
Faster Differentially Private Convex Optimization via Second-Order Methods

Arun Ganesh
Google Research

Mahdi Haghifam*
University of Toronto,
Vector Institute

Thomas Steinke
Google DeepMind

Abhradeep Thakurta
Google DeepMind

Abstract

Differentially private (stochastic) gradient descent is the workhorse of DP private machine learning in both the convex and non-convex settings. Without privacy constraints, second-order methods, like Newton’s method, converge faster than first-order methods like gradient descent. In this work, we investigate the prospect of using the second-order information from the loss function to accelerate DP convex optimization. We first develop a private variant of the regularized cubic Newton method of Nesterov and Polyak [NP06], and show that for the class of strongly convex loss functions, our algorithm has quadratic convergence and achieves the optimal excess loss. We then design a practical second-order DP algorithm for the unconstrained logistic regression problem. We theoretically and empirically study the performance of our algorithm. Empirical results show our algorithm consistently achieves the best excess loss compared to other baselines and is 10-40 \times faster than DP-GD/DP-SGD for challenging datasets.

1 Introduction

Many machine learning tasks reduce to a convex optimization problem. More precisely, given a dataset $S_n = (z_1, \dots, z_n) \in \mathcal{Z}^n$, a closed, convex set $\mathcal{W} \subseteq \mathbb{R}^d$, and a loss function $f : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}$ such that, for every $z \in \mathcal{Z}$, $f(w, z)$ is a convex function in w , our goal is to compute an approximation to $\arg \min_{w \in \mathcal{W}} \left(\ell(w, S_n) \triangleq \frac{1}{n} \sum_{i \in [n]} f(w, z_i) \right)$. In this paper, we are interested in the problem of designing optimization algorithms in the scenario that the dataset S_n contains private information. Differential privacy (DP) [DMNS06] is a formal standard for privacy-preserving data analysis that provides a framework for ensuring that the output of an analysis on the data does not leak this private information. This problem is known as *private convex optimization*: Design an algorithm $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$ that is both DP and ensures low *excess loss* $\triangleq \ell(\mathcal{A}(S_n), S_n) - \min_{w \in \mathcal{W}} \ell(w, S_n)$.

The predominant algorithm for private convex optimization is DP (stochastic) gradient descent (DP-GD/DP-SGD). This is a *first-order* iterative method. I.e., we start with an initial value w_0 and iteratively update it using the gradient of the loss $\nabla_{w_t} \ell(w_t, S_n)$ following the update rule $w_{t+1} = w_t - \eta \cdot (\nabla_{w_t} \ell(w_t, S_n) + \xi_t)$, where $\eta > 0$ is a constant and ξ_t is Gaussian noise to ensure privacy. The number of iterations T also determines the amount of noise at each iteration, i.e., the scale of ξ_t is proportional to \sqrt{T} due to the composition of DP. Note that we assume $\|\nabla_{w_t} \ell(w_t, S_n)\| \leq 1$.

One of the major drawbacks of DP-(S)GD is *slow convergence*. The choice of (η, T) exhibits a tradeoff in terms of the excess loss: if $\eta \cdot T$ is small, the algorithm cannot reach the optimal solution; on the other hand, the magnitude of noise at each iteration is $\eta \cdot \sqrt{T}$, which cannot be too large. Therefore, to maximize $\eta \cdot T$ and minimize $\eta \cdot \sqrt{T}$, implementations of DP-(S)GD err on the side of

*This work was carried out while the author was an intern at Google Research, Brain Team.
m.haghifam@northeastern.edu, {arunganesh, steinke, athakurta}@google.com

large T and small η , which results in a long, slow path to convergence. This fact has been shown theoretically as well: for the class of β -smooth convex functions, the optimal instantiations of DP-GD use a step size of $\max\{1/\sqrt{n}, \sqrt{d}/\varepsilon n\}$ [BFTG19] while in the non-private setting the stepsize for GD is set to $1/\beta$. Smaller step size requires more steps (i.e. more iterations) to converge. This slowness is exacerbated by the facts that (1) DP-SGD requires large batch sizes for good performance [PHKX+23] and (2) the hyperparameter tuning of DP-(S)GD, and generally DP algorithms, is a challenging task [PS22]. *Can we design a DP optimization algorithm which accelerates DP-(S)GD by choosing the step size dynamically based on the local geometry of the loss function?*

We draw inspiration from the non-private optimization literature: To address the slow convergence of GD and of first-order methods in general, a class of algorithms based on *preconditioning* the gradient using second-order information has been developed [Nes98; NW99]. This class of algorithms is based on successively minimizing a quadratic *approximation* of the function, i.e., $w_{t+1} = w_t + \Delta_t$ where $\Delta_t = \arg \min_{\Delta} \{\ell(w_t, S_n) + \langle \nabla \ell(w_t, S_n), \Delta \rangle + \frac{1}{2} \langle H_t \cdot \Delta, \Delta \rangle\} = - (H_t)^{-1} \nabla \ell(w_t, S_n)$. Here, H_t is a scaling matrix which provides curvature information about the loss $\ell(\cdot, S_n)$ at w_t . For instance, Newton’s method uses the Hessian $H_t = \nabla^2 \ell(w_t, S_n)$. Second-order algorithms significantly improve over the convergence speed of GD, and key to their success is that at each step they *automatically* tune the stepsize along each dimension based on the local curvature.

In this paper, our goal is to accelerate DP convex optimization. In particular, the current paper revolves around the following questions: Can the second-order information *accelerate* private convex optimization while achieving *optimal excess error*? What is the best way to *privatize second-order information*, e.g., the Hessian matrix? How does the achievable *privacy-utility-runtime tradeoff* compare with first-order methods such as DP-GD? We show that second-order information can accelerate DP optimization while achieving excess loss that matches or improves on DP-GD. Our main contributions are both theoretical and empirical:

1.1 Provably Optimal Algorithm for Strongly Convex Functions

Newton’s method is a second-order optimization technique that is well-known for its rapid convergence for strongly convex and smooth functions in non-private optimization. Specifically, to achieve an excess loss of α , the method only requires $O(\log \log(1/\alpha))$ iterations, which is provably faster than the convergence rate of *any* first-order method. One natural question is whether it is possible to design a second-order DP convex optimization algorithm that can achieve the *optimal minmax* excess error err^{opt} in $O(\log \log(1/\text{err}^{\text{opt}}))$ iterations? We provide an affirmative answer to this question in Section 4 by designing a second-order DP algorithm based on the cubic regularized Newton’s method of Nesterov and Polyak [NP06]. At each step t , we compute a *cubic* upper bound $\ell(w + \Delta, S_n) \leq \ell(w, S_n) + \langle \nabla_w \ell(w, S_n), \Delta \rangle + \frac{1}{2} \langle \nabla_w^2 \ell(w, S_n) \cdot \Delta, \Delta \rangle + O(\|\Delta\|^3)$. We can minimize this cubic upper bound using *any* DP convex optimization subroutine; the minimizer becomes the next iterate w_{t+1} . Since the cubic is a universal upper bound, our algorithm converges globally

1.2 Fast Practical Algorithms for DP Logistic Regression

DP logistic regression is a popular approach for private classification, with DP-GD/DP-SGD being the predominant class of algorithms for this task. As we numerically show, DP-GD/DP-SGD exhibit slow convergence for this task (See Figure 1). In Section 5, we develop a practical algorithm that injects carefully designed noise into Newton’s update rule as follows:

$$w_{t+1} = w_t - \Psi(\nabla_{w_t}^2 \ell(w_t, S_n))^{-1} \cdot (\nabla_{w_t} \ell(w_t, S_n) + \xi_{t,1}) + \xi_{t,2}. \quad (1)$$

In particular, we inject noise twice: $\xi_{t,1}$ privatizes the gradient and $\xi_{t,2}$ privatizes the direction. The function Ψ modifies the Hessian to ensure that the eigenvalues are not too small; this is essential for bounding the sensitivity and, hence, the scale of $\xi_{t,2}$. We consider two types of modification based on *eigenvalue clipping* and *eigenvalue adding*. For eigenvalue clipping, $\Psi(\nabla_{w_t}^2 \ell(w_t, S_n))$ replaces the eigenvalues λ_i of $\nabla_{w_t}^2 \ell(w_t, S_n)$ with $\max\{\lambda_i, \lambda_0\}$, where $\lambda_0 > 0$ is a carefully chosen constant. For eigenvalue adding, $\Psi(\nabla_{w_t}^2 \ell(w_t, S_n)) = \nabla_{w_t}^2 \ell(w_t, S_n) + \lambda_0 I$. Using Ψ we can control the sensitivity and still have fast convergence, since important curvature information is generally

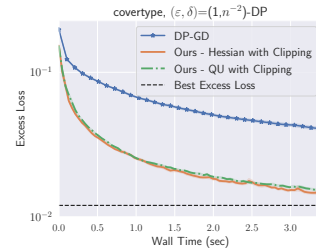


Figure 1: Excess loss versus runtime of DP-GD & our algorithms.

contained in the larger eigenvalues/vectors of the Hessian. We prove the local convergence of the update rule (1) in Section 5.3 and perform a thorough empirical evaluation Section 6. We demonstrate that our algorithm outperforms existing baselines on a variety of benchmarks.

Ensuring Global Convergence. One limitation of the update rule in Equation (1) is it does not converge globally (even without noise added for DP). That is, if the initial point w_0 is too far from the optimal solution, then the iterates may diverge. To address this problem, we propose a variant of Newton’s update rule where we replace the Hessian with a different form of second-order information which gives a *Quadratic Upperbound* (QU) on the logistic loss. This is *guaranteed to converge globally*, like the cubic Newton approach. And we show numerically that this algorithm converges almost as fast as the regular Newton’s method in the private setting. Figure 1 shows the convergence speed of our algorithms and DP-GD in terms of real wall time for the task of logistic regression on the Covertype dataset for $(\epsilon, \delta) = (1, (\text{num. samples})^{-2})$ -DP. Despite DP-GD having a lower per-iteration cost, our algorithm is $30\times$ faster than DP-GD and achieves better excess loss.

Stochastic Minibatch Variant. We also show that our algorithms naturally extend to the minibatch setting where gradient and second-order information are computed on a subset of samples. We numerically compare it with DP-SGD and show that it has faster convergence.

2 Related Work

DP optimization is a well-studied topic [e.g., SCS13; MRTZ17; ACGM+16; STU17; WLKC+17; INST+19; STT20; SSTT21; GTU22; GLL22; BFTG19; BST14]. Most similar to our work, Avella-Medina, Bradshaw, and Loh [ABL21] consider second-order methods for DP convex optimization. We provide a detailed comparison between our results and theirs in Remark 4.5 and Section 6 showing that our algorithms relax restrictive assumptions and provide better excess error for logistic regression.

There are numerous non-private second-order optimization methods in the literature. The choice of method depends primarily on the values of n and d . When n is large, several works consider various sampling techniques for constructing second-order information, see [RM19; XYRRM16; Erd15; EM15]. When d is large, various methods are proposed in the literature for efficient approximation of the Hessian matrix, see [ABH17; Erd15; EM15; XYRRM16; GKL19]. There is also a family of algorithms based on the estimation of the curvature from the change in gradients. These algorithms are generally known as quasi-Newton methods stemming from the seminal BFGS algorithm [JM23].

3 Preliminaries

Let $d \in \mathbb{N}$. For a vector $x \in \mathbb{R}^d$, $\|x\|$ denotes the ℓ_2 norm of x . Let $n, m \in \mathbb{N}$. For a matrix $A \in \mathbb{R}^{n \times m}$, $\|A\| = \sup_{x \in \mathbb{R}^m: \|x\| \leq 1} \|Ax\|$ denotes the operator norm, and $\|A\|_F \triangleq \sqrt{\text{trace}(A^T \cdot A)}$ denotes the Frobenius norm of A where trace denotes the trace operator. $I_d \in \mathbb{R}^{d \times d}$ denotes the identity matrix. $\langle \cdot, \cdot \rangle$ denotes the standard inner product in \mathbb{R}^d . For a convex and closed subset $\mathcal{W} \subseteq \mathbb{R}^d$, let $\Pi_{\mathcal{W}} : \mathbb{R}^d \rightarrow \mathcal{W}$ be the Euclidean projection operator, given by $\Pi_{\mathcal{W}}(x) = \arg \min_{y \in \mathcal{W}} \|y - x\|_2$. For a (measurable) space \mathcal{R} , $\mathcal{M}_1(\mathcal{R})$ denotes the set of all probability measures on \mathcal{R} . Note that the statements in the paper about random variables hold almost surely. We will skip such declarations to aid readability. Let \mathcal{Z} be the data and let $\mathcal{W} \subseteq \mathbb{R}^d$ be the parameter space. Let $f : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}$ be a loss function. Throughout the paper, we assume f is doubly continuous, a convex function in w , and \mathcal{W} is a closed and convex set. We say (1) f is L_0 -Lipschitz iff there exists $L_0 \in \mathbb{R}$ such that $\forall z \in \mathcal{Z}, \forall w, v \in \mathcal{W} : |f(w, z) - f(v, z)| \leq L_0 \|w - v\|$, (2) f is L_1 -smooth iff there exists $L_1 \in \mathbb{R}$ such that $\forall z \in \mathcal{Z}, \forall w, v \in \mathcal{W} : \|\nabla f(w, z) - \nabla f(v, z)\| \leq L_1 \|w - v\|$, (3) f has a L_2 -Lipschitz Hessian iff there exists $L_2 \in \mathbb{R}$ such that $\forall z \in \mathcal{Z}, \forall w, v \in \mathcal{W} : \|\nabla^2 f(w, z) - \nabla^2 f(v, z)\| \leq L_2 \|w - v\|$, (4) f is μ -strongly convex iff for all $w, v \in \mathcal{W}$ and $z \in \mathcal{Z}$ we have $f(v, z) \geq f(w, z) + \langle \nabla f(w, z), v - w \rangle + \frac{\mu}{2} \|v - w\|^2$.

3.1 Zero-Concentrated DP

For our privacy analysis, we use concentrated differential privacy [DR16; BS16], as it provides a simpler composition theorem – the privacy parameter ρ adds up when we compose.

Definition 3.1 ([BS16, Def. 1.1]). A randomized mechanism $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{M}_1(\mathcal{R})$ is ρ -zCDP, iff, for every neighbouring dataset (i.e., addition or removal) $S_n \in \mathcal{Z}^n$ and $S'_n \in \mathcal{Z}^n$, and for every $\alpha \in (1, \infty)$, it holds $D_\alpha(\mathcal{A}(S_n) \parallel \mathcal{A}(S'_n)) \leq \rho\alpha$, where $D_\alpha(\mathcal{A}_n(S_n) \parallel \mathcal{A}_n(S'_n))$ is the α -Renyi divergence between $\mathcal{A}_n(S_n)$ and $\mathcal{A}_n(S'_n)$.

We should think of $\rho \approx \varepsilon^2$: to attain (ε, δ) -DP, it suffices to set $\rho = \frac{\varepsilon^2}{4 \log(1/\delta) + 4\varepsilon}$ [BS16, Lem. 3.5].

Lemma 3.2 ([BS16, Prop. 1.3]). Assume we have a randomized mechanism $\mathcal{A} : \mathcal{Z} \rightarrow \mathcal{M}_1(\mathcal{R})$ that satisfies ρ -zCDP, then for every $\delta > 0$, \mathcal{A} is $(\rho + 2\sqrt{\rho \log(1/\delta)}, \delta)$ -DP.

4 Optimal Algorithm for the Class of Strongly Convex Functions

In this section, we present a DP variant of the cubic-regularized Newton's method of Nesterov and Polyak [NP06]. To motivate the idea behind our algorithm, we revisit DP gradient descent (DP-GD) for the class of L_0 -Lipschitz and L_1 -smooth convex loss functions.

Let $\{w_t^{\text{GD}}\}_{t \in [T]}$ be the iterates of DP-GD. The smoothness of ℓ lets us construct a global quadratic upper bound on the function [Nes98, Thm. 2.1.5] as follows $\forall w \in \mathcal{W}$ and $S_n \in \mathcal{Z}^n$:

$$\ell(w, S_n) \leq q_t(w) \triangleq \ell(w_t^{\text{GD}}, S_n) + \langle \nabla \ell(w_t^{\text{GD}}, S_n), w - w_t^{\text{GD}} \rangle + \frac{L_1}{2} \|w - w_t^{\text{GD}}\|^2. \quad (2)$$

Then, DP-GD can be seen as a two-step process:

$$(\text{Step I}) \quad v_{t+1} = \arg \min_v q_t(v) = w_t^{\text{GD}} - L_1^{-1} \nabla \ell(w_t^{\text{GD}}, S_n), \quad (\text{Step II}) \quad w_{t+1}^{\text{GD}} = \Pi_{\mathcal{W}}(v_{t+1} + L_1^{-1} \xi_t),$$

where $\xi_t = \mathcal{N}(0, \sigma^2 I_d)$ with $\sigma^2 = \frac{L_0^2}{2\rho n^2}$ so that w_{t+1}^{GD} satisfies ρ -zCDP [BS16, Lem. 2.5]. That is, in each iteration of DP-GD, we find a minimum of the quadratic upper bound $q_t(w)$ and then project back to \mathcal{W} . (In the unconstrained setting where $\mathcal{W} = \mathbb{R}^d$ we do not need the second projection step.)

Consider the class of L_2 -Lipschitz Hessian convex loss functions. Nesterov and Polyak [NP06, Lem. 1] show that we can construct a *global cubic upper bound* exploiting the second-order information (i.e., Hessian) as follows: for all w and w_t , $\ell(w, S_n) \leq \phi_t(w)$ where

$$\phi_t(w) \triangleq \ell(w_t, S_n) + \langle \nabla \ell(w_t, S_n), w - w_t \rangle + \frac{1}{2} \langle \nabla^2 \ell(w_t, S_n)(w - w_t), w - w_t \rangle + \frac{L_2}{6} \|w - w_t\|^3. \quad (3)$$

Their non-private algorithm is based on the *exact* minimization of $\phi_t(w)$, i.e., the next iterate is $w_{t+1} = \arg \min \phi_t(w)$. Note that $\arg \min \phi_t(w)$ does not admit a closed form solution, as opposed to the quadratic upper bound (2). Similar to the intuition for DP-GD on smooth loss functions (2), our algorithms in this section are based on *privately* minimizing $\phi_t(w)$ at each iteration. Our algorithm is shown in Algorithm 1. In each iteration the algorithm makes an oracle call to obtain $(\ell(w_t, S_n), \nabla \ell(w_t, S_n), \nabla^2 \ell(w_t, S_n))$. Then, the algorithm calls an efficient DPSolver for privately optimizing the cubic upper bound (3). The privacy analysis of Algorithm 1 is a direct application of the composition property of zCDP [BS16, Lemma 2.3]; the output of DPSolver at each iteration satisfies ρ/T -zCDP where ρ is the total privacy budget and T is the total number of iterations.

Remark 4.1. DPSolver in Algorithm 1 does not affect the *oracle complexity* of Algorithm 1, as it is applied to the proxy loss $\phi_t(w)$, rather than the underlying loss $\ell(w, S_n)$. \triangleleft

Algorithm 1 Meta Algorithm

- 1: Input: training set $S_n \in \mathcal{Z}^n$, privacy budget ρ -zCDP, initialization $w_0 \in \mathcal{W}$, number of iterations T .
 - 2: **for** $t = 0, \dots, T - 1$ **do**
 - 3: Query $\ell(w_t, S_n), \nabla \ell(w_t, S_n), \nabla^2 \ell(w_t, S_n)$
 - 4: Construct $\phi_t(w)$ from Equation (3)
 - 5: $w_{t+1} = \text{DPSolver}(\phi_t(w), \rho/T, w_t)$
 - 6: Output w_T .
-

Algorithm 2 DPSolver

- 1: Input: function $\phi : \mathcal{W} \rightarrow \mathbb{R} : \phi(\theta) = \ell + \langle g, \theta - \theta_0 \rangle + \frac{1}{2} \langle H(\theta - \theta_0), (\theta - \theta_0) \rangle + \frac{L_2}{6} \|\theta - \theta_0\|^3$, privacy budget $\tilde{\rho}$ -zCDP, initialization θ_0 .
 - 2: $N = \frac{2\tilde{\rho}(L_0 + L_1 D + L_2 D^2)^2 n^2}{(L_0 + L_1 D)^2 d}, \sigma^2 = \frac{N(L_0 + L_1 D)^2}{2\tilde{\rho}}$
 - 3: **for** $i = 0, \dots, N - 1$ **do**
 - 4: $\eta_i = \frac{2}{\mu(i+2)}$
 - 5: $\text{grad}_i = g + H(\theta_i - \theta_0) + \frac{L_2}{2} \|\theta_i - \theta_0\| (\theta_i - \theta_0)$
 - 6: $\theta_{i+1} = \Pi_{\mathcal{W}}(\theta_i - \eta_i(\text{grad}_i + \mathcal{N}(0, \sigma^2 I_d)))$
 - 7: Return $\sum_{i=0}^{N-1} \frac{2i}{N(N+1)} \theta_i$
-

Theorem 4.2. Let f be a L_0 -Lipschitz, L_1 -smooth, L_2 -Lipschitz Hessian, and μ -strongly convex function. Also, assume that $\mathcal{W} \subseteq \mathbb{R}^d$ has finite diameter D . Let $w^* = \arg \min_{w \in \mathcal{W}} \ell(w, S_n)$. Then, for every $\rho > 0$, $\beta \in (0, 1)$, and $S_n \in \mathcal{Z}^n$ for sufficiently large n , by setting the number of iterations in Algorithm 1 to

$$T = \Theta\left(\frac{\sqrt{L_2}}{\mu^{3/4}}(\ell(w_0, S_n) - \ell(w^*, S_n))^{\frac{1}{4}} + \log \log\left(\frac{n\sqrt{\rho}}{\sqrt{\log(1/\beta)d}}\right)\right),$$

and using Algorithm 2 as DPSolver, we have the following: The output of Algorithm 1, i.e., w_T , satisfies ρ -zCDP and with probability at least $1 - \beta$

$$\ell(w_T, S_n) - \ell(w^*, S_n) \leq \tilde{O}\left(\frac{d(L_0 + L_1 D)^2 \log(1/\beta)}{\mu \rho n^2} \cdot \left(\frac{L_2^2 L_0 D}{\mu^3}\right)^{\frac{1}{4}}\right)$$

Remark 4.3. The lower bound on the excess error of any DP algorithm for the class of strongly convex functions [BST14, Thm. 5.5] implies that the achievable excess error in Theorem 4.2 is *optimal* in terms of the dependence on d , ρ , and n . Also, the oracle complexity of our algorithm is an exponential improvement over the oracle complexity of first-order methods [STU17]. \triangleleft

Remark 4.4. The proof of Theorem 4.2 suggests that Algorithm 1 has two phases. First, while w_t is far from w^* , the convergence rate is $1/T^4$. Second, when w_t is close to w^* , the algorithm exhibits the convergence rate of $\exp(\exp(-T))$. Notice that Algorithm 1 is agnostic to this transition in the sense that we do not have an explicit switching step in Algorithm 1 and Algorithm 2. It is also interesting to note that the transition happens when $\|w_t - w^*\| \leq 3\mu/4L_2$. \triangleleft

Remark 4.5 (Comparison with [ABL21]). In [ABL21, §4], the authors propose a DP variant of Newton’s method. Their main idea is to add independent noise *directly* to the Hessian matrix and the gradient vector using the Gaussian mechanism. They also require that *the Hessian be a rank-1 matrix*. The issue with adding noise directly to a full-rank Hessian matrix is that the noise scales with the dimension d , which can lead to a suboptimal excess loss. In contrast, our algorithm has a global convergence without placing restrictions on the rank of the Hessian matrix or the initialization. \triangleleft

Remark 4.6. We showed in Theorem 4.2 that our algorithm has an exponentially smaller *oracle complexity* than the first-order methods in terms of the dependence to n . For the class of convex, smooth, Lipschitz, and strongly convex, [ZZMW17] proposes a first-order algorithm with an oracle complexity of $T_1 = \Theta(\sqrt{L_1}/\sqrt{\mu} + \log(n))$. It is important to note that the *constant* term in T_1 differs from our result, making a direct comparison challenging. It is an interesting question to develop a second-order DP algorithm with a smaller oracle complexity than both the algorithms proposed in [ZZMW17] and ours in Algorithm 1. \triangleleft

Remark 4.7. The cubic Newton method has a non-private convergence rate of T^{-2} for the class of convex (but not strongly convex) functions [NP06, Thm. 4]. We leave it as an open question whether there exists a DPSolver such that Algorithm 1 achieves an optimal excess error and oracle complexity for convex functions. However, this can be achieved by a DP variant of the first-order accelerated Nesterov’s method [Nes98; NJLS09; GL12]; see Appendix A.2. \triangleleft

5 DP Logistic Regression using Second-Order Information

The main limitation of our cubic Newton’s method (Algorithm 1) is that each iteration requires solving a nontrivial subproblem. So, despite low oracle complexity, it is computationally expensive. Moreover, many loss functions, such as logistic loss, are not strongly convex in the unconstrained setting. In this section, we aim to develop a fast second-order algorithm for unconstrained logistic regression avoiding this issue. In many real-world classification tasks, the logistic loss is the loss of choice. The logistic loss is a convex surrogate of the 0-1 loss, and satisfies many regularity conditions that give rise to various practical optimization algorithms [Bac10; Erd15; KSJ18]. Also, note that our results in this section can readily be extended to the class of smooth and convex GLMs.

First, we recall the logistic loss function. Let $d \in \mathbb{N}$ and $\mathcal{Z} = \mathcal{B}^d(1) \times \{-1, 1\}$ be the dimension and data space, where $\mathcal{B}^d(1) = \{x \in \mathbb{R}^d : \|x\| \leq 1\}$ is the unit ball in \mathbb{R}^d . Let $f_{\text{LL}} : \mathbb{R}^d \times \mathcal{Z} \rightarrow \mathbb{R}$ denote the logistic loss function defined as

$$f_{\text{LL}}(w, (x, y)) = \log(1 + \exp(-y \cdot \langle w, x \rangle)). \quad (4)$$

The gradient and Hessian of f_{LL} are given by

$$\nabla_w f_{\text{LL}}(w, (x, y)) = \frac{-xy}{1 + \exp(y \langle w, x \rangle)}, \quad \nabla_w^2 f_{\text{LL}}(w, (x, y)) = \frac{xx^\top}{(\exp(-\frac{\langle w, x \rangle}{2}) + \exp(\frac{\langle w, x \rangle}{2}))^2}. \quad (5)$$

Newton’s method [BV04, §9.5] is based on successively minimizing a *local* second-order Taylor approximation on the function. Newton’s method does not guarantee a global convergence [JT16]; the reason is that the second-order Taylor approximation of the logistic loss can greatly underestimate the function. Next we show that it is possible to obtain a quadratic *global upper bound* on the logistic loss function. We will use this to develop an algorithm that converges globally.

Lemma 5.1. *For every $v \in \mathbb{R}^d$, $x \in \mathbb{R}^d$, $w \in \mathbb{R}^d$, and $y \in \{-1, +1\}$, we have*

$$f_{\text{LL}}(w, (x, y)) \leq f_{\text{LL}}(v, (x, y)) + \langle \nabla f_{\text{LL}}(v, (x, y)), w - v \rangle + \frac{1}{2} \langle H_{\text{qu}}(v, (x, y))(w - v), w - v \rangle,$$

$$\text{where } H_{\text{qu}}(v, (x, y)) \triangleq \frac{\tanh(\langle x, v \rangle / 2)}{2 \langle x, v \rangle} x x^\top \in \mathbb{R}^{d \times d}.$$

Remark 5.2. Since f_{LL} is $\frac{1}{4}$ -smooth, we can construct a simpler global quadratic upper-bound as follows [Nes98, Thm. 2.1.5]: $f_{\text{LL}}(w, (x, y)) \leq f_{\text{LL}}(v, (x, y)) + \langle \nabla f_{\text{LL}}(v, (x, y)), w - v \rangle + \frac{1}{8} \|w - v\|^2$. Lemma 5.1 is tighter than this, since $H_{\text{qu}}(v, (x, y)) \preceq \frac{1}{4} I_d$; see Appendix B.2. \triangleleft

Remark 5.3. The second-order Taylor approximation and our upper bound in Lemma 5.1 both provide a quadratic approximation of the logistic loss. In the remainder of the paper, we write $H(v, (x, y))$ to refer to both $\nabla^2 f_{\text{LL}}(v, (x, y))$ and $H_{\text{qu}}(v, (x, y))$. We refer to $H(v, (x, y))$ as the second-order information (SOI) and to H_{qu} as *quadratic upperbound SOI*. Finally, notice both $\nabla^2 f_{\text{LL}}(v, (x, y))$ and $H_{\text{qu}}(v, (x, y))$ are PSD rank-1 matrices, with maximum eigenvalue $\leq \frac{1}{4} \|x\|^2 \leq \frac{1}{4}$. \triangleleft

5.1 Algorithm Description

We are given a dataset $S_n = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathcal{B}^d(1) \times \{-1, +1\})^n$ and we aim to minimize $\ell_{\text{LL}}(w, S_n) \triangleq \frac{1}{n} \sum_{i \in [n]} f_{\text{LL}}(w, (x_i, y_i))$. Our algorithm iteratively minimizes a quadratic approximation of $\ell_{\text{LL}}(w, S_n)$. Consider

$$q_t(w) \triangleq \ell_{\text{LL}}(w_t, S_n) + \langle \nabla \ell_{\text{LL}}(w_t, S_n), w - w_t \rangle + \frac{1}{2} \langle H(w_t, S_n)(w - w_t), (w - w_t) \rangle, \quad (6)$$

where $H(w_t, S_n) \triangleq \frac{1}{n} \sum_{i \in [n]} H(w_t, (x_i, y_i))$. In the non-private setting the next iterate is set to $w_{t+1} = \arg \min_w q_t(w) = w_t - H(w_t, S_n)^{-1} \nabla \ell_{\text{LL}}(w_t, S_n)$. To develop a private variant of Newton’s method, we need to characterize the sensitivity of this update rule. Our key observation is that *the directions corresponding to small eigenvalues of $H(w_t, S_n)$ are more sensitive than the directions corresponding to large eigenvalues*. To overcome this issue, we modify the eigenvalues of $H(w_t, S_n)$ to ensure a minimum eigenvalue $\geq \lambda_0$, where $\lambda_0 > 0$ is a carefully chosen constant. We show how to *adaptively* tune λ_0 in Section 5.2. This procedure yields the desired stability with respect to neighbouring datasets. Formally, the modification operator is defined as follows:

Definition 5.4. Let $A \in \mathbb{R}^{d \times d}$ be a positive semi-definite (PSD) matrix and $\lambda_0 \geq 0$. Define

$$\Psi_{\lambda_0}(A, \text{clip}) = \sum_{i=1}^d \max\{\lambda_0, \lambda_i\} u_i u_i^\top, \quad \Psi_{\lambda_0}(A, \text{add}) = \sum_{i=1}^d (\lambda_i + \lambda_0) u_i u_i^\top = A + \lambda_0 I_d.$$

where $A = \sum_{i=1}^d \lambda_i u_i u_i^\top$ is the eigendecomposition of A – i.e., $0 \leq \lambda_1 \leq \dots \leq \lambda_d$ are the eigenvalues and $u_1, \dots, u_d \in \mathbb{R}^d$ are the eigenvectors, which satisfy $\forall i \neq j \ \|u_i\| = 1 \wedge \langle u_i, u_j \rangle = 0$.

Algorithm 3 describes our algorithm. First, we state the privacy guarantee of Algorithm 3 whose proof can be found in Appendices B.3 and B.4.

Theorem 5.5. *Assume in Algorithm 3 we choose add for the SOI modification. Then, for every training set $S_n \in (\mathbb{R}^d \times \{-1, +1\})^n$, $w_0 \in \mathcal{W}$, $\lambda_0 > 0$, $T \in \mathbb{N}$, $\rho \in \mathbb{R}_+$, and $\theta \in (0, 1)$, by setting $\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}$ and $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 + \lambda_0)\sqrt{2\rho\theta}}$, w_T satisfies ρ -zCDP.*

Theorem 5.6. *Assume in Algorithm 3, we choose clip for the SOI modification. Then, for every training set $S_n \in (\mathbb{R}^d \times \{-1, +1\})^n$, $w_0 \in \mathcal{W}$, $\lambda_0 > 0$, $T \in \mathbb{N}$, $\rho \in \mathbb{R}_+$, and $\theta \in (0, 1)$ such that $n > \frac{1}{4\lambda_0}$, by setting $\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}$ and $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 - \lambda_0)\sqrt{2\rho\theta}}$, w_T satisfies ρ -zCDP.*

Algorithm 3 Newton Method with Double noise

```

1: Inputs: training set  $S_n \in \mathcal{Z}^n$ ,  $\lambda_0 > 0$ ,  $\theta \in (0, 1)$ ,
   privacy budget  $\rho$ -zCDP, initialization  $w_0$ , number of
   iterations  $T$ , SOI modification  $\in \{\text{clip}, \text{add}\}$ .
2: Set  $\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}$ 
3: if SOI modification = Add then
4:    $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 + \lambda_0)\sqrt{2\rho\theta}}$ 
5: else if SOI modification = Clip then
6:    $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 - \lambda_0)\sqrt{2\rho\theta}}$ 
7: for  $t = 0, \dots, T - 1$  do
8:   Query  $\nabla f(w_t, S_n)$  and  $H(w_t, S_n)$ 
9:    $\tilde{H}_t = \Psi_{\lambda_0}(H(w_t, S_n), \text{SOI modification})$ 
10:   $\tilde{g}_t = \nabla f_{\text{LL}}(w_t, S_n) + \mathcal{N}(0, \sigma_1^2 I_d)$ 
11:   $w_{t+1} = w_t - \tilde{H}_t^{-1} \tilde{g}_t + \mathcal{N}(0, \|\tilde{g}_t\|^2 \sigma_2^2 I_d)$ 
12: Output  $w_T$ .
```

Our DP algorithm differs from the non-private Newton’s method in three ways: (1) We first privatize the gradient by adding noise. (2) We modify $H(w_t, S_n)$ to ensure its eigenvalues are not too small. And (3) we add a second noise to the update computed using the noised gradient and modified second-order information (SOI).

Notice that Algorithm 3 has *four variations* based on the SOI and the modification of SOI, namely, Hess-clip, Hess-add, QU-clip, and QU-add which refer to using Hessian and clip, Hessian and add, quadratic upper bound (See Lemma 5.1) and clip, and quadratic upper bound and add, respectively.

Remark 5.7 (Generalization of Algorithm 3). In this section our main focus is on DP logistic regression, and the privacy guarantees hold for

the logistic loss. Nevertheless, in Appendix B.6, we present a generalization of Algorithm 3 whose privacy guarantee holds for *every* convex, doubly differentiable, Lipschitz, and smooth loss function *without any constraints on the rank of Hessian*. The main technical challenge for sensitivity analysis is proving the approximate Lipschitzness of Ψ in the operator norm (See Lemma B.7). This demonstrates that our algorithm is more general than objective perturbation [CMS11; KST12; INST+19] and the private damped Newton’s method [ABL21] which both require a low-rank Hessian. \triangleleft

5.2 Private and Adaptive Selection of Minimum Eigenvalue

One of the hyperparameters of Algorithm 3 is the minimum eigenvalue λ_0 . There exists a tradeoff for choosing λ_0 . We ideally want the modification to be as small as possible, so that the SOI is preserved. However, decreasing λ_0 increases σ_2 and we add more noise. To deal with this problem, we propose a heuristic rule for an adaptive, private, and time-varying selection of the minimum eigenvalue. We wish to find $\lambda_{0,t}$ that minimizes expected loss at the next iteration, for which we have the quadratic approximation (6). More formally, we compute $\lambda_{0,t}$ as $\arg \min_{\lambda} \mathbb{E}[q_t(w_t - \Psi_{\lambda}(H(w_t, S_n), \text{SOI modification})\tilde{g}_t + \|\tilde{g}_t\| \sigma_2(\lambda) \cdot \xi)]$ where q_t is given in (6) and $\xi \sim \mathcal{N}(0, I_d)$. We show in Appendix B.5 that an approximate minimizer is $\lambda_{0,t} \propto \left(\frac{\text{trace}(H_t(w_t, S_n))}{n^2 \times \text{privacy budget for the direction}} \right)^{\frac{1}{3}}$. Note that $\lambda_{0,t}$ depends on the data through $\text{trace}(H(w_t, S_n))$, which has sensitivity $1/4n$, so it can be estimated privately. In Appendix B.5, we provide the algorithmic description of a variant of Algorithm 3 with an adaptive and private minimum eigenvalue. In particular, we divide the privacy budget at each iteration into three parts: (1) privatizing the gradient; (2) estimating the trace of SOI; and (3) privatizing the direction. We use this variant for our numerical experiments in Section 6.

5.3 Convergence Results for Algorithm 3

In this section, we provide data-dependent convergence guarantees for Algorithm 3. We express these guarantees in terms of the conditional expectation $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \{w_i\}_{i \in [t]}]$ and they can be easily extended to obtain high probability bounds. Before presenting the results, we introduce a notation. For a dataset $S_n = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathbb{R}^d \times \{-1, +1\})^n$, let $V \in \mathbb{R}^{d \times d}$ denote the *orthogonal projection matrix* on the linear subspace spanned by $\{x_1, \dots, x_n\}$. For every vector $u \in \mathbb{R}^d$, define $\|u\|_V \triangleq \sqrt{u^\top V u}$. This norm naturally arises since for every $w \in \mathbb{R}^d$ we have $\ell_{\text{LL}}(w, S_n) - \ell_{\text{LL}}(w^*, S_n) \leq \frac{1}{8} \|w - w^*\|_V^2$ where $w^* = \arg \min \ell_{\text{LL}}(w, S_n)$ (See Appendix B.7).

5.3.1 Local Convergence Guarantee of Hess-clip and Hess-add

Theorem 5.8. *Let S_n denote the dataset and rank denote the dimension of the linear subspace spanned by $\{x_1, \dots, x_n\}$. Let $\lambda_{\min,t}$ be the smallest non-zero eigenvalue of $\nabla^2 \ell_{\text{LL}}(w_t, S_n)$ and ρ be*

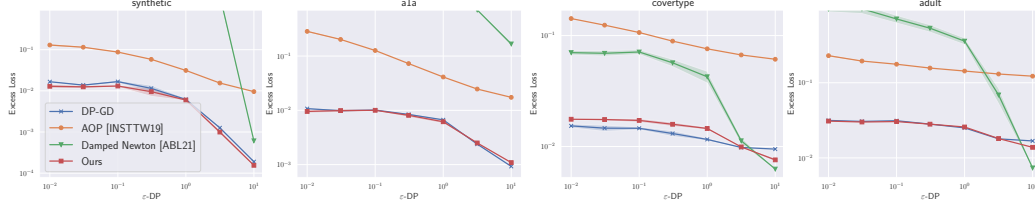


Figure 2: Privacy-Utility tradeoff on different datasets.

the privacy budget (in ϵ CDP) per iteration. Then,

$$\mathbb{E}_t \left[\|w_{t+1} - w^*\|_V^2 \right] \leq \nu_{1,t}^2 \|w_t - w^*\|_V^2 + 2\nu_{1,t}\nu_{2,t} \|w_t - w^*\|_V^3 + \nu_{2,t}^2 \|w_t - w^*\|_V^4 + \Delta,$$

where the coefficients are given by

$$\nu_{1,t} = 1 - \frac{\tilde{\lambda}_{\min,t}}{\lambda_0} + \frac{\sqrt{\text{rank}}}{(4n\lambda_0^2 - \lambda_0)\sqrt{2\rho\theta}}, \quad \nu_{2,t} = \frac{0.05}{\tilde{\lambda}_{\min,t}}, \quad \Delta = O\left(\frac{\text{rank}}{\rho(1-\theta)n^2} \frac{1}{(\tilde{\lambda}_{\min,t})^2}\right). \quad (7)$$

Here, $\tilde{\lambda}_{\min,t} = \begin{cases} \min\{\lambda_{\min,t}, \lambda_0\} & \text{for Hess-clip,} \\ \lambda_{\min,t} + \lambda_0 & \text{for Hess-add,} \end{cases}$ depends on the modification procedure.

This type of convergence is known as *composite convergence*, as it is a combination of linear and quadratic rates, and has been observed in the convergence analysis of several quasi-Newton's methods [EM15; Erd15; RM16; XYRRM16].

Remark 5.9. $\lambda_{\min,t}$ is the smallest *non-zero* eigenvalue of $\nabla^2 \ell_{\text{LL}}(w_t, S_n)$. Therefore, for sufficiently large n we have $0 < \nu_{1,t} < 1$. It shows Algorithm 3 with Hessian as SOI is, in-expectation, a descent algorithm locally given $\|w_t - w^*\|$ is sufficiently larger than Δ . Roughly speaking, Theorem 5.8 guarantees a linear convergence to a ball around the optimum whose radius is given by Δ . We also observe the linear rate in Figure 3. Moreover, the error due to the privacy, i.e., Δ in Equation (7), is proportional to the rank of the feature vectors which is always smaller than d . These interesting properties is due to the convergence analysis with respect to $\|\cdot\|_V$. \triangleleft

Remark 5.10. The coefficients of the convergence in Equation (7) depend on the iteration step which is an undesirable aspect of the results. In Lemma B.11, we prove that $|\lambda_{\min,t} - \lambda_{\min}^*| \leq 0.1 \|w_t - w^*\|_V$ where λ_{\min}^* is the smallest non-zero eigenvalue of $\nabla^2 \ell_{\text{LL}}(w^*, S_n)$. Therefore, the coefficients can be well-approximated by their analogous values evaluated at the optimum. \triangleleft

5.3.2 Global Convergence Guarantee of QU-clip and QU-add

We also establish a global convergence guarantee for QU-clip and QU-add. Due to the space the formal statement and proof are deferred to Appendix B.9. Roughly speaking, under the assumption of *local strong convexity at the optimum* [Bac14], QU-clip and QU-add converge globally: this is intuitive since QU-clip and QU-add are based on minimizing a global upper bound on the function.

6 Numerical Results

In this section, we evaluate the performance of our algorithm (Algorithm 3 with the adaptive minimum eigenvalue selection from Section 5.2) for the problem of *binary classification* using *logistic regression*. For brevity, many of the details behind our implementation and more experimental results are deferred to Appendix C.

6.1 Setup

The setup of the experiments is as follows: **Baseline1- DP-(S)GD:** The update rule is $w_{t+1} = w_t - \eta \nabla \ell(w_t, S_n) + \xi$ where ξ is a Gaussian noise [SCS13; BST14; ACGM+16]. Since the logistic loss is 1-Lipschitz, we do not need gradient clipping. The Lipschitzness parameter controls the variance of the Gaussian random vector. To draw a fair comparison and show the advantage of using second-order information, we chose the stepsize to be equal to the inverse smoothness. **Baseline2-**

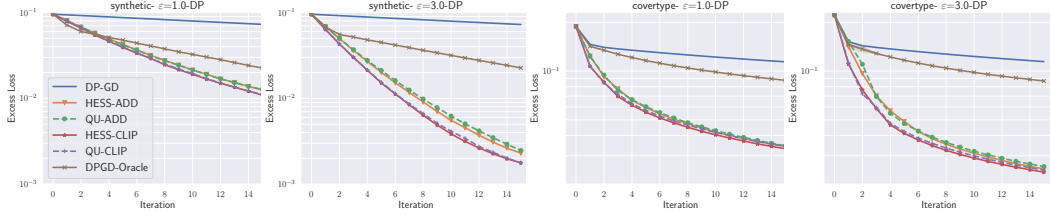


Figure 3: Comparison with DP-GD Oracle where at each iteration the stepsize tuned non-privately.

	$\frac{T_{\text{DP-GD}}^*}{T_{\text{ours}}^*}$				$T_{\text{ours}}^*(\text{sec})$	
	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 1$	$\epsilon = 10$	$\min(T_{\text{ours}}^*)(\text{sec.})$	$\max(T_{\text{ours}}^*)(\text{sec.})$
a1a	$4.87 \times$	$2.95 \times$	$5.09 \times$	$30.59 \times$	2.45	4.2
synthetic	$2.90 \times$	$2.90 \times$	$5.19 \times$	$11.61 \times$	0.18	0.21
adult	$12.08 \times$	$11.84 \times$	$22.17 \times$	$38.16 \times$	6.81	8.07
coartype	$24.19 \times$	$19.85 \times$	$35.70 \times$	$36.20 \times$	2.93	3.58

Table 1: Comparison between the run time of our algorithm and DP-GD in terms of the ratio $T_{\text{DP-GD}}^*/T_{\text{ours}}^*$. The last two columns show the minimum and maximum run time of our algorithm.

Approximate Objective Perturbation (AOP): AOP is built on objective perturbation [CMS11; KST12]. Objective perturbation consists of a two-stage process: (1) *perturbing* the objective function by adding a random linear term and (2) outputting the minimum of the perturbed objective. Releasing such a minimum is sufficient for achieving DP guarantees [CMS11; KST12], but only if we can find the exact minimum of the perturbed objective. AOP extends objective perturbation to permit using an *approximate* minimum of the perturbed objective [INST+19; INST+]. Notice AOP is not an iterative optimization algorithm. **Baseline3- Damped Newton Method [ABL21]:** The algorithm in [ABL21] is a variant of damped Newton’s method with the assumption that the Hessian of loss function is rank-1, which holds for the logistic loss. Their algorithm is based on adding two i.i.d. noises to the Hessian and the gradient: $w_{t+1} = w_t - \eta_t H_{\text{noisy},t}(w_t, S_n)^{-1} \tilde{g}_t$, where η_t is the stepsize, $H_{\text{noisy},t}(w_t, S_n) = \nabla^2 \ell_{\text{LL}}(w_t, S_n) + \Xi_t$ and $\tilde{g}_t = \nabla \ell_{\text{LL}}(w_t, S_n) + \xi_t$. Here Ξ_t and ξ_t are carefully chosen Gaussian noise. With $\eta_t = 1$, our experiments show that their algorithm is not converging. We use the strategy suggested in [ABL21, Page 22] and set $\eta_t = \log(1 + \beta_t)/\beta_t$ where $\beta_t = \|\nabla^2 \ell_{\text{LL}}(w_t, S_n)^{-1} \nabla \ell_{\text{LL}}(w_t, S_n)\|$. This stepsize selection makes the algorithm *non-private*, however, it serves as a good baseline. **Datasets:** We conducted experiments on six publicly available datasets: a1a, Adult, coartype, synthetic, fashion-MNIST, and protein (Appendix C includes fashion-MNIST and protein results). The synthetic dataset is generated as follows: Fix $d \in \mathbb{N}$ and $w^* \in \mathbb{R}^d$. Then, (1) the feature vectors $\{x_i \in \mathbb{R}^d : i \in [n]\}$ are independent and sampled uniformly at random from the unit sphere in \mathbb{R}^d , (2) for the i -th datapoint the label is $+1$ with probability $(1 + \exp(-\langle x_i, w^* \rangle))^{-1}$ and -1 otherwise. **Privacy Notion:** The privacy notion for our experiments is $(\epsilon, \delta = (\text{num. of samples})^{-2})$ -DP. Next, we present the results.

6.2 Privacy-Utility-Run Time Tradeoff

We study the tradeoff for our algorithm and compare it with other baselines for a broad range of $\epsilon \in \{0.01, \dots, 10\}$. We *non-privately tune* the total number of iterations of the iterative algorithms and report the best achievable excess error in Figure 2. As can be seen our algorithm almost always achieves the best excess loss for a broad range of ϵ . Also, Figure 2 shows that damped private Newton method of [ABL21] achieves a low excess loss only for large ϵ . Figure 2 indicates that DP-GD and our algorithm are the best in terms of excess loss. In Table 1, we compare the run time of DP-GD and our algorithm, i.e., the computational time in seconds for achieving the excess loss in Figure 2. As can be seen, for many challenging datasets, our algorithm is 10 - $40 \times$ faster than DP-GD. Our experiments are run on CPU. We also remark that each step of Algorithm 3, i.e., computing gradient and SOI, is heavily parallelizable implying that the run time of Algorithm 3 can be made much smaller by an efficient implementation. Also, the reported numbers in Figure 2 and Table 1 correspond to Hess-clip.

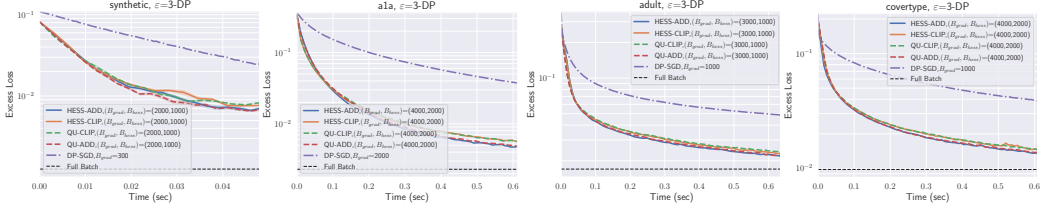


Figure 4: Minibatch Variant of Our Algorithm and Comparison with DP-SGD

6.3 Second Order Information vs Optimal Stepsize

In non-private convex optimization, the key to the success of second-order optimization algorithms is that the second-order information acts as a preconditioner, and the same performance *cannot* be attained by optimally tuning the stepsize for GD algorithm. To investigate whether the same holds for our algorithms, we consider the following variant of DP-GD. Let \tilde{g}_t denote the perturbed gradient obtained by adding a Gaussian random vector to $\nabla \ell_{LL}(w_t, S_n)$. Instead of a constant stepsize, the stepsize at iteration t is chosen based on $\eta_t = \arg \min_{\eta \geq 0} \ell_{LL}(w_t - \eta \tilde{g}_t)$. Notice this variant is obviously *not DP*. We refer to this variant as *DP-GD-Oracle*. The comparison with DP-GD-Oracle lets us answer the following question: *Could we have just computed a single number, i.e., stepsize, to achieve the same performance as our second-order optimization algorithms which require computing a $d \times d$ matrix?* In Figure 3, we compare the convergence speed of our algorithms with DP-GD-Oracle in low- and high-privacy regimes. Figure 3 shows our algorithms converge faster than DP-GD-Oracle which is not even a DP algorithm. Figure 3 confirms the expectation that as the privacy budget increases the difference between our algorithms and DP-GD-Oracle increases since we can use more curvature information.

6.4 Minibatch Variant of Our Algorithm and Comparison with DP-SGD

So far we have considered full-batch algorithms that compute first- and second-order information on the entire dataset. We extend Algorithm 3 to the minibatch setting, where, at each iteration, the gradient and SOI matrix are computed using a subsample of the data points. In Appendix C.1 we provide a formal algorithmic description of the minibatch version of Algorithm 3 along with its privacy proof. Then, we compare the convergence speed and excess loss with DP-SGD.

DP-SGD is faster than DP-GD, but to achieve good privacy and utility, we need large batches [PHKX+23, Fig. 2]. This is in stark contrast with non-private SGD, where larger batch sizes yield diminishing returns [ZLNM+19]. In particular, to achieve the best excess loss we need to select the batch size as large as possible. We select the batch size of DP-SGD so that the achievable excess loss will be close to the full batch versions. Specifically, we select $\frac{\text{batch size DP-SGD}}{\text{number of samples}} \approx 0.02$ and tune the number of iterations of DP-SGD to obtain the best result. Figure 4 shows the progress of different algorithm versus run time. Obviously, for a fixed run time DP-SGD performs more iterations compared to our algorithms. Nevertheless, our algorithms achieve the same excess error as DP-GD with 8-10 \times faster run time over all the datasets while *the batch sizes of our algorithms are larger than that of DP-SGD*. We observe that the variations of our algorithms based on the adding operator performs better in the minibatch setting. This can be attributed to the smaller σ_2 for the adding operator in Algorithm 3. In summary, the comparison between privacy-utility-wall time tradeoff of the subsampled variant of our algorithm and DP-SGD is similar to their full-batch counterparts.

7 Conclusion and Limitations

We showed that second-order methods can be used in the DP setting both for improving worst-case convergence guarantees and designing faster practical algorithms. We believe our results open up many directions: A limitation of our algorithms is that the cost of forming and inverting the Hessian can be prohibitive when d is large. In the non-private setting, a line of research tries to address this limitation by constructing an approximation to SOI such that the update is efficient, yet still provides sufficient SOI [EM15; Erd15; XYRRM16; ABH17]. It would be interesting to investigate how the ideas developed in our paper could be incorporated into these methods.

Acknowledgments

The authors would like to thank Murat Erdogdu, Jalaj Upadhyay, and Mohammad Yaghini for helpful discussions. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute www.vectorinstitute.ai/partners.

References

- [ABH17] N. Agarwal, B. Bullins, and E. Hazan. “Second-order stochastic optimization for machine learning in linear time”. *The Journal of Machine Learning Research* 18.1 (2017), pp. 4148–4187.
- [ABL21] M. Avella-Medina, C. Bradshaw, and P.-L. Loh. “Differentially private inference via noisy optimization”. *arXiv preprint arXiv:2103.11003* (2021).
- [ACGM+16] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, et al. “Deep learning with differential privacy”. In: *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*. 2016, pp. 308–318.
- [AMN20] J. Arbel, O. Marchal, and H. D. Nguyen. “On strict sub-Gaussianity, optimal proxy variance and symmetry for bounded random variables”. *ESAIM: Probability and Statistics* 24 (2020), pp. 39–55.
- [Bac10] F. Bach. “Self-concordant analysis for logistic regression”. *Electronic Journal of Statistics* 4 (2010), pp. 384–414.
- [Bac14] F. Bach. “Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression”. *The Journal of Machine Learning Research* 15.1 (2014), pp. 595–627.
- [BD99] J. A. Blackard and D. J. Dean. “Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables”. *Computers and electronics in agriculture* 24.3 (1999), pp. 131–151.
- [BFTG19] R. Bassily, V. Feldman, K. Talwar, and A. Guha Thakurta. “Private stochastic convex optimization with optimal rates”. *Advances in neural information processing systems* 32 (2019).
- [BS16] M. Bun and T. Steinke. “Concentrated differential privacy: Simplifications, extensions, and lower bounds”. In: *Theory of Cryptography Conference*. Springer. 2016, pp. 635–658.
- [BST14] R. Bassily, A. Smith, and A. Thakurta. “Private empirical risk minimization: Efficient algorithms and tight error bounds”. In: *2014 IEEE 55th annual symposium on foundations of computer science*. IEEE. 2014, pp. 464–473.
- [BV04] S. P. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [CJB04] R. Caruana, T. Joachims, and L. Backstrom. “KDD-Cup 2004: results and analysis”. *ACM SIGKDD Explorations Newsletter* 6.2 (2004), pp. 95–108.
- [CMS11] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. “Differentially Private Empirical Risk Minimization”. *Journal of Machine Learning Research* 12.29 (2011), pp. 1069–1109.
- [DG17] D. Dua and C. Graff. *UCI Machine Learning Repository*. 2017.
- [DMNS06] C. Dwork, F. McSherry, K. Nissim, and A. Smith. “Calibrating noise to sensitivity in private data analysis”. In: *Theory of cryptography conference*. Springer. 2006, pp. 265–284.
- [DR16] C. Dwork and G. N. Rothblum. “Concentrated differential privacy”. *arXiv preprint arXiv:1603.01887* (2016).
- [EM15] M. A. Erdogdu and A. Montanari. “Convergence rates of sub-sampled Newton methods”. *Advances in Neural Information Processing Systems* 28 (2015).

- [Erd15] M. A. Erdogdu. “Newton-Stein method: A second order method for GLMs via Stein’s lemma”. *Advances in Neural Information Processing Systems* 28 (2015).
- [GKLR19] R. Gower, D. Kovalev, F. Lieder, and P. Richtárik. “RSN: randomized subspace Newton”. *Advances in Neural Information Processing Systems* 32 (2019).
- [GL12] S. Ghadimi and G. Lan. “Optimal Stochastic Approximation Algorithms for Strongly Convex Stochastic Composite Optimization I: A Generic Algorithmic Framework”. *SIAM Journal on Optimization* 22.4 (2012), pp. 1469–1492. eprint: <https://doi.org/10.1137/110848864>.
- [GLL22] S. Gopi, Y. T. Lee, and D. Liu. “Private Convex Optimization via Exponential Mechanism”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, Feb. 2022, pp. 1948–1989.
- [GTU22] A. Ganesh, A. Thakurta, and J. Upadhyay. “Langevin diffusion: An almost universal algorithm for private Euclidean (convex) optimization”. *arXiv preprint arXiv:2204.01585* (2022).
- [GV13] G. H. Golub and C. F. Van Loan. *Matrix computations*. JHU press, 2013.
- [HLR19] N. J. Harvey, C. Liaw, and S. Randhawa. “Simple and optimal high-probability bounds for strongly-convex stochastic gradient descent”. *arXiv preprint arXiv:1909.00843* (2019).
- [INST+] R. Iyengar, J. P. Near, D. Song, O. Thakkar, et al. *Differentially Private Convex Optimization Benchmark*. URL: <https://github.com/sunblaze-ucb/dpml-benchmark>.
- [INST+19] R. Iyengar, J. P. Near, D. Song, O. Thakkar, et al. “Towards practical differentially private convex optimization”. In: *2019 IEEE Symposium on Security and Privacy (SP)*. 2019.
- [JM23] Q. Jin and A. Mokhtari. “Non-asymptotic superlinear convergence of standard quasi-Newton methods”. *Mathematical Programming* 200.1 (2023), pp. 425–473.
- [JNGKJ19] C. Jin, P. Netrapalli, R. Ge, S. M. Kakade, and M. I. Jordan. “A short note on concentration inequalities for random vectors with subgaussian norm”. *arXiv preprint arXiv:1902.03736* (2019).
- [JT16] F. Jarre and P. L. Toint. “Simple examples for the failure of Newton’s method with line search for strictly convex minimization”. *Mathematical Programming* 158.1 (2016), pp. 23–34.
- [Kat73] T. Kato. “Continuity of the map $S \rightarrow |S|$ for linear operators”. *Proceedings of the Japan Academy* 49.3 (1973), pp. 157–160.
- [KSJ18] S. P. Karimireddy, S. U. Stich, and M. Jaggi. “Global linear convergence of Newton’s method without strong-convexity or Lipschitz gradients”. *arXiv preprint arXiv:1806.00413* (2018).
- [KST12] D. Kifer, A. Smith, and A. Thakurta. “Private Convex Empirical Risk Minimization and High-dimensional Regression”. In: *Proceedings of the 25th Annual Conference on Learning Theory*. Ed. by S. Mannor, N. Srebro, and R. C. Williamson. Vol. 23. Proceedings of Machine Learning Research. Edinburgh, Scotland: PMLR, 25–27 Jun 2012, pp. 25.1–25.40.
- [Lyu22] X. Lyu. “Composition Theorems for Interactive Differential Privacy”. *arXiv preprint arXiv:2207.09397* (2022).
- [Mir75] L. Mirsky. “A trace inequality of John von Neumann”. *Monatshefte für mathematik* 79.4 (1975), pp. 303–306.
- [MRTZ17] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. “Learning differentially private recurrent language models”. *arXiv preprint arXiv:1710.06963* (2017).
- [MTZ19] I. Mironov, K. Talwar, and L. Zhang. “Renyi differential privacy of the sampled Gaussian mechanism”. *arXiv preprint arXiv:1908.10530* (2019).
- [Nes98] Y. Nesterov. “Introductory lectures on convex programming volume i: Basic course”. *Lecture notes* 3.4 (1998), p. 5.

- [NJLS09] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. “Robust Stochastic Approximation Approach to Stochastic Programming”. *SIAM Journal on Optimization* 19.4 (2009), pp. 1574–1609. eprint: <https://doi.org/10.1137/070704277>.
- [NP06] Y. Nesterov and B. T. Polyak. “Cubic regularization of Newton method and its global performance”. *Mathematical Programming* 108.1 (2006), pp. 177–205.
- [NW99] J. Nocedal and S. J. Wright. *Numerical optimization*. Springer, 1999.
- [PHKX+23] N. Ponomareva, H. Hazimeh, A. Kurakin, Z. Xu, et al. “How to dp-fy ml: A practical guide to machine learning with differential privacy”. *arXiv preprint arXiv:2303.00654* (2023).
- [PS22] N. Papernot and T. Steinke. “Hyperparameter Tuning with Renyi Differential Privacy”. In: *International Conference on Learning Representations*. 2022.
- [RM16] F. Roosta-Khorasani and M. W. Mahoney. “Sub-sampled newton methods ii: Local convergence rates”. *arXiv preprint arXiv:1601.04738* (2016).
- [RM19] F. Roosta-Khorasani and M. W. Mahoney. “Sub-sampled Newton methods”. *Mathematical Programming* 174 (2019), pp. 293–326.
- [SCS13] S. Song, K. Chaudhuri, and A. D. Sarwate. “Stochastic gradient descent with differentially private updates”. In: *2013 IEEE global conference on signal and information processing*. IEEE. 2013, pp. 245–248.
- [SSTT21] S. Song, T. Steinke, O. Thakkar, and A. Thakurta. “Evading the curse of dimensionality in unconstrained private glms”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2021, pp. 2638–2646.
- [Ste98] G. W. Stewart. *Matrix algorithms: volume 1: basic decompositions*. SIAM, 1998.
- [STT20] S. Song, O. Thakkar, and A. Thakurta. “Characterizing Private Clipped Gradient Descent on Convex Generalized Linear Problems”. *arXiv preprint arXiv:2006.06783* (2020).
- [STU17] A. Smith, A. Thakurta, and J. Upadhyay. “Is interaction necessary for distributed private learning?” In: *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE. 2017, pp. 58–77.
- [VZ22] S. Vadhan and W. Zhang. “Concurrent Composition Theorems for Differential Privacy”. *arXiv preprint arXiv:2207.08335* (2022).
- [WLKC+17] X. Wu, F. Li, A. Kumar, K. Chaudhuri, et al. “Bolt-on Differential Privacy for Scalable Stochastic Gradient Descent-based Analytics”. In: *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD*. Ed. by S. Salihoglu, W. Zhou, R. Chirkova, J. Yang, and D. Suciu. 2017.
- [XRV17] H. Xiao, K. Rasul, and R. Vollgraf. “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms”. *arXiv preprint arXiv:1708.07747* (2017).
- [XYRRM16] P. Xu, J. Yang, F. Roosta, C. Ré, and M. W. Mahoney. “Sub-sampled Newton methods with non-uniform sampling”. *Advances in Neural Information Processing Systems* 29 (2016).
- [ZLNM+19] G. Zhang, L. Li, Z. Nado, J. Martens, et al. “Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model”. *Advances in neural information processing systems* 32 (2019).
- [ZZMW17] J. Zhang, K. Zheng, W. Mou, and L. Wang. “Efficient private ERM for smooth objectives”. In: *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 2017, pp. 3922–3928.

A Appendix of Section 4

A.1 Proof of Theorem 4.2

Given a training set $S_n = (z_1, \dots, z_n) \in \mathcal{Z}^n$, our goal is to minimize

$$\ell(w, S_n) = \frac{1}{n} \sum_{i \in [n]} f(w, z_i).$$

Since f is a strongly convex function and \mathcal{W} is a closed and convex set, there exists a unique $w^* = \arg \min_{w \in \mathcal{W}} \ell(w, S_n)$.

Let $M \in \mathbb{R}$. In each step of the algorithm, we construct a cubic function $\phi : \mathcal{W} \times \mathcal{W} \rightarrow \mathbb{R}$ defined as

$$\phi_M(v; w) \triangleq \ell(w, S_n) + \langle \nabla \ell(w, S_n), v - w \rangle + \frac{1}{2} \langle \nabla^2 \ell(w, S_n)(v - w), v - w \rangle + \frac{M}{6} \|v - w\|^3. \quad (8)$$

We provide a lemma on the properties of $\phi_M(v; w)$.

Lemma A.1. *Let f be a L_2 -Lipschitz hessian function. Then, ϕ_M in Equation (8) satisfies the following properties:*

1. For every $M \geq 0$ and $w, v \in \mathcal{W}$ such that $v \neq w$,

$$\nabla_v^2 \phi_M(v; w) = \nabla^2 \ell(w, S_n) + \frac{M}{2} \|v - w\| I_d + \frac{M}{2 \|v - w\|} (v - w)(v - w)^T.$$

Therefore, $\nabla_v^2 \phi_M(v; w) \succcurlyeq \lambda_{\min}(\nabla^2 \ell(w, S_n)) I_d + M \|w - v\| I_d$ where $\lambda_{\min}(\nabla^2 \ell(w, S_n))$ denotes the minimum eigenvalue of $\nabla^2 \ell(w, S_n)$.

2. For every $M \geq L_2$ and $v, w \in \mathcal{W}$,

$$\ell(v, S_n) \leq \phi_M(v; w).$$

3. For every $M \in \mathbb{R}_+$ and $v, w \in \mathcal{W}$,

$$\phi_M(v; w) \leq \ell(v, S_n) + \frac{M + L_2}{6} \|v - w\|^3.$$

Proof. To show the first claim, consider

$$\nabla_v \phi_M(v; w) = \nabla \ell(w, S_n) + \nabla^2 \ell(w, S_n)(v - w) + \frac{M}{2} \|v - w\| (v - w).$$

Then, the hessian of $\phi_M(v; w)$ is given by

$$\nabla_v^2 (\phi_M(v; w)) = \nabla^2 \ell(w, S_n) + M \|w - v\| I_d + \frac{M}{\|w - v\|} (w - v)(w - v)^T.$$

Note that $(w - v)(w - v)^T$ is a PSD and rank-1 matrix whose non-zero eigenvalue is given by $\|w - v\|^2$.

The second and third parts follow from [NP06, Lemma 1] where it is shown for L_2 -Lipschitz hessian functions we have

$$\begin{aligned} & \left| \ell(v, S_n) - \left(\ell(w, S_n) + \langle \nabla \ell(w, S_n), v - w \rangle + \frac{1}{2} \langle \nabla^2 \ell(w, S_n)(v - w), (v - w) \rangle \right) \right| \\ & \leq \frac{L_2}{6} \|v - w\|^3. \end{aligned}$$

The claims are straightforward applications of this inequality. \square

We can rephrase Equation (8) as follows:

$$\begin{aligned} \phi_M(v; w) = & \frac{1}{n} \sum_{i \in [n]} f(w, z_i) + \frac{1}{n} \sum_{i \in [n]} \langle \nabla f(w, z_i), v - w \rangle + \frac{1}{2n} \sum_{i \in [n]} \langle \nabla^2 f(w, z_i), v - w \rangle + \frac{M}{6} \|w - v\|^3. \end{aligned} \quad (9)$$

Notice that $\phi_M(v; w)$ is a convex function as it is sum of a quadratic function and a cubic term, i.e., $\frac{M}{6} \|v - w\|^3$. Also, Part 1 of Lemma A.1 shows that $\phi_M(v; w)$ is a μ -strong convex function since f is a μ -strong convex function. Moreover, the ℓ_2 -sensitivity of $\nabla \phi_M(v; w)$ is $n^{-1}(L_0 + L_1 D)$, and $\phi_M(v; w)$ is $(L_0 + L_1 D + \frac{M}{2} D^2)$ -Lipschitz where D denotes the diameter of \mathcal{W} .

First we provide the performance guarantee of the solver of the cubic subproblem in Algorithm 2.

Lemma A.2. *For every $\beta \in (0, 1)$, $\rho > 0$, and $w \in \mathcal{W}$ the output of the subproblem solver, denoted by \hat{v} satisfies ρT^{-1} -zCDP and with probability at least $1 - \beta$*

$$\phi_M(\hat{v}; w) - \min_{v \in \mathcal{W}} \phi_M(v; w) = O\left(\frac{d(L_0 + L_1 D)^2 T}{\mu \rho n^2} \cdot \log(1/\beta)\right).$$

Proof. The privacy analysis is as follows: let the total privacy budget and the number of iterations of Meta Algorithm, i.e., Algorithm 1, denoted by ρ and T , respectively. We require that the output of the subproblem solver at each iteration satisfies ρ/T -zCDP which implies the output of Algorithm 1 satisfies ρ -zCDP since the zCDP constant increases linearly with the number of iterations.

Now we provide a detailed proof of the suboptimality gap. To ease the notations consider the following problem. Lets assume that we want to use DPSolver to minimize the function h which is μ -strongly convex and L -Lipschitz function whose ℓ_2 -gradient sensitivity is given by Δ . We are interested in analyzing the suboptimality gap of Algorithm 2 under the condition that the output satisfies $\tilde{\rho}$ -zCDP.

Lets assume we want to run the algorithm for N iterations where N will be determined later. Therefore, we need to make sure that each noisy gradient computation satisfies $\tilde{\rho}/N$ -zCDP. To do so, the variance of the noise needs to be $\sigma^2 = \frac{N\Delta^2}{2\tilde{\rho}}$ from ??.

Let ξ_t be the noise added to the gradient at iteration t . Then, in each step we consider $\text{grad}_t + \xi_t$ for noisy gradient. From [JNGKJ19, Lemma 1], we know that $\|\xi_t\|$ is a SubGaussian random variable with variance proxy of $c\sigma\sqrt{d}$ where c is a universal constant. We are now ready to use [HLR19, Thm. C.3]. [HLR19, Thm. C.3] shows that for every $\beta \in (0, 1]$ the suboptimality gap with probability at least $1 - \beta$ is given by

$$O\left(\frac{(L + \sigma\sqrt{d})^2 \log(1/\beta)}{\mu N}\right) = O\left(\frac{(L^2 + \sigma^2 d) \log(1/\beta)}{\mu N}\right),$$

where we simply use $(a + b)^2 \leq 2a^2 + 2b^2$ for every a, b . Finally, we need to plug in the value of σ to obtain that the suboptimality gap:

$$O\left(\frac{(L^2 + \sigma^2 d) \log(1/\beta)}{\mu N}\right) = O\left(\left(\frac{L^2}{\mu N} + \frac{\sigma^2 d}{\mu N}\right) \log(1/\beta)\right) = O\left(\left(\frac{L^2}{\mu N} + \frac{d\Delta^2}{\tilde{\rho}}\right) \log(1/\beta)\right).$$

Then, by setting the number of iterations to $N = \frac{2L^2\tilde{\rho}}{\mu d\Delta^2}$, we obtain that the suboptimality gap is given by $O\left(\frac{d\Delta^2}{\mu\tilde{\rho}} \log(1/\beta)\right)$.

In the context of our paper, $\Delta = n^{-1}(L_0 + L_1 D)$, $L = (L_0 + L_1 D + \frac{M}{2} D^2)$, and $\tilde{\rho} = \rho/T$. Setting these constants proves the lemma. \square

We drop the S_n argument from $\ell(w, S_n)$ to reduce notational clutter. Using Part 2 of Lemma A.1 we can write

$$\ell(w_{t+1}) - \ell(w^*) \leq \phi_M(w_{t+1}; w_t) - \min_{w \in \mathcal{W}} \phi_M(w; w_t) + \min_{w \in \mathcal{W}} \phi_M(w; w_t) - \ell(w^*). \quad (10)$$

Since $\phi_M(w; w_t)$ as a function of w is a strongly convex function and \mathcal{W} is a closed and convex set there exists a unique $w_{t+1}^* = \arg \min_{w \in \mathcal{W}} \phi_M(w; w_t)$.

Fix a $\beta \in (0, 1]$ and define the following event

$$\mathcal{G} = \{\forall t \in [T] : \phi_M(w_t; w_{t-1}) - \phi_M(w_t^*; w_{t-1}) \leq O\left(\frac{d(L_0 + L_1 D)^2 T}{\mu \rho n^2} \cdot \log(T/\beta)\right) \triangleq \Delta_0\}. \quad (11)$$

We claim that $\mathbb{P}(\mathcal{G}) \geq 1 - \beta$. In each step of the algorithm, we find an approximate minimizer of $\phi_M(w; w_t)$ using the subproblem solver in Algorithm 2. The performance guarantee of the subproblem solver is given in Lemma A.2 which shows that at each step of the algorithm the excess error in minimizing $\phi_M(w; w_t)$ is less than Δ_0 with probability greater than $1 - \beta/T$. Ergo, a union bound concludes the proof.

Next we provide an upperbound on $\phi_M(w_{t+1}^*; w_t) - \ell(w^*)$ in Equation (10). By the third part of Lemma A.1 we have

$$\phi_M(w_{t+1}^*; w_t) - \ell(w^*) \leq \min_{w \in \mathcal{W}} \{\ell(w) + \frac{M + L_2}{6} \|w - w_t\|^3 - \ell(w^*)\}.$$

Since \mathcal{W} is a convex set and $w_t, w^* \in \mathcal{W}$, for all $\alpha \in [0, 1]$, $(1 - \alpha)w_t + \alpha w^* \in \mathcal{W}$. Therefore,

$$\begin{aligned} & \min_{w \in \mathcal{W}} \{\ell(w) + \frac{M + L_2}{6} \|w - w_t\|^3 - \ell(w^*)\} \\ & \leq \min_{\alpha_t \in [0, 1]} \{\ell((1 - \alpha_t)w_t + \alpha_t w^*) + \frac{M + L_2}{6} \alpha_t^3 \|w_t - w^*\|^3 - \ell(w^*)\}. \end{aligned}$$

By the convexity of ℓ we have $\ell((1 - \alpha_t)w_t + \alpha_t w^*) - \ell(w^*) \leq \ell(w_t) - \ell(w^*) - \alpha_t(\ell(w_t) - \ell(w^*))$. Also, strong convexity implies that $(\frac{2}{\mu}(\ell(w_t) - \ell(w^*)))^{\frac{3}{2}} \geq \|w_t - w^*\|^3$ [Nes98]. Thus,

$$\begin{aligned} & \phi_M(w_{t+1}^*; w_t) - \ell(w^*) \\ & \leq \min_{\alpha_t \in [0, 1]} \{\ell(w_t) - \ell(w^*) - \alpha_t(\ell(w_t) - \ell(w^*)) + \alpha_t^3 \frac{M + L_2}{6} (\frac{2}{\mu}(\ell(w_t) - \ell(w^*)))^{\frac{3}{2}}\} \quad (12) \end{aligned}$$

In the rest of the proof, under the event \mathcal{G} , we provide a convergence analysis.

Let $\lambda = (\frac{3}{M + L_2})^2 (\frac{\mu}{2})^3$ and $q_t = \lambda^{-1}(\ell(w_t) - \ell(w^*))$. Then, under the event \mathcal{G} and by Equation (12), we can rephrase Equation (10) as

$$q_{t+1} \leq \lambda^{-1} \Delta_0 + \min_{\alpha_t \in [0, 1]} \{q_t - \alpha_t q_t + \frac{1}{2} \alpha_t^3 q_t^{\frac{3}{2}}\}. \quad (13)$$

Let $\alpha_t^* = \arg \min_{\alpha_t \in [0, 1]} \{q_t - \alpha_t q_t + \frac{1}{2} \alpha_t^3 q_t^{\frac{3}{2}}\} = \min\{\sqrt{\frac{2}{3\sqrt{q_t}}}, 1\}$.

First, consider the case that $q_t \geq 4/9$ so that $\alpha_t^* = \sqrt{\frac{2}{3\sqrt{q_t}}}$. We can rephrase Equation (13) as follows

$$q_{t+1} \leq \lambda^{-1} \Delta_0 + q_t - (\frac{2}{3})^{\frac{3}{2}} q_t^{\frac{3}{4}}. \quad (\text{Phase I}) \quad (14)$$

In the second case, i.e., $q_t < 4/9$, we have $\alpha_t^* = 1$ Equation (13) is given by

$$q_{t+1} \leq \lambda^{-1} \Delta_0 + \frac{1}{2} q_t^{\frac{3}{2}}. \quad (\text{Phase II}) \quad (15)$$

Assume that $q_0 \geq 4/9$. We will show that, under the event \mathcal{G} , $\{q_t\}_{t \in [T]}$ is a decreasing sequence, and as a result there exists $T_1 \in \mathbb{N}$, independent of n , such that $q_t < 4/9$ for $t \geq T_1^*$, and as a result $\alpha_t^* = 1$ for $t \geq T_1^*$.

To prove the convergence in Phase I (see Equation (14)), we follow the techniques of Nesterov and Polyak [NP06]. Let $\tilde{q}_t = \frac{9q_t}{4}$, and assume $\Delta_0 \leq \frac{4\lambda}{27}$. Then, we can rephrase the recursion for Phase I as follows:

$$\begin{aligned} \tilde{q}_{t+1} & \leq \frac{9\Delta_0}{4\lambda} + \tilde{q}_t - \frac{2}{3} \tilde{q}_t^{\frac{3}{4}} \\ & \leq \tilde{q}_t - \frac{1}{3} \tilde{q}_t^{\frac{3}{4}}, \end{aligned}$$

where the last step follows from $\tilde{q}_t \geq 1$ and $\frac{9\Delta_0}{4\lambda} \leq \frac{1}{3} \leq \frac{\tilde{q}_t^{\frac{3}{4}}}{3}$. It also shows that provided that $\tilde{q}_t \geq 1$, $\tilde{q}_{t+1} \leq \tilde{q}_t$.

Using induction, it is straightforward to show that in Phase I

$$\frac{9q_t}{4} \leq \left[\left(\frac{9q_0}{4} \right)^{\frac{1}{4}} - \frac{t}{12} \right]^4. \quad (\text{Phase I})$$

This result implies that after T_1^* iterations where

$$T_1^* \leq O\left(\frac{\sqrt{M} + L_2}{\mu^{3/4}}(\ell(w_0) - \ell(w^*))^{\frac{1}{4}}\right), \quad (16)$$

we have $q_{T_1^*} < \frac{4}{9}$, and we enter Phase II (see Equation (15)).

Next, we analyze Phase II in which the recursion is given by

$$q_{t+1} \leq \lambda^{-1}\Delta_0 + \frac{1}{2}q_t^{\frac{3}{2}}.$$

Using Lemma A.3, we obtain that after $\Theta(\log(\log(\frac{\lambda}{\Delta_0})))$ iterations we have $O(\lambda^{-1}\Delta_0)$. Therefore, the number of iterations to achieve the minimum excess error in Phase II

$$T_2^* = \tilde{\Theta}\left(\log \log\left(\frac{n}{\sqrt{\rho} \log(1/\beta)d}\right)\right). \quad (17)$$

Finally, the excess error is given by

$$\ell(w_T) - \ell(w^*) = \tilde{O}\left(\frac{d(L_0 + L_1 D)^2}{\mu \rho n^2} \cdot \log(1/\beta) \cdot (T_1^* + T_2^*)\right), \quad (18)$$

where $T = T_1^* + T_2^*$ and T_1^* and T_2^* are given by Equation (16) and Equation (17), respectively.

Lemma A.3. Let $\beta_0 > 0$ and define the sequence $a_{t+1} \leq \beta_0 + \frac{1}{2}a_t^{3/2}$ where $a_0 \leq \frac{16}{9}$. Then, after $T = \Theta(\log \log(\frac{1}{\beta_0}))$ we have $a_T = O(\beta_0)$.

Proof. Without loss of generality, assume $a_{t+1} = \beta_0 + \frac{1}{2}a_t^{3/2}$. We define another sequence $\{b_t\}_{t \in \mathbb{N}}$ as follows: $b_0 = a_0$ and $b_{t+1} = \frac{3}{4}(b_t)^{\frac{3}{2}}$. By induction one can easily prove that for every $t \in \mathbb{N}$ such that $\beta_0 \leq \frac{1}{4}(b_t)^{\frac{3}{2}}$, we have $b_{t+1} \geq a_{t+1}$. Then, we can write

$$b_{t+1} = \frac{3}{4}(b_t)^{\frac{3}{2}} \Leftrightarrow \frac{9}{16}b_{t+1} = \left(\frac{9}{16}b_t\right)^{\frac{3}{2}}.$$

Therefore, we obtain that $\log(\frac{9}{16}b_t) = (\frac{3}{2})^t \log(\frac{9}{16}b_0)$. We want to find T such that $\beta_0 \leq \frac{1}{4}(b_T)^{\frac{3}{2}} \leq 2\beta_0$ which is equivalent to $\log(\frac{8\beta_0}{3}) \leq \log(b_T) \leq \log(\frac{16\beta_0}{3})$. Then, by some simple manipulations we can see that $T = \Theta(\log(\log(\frac{1}{\beta_0})))$. Therefore, we have $b_{T+1} = O(\beta_0)$. Also, by the construction, $a_{T+1} \leq b_{T+1} = O(\beta_0)$. \square

A.2 Private Accelerated Nesterov's Method

In this section, we present a DP variant of the accelerated Nesterov's Method. The proof ideas are based on [Nes98; NJLS09; GL12].

Algorithm 4 Private Accelerated Nesterov's Method for L_0 -Lipschitz, L_1 -smooth convex function on a bounded feasible set \mathcal{W} with diameter D .

- 1: Input: $w_0 \in \mathcal{W}$, Privacy Guarantee ρ -zCDP.
 - 2: $T = \Theta\left(\left(\frac{D\sqrt{\rho n}}{L_0}\right)^{1/2}\right)$, $\sigma^2 = \frac{L_0^2 T}{2\rho n^2}$.
 - 3: $w_0^{\text{ag}} = w_0 \in \mathcal{W}$
 - 4: $\alpha_t = \frac{2}{t+1}$, $\gamma_t = \frac{4\gamma}{t(t+1)}$ where $\gamma = 2L_1$.
 - 5: **for** $t = 1, \dots, T$ **do**
 - 6: $w_t^{\text{md}} = w_{t-1}^{\text{ag}} + \alpha_t(w_t - w_{t-1}^{\text{ag}})$
 - 7: $G_t = \nabla \ell(w_t^{\text{md}}, S_n) + \mathcal{N}(0, \sigma^2 I_d)$
 - 8: $w_t = \Pi_{\mathcal{W}}(w_{t-1} - \frac{\alpha_t}{\gamma_t} G_t) = \arg \min_{v \in \mathcal{W}} \{\alpha_t \langle G_t, v \rangle + \frac{\gamma_t}{2} \|v - w_{t-1}\|^2\}$
 - 9: $w_t^{\text{ag}} = \alpha_t w_t + (1 - \alpha_t) w_{t-1}^{\text{ag}}$
 - 10: **Return** w_T^{ag}
-

Theorem A.4. Let f be a convex, L_0 -Lipschitz, and L_1 -smooth. Also, assume that $\mathcal{W} \subseteq \mathbb{R}^d$ is a convex set and has finite diameter D . Let $w^* \in \arg \min_{w \in \mathcal{W}} \ell(w, S_n)$. Then, for every $n \in \mathbb{N}$, $S_n \in \mathcal{Z}^n$, and $\rho > 0$, the output of Algorithm 4, i.e., w_T^{ag} , satisfies ρ -zCDP and

$$\mathbb{E}[\ell(w_T^{\text{ag}}, S_n) - \ell(w^*, S_n)] = O\left(\frac{L_0 D \sqrt{d}}{n \sqrt{\rho}}\right).$$

Also, the oracle complexity of Algorithm 4 is

$$T = \Theta\left(\sqrt{\frac{Dn\sqrt{\rho}}{L_0}}\right).$$

Proof. First, we start with the privacy proof. The ℓ_2 -sensitivity of $\nabla \ell(w, S_n)$ for every $w \in \mathcal{W}$ is given by $\frac{L_0}{n}$. Therefore, by the composition properties of zCDP in [BS16, Lem. 2.3] and the zCDP analysis of the Gaussian mechanism in [BS16, Lem. 2.5], it is straightforward to show that w_T^{ag} satisfies ρ -zCDP.

Then, we analyze the excess error. For every $v \in \mathcal{W}$, by the smoothness of ℓ we can write

$$\begin{aligned} \ell(w_t^{\text{ag}}) &\leq \ell(v) + \langle \nabla \ell(v), w_t^{\text{ag}} - v \rangle + \frac{L_1}{2} \|w_t^{\text{ag}} - v\|^2 \\ &= \ell(v) + \alpha_t \langle \nabla \ell(v), w_t - v \rangle + (1 - \alpha_t) \langle \nabla \ell(v), w_{t-1}^{\text{ag}} - v \rangle + \frac{L_1}{2} \|w_t^{\text{ag}} - v\|^2 \\ &\leq \alpha_t (\ell(v) + \langle \nabla \ell(v), w_t - v \rangle) + (1 - \alpha_t) \ell(w_{t-1}^{\text{ag}}) + \frac{L_1}{2} \|w_t^{\text{ag}} - v\|^2. \end{aligned} \quad (19)$$

Here, the second step is by definition of w_t^{ag} , and the last step is by convexity of ℓ which implies $\ell(v) + \langle \nabla \ell(v), w_{t-1}^{\text{ag}} - v \rangle \leq \ell(w_{t-1}^{\text{ag}})$.

Note that

$$\begin{aligned} w_t^{\text{ag}} - w_t^{\text{md}} &= \alpha_t w_t + (1 - \alpha_t) w_{t-1}^{\text{ag}} - (w_{t-1}^{\text{ag}} + \alpha_t (w_{t-1} - w_{t-1}^{\text{ag}})) \\ &= \alpha_t (w_t - w_{t-1}). \end{aligned}$$

In the next step, we substitute $v = w_t^{\text{md}}$ in Equation (19) to obtain

$$\begin{aligned} \ell(w_t^{\text{ag}}) &\leq \alpha_t (\ell(w_t^{\text{md}}) + \langle \nabla \ell(w_t^{\text{md}}), w_t - w_t^{\text{md}} \rangle) + (1 - \alpha_t) \ell(w_{t-1}^{\text{ag}}) + \frac{L_1}{2} \|w_t^{\text{ag}} - w_t^{\text{md}}\|^2 \\ &= (1 - \alpha_t) \ell(w_{t-1}^{\text{ag}}) + \alpha_t (\ell(w_t^{\text{md}}) + \langle \nabla \ell(w_t^{\text{md}}), w_t - w_t^{\text{md}} \rangle) + \frac{L_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2 \\ &= (1 - \alpha_t) \ell(w_{t-1}^{\text{ag}}) + \alpha_t (\ell(w_t^{\text{md}}) + \langle \nabla \ell(w_t^{\text{md}}), w_t - w_t^{\text{md}} \rangle) + \frac{\gamma_t}{2} \|w_t - w_{t-1}\|^2 \\ &\quad - \frac{\gamma_t - L_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2. \end{aligned} \quad (20)$$

Let $\xi_t = G_t - \nabla \ell(w_t^{\text{md}})$. Notice that the projection step can be written as $w_t = \arg \min_{v \in \mathcal{W}} \{\alpha_t \langle G_t, v - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|v - w_{t-1}\|^2\}$. The function $g : \mathcal{W} \rightarrow \mathbb{R}$, $g(v) = \alpha_t \langle G_t, v - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|v - w_{t-1}\|^2$ is a γ_t strongly convex function. By the optimality condition for strongly convex functions we can write, for every $v \in \mathcal{W}$

$$\begin{aligned} & \alpha_t \langle G_t, w_t - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|w_t - w_{t-1}\|^2 \\ & \leq \alpha_t \langle G_t, v - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|v - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|v - w_t\|^2 \\ & = \alpha_t \langle \nabla \ell(w_t^{\text{md}}), v - w_t^{\text{md}} \rangle + \alpha_t \langle \xi_t, v - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|v - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|v - w_t\|^2. \end{aligned} \quad (21)$$

By replacing $v = w^*$ where $w^* \in \arg \min \ell(w, S_n)$ in Equation (21), we obtain

$$\begin{aligned} & \alpha_t (\ell(w_t^{\text{md}}) + \langle \nabla \ell(w_t^{\text{md}}), w_t - w_t^{\text{md}} \rangle) + \frac{\gamma_t}{2} \|w_t - w_{t-1}\|^2 \\ & = \alpha_t \ell(w_t^{\text{md}}) + \alpha_t \langle G_t, w_t - w_t^{\text{md}} \rangle - \alpha_t \langle \xi_t, w_t - w_t^{\text{md}} \rangle + \frac{\gamma_t}{2} \|w_t - w_{t-1}\|^2 \\ & \leq \alpha_t (\ell(w_t^{\text{md}}) + \langle \nabla \ell(w_t^{\text{md}}), w^* - w_t^{\text{md}} \rangle) + \alpha_t \langle \xi_t, w^* - w_t \rangle + \frac{\gamma_t}{2} \|w^* - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|w^* - w_t\|^2 \\ & = \alpha_t \ell(w^*) + \alpha_t \langle \xi_t, w^* - w_t \rangle + \frac{\gamma_t}{2} \|w^* - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|w^* - w_t\|^2. \end{aligned} \quad (22)$$

From Equation (21) and Equation (22) we can write

$$\begin{aligned} \ell(w_t^{\text{ag}}) - \ell(w^*) & \leq (1 - \alpha_t)(\ell(w_{t-1}^{\text{ag}}) - \ell(w^*)) + \alpha_t \langle \xi_t, w^* - w_t \rangle \\ & \quad + \frac{\gamma_t}{2} \|w^* - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|w^* - w_t\|^2 - \frac{\gamma_t - \mathbf{L}_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2. \end{aligned}$$

Note that ξ_t is independent from the history up to time $t - 1$, i.e., $\{(w_i^{\text{md}}, w_i, w_i^{\text{ag}})\}_{i=1}^{t-1}$. Therefore, $\mathbb{E}[\langle \xi_t, w_{t-1} \rangle] = 0$ as $\xi_t \sim \mathcal{N}(0, \sigma^2 I_d)$. Using this observation, we can write

$$\begin{aligned} & \mathbb{E}[\alpha_t \langle \xi_t, w^* - w_t \rangle - \frac{\gamma_t - \mathbf{L}_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2] \\ & = \mathbb{E}[\alpha_t \langle \xi_t, w^* - w_{t-1} \rangle + \alpha_t \langle \xi_t, w_t - w_{t-1} \rangle - \frac{\gamma_t - \mathbf{L}_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2] \\ & \leq \mathbb{E}[\alpha_t \|\xi_t\| \|w_t - w_{t-1}\| - \frac{\gamma_t - \mathbf{L}_1 \alpha_t^2}{2} \|w_t - w_{t-1}\|^2] \\ & \leq \frac{\alpha_t^2}{\gamma_t - \mathbf{L}_1 \alpha_t^2} \mathbb{E}[\|\xi_t\|^2] \\ & = \frac{\alpha_t^2}{\gamma_t - \mathbf{L}_1 \alpha_t^2} \cdot \sigma^2 I_d, \end{aligned} \quad (23)$$

where the second step follows from Cauchy–Schwarz inequality. Therefore, we obtain that

$$\begin{aligned} & \mathbb{E}[\ell(w_t^{\text{ag}}) - \ell(w^*)] \\ & \leq (1 - \alpha_t) \mathbb{E}[\ell(w_{t-1}^{\text{ag}}) - \ell(w^*)] + \frac{\alpha_t^2 (\sigma^2 d)}{\gamma_t - \mathbf{L}_1 \alpha_t^2} + \frac{\gamma_t}{2} \|w^* - w_{t-1}\|^2 - \frac{\gamma_t}{2} \|w^* - w_t\|^2. \end{aligned} \quad (24)$$

Let $\Gamma_t = \begin{cases} 1 & t = 1 \\ (1 - \alpha_t) \Gamma_{t-1} & t \geq 2 \end{cases} = \frac{2}{t(t+1)}$. Note that since $\gamma = 2\mathbf{L}_1$, we have $\gamma_t - \mathbf{L}_1 \alpha_t^2 = \frac{4\gamma}{t(t+1)} - \frac{4\mathbf{L}_1}{(t+1)^2} = \frac{2\gamma}{(t+1)^2}$, and $\frac{\gamma_t}{\Gamma_t} = 2\gamma$. Consider dividing both side of Equation (24) by Γ_t and summing up from 1 to T to obtain

$$\frac{\mathbb{E}[\ell(w_T^{\text{ag}}) - \ell(w^*)]}{\Gamma_T} \leq \sum_{\tau=1}^T \frac{\gamma_\tau}{2\Gamma_\tau} (\|w^* - w_{\tau-1}\|^2 - \|w^* - w_\tau\|^2) + \sigma^2 d \cdot \sum_{\tau=1}^T \frac{1}{\Gamma_\tau} \cdot \frac{\alpha_\tau^2}{\gamma_\tau - \mathbf{L}_1 \alpha_\tau^2}.$$

Note that $\frac{\gamma_t}{\Gamma_t} = 2\gamma$. Therefore,

$$\sum_{\tau=1}^T \frac{\gamma_\tau}{2\Gamma_\tau} (\|w^\star - w_{\tau-1}\|^2 - \|w^\star - w_\tau\|^2) = \gamma \|w_0 - w^\star\|^2 \leq \gamma D^2.$$

Then, for the last term consider

$$\sum_{\tau=1}^T \frac{1}{\Gamma_\tau} \cdot \frac{\alpha_\tau^2}{\gamma_\tau - L_1 \alpha_\tau^2} = \frac{1}{3\gamma} T(T+1)(T+2).$$

By combining all the previous steps, we get the following bound on the expected excess error

$$\begin{aligned} \mathbb{E}[\ell(w_T^{\text{ag}}) - \ell(w^\star)] &\leq \frac{2}{T(T+1)} \gamma D^2 + 2\sigma^2 d \frac{T+2}{3\gamma} \\ &= \frac{2}{T(T+1)} \gamma D^2 + d \frac{2L_0^2 T(T+2)}{6\gamma \rho n^2}. \end{aligned} \quad (25)$$

Finally, optimizing Equation (25) over T , we conclude that with at most T oracle calls where

$$T = \Theta\left(\left(\frac{D^2 \rho n^2}{d L_0^2}\right)^{1/4}\right),$$

the achievable excess error is given by

$$\mathbb{E}[\ell(w_T^{\text{ag}}) - \ell(w^\star)] = O\left(\frac{L_0 D \sqrt{d}}{n \sqrt{\rho}}\right).$$

□

B Appendix of Section 5

B.1 Proof of Lemma 5.1

We begin the proof by a lemma from [AMN20].

Lemma B.1 ([AMN20, Prop. 4.1]). *For all $0 < \mu < 1$ and $\lambda \in \mathbb{R}$, we have*

$$\frac{2}{\lambda^2} (\log(\mu \exp(\lambda) + 1 - \mu) - \mu \lambda) \leq \frac{1/2 - \mu}{\log(1/\mu - 1)}.$$

We will use the following reformulation of Lemma B.1. Let $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}$ be two constants. Then, substitute $\mu = \frac{\exp(\alpha)}{1 + \exp(\alpha)}$ and $\lambda = \beta - \alpha$ in Lemma B.1. By some simple manipulations we obtain that $\forall \beta, \alpha \in \mathbb{R}$

$$\log(1 + \exp(\beta)) \leq \log(1 + \exp(\alpha)) + \frac{\beta - \alpha}{1 + \exp(-\alpha)} + \begin{cases} \frac{\exp(\alpha) - 1}{4\alpha(\exp(\alpha) + 1)} (\beta - \alpha)^2 & \alpha \neq 0 \\ \frac{(\beta - \alpha)^2}{4} & \alpha = 0 \end{cases}. \quad (26)$$

Finally, let $w, v, x \in \mathbb{R}^d$ and $y \in \{-1, +1\}$, by substituting $\alpha = \langle -yx, v \rangle$ and $\beta = \langle -yx, w \rangle$, we obtain the stated result in Lemma 5.1.

B.2 Comparison of the Approximations

Consider $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = \log(1 + \exp(x))$ which can be seen as a logistic loss in one dimension. Figure 5 compares the three approaches for the quadratic approximation: second-order Taylor approximation, our upper bound in Lemma 5.1, and the upper bound based on smoothness. As can be seen the upper bound in Lemma 5.1 provides tighter approximation compared to the upper bound based on smoothness. Also, the second-order Taylor approximation is not an upper bound on the function.

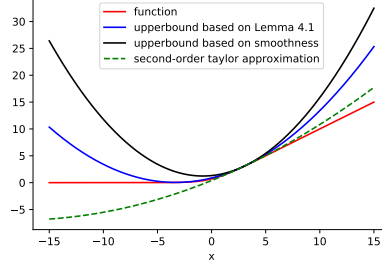


Figure 5: Comparison between the approximation of the logistic loss function

B.3 Privacy proof of Algorithm 3 for add

Theorem B.2. Assume in Algorithm 3 we choose add for the SOI modification. Then, for every training set $S_n \in (\mathbb{R}^d \times \{-1, +1\})^n$, $w_0 \in \mathcal{W}$, $\lambda_0 > 0$, $T \in \mathbb{N}$, $\rho \in \mathbb{R}_+$, and $\theta \in (0, 1)$, by setting

$$\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}, \quad \sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 + \lambda_0)\sqrt{2\rho\theta}},$$

w_T satisfies ρ -zCDP.

Proof. For the privacy analysis, we assume a two stage procedure. The loss function is 1-Lipschitz. Therefore, by setting

$$\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}},$$

the mechanism in Line 10 of Algorithm 3 satisfies $(1-\theta)\frac{\rho}{T}$ -zCDP by ??.

Notice that $\Psi_{\lambda_0}(H(w_t, S_n), \text{add}) = H(w_t, S_n) + \lambda_0 I_d$. We need to bound the ℓ_2 sensitivity of the search direction which is given by

$$\sup_{S_n \in (\mathbb{R}^d \times \{-1, 1\})^n} \sup_{z_{n+1} = (x_{n+1}, y_{n+1}) \in \mathbb{R}^d \times \{-1, 1\}: \|x_{n+1}\| \leq 1} \left\| [H(w_t, S_n) + \lambda_0 I_d]^{-1} \tilde{g}_t - [H(w_t, S_n) + \frac{1}{n} H(w_t, z_{n+1}) + \lambda_0 I_d]^{-1} \tilde{g}_t \right\|. \quad (27)$$

By the definition of the operator norm, we have

$$\begin{aligned} & \left\| [H(w_t, S_n) + \lambda_0 I_d]^{-1} \tilde{g}_t - [H(w_t, S_n) + \frac{1}{n} H(w_t, z_{n+1}) + \lambda_0 I_d]^{-1} \tilde{g}_t \right\| \\ & \leq \left\| [H(w_t, S_n) + \lambda_0 I_d]^{-1} - [H(w_t, S_n) + \frac{1}{n} H(w_t, z_{n+1}) + \lambda_0 I_d]^{-1} \right\| \|\tilde{g}_t\|. \end{aligned}$$

Let $A \triangleq H(w_t, S_n) + \lambda_0 I_d$. For both type of the SOI, we have $H(w_t, S_n) + \frac{1}{n} H(w_t, z_{n+1}) + \lambda_0 I_d = A + \beta x_{n+1} x_{n+1}^\top$, where β only depends on x_{n+1}, w_t and $\beta \leq \frac{1}{4n}$ (see Remark 5.3).

We can drop the subscript $n+1$ and rephrase the problem as follows

$$\sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \left\| (A + \beta x x^\top)^{-1} - A^{-1} \right\|. \quad (28)$$

We begin by applying the Sherman–Morrison formula [GV13] to $(A + \beta x x^\top)^{-1}$ to obtain

$$(A + \beta x x^\top)^{-1} - A^{-1} = -\frac{\beta A^{-1} x x^\top A^{-1}}{1 + \beta x^\top A^{-1} x}.$$

A is a PSD matrix. Let the eigenvalue decomposition of A be $A = \sum_{i \in [d]} \lambda_i u_i u_i^\top = U \Lambda U^\top$ where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$. Using this representation, we can write

$$\begin{aligned} \sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \|(A + \beta x x^\top)^{-1} - A^{-1}\| &= \sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \frac{\|\beta A^{-1} x x^\top A^{-1}\|}{1 + \beta x^\top A^{-1} x} \\ &= \sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \frac{\|\beta U \Lambda^{-1} U^\top x x^\top U \Lambda^{-1} U^\top\|}{1 + \beta x^\top U \Lambda^{-1} U^\top x}. \end{aligned}$$

Then, consider the change of variable to $v = U^\top x$:

$$\begin{aligned} \sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \|(A + \beta x x^\top)^{-1} - A^{-1}\| &= \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\|\beta U \Lambda^{-1} v v^\top \Lambda^{-1} U^\top\|}{1 + \beta v^\top \Lambda^{-1} v} \\ &= \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\beta \|U \Lambda^{-1} v v^\top \Lambda^{-1} U^\top\|}{1 + \beta v^\top \Lambda^{-1} v}. \end{aligned}$$

Notice that $U \Lambda^{-1} v v^\top \Lambda^{-1} U^\top$ is a rank-one matrix. For a rank-1 matrix, the operator norm is given by its non-zero eigenvalue. Thus,

$$\begin{aligned} \|U \Lambda^{-1} v v^\top \Lambda^{-1} U^\top\| &= \|U \Lambda^{-1} v\|_2^2 \\ &= \|\Lambda^{-1} v\|_2^2. \end{aligned}$$

The last step follows from the fact that U is an orthonormal matrix. Therefore, by combining the previous representations we obtain

$$\sup_{x \in \mathbb{R}^d: \|x\| \leq 1} \|(A + \beta x x^\top)^{-1} - A^{-1}\| = \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\beta \|\Lambda^{-1} v\|_2^2}{1 + \beta v^\top \Lambda^{-1} v}.$$

For every $a > 0$, define $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) \triangleq \frac{x}{1+ax}$ and notice that h is increasing for $x > 0$. Using this fact and by considering $v^\top \Lambda^{-1} v > 0$ and $0 < \beta \leq \frac{1}{4n}$, we obtain

$$\begin{aligned} \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\beta \|\Lambda^{-1} v\|_2^2}{1 + \beta v^\top \Lambda^{-1} v} &\leq \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \sup_{\beta \in [0, \frac{1}{4}]} \frac{\beta \|\Lambda^{-1} v\|_2^2}{1 + \beta v^\top \Lambda^{-1} v} \\ &\leq \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\frac{1}{4n} \|\Lambda^{-1} v\|_2^2}{1 + \frac{1}{4n} v^\top \Lambda^{-1} v} \\ &= \sup_{v \in \mathbb{R}^d: \|v\| \leq 1} \frac{\sum_{i=1}^d \frac{1}{\lambda_i^2} v_i^2}{4n + \sum_{i=1}^d \frac{1}{\lambda_i} v_i^2}. \end{aligned}$$

Notice that by the definition of A , for $i \in [d]$, $\lambda_0 \leq \lambda_i$. Therefore,

$$\frac{\sum_{i=1}^d \frac{1}{\lambda_i^2} v_i^2}{4n + \sum_{i=1}^d \frac{1}{\lambda_i} v_i^2} \leq \frac{1}{\lambda_0} \frac{\sum_{i=1}^d \frac{1}{\lambda_i} v_i^2}{4n + \sum_{i=1}^d \frac{1}{\lambda_i} v_i^2}.$$

Then, note that for every $v \in \mathbb{R}^d$ such that $\|v\| \leq 1$, we have $\sum_{i=1}^d \frac{1}{\lambda_i} v_i^2 \leq \frac{1}{\lambda_0}$. Also, $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) \triangleq \frac{x}{4n+x}$ is increasing for $x > 0$. Thus, using these two facts we obtain

$$\frac{1}{\lambda_0} \frac{\sum_{i=1}^d \frac{1}{\lambda_i} v_i^2}{4n + \sum_{i=1}^d \frac{1}{\lambda_i} v_i^2} \leq \frac{1}{4n\lambda_0^2 + \lambda_0}.$$

Therefore, we showed that

$$\begin{aligned} \sup_{\substack{z_{n+1} = (x_{n+1}, y_{n+1}): \\ \|x_{n+1}\| \leq 1}} \left\| [H(w_t, S_n) + \lambda_0 I_d]^{-1} \tilde{g}_t - [H(w_t, S_n) + \frac{1}{n} H(w_t, z_{n+1}) + \lambda_0 I_d]^{-1} \tilde{g}_t \right\| \\ \leq \frac{\|\tilde{g}_t\|}{4n\lambda_0^2 + \lambda_0}. \end{aligned}$$

This shows that by setting

$$\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 + \lambda_0)\sqrt{2\rho\theta}},$$

the mechanism in Line 11 of Algorithm 3 is $\frac{\theta\rho}{T}$ -zCDP using ??.

In each step of the algorithm we have two privatizing step that satisfy $\frac{(1-\theta)\rho}{T}$ and $\frac{\theta\rho}{T}$. By the composition property of zCDP [BS16, Lemma 2.3], we conclude that w_T satisfies ρ -zCDP. \square

B.4 Privacy Proof of Algorithm 3 for clip

Theorem B.3. Assume in Algorithm 3, we choose clip for the SOI modification. Then, for every training set $S_n \in (\mathbb{R}^d \times \{-1, +1\})^n$, $w_0 \in \mathcal{W}$, $\lambda_0 > 0$, $T \in \mathbb{N}$, $\rho \in \mathbb{R}_+$, and $\theta \in (0, 1)$ such that $n > \frac{1}{4\lambda_0}$, by setting

$$\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}, \quad \sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 - \lambda_0)\sqrt{2\rho\theta}},$$

w_T satisfies ρ -zCDP.

Proof. Similar to the proof of Theorem B.2, we use a two-stage approach. Since the logistic loss is a 1-Lipschitz function, by setting

$$\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}},$$

the mechanism in Line 10 of Algorithm 3 satisfies $(1-\theta)\frac{\rho}{T}$ -zCDP by ??.

For the second step, following the same line as in the proof of Theorem B.2, we need to upper bound

$$\sup_{S_n \in (\mathbb{R}^d \times \{-1, 1\})^n} \sup_{\substack{z_{n+1} = (x_{n+1}, y_{n+1}) \in \mathbb{R}^d \times \{-1, 1\}: \\ \|x_{n+1}\| \leq 1}} \left\| \left[\Psi_{\lambda_0}(H(w_t, S_n), \text{clip}) \right]^{-1} - \left[\Psi_{\lambda_0}(H(w_t, S_n) + \frac{1}{n}H(w_t, z_{n+1}), \text{clip}) \right]^{-1} \right\|.$$

Let $A = \Psi_{\lambda_0}(H(w_t, S_n), \text{clip})$ and $B = \Psi_{\lambda_0}(H(w_t, S_n) + \frac{1}{n}H(w_t, z_{n+1}), \text{clip})$. We need a lemma for the next step of the proof.

Lemma B.4. Let $A, B \in \mathbb{R}^{d \times d}$ be positive definite matrices. If $\|A - B\| \cdot \|A^{-1}\| < 1$, then

$$\|A^{-1} - B^{-1}\| \leq \frac{\|A - B\| \cdot \|A^{-1}\|^2}{1 - \|A - B\| \cdot \|A^{-1}\|}.$$

Proof. Let $B = A - C$. We have the identity $(A - C)^{-1} = A^{-1} \sum_{k=0}^{\infty} (CA^{-1})^k$, which holds as long as $\|CA^{-1}\| < 1$ [Ste98, Thm. 4.8]. Thus

$$\|A^{-1} - B^{-1}\| = \left\| A^{-1} \sum_{k=1}^{\infty} (CA^{-1})^k \right\| \leq \|A^{-1}\| \sum_{k=1}^{\infty} \|CA^{-1}\|^k = \frac{\|A^{-1}\| \cdot \|CA^{-1}\|}{1 - \|CA^{-1}\|}. \quad (29)$$

Now $\|CA^{-1}\| \leq \|C\| \cdot \|A^{-1}\| = \|A - B\| \cdot \|A^{-1}\|$, which gives the result. \square

Using Lemma B.4, we can write

$$\|A^{-1} - B^{-1}\| \leq \frac{\|A - B\| \cdot \|A^{-1}\|^2}{1 - \|A - B\| \cdot \|A^{-1}\|}, \quad (30)$$

provided that $\|A - B\| \cdot \|A^{-1}\| < 1$. Note that Frobenius norm of a matrix is not smaller than the operator norm, i.e., $\|A - B\| \leq \|A - B\|_F$. For the next step of the proof, we need a lemma.

Lemma B.5. For every $\lambda_0 \geq 0$ and $A \in \mathbb{R}^{d \times d}$ with $A^\top = A$, we have

$$\Psi_{\lambda_0}(A, \text{clip}) = \arg \min_{\hat{A} \in \mathbb{R}^{d \times d}: \hat{A}^\top = \hat{A}, \forall x \in \mathbb{R}^d \ x^\top \hat{A} x \geq \lambda_0 \|x\|_2^2} \|\hat{A} - A\|_F. \quad (31)$$

Moreover, for every $\lambda_0 \geq 0$ and every PSD matrices $A \in \mathbb{R}^{d \times d}$ and $B \in \mathbb{R}^{d \times d}$, we have

$$\|\Psi_{\lambda_0}(A, \text{clip}) - \Psi_{\lambda_0}(B, \text{clip})\|_F \leq \|A - B\|_F.$$

Proof. Consider the eigenvalue decomposition of A as $A = \sum_{i=1}^d \lambda_i u_i u_i^\top = U \Lambda U^\top$ where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ and $U \in \mathbb{R}^{d \times d}$ is matrix with u_i on its i -th column. We can represent every matrix in the feasible set of the optimization in Equation (31) as $\hat{A} = \sum_{i=1}^d v_i v_i^\top \tilde{\lambda}_i = V \tilde{\Lambda} V^\top$ where $\{v_i\}_{i \in [d]}$ are an orthonormal basis for \mathbb{R}^d and $\min_{i \in [d]} \tilde{\lambda}_i \geq \lambda_0$. By using simple facts about Frobenius norm and the eigenvalue decomposition, we can write

$$\begin{aligned} \|A - \hat{A}\|_F^2 &= \|U \Lambda U^\top - \hat{A}\|_F^2 \\ &= \text{trace}((\Lambda - U^\top \hat{A} U)^2). \end{aligned}$$

We have $(\Lambda - U^\top \hat{A} U)^2 = \Lambda^2 - U^\top \hat{A} U \Lambda - \Lambda U^\top \hat{A} U + U^\top \hat{A}^2 U$. Thus,

$$\begin{aligned} \|A - \hat{A}\|_F^2 &= \text{trace}(\Lambda^2) - 2\text{trace}(U^\top \hat{A} U \Lambda) + \text{trace}(U^\top \hat{A}^2 U) \\ &= \text{trace}(\Lambda^2) - 2\text{trace}(U^\top \hat{A} U \Lambda) + \text{trace}(\hat{A}^2), \end{aligned} \quad (32)$$

where the last step follows from $\text{trace}(U^\top \hat{A}^2 U) = \text{trace}(\hat{A}^2) = \text{trace}(V \tilde{\Lambda}^2 V^\top) = \text{trace}(\tilde{\Lambda}^2)$.

As the trace operator is invariant under the cyclic permutation, we have $\text{trace}(U^\top \hat{A} U \Lambda) = \text{trace}(U \Lambda U^\top \hat{A}) = \text{trace}(A \hat{A})$. Then, we invoke Von Neumann's trace inequality [Mir75] which states that

$$\text{trace}(A \hat{A}) \leq \sum_{i=1}^d \lambda_i \hat{\lambda}_i, \quad (33)$$

where the equality holds if A and \hat{A} share the same eigenvectors. Therefore, by Equation (32) and Equation (33), we have

$$\begin{aligned} \|A - \hat{A}\|_F^2 &\geq \text{trace}(\Lambda^2) - 2 \sum_{i=1}^n \lambda_i \hat{\lambda}_i + \text{trace}(\hat{A}^2) \\ &= \sum_{i=1}^d (\lambda_i - \hat{\lambda}_i)^2. \end{aligned}$$

It is straightforward to see that

$$\sum_{i=1}^d (\lambda_i - \hat{\lambda}_i)^2 \geq \sum_{i=1}^d (\lambda_i - \max\{\lambda_0, \lambda_i\})^2. \quad (34)$$

Thus, we obtain that for every \hat{A} in the feasible set of Equation (31) the following holds

$$\|A - \hat{A}\|_F^2 \geq \sum_{i=1}^d (\lambda_i - \max\{\lambda_0, \lambda_i\})^2.$$

For deriving this lower bound we used two inequalities in Equation (33) and Equation (34). The equality condition for Equation (33) is that A and \hat{A} share the same eigenvectors, and, for Equation (34), the equality condition is $\hat{\lambda}_i = \max\{\lambda_i, \lambda_0\}$ for every $i \in [d]$. Therefore, we conclude that $\Psi_{\lambda_0}(A, \text{clip})$ is a minimizer of Equation (31).

For the second part, notice that the feasible set in the optimization problem Equation (31) is a convex and closed subset of $\mathbb{R}^{d \times d}$. Also, Frobenius norm is a metric induced by an inner product over the vector space of the real symmetric matrices. Therefore, we conclude that for every $\lambda_0 > 0$, $\Psi_{\lambda_0}(\cdot, \text{clip})$ is a projection onto a convex and closed set, and the second claim follows. \square

In Lemma B.5, we show that $\Psi_{\lambda_0}(\cdot, \text{clip})$ is a Frobenius-norm projection onto a convex and closed set. Therefore, by the contraction property of the projection, we have

$$\begin{aligned}\|A - B\|_F &= \left\| \Psi_{\lambda_0}(H(w_t, S_n) + \frac{1}{n}H(w_t, z_{n+1}), \text{clip}) - \Psi_{\lambda_0}(H(w_t, S_n), \text{clip}) \right\|_F \\ &\leq \frac{1}{n} \|H(w_t, z_{n+1})\|_F.\end{aligned}$$

Since $H(w_t, z_{n+1})$ is a rank-1 matrix, we have $\frac{1}{n} \|H(w_t, z_{n+1})\|_F \leq \frac{1}{4n}$ (see Remark 5.3.).

For every $a > 0$, $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) = \frac{x}{1-ax}$ is increasing for $x < \frac{1}{a}$. Therefore, from Equation (30)

$$\begin{aligned}\|A^{-1} - B^{-1}\| &\leq \|A^{-1}\|^2 \frac{\|A - B\|}{1 - \|A - B\| \cdot \|A^{-1}\|} \\ &\leq \frac{\|A\|^{-2}}{4n - \|A\|^{-1}}.\end{aligned}$$

Consider $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) = \frac{x^2}{4n-x}$. This function is increasing in the interval $0 \leq x < 4n$. Using this observation, $\|A\|^{-1} \leq \frac{1}{\lambda_0}$, and $4n\lambda_0 \geq 1$ we obtain

$$\begin{aligned}\|A^{-1} - B^{-1}\| &\leq \frac{\|A\|^{-2}}{4n - \|A\|^{-1}} \\ &\leq \frac{1}{4n\lambda_0^2 - \lambda_0}.\end{aligned}$$

Therefore, we have shown that

$$\begin{aligned}&\sup_{S_n} \sup_{z_{n+1}=(x_{n+1}, y_{n+1}): \|x_{n+1}\| \leq 1} \\ &\left\| [\Psi_{\lambda_0}(H(w_t, S_n), \text{clip})]^{-1} \tilde{g}_t - [\Psi_{\lambda_0}(H(w_t, S_n) + \frac{1}{n}H(w_t, z_{n+1}), \text{clip})]^{-1} \tilde{g}_t \right\| \\ &\leq \frac{\|\tilde{g}_t\|}{4n\lambda_0^2 - \lambda_0}.\end{aligned}$$

This shows that by setting

$$\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_0^2 - \lambda_0)\sqrt{2\rho\theta}},$$

the mechanism in Line 11 of Algorithm 3 is $\frac{\theta\rho}{T}$ -zCDP.

In each step of the algorithm we have two privatizing step that satisfy $\frac{(1-\theta)\rho}{T}$ and $\frac{\theta\rho}{T}$. By the composition property of zCDP [BS16, Lemma 2.3], we conclude that w_T satisfies ρ -zCDP. \square

B.5 Description of Double noise with adaptive minimum eigenvalue selection

In Algorithm 5, we provide the detailed algorithmic description of the variant of our algorithm with an adaptive minimum eigenvalue selection. We use this variant in our numerical results.

B.5.1 Derivation

Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ be the second-order approximation of $\ell(v, S_n)$ at w_t given by

$$\phi(v) = \ell(w_t, S_n) + \langle \nabla \ell(w_t, S_n), v - w_t \rangle + \frac{1}{2} \langle H(w_t, S_n)(v - w_t), (v - w_t) \rangle, \quad (35)$$

where $H(w_t, S_n)$ can be either $H_{\text{qu}}(w_t, S_n)$ from Lemma 5.1 or $\nabla^2 \ell(w_t, S_n)$.

Let $\lambda > 0$, $\tilde{H}_t = \Psi_\lambda(H(w_t, S_n), \text{hessian modification})$, and $v_\lambda = w_t - \tilde{H}_t^{-1} \tilde{g}_t + \sigma_2 \|\tilde{g}_t\| \xi_t$ where $\sigma_2 > 0$ is a constant and $\xi_t \sim \mathcal{N}(0, I_d)$. Our goal here is to find $\lambda > 0$ as an approximate minimizer of $\mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)}[\phi(v_\lambda)]$. Note that we condition on the random variables w_t and \tilde{g}_t .

Algorithm 5 Newton Method with Double noise and adaptive min. eigenvalue

- 1: Input: training set $S_n \in \mathcal{Z}^n$, $\theta \in (0, 1)$ for dividing privacy budget for gradient vs SOI, $\gamma \in (0, 1)$ for dividing privacy budget for trace estimation, $\beta > 0$ as the coefficient for min. eig. value, privacy budget ρ -zCDP, initialization w_0 , number of iterations T , Hessian modification $\in \{\text{clip}, \text{add}\}$.
 - 2: Set $\sigma_1 = \frac{\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}$
 - 3: Set $\sigma_{\text{tr}} = \frac{\sqrt{T}}{4n\sqrt{2\theta\rho\gamma}}$
 - 4: **for** $t = 0, \dots, T-1$ **do**
 - 5: Query $\nabla \ell(w_t, S_n)$ and $H(w_t, S_n)$
 - 6: $\text{trace}_t = \max\{\text{trace}(H(w_t)) + \mathcal{N}(0, \sigma_{\text{tr}}^2 I_d), 0\}$
 - 7:
 - 8: $\lambda_{0,t} = \max\left\{\beta \cdot (\widetilde{\text{trace}_t})^{1/3} \left(\frac{T}{n^2(1-\gamma)\rho\theta}\right)^{1/3}, \frac{1}{n}\right\}$ \triangleright To prevent $\lambda_{0,t}$ makes σ_2 negative.
 - 9: **if** Hessian modification = Add **then**
 - 10: $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_{0,t}^2 + \lambda_{0,t})\sqrt{2(1-\gamma)\rho\theta}}$
 - 11: **else if** Hessian modification = Clip **then**
 - 12: $\sigma_2 = \frac{\sqrt{T}}{(4n\lambda_{0,t}^2 - \lambda_{0,t})\sqrt{2(1-\gamma)\rho\theta}}$
 - 13: $\tilde{H}_t = \Psi_{\lambda_0}(H(w_t, S_n), \text{Hessian modification})$
 - 14: $\tilde{g}_t = H(w_t) + \mathcal{N}(0, \sigma_1^2 I_d)$
 - 15: $w_{t+1} = w_t - \tilde{H}_t^{-1} \tilde{g}_t + \mathcal{N}(0, \|\tilde{g}_t\|^2 \sigma_2^2 I_d)$
 - 16: Output w_T .
-

We begin by expanding $\phi(v_\lambda)$ as follows:

$$\begin{aligned}
\mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)}[\phi(v_\lambda)] &= \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)}\left[\ell(w_t, S_n) + \left\langle \nabla \ell(w_t, S_n), -\tilde{H}_t^{-1} \tilde{g}_t + \sigma_2 \|\tilde{g}_t\| \xi_t \right\rangle \right. \\
&\quad \left. + \frac{1}{2}(-\tilde{H}_t^{-1} \tilde{g}_t + \sigma_2 \|\tilde{g}_t\| \xi_t)^\top H(w_t, S_n)(-\tilde{H}_t^{-1} \tilde{g}_t + \sigma_2 \|\tilde{g}_t\| \xi_t) \right] \\
&= \ell(w_t, S_n) + \left\langle \nabla \ell(w_t, S_n), -\tilde{H}_t^{-1} \tilde{g}_t \right\rangle + \frac{1}{2} \tilde{g}_t^\top (\tilde{H}_t^{-1})^\top H(w_t, S_n) \tilde{H}_t^{-1} \tilde{g}_t \\
&\quad + \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)}\left[\frac{1}{2} \sigma_2^2 \|\tilde{g}_t\|^2 \xi_t^\top H(w_t, S_n) \xi_t\right] \\
&= \ell(w_t, S_n) + \left\langle \nabla \ell(w_t, S_n), -\tilde{H}_t^{-1} \tilde{g}_t \right\rangle + \frac{1}{2} \tilde{g}_t^\top (\tilde{H}_t^{-1})^\top H(w_t, S_n) \tilde{H}_t^{-1} \tilde{g}_t \\
&\quad + \frac{\sigma_2^2 \|\tilde{g}_t\|^2}{2} \text{trace}(H(w_t, S_n)), \tag{36}
\end{aligned}$$

where we have used $\mathbb{E}[\xi_t] = 0$ and $\mathbb{E}[\xi_t^\top H(w_t, S_n) \xi_t] = \mathbb{E}[\text{trace}(\xi_t \xi_t^\top H(w_t, S_n))] = \text{trace}(H(w_t, S_n))$.

Consider the eigenvalue decomposition of $\tilde{H}_t \triangleq U \tilde{\Lambda} U^\top$ and $H(w_t, S_n) \triangleq U \Lambda U^\top$. Notice that by the definition of the adding and clipping operators in Definition 5.4, $H(w_t, S_n)$ and \tilde{H}_t share the same eigenvectors.

To approximate Equation (36), we assume $\tilde{g}_t \approx \nabla \ell(w_t, S_n)$. Then, by the change of variable $b = U^\top \tilde{g}_t$, we can rephrase Equation (36) as follows

$$\arg \min_{\lambda > 0} \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)}[\phi(v_\lambda)] \approx \arg \min_{\lambda > 0} \left\{ -b^\top \tilde{\Lambda}^{-1} b + \frac{1}{2} b^\top \tilde{\Lambda}^{-1} \Lambda \tilde{\Lambda}^{-1} b + \frac{\sigma_2^2 \|b\|^2}{2} \text{trace}(H(w_t, S_n)) \right\}. \tag{37}$$

Consider the eigenvalue modification using add operator. In this case $\tilde{H}_t = H(w_t, S_n) + \lambda I_d$. Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ and $\tilde{\Lambda} = \Lambda + \lambda I_d$. Also from Theorem B.2, $\sigma_2 = \frac{1}{(4n\lambda^2 + \lambda)\sqrt{2\rho_2}}$ where $\rho_2 > 0$

is the privacy budget. Setting these parameters in Equation (37), we get

$$h(\lambda) = \sum_{i=1}^d b_i^2 \left(\frac{-1}{\lambda_i + \lambda} + \frac{0.5\lambda_i}{(\lambda_i + \lambda)^2} \right) + \frac{\|b\|^2}{2} \left(\frac{1}{(4n\lambda^2 + \lambda)\sqrt{2\rho_2}} \right)^2 \text{trace}(H(w_t, S_n)).$$

By taking the derivative of $h(\lambda)$ and setting it to zero, we obtain

$$\frac{dh(\lambda^*)}{d\lambda^*} = 0 \Rightarrow \sum_{i=1}^d b_i^2 \frac{\lambda^*}{(\lambda_i + \lambda^*)^3} = \frac{\|\tilde{g}_t\|^2 \text{trace}(H(w_t, S_n))}{2\rho_2} \frac{1 + 8n\lambda^*}{(4n(\lambda^*)^2 + \lambda^*)^3}. \quad (38)$$

In many practical scenarios, the SOI matrix has zero eigenvalues. This observation motivates us to use the approximation $\frac{\lambda^*}{(\lambda_i + \lambda^*)^3} \approx \frac{1}{(\lambda^*)^2}$ for all $i \in [d]$. Let $\beta \in (0, 1)$ such that

$$\begin{aligned} \sum_{i=1}^d b_i^2 \frac{\lambda^*}{(\lambda_i + \lambda^*)^3} &= \frac{\beta}{(\lambda^*)^2} \sum_{i=1}^d b_i^2 \\ &= \frac{\beta}{(\lambda^*)^2} \|\tilde{g}_t\|^2. \end{aligned}$$

where we have used $\|b\| = \|U^\top \tilde{g}_t\| = \|\tilde{g}_t\|$. We can approximate Equation (38) by

$$\beta\lambda^* = \frac{\text{trace}(H(w_t, S_n))}{2\rho_2} \frac{1 + 8n\lambda^*}{(4n\lambda^* + 1)^3}.$$

Assume $n\lambda_0^* \gg 1$, we obtain

$$\lambda^* \approx \left(\frac{\text{trace}(H(w_t, S_n))}{n^2\rho_2} \right)^{\frac{1}{3}}. \quad (39)$$

The derivation for clip follows similarly using the smooth approximation of the max function: Let $m > 0$ be a constant. For all $i \in [d]$, we approximate $\max\{\lambda_i, \lambda\} \approx m^{-1} \log(\exp(m\lambda_i) + \exp(m\lambda))$.

B.6 Generalization of Algorithm 3 for convex, Lipschitz, and smooth loss functions

Algorithm 6 Generalization of Algorithm 3 for convex, L_0 -Lipschitz, and L_1 -smooth losses

- 1: Inputs: training set $S_n \in \mathcal{Z}^n$, $\lambda_0 > 0$, $\theta \in (0, 1)$, privacy budget ρ -zCDP, initialization w_0 , number of iterations T , hessian modification $\in \{\text{clip}, \text{add}\}$.
 - 2: Set $\sigma_1 = \frac{L_0\sqrt{T}}{n\sqrt{2\rho(1-\theta)}}$
 - 3: **if** hessian modification = Add **then**
 - 4: Condition: $n\lambda_0 > L_1$
 - 5: $\sigma_2 = \frac{L_1}{n\lambda_0^2 - \lambda_0 L_1} \cdot \frac{\sqrt{T}}{\sqrt{2\rho\theta}}$
 - 6: **else if** hessian modification = Clip **then**
 - 7: Condition: $n\lambda_0 > L_1 \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right)$ and $2\lambda_0 \leq L_1$
 - 8: $\sigma_2 = \frac{L_1 \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right)}{n\lambda_0^2 - \lambda_0 L_1 \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right)} \cdot \frac{\sqrt{T}}{\sqrt{2\rho\theta}}$.
 - 9: **for** $t = 0, \dots, T - 1$ **do**
 - 10: Query $\nabla \ell(w_t, S_n)$ and $\nabla^2 \ell(w_t, S_n)$
 - 11: $\tilde{H}_t = \Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n); \text{hessian modification})$
 - 12: $\tilde{g}_t = \nabla \ell(w_t, S_n) + \mathcal{N}(0, \sigma_1^2 I_d)$.
 - 13: $w_{t+1} = w_t - \tilde{H}_t^{-1} \tilde{g}_t + \mathcal{N}(0, \|\tilde{g}_t\|^2 \sigma_2^2 I_d)$.
 - 14: Output w_T .
-

Theorem B.6. For every convex, L_0 -Lipschitz, L_1 -smooth loss function $f(\cdot, \cdot)$, training set $S_n \in \mathcal{Z}^n$, initialization $w_0 \in \mathcal{W}$, $T \in \mathbb{N}$, $\rho \in \mathbb{R}_+$, and $\theta \in (0, 1)$, w_T in Algorithm 6 satisfies ρ -zCDP.

Proof. We follow the two-stage procedure of Theorem B.2 and Theorem B.3. Since the loss function is L_0 -Lipschitz, by setting,

$$\sigma_1 = \frac{L_0 \sqrt{T}}{n \sqrt{2\rho(1-\theta)}},$$

the mechanism in Line 12 of Algorithm 6 satisfies $\frac{\rho(1-\theta)}{T}$ -zCDP.

First consider the case with using add for the SOI modification. For the second step, following the same line as in the proof of Theorem B.2, we need to upper bound

$$\sup_{S_n \in \mathcal{Z}^n} \sup_{z_{n+1} \in \mathcal{Z}} \left\| \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{add}) \right]^{-1} - \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{add}) \right]^{-1} \right\|.$$

Let $A = \Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{add}) = \nabla^2 \ell(w_t, S_n) + \lambda_0 I_d$ and $B = \frac{1}{n} \nabla^2 f(w_t, z_{n+1})$. We need a lemma for the next step of the proof.

Lemma B.7. For every PSD matrix $A \in \mathbb{R}^{d \times d}$, $B \in \mathbb{R}^{d \times d}$, and $\lambda_0 \geq 0$ such that $\|A\|, \|B\| < \infty$ we have

$$\begin{aligned} \|\Psi_{\lambda_0}(A + B, \text{clip}) - \Psi_{\lambda_0}(A, \text{clip})\| &\leq \|B\| \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{\|A - \lambda_0 I_d\| + \|B\|}{\|B\|} \right) \right), \\ \|\Psi_{\lambda_0}(A + B, \text{add}) - \Psi_{\lambda_0}(A, \text{add})\| &\leq \|B\|. \end{aligned}$$

Proof. The Lipschitzness of $\Psi_{\lambda_0}(\cdot, \text{add})$ is obvious from Definition 5.4. We prove the result for $\Psi_{\lambda_0}(\cdot, \text{clip})$.

For a symmetric matrix $A \in \mathbb{R}^{d \times d}$, let $A = U \Lambda U^\top$ be the eigenvalue decomposition of A where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$. Then, define the absolute value of A as $|A| \triangleq U |\Lambda| U^\top \in \mathbb{R}^{d \times d}$ where $|\Lambda| = \text{diag}(|\lambda_1|, \dots, |\lambda_d|)$.

It is straightforward to see that

$$\Psi_{\lambda_0}(A, \text{clip}) = \frac{1}{2} (|A - \lambda_0 I_d| + A + \lambda_0 I_d).$$

Therefore,

$$\begin{aligned} &\|\Psi_{\lambda_0}(A + B, \text{clip}) - \Psi_{\lambda_0}(A, \text{clip})\| \\ &= \frac{1}{2} \left\| (|A + B - \lambda_0 I_d| + (A + B) + \lambda_0 I_d) - (|A - \lambda_0 I_d| + A + \lambda_0 I_d) \right\| \\ &= \frac{1}{2} \left\| |A + B - \lambda_0 I_d| - |A - \lambda_0 I_d| + B \right\| \\ &\leq \frac{1}{2} \left\| |A + B - \lambda_0 I_d| - |A - \lambda_0 I_d| \right\| + \frac{1}{2} \|B\|. \end{aligned} \tag{40}$$

Then, we invoke the result of [Kat73] which states that

$$\frac{1}{2} \left\| |A + B - \lambda_0 I_d| - |A - \lambda_0 I_d| \right\| \leq \frac{\|B\|}{\pi} \left(2 + \log \left(\frac{\|A - \lambda_0 I_d\| + \|B\|}{\|B\|} \right) \right) \tag{41}$$

Combining Equations (40) and (41) concludes the proof. \square

Using Lemma B.7 and Lemma B.4 we can write

$$\begin{aligned} &\left\| \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{add}) \right]^{-1} - \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{add}) \right]^{-1} \right\| \\ &= \left\| (A + B)^{-1} - A^{-1} \right\| \\ &\leq \frac{\|A^{-1}\|^2 \|B\|}{1 - \|A^{-1}\| \|B\|}. \end{aligned} \tag{42}$$

Let $n \geq L_1 \lambda_0^{-1}$. Notice that $\|A^{-1}\| \leq \lambda_0^{-1}$ and $\|B\| \leq L_1 n^{-1}$ because of the modification operator and the smoothness of the loss function. Using this observation, we can write

$$\begin{aligned} \frac{\|A^{-1}\|^2 \|B\|}{1 - \|A^{-1}\| \|B\|} &\leq \sup_{0 \leq x \leq L_1 n^{-1}} \frac{\|A^{-1}\|^2 x}{1 - \|A^{-1}\| x} \\ &= \|A^{-1}\|^2 \frac{L_1}{n - \|A^{-1}\| L_1}. \end{aligned}$$

Here the last step follows from the following fact: For every $a > 0$, $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) = \frac{x}{1-ax}$ is increasing for $x < \frac{1}{a}$. Then,

$$\begin{aligned} \|A^{-1}\|^2 \frac{L_1}{n - \|A^{-1}\| L_1} &\leq \sup_{0 \leq x \leq \lambda_0^{-1}} \frac{x^2 L_1}{n - x L_1} \\ &= \frac{L_1}{n \lambda_0^2 - \lambda_0 L_1}, \end{aligned}$$

where the last step follows from the following fact: $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) = \frac{x^2}{n-xL_1}$ is increasing in the interval $0 \leq x < nL_1^{-1}$. Therefore, we conclude that

$$\begin{aligned} &\sup_{S_n \in \mathcal{Z}^n} \sup_{z_{n+1} \in \mathcal{Z}} \left\| \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{add}) \right]^{-1} - \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{add}) \right]^{-1} \right\| \\ &\leq \frac{L_1}{n \lambda_0^2 - \lambda_0 L_1}. \end{aligned} \quad (43)$$

This shows that by setting

$$\sigma_2 = \frac{L_1 \sqrt{T}}{(n \lambda_0^2 - \lambda_0 L_1) \sqrt{2\rho\theta}},$$

the mechanism in Line 13 of Algorithm 6 is $\frac{\theta\rho}{T}$ -zCDP.

In each step of the algorithm we have two privatization step that satisfy $\frac{(1-\theta)\rho}{T}$ and $\frac{\theta\rho}{T}$. By the composition property of zCDP [BS16, Lemma 2.3], we conclude that w_T satisfies ρ -zCDP.

Next, we provide a privacy analysis for the clipping operator. We are interested in upper-bounding the following term

$$\sup_{S_n \in \mathcal{Z}^n} \sup_{z_{n+1} \in \mathcal{Z}} \left\| \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{clip}) \right]^{-1} - \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{clip}) \right]^{-1} \right\|.$$

Let

$$A = \Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{clip}), \quad B = \Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{clip}).$$

Then, using Lemma B.4 we can write

$$\begin{aligned} &\left\| \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{clip}) \right]^{-1} - \left[\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{clip}) \right]^{-1} \right\| \\ &\leq \frac{\|B - A\| \|A^{-1}\|^2}{1 - \|B - A\| \|A^{-1}\|}. \end{aligned} \quad (44)$$

Then, we invoke Lemma B.7 to write

$$\begin{aligned} &\|B - A\| \\ &\leq \frac{1}{n} \|\nabla^2 f(w_t, z_{n+1})\| \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n \|\nabla^2 \ell(w_t, S_n) - \lambda_0 I_d\| + \|\nabla^2 f(w_t, z_{n+1})\|}{\|\nabla^2 f(w_t, z_{n+1})\|} \right) \right) \\ &\leq \frac{1}{n} \|\nabla^2 f(w_t, z_{n+1})\| \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + \|f(w_t, z_{n+1})\|}{\|\nabla^2 f(w_t, z_{n+1})\|} \right) \right), \end{aligned}$$

where the last step follows from the smoothness of f and the assumption that $2\lambda_0 \leq L_1$. By the smoothness we have $\|\nabla^2 f(w_t, z_{n+1})\| \leq L_1$, therefore, to upper bound $\|B - A\|$ we can write

$$\begin{aligned} \|B - A\| &\leq \frac{1}{n} \|\nabla^2 f(w_t, z_{n+1})\| \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + \|\nabla^2 f(w_t, z_{n+1})\|}{\|\nabla^2 f(w_t, z_{n+1})\|} \right) \right) \\ &\leq \sup_{0 \leq y \leq L_1} \frac{y}{n} \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + y}{y} \right) \right) \\ &= \frac{L_1}{n} \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right) \triangleq \Delta. \end{aligned}$$

where the last step follows from the following technical lemma.

Lemma B.8. *For every $a > 0$, function $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = x \left(\log \frac{x+a}{x} \right)$ is increasing for $x > 0$.*

Proof. The derivative of f is given by $\frac{df(x)}{dx} = \log(1 + \frac{a}{x}) - \frac{a}{x+a}$. By using the inequality $\log(1+y) \geq \frac{y}{1+y}$ for $y > -1$, we can show that $\frac{df(x)}{dx} \geq 0$, as was to be shown. \square

Then, we can further upper bound Equation (44) as follows:

$$\begin{aligned} \frac{\|B - A\| \|A^{-1}\|^2}{1 - \|B - A\| \|A^{-1}\|} &\leq \frac{\Delta \|A^{-1}\|^2}{1 - \Delta \|A^{-1}\|} \\ &\leq \frac{\Delta}{\lambda_0^2 - \Delta \lambda_0}, \end{aligned}$$

where the last step follows from $\|A^{-1}\| \leq \lambda_0^{-1}$.

Therefore, we conclude that

$$\begin{aligned} &\sup_{S_n \in \mathcal{Z}^n} \sup_{z_{n+1} \in \mathcal{Z}} \left\| [\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n) + \frac{1}{n} \nabla^2 f(w_t, z_{n+1}), \text{clip})]^{-1} - [\Psi_{\lambda_0}(\nabla^2 \ell(w_t, S_n), \text{clip})]^{-1} \right\| \\ &\leq \frac{L_1 \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right)}{n\lambda_0^2 - L_1\lambda_0 \left(\frac{2}{\pi} + \frac{1}{2} + \frac{1}{\pi} \log \left(\frac{n(L_1 - \lambda_0) + L_1}{L_1} \right) \right)}. \end{aligned}$$

The rest of the proof is similar to the proof of the Hessian modification using the adding operator. \square

B.7 Suboptimality Gap for Logistic Loss and $\|\cdot\|_V$

From Lemma 5.1, since $\nabla \ell_{\text{LL}}(w^*, S_n) = 0$ we have

$$\ell_{\text{LL}}(w, S_n) \leq \ell_{\text{LL}}(w^*, S_n) + (w - w^*)^\top \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\tanh(\langle x_i, w^* \rangle / 2)}{4 \langle x_i, w^* \rangle} \right) (w - w^*).$$

By definition of V , $V x_i = x_i$. Therefore,

$$\ell_{\text{LL}}(w, S_n) \leq \ell_{\text{LL}}(w^*, S_n) + (w - w^*)^\top V \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\tanh(\langle x_i, w^* \rangle / 2)}{4 \langle x_i, w^* \rangle} \right) V (w - w^*).$$

Since $\left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\tanh(\langle x_i, w^* \rangle / 2)}{2 \langle x_i, w^* \rangle} \right) \preceq \frac{1}{4} I_d$ and $V^2 = V$, we have

$$\begin{aligned} \ell_{\text{LL}}(w, S_n) &\leq \ell_{\text{LL}}(w^*, S_n) + \frac{1}{8} (w - w^*)^\top V (w - w^*) \\ &= \frac{1}{8} \|w - w^*\|_V^2, \end{aligned}$$

which was to be shown.

B.8 Proof of Theorem 5.8

We start this section by recalling some of the well-known properties of *Mahalanobis semi-norm*.

Lemma B.9. *Let $A \in \mathbb{R}^{d \times d}$ be a positive semi-definite matrix. For every $x, y \in \mathbb{R}^d$ define $\langle x, y \rangle_A \triangleq x^\top A y$ and $\|x\|_A \triangleq \sqrt{\langle x, x \rangle_A}$. Then, the following holds:*

- For every $x \in \mathbb{R}^d$, we have $\|x\|_A \geq 0$.
- for every $\alpha \in \mathbb{R}$, we have $\|\alpha x\|_A = |\alpha| \|x\|_A$.
- For every $x, y \in \mathbb{R}^d$, $\|x + y\|_A \leq \|x\|_A + \|y\|_A$.
- For every $x, y \in \mathbb{R}^d$, we have $|\langle x, y \rangle_A| \leq \|x\|_A \|y\|_A$.

Lemma B.10. *Let $A \in \mathbb{R}^{d \times d}$ be a positive semi-definite matrix. Then, for every $M \in \mathbb{R}^{d \times d}$ define*

$$\|M\|_A \triangleq \sup_{x \in \mathbb{R}^d} \frac{\|Mx\|_A}{\|x\|_A}.$$

Then, $\|M\|_A \|x\|_A \geq \|Mx\|_A$. Also, for $M, M' \in \mathbb{R}^{d \times d}$, we have $\|M + M'\|_A \leq \|M\|_A + \|M'\|_A$.

The following lemma summarizes some of the properties of the logistic loss that will be used in the proof.

Lemma B.11. *Fix $n \in \mathbb{N}$ and data set $S_n = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathbb{R}^d \times \{-1, +1\})^n$. Let $V \in \mathbb{R}^{d \times d}$ denote the orthogonal projection matrix on the linear subspace spanned by $\{x_1, \dots, x_n\}$. Then, the following holds:*

1. For every $w \in \mathbb{R}^d$ and $w' \in \mathbb{R}^d$

$$\|\nabla^2 \ell_{\text{LL}}(w', S_n) - \nabla^2 \ell_{\text{LL}}(w, S_n)\|_V \leq 0.1 \cdot \|w' - w\|_V.$$

2. For every $w \in \mathbb{R}^d$, $u^\top \nabla^2 \ell_{\text{LL}}(w, S_n) u = 0 \Leftrightarrow Vu = 0$. In words, the directions of zero eigenvalue of $\nabla^2 \ell_{\text{LL}}(w, S_n)$ are orthogonal to the linear subspace spanned by $\{x_1, \dots, x_n\}$.
3. For every $w \in \mathbb{R}^d$, the eigenvectors of $\nabla^2 \ell_{\text{LL}}(w, S_n)$ corresponding to non-zero eigenvalue lie in the linear subspace spanned by $\{x_1, \dots, x_n\}$.
4. Fix $w \in \mathbb{R}^d$ and consider the eigenvalue decomposition of $\nabla^2 \ell_{\text{LL}}(w, S_n)$ as $\sum_{i=1}^d \lambda_i u_i u_i^\top$ where $\{\lambda_i \in \mathbb{R} : i \in [d]\}$ and $\{u_i \in \mathbb{R}^d : i \in [d]\}$ denote the eigenvalues and eigenvectors. Let $\lambda_{\min, w} = \min\{\lambda_i : \lambda_i > 0\}$. Then,

$$\lambda_{\min, w} = \min_{u \in \text{span}\{x_1, \dots, x_n\}, \|u\|=1} u^\top \nabla^2 \ell_{\text{LL}}(w, S_n) u.$$

Also,

$$|\lambda_{\min, w} - \lambda_{\min, w'}| \leq \|\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)\|_V.$$

5. Let $\lambda_0 > 0$. For every $w \in \mathbb{R}^d$,

$$\|\Psi_{\lambda_0}(\nabla^2 \ell_{\text{LL}}(w, S_n), \text{clip})\|_V = \frac{1}{\max\{\lambda_0, \lambda_{\min, w}\}}, \quad \|\Psi_{\lambda_0}(\nabla^2 \ell_{\text{LL}}(w, S_n), \text{add})\|_V = \frac{1}{\lambda_0 + \lambda_{\min, w}},$$

where $\Psi_{\lambda_0}(\cdot, \cdot)$ is defined in Definition 5.4.

Proof. For Part 1, from Equation (5), we know that for every w

$$\nabla^2 \ell_{\text{LL}}(w, S_n) = \frac{1}{n} \sum_{i=1}^n \frac{x_i x_i^\top}{(\exp(-\langle w, x_i \rangle / 2) + \exp(\langle w, x_i \rangle / 2))^2}$$

For every $i \in [n]$, let $g : \mathbb{R} \rightarrow \mathbb{R}$ be $g(t) = \frac{1}{(\exp(-t/2) + \exp(t/2))^2}$. Then

$$\begin{aligned} \|\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)\|_V &= \frac{1}{n} \left\| \sum_{i=1}^n x_i x_i^\top (g(\langle w, x_i \rangle) - g(\langle w', x_i \rangle)) \right\|_V \\ &\leq \frac{1}{n} \sum_{i=1}^n \|x_i x_i^\top\|_V \max_{i \in [n]} |g(\langle w, x_i \rangle) - g(\langle w', x_i \rangle)| \\ &\leq \max_{i \in [n]} \|g(\langle w, x_i \rangle) - g(\langle w', x_i \rangle)\|, \end{aligned}$$

where the second and third steps follow from Lemma B.10 and $\|x_i x_i^\top\|_V = \|x_i x_i^\top\|_2 \leq 1$. It is easy to show there exists $L_2 < 0.1$ such that g is L_2 -Lipschitz. Therefore,

$$\begin{aligned} |g(\langle w, x_i \rangle) - g(\langle w', x_i \rangle)| &\leq L_2 |\langle w - w', x_i \rangle| \\ &= L_2 |\langle w - w', x_i \rangle_V| \\ &\leq \|w - w'\|_V, \end{aligned}$$

where the second step follows from $\langle w - w', x_i \rangle = (w - w')^\top x_i = (w - w')^\top V x_i$ since $x_i = V x_i$ by definition. Also, the last step follows from Lemma B.9.

For Part 2, by Equation (5), we have

$$u^\top \nabla^2 \ell_{\text{LL}}(w, S_n) u = \frac{1}{n} \sum_{i=1}^n \frac{(x_i^\top u)^2}{(\exp(-\langle w, x_i \rangle/2) + \exp(\langle w, x_i \rangle/2))^2}$$

Notice that every summand is positive, therefore, given $u \in \mathbb{R}^d$ such that $u^\top \nabla^2 \ell_{\text{LL}}(w, S_n) u = 0$ implies that for every $i \in [n]$, $x_i^\top u = 0$. The other direction is obvious.

The proof of Part 3 follows from the definition of eigenvalues. Let $u \in \mathbb{R}^d$ be an eigenvector corresponding to eigenvalue of $\lambda > 0$, then

$$\nabla^2 \ell_{\text{LL}}(w, S_n) u = \lambda u \Rightarrow \sum_{i=1}^n \frac{x_i^\top u}{n \lambda (\exp(-\langle w, x_i \rangle/2) + \exp(\langle w, x_i \rangle/2))^2} x_i = u,$$

which shows that u is a linear combination of x_i s.

For Part 4, the first statement is a corollary of Part 3. For the second statement, let $u \in \mathbb{R}^d$ be the eigenvector corresponding to $\lambda_{\min, w}$. Then,

$$\begin{aligned} \lambda_{\min, w'} - \lambda_{\min, w} &= \min_{u' \in \text{span}\{x_1, \dots, x_n\}, \|u'\|=1} (u')^\top \nabla^2 \ell_{\text{LL}}(w', S_n) (u') - u^\top \nabla^2 \ell_{\text{LL}}(w', S_n) u \\ &\leq u^\top \nabla^2 \ell_{\text{LL}}(w, S_n) u - u^\top \nabla^2 \ell_{\text{LL}}(w', S_n) u \\ &= u^\top (\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)) u \\ &= u^\top V (\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)) u \\ &\leq \|u\|_V \|\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)\|_V \|u\|_V \\ &\leq \|\nabla^2 \ell_{\text{LL}}(w, S_n) - \nabla^2 \ell_{\text{LL}}(w', S_n)\|_V. \end{aligned}$$

Part 4 is based on the definition of the matrix norm in Lemma B.10 and Parts 2, 3. \square

Let $S_n = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathbb{R}^d \times \{-1, +1\})^n$. Let $V \in \mathbb{R}^{d \times d}$ denote the orthogonal projection matrix on the linear subspace spanned by $\{x_1, \dots, x_n\}$. Let w_t^* denote the minimizer of the empirical loss. We assume it exists and $\nabla \ell_{\text{LL}}(w_t^*, S_n) = 0$. To reduce notation clutter we drop S_n . Let $\xi_1 \sim \mathcal{N}(0, I_d)$ and $\xi_2 \sim \mathcal{N}(0, I_d)$. We can rephrase the update rule of Algorithm 3 as

$$\begin{aligned} \|w_{t+1} - w^*\|_V^2 &= \left\| w_t - w^* - \tilde{H}_t^{-1} (\nabla \ell_{\text{LL}}(w_t) + \sigma_1 \xi_1) + \|\nabla \ell_{\text{LL}}(w_t) + \sigma_1 \xi_1\|_2 \sigma_2 \xi_2 \right\|_V^2 \\ &= \left\| w_t - w^* - \tilde{H}_t^{-1} \nabla \ell_{\text{LL}}(w_t) \right\|_V^2 + \left\| \sigma_1 \tilde{H}_t^{-1} \xi_1 - \|\nabla \ell_{\text{LL}}(w_t) + \sigma_1 \xi_1\|_2 \sigma_2 \xi_2 \right\|_V^2 \\ &\quad - 2 \left\langle w_t - w^* - \tilde{H}_t^{-1} \nabla \ell_{\text{LL}}(w_t), \sigma_1 \tilde{H}_t^{-1} \xi_1 - \|\nabla \ell_{\text{LL}}(w_t) + \sigma_1 \xi_1\|_2 \sigma_2 \xi_2 \right\rangle_V. \end{aligned} \tag{45}$$

In the next step we analyze the first term in Equation (45):

$$\begin{aligned}
& \left\| w_t - \tilde{H}_t^{-1} \nabla \ell_{\text{LL}}(w_t) - w^* \right\|_V^2 \\
&= \left\| w_t - w^* - \tilde{H}_t^{-1} (\nabla \ell_{\text{LL}}(w_t) - \nabla \ell_{\text{LL}}(w^*)) \right\|_V^2 \\
&= \left\| w_t - w^* - \tilde{H}_t^{-1} \left(\int_0^1 \nabla^2 \ell_{\text{LL}}(w^* + \tau(w_t - w^*)) (w_t - w^*) d\tau \right) \right\|_V^2 \\
&= \left\| w_t - w^* - \tilde{H}_t^{-1} \left(\int_0^1 [\nabla^2 \ell_{\text{LL}}(w^* + \tau(w_t - w^*)) - \nabla^2 \ell_{\text{LL}}(w_t) + \nabla^2 \ell_{\text{LL}}(w_t)] (w_t - w^*) d\tau \right) \right\|_V^2.
\end{aligned}$$

For every $w \in \mathbb{R}^d$ and $\tau \in [0, 1]$, let $\Delta_\tau(w) = \nabla^2 \ell_{\text{LL}}(w^* + \tau(w - w^*)) - \nabla^2 \ell_{\text{LL}}(w)$. We write

$$\begin{aligned}
& \left\| w_t - \tilde{H}_t^{-1} \nabla \ell_{\text{LL}}(w_t) - w^* \right\|_V^2 \\
&\leq \left\| w_t - w^* - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) (w_t - w^*) \right\|_V^2 + \left\| \tilde{H}_t^{-1} \left(\int_0^1 \Delta_\tau(w_t) (w_t - w^*) d\tau \right) \right\|_V^2 \\
&\quad + 2 \left\| w_t - w^* - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) (w_t - w^*) \right\|_V \left\| \tilde{H}_t^{-1} \left(\int_0^1 \Delta_\tau(w_t) (w_t - w^*) d\tau \right) \right\|_V \\
&\leq \left\| I - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) \right\|_V^2 \|w_t - w^*\|_V^2 + \left\| \tilde{H}_t^{-1} \right\|_V^2 \left(\int_0^1 \|\Delta_\tau(w_t)\|_V d\tau \right)^2 \|w_t - w^*\|_V^2 \\
&\quad + 2 \left\| I - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) \right\|_V \left\| \tilde{H}_t^{-1} \right\|_V \left(\int_0^1 \|\Delta_\tau(w_t)\|_V d\tau \right) \|w_t - w^*\|_V^2.
\end{aligned}$$

Here, we repeatedly use the properties of $\|\cdot\|_V$ from Lemmas B.9 and B.10.

Consider the eigenvalue decomposition of $\nabla^2 \ell_{\text{LL}}(w_t) = \sum_{i=1}^d \lambda_i u_i u_i^\top$ where some λ_i may be zero since we do not assume that $\nabla^2 \ell_{\text{LL}}(w_t)$ is a full-rank matrix. Let $\lambda_{\min, t}$ be the smallest *non-zero* eigenvalue of $\nabla^2 \ell_{\text{LL}}(w_t)$. Then, by Definition 5.4

$$\tilde{H}_t^{-1} = \begin{cases} \sum_{i=1}^d \frac{1}{\max\{\lambda_i, \lambda_0\}} u_i u_i^\top & \text{if Hessian modification is clip,} \\ \sum_{i=1}^d \frac{1}{\lambda_i + \lambda_0} u_i u_i^\top & \text{if Hessian modification is add,} \end{cases} \quad (46)$$

and

$$I - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) = \begin{cases} \sum_{i: \lambda_i < \lambda_0} \left(1 - \frac{\lambda_i}{\lambda_0}\right) u_i u_i^\top & \text{if Hessian modification is clip,} \\ \sum_{i=1}^d \frac{\lambda_i}{\lambda_i + \lambda_0} u_i u_i^\top & \text{if Hessian modification is add,} \end{cases} \quad (47)$$

Therefore from Equation (46), Equation (47), and Lemma B.11,

$$\left\| \tilde{H}_t^{-1} \right\|_V = \begin{cases} \frac{1}{\max\{\lambda_0, \lambda_{\min, t}\}} & \text{if Hessian modification is clip,} \\ \frac{1}{\lambda_0 + \lambda_{\min, t}} & \text{if Hessian modification is add,} \end{cases}$$

and

$$\left\| I - \tilde{H}_t^{-1} \nabla^2 \ell_{\text{LL}}(w_t) \right\|_V = \begin{cases} 1 - \frac{\min\{\lambda_0, \lambda_{\min, t}\}}{\lambda_0} & \text{if Hessian modification is clip,} \\ 1 - \frac{\lambda_{\min, t}}{\lambda_0 + \lambda_{\min, t}} & \text{if Hessian modification is add.} \end{cases}$$

Also, from Lemma B.11,

$$\begin{aligned} \int_0^1 \|\Delta_\tau(w_t)\|_V d\tau &= \int_0^1 \|\nabla^2 \ell_{LL}(w^* + \tau(w_t - w^*)) - \nabla^2 \ell_{LL}(w_t)\|_V d\tau \\ &\leq \frac{0.1}{2} \|w_t - w^*\|_V, \end{aligned} \quad (48)$$

Therefore,

$$\begin{aligned} \|w_t - \tilde{H}_t^{-1} \nabla \ell_{LL}(w_t) - w^*\|_V^2 &\leq \|I - \tilde{H}_t^{-1} \nabla^2 \ell_{LL}(w_t)\|_V^2 \|w_t - w^*\|_V^2 \\ &+ 0.1 \cdot \|I - \tilde{H}_t^{-1} \nabla^2 \ell_{LL}(w_t)\|_V \|\tilde{H}_t^{-1}\|_V \|w_t - w^*\|_V^3 + \frac{(0.1)^2}{4} \|\tilde{H}_t^{-1}\|_V^2 \|w_t - w^*\|_V^4. \end{aligned} \quad (49)$$

Consider the second term in Equation (45). Using the facts that $\mathbb{E}[\xi_1] = \mathbb{E}[\xi_2] = 0$, $\xi_1 \perp \xi_2$, $(\xi_1, \xi_2) \perp w_t$, and $\mathbb{E}[\|\xi_1\|_V^2] = \mathbb{E}[\|\xi_2\|_V^2] = \text{rank}$, we obtain

$$\begin{aligned} &\mathbb{E}_t \left[\left\| \sigma_1 \tilde{H}_t^{-1} \xi_1 - \|\nabla \ell_{LL}(w_t) + \sigma_1 \xi_1\|_2 \sigma_2 \xi_2 \right\|_V^2 \right] \\ &= \sigma_1^2 \mathbb{E}_t \left[\left\| \tilde{H}_t^{-1} \xi_1 \right\|_V^2 \right] + \sigma_2^2 \mathbb{E}_t \left[\|\nabla \ell_{LL}(w_t) + \sigma_1 \xi_1\|_2^2 \|\xi_2\|_V^2 \right] \\ &= \sigma_1^2 \left\| \tilde{H}_t^{-1} \right\|_V^2 \mathbb{E}_t \left[\|\xi_1\|_V^2 \right] + \sigma_2^2 \mathbb{E}_t \left[\|\nabla \ell_{LL}(w_t) + \sigma_1 \xi_1\|_2^2 \right] \mathbb{E}_t \left[\|\xi_2\|_V^2 \right] \\ &\leq \sigma_1^2 \left\| \tilde{H}_t^{-1} \right\|_V^2 \mathbb{E}_t \left[\|\xi_1\|_V^2 \right] + \sigma_2^2 \|\nabla \ell_{LL}(w_t)\|_2^2 \mathbb{E}_t \left[\|\xi_2\|_V^2 \right] \\ &+ \sigma_2^2 \sigma_1^2 \mathbb{E}_t \left[\|\xi_1\|^2 \right] \mathbb{E}_t \left[\|\xi_2\|_V^2 \right] \\ &= \sigma_1^2 \left\| \tilde{H}_t^{-1} \right\|_V^2 \text{rank} + \sigma_2^2 \|w_t - w^*\|_V^2 \text{rank} + \sigma_2^2 \sigma_1^2 d \text{rank}. \end{aligned} \quad (50)$$

Notice that the expectation of the third term in Equation (45) is zero. By combining Equation (45), Equation (49), and Equation (50) we obtain

$$\begin{aligned} \mathbb{E}_t \left[\|w_{t+1} - w^*\|_V^2 \right] &\leq \left(\left\| I - \tilde{H}_t^{-1} \nabla^2 \ell_{LL}(w_t) \right\|_V^2 + \sigma_2^2 \cdot \text{rank} \right) \|w_t - w^*\|_V^2 \\ &+ 0.1 \cdot \left\| I - \tilde{H}_t^{-1} \nabla^2 \ell_{LL}(w_t) \right\|_V \|\tilde{H}_t^{-1}\|_V \|w_t - w^*\|_V^3 + \frac{(0.1)^2}{4} \|\tilde{H}_t^{-1}\|_V^2 \|w_t - w^*\|_V^4 \\ &+ \sigma_1^2 \left\| \tilde{H}_t^{-1} \right\|_V^2 \text{rank} + \sigma_2^2 \sigma_1^2 d \text{rank}. \end{aligned} \quad (51)$$

Finally, setting the values of σ_1 and σ_2 from Algorithm 3 completes the proof.

B.9 Global Convergence of QU-clip and QU-add

Theorem B.12 (Global Convergence Guarantee of QU-clip and QU-add). *Let λ_{\min}^* be the minimum non-zero eigenvalue of $\nabla^2 \ell_{LL}(w^*, S_n)$, ρ be the privacy budget (in zCDP) per iteration, $\delta_t = \ell_{LL}(w_t, S_n) - \ell_{LL}(w^*, S_n)$ be the suboptimality gap at iteration t , and $\lambda_{\max, t}$ be the maximum eigenvalue of $H_{qu}(w_t, S_n)$ from Lemma 5.1. Let*

$$\tilde{\lambda}_{\max, t} = \begin{cases} \max\{\lambda_0, \lambda_{\max, t}\} & \text{if SOI modification is clip,} \\ (\lambda_0 + \lambda_{\max, t}) & \text{if SOI modification is add.} \end{cases}$$

Then, if $\|\nabla \ell_{LL}(w_t, S_n)\| \geq \frac{3\lambda_{\min}^*}{4}$

$$\mathbb{E}_t [\delta_{t+1}] \leq \delta_t - \frac{9}{8} \lambda_{\min}^* \cdot \nu + \Delta, \quad (52)$$

Also, if $\|\nabla \ell_{LL}(w_t, S_n)\| < \frac{3\lambda_{\min}^*}{4}$, we have

$$\mathbb{E}_t [\delta_{t+1}] \leq (1 - \nu) \delta_t + \Delta, \quad (53)$$

where

$$\begin{aligned} \nu &= \frac{\lambda_{\min}^*}{4\tilde{\lambda}_{\max,t}} - \frac{\lambda_{\max}\lambda_{\min}^*\text{rank}}{8\rho\theta(4n\lambda_0^2 - \lambda_0)^2}, \quad \Delta = O\left(\frac{\text{rank}}{4\lambda_0\rho\theta(1-\theta)n^2}\right), \quad \text{if SOI modification is clip.} \\ \nu &= \frac{\lambda_{\min}^*}{4\tilde{\lambda}_{\max,t}} - \frac{\lambda_{\max}\lambda_{\min}^*\text{rank}}{8\rho\theta(4n\lambda_0^2 + \lambda_0)^2}, \quad \Delta = O\left(\frac{\text{rank}}{4\lambda_0\rho\theta(1-\theta)n^2}\right), \quad \text{if SOI modification is add.} \end{aligned}$$

Proof. We begin the proof by quoting a result from [Bac14]. Note that the statement in [Bac14] is stated in terms of $\|\cdot\|_2$, but the extension to the norm induced by V is straightforward.

Lemma B.13 ([Bac14, Lemma 9]). *Let S_n be a training set and $w^* = \arg \min \ell_{\text{LL}}(w, S_n)$. Let λ_{\min}^* be the minimum non-zero eigenvalue of $\nabla^2 \ell_{\text{LL}}(w^*, S_n)$. Then, for every w such that $\|\nabla \ell_{\text{LL}}(w, S_n)\| \leq \frac{3}{4}\lambda_{\min}^*$, we have*

$$\ell_{\text{LL}}(w, S_n) - \ell_{\text{LL}}(w^*, S_n) \leq 2 \frac{\|\nabla \ell_{\text{LL}}(w, S_n)\|_V^2}{\lambda_{\min}^*}.$$

Define the $\mathbb{E}_t[\cdot]$ as the conditional expectation conditioned on the history up to time t , i.e., $\{w_0, \dots, w_t\}$. We can write by Lemma 5.1

$$\ell_{\text{LL}}(w_{t+1}) \leq \ell_{\text{LL}}(w_t) + \langle \nabla \ell_{\text{LL}}(w_t), w_{t+1} - w_t \rangle + \frac{1}{2}(w_{t+1} - w_t)^\top H_{\text{qu},t}(w_{t+1} - w_t). \quad (54)$$

Let $g_t = \nabla \ell_{\text{LL}}(w_t)$, $\xi_1 \sim \mathcal{N}(0, I_d)$, $\xi_2 \sim \mathcal{N}(0, I_d)$. Then, by the definition of the update rule,

$$w_{t+1} - w_t = -\tilde{H}_t^{-1}(g_t + \sigma_1 \xi_1) + \|g_t + \sigma_1 \xi_1\| \sigma_2 \xi_2.$$

We use \tilde{H}_t to denote $\Psi_{\lambda_0}(H_{\text{qu},t}, \text{SOI modification})$. Then,

$$\begin{aligned} \mathbb{E}_t[\langle \nabla \ell_{\text{LL}}(w_t), w_{t+1} - w_t \rangle] &= \mathbb{E}_t\left[\left\langle g_t, -\tilde{H}_t^{-1}(g_t + \sigma_1 \xi_1) + \|g_t + \sigma_1 \xi_1\| \sigma_2 \xi_2 \right\rangle\right] \\ &= -g_t^\top \tilde{H}_t^{-1} g_t, \end{aligned} \quad (55)$$

where the last step follows from ξ_1 and ξ_2 being independent of w_t . For the third term on RHS of Equation (54), using $\mathbb{E}[\xi_1] = \mathbb{E}[\xi_2] = 0$ and $\xi_1 \perp \xi_2$, we can write

$$\begin{aligned} &\mathbb{E}_t[(w_{t+1} - w_t)^\top H_{\text{qu},t}(w_{t+1} - w_t)] \\ &= \mathbb{E}_t\left[(g_t + \sigma_1 \xi_1)^\top \tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1} (g_t + \sigma_1 \xi_1)\right] + \sigma_2^2 \mathbb{E}_t\left[\|g_t + \sigma_1 \xi_1\|^2 \xi_2^\top H_{\text{qu},t} \xi_2\right] \\ &= g_t^\top \tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1} g_t + \sigma_1^2 \mathbb{E}_t\left[\xi_1^\top \tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1} \xi_1\right] \\ &\quad + \left(\sigma_2^2 \|g_t\|^2 + \sigma_1^2 \sigma_2^2 \mathbb{E}_t[\|\xi_1\|^2]\right) \mathbb{E}_t[\xi_2^\top H_{\text{qu},t} \xi_2]. \end{aligned}$$

By the definition of the modification operators in Definition 5.4 we have $\tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1} \preceq \tilde{H}_t^{-1}$. Also, by the fact that for a symmetric matrix A and $\xi \sim \mathcal{N}(0, I_d)$, it holds $\mathbb{E}[\xi^\top A \xi] = \text{trace}(A)$, we can write

$$\begin{aligned} &\frac{1}{2} \mathbb{E}_t[(w_{t+1} - w_t)^\top H_{\text{qu},t}(w_{t+1} - w_t)] \\ &\leq -\frac{1}{2} g_t^\top \tilde{H}_t^{-1} g_t + \frac{1}{2} \sigma_1^2 \text{trace}(\tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1}) + \frac{1}{2} \sigma_2^2 \|g_t\|^2 \text{trace}(H_{\text{qu},t}) + \frac{1}{2} \sigma_1^2 \sigma_2^2 d \text{trace}(H_{\text{qu},t}) \\ &\leq -\frac{1}{2} g_t^\top \tilde{H}_t^{-1} g_t + \frac{1}{2} \sigma_1^2 \frac{\text{rank}}{\lambda_0} + \frac{\lambda_{\max,t}}{2} \sigma_2^2 \|g_t\|^2 \text{rank} + \frac{\lambda_{\max,t}}{2} \sigma_1^2 \sigma_2^2 \cdot d \cdot \text{rank}, \end{aligned} \quad (56)$$

where the last line follows from $\text{trace}(\tilde{H}_t^{-1} H_{\text{qu},t} \tilde{H}_t^{-1}) \leq \frac{\text{rank}}{\lambda_0}$ and $\text{trace}(H_{\text{qu},t}) \leq \text{rank} \cdot \lambda_{\max,t}$ where the maximum eigenvalue of $H_{\text{qu},t}$ is denoted by $\lambda_{\max,t}$. Also,

$$\tilde{H}_t \preceq \tilde{\lambda}_{\max,t} I_d \triangleq \begin{cases} \max\{\lambda_0, \lambda_{\max,t}\} I_d & \text{if Hessian modification is clip,} \\ (\lambda_0 + \lambda_{\max,t}) I_d & \text{if Hessian modification is add.} \end{cases} \quad (57)$$

Therefore,

$$\begin{aligned}\mathbb{E}_t [\ell_{\text{LL}}(w_{t+1}) - \ell_{\text{LL}}(w_t)] &\leq -\frac{1}{2}g_t^\top \tilde{H}_t^{-1}g_t + \frac{1}{2}\sigma_1^2 \frac{\text{rank}}{\lambda_0} + \frac{\lambda_{\max,t}}{2}\sigma_2^2 \|g_t\|^2 \text{rank} + \frac{\lambda_{\max,t}}{2}\sigma_1^2\sigma_2^2 \cdot d \cdot \text{rank} \\ &\leq -\frac{1}{2}\|g_t\|^2 \left(\frac{1}{\tilde{\lambda}_{\max,t}} - \sigma_2^2 \cdot \text{rank} \cdot \lambda_{\max,t} \right) + \frac{1}{2}\sigma_1^2 \frac{\text{rank}}{\lambda_0} + \frac{\lambda_{\max,t}}{2}\sigma_1^2\sigma_2^2 \cdot d \cdot \text{rank},\end{aligned}\quad (58)$$

where the last step follows from the fact that for every $u \in \mathbb{R}^d$, $u^\top \tilde{H}_t^{-1}u \geq \frac{1}{\tilde{\lambda}_{\max,t}} \|u\|^2$.

In the last step, we will use Lemma B.13. Let λ_{\min}^* be the minimum non-zero eigenvalue of $\nabla^2 \ell_{\text{LL}}(w^*, S_n)$. Since g_t is a linear combination of x_i s (See Equation (5)), we have $\|g_t\|_2 = \|g_t\|_V$. Consider two cases: Case 1) $\|g_t\|_V > \frac{3}{4}\lambda_{\min}^*$, Case 2) $\|g_t\|_V \leq \frac{3}{4}\lambda_{\min}^*$.

For Case 1, we can simplify Equation (58) as follows

$$\begin{aligned}\mathbb{E}_t [\ell_{\text{LL}}(w_{t+1}) - \ell_{\text{LL}}(w^*)] &\leq \ell_{\text{LL}}(w_t) - \ell_{\text{LL}}(w^*) \\ &\quad - \frac{9}{32}(\lambda_{\min}^*)^2 \left(\frac{1}{\tilde{\lambda}_{\max,t}} - \sigma_2^2 \cdot \text{rank} \cdot \lambda_{\max,t} \right) + \frac{1}{2}\sigma_1^2 \frac{\text{rank}}{\lambda_0} + \frac{\lambda_{\max,t}}{2}\sigma_1^2\sigma_2^2 d \text{rank}.\end{aligned}$$

For the second case, from Lemma B.13 we have

$$\begin{aligned}\mathbb{E}_t [\ell_{\text{LL}}(w_{t+1}) - \ell_{\text{LL}}(w^*)] &\leq \left[1 - \frac{\lambda_{\min}^*}{4\tilde{\lambda}_{\max,t}} + \frac{\sigma_2^2 \text{rank} \lambda_{\max,t} \lambda_{\min}^*}{4} \right] (\ell_{\text{LL}}(w_t) - \ell_{\text{LL}}(w^*)) + \sigma_1^2 \frac{\text{rank}}{2\lambda_0} + \frac{\lambda_{\max,t}}{2}\sigma_1^2\sigma_2^2 d \text{rank}.\end{aligned}\quad (59)$$

The stated results follow from setting σ_1 and σ_2 . \square

C Appendix of Section 6

In this section, we present the details of the implementation and additional experiment results. An implementation of our proposed optimization algorithms can be found in github.com/tensorflow/privacy/tree/master/research/dp_newton.

C.1 Subsampled variant of Our Algorithm

In this section, we show how to extend Algorithm 3 to the minibatch version.

Let's assume we have m queries, denoted as $q_i : \mathcal{Z}^* \rightarrow \mathbb{R}^d$, where $i \in [m]$, and each query has an ℓ_2 sensitivity of one. We want to sequentially compose these queries using the Sampled Gaussian Mechanism (SGM), which combines subsampling and additive Gaussian noise [MTZ19]. To determine the appropriate noise level for achieving the desired privacy, we assume we have an access to function `get_noise_multiplier` which takes as input the total privacy budget, m , and the subsampling probability and outputs the minimum standard deviation of noise for Gaussian Mechanism to achieve the required privacy. Such a function can be found in various publicly available privacy libraries.

Theorem C.1. *For every training set $S_n \in \mathcal{Z}^n$, $\lambda_0 > 0$, $\theta \in (0, 1)$, privacy budget (ε, δ) -DP, initialization w_0 , number of iterations T , SOI modification $\in \{\text{clip}, \text{add}\}$, and sampling rates $p_g, p_H \in (0, 1)$ for gradient and SOI, the output of Algorithm 7, i.e., w_T satisfies (ε, δ) -DP.*

Proof. In Algorithm 7, we have two types of SGMs, 1) gradient SGM, 2) SOI SGM. The result from [Lyu22; VZ22] indicate that we can *interleave* these mechanisms in a way that we have T composition of only gradient SGM followed by T composition of SOI SGM. Using this observation, the privacy proof is a straightforward extension of the sensitivity analysis in Theorem B.2 and Theorem B.3. \square

Algorithm 7 Newton Method with Double noise - Minibatch Version

Inputs: training set $S_n \in \mathcal{Z}^n$, $\lambda_0 > 0$, $\theta \in (0, 1)$, privacy budget (ϵ, δ) -DP, initialization w_0 , number of iterations T , SOI modification $\in \{\text{clip}, \text{add}\}$, sampling rates $p_g, p_H \in (0, 1)$ for gradient and SOI.
Set $\sigma_1 = \text{get_noise_multiplier}$ (privacy budget = $((1 - \theta)\epsilon, (1 - \theta)\delta)$, sampling rate = p_g , steps = T)
if SOI modification = Add **then**
 $\sigma_2 = \frac{1}{(4np_H\lambda_0^2 + \lambda_0)} \cdot \text{get_noise_multiplier}$ (privacy budget = $(\theta\epsilon, \theta\delta)$, sampling rate = p_H , steps = T)
else if SOI modification = Clip **then**
 $\sigma_2 = \frac{1}{(4np_H\lambda_0^2 - \lambda_0)} \cdot \text{get_noise_multiplier}$ (privacy budget = $(\theta\epsilon, \theta\delta)$, sampling rate = p_H , steps = T)
for $t = 0, \dots, T - 1$ **do**
Take a Poisson subsample $\mathcal{I}_{t,g} \subseteq [n]$ with sampling probability p_g
Take a Poisson subsample $\mathcal{I}_{t,H} \subseteq [n]$ with sampling probability p_H
Query $g_t = \frac{1}{np_g} \sum_{i \in \mathcal{I}_{t,g}} \nabla f_{\text{LL}}(w_t, z_i)$ and $H_t = \frac{1}{np_H} \sum_{j \in \mathcal{I}_{t,H}} H(w_t, z_j)$
 $\tilde{H}_t = \Psi_{\lambda_0}(H_t, \text{SOI modification})$
 $\tilde{g}_t = g_t + \frac{1}{np_g} \mathcal{N}(0, \sigma_1^2 I_d)$
 $w_{t+1} = w_t - \tilde{H}_t^{-1} \tilde{g}_t + \mathcal{N}(0, \|\tilde{g}_t\|^2 \sigma_2^2 I_d)$
Output w_T .

Parameter	Value
θ : fraction of the privacy budget for the search direction in Algorithm 5	0.3
γ : fraction of the privacy budget for computing trace in Algorithm 5	0.1
β : the coefficient for minimum eigenvalue in Algorithm 5	$\{0.5, 1, 2\}$
number of independent runs	15

Table 2: Hyperparameters of Algorithm 5

C.2 Details of the experiments

Table 2 summarizes the hyperparameters of Algorithm 5 used for the experiments. Notice that these parameters are *data-independent*. Table 3 lists the public datasets used in our experimental evaluation.

dataset name	number of samples	dimension	Reference
ala	30956	134	[DG17]
adult	45220	104	[DG17]
(binary) coverytype	53121	55	[BD99; DG17]
synthetic	10000	100	See Section 6
(binary) FMNIST	12000	784	[XRV17]
protein	50000	74	[CJB04]

Table 3: Datasets used in the experiments

C.3 Privacy-Utility-Run Time tradeoff

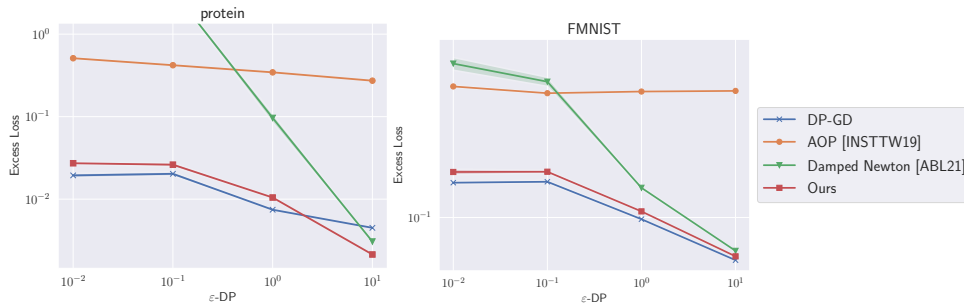


Figure 6: Privacy-Excess Loss Tradeoff for FMNIST and protein

	$\frac{T_{\text{DP-GD}}^*}{T_{\text{ours}}^*}$				$T_{\text{ours}}^* (\text{sec})$	
	$\varepsilon = 0.01$	$\varepsilon = 0.1$	$\varepsilon = 1$	$\varepsilon = 10$	$\min(T_{\text{ours}}^*) (\text{sec.})$	$\max(T_{\text{ours}}^*) (\text{sec.})$
FMNIST	$3.44 \times$	$2.79 \times$	$2.77 \times$	$8.74 \times$	11.36	25.61
protein	$6.65 \times$	$9.62 \times$	$24.16 \times$	$26.99 \times$	3.99	4.66

Table 4: Comparison between the run time of our algorithm and DP-GD in terms of the ratio $T_{\text{DP-GD}}^*/T_{\text{ours}}^*$. The last two columns show the minimum and maximum run time of our algorithm.

C.4 Second Order Information vs Optimal Stepsize

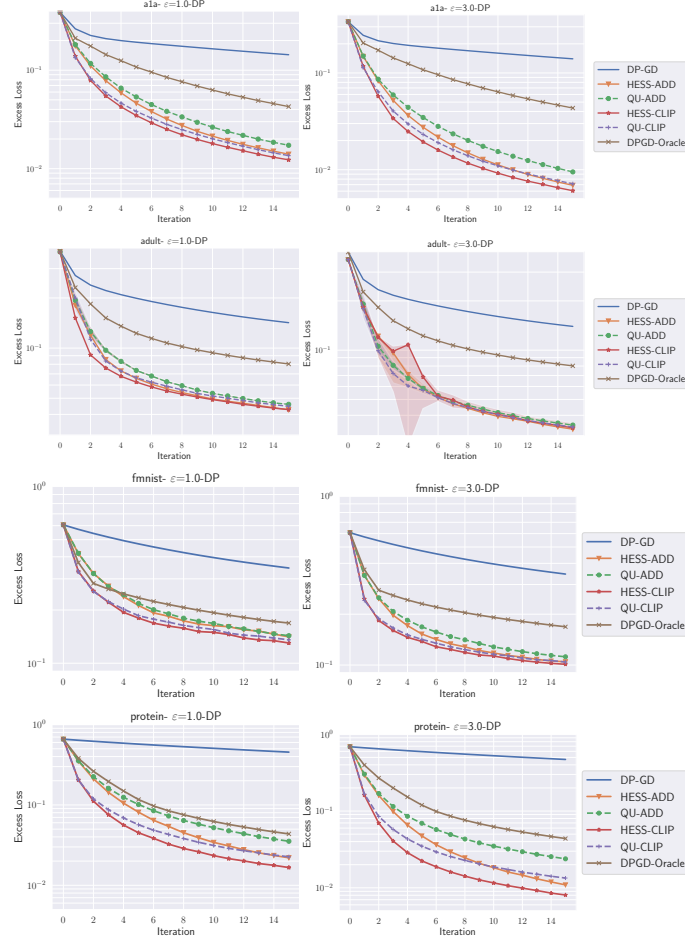
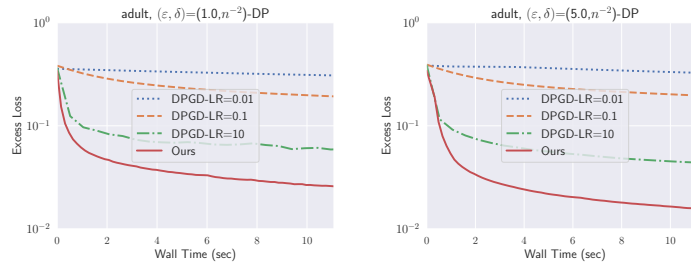


Figure 7: Comparison with DP-GD-Oracle (left $(\varepsilon, \delta) = (1, n^{-2})$, right $(\varepsilon, \delta) = (3, n^{-2})$)

C.5 Experiment with different learning rates for DPGD



C.6 Minibatch Variant

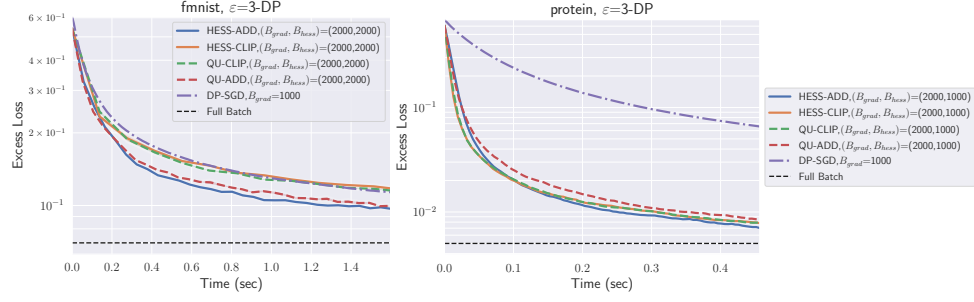


Figure 8: Minibatch variants for FMNIST and protein

C.7 Memory Usage of SOI modification

The proposed method for SOI modification may need a large space. The memory usage varies with the chosen SVD implementation; we utilized numpy's `linalg.eigh` in our experiments. For Adult dataset with dimension 100, the memory usage is 52 mebibyte (MiB), and for Fashion-MNIST with dimension 784, the memory usage is 67 MiB.