 - 6: Draw privacy noise $v_t \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$.
 - 7: $g_t = \alpha \left(\frac{1}{K_{\text{priv}}} \sum_{x \in B_t^{\text{priv}}} \text{clip}_C \int \langle w_t, x \rangle f(x, w_t) + \frac{1}{K_{\text{pub}}} \sum_{x \in B_t^{\text{pub}}} \int \langle w_t, x \rangle f(x, w_t) \right) + v_t$.
 - 8: Update $w_{t+1} = w_t - \eta_t g_t$.
 - 9: **end for**
 - 10: **Output:** w_T or an average of the iterates $\{w_t\}_{t \in \mathbb{T}_S}$.
-

3. Optimal Locally Private Model Training with Public Data

We now turn to a stronger privacy notion that we refer to as *semi-local DP* (semi-LDP). Semi-LDP guarantees *privacy for each private x_i , without requiring person i to trust others* (e.g. central server). Semi-LDP generalizes LDP (Kasiviswanathan et al., 2011; Duchi et al., 2013), which has been deployed in industry (Apple, 2016; Úlfar Erlingsson et al., 2014; Ding et al., 2017).

Following Duchi & Rogers (2019), we permit algorithms to be *fully interactive*: algorithms may adaptively query the same person i multiple times over the course of T communication rounds. For example, n cell phone users send messages to a server over T rounds, and message $Z_{i,t} \in \mathcal{Z}$ sent by user i in round t can depend on the previous messages $\{Z_{j,t'}\}_{j \in [n], t' \leq t}$. Semi-LDP requires the messages

⁶In Algorithm 1, we re-parameterize by setting $\eta = \frac{\alpha}{K_{\text{priv}}}$ for $\alpha \in \mathbb{R}^+$.

$\{Z_i; t_{Pr} T_S\}$ to be semi-DP:

Definition 13 (Semi-Local Differential Privacy). *The T -round interactive algorithm A is (ϵ, δ) -semi-LDP if the transcript $Z = \{Z_i; t_{Pr} T_S\}$ is (ϵ, δ) -semi-DP: i.e. for all $X_{pub} \in \mathcal{B}^n$, all adjacent $X_{priv} = X_{priv}^1$ and all $S \in \mathcal{Z}^{nT}$, $\mathbb{P}[Z \in S | X_{priv}, X_{pub}] \leq e^{\epsilon} \mathbb{P}[Z \in S | X_{priv}^1, X_{pub}] + \delta$.*

Definition 13 is stronger than Definition 3, since the latter only requires the final output of A to be semi-DP.

3.1. Optimal Semi-LDP Mean Estimation

We will characterize the minimax squared error of ϵ -semi-LDP d -dimensional mean estimation:

$$\mathcal{M}_{pop}^{loc}(\epsilon, n_{priv}, n, d) := \inf_{A \in \mathcal{A}^{loc}} \sup_{P \in \mathcal{P}} \mathbb{E}_{A; X \sim P} \{ \text{Ap}X \}_{X \sim P} \{ \text{r}xs \}^2, \quad (9)$$

where \mathcal{A}^{loc} contains all (fully interactive) ϵ -semi-LDP estimators $A: \mathcal{B}^n \rightarrow \mathcal{B}$, and $|X_{priv}| = n_{priv}$.

Theorem 14. *Let $\epsilon \in (0, 1]$. There are absolute constants $0 < c \leq C$ s.t.*

$$c \min \left\{ \frac{1}{n_{pub}}, \frac{d}{n\epsilon^2} \right\} \leq \mathcal{M}_{pop}^{loc}(\epsilon, n_{priv}, n, d) \leq C \min \left\{ \frac{1}{n_{pub}}, \frac{d}{n\epsilon^2} \right\}. \quad (10)$$

Remark 15 (Approximate Semi-LDP). *Theorem 14 still holds if we replace \mathcal{A}^{loc} in the definition of $\mathcal{M}_{pop}^{loc}(\epsilon, n_{priv}, n, d)$ by the set of all (ϵ, δ) -semi-LDP estimators A for which either $\delta \leq \frac{1}{2}$ and A is “compositional” (Duchi & Rogers, 2019) (e.g. sequentially interactive (Duchi et al., 2013)) or $\delta \leq \frac{1}{2}d$.*

The upper bound in Theorem 14 is the minimum of the error of the throw-away estimator and the optimal ϵ -LDP estimator of Duchi et al. (2013). The LDP estimator of Duchi et al. (2013) takes the form

$$\mathbb{E}_{\mathcal{M}_{Duchi}(\mathbb{P}X)} = \frac{1}{n} \sum_{i=1}^n \mathbb{M}_{Duchi}(\mathbb{P}x_i),$$

where $\mathbb{M}_{Duchi}(\mathbb{P}x_i)$ samples a vector uniformly from a carefully chosen subset of \mathcal{B} , depending on x_i .

3.2. An “Even More Optimal” Semi-LDP Estimator

By applying \mathcal{M}_{Duchi} only to the private samples, we obtain an ϵ -semi-LDP algorithm with smaller error than the asymptotically optimal \mathcal{M}_{Duchi} . Define the *Semi-LDP* $A_{Semi-Duchi}$:

$$A_{Semi-Duchi}(\mathbb{P}X) = \frac{1}{n} \sum_{x \in X_{priv}} \mathbb{M}_{Duchi}(\mathbb{P}x) + \sum_{x \in X_{pub}} x^1. \quad (11)$$

Let P be a distribution on \mathcal{B} with $V^2 = \mathbb{E}_{X \sim P} \{ \text{r}xs \}^2$.

Lemma 16. *Let $c > 0$ such that $\mathbb{E}_{X \sim P} \{ \mathcal{M}_{Duchi}(\mathbb{P}x) \}_{X \sim P} \{ \text{r}xs \}^2 \leq \frac{cd}{n^2}$, so that $\mathbb{E}_{X \sim P} \{ \mathcal{M}_{Duchi}(\mathbb{P}X) \}_{X \sim P} \{ \text{r}xs \}^2 \leq \frac{cd}{n^2} \frac{V^2}{n}$. Then,*

$$\mathbb{E}_{X \sim P} \{ A_{Semi-Duchi}(\mathbb{P}X) \}_{X \sim P} \{ \text{r}xs \}^2 \leq \frac{n_{priv}}{n} \frac{cd}{n\epsilon^2} + \frac{n_{pub}}{n} \frac{V^2}{n}.$$

The constant c in Lemma 16 that bounds the error of \mathcal{M}_{Duchi} may depend on the distribution P . Lemma 16 shows that given any P , the error of our semi-LDP estimator is smaller than the error of \mathcal{M}_{Duchi} . Quantitatively, the MSE of our estimator is smaller than the MSE of \mathcal{M}_{Duchi} by a factor of n_{priv}/n if the privacy noise error term is dominant (e.g. if $d \ll \epsilon^2$).

3.3. A Semi-LDP Estimator with Optimal Constants

In this subsection, we consider the task of estimating the average of data X on the unit sphere $S^{d-1} \in \mathbb{R}^d$. We give a semi-LDP estimator that is *truly optimal*—i.e. *our estimator has the smallest MSE, including constants*—among a large class of unbiased semi-LDP estimators of X . Our semi-LDP estimator, $A_{Semi-PrivU}$ takes a similar shape to $\mathcal{A}_{Semi-Duchi}$, but uses *PrivUnit* (Bhowmick et al., 2018) instead of \mathcal{M}_{Duchi} as the LDP randomizer in (11). We recall *PrivUnit* in Algorithm 3 in Appendix F.

Proposition 17. *Let $R: S^{d-1} \rightarrow \mathcal{Z}$ be an ϵ -LDP randomizer, M_{priv} and M_{pub} be aggregation protocols, and $\mathbb{A}(\mathbb{P}X) = \frac{1}{n} \{ M_{priv}(\mathbb{R}(\mathbb{P}x_1), \dots, \mathbb{R}(\mathbb{P}x_{n_{priv}})) + M_{pub}(\mathbb{P}X_{pub}) \}$. Assume $\mathbb{E} \{ M_{priv}(\mathbb{R}(\mathbb{P}x_1), \dots, \mathbb{R}(\mathbb{P}x_{n_{priv}})) | X_{priv} \} = \sum_{x \in X_{priv}} x$ and $\mathbb{E} \{ M_{pub}(\mathbb{P}X_{pub}) | X_{pub} \} = \sum_{x \in X_{pub}} x$. Then,*

$$\sup_{X \in \mathcal{P}S^{d-1}} \mathbb{E}_{A_{Semi-PrivU}} \{ A_{Semi-PrivU}(\mathbb{P}X) \}_{X \sim P} \{ \text{r}xs \}^2 \leq \sup_{X \in \mathcal{P}S^{d-1}} \mathbb{E}_A \{ \mathbb{A}(\mathbb{P}X) \}_{X \sim P} \{ \text{r}xs \}^2.$$

Proposition 17 is proved by extending the analysis of Asi et al. (2022) to the semi-LDP setting.

3.4. Optimal Semi-LDP Stochastic Convex Optimization

We will characterize the minimax optimal excess population risk of semi-LDP SCO

$$\mathcal{R}_{SCO}^{loc}(\epsilon, n_{priv}, n, d, L, D, \mu) := \inf_{A \in \mathcal{A}^{loc}} \sup_{F \in \mathcal{F}} \sup_{P \in \mathcal{P}} \mathbb{E}_{A; X \sim P} \{ F(\mathbb{A}(\mathbb{P}X)) - F \}, \quad (12)$$

where \mathcal{A}^{loc} denotes the set of all algorithms $A : X^n \rightarrow W$ that are ε -semi-LDP for some X and W , and exactly n_{priv} samples in X are private.

Theorem 18. Let $\varepsilon \in (0, 1]$ and $h(\varepsilon, n_{\text{priv}}, n, d, L, D) \leq \frac{1}{n_{\text{pub}}}$. There are absolute constants c and C with $0 < c \leq C$, such that

$$c h(\varepsilon, n_{\text{priv}}, n, d, L, D) \leq F_{\text{SCO}}^{\text{loc}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \leq C h(\varepsilon, n_{\text{priv}}, n, d, L, D).$$

See Appendix F.2 for the $\mu \leq 0$ result and proofs. The first term in the upper bound is achieved by throwing away X_{priv} and running SGD on X_{pub} (Nemirovski & Yudin, 1983). The second term in the upper bound is achieved by the one-pass LDP-SGD of Duchi et al. (2013). Remark 15 also applies to Theorem 18.

3.5. “Even More Optimal” Semi-LDP SCO Algorithm

We give a semi-LDP algorithm, called *Semi-LDP-SGD*, with smaller excess risk than the optimal LDP-SGD of Duchi et al. (2013). Essentially, Semi-LDP-SGD runs as follows: In each iteration $t \in [1, n]$, we draw a random sample $x_t \in X$ without replacement. If $x_t \in X_{\text{priv}}$, update $w_{t+1} = \text{Proj}_{\mathcal{W}}(w_t - \eta \nabla f(w_t, x_t))$; if $x_t \in X_{\text{pub}}$, instead update $w_{t+1} = \text{Proj}_{\mathcal{W}}(w_t - \eta \nabla f(w_t, x_t))$. See Algorithm 4 in Appendix F.2.1 for pseudocode.

Proposition 19. Let $f \in \mathcal{F}_{0:L;D}$, let P be any distribution and $\varepsilon \leq d$. Algorithm 4 is ε -semi-LDP. Further, there is an absolute constant c such that the output $\text{Ap}X$ of Algorithm 4 satisfies

$$\mathbb{E}_{A;X} \text{Pr}[\text{Ap}X] \leq F_{\text{SCO}}^s \leq \frac{c}{\varepsilon^2} \frac{D}{n_{\text{priv}}} + \frac{c}{n_{\text{pub}}}. \quad (13)$$

Thus, Algorithm 4 has smaller excess risk than LDP-SGD, roughly by a factor of $\frac{1}{n_{\text{priv}}}$.

4. Numerical Experiments

In this section, we empirically evaluate the performance of four different semi-DP algorithms: 1. *Throw-away* (i.e. minimize the public loss). 2. *DP-SGD* (Abadi et al., 2016; De et al., 2022). 3. *PDA-MD* (Amid et al., 2022), which is the state-of-the-art semi-DP algorithm for training convex models. 4. *Our Algorithm 1*. Unless otherwise noted, we evaluate all algorithms with “warm start,” which means finding a minimizer w_{pub} of the public loss and initializing training at w_{pub} . The hyperparameters of each algorithm were carefully tuned. See Appendix G for details on the experimental setups and additional results.

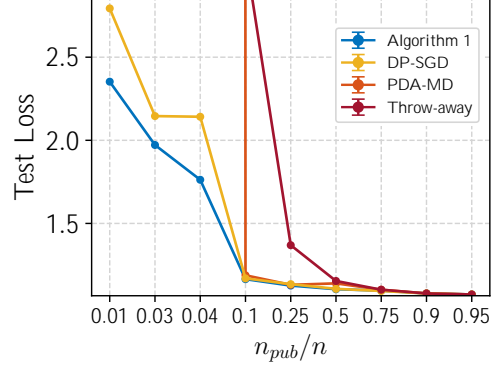


Figure 3. Test loss vs. n_{pub}/n ($\varepsilon = 2$).

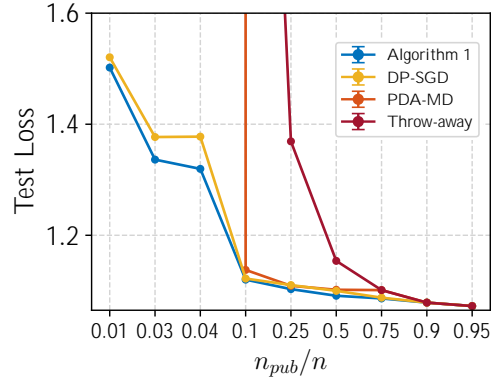


Figure 4. Test loss vs. n_{pub}/n ($\varepsilon = 4$).

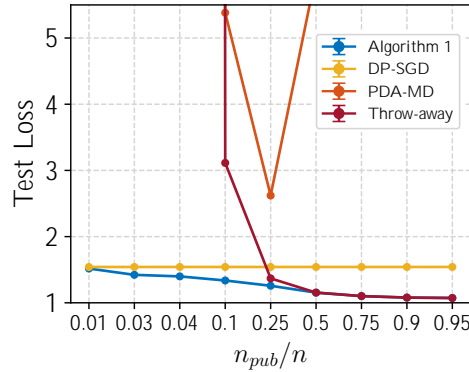


Figure 5. Test loss vs. n_{pub}/n ($\varepsilon = 4$, without warm-start).

Our Algorithm 1 achieves the smallest test loss among the semi-DP baselines across different levels of ε (privacy) and n_{pub} : Figures 3-4 show results for $\varepsilon, \delta = 10^{-5}$ q-semi-DP linear regression with synthetic Gaussian data. In the Appendix, we evaluate the algorithms in several other tasks: e.g., logistic regression and Wide-ResNet: see Figures 15-18 and 19-20. Our results indicate that *Algorithm 1 consistently outperforms all baselines*.

Our Algorithm 1 can converge even when DP-SGD diverges: Figure 6 gives an example in which DP-SGD diverges but Algorithm 1 converges. We used the following

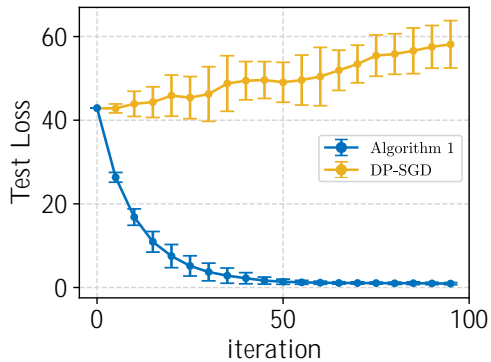


Figure 6. Test loss vs. iterations. $\frac{n_{\text{pub}}}{n} = 0.1$

parameters: $d = 50$, $n = 1000$, $\epsilon = 0.01$, and $\frac{n_{\text{pub}}}{n} = 0.1$.

Algorithm 1 performs well even without “warm start”, whereas PDA-MD performs poorly: e.g., Figures 5 and (in Appendix) 11 show that PDA-MD even performs worse than throw-away. By contrast, Algorithm 1 is resilient to the “cold start” condition and outperforms all baselines. In certain practical applications, such as advertising and health-care, samples are often obtained in an online/streaming fashion, precluding the possibility of “warm start”.

More public data always improves the performance of Algorithm 1, but does not always benefit PDA-MD. The main reason for this is that our algorithm more effectively handles the increasing privacy noise that is needed to maintain semi-DP with increasing n_{pub}/n . We give details and numerical evidence of this explanation in Appendix G.1.1.

Appendix G contains other findings, too. For example, Figure 8 shows that PDA-MD is sensitive to Hessian regularization parameter, which requires extra tuning on complicated tasks. By contrast, Algorithm 1 does not use any Hessian information.

5. Conclusion

We considered training DP models with side access to public data. Theoretically, we characterized the optimal error bounds (up to constants) for three fundamental problems: mean estimation, empirical risk minimization, and stochastic convex optimization. We show that it is impossible to improve over the naïve semi-DP algorithms asymptotically, in the worst case. Algorithmically, we developed new optimal methods for semi-DP learning that have smaller error than the asymptotically optimal algorithms. Empirically, we showed that our algorithms are effective in training semi-DP models. Our work raises interesting open questions. For instance, why do certain learning/optimization problems benefit more from public data than others? Is there some general underlying property that (don’t) permit asymptotic benefits over the naïve baselines? Also, what can be said

about semi-DP learning with *out-of-distribution* public data? Lastly, it would be useful to have an extensive empirical study that evaluates the efficacy of combining our semi-DP algorithms with various other techniques, such as dimensionality reduction (Yu et al., 2021b; Pinto et al., 2024).

Impact Statement

Our work provides algorithms for protecting the privacy of individuals who contribute training data. Privacy is commonly regarded in a positive light and is even enshrined as a fundamental right in various legal systems. However, there is a risk that corporations or governments could exploit our algorithms for nefarious purposes, such as unauthorized collection of personal data. Furthermore, the use of semi-privately trained models may result in decreased accuracy compared to non-private models, which can have adverse consequences. For instance, if a semi-DP model is utilized to forecast the effects of pollution, but yields less precise and overly optimistic outcomes, it could provide pretext for a government to unjustly dismantle environmental safeguards. Nonetheless, we firmly believe that the dissemination of privacy-preserving machine learning algorithms, coupled with enhanced understanding of these algorithms, ultimately offers a net benefit to society.

Another potential misuse of our work would be using the accuracy benefits of public data to argue for less stringent data privacy policies, laws, or regulations. However, we want to emphasize that even though public data can enhance the accuracy of models, we firmly believe that privacy laws and corporate policies should not be weakened. Differentially private synthetic data generation (Torkzadehmahani et al., 2019; Vietri et al., 2020; Boedihardjo et al., 2022; He et al., 2023) is one possible avenue for generating public data in an ethical, privacy-preserving manner.

Acknowledgements

The authors would like to thank Gautam Kamath, Adam Smith, Thomas Steinke, and Jonathan Ullman for very helpful pointers and explanations related to existing lower bound proof techniques. We thank Arnold Pereira for discussions and providing feedback at various stages of this project. We thank Michael Menart for helpful feedback on an earlier version of this manuscript. We thank the TPDP and ICML reviewers for their helpful feedback. This work was supported in part by a gift from Meta and the USC-Meta Center for Research and Education in AI & Learning. AL’s work was supported in part by NSF award DMS-2023239 and AFOSR award FA9550-21-1-0084.

References

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, Oct 2016. doi: 10.1145/2976749.2978318. URL <http://dx.doi.org/10.1145/2976749.2978318>.
- Acharya, J., Sun, Z., and Zhang, H. Differentially private assouad, fano, and le cam. In *Algorithmic Learning Theory*, pp. 48–78. PMLR, 2021.
- Alon, N., Bassily, R., and Moran, S. Limits of private learning with access to public data. *Advances in Neural Information Processing Systems*, 32, 2019.
- Amid, E., Ganesh, A., Mathews, R., Ramaswamy, S., Song, S., Steinke, T., Suriyakumar, V. M., Thakkar, O., and Thakurta, A. Public data-assisted mirror descent for private model training. In *International Conference on Machine Learning*, pp. 517–535. PMLR, 2022.
- Apple. Differential privacy overview, 2016.
- Asi, H. and Duchi, J. C. Near instance-optimality in differential privacy. *arXiv preprint arXiv:2005.10630*, 2020a.
- Asi, H. and Duchi, J. C. Instance-optimality in differential privacy via approximate inverse sensitivity mechanisms. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 14106–14117. Curran Associates, Inc., 2020b. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/a267f936e54d7c10a2bb70dbe6ad7a89-Paper.pdf.
- Asi, H., Feldman, V., and Talwar, K. Optimal algorithms for mean estimation under local differential privacy. In *International Conference on Machine Learning*, pp. 1046–1056. PMLR, 2022.
- Avent, B., Korolova, A., Zeber, D., Hovden, T., and Livshits, B. Blender: Enabling local search with a hybrid differential privacy model. In *USENIX Security Symposium*, pp. 747–764, 2017.
- Barber, R. F. and Duchi, J. C. Privacy and statistical risk: Formalisms and minimax bounds. *arXiv preprint arXiv:1412.4451*, 2014.
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 464–473. IEEE, 2014.
- Bassily, R., Thakkar, O., and Guha Thakurta, A. Model-agnostic private learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Bassily, R., Feldman, V., Talwar, K., and Thakurta, A. Private stochastic convex optimization with optimal rates. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Bassily, R., Cheu, A., Moran, S., Nikolov, A., Ullman, J., and Wu, S. Private query release assisted by public data. In *International Conference on Machine Learning*, pp. 695–703. PMLR, 2020.
- Beimel, A., Nissim, K., and Stemmer, U. Private learning and sanitization: Pure vs. approximate differential privacy. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques: 16th International Workshop, APPROX 2013, and 17th International Workshop, RANDOM 2013, Berkeley, CA, USA, August 21-23, 2013. Proceedings*, pp. 363–378. Springer, 2013.
- Ben-David, S., Bie, A., Canonne, C. L., Kamath, G., and Singhal, V. Private distribution learning with public data: The view from sample compression. *arXiv preprint arXiv:2308.06239*, 2023.
- Bhowmick, A., Duchi, J., Freudiger, J., Kapoor, G., and Rogers, R. Protection against reconstruction and its applications in private federated learning. *arXiv preprint arXiv:1812.00984*, 2018.
- Bie, A., Kamath, G., and Singhal, V. Private estimation with public data. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 18653–18666. Curran Associates, Inc., 2022.
- Boedihardjo, M., Strohmer, T., and Vershynin, R. Covariance’s loss is privacy’s gain: Computationally efficient, private and accurate synthetic data. *Foundations of Computational Mathematics*, pp. 1–48, 2022.
- Bubeck, S. et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Bun, M. and Steinke, T. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Proceedings, Part I, of the 14th International Conference on Theory of Cryptography - Volume 9985*, pp. 635–658, Berlin, Heidelberg, 2016. Springer-Verlag. ISBN 9783662536407. doi: 10.1007/978-3-662-53641-4_24. URL https://doi.org/10.1007/978-3-662-53641-4_24.

- Bun, M., Ullman, J., and Vadhan, S. Fingerprinting codes and the price of approximate differential privacy. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pp. 1–10, 2014.
- Bun, M., Steinke, T., and Ullman, J. Make up your mind: The price of online queries in differential privacy. In *Proceedings of the twenty-eighth annual ACM-SIAM symposium on discrete algorithms*, pp. 1306–1325. SIAM, 2017.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T. B., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *USENIX Security Symposium*, volume 6, 2021.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- Church, G. M. The personal genome project, 2005.
- De, S., Berrada, L., Hayes, J., Smith, S. L., and Balle, B. Unlocking high-accuracy differentially private image classification through scale. *arXiv preprint arXiv:2204.13650*, 2022.
- Ding, B., Kulkarni, J., and Yekhanin, S. Collecting telemetry data privately. *Advances in Neural Information Processing Systems*, 30, 2017.
- Duchi, J. Lecture notes for statistics 311/electrical engineering 377. URL: https://stanford.edu/class/stats311/Lectures/full_notes.pdf, 2021.
- Duchi, J. and Rogers, R. Lower bounds for locally private estimation via communication complexity. In *Conference on Learning Theory*, pp. 1161–1191. PMLR, 2019.
- Duchi, J. C., Jordan, M. I., and Wainwright, M. J. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pp. 429–438, 2013. doi: 10.1109/FOCS.2013.53.
- Duchi, J. C., Jordan, M. I., and Wainwright, M. J. Minimax optimal procedures for locally private estimation. *Journal of the American Statistical Association*, 113(521):182–201, 2018.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pp. 265–284. Springer, 2006.
- Dwork, C., Roth, A., et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Dwork, C., Smith, A., Steinke, T., Ullman, J., and Vadhan, S. Robust traceability from trace amounts. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pp. 650–669. IEEE, 2015.
- Fallah, A., Makhdoumi, A., Malekian, A., and Ozdaglar, A. Optimal and differentially private data acquisition: Central and local mechanisms. In *Proceedings of the 23rd ACM Conference on Economics and Computation*, EC ’22, pp. 1141, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450391504. doi: 10.1145/3490486.3538329. URL <https://doi.org/10.1145/3490486.3538329>.
- Feldman, V., Mironov, I., Talwar, K., and Thakurta, A. Privacy amplification by iteration. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 521–532. IEEE, 2018.
- Ferrando, C., Gillenwater, J., and Kulesza, A. Combining public and private data. *arXiv preprint arXiv:2111.00115*, 2021.
- Fredrikson, M., Jha, S., and Ristenpart, T. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 1322–1333, 2015.
- Ganesh, A., Thakurta, A., and Upadhyay, J. Langevin diffusion: An almost universal algorithm for private euclidean (convex) optimization. *arXiv preprint arXiv:2204.01585*, 2022.
- Ganesh, A., Haghifam, M., Nasr, M., Oh, S., Steinke, T., Thakkar, O., Thakurta, A., and Wang, L. Why is public pretraining necessary for private model training?, 2023.
- Ghadimi, S. and Lan, G. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on Optimization*, 22(4):1469–1492, 2012. doi: 10.1137/110848864. URL <https://doi.org/10.1137/110848864>.
- Golatkar, A., Achille, A., Wang, Y.-X., Roth, A., Kearns, M., and Soatto, S. Mixed differential privacy in computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8376–8386, 2022.
- Gu, X., Kamath, G., and Wu, Z. S. Choosing public datasets for private machine learning via gradient subspace distance, 2023.
- Hardt, M. and Talwar, K. On the geometry of differential privacy. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pp. 705–714, 2010.

- He, Y., Vershynin, R., and Zhu, Y. Algorithmically effective differentially private synthetic data. *arXiv preprint arXiv:2302.05552*, 2023.
- Jorgensen, Z., Yu, T., and Cormode, G. Conservative or liberal? personalized differential privacy. In *2015 IEEE 31st international conference on data engineering*, pp. 1023–1034. IEEE, 2015.
- Kairouz, P., Diaz, M. R., Rush, K., and Thakurta, A. (nearly) dimension independent private erm with adagrad rates via publicly estimated subspaces. In *Conference on Learning Theory*, pp. 2717–2746. PMLR, 2021.
- Kamath, G., Liu, X., and Zhang, H. Improved rates for differentially private stochastic convex optimization with heavy-tailed data. In *International Conference on Machine Learning*, pp. 10633–10660. PMLR, 2022a.
- Kamath, G., Mouzakis, A., and Singhal, V. New lower bounds for private estimation and a generalized fingerprinting lemma. *Advances in Neural Information Processing Systems*, 35:24405–24418, 2022b.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811. Springer, 2016.
- Kasiviswanathan, S. P., Lee, H. K., Nissim, K., Raskhodnikova, S., and Smith, A. What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826, 2011.
- Kerrigan, G., Slack, D., and Tuyls, J. Differentially private language models benefit from public pre-training. *arXiv preprint arXiv:2009.05886*, 2020a.
- Kerrigan, G., Slack, D., and Tuyls, J. Differentially private language models benefit from public pre-training. *arXiv preprint arXiv:2009.05886*, 2020b.
- Klimt, B. and Yang, Y. The enron corpus: A new dataset for email classification research. In *Machine Learning: ECML 2004: 15th European Conference on Machine Learning, Pisa, Italy, September 20-24, 2004. Proceedings 15*, pp. 217–226. Springer, 2004.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Li, X., Tramer, F., Liang, P., and Hashimoto, T. Large language models can be strong differentially private learners. *arXiv preprint arXiv:2110.05679*, 2021a.
- Li, X., Tramer, F., Liang, P., and Hashimoto, T. Large language models can be strong differentially private learners. In *International Conference on Learning Representations*, 2021b.
- Liu, J. and Talwar, K. Private selection from private candidates. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 298–309, 2019.
- Liu, R., Bu, Z., Wang, Y.-x., Zha, S., and Karypis, G. Coupling public and private gradient provably helps optimization. *arXiv preprint arXiv:2310.01304*, 2023.
- Liu, T., Vietri, G., Steinke, T., Ullman, J., and Wu, S. Leveraging public data for practical private query release. In *International Conference on Machine Learning*, pp. 6968–6977. PMLR, 2021.
- Lowy, A. and Razaviyayn, M. Output perturbation for differentially private convex optimization with improved population loss bounds, runtimes and applications to private adversarial training. *arXiv preprint:2102.04704*, 2021.
- Lowy, A. and Razaviyayn, M. Private stochastic optimization with large worst-case lipschitz parameter: Optimal rates for (non-smooth) convex losses and extension to non-convex losses. In *International Conference on Algorithmic Learning Theory*, pp. 986–1054. PMLR, 2023a.
- Lowy, A. and Razaviyayn, M. Private federated learning without a trusted server: Optimal algorithms for convex losses. In *International Conference on Learning Representations*, 2023b. URL <https://openreview.net/pdf?id=TVY6GoURrw>.
- Lowy, A., Ghafelebashi, A., and Razaviyayn, M. Private non-convex federated learning without a trusted server. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, 2023a.
- Lowy, A., Gupta, D., and Razaviyayn, M. Stochastic differentially private and fair learning. In *International Conference on Learning Representations*, 2023b. URL <https://openreview.net/pdf?id=3nM5uhPIfv6>.
- McSherry, F. D. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pp. 19–30, 2009.
- Mehta, H., Thakurta, A., Kurakin, A., and Cutkosky, A. Large scale transfer learning for differentially private image classification. *arXiv preprint arXiv:2205.02973*, 2022.
- Mühl, C. and Boenisch, F. Personalized pate: Differential privacy for machine learning with individual privacy guarantees. *arXiv preprint arXiv:2202.10517*, 2022.
- Nasr, M., Mahloujifar, S., Tang, X., Mittal, P., and Houmansadr, A. Effectively using public data in privacy preserving machine learning. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.),

- Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 25718–25732. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/nasr23a.html>.
- Nemirovski, A. and Yudin, D. *Problem complexity and method efficiency in optimization*. Chichester, 1983.
- Papernot, N. and Steinke, T. Hyperparameter tuning with renyi differential privacy. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=-70L8lpp9DF>.
- Papernot, N., Abadi, M., Erlingsson, Ú., Goodfellow, I., and Talwar, K. Semi-supervised knowledge transfer for deep learning from private training data. In *International Conference on Learning Representations*, 2017.
- Papernot, N., Song, S., Mironov, I., Raghunathan, A., Talwar, K., and Erlingsson, U. Scalable private learning with pate. In *International Conference on Learning Representations*, 2018.
- Pinto, F., Hu, Y., Yang, F., and Sanyal, A. Pillar: How to make semi-private learning more effective. In *2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 110–139, 2024. doi: 10.1109/SaTML59370.2024.00014.
- Polyak, B. T. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- Rakhlin, A., Shamir, O., and Sridharan, K. Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 1571–1578, 2012.
- Rényi, A. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, volume 4, pp. 547–562. University of California Press, 1961.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18. IEEE, 2017.
- Tang, X., Panda, A., Sehwal, V., and Mittal, P. Differentially private image classification by learning priors from random processes, 2023.
- Thakurta, A. G., Vyrros, A. H., Vaishampayan, U. S., Kapoor, G., Freudiger, J., Sridhar, V. R., and Davidson, D. Learning new words, March 14 2017. US Patent 9,594,741.
- Torkzadehmahani, R., Kairouz, P., and Paten, B. Dp-cgan: Differentially private synthetic data and label generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019.
- Úlfar Erlingsson, Pihur, V., and Korolova, A. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 21st ACM Conference on Computer and Communications Security*, 2014.
- Ullah, E., Menart, M., Bassily, R., Guzmán, C., and Arora, R. Public-data assisted private stochastic optimization: Power and limitations. *arXiv preprint arXiv:2403.03856*, 2024.
- U.S. Census Bureau. Differential privacy and the 2020 census, 2020.
- Vietri, G., Tian, G., Bun, M., Steinke, T., and Wu, S. New oracle-efficient algorithms for private synthetic data release. In *International Conference on Machine Learning*, pp. 9765–9774. PMLR, 2020.
- Wang, J. and Zhou, Z.-H. Differentially private learning with small public data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 6219–6226, 2020.
- Yu, B. Assouad, fano, and le cam. *Festschrift for Lucien Le Cam: research papers in probability and statistics*, pp. 423–435, 1997.
- Yu, D., Naik, S., Backurs, A., Gopi, S., Inan, H. A., Kamath, G., Kulkarni, J., Lee, Y. T., Manoel, A., Wutschitz, L., et al. Differentially private fine-tuning of language models. In *International Conference on Learning Representations*, 2021a.
- Yu, D., Zhang, H., Chen, W., and Liu, T.-Y. Do not let privacy overbill utility: Gradient embedding perturbation for private learning. *arXiv preprint arXiv:2102.12677*, 2021b.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2017.
- Zhou, Y., Wu, Z. S., and Banerjee, A. Bypassing the ambient dimension: Private sgd with gradient subspace identification. *arXiv preprint arXiv:2007.03813*, 2020.

Appendix

A. Related Work

Many works have considered a variety of semi-DP learning problems, empirically and theoretically (Bassily et al., 2018; Feldman et al., 2018; Bassily et al., 2020; Kairouz et al., 2021; Wang & Zhou, 2020; Zhou et al., 2020; Li et al., 2021a; Liu et al., 2021; Papernot et al., 2017; 2018; Yu et al., 2021b; Alon et al., 2019; Bie et al., 2022; Ferrando et al., 2021; Amid et al., 2022). Here we discuss the works that are most closely related to our own.

Sample complexity bounds for semi-DP learning and estimation: The works of Alon et al. (2019); Bassily et al. (2020) give sample complexity bounds for semi-DP PAC learning and query release. Their upper bounds show that for hypothesis classes with finite VC-dimension, asymptotic improvements over the naïve approaches are possible. These results stand in contrast with our lower bounds for model training/optimization with Lipschitz loss functions. Both of these works also provide lower bounds. Below, we compare their results with our own lower bounds.

The negative result of Theorem 4.2 in (Alon et al., 2019) is similar in spirit to our lower bounds: no semi-DP algorithm can achieve error better than the minimum of the optimal DP error and a term that depends on n_{pub} . That being said, there are some significant and consequential differences between our lower bounds and (Alon et al., 2019, Theorem 4.2):

1. *Different learning problems:* (Alon et al., 2019) considers agnostic PAC learning/binary classification, whereas we consider model training (mean estimation and optimization). We are not aware of any way to obtain our lower bounds from the results of (Alon et al., 2019).
2. *Quantitative differences in the lower bounds:* The lower bound implied by (Alon et al., 2019, Theorem 4.2) is of the form $\min\{n_{\text{pub}}, \text{optimal DP error}\}$. Our lower bounds do not always take this form. For example, consider Theorem 10: the first term in our lower bound is $n_{\text{priv}}\{n$, which can imply a much larger error than the bound in (Alon et al., 2019) (e.g., if $d = \epsilon n$ and $n_{\text{priv}} = n/2$). This illustrates how different learning problems can benefit more from public data than others.
3. *Central vs. Local Semi-DP:* (Alon et al., 2019) only covers central semi-DP, whereas we cover both central and local semi-DP.
4. *Pure vs. Approximate Semi-DP:* (Alon et al., 2019)’s proof technique cannot handle approximate semi-DP because their Lemma 2.6 is limited to pure DP. By contrast, we give lower bounds for both pure and approximate semi-DP.
5. *New Techniques:* Our techniques differ substantially from those in (Alon et al., 2019). For example, we develop a novel semi-DP Fano’s inequality and a novel semi-DP packing argument. For approximate semi-DP, we utilize fingerprinting proofs. Also, our semi-LDP lower bound techniques are completely different from (Alon et al., 2019)’s techniques. We hope that our novel techniques to find applications beyond those in our paper.

The result of Bassily et al. (2020, Theorem 13) showed (up to logarithmic factors) that no improvement over the naïve approaches is possible for *approximate* ϵ_1, δ -semi-DP releasing decision stumps. This implies a lower bound for ϵ_1, δ -semi-DP mean estimation in the ℓ_∞ norm, but does not imply the tight lower bound for the ℓ_2 setting that we provide in Theorem 4. Moreover, (Bassily et al., 2018)’s results and techniques do not lead to tight lower bounds for *pure* ϵ -semi-DP decision stumps or mean estimation. Thus, we develop novel techniques (e.g. semi-DP Fano and semi-DP packing arguments) for pure semi-DP estimation and model training.

For semi-DP d -dimensional Gaussian mean estimation, (Bie et al., 2022) gave sample complexity bounds that do not depend on the range parameters of the distribution if $n_{\text{pub}} \asymp d^{-1}$; this is known to be impossible without public data. The concurrent and independent work of Ben-David et al. (2023) established a lower bound for this problem. (Ben-David et al., 2023) also explored the connection between semi-DP distribution learning of a class and the existence of a sample compression scheme for that class.

DP model training (ERM and SCO) with public data: The works of Kairouz et al. (2021); Zhou et al. (2020) considered DP ERM with public data and additional assumptions on the gradients lying in a certain low-dimensional subspace. Under these additional assumptions, (Kairouz et al., 2021; Zhou et al., 2020) show that nearly dimension-independent excess empirical risk bounds are possible (e.g. by using the public data to estimate the low-dimensional subspace and projecting

noisy gradients onto this subspace). Our lower bounds show that these additional assumptions are strictly necessary: in general, polynomial dependence on the dimension is necessary for semi-DP ERM and SCO. The work of Wang & Zhou (2020) used public data to adjust the parameters of DP-SGD. Empirically, pre-training on public data sets and privately fine-tuning the model (Li et al., 2021a; Kerrigan et al., 2020a; Mehta et al., 2022) has shown great promise for large-scale ML.

The work of Amid et al. (2022) developed a public data-assisted DP mirror descent (PDA-MD) algorithm that sometimes outperforms DP-SGD empirically in training ML models, and theoretically in terms of excess risk for linear regression under certain distributional assumptions. We use the PDA-MD of Amid et al. (2022) as a baseline in our experiments. (Amid et al., 2022) also gave an “efficient approximation” of their PDA-MD in (Amid et al., 2022, Equation 1), which they used for training non-convex models. While finalizing this manuscript, we became aware that this “efficient approximation” is nearly equivalent to our Algorithm 1. However, there are differences in the implementation of our algorithm. For example, we use a constant weight parameter α , whereas (Amid et al., 2022) uses decaying weights $\alpha_t u_t^T$. Also, we clip both the private and public gradients in our implementation of Algorithm 1, which empirically improves performance: see Appendix G.1.1 for further discussion. Our Algorithm 1 was derived in a different fashion from the algorithm in (Amid et al., 2022, Equation 1): we derived our algorithm as an application of our “even more optimal” mean estimator, while theirs was derived as an approximation of their mirror descent method. Moreover, no theoretical analysis was provided for the “efficient approximation” in (Amid et al., 2022). In the linear regression setting, the PDA-MD algorithm can be viewed as Newton’s algorithm where the Hessian is estimated via public data. Then the Hessian is inverted and multiplied by privatized gradients at each iteration. In other words, the update rule of the algorithm is given by $w^{t+1} = w^t - \alpha_t (X_{pub}^T X_{pub} + \eta I)^{-1} p g_t + n_t$ where X_{pub} is the public data, g_t is the sampled gradient from private data, and n_t is the added noise. When the number of public data samples is small, the estimate of the Hessian becomes low rank (or inaccurate) and the (pseudo)inverse of it may introduce additional error (even after proper regularization). On the other hand, when the number of public data samples is large, but the number of private data is small, the PDA-MD algorithm can still suffer if it is not warm-started. This is because, although the Hessian $X^T X$ is estimated accurately in this case, but there is not enough private data to generate enough gradients to converge to optimal solution. This poor performance of PDA-MD when it is not warm started is also observed in our experiments.

The work of Nasr et al. (2023) used public data to train a generative model for data augmentation and to estimate the center of the clipping balls in DP-SGD. They did not provide any code for their experiments and thus we do not compare against them as a baseline in our experiments. However, combining our algorithm with the tricks used in (Nasr et al., 2023) could be a promising avenue for future empirical work.

Personalized DP: A related line of work is that of *personalized DP* (PDP) (Jorgensen et al., 2015; Golatkar et al., 2022; Mühl & Boenisch, 2022; Fallah et al., 2022), a generalization of DP in which each person may have different privacy parameters ϵ_i, δ_i . By letting $\epsilon_i = \epsilon$ for some person i , PDP also generalizes semi-DP. We leverage this connection to borrow techniques from the work of Fallah et al. (2022), which considers pure ($\delta_i = 0$) PDP estimation in one dimension ($d = 1$). We also note that the 1-dimensional pure (central) PDP mean estimation bound of Fallah et al. (2022) extends easily to a 1-dimensional ϵ -semi-DP bound. However, our d -dimensional semi-DP lower bounds require a different set of techniques. Additionally, the personalized LDP bound of Fallah et al. (2022) relies on the assumption $\epsilon_i \asymp 1$ and does not seem to extend to the ϵ -semi-LDP setting.

Concurrent and Subsequent Work: The work of Ullah et al. (2024) first appeared on arXiv a few weeks after v2 of our paper appeared.⁷ (Ullah et al., 2024) proves results that are very similar to Theorems 4 and 12. However, their lower bound proofs do not require the (mild) symmetry assumption that our proof requires: see Remark 5. Moreover, in a restricted parameter regime, Ullah et al. (2024) also provide lower bounds that are tighter by a $\log_2 \frac{1}{\delta}$ factor. Ullah et al. (2024) complement these lower bounds by giving novel algorithms for leveraging unlabeled public data in training private generalized linear models (GLM) with dimension-independent rates.

The work of Liu et al. (2023), which appeared on arXiv 3 months after v1 of our paper, couples private and public gradients in non-convex optimization via a similar weighting scheme to our own. They show benefits of their algorithm over standard DP approaches both theoretically and empirically.

⁷The first version of our paper appeared on arXiv more than eight months before (Ullah et al., 2024). However, v1 of our paper did not contain our tight high-dimensional approximate semi-DP population mean estimation and SCO lower bounds.

Tang et al. (2023), which appeared on arXiv 18 days before v1 of our paper, explores how to improve the privacy-utility tradeoff of DP-SGD by learning priors from images generated by random processes and transferring these priors to private data.

Other Related Works: The work of Gu et al. (2023) gave an algorithm for selecting an appropriate public dataset that can be used to enhance private optimization by projecting gradients onto a subspace prescribed by the this public dataset.

Ganesh et al. (2023) provided an explanation for the empirically reported benefits of pre-training on public data, arguing that non-convex optimization algorithms must go through two phases: (i) selecting a good “basin” in the loss landscape; (ii) solving an easy optimization problem within that basin. They hypothesize that public pre-training can be helpful in selecting a good basin. They also demonstrated a separation between pretrained and non-pretrained models by constructing a non-convex optimization problem for which public pretraining is necessary to achieve non-trivial error.

B. Zero-Concentrated DP vs. Pure and Approximate DP

Proposition 20. (Bun & Steinke, 2016) *If A is ρ -zCDP, then A is $\rho \log \frac{1}{1-\rho}$ -DP. Moreover, if A is ε -DP, then A is ε^2 -zCDP.*

C. Summary table of pure semi-DP and semi-LDP results

Learning problem	Semi-DP error	When is semi-DP error less than DP?	Learning problem	Semi-LDP error	When is semi-LDP error less than LDP?
Mean Estimation (Pop. MSE)	$\min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2}{n^2 \varepsilon^2} + \frac{1}{n} \right\}$ (Theorem 32)	$n_{\text{pub}} > \frac{n^2 \varepsilon^2}{d^2}$ or $n_{\text{pub}} = \Theta(n)$	Mean Estimation (Pop. MSE)	$\min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n \varepsilon^2} \right\}$ (Theorem 14)	$n_{\text{pub}} > \frac{\varepsilon^2 n}{d}$ or $n_{\text{pub}} = \Theta(n)$
ERM (Excess emp. Risk)	$\min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n \varepsilon} \right\}$ (Theorem 10)	$n_{\text{pub}} > n - d/\varepsilon$	SCO (Excess pop. risk)	$\min \left\{ \frac{1}{\sqrt{n_{\text{pub}}}}, \sqrt{\frac{d}{n \varepsilon^2}} \right\}$ (Theorem 18)	$n_{\text{pub}} > \frac{\varepsilon^2 n}{d}$ or $n_{\text{pub}} = \Theta(n)$
SCO (Excess pop. risk)	$\min \left\{ \frac{1}{\sqrt{n_{\text{pub}}}}, \frac{d}{n \varepsilon} + \frac{1}{\sqrt{n}} \right\}$ (Theorem 42)	$n_{\text{pub}} > \frac{n^2 \varepsilon^2}{d^2}$ or $n_{\text{pub}} = \Theta(n)$			

Figure 7. Minimax optimal error rates for central ε -semi-DP and (local) ε -semi-LDP SCO and mean estimation results. Dependence on range and Lipschitz parameters, and constraint set diameter omitted. Mean estimation and SCO lower bounds are only tight if $n_{\text{pub}} \geq \frac{1}{\rho} \log \frac{1}{1-\rho}$. Strongly convex ERM and SCO results are included later in this Appendix, but excluded from this table.

D. Notation

We recall notation from the main body and include some additional basic definitions below for convenience.

Let $\|\cdot\|$ be the ℓ_2 norm. W denotes a convex, compact subset of \mathbb{R}^d with ℓ_2 diameter D . X denotes a data universe. Function $g : W \rightarrow \mathbb{R}$ is μ -strongly convex if $g(\alpha w + (1-\alpha)w^1) \leq \alpha g(w) + (1-\alpha)g(w^1) - \frac{\mu}{2} \|w - w^1\|^2$ for all $\alpha \in [0, 1]$ and all $w, w^1 \in W$. If $\mu = 0$, we say g is convex. For convex $f : W \rightarrow \mathbb{R}$, denote any subgradient of f at w, x w.r.t. w by $\gamma \in \partial_w f(w, x) \in \mathbb{B}_w(f(w, x))$: i.e. $f(w^1, x) \geq f(w, x) + \langle \gamma, w^1 - w \rangle$ for all $w^1 \in W$. Function $f : W \times X \rightarrow \mathbb{R}$ is uniformly L -Lipschitz in w if $\sup_{x \in X} |f(w, x) - f(w^1, x)| \leq L \|w - w^1\|$. Let $\mathbb{B}(x, r) \subset \mathbb{R}^d$ denote the unit ℓ_2 -ball. For functions $a : \Theta \rightarrow \mathbb{R}$ and $b : \Phi \rightarrow \mathbb{R}$ of input parameter vectors θ and ϕ , we write $a \lesssim b$ or $a = O(b)$ if there is an absolute constant $C > 0$ such that $a \leq Cb$ for all values of input parameter vectors θ and ϕ .

E. Optimal Centrally Private Model Training with Public Data

E.1. Optimal Semi-DP Mean Estimation

We begin in Appendix E.1.1 with empirical mean estimation. This subsection was omitted from the main body due to space constraints. Theorem 21 will be useful for proving our ERM bounds (Theorem 10). Then, in Appendix E.1.2, we turn to population mean estimation (i.e. the proof of Theorem 4). Appendix E.1.3 contains proofs for Section 2.2.

E.1.1. ESTIMATING THE EMPIRICAL MEAN

For a given data set $X \in \mathcal{B} : \{x\} \times \mathbb{R}^d : \{x\} \times \mathbb{R}^d$, consider the problem of estimating $X = \frac{1}{n} \sum_{i=1}^n x_i$ subject to the constraint that the estimator satisfies semi-DP. We will characterize the minimax squared error of d -dimensional empirical mean estimation under ε -semi-DP:

$$\mathcal{M}_{\text{emp}}(\varepsilon, n_{\text{priv}}, n, d) : \inf_{A \in \mathcal{A}(\varepsilon, n_{\text{priv}}, n, d)} \sup_{X \in \mathcal{B}} \mathbb{E}_A \left[\sum_{i=1}^d (A_i(X) - X_i)^2 \right], \quad (14)$$

where $\mathcal{A}(\varepsilon, n_{\text{priv}}, n, d)$ denotes the set of all ε -semi-DP estimators $A : \mathcal{B} \rightarrow \mathbb{R}^d$.

Theorem 21. *Let $\varepsilon \in (0, 1]$, $n, d \in \mathbb{N}$, $n_{\text{priv}} \in \mathbb{N}$. There exist absolute constants c and C with $0 < c \leq C$ such that*

$$c \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\} \leq \mathcal{M}_{\text{emp}}(\varepsilon, n_{\text{priv}}, n, d) \leq C \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\}.$$

Proof. Lower bound: We use a packing argument to prove our lower bound. Denote $X = \frac{1}{n} \sum_{i=1}^n x_i$. Let $A : \mathcal{B} \rightarrow \mathbb{R}^d$ be ε -semi-DP. Choose $K = 2^{d/2}$ private data points $\{x_i\}_{i=1}^K \in \mathbb{R}^d$ such that $\|x_i - x_j\| \geq \frac{1}{8}$ for all $i \neq j$. The existence of such a set of points is well-known (see e.g. the Gilbert-Varshamov construction). Let $n = \frac{nd}{3n_{\text{priv}}}$.

Case 1: $n \leq n$ (i.e. $d \leq 3\varepsilon n_{\text{priv}}$). In this case, we'll show that $\mathbb{E} \sum_{i=1}^d (A_i(X) - X_i)^2 \geq \frac{n_{\text{priv}}}{n} \varepsilon^2$ for some $X \in \mathcal{B}$. For $i \in [K]$, let $X_i = \frac{1}{n} \sum_{j=1}^n x_j$, $\mathbf{0}_{n-n_{\text{priv}}}$ consist of n_{priv} copies of x_i followed by $n - n_{\text{priv}}$ copies of $\mathbf{0} \in \mathbb{R}^d$. Suppose for the sake of contradiction that for every $i \in [K]$, with probability $\geq \frac{1}{3}$ we have $\sum_{i=1}^d (A_i(X) - X_i)^2 \leq \frac{1}{32} \frac{n_{\text{priv}}}{n}$. That is, we are supposing $\mathbb{P}(\sum_{i=1}^d (A_i(X) - X_i)^2 \leq \frac{1}{32} \frac{n_{\text{priv}}}{n}) \geq \frac{1}{3}$ for all i , where $B_i = \{x \in \mathcal{B} : \sum_{i=1}^d (x_i - X_i)^2 \leq \frac{1}{32} \frac{n_{\text{priv}}}{n}\}$. Note that the sets $\{B_i\}_{i=1}^K$ are disjoint by construction. Since A is ε -semi-DP, group privacy implies that $\mathbb{P}(\sum_{i=1}^d (A_i(X) - X_i)^2 \leq \frac{1}{32} \frac{n_{\text{priv}}}{n}) \leq e^{-\frac{1}{3} \varepsilon n_{\text{priv}}}$ for all $i \in [K]$. Thus,

$$K \frac{1}{3} e^{-\frac{1}{3} \varepsilon n_{\text{priv}}} \leq \sum_{i=1}^K \mathbb{P}(\sum_{i=1}^d (A_i(X) - X_i)^2 \leq \frac{1}{32} \frac{n_{\text{priv}}}{n}) \leq 1,$$

where the last inequality follows from disjointness of the balls B_i . Thus, we obtain $\ln(K \frac{1}{3}) \leq \varepsilon n_{\text{priv}}$. Assume for now that $d \leq 8$. (A 1-dimensional lower bound that is tight up to constant factors can be shown easily by following the proof of the d -dimensional case but choosing $K = 16$ instead of $K = 2^{d/2}$.) Then $d/2 \leq d - 4 \leq \varepsilon n_{\text{priv}}$ implies $d \leq 2\varepsilon n_{\text{priv}}$, contradicting the assumption made in Case 1. Thus, we conclude that there exists a data set X_i such that with probability at least $\frac{2}{3}$, $\sum_{i=1}^d (A_i(X) - X_i)^2 \geq \frac{1}{32} \frac{n_{\text{priv}}}{n}$. Squaring both sides of this inequality and applying Markov's inequality yields the desired lower bound.

Case 2: $n > n$.

Additionally, suppose for now that $n \leq n_{\text{priv}}$. In this case, we'll show that $\mathbb{E} \sum_{i=1}^d (A_i(X) - X_i)^2 \geq \frac{d}{n} \varepsilon^2$ for some $X \in \mathcal{B}$. Let $X_i = \frac{1}{n} \sum_{j=1}^n x_j$, $\mathbf{0}_{n-n}$ consist of n copies of x_i followed by $n - n$ copies of $\mathbf{0} \in \mathbb{R}^d$.⁸ Denoting the mean of a dataset X by $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$ for convenience, we see that $\bar{X}_i = \frac{1}{n} \sum_{j=1}^n x_j = \frac{1}{n} x_i$. Define the algorithm $\hat{A} : \mathcal{B} \rightarrow \mathbb{R}^d$ by $\hat{A}_i(X) = \frac{1}{n} \sum_{j=1}^n A_j(X)$. Since A is ε -semi-DP, we see that \hat{A} is ε -semi-DP by post-processing. Also, the domain of \hat{A} is \mathcal{B} and $n \leq n$, so the argument in Case 1 applies to \hat{A} . Thus, by applying the result in Case 1, there exists $i \in [K]$ such that with probability at least $\frac{2}{3}$, $\sum_{i=1}^d (\hat{A}_i(X) - \bar{X}_i)^2 \geq \frac{1}{32} \frac{n_{\text{priv}}}{n}$. (Here $\bar{X}_i = \frac{1}{n} x_i$.) But this implies $\sum_{i=1}^d (A_i(X) - X_i)^2 \geq \frac{1}{32} \frac{n_{\text{priv}}}{n} \frac{n}{n} = \frac{1}{96} \frac{d}{n}$ with probability at least $\frac{2}{3}$. Again, squaring both sides and applying Markov yields the desired lower bound.

Next, consider the complementary subcase where $n > n_{\text{priv}}$. Define the algorithm $\hat{A} : \mathcal{B} \rightarrow \mathbb{R}^d$ by $\hat{A}_i(X) = \frac{1}{n} \sum_{j=1}^n A_j(X)$. Since A is ε -semi-DP, we see that \hat{A} is ε -semi-DP by post-processing. Also, the domain of \hat{A} is \mathcal{B} and $n \leq n$, so the argument in Case 1 applies to \hat{A} . Thus, by applying the result in Case 1, there exists $i \in [K]$ such that with probability at least $\frac{2}{3}$, $\sum_{i=1}^d (\hat{A}_i(X) - \bar{X}_i)^2 \geq \frac{1}{32} \frac{n_{\text{priv}}}{n}$. (Here $\bar{X}_i = \frac{1}{n} x_i$.) But this implies that $\sum_{i=1}^d (A_i(X) - X_i)^2 \geq \frac{1}{32} \frac{n_{\text{priv}}}{n}$.

⁸We assume without loss of generality that $n \in \mathbb{N}$. If n is not an integer, then choosing $\lceil n \rceil$ instead yields the same bound up to constant factors.

with probability at least $2/3$. Again, squaring both sides and applying Markov yields the desired lower bound. Thus, the lower bound holds in all cases.

Upper bound: For the first term in the minimum, consider the algorithm which throws away the private data and outputs $\frac{1}{n} \sum_{x \in X_{\text{pub}}} x$. Clearly, A is semi-DP. Moreover,

$$\mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X_{\text{pub}}} x - \frac{1}{n} \sum_{x \in X_{\text{priv}}} x \right)^2 \right\} \leq \frac{n_{\text{priv}}}{n^2}.$$

For the second term, consider the Laplace mechanism $\frac{1}{n} \sum_{x \in X} x + \frac{1}{n} \sum_{i=1}^d L_i$, where $L_i \sim \text{Lap}(2/\epsilon)$ are i.i.d. mean-zero Laplace random variables. We know A is ϵ -DP by (Dwork et al., 2014), since the ℓ_1 -sensitivity is $\frac{1}{n}$. Hence A is ϵ -semi-DP. Moreover, A has error

$$\mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x - \frac{1}{n} \sum_{i=1}^d L_i \right)^2 \right\} \leq \frac{8d^2}{n^2 \epsilon^2}.$$

Combining the two upper bounds completes the proof. \square

E.1.2. ESTIMATING THE POPULATION MEAN

Approximate (ϵ, δ) -Semi-DP Mean Estimation

Theorem 22 (Formal statement of Theorem 4). *Let $\epsilon \in (0, 1]$, $\delta \in (0, 1]$. Then, there is an absolute constant $C > 0$ such that*

$$\mathcal{M}_{\text{pop}}(\epsilon, \delta, n_{\text{priv}}, n, d) \geq C \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d \ln(1/\delta)}{n^2 \epsilon^2}, \frac{1}{n} \right\}, \quad (15)$$

where $\mathcal{M}_{\text{pop}}(\epsilon, \delta, n_{\text{priv}}, n, d)$ is logarithmic in d and n . The lower bound holds for symmetric algorithms A^1, \dots, A^d such that $A^j = A^l$ for all $j, l \in [d]$.

For the lower bound, we will first prove a stronger result in Theorem 23, in which we construct a hard distribution whose mean is small—scaling with the accuracy lower bound that we aim to prove. This “small mean” property will be needed for our semi-DP SCO lower bound (Theorem 12), even though it is not necessary for the proof of Theorem 22.

Theorem 23. *Let $X = \{x^1, \dots, x^d\}$, $\epsilon \in (0, 1]$, and $\delta \in (0, 1]$. Then, for any symmetric (ϵ, δ) -semi-DP A , there exists a product distribution P on X with $\mathbb{E}_{x \sim P} \{x\} = \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n} \right\}$ such that*

$$\mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x - \mathbb{E}_{x \sim P} \{x\} \right)^2 \right\} \geq \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{\epsilon^2 n_{\text{priv}}}, \frac{1}{n} \right\}.$$

For any symmetric A and any distribution P with $X \sim P^n$, we have

$$\begin{aligned} \mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x - \mathbb{E}_{x \sim P} \{x\} \right)^2 \right\} &= \mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x - \mathbb{E}_{x \sim P} \{x\} \right)^2 \right\} \\ &= \mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x - \mathbb{E}_{x \sim P} \{x\} \right)^2 \right\} \\ &= d \mathbb{E} \left\{ \left(\frac{1}{n} \sum_{x \in X} x^j - \mathbb{E}_{x \sim P} \{x^j\} \right)^2 \right\}, \end{aligned} \quad (16)$$

where the last equality used assumption that A is symmetric. For a $P \in \mathcal{P}(\mathbb{R})$, scalar random variable $x^j \sim P_a$ with mean $\mathbb{E} \{x^j\} = a$ and $X^j \sim P_a^n$ denote

$$\text{Bias}_a(P_a) := \left| \mathbb{E} \{x^j\} - a \right|.$$

Definition 24 (Low bias algorithms). *We say symmetric A is low bias if for every $a \in \mathbb{R}$,*

$$\text{Bias}_a(P_a)^2 \leq \frac{1}{d} \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{1}{n}, \frac{d}{n^2 \epsilon^2} \right\}.$$

We can assume without loss of generality that A is low bias when we are proving the lower bound in Theorem 22: if A is not low bias, then inequality (16) implies that the worst-case MSE of A is lower bounded as in (15).

To prove Theorem 23 for low bias and symmetric A , we will use Theorem 25.⁹ This result shows that any sufficiently accurate A is vulnerable to an attack that traces many individuals in the data set with high probability.

Theorem 25. *Let $a \in [0, 1]$. Consider the product distribution on $\mathbb{t} \subseteq \mathbb{R}^d$ defined in the following way: for $j \in [d]$, independently draw $\theta^j \sim \text{Unif}[a, a + \delta]$ and $x_j^i \sim P$ such that $x_j^i \in \mathbb{t}$ with mean $\mathbb{E}_{x_j^i \sim P} x_j^i = \theta^j$ for $i \in [n]$. Denote $X = (x_1, \dots, x_n)$ where $\theta = (\theta^1, \dots, \theta^d)$ and $P = \prod_{j=1}^d P_j$. Let $A : \mathbb{t} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^d$ satisfy $\mathbb{E}_{X \sim P^n} \mathbb{E}_{\mathcal{A} \sim \mathcal{A}}$ and $\mathbb{E}_{\mathcal{A} \sim \mathcal{A}} \|\mathcal{A}(X) - \theta\|^2 \leq \alpha^2$ for $\frac{d}{n} \leq \alpha \leq \frac{a + \delta}{100}$. Assume $d \leq 400\alpha n / \ln 2 \ln \frac{1}{\delta}$. Moreover, assume that $\mathcal{A}(X) = (A^1(X), \dots, A^d(X))$ with $A^l = A^l$ for all $j, l \in [d]$. Then, the attack $l : \mathbb{t} \subseteq \mathbb{R}^d \rightarrow \{a, a + \delta\} \subseteq \mathbb{R}$, OUT if described in Algorithm 2 satisfies the following properties: a) if $y \sim P$ independently of X , then $\Pr[l(y) = \mathcal{A}(X)] \leq \delta$; and b) $\Pr[\text{IN}] \leq \delta$, where \mathcal{A} is (ϵ, δ) -semi-DP if and only if A is (ϵ, ζ) -semi-DP, and $\mathbb{E}_{X \sim P^n} \mathcal{A}(X) = \theta$.*

The proof of Theorem 25 uses a convenient reduction in Lemma 27: for any low bias and symmetric A that has small expected mean squared error (in ℓ_2 -norm), there exists another low bias mechanism \mathcal{A} that has even smaller ℓ_2 accuracy with high probability. Moreover, A is (ϵ, δ) -semi-DP if and only if \mathcal{A} is (ϵ, δ) -semi-DP. This lemma allows us to: construct a new product distribution whose mean is small, modify the attack of Dwork et al. (2015), and derive a generalized fingerprinting lemma (Lemma 26) in order to prove Theorem 25.

With Theorem 25 in hand, we leverage the proof technique of Bassily et al. (2020) to show that any sufficiently accurate \mathcal{A} must leak the data of more than n_{pub} individuals and thus cannot be (ϵ, δ) -semi-DP. But by Lemma 27, this means that A cannot be (ϵ, δ) -semi-DP. Below, we discuss the attack that we use and then fill in the details of the proof.

The attack in Algorithm 2 is a modification of the robust tracing attack of Dwork et al. (2015). The attacker in Algorithm 2 receives as input the output $q = \mathcal{A}(X)$ of a mechanism, a target point y , and the mean of the data distribution P , $\theta = \mathbb{E}_{x \sim P} x$. (An estimate of the true mean or access to sufficiently many independent draws from P would also suffice in lieu of θ .) The target y is either a data point used by $\mathcal{A}(X)$ (i.e. $y \in X$) or an independent draw from the distribution P that X was drawn from. The attacker aims to infer whether or not y was in X . If the attacker outputs IN when y is in X and OUT when y is not in X (i.e. the attack succeeds) with high probability, then the algorithm A is not private. The truncation parameter is $\eta = 2\alpha \sqrt{\frac{d}{\delta}}$, where α is the expected ℓ_2 -error of the mechanism A . The parameter δ can be chosen by the attacker.

Algorithm 2 Tracing Attack Against ℓ_2 -Accurate Mechanisms

- 1: **Input:** Target $y \in \mathbb{t} \subseteq \mathbb{R}^d$, mean $\theta \in \mathbb{R}^d$, output of mechanism $q = \mathcal{A}(X)$.
 - 2: Let $\eta = 2\alpha \sqrt{\frac{d}{\delta}}$.
 - 3: Let $t_q = \theta \vee \eta$ denote the entrywise truncation of q .
 - 4: Compute $\chi_y = \theta, t_q \vee y = \frac{1}{d} \sum_{j=1}^d y^j \theta^j t_q^j$.
 - 5: **if** $\chi_y = \theta, t_q \vee y \geq \tau : 2\eta \sqrt{\frac{d}{\delta}}$ **then**
 - 6: **Output:** IN.
 - 7: **else**
 - 8: **Output:** OUT.
 - 9: **end if**
-

Compared to (Dwork et al., 2015), we choose a smaller truncation parameter to handle the ℓ_2 case ($\eta = 2\alpha \sqrt{\frac{d}{\delta}}$ instead of 2α). We also scale the “IN/OUT” decision threshold τ accordingly. Moreover, since we are not concerned with the size of the reference sample, we assume that the attacker knows the mean of the data distribution. Thus, we use θ in place of the average of reference samples to give a simpler attack than the one in (Dwork et al., 2015). This is not necessary for our attack to work: the attacker just needs a sufficiently close approximation to the true mean (which can be attained with access to enough reference samples) for our proof to go through.

Theorem 25 builds on (Dwork et al., 2015, Theorem 17), which gave a similar result for ℓ_2 -accurate mechanisms and θ drawn from a so-called *strong distribution* (see (Dwork et al., 2015, Definition 5)). By contrast, Theorem 25 only requires a

⁹For convenience, we work with data drawn from $\mathbb{t} \subseteq \mathbb{R}^d$ and then re-scale to obtain the final mean estimation lower bound.

weaker ℓ_2 -accuracy guarantee and uses a distribution which is not “strong” in the sense required by (Dwork et al., 2015, Theorem 17). Our generalization of the *fingerprinting lemma* (Bun et al., 2017), given below, allows us to handle such a distribution, and is the first step towards proving Theorem 25:

Lemma 26 (Generalized Fingerprinting Lemma). *Let $a \in \mathbb{R}, \theta \in \mathbb{R}$. Draw $\theta \sim \text{Unif}(a, a+s)$ and $x_i \sim P$ be such that $x_i \in \mathbb{R}$ with mean $E_{x_i \sim P} \{x_i\} = \theta$ for $i = 1, \dots, n$. Denote $X = (x_1, \dots, x_n) \in \mathbb{R}^n$. Then, for any $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that does not depend directly on θ , we have*

$$E_{\theta \sim \text{Unif}(a, a+s)} \left[\sum_{i=1}^n f(x_i, \theta) \right] \leq \frac{a^2}{3} E_{\theta \sim \text{Unif}(a, a+s)} \left[\sum_{i=1}^n f(x_i, \theta)^2 \right] + \frac{1-a^2}{2a} E_{X \sim P^n} \left[\sum_{i=1}^n f(x_i, \theta) \right] - E_{X \sim P^n} \left[\sum_{i=1}^n f(x_i, \theta) \right],$$

where all expectations are over the random draw of θ and $X \sim P^n$ unless otherwise indicated.

Proof. Define $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(\theta) = E_{X \sim P^n} \left[\sum_{i=1}^n f(x_i, \theta) \right]$. Our first claim is

$$E_{X \sim P^n} \left[\sum_{i=1}^n f(x_i, \theta) \right] = g(\theta) \quad (17)$$

To prove (17), write

$$g(\theta) = E_{X \sim P^n} \left[\sum_{i=1}^n f(x_i, \theta) \right]$$

and

$$\begin{aligned} g'(\theta) &= E_{X \sim P^n} \left[\sum_{i=1}^n \frac{d}{d\theta} f(x_i, \theta) \right] \\ &= E_{X \sim P^n} \left[\sum_{i=1}^n \frac{d}{d\theta} \left(\frac{1}{2} \rho\theta\{x_i\} \right) \right] \\ &= E_{X \sim P^n} \left[\sum_{i=1}^n \frac{1}{2} \rho\theta\{x_i\} \right] \\ &= E_{X \sim P^n} \left[\sum_{i=1}^n \frac{1}{2} \frac{x_i}{\theta} \right] \\ &= E_{X \sim P^n} \left[\sum_{i=1}^n \frac{x_i}{\theta} \right] \cdot \frac{1}{2} \end{aligned}$$

Since $E_{X \sim P^n} \left[\sum_{i=1}^n x_i \right] = 0$, (17) is proved.

Next, we claim: for any differentiable function $g : \mathbb{R} \rightarrow \mathbb{R}$, we have

$$E_{\theta \sim \text{Unif}(a, a+s)} \left[g'(\theta) \right] = 2 E_{\theta \sim \text{Unif}(a, a+s)} \left[g(\theta) \right] - \frac{1-a^2}{2a} g(a) - g(a+s). \quad (18)$$

To prove (18), let $u(\theta) = 1 - \theta^2$ and use the product rule and fundamental theorem of calculus to write

$$\begin{aligned} E_{\theta \sim \text{Unif}(a, a+s)} \left[g'(\theta) \right] &= \frac{1}{2a} \int_a^{a+s} g'(\theta) u(\theta) d\theta \\ &= \frac{1}{2a} \int_a^{a+s} \frac{d}{d\theta} (g(\theta) u(\theta)) d\theta \\ &= \frac{1}{2a} (g(a+s) u(a+s) - g(a) u(a)) + \int_a^{a+s} g(\theta) u'(\theta) d\theta \\ &= \frac{1}{2a} (g(a+s) (1 - (a+s)^2) - g(a) (1 - a^2)) + \int_a^{a+s} g(\theta) (-2\theta) d\theta \\ &= \frac{1-a^2}{2a} (g(a+s) - g(a)) - 2 E_{\theta \sim \text{Unif}(a, a+s)} \left[g(\theta) \right]. \end{aligned}$$

Now we will apply (17) and (18) with the differentiable function $g : \mathbb{R}^d \rightarrow \mathbb{R}$, $g(\theta) = \mathbb{E}_X \sum_{i=1}^n f(\theta; X_i)$ to obtain

$$\begin{aligned} \mathbb{E} \sum_{i=1}^n f(\theta; X_i) - \theta^T \mathbb{E} \sum_{i=1}^n \nabla f(\theta; X_i) &= \mathbb{E} \sum_{i=1}^n \text{Unif}_{a, a^d} \text{r} g'(\theta) \cdot \theta^T \\ &= 2 \mathbb{E} \sum_{i=1}^n \text{Unif}_{a, a^d} \text{r} \theta g'(\theta) - \frac{1}{2a} \text{r} g''(\theta) \theta^T. \end{aligned}$$

Moreover,

$$\begin{aligned} \mathbb{E} \sum_{i=1}^n f(\theta; X_i) - \theta^T \mathbb{E} \sum_{i=1}^n \nabla f(\theta; X_i) &\leq 2 \mathbb{E} \sum_{i=1}^n \text{r} \theta g'(\theta) - \mathbb{E} \sum_{i=1}^n \theta^T \\ &= 2 \mathbb{E} \sum_{i=1}^n \text{r} \theta g'(\theta) - \frac{a^2}{3}. \end{aligned}$$

Combining the above pieces completes the proof. \square

We state a lemma that will allow us to conveniently assume without loss of generality that the given mechanism A is ℓ_2 -accurate:

Lemma 27. Consider $\theta \sim \text{Unif}_{a, a^d}$ and $X \sim P^n$ be distributed as described above. Suppose $\mathbb{E} \sum_{i=1}^n A(\theta; X_i) \leq \alpha^2$, where $A : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $A(\theta; X_i) = (A^1(\theta; X_i), \dots, A^d(\theta; X_i))$. Assume that $A(\theta; X_i)$ is a low bias estimator of θ and that $A^j = A^l$ for all j, l . Then, for any $0 < \beta \leq \min\{\alpha^2/4d, 1/(4n^2d^2)\}$, there exists a low bias \mathcal{A} such that $\mathbb{P}(|A^j(\theta; X_i) - \theta^j| \leq \beta) \geq \beta$, $A^j = A^l$ for all $j, l \in [d]$, and $\mathbb{E} \sum_{i=1}^n \mathcal{A}(\theta; X_i) \leq 5\alpha^2$. Also, A is (ϵ, δ) -semi-DP if and only if \mathcal{A} is (ϵ, δ) -semi-DP.

Proof. By the assumption that $A^j = A^l$ for all j, l and the fact that X is drawn from a product distribution, we have

$$\mathbb{E} \sum_{i=1}^n A(\theta; X_i) \leq \alpha^2 \implies \mathbb{E} |A^1(\theta; X_i) - \theta^1|^2 \leq \frac{\alpha^2}{d}.$$

Thus, $\mathbb{E} |A^j(\theta; X_i) - \theta^j|^2 \leq \frac{\alpha^2}{d}$ for all j . (We are also assuming without loss of generality here that A^j depends only on X^j : if this is not the case, then it's easy to see by the choice of distribution that there exists another algorithm with smaller error than A satisfying this assumption.) Also,

$$\mathbb{P}(|A^j(\theta; X_i) - \theta^j| \leq \frac{2\alpha}{d}) \geq \frac{\mathbb{E} |A^j(\theta; X_i) - \theta^j|^2}{\frac{2\alpha^2}{d}} \geq \frac{1}{4}, \quad (19)$$

by Chebyshev's inequality. To construct \mathcal{A} we will use the well-known ‘‘median trick’’: run each A^j for $m = \lceil \frac{1}{\beta} \log \frac{1}{\beta} \rceil$ times, each time using a batch X_i of n independent samples in X . Then take the median of the m outputs: $\mathcal{A}^j(\theta; X_i) = \text{median}(A^j(\theta; X_{i,1}), \dots, A^j(\theta; X_{i,m}))$. Then, $\mathbb{P}(|\mathcal{A}^j(\theta; X_i) - \theta^j| \leq \beta) \geq \beta$ by a Chernoff/Hoeffding bound. Applying the law of total expectation with $\beta \leq \min\{\alpha^2/4d, 1/(4n^2d^2)\}$ establishes the expected ℓ_2 -accuracy claim for \mathcal{A} . Moreover, \mathcal{A} is low bias since for all $j \in [d]$, we have $\mathcal{A}^j(\theta; X_i) = A^j(\theta; X_{i,i})$ for some $i \in [n]$ with probability 1 and A is low bias.

We now prove the privacy claim. First note that A is (ϵ, δ) -semi-DP if and only if A^j is (ϵ, δ) -semi-DP by the advanced composition theorem (Dwork et al., 2014), since $A^j = A^l$ for all $j, l \in [d]$. Also, A^j is (ϵ^1, δ^1) -semi-DP if and only if \mathcal{A}^j is (ϵ^1, δ^1) -semi-DP by parallel composition of (semi) DP (McSherry, 2009) and the fact that $m = \lceil \frac{1}{\beta} \log \frac{1}{\beta} \rceil$. Since $\mathcal{A}^j = \mathcal{A}^l$ for all j, l by construction, we can conclude that A is (ϵ, δ) -semi-DP if and only if \mathcal{A} is (ϵ, δ) -semi-DP. \square

By Lemma 27 and the preceding discussion, it suffices to assume that the given mechanism A is low bias, $2\alpha \sqrt{d}$ - ℓ_2 -accurate with probability at least $1 - \beta$ for any $\beta \in (0, 1]$, $A^j = A^l$ for all j, l , and that A has expected ℓ_2 -mean-squared-error bounded by α^2 , for the remainder of the proof of Theorem 25. We also assume in what follows that θ is drawn uniformly at random from \mathbb{R}^d and that conditional on θ , the data X is drawn i.i.d. from the product distribution P , as described in the statement of Theorem 25.

Now, we will use Lemma 26 to lower bound the sum of expected scores $\sum_{i=1}^n \mathbb{E} x_i \cdot \theta + \sum_{i=1}^n \mathbb{E} y_i$:

Lemma 28. Let $1 \leq a \leq 100\alpha \frac{1}{d}$. Suppose $E\{A^j p_{X^j}(\theta)^2\} \leq \alpha^2$, $A^j = A^l$ is low bias for all $j, l \in \mathcal{P}$ rds, and $P\{A^j p_{X^j}(\theta^j) \mid \eta\} \leq 1/48n$. Then,

$$E \sum_{i=1}^n p_{x_i^j}(\theta^j) t_{A^j p_{X^j}}(\theta^j) \leq \frac{1}{24}$$

for all $j \in \mathcal{P}$ rds.

Proof. Note that $E\{A^j p_{X^j}(\theta^j)^2\} \leq \frac{2}{d}$ because $A^j = A^l$ for all j, l and $E\{A^j p_{X^j}(\theta)^2\} \leq \alpha^2$. By Lemma 26, we have

$$\begin{aligned} E \sum_{i=1}^n p_{x_i^j}(\theta^j) t_{A^j p_{X^j}}(\theta^j) &\leq \frac{a^2}{3} E\{A^j p_{X^j}(\theta^j)^2\} \\ &\leq \frac{1}{2a} E_{X^j} p_{A^j} p_{X^j}(\theta^j) \leq \frac{1}{2a} E_{X^j} p_{A^j} p_{X^j}(\theta^j) \leq \frac{1}{2a} \alpha^2 \\ &\leq \frac{2}{3} \alpha^2 \frac{1}{d} \\ &\leq 1/6, \end{aligned}$$

since A^j is low bias and using the assumptions on the values of a and α .

Also, by the law of total expectation, we have

$$\begin{aligned} E \sum_{i=1}^n p_{A^j p_{X^j}}(\theta^j) t_{A^j p_{X^j}}(\theta^j) &\leq 2n E \sum_{i=1}^n p_{A^j p_{X^j}}(\theta^j) t_{A^j p_{X^j}}(\theta^j) \mid \eta \\ &\leq \frac{1}{12}. \end{aligned}$$

Combining the above inequalities completes the proof. \square

As mentioned earlier, $P\{A^j p_{X^j}(\theta^j) \mid \eta\} \leq 1/48n$ (and the other assumptions in the lemma) can be ensured by Lemma 27.

Below we obtain a high probability lower bound on the sum of the scores:

Proposition 29. Let $1 \leq a \leq 100\alpha \frac{1}{d}$. Suppose $E\{A^j p_{X^j}(\theta)^2\} \leq \alpha^2$, $A^j = A^l$ is low bias for all $j, l \in \mathcal{P}$ rds, and $P\{A^j p_{X^j}(\theta^j) \mid \eta\} \leq 1/48n$. If $d \geq 200\alpha n \ln(1/\delta)$, then

$$P \sum_{i=1}^n x_{x_i}(\theta, t_{A^j p_{X^j}}(\theta^j)) \geq \frac{d}{48} \geq 1 - \delta.$$

Proof. We use the concentration inequality of Dwork et al. (2015, Theorem 36) and Lemma 28. Specifically, Lemma 28 implies that the assumptions of Dwork et al. (2015, Theorem 36) are satisfied with $X_{i,j} = x_i^j(\theta^j)$, $c = 1/2$, $Y_j = t_{A^j p_{X^j}}(\theta^j)$, $\gamma_j = 1/24$, $\alpha \sqrt{n}$, and $\mathbf{a} = \mathbf{1}_n$ as the vector of all 1's. Thus, for any $\lambda \geq 0$, we have

$$\begin{aligned} P \sum_{i=1}^n x_{x_i}(\theta, t_{A^j p_{X^j}}(\theta^j)) \geq \frac{d}{24} - \lambda &\leq \exp\left(-\frac{\lambda^2}{8d\eta^2 n^2}\right) \\ &\leq \exp\left(-\frac{\lambda^2}{8\alpha^2 n^2}\right). \end{aligned}$$

Choosing $\lambda = \frac{a}{8 \ln(1/\delta)} \geq d/48$ completes the proof. \square

Next, we provide a high probability upper bound on the sum of the squared scores:

Proposition 30. (Dwork et al., 2015) Fix $\theta \in \mathbb{R}^d$ and let $X \sim P^n$. Assume $d \geq 64pn \frac{a}{\ln p1\{\delta\}}$. Then,

$$P \left[\sum_{i=1}^n x_i \cdot \theta, \text{tAp}Xq \leq \theta \cdot y \mid 4\eta^2 d^2 \leq \delta \right] \leq \delta.$$

Proof. This is immediate from Dwork et al. (2015, Lemma 22 and Proposition 23) and the proofs of these results. \square

The next elementary lemma is taken verbatim from (Dwork et al., 2015, Lemma 24):

Lemma 31. (Dwork et al., 2015) Let $\sigma \in \mathbb{R}^n$ satisfy $\sum_{i=1}^n \sigma_i \leq A$ and $\sum_{i=1}^n \sigma_i^2 \leq B^2$. Then,

$$P \left[\sum_{i=1}^n \sigma_i \geq \frac{A}{2n} \right] \leq \frac{A^2}{2B}.$$

We are now prepared to prove Theorem 25:

Proof of Theorem 25. For notational convenience, we will assume that $\mathcal{A} = \mathcal{A}$ is low bias, symmetric ($\mathcal{A}^j = \mathcal{A}^j$), and accurate in both expected ℓ_2 norm and high probability ℓ_∞ norm. This is without loss of generality, by Lemma and Lemma 27. We first prove a): Assume y is independent of X (drawn from the same distribution). By Hoeffding's inequality,

$$P \left[\sum_{i=1}^n p_i y_i, \text{Ap}Xq \leq \sum_{i=1}^n p_i x_i \cdot \theta, \text{tAp}Xq \leq \theta \cdot y \mid \tau \geq 2\eta \frac{a}{d \ln p1\{\delta\}} \right] \leq \exp \left[-\frac{2\tau^2}{d p 4\eta^2} \right] \leq \delta.$$

Next, we prove b): Note that the assumptions on d and α imply that $d \geq 64pn \frac{a}{\ln p1\{\delta\}}$. Proposition 29 implies that

$$\sum_{i=1}^n x_i \cdot \theta, \text{tAp}Xq \leq \theta \cdot y \leq \frac{d}{48} : A$$

with probability at least $1 - \delta$. Proposition 30 implies that

$$\sum_{i=1}^n x_i \cdot \theta, \text{tAp}Xq \leq \theta \cdot y^2 \leq 4\eta^2 : B^2$$

with probability at least $1 - \delta$. By a union bound, both of the above events occur with probability at least $1 - 2\delta$. Then, Lemma 31 implies that

$$P \left[\sum_{i=1}^n x_i \cdot \theta, \text{tAp}Xq \leq \theta \cdot y \leq \frac{d}{96n} \right] \leq \frac{A^2}{2B} \leq \frac{d}{10^6 \alpha^2}.$$

Moreover, $\frac{d}{96n} \leq \tau \frac{a}{\ln p1\{\delta\}}$ by the assumption $d \geq 400\alpha n \frac{a}{\ln p1\{\delta\}}$. This completes the proof. \square

Now we are ready to prove Theorem 23:

Proof of Theorem 23. We will first prove a lower bound that is larger than the one stated in Theorem 23 by a factor of d for distributions on $\mathcal{X}^1 : \text{t} \leq 1u^d$ and $\sum_{x \in \mathcal{X}} p(x) \leq \frac{1}{a}$, where $a = \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}}} \right\}$. We will re-scale at the end to obtain Theorem 23.

By a standard reduction (see e.g. (Bun et al., 2014, Lemma 2.5)), it suffices to prove the lower bound for $\epsilon = \frac{1}{2}$. Let P be the product distribution described in Theorem 25 with $\theta^j = \text{Unif}_{[a, aq]}$ for $j \in [d]$ and $a = \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}}} \right\}$. Let $\mathcal{A} = \{A^1, \dots, A^d\}$ be $(1, \frac{1}{n_{\text{priv}}}$ -semi-DP with $\sum_{x \in \mathcal{X}} p(x) \leq \frac{1}{a}$ and $\sum_{x \in \mathcal{X}} p(x) \leq \frac{1}{a}$, where the supremum is taken over all distributions on \mathcal{X}^1 . By Lemma 27 and our earlier discussions, we assume without loss of

generality that A is low bias, $A^j = A^l$ for all $j, l \in [r]$, and that $\mathbb{P}(\exists x \in \mathcal{X} : \mathbb{E}_x \sum_{j \in [r]} \mathbb{1}_{\{A^j(x) \neq \bar{d}\}}) \leq 2\alpha \frac{r}{n}$ with arbitrarily high probability.

Note that $\alpha^2 \frac{r}{n}$ always holds by the non-private lower bound. Moreover, if $\alpha^2 \geq \frac{r}{n}$, then we are done. Thus, we may assume $\alpha^2 \leq \frac{r}{n}$. We will derive a contradiction under this assumption.

Let $X_{\text{priv}} = \{x_1, \dots, x_{n_{\text{priv}}}\}$ denote the private samples in X and $X_{\text{pub}} = \{w_1, \dots, w_{n_{\text{pub}}}\}$ denote the public samples in $X = \{X_{\text{priv}}, X_{\text{pub}}\} \sim P^n$. Let $r = \frac{2d}{400 \ln 2} \frac{1}{\epsilon}$, $t = \frac{d}{10^6 \epsilon}$, and $\gamma = \frac{1}{\rho r t \epsilon^2}$. We will show that if $\alpha \leq \frac{d}{n_{\text{priv}}}$ and $\alpha \leq \frac{d}{n_{\text{pub}}}$, then A is not $(\epsilon, \frac{1}{n_{\text{priv}}})$ -semi-DP. Assume $n_{\text{priv}} \geq r$, $n_{\text{pub}} \geq t$, and $t \geq n_{\text{priv}} \geq \frac{1}{\epsilon}$. Thus, the assumptions in Theorem 25 hold. We will show that the attack given in Algorithm 2 succeeds at identifying at least $n_{\text{pub}} - 1$ people in X with high probability: By part b) of Theorem 25 with $\delta = \gamma$, we have

$$\begin{aligned} \mathbb{P}(\exists i \in [n_{\text{pub}}] : \mathbb{1}_{\{A(x_i) = \text{IN}\}} \neq \mathbb{1}_{\{A(w_i) = \text{IN}\}}) &\leq \mathbb{P}(\exists i \in [n_{\text{pub}}] : \mathbb{1}_{\{A(x_i) = \text{IN}\}} \neq \mathbb{1}_{\{A(w_i) = \text{IN}\}}) \\ &\leq 1 - \frac{1}{\rho r t \epsilon^2} \\ &\leq 1 - \frac{1}{n_{\text{priv}}^2}. \end{aligned}$$

That is, A identifies at least $n_{\text{pub}} - 1$ individuals in the data set with high probability $\geq 1 - \frac{1}{n_{\text{priv}}^2}$. Let $v_i = \mathbb{1}_{\{A(x_i) = \text{IN}\}}$ be the indicator of the event $\mathbb{1}_{\{A(x_i) = \text{IN}\}}$. By Markov's inequality, we have

$$\mathbb{E} \sum_{i=1}^n v_i \leq \sum_{i=1}^{n_{\text{priv}}} \mathbb{P}(A(x_i) = \text{IN}) + \sum_{i=1}^{n_{\text{pub}}} \mathbb{P}(A(w_i) = \text{IN}) \leq n_{\text{priv}} \epsilon + n_{\text{pub}} \frac{1}{n_{\text{priv}}^2}.$$

Now, $\sum_{i=1}^{n_{\text{pub}}} \mathbb{1}_{\{A(w_i) = \text{IN}\}} \geq n_{\text{pub}}$, which implies that

$$\begin{aligned} \sum_{i=1}^n v_i &\geq \sum_{i=1}^{n_{\text{pub}}} \mathbb{1}_{\{A(w_i) = \text{IN}\}} \geq n_{\text{pub}} \\ &\leq n_{\text{priv}} \epsilon + n_{\text{pub}} \frac{1}{n_{\text{priv}}^2} \\ &\leq 1 + \frac{1}{n_{\text{priv}}^2} \\ &\leq 2. \end{aligned}$$

Thus, there exists a private sample $x_j \in X_{\text{priv}}$ such that

$$\mathbb{P}(A(x_j) = \text{IN}) \geq \frac{1}{2n_{\text{priv}}}.$$

Now consider the adjacent data set X^1 obtained by replacing x_j with an independent sample $y \in P$. Then, part a) of Theorem 25 implies

$$\mathbb{P}(A(x_j) = \text{IN}) \leq \frac{1}{\rho r t \epsilon^2} \leq \frac{1}{n^2} \leq \frac{1}{n_{\text{priv}}^2},$$

where the probability is taken over the random independent draws of all the data points (including y) and the mechanism A . Thus, A cannot satisfy $(\epsilon, \frac{1}{4n_{\text{priv}}})$ -semi-DP unless

$$\frac{1}{2n_{\text{priv}}} \leq e^{-\frac{1}{n_{\text{priv}}^2}} \leq \frac{1}{4n_{\text{priv}}} \leq \ln n \frac{1}{4} \leq \epsilon.$$

In particular, if $n_{\text{priv}} \geq 11$, then A cannot be $(\epsilon, \frac{1}{4n_{\text{priv}}})$ -semi-DP, which contradicts our earlier hypothesis. This proves the unscaled lower bound on \mathcal{X}^1 : $\alpha^2 \geq \frac{1}{d} \min\{n_{\text{pub}}, d\} \frac{1}{\epsilon n_{\text{priv}}^2} \geq \frac{1}{n}$.

Now we will scale the lower bound: Let $X = \{x_1, \dots, x_n\}$. Draw θ^j uniformly from \mathcal{X} for all $j \in [r]$ and then let P_θ be the distribution on X that has mean $\theta = \frac{1}{r} \sum_{j \in [r]} \theta^j$. Note that $\mathbb{E}_x \sum_{j \in [r]} \mathbb{1}_{\{A^j(x) \neq \bar{d}\}} \leq \min\{n_{\text{pub}}, d\} \frac{1}{n_{\text{priv}}}$ for any θ . Moreover, for any (ϵ, δ) -semi-DP $A : \mathcal{X}^n \rightarrow \mathcal{Y}$, there exists (ϵ, δ) -semi-DP $A^1 : \mathcal{X}^1 \rightarrow \mathcal{Y}$ such that $\mathbb{P}(A^1(x) = \bar{d}) \geq \frac{1}{n}$.

such that $\mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X)\} = \theta^1$ for $X \sim \mathcal{P}(X)$. Applying our lower bound for the MSE of A^1 , we have

$$\begin{aligned} \sup_{\mathcal{P}(X)} \inf_{\{A^1\}} \mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X) - \theta^1\}^2 &\geq \sup_{\mathcal{P}(X)} \inf_{\{A^1\}} \mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X) - \theta^1\}^2 \\ &\geq \frac{1}{d} \sup_{\mathcal{P}(X)} \mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X) - \theta^1\}^2 \\ &\geq \frac{1}{d} \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{\varepsilon^2 n_{\text{priv}}} \right\} \frac{1}{n}, \end{aligned}$$

as desired. \square

With Theorem 23 in hand, we can easily complete the proof of Theorem 22 by giving a matching upper bound (up to log factors):

Proof of Theorem 22. Lower bound: This was proved in Theorem 23.

Upper bound: First, the throw-away estimator $\mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{x\}$ is clearly $(0, 0)$ -semi-DP and has MSE

$$\begin{aligned} \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{ \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{x\} - \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{x\} \}^2 &= \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \left\{ \frac{1}{n_{\text{pub}}} \sum_{x \in \mathcal{X}_{\text{pub}}} p(x) - \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{x\} \right\}^2 \\ &= \frac{1}{n_{\text{pub}}^2} \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{ \sum_{x \in \mathcal{X}_{\text{pub}}} p(x) - n_{\text{pub}} \mathbb{E}_{\mathcal{P}(X_{\text{pub}})} \{x\} \}^2 \\ &\leq \frac{1}{n_{\text{pub}}}. \end{aligned}$$

To get the second term in the minimum in (15), consider the Gaussian mechanism $\mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X)\} = \theta^1$, where $\sigma^2 = \frac{8 \ln 2 \{ \delta q \}}{\varepsilon^2 n^2}$. We know A^1 is $(\rho \varepsilon, \delta q)$ -DP (by e.g. (Dwork et al., 2014)) since the ℓ_2 -sensitivity is $\sup_{X, X'} \|X - X'\|_2 \leq \frac{2}{n}$. Hence A^1 is $(\rho \varepsilon, \delta q)$ -semi-DP. Moreover, the MSE of A^1 is

$$\begin{aligned} \mathbb{E}_{\mathcal{P}(X)} \{ \mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X)\} - \theta^1 \}^2 &= \mathbb{E}_{\mathcal{P}(X)} \{ \mathbb{E}_{\mathcal{P}(X)} \{A^1 p(X)\} - \theta^1 \}^2 \\ &\leq \frac{8d \ln 2 \{ \delta q \}}{\varepsilon^2 n^2} \frac{1}{n}. \end{aligned}$$

This completes the proof of Theorem 22. \square

Next, we turn to the pure ε -semi-DP case.

Pure ε -Semi-DP Mean Estimation

Theorem 32 (Pure ε -semi-DP mean estimation). *Let $\varepsilon \leq d \{ 8 \}$ and either $n_{\text{pub}} \geq \frac{n^2}{d}$ or $d \geq 1$. Then, there exist absolute constants c and C , with $0 < c \leq C$, such that*

$$c \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2}{n^2 \varepsilon^2} \right\} \frac{1}{n} \leq \mathcal{M}_{\text{pop}}(\rho \varepsilon, \delta, 0, n_{\text{priv}}, n, d) \leq C \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2}{n^2 \varepsilon^2} \right\} \frac{1}{n}. \quad (20)$$

Moreover, the upper bound in (20) holds for any n_{pub}, d .

The restriction on n_{pub} is needed for our lower bound proof—via Theorem 33—to work. It seems challenging to remove this restriction, but we do believe the same lower bound holds in the complementary parameter regime. Proving this may require the invention of new techniques, making it an interesting direction for future work.

The proof of (20) will require the following intermediate result, Theorem 33, which can be viewed as a “semi-DP Fano’s inequality.”

Theorem 33. Let $\{P_v\}_{v \in \mathcal{V}} \in \mathcal{P}$ be a family of distributions on X . Let P_0 be a distribution and $\rho \in [0, 1]$ such that $P_v = \rho P_0 + (1-\rho)P_v$ for all $v \in \mathcal{V}$. Denote $\theta_v = \mathbb{E}_{x \sim P_v} f(x)$ and $\rho \in [0, 1]$: $\min_{\theta} \sum_{v \in \mathcal{V}} \theta_v |v, v^1 \in \mathcal{V}, v \neq v^1$. Let $\mathcal{D}: X^n \rightarrow \mathbb{R}^d$ be any ϵ -semi-DP estimator. Draw $V \sim \text{Unif}(\mathcal{V})$; then conditional on $V = v$, draw an i.i.d. sample $X|V=v \sim P_v^n$ containing n_{priv} private samples and n_{pub} samples. Then,

$$\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P_v(\mathcal{D}(X) \in \theta_v) \leq \rho + \frac{\epsilon \sum_{v \in \mathcal{V}} \theta_v}{2} e^{-\epsilon n_{\text{priv}}} + \frac{\epsilon \sum_{v \in \mathcal{V}} \theta_v}{2} e^{-\epsilon n_{\text{pub}}} \quad (21)$$

Remark 34. Note that the right-hand-side of (21) is similar to the second term on the right-hand-side of DP Fano's inequality (Acharya et al., 2021, Equation 5), after aligning notation. The main differences are that (21) has an extra factor of $\rho + \frac{\epsilon \sum_{v \in \mathcal{V}} \theta_v}{2}$, and n_{priv} in place of n .

Proof of Theorem 33. Our proof builds on the techniques of Barber & Duchi (2014). A key step in the proof is (22): if A is a measurable set and $v, v^1 \in \mathcal{V}$, then

$$P_v(\mathcal{D}(X) \in A) \leq e^{-\epsilon n_{\text{priv}}} P_{v^1}(\mathcal{D}(X) \in A) + \frac{\epsilon \sum_{v \in \mathcal{V}} \theta_v}{2} \quad (22)$$

Let us now prove (22): We will use upper case letters to denote random variables and lower case letters to denote the values that the random variables take. Let $B = \{B_i\}_{i=1}^n$ be i.i.d. Bernoulli random variables. Assume that the random variables $X = \{X_i\}_{i=1}^n$ are generated in the following way: first draw $W_1^0, \dots, W_n^0 \sim P_0$ i.i.d. and draw $W_1^v, \dots, W_n^v \sim P_v$ i.i.d. For each i , if $B_i = 0$, set $X_i = W_i^0$; if $B_i = 1$, set $X_i = W_i^v$. Thus, conditional on $V = v$, the random variables X_i are each distributed according to $P_v = \rho P_0 + (1-\rho)P_v$. For fixed $v^1 \in \mathcal{V}$, generate a different sample $X^1 = \{X_i^1\}_{i=1}^n$ by drawing $W_i^{v^1} \sim P_{v^1}$ i.i.d. and setting $X_i^1 = W_i^0$ if $B_i = 0$, and $X_i^1 = W_i^{v^1}$ if $B_i = 1$. Note that if $B_i = 0$, then $X_i = X_i^1$. Thus, the hamming distance between X and X^1 is

$$d_{\text{ham}}(X, X^1) = \sum_{i=1}^n B_i.$$

Now let Q denote the conditional distribution of the ϵ -semi-DP estimator \mathcal{D} given input data $(X$ or $X^1)$. For notational convenience, assume without loss of generality that $X = \{X_1, \dots, X_{n_{\text{priv}}}, X_{n_{\text{pub}}}\}$ and $X^1 = \{X_1^1, \dots, X_{n_{\text{priv}}}^1, X_{n_{\text{pub}}}^1\}$. Then, for any fixed sequence $b = \{b_1, \dots, b_{n_{\text{priv}}}, \mathbf{0}_{n_{\text{pub}}}\} \in \{0, 1\}^{n_{\text{priv}}} \times \{0, 1\}^{n_{\text{pub}}}$, we have

$$Q(\mathcal{D}(X) \in A | X_i = W_i^0 + b_i) - Q(\mathcal{D}(X) \in A | X_i = W_i^1 + b_i) \leq e^{-\epsilon b_i} Q(\mathcal{D}(X) \in A | X_i = W_i^0 + b_i) - Q(\mathcal{D}(X) \in A | X_i = W_i^1 + b_i) \quad (23)$$

Then by a union bound and disjointness of the balls $\frac{1}{2}B_{\rho}(\theta_{v^1}, \theta_{v^2})$, we have

$$P_{\text{succ}} \leq 1 - \frac{1}{|V|} \sum_{v, v^1 \in V, v^2 \in V} P_{v^1} \cdot \mathbb{P} \left[B_{\rho}(\theta_{v^1}, \theta_{v^2}) \right].$$

An application of (22) yields

$$\begin{aligned} P_{\text{succ}} &\leq 1 - \frac{1}{|V|} \sum_{v, v^1 \in V, v^2 \in V} e^{-\epsilon n_{\text{priv}} \rho S} P_{v^1} \cdot \mathbb{P} \left[B_{\rho}(\theta_{v^1}, \theta_{v^2}) \right] \leq 1 - \frac{\rho}{2} \frac{pq^{\eta_{\text{pub}}}}{2} \\ &\leq 1 - e^{-\epsilon n_{\text{priv}} \rho S} \frac{\rho}{2} \frac{pq^{\eta_{\text{pub}}}}{2} \leq 1 - \frac{\rho}{2} \frac{pq^{\eta_{\text{pub}}}}{2}. \end{aligned}$$

Re-arranging this inequality leads to

$$P_{\text{succ}} \leq \frac{1 - \frac{\rho}{2} \frac{pq^{\eta_{\text{pub}}}}{2}}{1 - e^{-\epsilon n_{\text{priv}} \rho S}}$$

and hence

$$1 - P_{\text{succ}} \geq \frac{\frac{\rho}{2} \frac{pq^{\eta_{\text{pub}}}}{2}}{1 - e^{-\epsilon n_{\text{priv}} \rho S}}.$$

This last inequality is equivalent to the inequality stated in Theorem 33. \square

While we state Theorem 33 for mean estimation, it holds more generally for estimating any population statistic $\theta : \mathcal{P} \rightarrow \mathbb{R}^d$. However, this additional generality will not be necessary for our purposes. With Theorem 33 in hand, we now turn to the proof of Theorem 22.

Proof of Theorem 32. Lower Bounds: We begin by proving (20). First suppose $n_{\text{pub}} \geq \frac{n}{d}$ and $d \geq 8$.

Choose a finite subset $V \subseteq \mathbb{R}^d$ such that $|V| \geq 2^{d/2}$, $\|v\| \leq 1$, and $\|v - v^1\| \geq \frac{1}{8}$ for all $v, v^1 \in V, v \neq v^1$. The existence of such a set of points is well-known (see e.g. the Gilbert-Varshamov construction). Define P_0 to be the point mass distribution on $\{X = 0\}$ and P_v to be point mass on $\{X = v\}$ for $v \in V$. For $v \in V$, let $P_v : p_1 - pqP_0 - pP_v$ for some $p \in (0, 1)$ to be specified later. Note that if $X \sim P_v$, then $\|X\| \leq 1$ with probability 1. Thus, P_v is a valid distribution in the class \mathcal{P} of bounded (by 1) distributions on \mathcal{B} that we are considering. Also, note that $\theta_v : \mathbb{E}_{P_v} [X] = pv$. Further,

$$\rho \leq \frac{1}{2} \min_{v, v^1 \in V, v \neq v^1} \|\theta_v - \theta_{v^1}\| \geq \frac{p}{8}$$

by construction.

Now we use the classical reduction from estimation to testing (see (Barber & Duchi, 2014) for details) to lower bound the MSE of any ϵ -semi-DP estimator $\hat{\theta}$ by

$$\begin{aligned} \sup_{P \in \mathcal{P}} \mathbb{E}_{X \sim P} \|\hat{\theta} - X\|^2 &\geq \mathbb{E}_{X \sim P} \|\hat{\theta} - X\|^2 \geq \frac{1}{|V|} \sum_{v \in V} \mathbb{E}_{P_v} \|\hat{\theta} - X\|^2 \geq \frac{\rho}{2} \sum_{v \in V} \mathbb{E}_{P_v} \|\hat{\theta} - X\|^2 \\ &\geq \frac{p^2}{8} \frac{1}{2} \frac{\sum_{v \in V} \|v\|^2}{|V|} \frac{1 - e^{-\epsilon n_{\text{priv}} \rho S}}{1 - e^{-\epsilon n_{\text{priv}} \rho S}} \frac{pq^{\eta_{\text{pub}}}}{2} \\ &\geq \frac{p^2}{64} \frac{2^{d/2}}{2} \frac{1 - e^{-\epsilon n_{\text{priv}} \rho S}}{2^{d/2}} \frac{pq^{\eta_{\text{pub}}}}{2} \frac{1}{1 - e^{-\epsilon n_{\text{priv}} \rho S}} \end{aligned}$$

where we used Theorem 33 in the second inequality. Since we assumed $d \geq 8$, we have $2^{d/2} \geq 1 - e^{-\epsilon n_{\text{priv}} \rho S}$ and hence

$$\begin{aligned} \sup_{P \in \mathcal{P}} \mathbb{E}_{X \sim P} \|\hat{\theta} - X\|^2 &\geq \frac{p^2}{64} \frac{pq^{\eta_{\text{pub}}}}{2} \frac{e^{-\epsilon n_{\text{priv}} \rho S}}{1 - e^{-\epsilon n_{\text{priv}} \rho S}} \\ &\geq \frac{p^2}{64} \frac{pq^{\eta_{\text{pub}}}}{4}. \end{aligned}$$

for any $p \geq \frac{d}{4n^\alpha} - \frac{1}{n}$. Now choose

$$p = \min \left\{ \frac{d}{4n^\alpha}, \frac{1}{n}, \frac{1}{2^{\frac{1}{\alpha}} n_{\text{pub}}^\alpha} \right\}.$$

By assumption, there exists an absolute constant k such that $n_{\text{pub}} \geq k \frac{n^\alpha}{d}$. Thus, if $n^\alpha \geq 2$, then

$$\begin{aligned} p &\geq \frac{d}{4n^\alpha} \geq \frac{1}{4} \frac{d}{n^\alpha} \geq \frac{1}{4} \frac{d}{n^\alpha} \frac{k n^\alpha}{d} \\ &\geq \frac{1}{4} \frac{d}{n^\alpha} \frac{n^\alpha}{d} \frac{1}{k} \\ &\geq \frac{1}{4} \frac{1}{2^k} \\ &\geq \frac{1}{4^k}. \end{aligned}$$

On the other hand, if $n^\alpha < 2$, then $n_{\text{pub}} \geq 2k \frac{n^\alpha}{d}$. Also, note that $p \geq \min \left\{ \frac{d}{8n^\alpha}, \frac{1}{2^{\frac{1}{\alpha}} n_{\text{pub}}^\alpha} \right\}$ by the assumption that $d \geq 8\epsilon$. Therefore,

$$\sup_{P, P'} \mathbb{E}_X \left[\sum_{i=1}^n \ell_i(x_i) \right] - \mathbb{E}_X \left[\sum_{i=1}^n \ell_i(x_i) \right]^2 \leq \frac{1}{256} \frac{1}{4^k} \min \left\{ \frac{d}{8n^\alpha}, \frac{1}{2^{\frac{1}{\alpha}} n_{\text{pub}}^\alpha} \right\}^2.$$

Combining the above inequality with the non-private lower bound of $\frac{1}{2} \log \frac{1}{\epsilon}$ for mean estimation proves the lower bound in (20).

Now consider the alternative case in which $d \leq 1$ (i.e. $d \leq k$ for some absolute constant $k \in \mathbb{N}$), but n_{pub} is arbitrary. Then we will prove that the lower bound in (20) holds for $d \leq 1$, for any $\delta \in (0, \epsilon]$ and $\epsilon \geq 1$. By taking a k -fold product distribution, this will suffice to complete the proof of the lower bound in (20). To that end, we will use Le Cam's method and build on the techniques in (Barber & Duchi, 2014; Fallah et al., 2022). The key novel ingredient is the following extension of Fallah et al. (2022, Lemma 3) to the ϵ, δ -semi-DP setting:

Lemma 35. Let $\mathcal{D} : X^n \rightarrow \mathbb{R}^n$ be ϵ, δ -semi-DP and let P_1, P_2 be distributions on X such that P_1 is absolutely continuous w.r.t. P_2 . Denote the conditional distribution of \mathcal{D} given X by Q and let $Q_j|A$ denote $Q_j|_{\mathcal{D}^{-1}(A)}$ for any measurable set A . Then

$$\|Q_1 - Q_2\|_{TV} \geq \min \left\{ \frac{\epsilon}{2} \frac{D_{KL}(P_1, P_2)}{D_{KL}(P_1, P_2)}, \frac{\delta}{2} \frac{D_{KL}(P_1, P_2)}{D_{KL}(P_1, P_2)} \right\}.$$

Let us defer the proof of Lemma 35 for now. We will now use Lemma 35 to prove the lower bound in (20) for $d \leq 1$ and $\epsilon \geq 1$. Define distributions P_1, P_2 on $\{0, 1\}$ as follows:

$$P_1(0) = \frac{1}{2} - \frac{\gamma}{2}, \quad P_1(1) = \frac{1}{2} + \frac{\gamma}{2}$$

for some $\gamma \in (0, 1/2]$ to be chosen later. Clearly $P_1, P_2 \in \mathcal{P}(\{0, 1\})$ (i.e. they are bounded by 1 with probability 1). Also, $\mathbb{E}_{P_1}(x) = \frac{1}{2} + \gamma$ and $\mathbb{E}_{P_2}(x) = \frac{1}{2} - \gamma$, so (P_1, P_2) is a γ -packing of $\{0, 1\}$. Thus, by Le Cam's method (see (Barber & Duchi, 2014) for details), for any ϵ, δ -semi-DP \mathcal{D} , we have

$$\sup_{P, P'} \mathbb{E}_X \left[\sum_{i=1}^n \ell_i(x_i) \right] - \mathbb{E}_X \left[\sum_{i=1}^n \ell_i(x_i) \right]^2 \geq \frac{\gamma^2}{8} \|Q_1 - Q_2\|_{TV}.$$

Now, applying Lemma 35 and the assumption $\delta \asymp \varepsilon \asymp 1$ yields

$$\begin{aligned} & \sup_{PPP} E_{X \sim P^{n,p}} \mathbb{E} \left[\sum_{i=1}^k \ell_i(x_i) \right] - E_{X \sim P^{n,p}} \left[\sum_{i=1}^k \ell_i(x_i) \right]^2 \\ & \leq \frac{\gamma^2}{8} \frac{1}{\min \left\{ \frac{n}{2} D_{\text{KL}}(P_1, P_2), 6\gamma P_1, P_2 \right\}_{\text{TV}}} \frac{1}{n_{\text{priv}} \varepsilon} \frac{1}{\min \left\{ \frac{n}{2} D_{\text{KL}}(P_1, P_2), 6\gamma n_{\text{priv}} \varepsilon \right\}} \\ & \leq \frac{\gamma^2}{8} \frac{1}{\min \left\{ \frac{n}{2} D_{\text{KL}}(P_1, P_2), 6\gamma n_{\text{priv}} \varepsilon \right\}} \frac{1}{\min \left\{ \frac{n_{\text{pub}}}{2} D_{\text{KL}}(P_1, P_2), 6\gamma n_{\text{pub}} \varepsilon \right\}} \\ & \leq \frac{\gamma^2}{8} \frac{1}{\gamma \min \left\{ \frac{3n}{2}, 6n_{\text{priv}} \varepsilon \right\}} \frac{1}{\min \left\{ \frac{3n_{\text{pub}}}{2}, 6n_{\text{pub}} \varepsilon \right\}}. \end{aligned} \quad (24)$$

In the second inequality we used the fact that $\|P_1 - P_2\|_{\text{TV}} \leq \frac{1}{2} \frac{1}{2} \frac{1}{2} \leq \frac{1}{2}$ and $D_{\text{KL}}(P_1, P_2) \leq 3\gamma^2$ for $\gamma \leq 1/2$.

Now we will choose γ to (approximately) maximize the right-hand side of (24). Suppose $\min \left\{ \frac{3n}{2}, 6n_{\text{priv}} \varepsilon \right\} \leq \frac{3n_{\text{pub}}}{2}$. Then choosing $\gamma = \frac{1}{3} \frac{2}{3n}$ yields

$$\sup_{PPP} E_{X \sim P^{n,p}} \mathbb{E} \left[\sum_{i=1}^k \ell_i(x_i) \right] - E_{X \sim P^{n,p}} \left[\sum_{i=1}^k \ell_i(x_i) \right]^2 \leq \frac{k}{n}$$

for some absolute constant $k \geq 0$. Our assumption that $\min \left\{ \frac{3n}{2}, 6n_{\text{priv}} \varepsilon \right\} \leq \frac{3n_{\text{pub}}}{2}$ implies that there exists

$k' \geq 0$ such that $n \leq k' \max \left\{ n_{\text{priv}}^2 \varepsilon^2, n_{\text{pub}} \right\}$. Thus, $k/n \leq \frac{k'}{k'} \min \left\{ \frac{1}{n_{\text{priv}}^2 \varepsilon^2}, \frac{1}{n_{\text{pub}}} \right\}$, which gives the desired lower bound in (20).

Suppose instead that $\min \left\{ \frac{3n}{2}, 6n_{\text{priv}} \varepsilon \right\} \geq \frac{3n_{\text{pub}}}{2}$. Then choose $\gamma = \frac{2}{3} \frac{1}{6n_{\text{priv}} \varepsilon}$. Then, there are constants $k, c \geq 0$ such that

$$\begin{aligned} & \sup_{PPP} E_{X \sim P^{n,p}} \mathbb{E} \left[\sum_{i=1}^k \ell_i(x_i) \right] - E_{X \sim P^{n,p}} \left[\sum_{i=1}^k \ell_i(x_i) \right]^2 \leq \frac{\gamma^2}{8} \frac{1}{\gamma \min \left\{ \frac{3n}{2}, 6n_{\text{priv}} \varepsilon \right\}} \frac{1}{\min \left\{ \frac{3n_{\text{pub}}}{2}, 6n_{\text{pub}} \varepsilon \right\}} \\ & \leq \frac{k}{n_{\text{priv}}^2 \varepsilon^2} \frac{1}{n_{\text{pub}}} \\ & \leq c \min \left\{ \frac{1}{n_{\text{priv}}^2 \varepsilon^2}, \frac{1}{n_{\text{pub}}} \right\}. \end{aligned}$$

Combining the above inequality with the non-private lower bound $\frac{1}{n} \sum_{i=1}^k \ell_i(x_i) \geq \frac{1}{n} \sum_{i=1}^k \ell_i(x_i) - \frac{1}{n}$ completes the proof of the lower bound in (20), assuming the truth of Lemma 35.

It remains to prove Lemma 35. To that end, fix $k \in \mathbb{N}$, $\varepsilon \in (0, 1]$ and denote by \mathcal{Q} the marginal distribution of \mathcal{P} given $X_1, \dots, X_k \sim P_1$ (i.i.d.) and $X_{k+1}, \dots, X_n \sim P_2$ (i.i.d.); i.e. for measurable A ,

$$\mathcal{Q}(A) = \int \int \mathbb{1}_A(x_1, \dots, x_n) dP_1^k \otimes dP_2^{n-k}.$$

Note that if $k = 0$, then $\mathcal{Q} = Q_2$. We have

$$\|Q_1 - Q_2\|_{\text{TV}} \leq \frac{1}{2} \left(\frac{1}{2} + \frac{1}{2} \right) = 1.$$

Also,

$$\min_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \frac{1}{n} \sum_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \left[\frac{n}{2} D_{\text{KL}}(P_1, P_2) + 2 \mathbb{E}_{P_1, P_2} \text{TV} \right] \min_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \left[k e^{\epsilon} - 1 - \delta \right] \frac{n_{\text{pub}}}{2} D_{\text{KL}}(P_1, P_2), \quad (26)$$

so it suffices to upper bound the sum of the terms of $a + b$ by the left-hand-side of (26). First, we deal with a : for any k , we have

$$\mathbb{E}_{Q_1, Q_2} \mathbb{E}_{\text{TV}}^2 \leq \mathbb{E}_{P_1^n, P_1^k P_2^n} \mathbb{E}_{\text{TV}}^2 \quad (27)$$

$$\leq \frac{1}{2} D_{\text{KL}}(P_1^n, P_1^k P_2^n) \quad (28)$$

$$\leq \frac{n}{2} D_{\text{KL}}(P_1, P_2), \quad (29)$$

by the data processing inequality for f -divergences, Pinsker's inequality, and the chain-rule for KL-divergences (see, e.g. (Duchi, 2021) for a reference on these facts). Thus, it remains to show

$$\mathbb{E}_{Q_1, Q_2} \mathbb{E}_{\text{TV}} \leq 2 \mathbb{E}_{P_1, P_2} \text{TV} \min_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \left[k e^{\epsilon} - 1 - \delta \right] \quad (30)$$

for $k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u$. If $k = 0$, (30) is trivial. Assume $k > n_{\text{priv}}$. Now, for any measurable A , we may write

$$\mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} = \int_{\mathbb{R}^n} \mathbb{P}_{n_{\text{priv}}} dx \int_{\mathbb{R}^n} dP_2^{\text{pub}} \mathbb{P}_{n_{\text{priv}}} dx, \quad (31)$$

where

$$\mathbb{P}_{n_{\text{priv}}} dx = \int_{\mathbb{R}^n} Q_1 \mathbb{P} A | X_{1:n} = x_{1:n} dP_1^{n_{\text{priv}}} \mathbb{P}_{x_{1:n_{\text{priv}}}} dP_2^{n_{\text{priv}}} \mathbb{P}_{x_{1:n_{\text{priv}}}} dx.$$

By (31), it suffices to show that $\mathbb{E}_{\mathbb{P}_{n_{\text{priv}}}} \mathbb{E}_{\text{TV}} \leq 2 \mathbb{E}_{P_1, P_2} \text{TV} \min_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \left[k e^{\epsilon} - 1 - \delta \right]$ for all $x_{n_{\text{priv}}}$. To do so, let $x_{1:n}^i = \mathbb{P}_{x_1, \dots, x_{i-1}, x_i^i, x_{i+1}, \dots, x_n}$ for some $i \in n_{\text{priv}}$. Then

$$\left| \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n} - \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n}^i \right| \leq e^{-\epsilon} \int Q_1 \mathbb{P} A | X_{1:n} = x_{1:n}^i dx \quad (32)$$

since \mathbb{P} is ϵ, δ -semi-DP. Moreover, by the proof of Fallah et al. (2022, Lemma 3), we have

$$\mathbb{E}_{\mathbb{P}_{n_{\text{priv}}}} \mathbb{E}_{\text{TV}} \leq \int_{i=1}^{n_{\text{priv}}} \int_{\mathbb{R}^n} \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n} - \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n}^i dx$$

$$\leq \int_{i=1}^{n_{\text{priv}}} \int_{\mathbb{R}^n} \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n} - \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n}^i dx$$

Applying the triangle inequality and (32), we get

$$\left| \mathbb{E}_{\mathbb{P}_{n_{\text{priv}}}} \mathbb{E}_{\text{TV}} \right|$$

$$\leq \int_{i=1}^{n_{\text{priv}}} \int_{\mathbb{R}^n} e^{-\epsilon} \int Q_1 \mathbb{P} A | X_{1:n} = x_{1:n}^i dx \quad \mathbb{E}_{Q_1} \mathbb{E}_{Q_2} \mathbb{E}_{A} | X_{1:n} = x_{1:n}^i dx$$

$$\leq 2 \mathbb{E}_{P_1, P_2} \text{TV} \min_{k \in \mathbb{P}^{\text{t0}}, n_{\text{priv}}, u} \left[k e^{\epsilon} - 1 - \delta \right],$$

as desired. This completes the proof of Lemma 35 and hence the proof of the lower bound in (20).

Upper bound: First, the throw-away estimator $\mathbb{E}_{\mathbb{P}_{n_{\text{pub}}}} \mathbb{E}_{X_{n_{\text{pub}}}}$ is clearly $(0, 0)$ -semi-DP and has MSE

$$\mathbb{E}_{X_{n_{\text{pub}}}} \left\{ \mathbb{E}_{\mathbb{P}_{n_{\text{pub}}}} \left[\int_{\mathbb{R}^n} \mathbb{P} x \right] - \mathbb{E}_{X_{n_{\text{pub}}}} \left[\int_{\mathbb{R}^n} \mathbb{P} x \right] \right\}^2$$

$$\leq \frac{1}{n_{\text{pub}}} \mathbb{E}_{X_{n_{\text{pub}}}} \left\{ \int_{\mathbb{R}^n} \mathbb{P} x \right\}^2$$

$$\leq \frac{1}{n_{\text{pub}}^2} \mathbb{E}_{X_{n_{\text{pub}}}} \left\{ \int_{\mathbb{R}^n} \mathbb{P} x \right\}^2$$

$$\leq \frac{1}{n_{\text{pub}}}.$$

To get the second term in the minimum in (20), consider the Laplace mechanism $A_{\rho}X_{\mathcal{Q}} = X_{\mathcal{Q}} + \rho(L_1, \dots, L_d)\mathbf{q}$, where $L_i \sim \text{Lap}(2/\epsilon, \mathbf{0})$ are i.i.d. mean-zero Laplace random variables. We know A is ϵ -DP by (Dwork et al., 2014), since the ℓ_1 -sensitivity is $\sup_{X, X'} \|X - X'\|_1 = \frac{1}{n} \sup_{x, x'} \|x - x'\|_1 \leq \frac{2}{n}$. Hence A is ϵ -semi-DP. Moreover, A has MSE

$$\begin{aligned} \mathbb{E}_{X \sim P^n} \{A_{\rho}X_{\mathcal{Q}} - \mathbb{E}_{X \sim P} \{x\}\}^2 &\leq 2\mathbb{E} \{A_{\rho}X_{\mathcal{Q}} - X\}^2 = 2\mathbb{E} \{X - \mathbb{E}_{X \sim P} \{x\}\}^2 \\ &\leq 2d \text{Var} \text{Lap}(2/\epsilon, \mathbf{0}) = \frac{2}{n} \\ &= \frac{16d^2}{n^2\epsilon^2} = \frac{2}{n}. \end{aligned}$$

This completes the proof of (20). \square

E.1.3. AN ‘‘EVEN MORE OPTIMAL’’ SEMI-DP ALGORITHM FOR MEAN ESTIMATION

Lemma 36 (Re-statement of Lemma 7). *Recall the definition of $\mathcal{P}_{\rho B, V, \mathcal{Q}}$ (Definition 6): $\mathcal{P}_{\rho B, V, \mathcal{Q}}$ denotes the collection of all distributions P on \mathbb{R}^d such that for any $x \sim P$, we have $\|x\| \leq B$ P -almost surely and $\text{Var}x_{\mathcal{Q}} = V^2$. Then, the error of the ρ -semi-zCDP throw-away algorithm $A_{\rho}X_{\mathcal{Q}} = \frac{1}{n_{\text{pub}}} \sum_{x \sim P_{\text{pub}}} x$ is*

$$\sup_{P \in \mathcal{P}_{\rho B, V, \mathcal{Q}}} \mathbb{E}_{X \sim P^n} \{A_{\rho}X_{\mathcal{Q}} - \mathbb{E}_{X \sim P} \{x\}\}^2 \leq \frac{V^2}{n_{\text{pub}}}.$$

The minimax error of the ρ -semi-zCDP Gaussian mechanism $G_{\rho}X_{\mathcal{Q}} = X \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ is

$$\inf_{\text{-zCDP } G} \sup_{P \in \mathcal{P}_{\rho B, V, \mathcal{Q}}} \mathbb{E}_{G; X \sim P^n} \{G_{\rho}X_{\mathcal{Q}} - \mathbb{E}_{X \sim P} \{x\}\}^2 = \frac{2dB^2}{\rho n^2} = \frac{V^2}{n}. \quad (33)$$

Proof. For throw-away, the i.i.d. data assumption implies

$$\mathbb{E} \{A_{\rho}X_{\mathcal{Q}} - \mathbb{E}x\}^2 = \frac{1}{n_{\text{pub}}^2} \mathbb{E} \left\{ \sum_{x \sim P_{\text{pub}}} (x - \mathbb{E}x) \right\}^2 = \frac{V^2}{n_{\text{pub}}}.$$

The Gaussian mechanism $G_{\rho}X_{\mathcal{Q}} = X \sim \mathcal{N}(\mathbf{0}, \frac{2B^2}{\rho^2} \mathbf{I}_d)$ is ρ -zCDP by (Bun & Steinke, 2016, Proposition 1.6) since the ℓ_2 -sensitivity is bounded by $\|x - x'\|_2 \leq \frac{2B}{\rho}$. Moreover, this sensitivity bound is tight: consider any P such that $x \sim \rho B, \mathbf{0}_d$ and $x' \sim \rho B, \mathbf{0}_d$ are in the support of P . Then fix any y in the support of P and consider the adjacent data sets $X = \rho x, y, \dots, y$ and $X' = \rho x', y, \dots, y$. We have $\|X - X'\|_2 = \frac{1}{n} \|x - x'\|_2 = \frac{2B}{n}$. Additionally, if the variance of the additive isotropic Gaussian noise σ^2 is smaller than $\frac{2}{\rho^2} = \frac{2B^2}{\rho^2}$, then the Gaussian mechanism is not ρ -zCDP (Bun & Steinke, 2016). Thus, $G_{\rho}X_{\mathcal{Q}}$ is the ρ -zCDP Gaussian mechanism with the smallest noise variance σ^2 . Hence the infimum in (33) is attained by G . Finally, for any $P \in \mathcal{P}_{\rho B, V, \mathcal{Q}}$, the MSE of G is

$$\begin{aligned} \mathbb{E}_{G; X \sim P^n} \{G_{\rho}X_{\mathcal{Q}} - \mathbb{E}_{X \sim P} \{x\}\}^2 &= \mathbb{E} \{G_{\rho}X_{\mathcal{Q}} - X\}^2 = \mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E}x_i \right\}^2 \\ &= \frac{2dB^2}{\rho n^2} = \frac{V^2}{n}. \end{aligned}$$

\square

Proposition 37 (Re-statement of Proposition 8). *Recall the definition of A_r from (3). A_r is ρ -semi-zCDP. Also, there exists $r \geq 0$ such that*

$$\sup_{P \in \mathcal{P}_{\rho B, V, \mathcal{Q}}} \mathbb{E}_{X \sim P^n} \{A_r X_{\mathcal{Q}} - \mathbb{E}_{X \sim P} \{x\}\}^2 = \min \left\{ \frac{V^2}{n_{\text{pub}}}, \frac{2dB^2}{\rho n^2}, \frac{V^2}{n} \right\}. \quad (34)$$

Further, if $\frac{V^2}{n_{\text{pub}}} \asymp \frac{2dB^2}{r^2}$, then the quantitative advantage of A_r is

$$\sup_{P \in \mathcal{P}(B, V)} \mathbb{E}_X \{A_r(p, X) - \mathbb{E}_X \{A_r(x, S)\}^2 \asymp \frac{q}{s^2} \min \left\{ \frac{V^2}{n_{\text{pub}}}, \frac{2dB^2}{\rho n^2}, \frac{V^2}{n} \right\}, \quad (35)$$

where $q = 2 \frac{n_{\text{priv}} V^2}{dB^2}$ and $s = \frac{V^2 n_{\text{priv}}}{B^2 n_{\text{pub}}}$.

Proof. Privacy: Note that the ℓ_2 -sensitivity of $M(p, X)$ is

$$\sup_{x, x'} \left\| \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} r x_i - \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} r x'_i \right\|_2 \leq \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} r \|x_i - x'_i\|_2 \leq \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} r \cdot 2B \leq 2rB.$$

Recall that the Gaussian mechanism guarantees ρ -zCDP whenever $\sigma^2 \geq \frac{2}{\rho} \sum_{i \in \mathcal{X}_{\text{pub}}} r^2 B^2$ (Bun & Steinke, 2016, Proposition 1.6). Thus, A_r is ρ -semi-zCDP for $\sigma_r^2 \geq \frac{2rBq^2}{\rho} = \frac{2B^2 r^2}{\rho}$.

Error bounds: Let $P \in \mathcal{P}(B, V)$. We have

$$\begin{aligned} \mathbb{E}_X \{A_r(p, X) - \mathbb{E}_X \{A_r(x, S)\}^2 &= \frac{2dB^2 r^2}{\rho} \mathbb{E}_{\mathcal{X}_{\text{priv}}} \{r p x\} - \mathbb{E}_{\mathcal{X}_{\text{pub}}} \left\{ \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} p x_i \right\}^2 \\ &= \frac{2dB^2 r^2}{\rho} \mathbb{E}_{\mathcal{X}_{\text{priv}}} \{r^2 V^2\} - \mathbb{E}_{\mathcal{X}_{\text{pub}}} \left\{ \frac{1}{n_{\text{pub}}} \sum_{i \in \mathcal{X}_{\text{pub}}} r^2 V^2 \right\}^2 \\ &= \frac{2dB^2 r^2}{\rho} n_{\text{priv}} r^2 V^2 - \frac{\rho}{n_{\text{pub}}} \frac{n_{\text{priv}} r^2 q^2}{n_{\text{pub}}} V^2, \end{aligned} \quad (36)$$

using independence of the Gaussian noise and the data, basic properties of variance, and the fact that the data is i.i.d.

To prove (34), let

$$J(p, r) = \frac{2dB^2 r^2}{\rho} n_{\text{priv}} r^2 V^2 - \frac{\rho}{n_{\text{pub}}} \frac{n_{\text{priv}} r^2 q^2}{n_{\text{pub}}} V^2.$$

We compute first and second derivatives of J :

$$\frac{d}{dr} J(p, r) = 2r \frac{2dB^2}{\rho} n_{\text{priv}} V^2 - 2n_{\text{priv}} \frac{V^2}{n_{\text{pub}}} \rho = 2n_{\text{priv}} r q$$

and

$$\frac{d^2}{dr^2} J(p, r) = 2 \frac{2dB^2}{\rho} n_{\text{priv}} V^2 - 2n_{\text{priv}}^2 \frac{V^2}{n_{\text{pub}}}.$$

Since J is strongly convex, it has a unique minimizer r^* which satisfies $\frac{d}{dr} J(p, r) = 0$. We find

$$r^* = \frac{n_{\text{priv}} V^2}{n_{\text{pub}}} \frac{2dB^2}{\rho} n_{\text{priv}} V^2 = \frac{n_{\text{priv}}^2 V^2}{n_{\text{pub}}}.$$

One can verify that $r^* \geq 0$ and $r^* \leq \frac{1}{n}$, since $1 \geq n_{\text{priv}} \geq n$ and $\rho \geq 8$ by assumption. Thus, $J(p, r^*) = \min_p J(p, 0), J(p, \frac{1}{n})$, which yields (34) by Lemma 7.

To prove (35), we will choose a different r : $r = \frac{KV}{B} \frac{\partial \rho}{\partial n_{\text{pub}}}$ for $K \geq 0$ to be determined. Then by (36), we have

$$\begin{aligned} \mathbb{E}_{X, P^n} \{A_r \rho(X, q) - \mathbb{E}_{X, P} \rho(x, S)\}^2 &= \frac{2dB^2 r^2}{\rho} + n_{\text{priv}} r^2 V^2 + \frac{\rho}{n_{\text{pub}}} \frac{q^2}{V^2} \\ &= K^2 \frac{V^2}{n_{\text{pub}}} + 2 \frac{n_{\text{priv}} \rho V^2}{dB^2} + \frac{V^2}{n_{\text{pub}}} \left(1 + \frac{KV}{B} \frac{\partial \rho}{\partial n_{\text{pub}}}\right)^2 \\ &= \frac{V^2}{n_{\text{pub}}} \left(qK^2 + 1 + \frac{KV}{B} \frac{\partial \rho}{\partial n_{\text{pub}}} \right)^2, \end{aligned}$$

where $q = 2 \frac{n_{\text{priv}} V^2}{dB^2}$. Now, letting $s = \frac{V}{B} \frac{\partial \rho}{\partial n_{\text{pub}}}$ and choosing $K = \frac{s}{q s^2}$ gives

$$\begin{aligned} \mathbb{E}_{X, P^n} \{A_r \rho(X, q) - \mathbb{E}_{X, P} \rho(x, S)\}^2 &\leq \frac{V^2}{n_{\text{pub}}} \left(qK^2 + 1 + \frac{KV}{B} \frac{\partial \rho}{\partial n_{\text{pub}}} \right)^2 \\ &\leq \frac{V^2}{n_{\text{pub}}} \left(\frac{q^2}{q^2} + \frac{qs^2}{2qs^2} + \frac{s^2}{s^4} \right). \end{aligned}$$

Finally, the assumption $\frac{V^2}{n_{\text{pub}}} \leq \frac{2dB^2}{n^2}$ implies $\frac{V^2}{n_{\text{pub}}} \leq \min \left\{ \frac{V^2}{n_{\text{pub}}}, \frac{2dB^2}{n^2}, \frac{V^2}{n} \right\}$, completing the proof. \square

E.2. Optimal Semi-DP Empirical Risk Minimization

Practical Applications of Semi-DP ERM Beyond ML: Semi-DP ERM has numerous applications beyond training ML models. For example, consider semi-DP optimization of energy consumption in smart grids or semi-DP optimization of the total capacity of a multi-user wireless communication system. In these systems, the goal is to optimize the current performance of the system given existing users (e.g., optimize current beamforming strategies in wireless communications). Some users may opt-in to share their data (e.g. electricity consumption pattern) and some users may not. Thus, the problem is naturally a semi-DP ERM problem.

Theorem 38 (Complete statement of Theorem 10). *There exist absolute constants c_0 and C_0 , with $0 < c_0 \leq C_0$, such that*

$$c_0 LD \min_{\frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon}} R_{\text{ERM}}(\rho, \varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \leq C_0 LD \min_{\frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon}} R_{\text{ERM}}^*.$$

Further, if $\mu \geq 0$, then there exist absolute constants $0 < c_1 \leq C_1$ such that

$$c_1 LD \min_{\frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon}} R_{\text{ERM}}^* \leq C_1 \frac{L^2}{\mu} \min_{\frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon}} R_{\text{ERM}}^{\#} + 2.$$

Proof. Lower Bounds: Given a lower bound for empirical mean estimation, Bassily et al. (Bassily et al., 2014) show how to prove excess risk lower bounds for convex and strongly convex ERM by reducing these problems to mean estimation. Thus, our lower bounds follow immediately by combining the lower bound in Theorem 21 with the reduction in (Bassily et al., 2014). Roughly, the reduction works as follows:

In the strongly convex case, we simply take $f(w, x, q) = \frac{1}{2} \|w - x\|^2$ on $W \times B$; $f(\cdot, x, q)$ is 1-uniformly-Lipschitz and 1-strongly convex. Moreover, for any ε -semi-DP A with output $w_{\text{priv}} = A(\rho, X, q)$, we have

$$\mathbb{E} \mathbb{P}_{X, P} \rho(w_{\text{priv}}, q) - \mathbb{P}_X \left\{ \frac{1}{2} \mathbb{E} \rho(A(\rho, X, q), X) \right\}^2.$$

Applying Theorem 21 and then scaling $f \leftarrow \frac{L}{D} f$ and $W \leftarrow DW$ and $X \leftarrow DX$ completes the proof.

For the convex case, we take $f(w, x) = \langle w, x \rangle$ on $W \times \mathcal{B}$. Then $w^* = \arg\min_{w \in W} \mathbb{P}_X \rho(w)$ and

$$\mathbb{P}_X \rho(w_{\text{priv}}) - \mathbb{P}_X \rho(w^*) \leq \frac{1}{2} \mathbb{E} \{ \langle X, w_{\text{priv}} - w^* \rangle \}^2.$$

Also, the proof of Theorem 21 shows that there exists a dataset $X \in \mathcal{X}^n$ such that $\mathbb{E} \{ \langle X, w_{\text{priv}} - w^* \rangle \} \leq M \{ n : \min \frac{n_{\text{priv}}}{n}, \frac{d}{3n} \}$ and $\mathbb{E} \{ \langle X, w_{\text{priv}} - w^* \rangle^2 \} \leq \min \frac{n_{\text{priv}}}{n}, \frac{d}{n} \}^2$ for any ε -semi-DP A^1 . Note that $A^1 : \frac{M}{n} w_{\text{priv}}$ is ε -semi-DP by post-processing. Thus,

$$\begin{aligned} \mathbb{E} \mathbb{P}_X \rho(w_{\text{priv}}) - \mathbb{P}_X \rho(w^*) &\leq \frac{M}{2n} \mathbb{E} \{ \langle X, w_{\text{priv}} - w^* \rangle \}^2 \\ &\leq \frac{M}{2n} \frac{n}{M} \mathbb{E} \{ \langle X, w_{\text{priv}} - w^* \rangle \}^2 \\ &\leq \frac{n}{M} \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\}^2 \\ &\leq \min \left\{ \frac{n_{\text{priv}}}{n}, \frac{d}{n\varepsilon} \right\}. \end{aligned}$$

A standard scaling argument (see (Bassily et al., 2014) for details) completes the proof.

Upper Bounds: The second terms in each minimum follows by running the ε -DP algorithms in (Bassily et al., 2014): these achieve the desired excess empirical risk bounds and are automatically ε -semi-DP.

We now prove the first term in each respective minimum. Denote $\mathbb{P}_{\text{pub}} \rho(w) = \frac{1}{n} \sum_{x \in X_{\text{pub}}} f(w, x)$ and $\mathbb{P}_{\text{priv}} \rho(w) = \frac{1}{n} \sum_{x \in X_{\text{priv}}} f(w, x)$, so that $\mathbb{P}_X = \mathbb{P}_{\text{pub}} + \mathbb{P}_{\text{priv}}$. The algorithm we will use simply returns any minimizer of the public empirical loss: $w_{\text{pub}} = \arg\min_{w \in W} \mathbb{P}_{\text{pub}} \rho(w)$. (It will be easy to see from the proof that any approximate minimizer would also suffice.) A is clearly ε -semi-DP. Next, we bound the excess risk of A . Let $w^* = \arg\min_{w \in W} \mathbb{P}_X \rho(w)$.

Convex Upper Bound: We have

$$\begin{aligned} \mathbb{P}_X \rho(w_{\text{pub}}) - \mathbb{P}_X \rho(w^*) &= \mathbb{P}_{\text{pub}} \rho(w_{\text{pub}}) - \mathbb{P}_{\text{pub}} \rho(w_{\text{pub}}) + \mathbb{P}_{\text{pub}} \rho(w_{\text{pub}}) - \mathbb{P}_{\text{pub}} \rho(w^*) + \mathbb{P}_{\text{pub}} \rho(w^*) - \mathbb{P}_{\text{priv}} \rho(w^*) \\ &\quad + \mathbb{P}_{\text{priv}} \rho(w^*) - \mathbb{P}_X \rho(w^*) \\ &\leq \frac{1}{n} \sum_{x \in X_{\text{pub}}} f(w_{\text{pub}}, x) - 0 + \frac{1}{n} \sum_{x \in X_{\text{priv}}} f(w_{\text{pub}}, x) \\ &\leq \frac{1}{n} \sum_{x \in X_{\text{priv}}} L \langle w_{\text{pub}} - w^*, x \rangle \\ &\leq L \mathbb{E} \{ \langle w_{\text{pub}} - w^*, X_{\text{priv}} \rangle \} \\ &\leq LD \frac{n_{\text{priv}}}{n}. \end{aligned}$$

Strongly Convex Upper Bound: By the above, $\mathbb{P}_X \rho(w_{\text{pub}}) - \mathbb{P}_X \rho(w^*) \leq L \mathbb{E} \{ \langle w_{\text{pub}} - w^*, X_{\text{priv}} \rangle \}$. Now we will use strong convexity to bound $\mathbb{E} \{ \langle w_{\text{pub}} - w^*, X_{\text{priv}} \rangle \}$. To do so, we use the following lemma, versions of which have appeared, e.g. in (Lowy & Razaviyayn, 2021; Chaudhuri et al., 2011):

Lemma 39. (Lowy & Razaviyayn, 2021) *Let $H(w), h(w)$ be convex functions on some convex closed set $W \subseteq \mathbb{R}^d$ and suppose that $H(w)$ is μ_H -strongly convex. Assume further that h is L_H -Lipschitz. Define $w_1 = \arg\min_{w \in W} H(w)$ and $w_2 = \arg\min_{w \in W} H(w) + h(w)$. Then $\|w_1 - w_2\| \leq \frac{L_H}{\mu_H}$.*

We apply the lemma with $h(w) = \mathbb{P}_{\text{priv}} \rho(w)$ and $H(w) = \mathbb{P}_{\text{pub}} \rho(w)$. Then the conditions of the lemma are satisfied with $L_H = \frac{n_{\text{priv}}}{n} L$ and $\mu_H = \frac{n_{\text{pub}}}{n} \mu$. Thus,

$$\|w_{\text{pub}} - w^*\| \leq \frac{L_H}{\mu_H} \leq \frac{L n_{\text{priv}}}{\mu n_{\text{pub}}}$$

This leads to

$$\mathbb{P}_{X^p w_{pub}^q} \mathbb{P}_{X^p w^q} \propto \frac{L^2 n_{priv}^2}{\mu n n_{pub}}$$

Combining the two strongly convex upper bounds with the upper bound $LD \frac{n_{priv}}{n} \propto \frac{L^2 n_{priv}}{n}$ (which holds for any convex function), we have an algorithm A with the following excess risk:

$$\mathbb{E} \mathbb{P}_{X^p A^p X^q} \mathbb{P}_X \leq \frac{L^2}{\mu} \min \left\{ \frac{n_{priv}}{n}, \frac{n_{priv}^2}{n n_{pub}}, \frac{d^2 \ln p n q}{n^2 \varepsilon^2} \right\}. \quad (37)$$

We will show that (37) is equal to the strongly convex upper bound stated in Theorem 10 up to constant factors. First, suppose $n_{pub} \hat{A} n$: i.e. there is a constant $k \geq 0$ such that $n_{pub} \leq kn$ for all $n \geq 1$. Then, clearly (37) and the strongly convex upper bound stated in Theorem 10 are both equal to $\frac{L^2}{\mu} \min \left\{ \frac{n_{priv}^2}{n^2}, \frac{d^2}{n^2} \ln p n q \right\}$.

Next, suppose $n_{pub} \leq n$: i.e., for any $k \geq 0$, there exists $n \geq 1$ such that $n_{pub} \leq kn$. Then we claim that $\min \left\{ \frac{n_{priv}}{n}, \frac{d^2 \ln p n q}{n^2} \right\} \hat{A} \min \left\{ 1, \frac{d^2 \ln p n q}{n^2} \right\}$. If we prove this claim, then we are done. There are two subcases to consider: A) $n_{priv} \leq \frac{d^2 \ln p n q}{n}$; and B) $n_{priv} > \frac{d^2 \ln p n q}{n}$. In subcase B), the claim is immediate. Consider subcase A): if $n_{priv} \hat{A} n$, then we're done. If not, then we have $n_{priv} \leq n$ and $n_{pub} \leq n$, so $n = n_{priv} = n_{pub} \leq n$, a contradiction. This completes the proof. \square

Remark 40 (Details of Remark 11). *The same minimax risk bound (7) holds up to a logarithmic factor if we replace $F_{0:L;D}$ by the larger class of all Lipschitz non-convex (or convex) loss functions in the definition (6): First, the lower bound in Theorem 10 clearly still holds for non-convex loss functions. For the upper bound, the ε -DP (hence semi-DP) exponential mechanism achieves error $OpLD \frac{d}{n^2} \ln p n q$ (Bassily et al., 2014; Ganesh et al., 2022). Further, the proof of Theorem 10 reveals that convexity is not necessary for the throw-away algorithm to achieve error $OpLD n_{priv} \ln q$. However, the optimal algorithms are inefficient for non-convex loss functions: to the best of our knowledge, all existing polynomial time implementations of the exponential mechanism require convexity for their runtime guarantees to hold. Further, computing $\arg \min_{w \in W} \mathbb{P}_{pub} w^q$ in the implementation of throw-away may not be tractable in polynomial time for non-convex \mathbb{P}_{pub} .*

E.3. Optimal Semi-DP Stochastic Convex Optimization

Approximate $p\varepsilon, \delta q$ -Semi-DP SCO

Theorem 41 (Complete Version of Theorem 12). *Let $\varepsilon \hat{A} 1 \{ \log p n d q$ and $\delta \leq 1 \{ n$. Then, there is a constant $C \geq 0$ such that*

$$\ell_{p,d,n} LD \min \left\{ \frac{1}{n_{pub}}, \frac{d}{n\varepsilon} \right\} \frac{1}{n} \propto R_{SCO}(p\varepsilon, \delta, n_{priv}, n, d, L, D, \mu) \quad \text{or} \quad CLD \min \left\{ \frac{1}{n_{pub}}, \frac{d \ln p 1 \{ \delta q}{n\varepsilon} \right\} \frac{1}{n},$$

and

$$\ell_{p,d,n} LD \min \left\{ \frac{1}{n_{pub}}, \frac{d}{n\varepsilon} \right\} \frac{1}{n^2} \propto R_{SCO}(p\varepsilon, \delta, n_{priv}, n, d, L, D, \mu) \propto C \frac{L^2}{\mu} \min \left\{ \frac{1}{n_{pub}}, \frac{d \ln p 1 \{ \delta q}{n\varepsilon} \right\} \frac{1}{n^2},$$

where $1 \{ \ell_{p,d,n}$ is logarithmic in d and n . Our lower bounds hold for symmetric $A = pA^1, \dots, A^d q$.

Proof. Lower bounds: Let A be $p\varepsilon, \delta q$ -semi-DP and symmetric, and denote $w_{priv} = ApXq$.

Strongly convex lower bounds: We begin with the strongly convex lower bounds, which can be proved by reducing strongly convex SCO to mean estimation and applying Theorem 4. In a bit more detail, let $f : W \times \mathbb{N} \rightarrow \mathbb{R}$ be given by

$$f(w, x) = \frac{L}{2D} \|w - x\|^2,$$

where $W \times DB$. Note that $f(w, x)$ is L -uniformly Lipschitz and $\frac{L}{D}$ -strongly convex in w for all x . Further, $w = \arg \min_{w \in W} \mathbb{E}_X f(w, x) = \mathbb{E}_X x$. By a direct calculation (see e.g. (Kamath et al., 2022a, Lemma 6.2)), we have

$$\mathbb{E} F(w_{priv}) = F(w) + \frac{L}{2D} \mathbb{E} \|w_{priv} - w\|^2. \quad (38)$$

We can lower bound $\mathbb{E}\{w_{\text{priv}} - w\}^2 = \mathbb{E}\{A_{\text{priv}} - \mathbb{E}_X \{P_{\text{priv}}\}\}^2$ via Theorem 4 (and its proof, to account for the re-scaling). Specifically,

$$\mathbb{E}\{w_{\text{priv}} - w\}^2 \leq \ell p d, n q D^2 \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n^2 \varepsilon^2}, \frac{1}{n} \right\},$$

for a logarithmic function $\ell p d, n q$ of d and n , by Theorem 4. Applying (38) yields the desired excess risk lower bound for $\delta \geq 0$.

Convex lower bounds: We will begin by proving the lower bounds for the case in which $L = D = 1$, and then scale our construction to get the lower bounds for arbitrary L, D .

Let $X = \left\{ \frac{1}{d} \right\}^d \in \mathbb{R}^d$ and $W = \mathbb{B}$. Define

$$f(w, x) = \langle w, x \rangle,$$

which is convex and 1-uniformly-Lipschitz in w on X . Let P be the hard distribution used to prove Theorem 23, which satisfies $\mathbb{E}_X \{P_{\text{priv}}\} = \theta P_{\text{priv}} = \frac{1}{n} \left\{ \frac{1}{n_{\text{pub}}}, \frac{1}{n_{\text{priv}}} \right\}$. Further, $w = \arg \min_{w \in W} \mathbb{E} \{f(w, x)\} = \frac{1}{n} \left\{ \frac{1}{n_{\text{pub}}}, \frac{1}{n_{\text{priv}}} \right\}$. A direct calculation (see e.g. (Kamath et al., 2022a, Equation 14)) shows

$$\sup \mathbb{E} \{F_{\text{priv}}\} - F \leq \frac{1}{2} \mathbb{E} \{ \theta \} \{w_{\text{priv}} - w\}^2 \quad (39)$$

$$\frac{1}{2} \sup \mathbb{E} \left\{ \frac{1}{\theta} \right\} \{w_{\text{priv}} - \theta\}^2, \quad (40)$$

where $w_{\text{priv}} = A_{\text{priv}}(X)$ is the output of the algorithm $A : X^n \rightarrow W$ defined by $A_{\text{priv}}(X) = \theta A_{\text{priv}}(X)$. Note that A is $\rho \varepsilon, \delta q$ -semi-DP by post-processing, for any θ . Now, we invoke Theorem 23 to obtain

$$\sup \mathbb{E} \{w_{\text{priv}} - \theta\}^2 \leq \frac{1}{n} \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}}^2 \varepsilon^2}, \frac{1}{n} \right\}$$

for $\theta = \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{1}{n_{\text{priv}}} \right\}$. This implies the desired lower bound when $L = D = 1$.

For general L and D , we scale the problem instance as follows: let $\tilde{W} = DW$, $\tilde{X} = LX$, and $\tilde{f}(w, x) = Lx$ for $x \in P$. Define $\tilde{f} : \tilde{W} \times \tilde{X} \rightarrow \mathbb{R}$ by $\tilde{f}(w, x) = f(w, x)$. Then \tilde{f} is L -Lipschitz and convex. Moreover, if $F_{\text{priv}} = \mathbb{E}_X \{P_{\text{priv}}\}$, $\tilde{F}_{\text{priv}} = \mathbb{E}_{\tilde{X}} \{P_{\text{priv}}\}$, $w = Dw$, and $\theta = \mathbb{E}_{\tilde{X}} \{P_{\text{priv}}\} = L\theta$, then $Dw = P \arg \min_{w \in \tilde{W}} \tilde{F}_{\text{priv}}(w)$ and

$$\begin{aligned} \tilde{F}_{\text{priv}}(w) &= \mathbb{E} \{ \langle w, \tilde{\theta} \rangle \} = \mathbb{E} \{ \langle w, \theta \rangle \} \\ &= D \mathbb{E} \{ \langle w, \theta \rangle \} \\ &= LD \mathbb{E} \{ \langle w, \theta \rangle \} \\ &= LD \mathbb{E} \{ f(w, \theta) \} = F. \end{aligned}$$

This shows that excess risk scales by LD , completing the lower bound proofs.

Upper bounds: Convex upper bounds: Consider the 0-semi-DP throw-away algorithm that discards X_{priv} and runs n_{priv} steps of one-pass SGD (stochastic approximation) using X_{pub} . This algorithm has excess risk $O \left(\frac{L^2}{n_{\text{pub}}} \right)$ (Nemirovski & Yudin, 1983). To obtain the second term in the convex $\rho \varepsilon, \delta q$ -semi-DP upper bound, one can use, e.g. $\rho \varepsilon, \delta q$ -DP-SGD (Bassily et al., 2019).

Strongly convex upper bounds: Consider the 0-semi-DP throw-away algorithm that discards X_{priv} and runs n_{priv} steps of one-pass SGD (stochastic approximation) using X_{pub} . This algorithm has excess risk $O \left(\frac{L^2}{n_{\text{pub}}} \right)$ (Nemirovski & Yudin, 1983). The second term in the strongly convex $\rho \varepsilon, \delta q$ -semi-DP upper bound can be attained, e.g. by $\rho \varepsilon, \delta q$ -DP-SGD (Lowy & Razaviyayn, 2023b). \square

Next, we provide minimax optimal excess risk bounds for pure ε -semi-DP SCO.

Pure ε -Semi-DP SCO

Theorem 42 (Pure ε -Semi-DP SCO). *Suppose $\varepsilon \leq d/8$, and either $n_{\text{pub}} \geq \frac{n}{d}$ or $d \geq 1$. If $\mu = 0$ (convex case), then there exist absolute constants $0 < c_0 \leq C_0$ such that*

$$c_0 LD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon} \right\} \leq R_{\text{SCO}}(\varepsilon, \delta, 0, n_{\text{priv}}, n, d, L, D, \mu) \leq C_0 LD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon} \right\}.$$

If $\mu > 0$, there are constants $0 < c_1 \leq C_1$ such that

$$c_1 LD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2}{n^2\varepsilon^2} \right\} \leq R_{\text{SCO}}(\varepsilon, \delta, 0, n_{\text{priv}}, n, d, L, D, \mu) \leq C_1 \frac{L^2}{\mu} \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2 \ln pn}{n^2\varepsilon^2} \right\}.$$

The above upper bounds hold for any n_{pub}, d .

Proof of Theorem 42. Lower bounds: Let A be ε -semi-DP and denote $w_{\text{priv}} = \text{ApXq}$.

Strongly convex lower bound: We begin with the strongly convex lower bound, which can be proved by reducing strongly convex SCO to mean estimation and applying Theorem 32. In a bit more detail, let $f : W \times \mathbb{N} \rightarrow \mathbb{R}$ be given by

$$f(w, x) = \frac{L}{2D} \|w - x\|^2,$$

where $W \times \mathbb{N} = DB$. Note that f is L -uniformly Lipschitz and $\frac{L}{D}$ -strongly convex in w for all x . Further, $w^* = \text{argmin}_{w \in W} \mathbb{E}_x [f(w, x)] = \mathbb{E}_x [x]$. By a direct calculation (see e.g. (Kamath et al., 2022a, Lemma 6.2)), we have

$$\mathbb{E} [f(w_{\text{priv}})] - \mathbb{E} [f(w^*)] \leq \frac{L}{2D} \mathbb{E} \|w_{\text{priv}} - w^*\|^2. \quad (41)$$

We can lower bound $\mathbb{E} \|w_{\text{priv}} - w^*\|^2 = \mathbb{E} \|\text{ApXq} - \mathbb{E}_x [x]\|^2$ via Theorem 32 (and its proof, to account for the re-scaling). Specifically, if $\delta = 0$, $\varepsilon \leq \max\{1, d/8\}$, and either $d \geq 0.1$ or $n_{\text{pub}} \geq 0.1n\varepsilon/d$, then Theorem 32 and its proof imply

$$\mathbb{E} \|w_{\text{priv}} - w^*\|^2 \geq cD^2 \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d^2}{n^2\varepsilon^2} \right\} \frac{1}{n}.$$

Combining this with (41) leads to the desired excess risk lower bound for $\delta = 0$.

Convex lower bound: We will begin by proving the lower bound for the case in which $L = D = 1$, and then scale our construction to get the lower bounds for arbitrary L, D .

Assume $\delta = 0$, $\varepsilon \leq d/8$, and either $n_{\text{pub}} \geq n\varepsilon/d$ or $d \geq 1$. Let $X = \{0, \dots, \frac{1}{d}\}^d \in \mathbb{R}^d$ and $W = B$. Define

$$f(w, x) = \|x - w\|,$$

which is convex and 1-uniformly-Lipschitz in w on X . Choose V to be a finite subset of \mathbb{R}^d such that $|V| \geq 2^{d/2}$, $\|v\| = 1$ for all v , and $\|v - v^1\| \geq 1/8$ for all $v \in V$ (see e.g. the Gilbert-Varshamov construction). Following the proof of Theorem 4, we define $P_v = (1-p)P_0 + pP_v$ for all $v \in V$, where $p \in (0, 1)$ will be chosen later, P_0 is point mass on $\{0\}$ and P_v is point mass on $\{v\}$. Denote the mean $\theta_v = \mathbb{E}_x [x] = pv$. Note that $\|\theta_v - \theta_{v^1}\| \geq p$ for all v . Let $F_v(w) = \mathbb{E}_x [f(w, x)]$ and $w_v = \text{argmin}_{w \in W} F_v(w) = \frac{v}{|V|}$. A direct calculation (see e.g. (Kamath et al., 2022a, Equation 14)) shows

$$\mathbb{E} [F_v(w)] - F_v(w_v) \leq \frac{1}{2} \mathbb{E} \|\theta_v - w_v\|^2 \quad (42)$$

for any $w \in W, v \in V$. Also,

$$\rho(p, V) = \min_{w \in W} \max_{v \in V} \|w - w_v\| : v \in V, \|v - v^1\| \geq 1/8 \leq \min_{v \in V} \|v - v^1\| : v \in V, \|v - v^1\| \geq 1/8.$$

Thus, by combining (42) with the reduction from estimation to testing and Theorem 33 (see the proof of Theorem 4 for details), we have

$$\begin{aligned}
 \sup_{\mathcal{P}^V} \mathbb{E} \left[r_{F_V} \rho w_{\text{priv}} \right] &\leq \frac{1}{2} \sup_{\mathcal{P}^V} \mathbb{E} \left[\|\theta_V\| w_{\text{priv}} \right]^2 \\
 &\leq \frac{p}{2} \sup_{\mathcal{P}^V} \mathbb{E} \left[w_{\text{priv}} \right]^2 \\
 &\leq \frac{p}{2} \rho \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E} \left[\|\theta_v\| \right]^2 \\
 &\leq \frac{p}{128} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{1}{\rho} \mathbb{E} \left[\|\theta_v\| \right]^2 \\
 &\leq \frac{p}{128} \frac{2^{d/2}}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{1}{\rho} \mathbb{E} \left[\|\theta_v\| \right]^2 \\
 &\leq \frac{p}{512} \frac{1}{\rho} \sum_{v \in \mathcal{V}} \mathbb{E} \left[\|\theta_v\| \right]^2 \min \left(1, \frac{2^{d/2}}{e^{\rho n_{\text{priv}}}} \right).
 \end{aligned}$$

Now, assume $d \geq 4$ so that $2^{d/2} \geq e^{d/4}$. Then, as detailed in the proof of Theorem 4, choosing

$$\rho \leq \min \left\{ \frac{d}{4n\varepsilon}, \frac{1}{n}, \frac{1}{2n_{\text{pub}}} \right\}$$

and assuming $n_{\text{pub}} \geq kn\varepsilon$ for some absolute constant k implies

$$\sup_{\mathcal{P}^V} \mathbb{E} \left[r_{F_V} \rho w_{\text{priv}} \right] \leq c \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon} \right\}$$

for some absolute constant $c \geq 0$. Combining this with the non-private SCO lower bound (Nemirovski & Yudin, 1983) yields

$$\sup_{\mathcal{P}} \mathbb{E} \left[r_{F_V} \rho w_{\text{priv}} \right] \geq c' \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon}, \frac{1}{n} \right\},$$

where $r_{F_V} \rho w_{\text{priv}} : \mathbb{E}_X \rho f_{\text{priv}}(w, x)$.

Suppose instead that $0 \leq \delta \leq \varepsilon$ and $d \geq 1$ (i.e. $d \geq k$ for some absolute constant $k \geq 1$), but $n_{\text{pub}} \geq n\varepsilon$ is arbitrary. We will prove the lower bound for $d \geq 1$; by taking the k -fold product distribution, this is sufficient to complete the proof of the unscaled ε -semi-DP lower bound. Define distributions P_1, P_2 on \mathbb{R}^d as follows:

$$P_1 \text{ is } \frac{1}{2} \delta_{\theta_1} + \frac{1}{2} \delta_{\theta_2}, \quad P_2 \text{ is } \frac{1}{2} \delta_{\theta_1} + \frac{1}{2} \delta_{-\theta_2}$$

for some $\gamma \in (0, 1/2]$ to be chosen later. Note $\theta_1 : \mathbb{E}_{P_1} \theta = \gamma$ and $\theta_2 : \mathbb{E}_{P_2} \theta = \gamma$, so $|\theta_j| \leq \gamma$ for $j = 1, 2$. Let $F_j \rho w_{\text{priv}} : \mathbb{E}_X \rho f_{\text{priv}}(w, x)$ and $w_j = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \rho \arg \min_{w \in \mathcal{W}} F_j \rho w_{\text{priv}}$. Then by (42), we have

$$\begin{aligned}
 \max_{j \in \{1, 2\}} \mathbb{E} \left[F_j \rho w_{\text{priv}} \right] &\geq \frac{1}{2} \max_{j \in \{1, 2\}} \mathbb{E} \left[\|\theta_j\| w_{\text{priv}} \right]^2 \\
 &\geq \frac{\gamma}{2} \max_{j \in \{1, 2\}} \mathbb{E} \left[w_{\text{priv}} \right]^2 \\
 &\geq \frac{1}{2\gamma} \max_{j \in \{1, 2\}} \mathbb{E} \left[w_{\text{priv}}^1 \right]^2,
 \end{aligned}$$

where $w_{\text{priv}}^1 : \gamma w_{\text{priv}}$ is semi-DP iff w_{priv} is semi-DP (by post-processing). Thus, by applying Le Cam's method and Lemma 35 (see the proof of Theorem 4 for details), we get

$$\max_{j \in \{1, 2\}} \mathbb{E} \left[F_j \rho w_{\text{priv}} \right] \geq \frac{1}{2\gamma} \frac{\gamma^2}{8} \geq \frac{1}{8} \gamma \min \left\{ \frac{c}{3n}, \frac{c}{6n_{\text{priv}}\varepsilon}, \frac{c}{3n_{\text{pub}}} \right\}.$$

Now, we will choose γ to (approximately) maximize the right-hand side of the above inequality. If $\min \left\{ \frac{3n}{2}, 6n_{\text{priv}}\varepsilon \right\} \geq \frac{3n_{\text{pub}}}{2}$, then choosing $\gamma = \frac{1}{3} \frac{2}{3n}$ yields

$$\max_{j \in \text{Pt}1:2u} \mathbb{E} \|F_j(\rho w_{\text{priv}}) - F_j\| \leq \frac{k}{n}$$

for some absolute constant $k \geq 0$. If instead $\min \left\{ \frac{3n}{2}, 6n_{\text{priv}}\varepsilon \right\} < \frac{3n_{\text{pub}}}{2}$, then we choose $\gamma = \frac{2}{3} \frac{6n_{\text{priv}}\varepsilon}{\frac{3n_{\text{pub}}}{2}}$. This choice implies

$$\max_{j \in \text{Pt}1:2u} \mathbb{E} \|F_j(\rho w_{\text{priv}}) - F_j\| \leq k^1 \min \left\{ \frac{1}{n_{\text{priv}}\varepsilon}, \frac{1}{n_{\text{pub}}} \right\}$$

for some absolute constant $k^1 \geq 0$. Combining the pieces above with the non-private SCO lower bound (Nemirovski & Yudin, 1983) yields

$$\sup_P \mathbb{E} \|rF(\rho w_{\text{priv}}) - F\| \leq c \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{1}{n\varepsilon}, \frac{1}{n} \right\},$$

where $rF(\rho w) = \mathbb{E}_x \rho r f(\rho w, x)$.

A standard scaling argument completes the lower bound proofs (see e.g. the proof of Theorem 41 for details).

Upper bounds: Convex upper bounds: Consider the 0-semi-DP throw-away algorithm that discards X_{priv} and runs n_{priv} steps of one-pass SGD (stochastic approximation) using X_{pub} . This algorithm has excess risk $O \left(\frac{L^2}{n_{\text{pub}}} \right)$ (Nemirovski & Yudin, 1983). To obtain the second term in the convex ε -semi-DP upper bound, use the ε -DP (hence semi-DP) regularized exponential mechanism of Ganesh et al. (2022).

Strongly convex upper bounds: Consider the 0-semi-DP throw-away algorithm that discards X_{priv} and runs n_{priv} steps of one-pass SGD (stochastic approximation) using X_{pub} . This algorithm has excess risk $O \left(\frac{L^2}{n_{\text{pub}}} \right)$ (Nemirovski & Yudin, 1983). To obtain the second term in the strongly convex ε -semi-DP upper bound, one can use, e.g. the ε -DP (hence semi-DP) iterated exponential mechanism of Ganesh et al. (2022). \square

E.3.1. SEMI-DP SCO WITH AN ‘‘EVEN MORE OPTIMAL’’ GRADIENT ESTIMATOR

Proposition 43. *We provide privacy guarantees for Algorithm 1:*

1. Suppose we sample with replacement in line 4 of Algorithm 1. Then, there exist constants c_1, c_2 such that for any $\varepsilon \leq c_1 \frac{K_{\text{priv}}}{n_{\text{priv}}}^2 T$, Algorithm 1 is $\rho\varepsilon, \delta\varepsilon$ -semi-DP for any $\delta \geq 0$ if we choose $\sigma^2 \leq c_2 \frac{C^2 \ln p \{ \frac{qT}{\sigma^2 n_{\text{priv}}^2} \}}$.
2. Suppose we sample without replacement in line 4 and choose $T \leq \frac{n}{K_{\text{priv}}}$. Then Algorithm 1 is ρ -semi-zCDP if $\sigma^2 \leq \frac{2C^2}{K_{\text{priv}}^2}$.

Proof. Note that the ℓ_2 -sensitivity of the private stochastic gradient query is

$$\sup_{X_{\text{priv}}} \sup_{X_{\text{priv}}^1} \frac{\alpha}{K_{\text{priv}}} \cdot \text{clip}_{C\rho r} f(\rho w_t, x) \leq \frac{\alpha}{K_{\text{priv}}} \cdot \text{clip}_{C\rho r} f(\rho w_t, x^1) \leq \frac{2\alpha C}{K_{\text{priv}}}.$$

1. Consider sampling with replacement. Then we are randomly subsampling from the private data uniformly with sampling ratio K_{priv}/a . Thus, the theorem follows from (Abadi et al., 2016, Theorem 1).

2. Consider sampling without replacement. Then by the ρ -zCDP guarantee of the Gaussian mechanism (Bun & Steinke, 2016, Proposition 1.6) and the sensitivity bound above, \mathfrak{g}_t is ρ -semi-zCDP for every t . Moreover, since we are sampling without replacement, the privacy of every $x \in X_{\text{priv}}$ is only affected by \mathfrak{g}_t for a single $t \in [T]$. Thus, semi-zCDP of Algorithm 1 follows by parallel composition (McSherry, 2009). \square

Excess risk of By Proposition 8, there exists a choice of α such that the variance of our unbiased estimator in line 7 is always less than the variance of both the throw-away gradient estimator $\frac{1}{K_{\text{pub}}} \sum_{x \in \mathcal{X}} \nabla f(w_t, x)$ and the DP-SGD estimator $\frac{1}{K_{\text{priv}} K_{\text{pub}}} \sum_{x \in \mathcal{X}} \nabla f(w_t, x) + u_t$, where u_t is appropriately scaled (to ensure DP) Gaussian noise. Consequently, if we choose T and $K = K_{\text{priv}} K_{\text{pub}}$ such that $n = TK$ and sample without replacement (i.e. one pass), then Algorithm 1 always has smaller excess risk than both the throw-away SCO algorithm and one-pass DP-SGD. Moreover, if the loss function has Lipschitz continuous gradient, then one can combine the stochastic gradient estimator of line 7 with acceleration (Ghadimi & Lan, 2012) to obtain a linear-time semi-DP algorithm that always outperforms the accelerated DP algorithm of Lowy & Razaviyayn (2023b). This is because the variance of our gradient estimator (hence our excess risk) is strictly smaller than that of Lowy & Razaviyayn (2023b), by Theorem 8. For example, for β -smooth, μ -strongly convex loss functions, one-pass Algorithm 1 achieves excess risk that is optimal up to a factor of $O(\beta/\mu)$ and improves over (Lowy & Razaviyayn, 2023b). Moreover, (Lowy & Razaviyayn, 2023b) has the smallest excess risk among linear-time (one-pass) DP algorithms whose privacy analysis does not require convexity. Thus, our algorithm can be used for deep learning. In our numerical experiments, we implement the with-replacement sampling version of Algorithm 1.

We also note that near-optimal excess risk bounds for non-convex loss functions that satisfy the (Proximal) PL inequality (Polyak, 1963; Karimi et al., 2016) can be derived by combining a proximal variation of Algorithm 1 with the techniques of Lowy et al. (2023a). Further, if $\nabla f(w, x)$ is not uniformly Lipschitz, but has stochastic gradients with bounded k -th order moment for some $k \geq 2$, then excess risk bounds can still be derived for Algorithm 1 via techniques in (Lowy & Razaviyayn, 2023a). Our algorithm can also be extended to a variation of noisy stochastic gradient descent ascent, which could be used, e.g. for fair semi-DP model training (Lowy et al., 2023b). We leave it as future work to explore these and other potential applications of our gradient estimator in efficiently training private ML models with public data.

F. Optimal Locally Private Model Training with Public Data

Notation and Setup: Following Duchi & Rogers (2019), we permit algorithms to be *fully interactive*. That is, algorithms may adaptively query the same individual i multiple times over the course of T “communication rounds.” We denote i ’s message in round t by $Z_{i,t} \in \mathcal{Z}$. Person i ’s message $Z_{i,t} \in \mathcal{Z}$ in round t may depend on all previous communications $B^{pt} = \{Z_{i,r}, B^{pr}\}_{r=1}^{t-1}$ and on i ’s own data: $Z_{i,t} = Q_{i,t}(x_i, Z_{i,t}, B^{pt})$. If i ’s data is private, then $Z_{i,t}$ is a randomized view of x_i distributed (conditionally) according to $Q_{i,t}$. If i ’s data is public, then $Z_{i,t}$ may be deterministic. Full interactivity is the most general notion of interactivity. If $T = 1$, then we say the algorithm is *sequentially interactive*. If, in addition, each person’s message $Z_{i,1}$ depends only on x_i and not on $x_{j \neq i}$, then we say the algorithm is *non-interactive*. Semi-LDP (Definition 13) essentially requires that the messages $\{Z_{i,t}\}_{t=1}^T$ be DP for all private $x_i \in \mathcal{X}_{\text{priv}}$.

F.1. Optimal Semi-LDP Mean Estimation

Theorem 44 (Re-statement of Theorem 14). *Let $\epsilon \in (0, 1]$. There are absolute constants $0 < c \leq C$ s.t.*

$$c \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\epsilon^2} \right\} \leq \mathcal{N}_{\text{pop}}^{\text{loc}}(\epsilon, n_{\text{priv}}, n, d) \leq C \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\epsilon^2} \right\}.$$

Proof. Lower bound: We will actually prove a more general lower bound than the one in Theorem 14; namely, we will show a lower bound on the minimax ℓ_1 -error for estimation of distributions on $\mathcal{X}_r = \{ru^d \mid r \in [0, 1]\}$. To that end, let $\gamma \in (0, 1]$ and

$$P_1 : \begin{cases} r & \text{with probability } \frac{1}{2} \\ r & \text{with probability } \frac{1}{2} \end{cases}$$

and

$$P_{-1} : \begin{cases} r & \text{with probability } \frac{1}{2} \\ r & \text{with probability } \frac{1}{2} \end{cases}.$$

We define our hard distribution on \mathcal{X}_r by first drawing $V \sim \text{Unif}([0, 1])$ and then—conditional on $V = v$ —drawing $X_{i,j} \sim P_v = \prod_{j=1}^d P_{v_j}$ for $i \in [n_s], j \in [d]$, where P_v denotes the product distribution. We have Markov chains $V_j \stackrel{\text{d}}{\sim} X_{i,j} \stackrel{\text{d}}{\sim} Z$ for all $j \in [d], i \in [n_s]$, where Z is the semi-LDP transcript. Note that $\ln \frac{dP_1}{dP_{-1}} \leq \ln \frac{1}{1-b}$ and $e^b \leq 3$ for any $\gamma \in (0, 1/2]$. Now we will use the following lemma from Duchi & Rogers (2019):

Lemma 45. (Duchi & Rogers, 2019, Lemma 24) Let $V \tilde{N} X \tilde{N} Z$ be a Markov chain, where $X \sim P_V$ conditional on $V = v$. If $\ln \frac{dP_v}{dP_{v^1}} \leq \alpha$ for all v, v^1 , then

$$I_p(V; Z) \leq 2pe^{-\alpha} + \gamma^2 I_p(X; Z).$$

Thus, for $V_j \sim \text{Unif}(\mathcal{T})$, Lemma 45 implies $I_p(V_j; Z) \leq \frac{8}{\rho^2} I_p(X_{1:j}; Z)$. Hence the strong data processing constant (Duchi & Rogers, 2019, Definition 9) is $\beta = \beta(P_1, P_{1^q}) \leq \frac{8}{\rho^2}$.

Now, $\theta_{V_j} : \mathbb{E}_X P_{V_j}(r x_S) = \gamma r v_j$ for any $v_j \in \mathcal{T}$. Moreover, letting $\theta_{V_j} = \rho \theta_{v_j}$, $\theta_{v_j} \in \mathbb{R}^d$ for $v_j \in \mathcal{T}$ and $\theta \in \mathbb{R}^d$, we have

$$\|\theta - \theta_{V_j}\|_1 \leq \sum_{j=1}^d |\theta_j - \gamma r v_j| \leq r \gamma \sum_{j=1}^d \frac{\theta_j}{r \gamma} v_j \leq r \gamma \sum_{j=1}^d \mathbb{1}_{\{\text{sign}(\theta_j) \neq v_j\}}.$$

Thus, $\mathcal{T} = \mathcal{U}^d$ induces an $r\gamma$ -Hamming separation, so Assouad's lemma (Duchi et al., 2018, Lemma 1) yields

$$\inf_{A \in \mathcal{A}^{\text{loc}}} \sup_{P \in \mathcal{P}_r} \mathbb{E} \{A_p(X) - \theta_p(P)\}_1 \leq r \gamma \sum_{j=1}^d \inf_{\hat{V}} \mathbb{P}(\hat{V}_j \neq P_j),$$

where Z is the communication transcript of A , the infimum on the RHS is over all estimators of V , $\theta_p(P) = \mathbb{E}_X P(r x_S)$, and \mathcal{P}_r is the set of distributions on X_r .

Assume WLOG that the private samples are the first n_{priv} samples of X : $X_{\text{priv}} = (x_1, \dots, x_{n_{\text{priv}}})$. To lower bound $\sum_{j=1}^d \inf_{\hat{V}} \mathbb{P}(\hat{V}_j \neq P_j)$, we use a slight extension of Duchi & Rogers (2019, Theorem 10):

$$\sum_{j=1}^d \inf_{\hat{V}} \mathbb{P}(\hat{V}_j \neq P_j) \geq \frac{d}{2} \left(1 - \frac{C}{d} \frac{7pe^b + 1}{\beta} \frac{I_p(X_{\text{priv}}; Z|V) + I_p(X_{\text{pub}}; Z|V)}{\gamma} \right).$$

This follows since $V \tilde{N} X_{\text{priv}} \tilde{N} Z$ and $V \tilde{N} X_{\text{pub}} \tilde{N} Z$ are both Markov chains and the other assumptions in (Duchi & Rogers, 2019, Theorem 10) all hold. Combining this bound with Assouad's lemma (Duchi et al., 2018, Lemma 1) and substituting the definitions of b and β given above gives us

$$\inf_{A \in \mathcal{A}^{\text{loc}}} \sup_{P \in \mathcal{P}_r} \mathbb{E} \{A_p(X) - \theta_p(P)\}_1 \leq \frac{r\gamma d}{2} \left(1 + \frac{C}{d} \frac{896}{\gamma^2} \frac{I_p(X_{\text{priv}}; Z|V) + I_p(X_{\text{pub}}; Z|V)}{\gamma} \right)$$

for any $\gamma \in (0, 1/2]$. It remains to upper bound the conditional mutual information $I_p(X_{\text{priv}}; Z|V)$ and $I_p(X_{\text{pub}}; Z|V)$.

Now for any ε -semi-LDP algorithm with communication transcript Z , we have $I_p(X_{\text{priv}}; Z|V) \leq n_{\text{priv}} \min\{\varepsilon, 4\varepsilon^2\}$, by an easy extension of Duchi & Rogers (2019, Lemma 12). Also, $I_p(X_{\text{pub}}; Z|V) \leq H_p(X_{\text{pub}}|V) \leq \log \rho |X_{\text{pub}}| = d n_{\text{pub}}$, where $H_p(\cdot|q)$ denotes conditional entropy. Thus,

$$\inf_{A \in \mathcal{A}^{\text{loc}}} \sup_{P \in \mathcal{P}_r} \mathbb{E} \{A_p(X) - \theta_p(P)\}_1 \leq \frac{r\gamma d}{2} \left(1 + \frac{C}{d} \frac{4000}{\gamma^2} \frac{\rho n_{\text{priv}} \min\{\varepsilon, \varepsilon^2\} + d n_{\text{pub}}}{\gamma} \right).$$

Choosing $\gamma^2 = c \min\left\{\frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}} \min\{\varepsilon, \varepsilon^2\}}\right\}$ for some small constant $c > 0$ yields

$$\inf_{A \in \mathcal{A}^{\text{loc}}} \sup_{P \in \mathcal{P}_r} \mathbb{E} \{A_p(X) - \theta_p(P)\}_1 \leq r d \min\left\{\frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}} \min\{\varepsilon, \varepsilon^2\}}\right\},$$

whence

$$\inf_{A \in \mathcal{A}^{\text{loc}}} \sup_{P \in \mathcal{P}_r} \mathbb{E} \{A_p(X) - \theta_p(P)\}_2 \leq r d \min\left\{\frac{1}{n_{\text{pub}}}, \frac{d}{n_{\text{priv}} \min\{\varepsilon, \varepsilon^2\}}\right\}.$$

By applying the non-private mean estimation lower bound, we get

$$\inf_{A_{PA}^{loc}} \sup_{P_r} \mathbb{E} \{ \mathbb{A}_{P_r} \mathbb{X} \} \geq \frac{1}{n_{pub}} \min \left\{ \frac{d}{n_{priv} \min\{\varepsilon, \varepsilon^2\}}, \frac{1}{n} \right\}.$$

Choosing $r = \frac{1}{\sqrt{d}}$ ensures that $P_r \in \mathcal{P}_{PB}$ and yields

$$\inf_{A_{PA}^{loc}} \sup_{P_r \in \mathcal{P}_{PB}} \mathbb{E} \{ \mathbb{A}_{P_r} \mathbb{X} \} \geq \frac{1}{n_{pub}} \min \left\{ \frac{1}{n_{priv} \min\{\varepsilon, \varepsilon^2\}}, \frac{1}{n} \right\}.$$

Applying Jensen’s inequality completes the proof of the lower bound in Theorem 14. Since we assumed $\varepsilon \propto 1 \propto d$, the minimum in the denominator simplifies to ε^2 and the $\frac{1}{n}$ term is non-dominant.

Upper bound: The first term in the minimum can be realized by the algorithm that throws away the private data and returns $\frac{1}{n_{pub}} \sum_{x \in X_{pub}} x$, which is 0-semi-LDP. Also,

$$\mathbb{E} \{ \mathbb{A}_{P_r} \mathbb{X} \} - \mathbb{E} \{ x \}^2 \leq \frac{1}{n_{pub}^2} \sum_{x \in X_{pub}} x^2 - \mathbb{E} \{ x \}^2 \leq \frac{1}{n_{pub}}.$$

The second term in the upper bound can be realized by the ε -LDP (hence ε -semi-LDP) estimator of [Duchi et al. \(2013\)](#), which has worst-case MSE upper bounded by $O\left(\frac{d}{n^2}\right)$. \square

Algorithm 3 PrivUnit $_{pp, \gamma}$ ([Bhowmick et al., 2018](#))

- 1: **Input:** $v \in \mathbb{R}^{d-1}, \gamma \in \mathbb{R}, 1s, p \in \mathbb{R}, 1s, B_p; \cdot, \cdot, q$ below is the incomplete Beta function $B_p(x; a, b) = \int_0^x t^{a-1} (1-t)^{b-1} dt$ and $B_p(a, b) = B_p(1; a, b)$.
- 2: Draw $z \sim \text{Bernoulli}(p)$
- 3: **if** $z = 1$ **then**
- 4: Draw $V \sim \text{Unif}(u \in \mathbb{R}^{d-1} : \|u\| \leq \gamma)$
- 5: **else**
- 6: Draw $V \sim \text{Unif}(u \in \mathbb{R}^{d-1} : \|u\| \leq \gamma)$
- 7: **end if**
- 8: Set $\alpha = \frac{d-1}{2}$ and $\tau = \frac{1}{2}$
- 9: Calculate normalization constant

$$m = \frac{p^{1-\gamma^2} q}{2^{d-2} p d^{1-q}} \frac{p}{B_p(\alpha, \alpha)} \frac{1-p}{B_p(\tau; \alpha, \alpha)}$$

- 10: Return $\frac{1}{m} V$

F.1.1. AN “EVEN MORE OPTIMAL” SEMI-LDP ESTIMATOR

Lemma 46 (Re-statement of Lemma 16). *Let P be a distribution on \mathbb{B} with $V^2 = \mathbb{E} \{ \|x\|^2 \} - \mathbb{E} \{ \|x\| \}^2$. Let $c \geq 0$ such that $\mathbb{E} \{ \|x\| \} \leq \frac{cd}{n^2}$, so that $\mathbb{E} \{ \|x\| \} \leq \frac{cd}{n^2} \leq \frac{V^2}{n}$. Then,*

$$\mathbb{E} \{ \mathbb{A}_{Semi-Duchi} \mathbb{X} \} - \mathbb{E} \{ \|x\| \}^2 \leq \frac{n_{priv}}{n} \frac{cd}{n\varepsilon^2} + \frac{n_{pub}}{n} \frac{V^2}{n}.$$

Proof. We have

$$\mathbb{E} \{ \mathbb{A}_{Semi-Duchi} \mathbb{X} \} - \mathbb{E} \{ \|x\| \}^2 \leq \frac{1}{n^2} \sum_{x \in X_{priv}} \|x\|^2 - \mathbb{E} \{ \|x\| \}^2 + \frac{1}{n^2} \sum_{x \in X_{pub}} \|x\|^2 - \mathbb{E} \{ \|x\| \}^2 \leq \frac{ncd}{\varepsilon^2 n^2} + \frac{n_{pub} V^2}{n^2},$$

by independence of the data and the assumptions in the statement of the lemma. \square

F.1.2. A SEMI-LDP ESTIMATOR WITH OPTIMAL CONSTANTS

Proposition 47 (Re-statement of Proposition 17). *Let $A_{\text{priv}}: S^{d-1} \rightarrow \mathbb{R}^n$ and $A_{\text{pub}}: S^{d-1} \rightarrow \mathbb{R}^n$ be a ε -semi-LDP algorithm, where $R: S^{d-1} \rightarrow \mathbb{R}^n$ is an ε -LDP randomizer and $M_{\text{priv}}: Z^{n_{\text{priv}}} \rightarrow \mathbb{R}^d$ and $M_{\text{pub}}: Z^{n_{\text{pub}}} \rightarrow \mathbb{R}^d$ are aggregation protocols such that $\mathbb{E}_{M_{\text{priv}}; R} M_{\text{priv}}(R(x_1), \dots, R(x_{n_{\text{priv}}})) = x$ and $\mathbb{E}_{M_{\text{pub}}} M_{\text{pub}}(X_{\text{pub}}) = x$ for all $x \in S^{d-1}$. Then,*

$$\sup_{X \in S^{d-1}} \mathbb{E}_{A_{\text{semi-PrivU}}(X)} \|A_{\text{semi-PrivU}}(X) - x\|^2 \leq \sup_{X \in S^{d-1}} \mathbb{E}_{A_{\text{priv}}(X)} \|A_{\text{priv}}(X) - x\|^2.$$

Proof. First, Asi et al. (Asi et al., 2022, Proposition 3.4) showed that PrivUnit (with a proper choice of ρ, γ) has the smallest worst-case variance among all unbiased ε -LDP randomizers:

$$\sup_{X \in S^{d-1}} \mathbb{E} \|R(x) - x\|^2 \leq \sup_{X \in S^{d-1}} \mathbb{E} \|\text{PrivUnit}(x) - x\|^2 \quad (43)$$

for all ε -LDP randomizers R such that $\mathbb{E} R(x) = x$ for all $x \in S^{d-1}$.

Now, let R be a ε -LDP randomizer and M_{priv} and M_{pub} be aggregation protocols such that the assumptions in Proposition 17 are satisfied. We claim that there exists an unbiased ε -LDP randomizer $R^1: S^{d-1} \rightarrow \mathbb{R}^n$ such that

$$\sup_{X \in S^{d-1}} \mathbb{E}_{A_{\text{priv}}(X)} \|A_{\text{priv}}(X) - x\|^2 \leq \sup_{X \in S^{d-1}} \mathbb{E}_{R^1} \|R^1(x) - x\|^2. \quad (44)$$

To prove (44), we follow the idea in the proof of Asi et al. (2022, Proposition 3.3). Let P denote the uniform distribution on S^{d-1} . We have

$$\begin{aligned} \sup_{X \in S^{d-1}} \mathbb{E}_{A_{\text{priv}}(X)} \|A_{\text{priv}}(X) - x\|^2 &\leq \mathbb{E}_{X \sim P} \mathbb{E}_{A_{\text{priv}}(X)} \|A_{\text{priv}}(X) - x\|^2 \\ &\leq \mathbb{E}_{X \sim P} \mathbb{E}_{M_{\text{priv}}; R} \|M_{\text{priv}}(R(x_1), \dots, R(x_{n_{\text{priv}}})) - x\|^2 \\ &\quad + \mathbb{E}_{M_{\text{pub}}} \|M_{\text{pub}}(X_{\text{pub}}) - x\|^2, \end{aligned}$$

since the cross-term (inner product) vanishes by independence of X_{pub} and X_{priv} , and unbiasedness of M_{pub} . Now, (Asi et al., 2022, Lemma A.1) shows that there exist ε -LDP randomizers $\hat{R}_x: X_{\text{priv}} \rightarrow S^{d-1}$ such that $\mathbb{E} \hat{R}_x(v) = v$ for all $v \in S^{d-1}$ and

$$\mathbb{E}_{X_{\text{priv}} \sim P^{n_{\text{priv}}}} M_{\text{priv}}(R(x_1), \dots, R(x_{n_{\text{priv}}})) - x \leq \mathbb{E}_{X_{\text{priv}}} \mathbb{E}_{\rho} \|\hat{R}_x(\rho) - v\|^2.$$

Hence

$$\sup_{X \in S^{d-1}} \mathbb{E}_{A_{\text{priv}}(X)} \|A_{\text{priv}}(X) - x\|^2 \leq \mathbb{E}_{X_{\text{priv}}} \mathbb{E}_{\rho} \|\hat{R}_x(\rho) - v\|^2.$$

Define $R_x^1(v) := U^T \hat{R}_x(Uv)$ for $v \in S^{d-1}$, where U is a uniformly random rotation matrix such that $U^T U = \mathbf{I}_d$. Note that R_x^1 is an ε -LDP randomizer such that $\mathbb{E} R_x^1(v) = v$ for all $v \in S^{d-1}, x \in X_{\text{priv}}$. Moreover, for any fixed $v \in S^{d-1}, x \in X_{\text{priv}}$, we have

$$\begin{aligned} \mathbb{E} \|R_x^1(v) - v\|^2 &= \mathbb{E}_U \|\hat{R}_x(Uv) - Uv\|^2 \\ &= \mathbb{E}_{U, \rho} \|\hat{R}_x(\rho) - v\|^2 \end{aligned}$$

Let $x^* := \operatorname{argmin}_{x \in \mathcal{X}} \sup_{\nu \in \mathcal{P}^{S^d}} \mathbb{E} \{ R_x^1(\nu) \}$ and $R^1(\nu) := \mathbb{E} \{ R_x^1(\nu) \}$. Then putting the pieces together, we have

$$\begin{aligned} \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_A \{ n \mathcal{A}_P(X) \} &\leq \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E} \{ R_x^1(\nu) \} \\ &\leq n_{\text{priv}} \sup_{\mathcal{P}^{S^d}} \mathbb{E} \{ R^1(\nu) \} \\ &\leq \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_{R^1} \{ R^1(x) \}, \end{aligned}$$

by conditional independence of $\mathcal{R}^1(x)$ given X . This establishes (44). Thus,

$$\begin{aligned} n^2 \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_A \{ \mathcal{A}_P(X) \} &\leq \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_A \{ n \mathcal{A}_P(X) \} \\ &\leq \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_{R^1} \{ R^1(x) \} \\ &\leq \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E} \{ \text{PrivUnit}(x) \} \\ &\leq n^2 \sup_{\mathcal{X} \times \mathcal{P}^{S^d}} \mathbb{E}_{A_{\text{semi-PrivU}}} \{ \mathcal{A}_{\text{semi-PrivU}}(X) \}^2, \end{aligned}$$

where we used (43) in the last inequality. Dividing both sides of the above inequality by n^2 completes the proof. \square

F.2. Optimal Semi-LDP Stochastic Convex Optimization

If $\mu = 0$ (convex case), we denote $\mathcal{R}_{\text{SCO}}^{\text{loc}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) := \mathcal{R}_{\text{SCO}}^{\text{loc}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu = 0)$.

Theorem 48 (Complete statement of Theorem 18). *Let $\varepsilon \in (0, 1]$. There exist absolute constants c and C , with $0 < c \leq C$, such that*

$$cLD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon^2} \right\} \leq \mathcal{R}_{\text{SCO}}^{\text{loc}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \leq CLD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon^2} \right\},$$

and

$$cLD \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon^2} \right\} \leq \mathcal{R}_{\text{SCO}}^{\text{loc}}(\varepsilon, n_{\text{priv}}, n, d, L, D, \mu) \leq C \frac{L^2}{\mu} \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon^2} \right\}.$$

Proof. Lower bounds: Let A be ε -semi-LDP and denote $w_{\text{priv}} = \mathcal{A}_P(X)$.

Strongly convex lower bound: We begin with the strongly convex lower bounds, which can be proved straightforwardly by reducing strongly convex SCO to mean estimation and applying Theorem 14. In a bit more detail, let $f : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}$ be given by

$$f(w, x) = \frac{L}{2D} \|w - x\|^2,$$

where $\mathcal{W} = \mathcal{X} = \mathcal{D}B$. Note that f is L -uniformly Lipschitz and $\frac{L}{D}$ -strongly convex in w for all x . Further, $w^* := \operatorname{argmin}_{w \in \mathcal{W}} \mathbb{E} \{ f(w, x) \} = \mathbb{E} \{ x \}$. By a direct calculation (see e.g. (Kamath et al., 2022a, Lemma 6.2)), we have

$$\mathbb{E} \{ f(w_{\text{priv}}) \} - \mathbb{E} \{ f(w^*) \} = \frac{L}{2D} \mathbb{E} \{ \|w_{\text{priv}} - w^*\|^2 \}. \quad (45)$$

We can lower bound $\mathbb{E} \{ \|w_{\text{priv}} - w^*\|^2 \} = \mathbb{E} \{ \mathcal{A}_P(X) - \mathbb{E} \{ x \} \|^2 \}$ via Theorem 14 (and its proof, to account for the re-scaling). Specifically, there is a distribution P on \mathcal{X} such that

$$\mathbb{E} \{ \mathbb{E} \{ \mathcal{A}_P(X) - \mathbb{E} \{ x \} \|^2 \} \} \geq cD^2 \min \left\{ \frac{1}{n_{\text{pub}}}, \frac{d}{n\varepsilon^2} \right\}.$$

Combining this with (45) leads to the desired excess risk lower bound.

Convex lower bound: We will begin by proving the lower bounds for the case in which $L = D = 1$, and then scale our construction to get the lower bounds for arbitrary L, D .

Let $W = \mathcal{B}, X = \mathcal{X}^d$, and

$$f(w, x) = \langle w, xy \rangle,$$

which is convex and 1-uniformly-Lipschitz in w on X . For any ε -semi-LDP $A^1(p, X, \mathcal{W})$ and any $\gamma \in (0, 1]$, the proof of Theorem 14 constructs a distribution P on X with mean $E_X P \langle w, x \rangle = \theta P \langle w, x \rangle$ such that

$$E_{w^1 \sim X, P} \langle w^1, x \rangle - \theta^2 \leq \frac{\gamma^2}{4} \left(1 + \frac{d}{4000\gamma^2} \frac{n_{\text{priv}}\varepsilon^2}{n_{\text{pub}}} \right). \quad (46)$$

Now, let $F(w) = E_X P \langle w, x \rangle$ and $w = \arg\min_{w \in W} F(w)$. A direct calculation (see e.g. (Kamath et al., 2022a, Equation 14)) shows

$$\begin{aligned} E F(w) - F &\leq \frac{1}{2} E \langle \theta \rangle \langle w - w^1 \rangle^2 \\ &\leq \frac{1}{2} E \langle w^1 - \theta \rangle^2 \\ &\leq \frac{1}{2\gamma} E \langle w^1 - \theta \rangle^2 \end{aligned} \quad (47)$$

for any $w \in W, w^1 \in \mathcal{X}$. Note that $A^1(p, X, \mathcal{W})$ is ε -semi-DP if and only if $A^1(p, X, \mathcal{W}) : \langle \theta \rangle_{w_{\text{priv}}} \rightarrow \langle \theta \rangle_{w_{\text{priv}}}$ is ε -semi-DP, by post-processing. Thus, (46) and (47) together imply that any ε -semi-DP w^1_{priv} has worst-case excess risk that is lower bounded by

$$E F(w^1_{\text{priv}}) - F \leq \frac{\gamma}{8} \left(1 + \frac{d}{4000\gamma^2} \frac{n_{\text{priv}}\varepsilon^2}{n_{\text{pub}}} \right).$$

Choosing $\gamma^2 = c \min\left\{\frac{d}{n_{\text{priv}}}, \frac{1}{n_{\text{pub}}}\right\}$ for some small $c \in (0, 1]$ implies

$$E F(w^1_{\text{priv}}) - F \leq c^1 \min\left\{\frac{d}{n_{\text{priv}}\varepsilon^2}, \frac{1}{n_{\text{pub}}}\right\}$$

for some $c^1 \in (0, 1]$. This proves the desired lower bound for the case when $L = D = 1$. In the general case, we scale our hard instance, as in the proof of Theorem 4: Let $\mathcal{W} = DW, \mathcal{X} = LX$, and $x = \tilde{P} \tilde{\theta} \tilde{n} x = Lx$ for $x \in P$. Define $f : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}$ by $f(w, x) = \langle w, x \rangle$. Then f is L -Lipschitz and convex. Moreover, if $F(w) = E_X P \langle w, x \rangle, \theta = E_X P \langle x, x \rangle, w = \arg\min_{w \in W} F(w)$, $F(w) = E_X P \langle w, x \rangle, w = Dw$, and $\theta = E_X P \langle x, x \rangle = L\theta$, then $w = Dw \in \arg\min_{w \in W} F(w)$ and

$$\begin{aligned} F(w) - F &= \langle w, \theta \rangle - \langle w, \theta \rangle \\ &= D \langle w, \theta \rangle - \langle w, \theta \rangle \\ &= LD \langle w, \theta \rangle - \langle w, \theta \rangle \\ &= LD (F(w) - F). \end{aligned}$$

This shows that the excess risk of the scaled instance scales by LD , completing the lower bound proofs.

Upper bounds: The first term in each of the upper bounds ($LD \frac{1}{n_{\text{pub}}}$ for convex, and $L^2 \frac{1}{n_{\text{pub}}}$ for strongly convex) is attained by the throw-away algorithm that runs n_{pub} steps of (one-pass) SGD on X_{pub} (Nemirovski & Yudin, 1983).

The second term in the convex upper bound follows from the ε -LDP (hence semi-LDP) upper bound of Duchi et al. (2013, Proposition 3).

For the second term in the strongly convex upper bound, we run ε -LDP-SGD as in (Duchi et al., 2013). We also return a non-uniform weighted average \hat{w}_n of the iterates w_1, \dots, w_n as in (Rakhlin et al., 2012) to obtain

$$\mathbb{E}F(\hat{w}_n) - F \leq C \frac{G^2}{\mu n},$$

where $G^2 = \sup_{t \in [1, n]} \mathbb{E} \|\mathcal{M}_{\text{Duchi}}(\nabla f(w_t, x_t))\|^2 \leq L^2 \frac{d}{\varepsilon^2}$ (Duchi et al., 2013). Thus,

$$\mathbb{E}F(\hat{w}_n) - F \leq C^1 \frac{L^2}{\mu} \frac{d}{n\varepsilon^2},$$

since $d \leq 1 \leq \varepsilon^2$. This completes the proof. \square

F.2.1. AN “EVEN MORE OPTIMAL” SEMI-LDP ALGORITHM FOR SCO

Algorithm 4 Semi-LDP-SGD

- 1: **Input:** clip threshold $C \geq 0$, stepsize η , privacy parameter $\varepsilon \in (0, 1]$.
 - 2: Initialize $w_0 \in \mathcal{W}$.
 - 3: **for** $t \in \{0, 1, \dots, n-1\}$ **do**
 - 4: Draw random sample x_t from X without replacement.
 - 5: **if** $x_t \in X_{\text{priv}}$ **then**
 - 6: $\tilde{g}_t \in \mathcal{M}_{\text{Duchi}}(\text{clip}_{C, \varepsilon}(\nabla f(w_t, x_t)))$.
 - 7: **else**
 - 8: $\tilde{g}_t \in \nabla f(w_t, x_t)$
 - 9: **end if**
 - 10: Update $w_{t+1} = w_t - \eta \tilde{g}_t$.
 - 11: **end for**
 - 12: **Output:** $w_n = \frac{1}{n} \sum_{i=1}^n w_i$.
-

Proposition 49 (Re-statement of Proposition 19). *Let $f \in \mathcal{F}_{0, L; D}$ and P be any distribution and $\varepsilon \leq d$. Algorithm 4 is ε -semi-LDP. Further, there is an absolute constant c such that the output $\mathbb{E}F(w_n)$ of Algorithm 4 satisfies*

$$\mathbb{E}_{A; X} F(w_n) - F \leq c \frac{LD}{n} \max \left\{ \frac{d}{\varepsilon^2}, \frac{C}{n_{\text{priv}}}, \frac{C}{n_{\text{pub}}} \right\}.$$

Proof. Privacy: Since $\mathcal{M}_{\text{Duchi}}$ is an ε -LDP randomizer and we are applying $\mathcal{M}_{\text{Duchi}}$ to the gradients of all the private samples $x \in X_{\text{priv}}$, Algorithm 4 is ε -semi-LDP.

Excess risk: Choose $C \leq L$: i.e. we don't clip, since stochastic subgradients are already uniformly bounded by the L -Lipschitz assumption. By the classical analysis of the stochastic subgradient method (see e.g. (Bubeck et al., 2015)), we can obtain

$$\begin{aligned} \mathbb{E}F(w_n) - F &\leq \frac{1}{n} \sum_{t=1}^n \mathbb{E} \|\tilde{g}_t\|^2 \eta^2 \\ &\leq \frac{D^2}{\eta n} \left(\frac{\eta}{n} n_{\text{priv}} G_a^2 + n_{\text{pub}} G_b^2 \right), \end{aligned}$$

where $G_a^2 = \sup_t \mathbb{E} \|\mathcal{M}_{\text{Duchi}}(\nabla f(w_t, x_t))\|^2$ and $G_b^2 = \sup_t \mathbb{E} \|\nabla f(w_t, x_t)\|^2$. By the uniform Lipschitz assumption, we have $G_b^2 \leq L^2$. By (Duchi et al., 2013), we have $G_a^2 \leq c^2 L^2 \frac{d}{\varepsilon^2}$ for some absolute constant $c \geq 0$. Thus, choosing $\eta = \frac{D}{L} \min \left\{ \frac{\varepsilon^2}{n_{\text{priv}} c^2 d}, \frac{1}{n_{\text{pub}}} \right\}$ yields

$$\mathbb{E}F(w_n) - F \leq 3 \frac{LD}{n} \max \left\{ \frac{C}{\varepsilon^2}, \frac{C}{n_{\text{priv}}}, \frac{C}{n_{\text{pub}}} \right\}.$$

This completes the proof. \square

G. Numerical Experiments

Code for all of the experiments is available here: <https://github.com/optimizati-on-for-data-driven-science/DP-with-public-data>.

G.1. Central Semi-DP Experiments

G.1.1. SEMI-DP LINEAR REGRESSION WITH GAUSSIAN DATA

Data Generation: We implement Algorithm 1 on synthetic data designed for a linear regression problem of dimension 2,000, using the squared error loss: $f(w, x, y) = \|xw - y\|^2$, where $x \in \mathbb{R}^d$ denotes the feature vector and $y \in \mathbb{R}$ is the target. Here, $d = 2,000$. Our synthetic dataset consists of $n = 30,000$ training samples, 7500 validation samples, and 37,500 test samples. The feature vectors $x_i \in \mathbb{R}^d$ and the optimal parameter vector w are drawn i.i.d. from a multivariate Gaussian distribution $\mathcal{N}(0, \mathbf{I}_{2000})$. We generate predicted values \hat{y}_i from a Gaussian distribution $\mathcal{N}(x_i w, \sigma^2)$. Thus, the optimal linear regression model has which ensures an optimal mean squared error of 1.

Experimental Setup: Our experiments investigate two phenomena: 1) the effect of the ratio $\frac{n_{\text{pub}}}{n}$ on test loss when $\epsilon \in \{2, 4\}$ is fixed, for values of $\frac{n_{\text{pub}}}{n}$ ranging from 0.01 to 0.95, see Figures 3-4; and 2) the effect of privacy (quantified by ϵ) on test loss, for fixed $\frac{n_{\text{pub}}}{n} \in \{0.1, 0.25\}$, and varying $\epsilon \in \{0.1, 0.5, 1.0, 2.0, 4.0, 8.0\}$, see Figures 9-10. We maintain the privacy parameter δ at a constant value of 10^{-5} throughout our experiments. We set the private batch size $K_{\text{priv}} = 500$, public batch size $K_{\text{pub}} = 200$, and iterations $T = 5000$. All algorithms undergo extensive hyperparameter tuning using the validation dataset, and the performance of each tuned algorithm is subsequently assessed using the test dataset. (See “Hyperparameter Tuning” paragraph below for details on the tuning process.)

Details on Implementations of Algorithms: We compare four different semi-DP algorithms: 1. *Throw-away*. 2. *DP-SGD* (Abadi et al., 2016; De et al., 2022). 3. *PDA-MD* (Amid et al., 2022, Algorithm 1). 4. *Our Algorithm 1*—specifically, the *sample-with-replacement* version of our algorithm. If $n_{\text{pub}} \neq d$, the *Throw-away* algorithm simply returns a minimizer w_{pub} of the public loss: $w_{\text{pub}} = (X_{\text{pub}}^T X_{\text{pub}})^{-1} X_{\text{pub}}^T y_{\text{pub}}$. Otherwise, we used pretrained warm-start models with all public samples. (See “warm-start” paragraph below for details.) *DP-SGD* adds noise to all (public or private) gradients. We use the state-of-the-art (for image classification) implementation of DP-SGD of (De et al., 2022). We adopt the re-parameterization of DP-SGD in (De et al., 2022, Equation 3) to ease the hyperparameter tuning. For *PDA-MD*, we implement (Amid et al., 2022, Algorithm). We use their exact Mirror descent form by multiplying the inverse of the hessian $X_{\text{pub}}^T X_{\text{pub}}$ by the private gradient (Amid et al., 2022). In their original implementation, they added a small constant, Hessian Regularization, times the identity matrix to the Hessian before calculating the inverse for numerically stable. We choose the hessian regularization constant as 0.01, the same value in their original implementation.

Effect of increasing ratio $n_{\text{pub}}/n_{\text{priv}}$: More public data always improves the performance of Algorithm 1, but does not always benefit PDA-MD. The primary reason for this is that our algorithm more effectively handles the increasing privacy noise that is needed to maintain semi-DP with increasing ratio n_{pub}/n . Our algorithm achieves this by using the weight parameter to reduce the variance of the increasingly noisy private gradients and leverage public gradients. By contrast, as n_{pub}/n ratio rises, the efficacy of PDA-MD may diminish due to its over-reliance on increasingly noisy private gradients. We verify our reasoning numerically in Appendix G.3: Tables 11 and 12 record the standard deviation σ of the privacy noise for different ratios.

Effect of hessian regularization parameter on PDA-MD: Upon reproducing the results of PDA-MD, we discovered a high sensitivity to the choices of Hessian Regularization, especially when the ratio of the largest and the smallest eigenvalue of the data matrix is large. We test PDA-MD on the dataset proposed in the original study as well as on our own dataset. The results of these tests are displayed in Fig. 8. We see PDA-MD is sensitive to hessian regularization parameter, which requires extra tuning on complicated tasks. We implement PDA-MD with the optimal regularization value of 0.01.

Details on warm-start and cold-start: In our evaluation, we adopt a “warm start” strategy for all algorithms: we first find a minimizer w_{pub} of the public loss, and then initialize the training process with w_{pub} . Note that there are two cases: $n_{\text{pub}} \neq n$ and $n_{\text{pub}} = n$. In the case of $n_{\text{pub}} \neq n$, the minimizer w_{pub} of the public loss can be obtained via $w_{\text{pub}} = (X_{\text{pub}}^T X_{\text{pub}})^{-1} X_{\text{pub}}^T y_{\text{pub}}$. In the case of $n_{\text{pub}} = n$, we run SGD on linear regression with all n_{pub} . Specifically, we minimize w_{pub} of the public loss by running the SGD with public batch size $K_{\text{pub}} = 200$, stepsize $\eta_t = 0.5$, and iterations $T = 50000$ to allow the models fully converge. For cold start scenarios, we initialize all model parameters w to 0. The cold

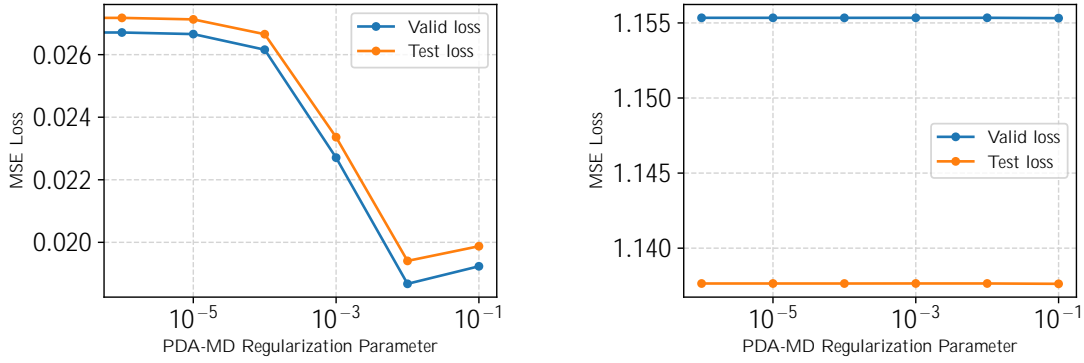


Figure 8. Loss v.s. PDA-MD Regularization Parameter. Left: Results on proposed dataset in (Amid et al., 2022). Right: Results on our dataset. PDA-MD is sensitive to hessian regularization parameter, which requires extra tuning on complicated tasks.

start experiment results can be seen in Figure 5 and Figure 11.

Clipping public gradients improves performance: We have empirically found that a slight modification of Algorithm 1 offers performance benefits. Namely, it is beneficial to project the public gradients onto the ℓ_2 -sphere of radius C ; that is, we re-scale the public gradients to have ℓ_2 -norm equal to C . Note that semi-DP still holds regardless of whether or not the public gradients are re-scaled. Re-scaling the public gradients helps balance the effects of the public and private gradients on the optimization trajectory. In the original method stated in Algorithm 1, if unclipped public gradients and the private gradient are of very different magnitudes, then one gradient direction might dominate the optimization procedure, leading to a sub-optimal model. Thus, our public gradient re-scaling technique promotes a more balanced update, which gracefully combines the public and private data in each iteration. To the best of our knowledge, this technique is novel.

Privacy accounting: We compute the privacy loss of each algorithm by using the moments accountant of Abadi et al. (Abadi et al., 2016). For a fixed clip threshold C , privacy level $\rho_\epsilon, \delta q$ is determined by three parameters: the variance of privacy noise σ^2 , the private sampling ratio $q : K_{\text{priv}}\{n_{\text{priv}}\}$, and the total number of iterations T . In our setting, the privacy parameters $\rho_\epsilon, \delta q$ are given, and we use the moments accountant to compute an approximation of σ^2 for any choice of hyperparameters T and q . We utilize the implementation of the privacy accountant provided by the Pytorch privacy framework, Opacus.

Hyperparameter Tuning: The results reported are for each algorithm with the hyperparameters (step size and α in Semi-DP) that attain the best performance for a given experiment. For simplicity and computation efficiency, we keep clip threshold $C = 1$ for all of our experiments. Preliminary experiments found that a clip threshold of $C = 1$ worked well for all algorithms. To tune all algorithms, we use grid search. See Tables 1 and 2 in Appendix G.2 for detailed descriptions of the hyperparameter search grids.

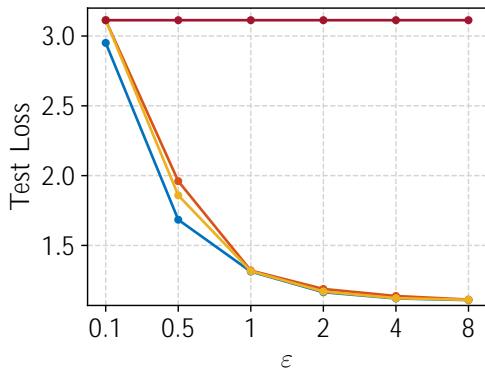


Figure 9. Test loss vs. ϵ . $\frac{n_{\text{pub}}}{n} = 0.1$.

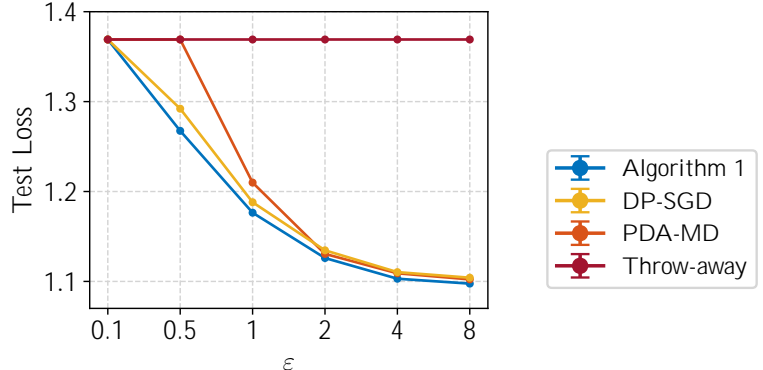


Figure 10. Test loss vs. ϵ . $\frac{n_{\text{pub}}}{n} = 0.25$.

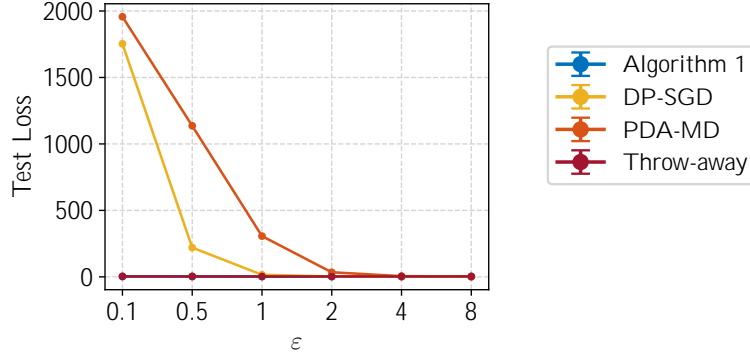


Figure 11. Test loss vs. ϵ . $\frac{n_{\text{pub}}}{n} = 0.1$, without warm-start.

G.1.2. SEMI-DP MEAN ESTIMATION

We present an experiment with mean estimation: We fix ρ -semi-zCDP with privacy parameter $\rho = 0.5$. We draw n i.i.d. samples from a $pd = 1000q$ -dimensional **Bernoulli** $p \in [0, 1]$ product distribution. We investigated the effect of the ratio $\frac{n_{\text{pub}}}{n}$ on mean ℓ_2 error, for values of $\frac{n_{\text{pub}}}{n}$ ranging from 0.05 to 0.95. We presented the experiment in three different high-dimensional and low-dimensional settings: $n = d$, $n < d$, and $n \gg d$. Fig. 12 shows that when $n = d$, throw-away outperforms Gaussian mechanism. Fig. 13 and Fig. 14 show that when $d \propto n$, throw-away outperforms Gaussian mechanism except when $\frac{n_{\text{pub}}}{n} = 0.05$. Moreover, **our Weighted Gaussian estimator outperforms both throw-away and the Gaussian mechanism in every case.**

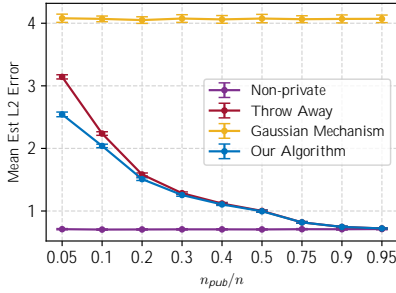


Figure 12. $d = 1000, n = 500$

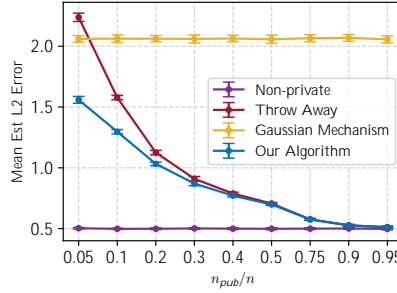


Figure 13. $d = 1000, n = 1000$

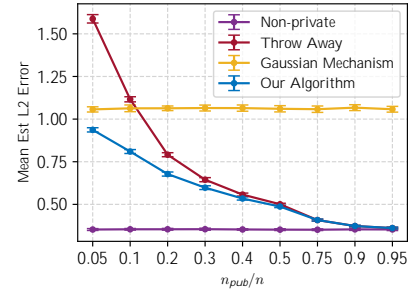


Figure 14. $d = 1000, n = 2000$

G.1.3. SEMI-DP LOGISTIC REGRESSION WITH CIFAR-10

We evaluate the performance of Algorithm 1 in training a logistic regression model to classify digits in the CIFAR-10 dataset (Krizhevsky et al., 2009). We compare Algorithm 1 against DP-SGD and throw-away. Note that PDA-MD does not have an efficient implementation for logistic regression. This is because there does not exist a closed form update rule for the mirror descent step. Thus, we do not compare against PDA-MD.

We flatten the images and feed them to the logistic (softmax) model. Cross-entropy loss is used here; therefore, the model is convex. Implementations of the algorithms are similar to the linear regression case. However, in this case, throw-away consists of running non-private SGD on X_{pub} to find an approximate minimizer w_{pub} . For all three algorithms, we fixed batch-size 256 and privacy parameter $\delta = 10^{-6}$. The remaining hyperparameters are tuned by grid search, using the same grid for each algorithm. Also, in contrast to the linear regression experiments, we do not use warm-start for any of the algorithms and we use SGD for all algorithms in these experiments. See Table 3 in Appendix G.2 for detailed descriptions of the hyperparameter search grids.

Results are reported in Figs. 15 to 18. *Our Algorithm 1 always outperforms the baselines* in terms of minimizing test accuracy/error. The advantage of Algorithm 1 over these baselines is even more pronounced than it was for linear regression. This might be partially due to the fact that we did not use warm-start.

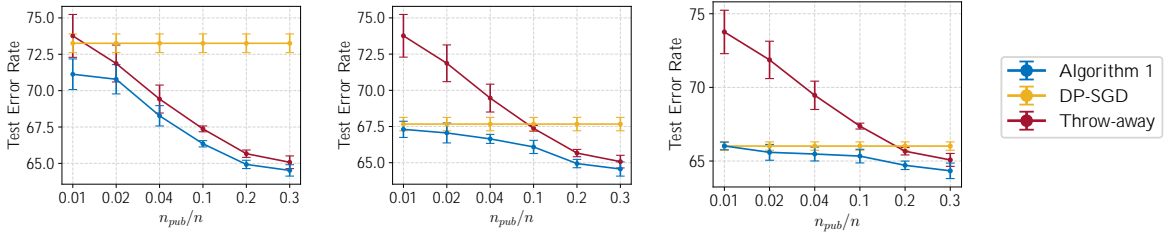


Figure 15. Test error rate vs. $\frac{n_{pub}}{n}$. Left: $\epsilon = 0.1$. Middle: $\epsilon = 0.5$. Right: $\epsilon = 1.0$

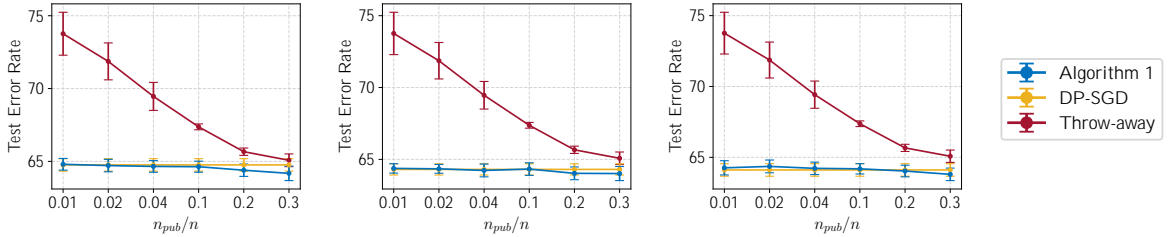


Figure 16. Test error rate vs. $\frac{n_{pub}}{n}$. Left: $\epsilon = 2.0$. Middle: $\epsilon = 4.0$. Right: $\epsilon = 8.0$

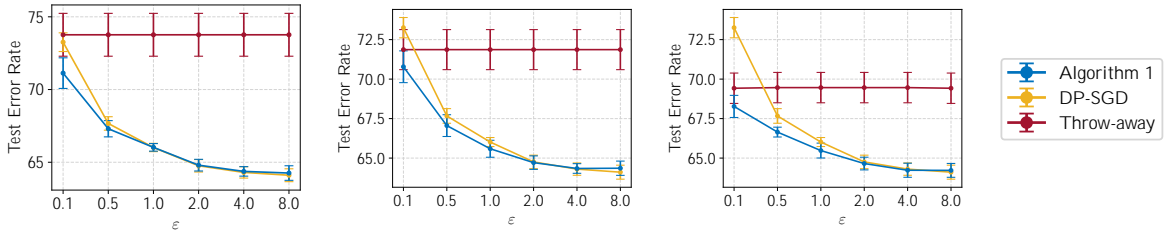


Figure 17. Test error rate vs. ϵ . Left: $\frac{n_{pub}}{n} = 0.01$. Middle: $\frac{n_{pub}}{n} = 0.02$. Right: $\frac{n_{pub}}{n} = 0.04$

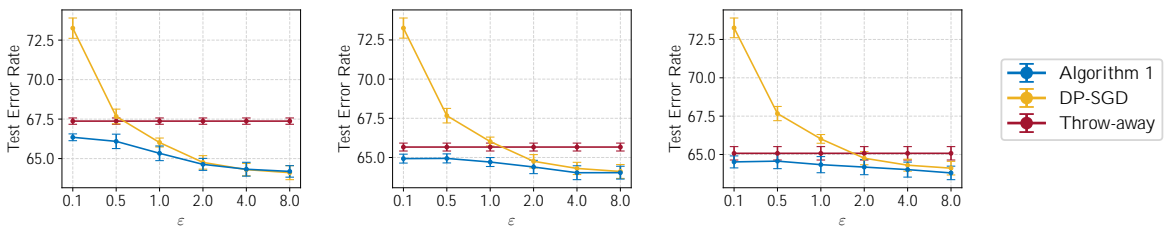


Figure 18. Test error rate vs. ϵ . Left: $\frac{n_{pub}}{n} = 0.1$. Middle: $\frac{n_{pub}}{n} = 0.2$. Right: $\frac{n_{pub}}{n} = 0.3$

G.1.4. SEMI-DP WIDE-RESNET-16-4 WITH CIFAR-10

To evaluate the performance of Algorithm 1 in training non-convex (neural) models, we use the Wide-ResNet with 16 layers and a width factor of 4 (WRN-16-4) (Zagoruyko & Komodakis, 2017). When trained non-privately on CIFAR-10 dataset, this model achieves an error rate of 5.02 (Zagoruyko & Komodakis, 2017).

Experimental Setup: We followed the same experimental setup as (De et al., 2022). For all algorithms, we used a constant learning rate and did not use momentum, weight decay, or dropout. We fixed the batch size of 256 and privacy parameter $\delta = 10^{-5}$. For simplicity and to isolate the effects of the weighted gradient estimator, we used the WRN-16-4

without batch/group normalization or data augmentation.¹⁰We split the CIFAR-10 dataset (Krizhevsky et al., 2009) of 50,000 training examples into 40,000 training samples and 10,000 validation samples. We used their 10,000 test images as our test set. Same as Appendix G.1.3, we did not use warm-start here.

Hyperparameter Tuning: We selected the hyperparameters settings with the highest validation accuracy for all algorithms and reported their test accuracy on the official test set. To tune the hyperparameters, we used the bayesian hyperparameter optimization technique. That is, we build a probability model of the objective function and use it to select the most promising hyperparameters to evaluate in the true objective function. See Table 4 in Appendix G.2 for detailed descriptions of the hyperparameter search range.

Results: We investigate two cases: 1) the effect of the ratio of public samples $\frac{n_{pub}}{n}$ on accuracy when $\epsilon = 8$ is fixed. We test values of $\frac{n_{pub}}{n}$ ranging from 0.01 to 0.3; and 2) the effect of privacy (quantified by ϵ) on accuracy, for fixed $\frac{n_{pub}}{n} = 0.04$. We vary $\epsilon \in \{0.1, 0.5, 1.0, 2.0, 4.0, 8.0\}$. Results are reported in Figs. 19 and 20. *Our Algorithm 1 always outperforms the baselines (DP-SGD and Throw-away) in terms of minimizing test error, across all experimental setups.*

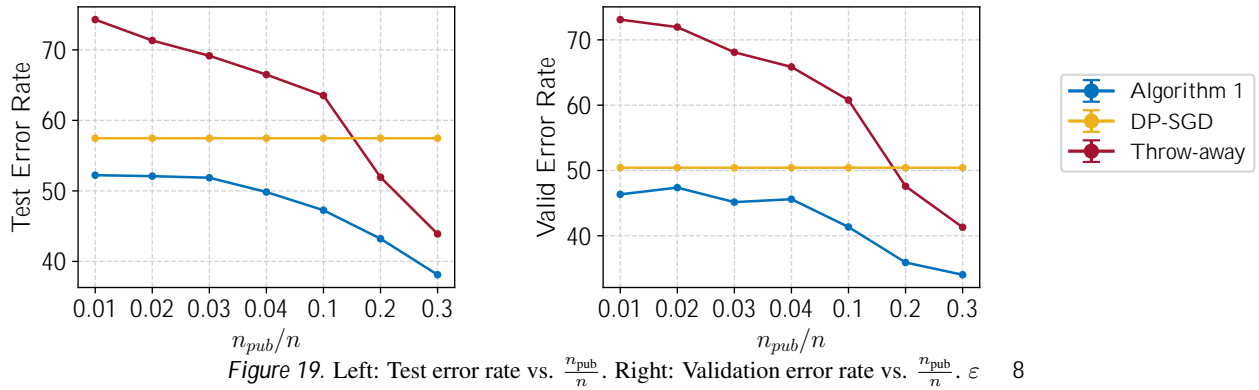


Figure 19. Left: Test error rate vs. $\frac{n_{pub}}{n}$. Right: Validation error rate vs. $\frac{n_{pub}}{n}$. $\epsilon = 8$

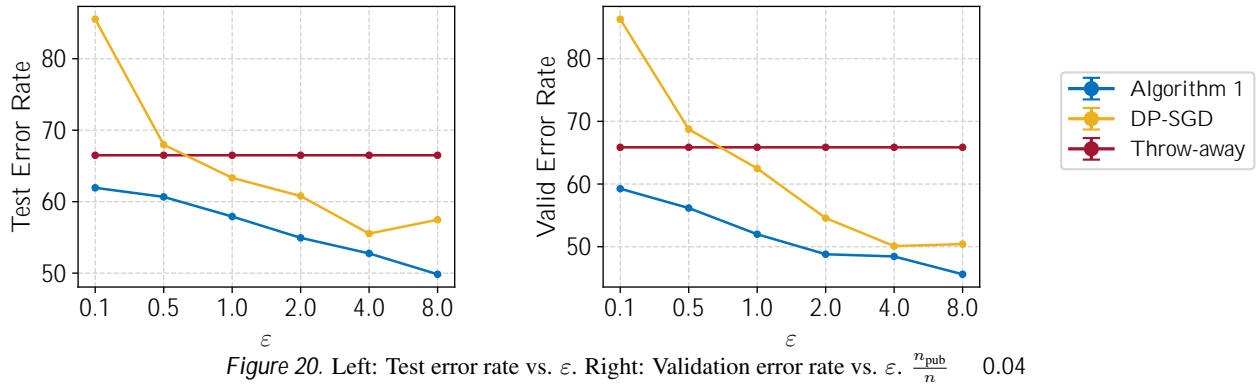


Figure 20. Left: Test error rate vs. ϵ . Right: Validation error rate vs. ϵ . $\frac{n_{pub}}{n} = 0.04$

G.2. Hyperparameters Search Grids

hyperparameter	learning-rate
value	0, 0.01, 0.03, 0.05, 0.07, 0.09, 0.1, 0.3, 0.5, 0.7, 0.9, 1.1, 1.3, 1.5, 1.7, 1.9

Table 1. Grid used for hyperparameter search for PDA-MD and DP-SGD in Appendix G.1.1

¹⁰We believe that accuracy could be improved by combining these tricks with our semi-DP gradient estimator (De et al., 2022; Nasr et al., 2023).

Optimal Differentially Private Model Training with Public Data

hyperparameter	learning-rate	α
value	0, 0.01, 0.03, 0.05, 0.07, 0.09, 0.1, 0.3, 0.5, 0.7, 0.9, 1.1, 1.3, 1.5, 1.7, 1.9	0, 0.1, 0.2 . . . 1

Table 2. Grid used for hyperparameter search for Algorithm 1 in Appendix G.1.1

hyperparameter	value
iterations	2000, 4000, 6000, . . . , 20000
learning-rate	0.005, 0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.5, 0.75, 1, 1.5, 2, 3, 4, 5
α	0, 0.1, 0.2, . . . , 1

Table 3. Grid used for hyperparameter search for DP-SGD and Semi-DP in Appendix G.1.3

hyperparameter	iterations	learning-rate	α
value	r3000, 7000s	r0.1, 3s	r0, 1s

Table 4. Range used for Bayesian hyperparameter search for DP-SGD and Semi-DP in Appendix G.1.4

G.3. Exact Numerical Results

$n_{\text{pub}} \{n$	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.01	2.3526	1689.5905	2.7935	1689.5905
0.03	1.9719	1084.2257	2.1457	1084.2256
0.04	1.7626	776.1302	2.1417	776.1301
0.1	1.1648	1.1880	1.1695	3.1133
0.25	1.1260	1.1306	1.1346	1.3691
0.5	1.1043	1.1394	1.1065	1.1539
0.75	1.0946	1.1013	1.0937	1.1013
0.9	1.0787	1.0787	1.0787	1.0788
0.95	1.0725	1.0725	1.0725	1.0725

Table 5. Exact training results of curves reported in Figure 3.

$n_{\text{pub}} \{n$	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.01	1.5020	1689.5905	1.5206	1689.5905
0.03	1.3363	1084.2257	1.3769	1084.2256
0.04	1.3196	776.1302	1.3776	776.1302
0.1	1.1201	1.1376	1.1221	3.1133
0.25	1.1030	1.1090	1.1103	1.3691
0.5	1.0911	1.1019	1.0999	1.1539
0.75	1.0863	1.1013	1.0877	1.1013
0.9	1.0787	1.0787	1.0787	1.0788
0.95	1.0725	1.0725	1.0725	1.0725

Table 6. Exact training results of curves reported in Figure 4.

ϵ	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.10	2.9509	3.1133	3.1133	3.1133
0.50	1.6841	1.9602	1.8588	3.1133
1.00	1.3125	1.3201	1.3158	3.1133
2.00	1.1648	1.1880	1.1695	3.1133
4.00	1.1201	1.1376	1.1221	3.1133
8.00	1.1088	1.1120	1.1104	3.1133

Table 7. Exact training results of curves reported in Figure 9.

ϵ	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.10	1.3691	1.3691	1.3691	1.3691
0.50	1.2647	1.3691	1.2679	1.3691
1.00	1.1764	1.2099	1.1880	1.3691
2.00	1.1260	1.1306	1.1346	1.3691
4.00	1.1030	1.1090	1.1103	1.3691
8.00	1.0976	1.1023	1.1041	1.3691

Table 8. Exact training results of curves reported in Figure 10.

$n_{\text{pub}} \setminus n$	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.01	1.5175	2010.2422	1.5411	1689.5905
0.03	1.4225	2010.2422	1.5411	1084.2256
0.04	1.3989	2010.2422	1.5413	776.1302
0.1	1.3371	5.3822	1.5413	3.1133
0.25	1.2574	2.6205	1.5411	1.3691
0.5	1.1539	5.9698	1.5411	1.1539
0.75	1.1013	78.9970	1.5411	1.1013
0.9	1.0787	1053.6185	1.5411	1.0788
0.95	1.0725	1662.3572	1.5413	1.0725

Table 9. Exact training results of curves reported in Figure 5.

ϵ	Algorithm 1	PDA-MD	DP-SGD	Throw-away
0.10	3.1133	1956.6069	1830.1820	3.1133
0.50	3.1133	1092.1338	213.6728	3.1133
1.00	2.1843	306.5518	15.8089	3.1133
2.00	1.6313	34.4665	2.7226	3.1133
4.00	1.3359	5.3648	1.4746	3.1133
8.00	1.1986	2.4319	1.2472	3.1133

Table 10. Exact training results of curves reported in Figure 11.

$n_{\text{pub}} \setminus n$	0.01	0.03	0.04	0.1	0.25	0.5	0.75	0.9	0.95
σ	2.490	2.529	2.568	2.744	3.252	4.805	9.531	23.672	47.344

Table 11. Standard Deviation σ of the private noise v_t used in experiment shown in Figure 3.

$n_{\text{pub}}\{n$	0.01	0.03	0.04	0.1	0.25	0.5	0.75	0.9	0.95
σ	1.470	1.489	1.509	1.597	1.860	2.671	5.176	12.812	25.586

Table 12. Standard Deviation σ of the private noise v_t used in experiment shown in Figure 4 and 5.

H. Limitations

Limitations of Theoretical Results: Our theoretical results rely on certain assumptions (e.g. convex, Lipschitz loss, i.i.d. data for SCO), that may be violated in certain applications. We leave it as future work to investigate the questions considered in this work under different assumptions (e.g. non-convexity, semi-DP SCO with *out-of-distribution* public data). Also, we reiterate that our theoretical results describe the optimal *worst-case* error. It might be the case that the worst-case distributions we construct in our lower bound proofs are unlikely to appear in practice. Thus, another interesting direction for future work would be to analyze “instance-optimal” (Asi & Duchi, 2020a;b) semi-DP error rates.

Limitations of Experiments: It is important to note that pre-processing and hyperparameter tuning were not done in a DP manner, since we did not want to detract focus from evaluation of the (fully tuned) semi-DP algorithms.¹¹ As a consequence, the overall privacy loss for the entire experimental process is higher than the ϵ indicated in the plots: the ϵ indicated in the plots solely reflects the privacy loss from running the algorithms with fixed hyperparameters and (pre-processed) data.

¹¹See, e.g. (Liu & Talwar, 2019; Papernot & Steinke, 2022) and the references therein for discussion of DP hyperparameter tuning.