
Private (Stochastic) Non-Convex Optimization Revisited: Second-Order Stationary Points and Excess Risks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We reconsider the challenge of non-convex optimization under differential privacy
2 constraint. Building upon the previous variance-reduced algorithm SpiderBoost,
3 we propose a novel framework that employs two types of gradient oracles: one
4 that estimates the gradient at a single point and a more cost-effective option
5 that calculates the gradient difference between two points. Our framework can
6 ensure continuous accuracy of gradient estimations and subsequently enhances
7 the rates of identifying second-order stationary points. Additionally, we consider
8 a more challenging task by attempting to locate the global minima of a non-
9 convex objective via the exponential mechanism without almost any assumptions.
10 Our preliminary results suggest that the regularized exponential mechanism can
11 effectively emulate previous empirical and population risk bounds, negating the
12 need for smoothness assumptions for algorithms with polynomial running time.
13 Furthermore, with running time factors excluded, the exponential mechanism
14 demonstrates promising population risk bound performance, and we provide a
15 nearly matching lower bound.

16 1 Introduction

17 Differential privacy [18] is a standard privacy guarantee for training machine learning models. Given
18 a randomized algorithm $\mathcal{A} : P^* \rightarrow R$, where P is a data domain and R is a range of outputs, we say
19 \mathcal{A} is (ε, δ) -differentially private (DP) for some $\varepsilon \geq 0$ and $\delta \in [0, 1]$ if for any neighboring datasets
20 $\mathcal{D}, \mathcal{D}' \in P^*$ that differ in at most one element and any $\mathcal{R} \subseteq R$, the distribution of the outcome of the
21 algorithm, e.g., pair of models trained on the respective datasets, are similar:

$$\Pr_{x \sim \mathcal{A}(\mathcal{D})} [x \in \mathcal{R}] \leq e^\varepsilon \Pr_{x \sim \mathcal{A}(\mathcal{D}')} [x \in \mathcal{R}] + \delta.$$

22 Smaller ε and δ imply the distributions are closer; hence, an adversary accessing the trained model
23 cannot tell with high confidence whether an example x was in the training dataset. Given this measure
24 of privacy, we consider the problem of optimizing a non-convex loss while ensuring a desired level of
25 privacy. In particular, suppose we are given a dataset $\mathcal{D} = \{z_1, \dots, z_n\}$ drawn i.i.d. from underlying
26 distribution \mathcal{P} . Each loss function $f(\cdot; z) : \mathcal{K} \rightarrow \mathbb{R}$ is G -Lipschitz over the convex set $\mathcal{K} \subset \mathbb{R}^d$ of
27 diameter D . Let the population risk function be $F_{\mathcal{P}}(x) := \mathbb{E}_{z \sim \mathcal{P}} [f(x; z)]$ and the empirical risk
28 function be $F_{\mathcal{D}}(x) := \frac{1}{n} \sum_{z \in \mathcal{D}} f(x; z)$. We also denote $F_S(x) := \frac{1}{|S|} \sum_{z \in S} f(x; z)$ for $S \subseteq \mathcal{D}$.

29 Our focus is in minimizing non-convex (empirical and population) risk functions, which may have
30 multiple local minima. Since finding the global optimum of a non-convex function can be challenging,
31 an alternative goal in the field is to find stationary points: A first-order stationary point is a point
32 with a small gradient of the function, and a second-order stationary point is a first-order stationary

33 point where additionally the function has a positive or nearly positive semi-definite Hessian. As first
 34 order stationary points can be saddle points or even a local maximum, we focus on the problem of
 35 finding a second order stationary point, i.e., a local minimum, privately. Existing works in finding
 36 approximate SOSP privately only give guarantees for the empirical function $F_{\mathcal{D}}$. We improve upon
 37 the state-of-the-art result for empirical risk minimization and give the first guarantee for the population
 38 function $F_{\mathcal{P}}$. This requires standard assumptions on bounded Lipschitzness, smoothness, and Hessian
 39 Lipschitzness, which we make precise in Section 2 and in Assumption 3.1.

40 Compared to finding a local minimum, finding a global minimum can be extremely challenging. We
 41 also present two methods, polynomial and exponential time, that outperform existing guarantees
 42 measured in excess risks for respective computational complexities. Our primary results are succinctly
 43 summarized in Table 1.

44 **Related Work.** We propose a novel and simple framework based on SpiderBoost [51], and its
 45 private version [2] that achieves the current best rate for finding the first order stationary point privately.
 46 We discuss the primary difference between our framework and theirs, that is their algorithms only
 47 promise small gradient estimation errors on average, but our framework can ensure small estimation
 48 errors consistently throughout all the iterations, and the motivation behind this briefly.

49 In SGD and its variants, the typical approach involves obtaining an estimation Δ_t of the gradient
 50 $\nabla f(x_t)$. In the stochastic variance-reduced algorithm SpiderBoost [51, 2], it queries the gradient
 51 $\mathcal{O}_1(x_t) \approx \nabla f(x_t)$ directly every q steps with some oracle \mathcal{O}_1 , and for the other $q - 1$ steps
 52 within each period, it queries the gradient difference between two steps, that is $\mathcal{O}_2(x_t, x_{t-1}) \approx$
 53 $\nabla f(x_t) - \nabla f(x_{t-1})$, and maintain $\Delta_t = \Delta_{t-1} + \mathcal{O}_2(x_t, x_{t-1})$. The contrast between these two
 54 types of oracles can be perceived as \mathcal{O}_1 being more accurate but also more costly, in terms of
 55 computation or privacy budget, although our framework does not strictly necessitate this assumption.

56 As SpiderBoost queries \mathcal{O}_1 every q steps, the error on the estimation may accumulate and $\|\Delta_t -$
 57 $\nabla f(x_t)\|$ can become large. Despite this, as demonstrated in [2], these estimations can, on average,
 58 suffice to find a private FOSP. However, such large deviations pose a challenge when scrutinizing
 59 behavior near a saddle point. For instance, when the current point is a saddle point, but the current
 60 estimation is unsatisfactory, it becomes uncertain whether the algorithm can escape the saddle point. It
 61 could be argued that average good estimations could achieve a SOSP, but to the best of our knowledge,
 62 there is no existing result addressing this concern.

63 A plausible solution to this challenge is to maintain high-quality gradient estimations throughout
 64 all iterations, a feat accomplished by our framework. We believe this feature holds promise for
 65 improving the outcomes of various other optimization problems, thus enhancing the overall appeal
 66 and significance of our work.

67 1.1 Main Results

68 **SOSP.** One of our main contributions is a refined optimization framework (Algorithm 1), predi-
 69 cated on the variance-reduced SpiderBoost [51], which guarantees consistently accurate gradient
 70 estimations. By integrating this framework with private gradient oracles, we achieve improved error
 71 rates for privately identifying SOSP of both empirical and population risks.

72 Advances in private non-convex optimization have focused on finding a first-order stationary point
 73 (FOSP), whose performance is measured in (i) the norm of the empirical gradient at the solution x ,
 74 i.e., $\|\nabla F_{\mathcal{D}}(x)\|$, and (ii) the norm of the population gradient, i.e., $\|\nabla F_{\mathcal{P}}(x)\|$. We survey the recent
 75 progress in the appendix in detail.

76 **Definition 1.1** (First-order stationary point). *We say $x \in \mathbb{R}^d$ is a First-Order Stationary Point (FOSP)*
 77 *of $g : \mathbb{R}^d \rightarrow \mathbb{R}$ iff $\nabla g(x) = 0$. x is an α -FOSP of g , if $\|\nabla g(x)\|_2 \leq \alpha$.*

78 Since FOSP can be a saddle point or a local maxima, finding a second-order stationary point is
 79 desired. Exact second-order stationary points can be extremely challenging to find [24]. Instead,
 80 progress is commonly measured in terms of how well the solution approximates an SOSP.

81 **Definition 1.2** (Second-order stationary point, [1]). *We say a point $x \in \mathbb{R}^d$ is a Second-Order*
 82 *Stationary Point (SOSP) of a twice differentiable function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ iff $\|\nabla g(x)\|_2 = 0$ and*
 83 *$\nabla^2 g(x) \succeq 0$. We say $x \in \mathbb{R}^d$ is an α -SOSP for ρ -Hessian Lipschitz function g , if $\|\nabla g(x)\|_2 \leq$
 84 $\alpha \wedge \nabla^2 g(x) \succeq -\sqrt{\rho}\alpha I$.*

	α -SOSP		Excess population risk	
	empirical	population	poly-time	exp-time
SOTA	$\min\left(\frac{d^{\frac{1}{4}}}{n^{\frac{1}{2}}\varepsilon^{\frac{1}{2}}}, \frac{d^{\frac{4}{7}}}{n^{\frac{4}{7}}\varepsilon^{\frac{4}{7}}}\right)$	N/A	$\frac{d}{\varepsilon^2 \log n}$ ♣	N/A
Ours	$\frac{d^{\frac{1}{3}}}{n^{\frac{2}{3}}\varepsilon^{\frac{2}{3}}}$	$\frac{1}{n^{\frac{1}{3}}} + \left(\frac{\sqrt{d}}{n\varepsilon}\right)^{\frac{3}{7}}$	$\frac{d \log \log n}{\varepsilon \log(n)}$	$\frac{d}{n\varepsilon} + \sqrt{\frac{d}{n}}$
LB	$\frac{\sqrt{d}}{n\varepsilon}$	$\frac{1}{\sqrt{n}} + \frac{\sqrt{d}}{n\varepsilon}$	$\frac{d}{n\varepsilon} + \sqrt{\frac{d}{n}}$	$\frac{d}{n\varepsilon} + \sqrt{\frac{d}{n}}$

Table 1: SOTA refers to the best previously known bounds on α for α -SOSP by [45, 47] and on the excess population risk by [45]. We introduce algorithm 1 that finds an α -SOSP (columns 2–3) with an improved rate. We show exponential mechanism can minimize the excess risk in polynomial time and exponential time, respectively (columns 4 and 5). ♣ requires extra assumption on bounded smoothness. The lower bounds for SOSP are from [2], and the lower bound on excess population risk is from Theorem 5.11. We omit logarithmic factors in n and d .

85 On the empirical risk $F_{\mathcal{D}}$, the SOTA on privately finding α -SOSP is by [45, 47], which achieves $\alpha =$
86 $\tilde{O}(\min\{(\sqrt{d}/n)^{1/2}, (d/n)^{4/7}\})$. In Theorem 4.2, we show that applying the proposed Algorithm 1
87 achieves a rate bounded by $\alpha = \tilde{O}((\sqrt{d}/n)^{2/3})$, which improves over the SOTA in all regime.¹ There
88 remains a factor $(\sqrt{d}/n)^{-1/6}$ gap to a known lower bound of $\alpha = \Omega(\sqrt{d}/n)$ that holds even if finding
89 only an α -FOSP [2]. On the population risk $F_{\mathcal{P}}$, applying Algorithm 1 with appropriate private
90 gradient oracles is the first private algorithm to guarantee finding an α -SOSP with $\alpha = \tilde{O}(n^{-1/3} +$
91 $(\sqrt{d}/n)^{3/7})$ in Theorem 4.6. There is a gap to a known lower bound of $\alpha = \Omega(1/\sqrt{n} + \sqrt{d}/n\varepsilon)$ that
92 holds even if finding only an α -FOSP [2].

93 **Minimizing Excess Risk.** In addition to the optimization framework, we present sampling-based
94 algorithms designed to identify a private solution $x^{priv} \in \mathbb{R}^d$ that minimizes both the excess empirical
95 risk: $\mathbb{E}[F_{\mathcal{D}}(x^{priv})] - \min_{x \in \mathcal{K}} F_{\mathcal{D}}(x)$, and the excess population risk: $\mathbb{E}[F_{\mathcal{P}}(x^{priv})] - \min_{x \in \mathcal{K}} F_{\mathcal{P}}(x)$.
96 Here, the expectation is over the randomness of the solution x^{priv} and the drawing of the training
97 date over \mathcal{P} . Our method is different from [45], which Gradient Langevin Dynamics and achieves
98 in polynomial time a bound of $O(d\sqrt{\log(1/\delta)}/(\varepsilon^2 \log n))$ for both excess empirical and population
99 risks with a need for the smoothness assumption. In Table 1 we omit excess empirical risk, as the
100 bounds align with those of the population risk. We introduce a sampling-based algorithm from the
101 exponential mechanism, which runs in polynomial time and achieves excess empirical and population
102 risks bounded by $O(d\sqrt{\log(1/\delta)}/(\varepsilon \log(nd)))$ with improved dependence on ε (Theorem 5.6).
103 Crucially, it achieves these results without the need for the smoothness assumption required by [45].

104 In the case of permitting an exponential running time, [22] demonstrated $\tilde{O}(d/(\varepsilon n))$ upper bound for
105 non-convex excess empirical risks alongside a nearly matching lower bound. However, establishing a
106 tight bound for the excess population risk remained an unresolved problem. We address this open
107 question by providing nearly matching upper and lower bounds of $\tilde{\Theta}(d/(\varepsilon n) + \sqrt{d}/n)$ for the excess
108 population risk (Theorem 5.8).

109 1.2 Our Techniques

110 **Stationary Points.** In our framework, we deviate from the traditional approach of querying \mathcal{O}_1
111 once every q steps. Instead, we introduce a novel but simple method of monitoring the total drift we
112 make, that is $\text{drift}_t = \sum_{i=\tau_t}^t \|x_i - x_{i-1}\|_2^2$, where τ_t represents the last timestamp when we employed
113 \mathcal{O}_1 . As we are considering smooth functions, the maximum error to estimate $\nabla f(x_t) - \nabla f(x_{t-1})$
114 is proportional to $\|x_t - x_{t-1}\|_2$. If the value drift_t is small, we know the current estimation should
115 still be good enough, eliminating the need for an expensive fresh estimation from \mathcal{O}_1 . Conversely,
116 when drift_t is large, the gradient estimation error may be substantial, necessitating a query to \mathcal{O}_1 and
117 thus obtaining $\Delta_t = \mathcal{O}_1(x_t)$. To effectively manage the total cost, it is crucial to set an appropriate
118 threshold to decide when the drift is significant. A smaller threshold would ensure more accurate
119 estimations but might incur higher costs due to more frequent queries to \mathcal{O}_1 .

¹We want $\alpha = o(1)$ and hence can assume $d \leq n^2$.

120 Our aim is to bound the total occurrences of the event that drift_t is large, which leads to querying \mathcal{O}_1 .
 121 A crucial observation is that, if drift_t increases rapidly, then the gradient norms are large and hence
 122 function values decrease quickly, which we know does not happen frequently under the standard
 123 assumption that the function is bounded.

124 In our framework, we assume $\mathcal{O}_1(x)$ is an unbiased estimation of $\nabla f(x)$, and $\mathcal{O}_1(x) - \nabla f(x)$ is
 125 Norm-SubGaussian (Definition 2.2), and similarly $\mathcal{O}_2(x, y)$ is an unbiased estimation of $\nabla f(x) -$
 126 $\nabla f(y)$ whose error is also Norm-SubGaussian. In the empirical case, we can simply add Gaussian
 127 noises with appropriately chosen variances to the gradients of the empirical function $\nabla F_{\mathcal{D}}$ for
 128 simplicity, and one can choose a smaller batch size to reduce the computational complexity. In
 129 the population case, we draw samples from the dataset without replacement to avoid dependence
 130 issues, and add the Gaussian noises to the sampled gradients. Hence we only need the gradient oracle
 131 complexity to be linear in the size of dataset for the population case.

132 **Minimizing Excess Risk.** Our polynomial time approach harnesses the power of the Log-Sobolev
 133 Inequality (LSI) and the classic Stroock perturbation lemma. The previous work of [38] shows that if
 134 the density $\exp(-\beta F_{\mathcal{D}}(x) - r(x))$ satisfies the LSI for some regularizer r , then sampling a model
 135 x from this density is DP with an appropriate (ε, δ) . If r is a μ strongly convex function, then the
 136 density proportional to $\exp(-r)$ satisfies LSI with constant $1/\mu$, and $\exp(-\beta F_{\mathcal{D}}(x) - r(x))$ satisfies
 137 LSI with constant $\exp(\max_{x,y} |F_{\mathcal{D}}(x) - F_{\mathcal{D}}(y)|)/\mu$ by the Stroock perturbation lemma. Our bound
 138 on the empirical risk follows from choosing the appropriate inverse temperature β and regularizer r
 139 to satisfy (ε, δ) -DP. The final bound on the population risk also follows from LSI, which bounds the
 140 stability of the sample drawn from the respective distribution.

141 When running time is not a priority, we employ an exponential mechanism over a discretization of
 142 \mathcal{K} to establish the upper bound. The empirical risk bound derives from [9], and we leverage the
 143 concentration of sums of bounded random variables to bound the maximum difference over the
 144 discretizations between the empirical and population risk. We show this is nearly tight by reductions
 145 from selection to non-convex Lipschitz optimization of [22].

146 1.3 Organization

147 In Section 2, we present necessary definitions and backgrounds for our work. In Section 3, we
 148 construct the optimization framework, with guarantees on finding the SOSGP with two different
 149 kinds of SubGaussian gradient oracles. It's crucial to note that this framework focuses solely on
 150 optimization and does not pertain to privacy. Section 4 explores the pursuits of finding the SOSGP
 151 privately by constructing private SubGaussian gradient oracles and seamlessly integrating them into
 152 the existing framework. We bound the private excess bounds in Section 5. For other preliminaries, all
 153 omitted proofs and some further discussions on related work can be found in the Appendix.

154 2 Preliminaries

155 Throughout the paper, if not stated explicitly, the norm $\|\cdot\|$ means the ℓ_2 norm.

156 **Definition 2.1** (Lipschitz, Smoothness and Hessian Lipschitz). *Given a function $f : \mathcal{K} \rightarrow \mathbb{R}$, we
 157 say f is G -Lipschitz, if for all $x_1, x_2 \in \mathcal{K}$, $|f(x_1) - f(x_2)| \leq G\|x_1 - x_2\|$, we say a function f is
 158 M -smooth, if for all $x_1, x_2 \in \mathcal{K}$, $\|\nabla f(x_1) - \nabla f(x_2)\| \leq M\|x_1 - x_2\|$. and we say the function f
 159 is ρ -Hessian Lipschitz, if for all $x_1, x_2 \in \mathcal{K}$, we have $\|\nabla^2 f(x_1) - \nabla^2 f(x_2)\| \leq \rho\|x_1 - x_2\|$.*

160 **Definition 2.2** (SubGaussian, and Norm-SubGaussian). *A random vector $x \in \mathbb{R}^d$ is SubGaussian
 161 (SG(ζ)) if there exists a positive constant ζ such that $\mathbb{E} e^{\langle v, x - \mathbb{E}x \rangle} \leq e^{\|v\|^2 \zeta^2 / 2}$, $\forall v \in \mathbb{R}^d$. $x \in \mathbb{R}^d$
 162 is norm-SubGaussian (nSG(ζ)) if there exists ζ such that $\Pr[\|x - \mathbb{E}x\| \geq t] \leq 2e^{-\frac{t^2}{2\zeta^2}}$, $\forall t \in \mathbb{R}$.*

163 **Fact 2.3.** *For a Gaussian $\theta \sim \mathcal{N}(0, \sigma^2 I_d)$, θ is SG(σ) and nSG($\sigma\sqrt{d}$).*

164 **Lemma 2.4** (Hoeffding type inequality for norm-subGaussian, [29]). *Let $x_1, \dots, x_k \in \mathbb{R}^d$ be
 165 random vectors, and for each $i \in [k]$, $x_i | \mathcal{F}_{i-1}$ is zero-mean nSG(ζ_i) where \mathcal{F}_i is the corresponding
 166 filtration. Then there exists an absolute constant c such that for any $\delta > 0$, with probability at least
 167 $1 - \omega$, $\|\sum_{i=1}^k x_i\| \leq c \cdot \sqrt{\sum_{i=1}^k \zeta_i^2 \log(2d/\omega)}$, which means $\sum_{i=1}^k x_i$ is nSG($\sqrt{c \log(d) \sum_{i=1}^k \zeta_i^2}$).*

168 3 Convergence to Stationary Points: Framework

169 We present the optimization framework for finding SOSP in this section. It's important to emphasize
 170 that this framework is dedicated exclusively to optimization concerns, with privacy considerations
 171 being outside of its purview. The results about SOSP throughout the paper follows the assumptions
 172 of [45].

173 **Assumption 3.1.** Any function drawn from \mathcal{P} is G -Lipschitz, ρ -Hessian Lipschitz, and M -smooth,
 174 almost surely, and the risk is upper bounded by B .

175 As discussed before, we define two different kinds of gradient oracles, one for estimating the gradient
 176 at one point and the other for estimating the gradient difference at two points.

177 **Definition 3.2** (SubGaussian gradient oracles). For a G -Lipschitz and M -smooth function F :

178 (1) We say \mathcal{O}_1 is a first kind of ζ_1 norm-subGaussian Gradient oracle if given $x \in \mathbb{R}^d$, $\mathcal{O}(x)$ satisfies
 179 $\mathbb{E} \mathcal{O}_1(x) = \nabla F(x)$ and $\mathcal{O}_1(x) - \nabla F(x)$ is $\text{nSG}(\zeta_1)$.

180 (2) We say \mathcal{O}_2 is a second kind of ζ_2 norm-subGaussian stochastic Gradient oracle if given $x, y \in$
 181 \mathbb{R}^d , $\mathcal{O}_2(x, y)$ satisfies that $\mathbb{E} \mathcal{O}_2(x, y) = \nabla F(x) - \nabla F(y)$ and $\mathcal{O}_2(x, y) - (\nabla F(x) - \nabla F(y))$ is
 182 $\text{nSG}(\zeta_2 \|x - y\|)$.

183 Note that we should assume $M \geq \sqrt{\rho\alpha}$ to make finding a second-order stationary point strictly
 184 more challenging than finding a first-order stationary point. We use $\text{smin}(\cdot)$ to denote the smallest
 185 eigenvalue of a matrix.

Algorithm 1 Stochastic Spider

1: **Input:** Objective function F , Gradient Oracle $\mathcal{O}_1, \mathcal{O}_2$ with SubGaussian parameters ζ_1 and ζ_2 ,
 parameters of objective function B, M, G, ρ , parameter κ , failure probability ω
 2: Set $\gamma = \sqrt{4C(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(BMd/\rho\omega)}$, $\Gamma = \frac{M \log(\frac{dMB}{\rho\gamma\omega})}{\sqrt{\rho\gamma}}$
 3: Set $\eta = 1/M, t = 0, T = BM \log^4(\frac{dMB}{\rho\gamma\omega})/\gamma^2$
 4: Set $\text{drift}_0 = \kappa, \text{frozen} = 1, \nabla_{-1} = 0$
 5: **while** $t \leq T$ **do**
 6: **if** $\|\nabla_{t-1}\| \leq \gamma \log^3(BMd/\rho\omega) \wedge \text{frozen}_{t-1} \leq 0$ **then**
 7: $\text{frozen}_t = \Gamma, \text{drift}_t = 0$
 8: $\nabla_t = \mathcal{O}_1(x_t) + g_t$, where $g_t \sim \mathcal{N}(0, \frac{\zeta_1^2}{d} I_d)$
 9: **else if** $\text{drift}_{t-1} \geq \kappa$ **then**
 10: $\nabla_t = \mathcal{O}_1(x_t), \text{drift}_t = 0, \text{frozen}_t = \text{frozen}_{t-1} - 1$
 11: **else**
 12: $\Delta_t = \mathcal{O}_2(x_t, x_{t-1}), \nabla_t = \nabla_{t-1} + \Delta_t, \text{frozen}_t = \text{frozen}_{t-1} - 1$
 13: **end if**
 14: $x_{t+1} = x_t - \eta \nabla_t, \text{drift}_t = \text{drift}_{t-1} + \eta^2 \|\nabla_t\|_2^2, t = t + 1$
 15: **end while**
 16: **Return:** $\{x_1, \dots, x_T\}$

186 We demonstrate a framework based on the SpiderBoost in Algorithm 1. Our analysis of Algorithm 1
 187 hinges on three key properties we establish in this section: (i) ∇_t remains consistently close to the
 188 true gradient $\nabla F(x_t)$ with high probability; (ii) the algorithm is capable of escaping the saddle point
 189 with high probability, and (iii) a large drift implies significant decrease in the function value, which
 190 enables us to limit the number of queries to the more accurate but costlier first kind of gradient oracle
 191 \mathcal{O}_1 .

192 **Lemma 3.3.** For any $0 \leq t \leq T$ and letting $\tau_t \leq t$ be the largest integer such that drift_{τ_t} is set to
 193 be 0, with probability at least $1 - \omega/T$, for some universal constant $C > 0$, we have

$$\|\nabla_t - \nabla F(x_t)\|^2 \leq (\zeta_2^2 \cdot \sum_{i=\tau_t+1}^t \|x_i - x_{i-1}\|^2 + 4\zeta_1^2) \cdot C \cdot \log(Td/\omega). \quad (1)$$

194 Hence with probability at least $1 - \omega$, we know for each $t \leq T$, $\|\nabla_t - \nabla F(x_t)\|^2 \leq \gamma^2/16$, where
 195 $\gamma^2 := 16C(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(Td/\omega)$ and κ is a parameter we can choose in the algorithm.

196 As shown in Lemma 3.3, the error on the gradient estimation for each step is bounded with high
 197 probability. Then we can show the algorithm can escape the saddle point efficiently based on previous
 198 results.

199 **Lemma 3.4** (Essentially from [45]). *Under Assumption 3.1, run SGD iterations $x_{t+1} = x_t -$
 200 $\eta \nabla_t$, with step size $\eta = 1/M$. Suppose x_0 is a saddle point satisfying $\|\nabla F(x_0)\| \leq \alpha$ and
 201 $\text{smin}(\nabla^2 F(x_0)) \leq -\sqrt{\rho\alpha}$, $\alpha = \gamma \log^3(dBM/\rho\omega)$. If $\nabla_0 = \nabla F(x_0) + \zeta_1 + \zeta_2$ where $\|\zeta_1\| \leq \gamma$,
 202 $\zeta_2 \sim \mathcal{N}(0, \frac{\gamma^2}{d \log(d/\omega)} I_d)$, and $\|\nabla_t - \nabla F(x_t)\| \leq \gamma$ for all $t \in [\Gamma]$, with probability at least
 203 $1 - \omega \cdot \log(1/\omega)$, one has $F(x_\Gamma) - F(x_0) \leq -\Omega(\frac{\gamma^{3/2}}{\sqrt{\rho} \log^3(\frac{dMB}{\rho\gamma\omega})})$, where $\Gamma = \frac{M \log(\frac{dMB}{\rho\gamma\omega})}{\sqrt{\rho\gamma}}$.*

204 We discuss this lemma in the Appendix in more details. The next lemma is standard, showing how
 205 large the function values can decrease in each step.

206 **Lemma 3.5.** *By setting $\eta = 1/M$, we have $F(x_{t+1}) \leq F(x_t) + \eta \|\nabla_t\| \cdot \|\nabla F(x_t) - \nabla_t\| - \frac{\eta}{2} \|\nabla_t\|^2$.
 207 Moreover, with probability at least $1 - \omega$, for each $t \leq T$ such that $\|\nabla F(x_t)\| \geq \gamma$, we have*

$$F(x_{t+1}) - F(x_t) \leq -\eta \|\nabla_t\|^2 / 6 \leq -\eta \gamma^2 / 6.$$

208 With the algorithm designed to control the drift term, the guarantee for Stochastic Spider to find the
 209 second order stationary point is stated below:

210 **Lemma 3.6.** *Suppose \mathcal{O}_1 and \mathcal{O}_2 are ζ_1 and ζ_2 norm-subGaussian respectively. If one sets $\gamma =$
 211 $O(1) \sqrt{(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(Td/\omega)}$, with probability at least $1 - \omega$, at least one point in the output set
 212 $\{x_1, \dots, x_T\}$ of Algorithm 1 is α -SOSP, where*

$$\alpha = \gamma \log^3(BMd/\rho\omega\gamma) = \sqrt{(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(\frac{d/\omega}{\zeta_2^2 \kappa + \zeta_1^2}) \cdot \log^3(\frac{BMd}{\rho\omega(\zeta_2^2 \kappa + \zeta_1^2)})}.$$

213 As mentioned before, we can bound the number of occurrences where the drift gets large and hence
 214 bound the total time we query the oracle of the first kind.

215 **Lemma 3.7.** *Under the event that $\|\nabla_t - \nabla F(x_t)\| \leq \gamma/4$ for all $t \in [T]$ and our parameter settings,
 216 letting $K = \{t \in [T] : \text{drift}_t \geq \kappa\}$ be the set of iterations where the drift is large, we know
 217 $|K| \leq O(\frac{B\eta}{\kappa} + T\gamma^2\eta^2/\kappa) = O(B\eta \log^4(\frac{dMB}{\rho\gamma\omega})/\kappa)$.*

218 4 Private SOSP

219 We adopt the framework before and get our main results on finding SOSP privately by constructing
 220 private gradient oracles in this section. Finding SOSP for empirical risk function $F_{\mathcal{D}}$ and for
 221 population risk function $F_{\mathcal{P}}$ are discussed in Subsection 4.1 and Subsection 4.2 respectively.

222 4.1 Convergence to the SOSP of the Empirical Risk

223 We use Stochastic Spider to improve the convergence to α -SOSP of the empirical risk, and aim at
 224 getting $\alpha = \tilde{O}(d^{1/3}/n^{2/3})$. We use the full-batch size for simplicity, and use the gradient oracles

$$\mathcal{O}_1(x) := \nabla F_{\mathcal{D}}(x) + g_1, \text{ and } \mathcal{O}_2(x, y) := \nabla F_{\mathcal{D}}(x) - \nabla F_{\mathcal{D}}(y) + g_2, \quad (2)$$

225 where $g_1 \sim \mathcal{N}(0, \sigma_1^2 I_d)$ and $g_2 \sim \mathcal{N}(0, \sigma_2^2 \|x - y\|_2^2 I_d)$ are added to ensure privacy by Gaussian
 226 mechanism (in Appendix).

227 Before stating the formal results, note that by Lemma 3.6, the framework can only guarantee the
 228 existence of an α -SOSP in the outputted set. In order to find the SOSP privately from the set, we
 229 adopt the well-known AboveThreshold algorithm, whose pseudo-code can be found in Algorithm 2
 230 in the Appendix. Algorithm 2 is a slight modification of the well-known AboveThreshold algorithm
 231 in [19], and we get the following guarantee immediately.

232 **Lemma 4.1.** *Algorithm 2 is $(\varepsilon, 0)$ -DP. Given the point set $\{x_1, \dots, x_T\}$ and S of size n as the input,
 233 (i) if it outputs any point x_i , then with probability at least $1 - \omega$, we know*

$$\|\nabla F_S(x_i)\| \leq \alpha + \frac{32 \log(2T/\omega)G}{n\varepsilon}, \text{ and } \text{smin}(\nabla^2 F_S(x_i)) \geq -\sqrt{\rho\alpha} - \frac{32 \log(2T/\omega)M}{n\varepsilon}$$

234

235 (ii) if there exists a α -SOSP point $x \in \{x_i\}_{i \in [T]}$, then with probability at least $1 - \omega$, Algorithm 2
 236 will output one point.

237 Choosing the appropriate noise scales for the Gaussian added in Equation (2) and running Algorithm 1
 238 can get a private set of points which contains at least one good SOSP. Then we can run Algorithm 2
 239 to find the good SOSP in the set privately. The formal guarantee is stated below:

240 **Theorem 4.2** (Empirical). For $\varepsilon \leq 10, \delta \in (0, 1/2)$, use Equation (2) as gradient or-
 241 acles with $\kappa = \frac{G^{4/3} B^{1/3}}{M^{5/3}} \left(\frac{\sqrt{d \log(1/\delta)}}{n\varepsilon} \right)^{2/3}$, $\sigma_1 = \frac{G \sqrt{B \eta \log^2(n/\delta) / \kappa \log^2(ndMB/\omega)}}{n\varepsilon}$, $\sigma_2 =$
 242 $\frac{M \sqrt{\log^2(n/\delta) B M / \alpha_1^2 \log^5(ndMB/\omega)}}{n\varepsilon}$. Running Algorithm 1, outputting the set $\{x_i\}_{i \in [T]}$ if the total
 243 time to query \mathcal{O}_1 is bounded by $O(B \eta \log^4(\frac{dMB}{\rho \gamma \omega}) / \kappa)$, otherwise outputting a set of T arbitrary
 244 points is $(\varepsilon/2, \delta)$ -DP. With probability at least $1 - \omega$, at least one point in the output set is α_1 -SOSP
 245 of $F_{\mathcal{D}}$ with

$$\alpha_1 = O \left(\left(\frac{\sqrt{dBGM \log^2(1/\delta)}}{n\varepsilon} \right)^{2/3} \cdot \log^6 \left(\frac{nBMd}{\rho\omega} \right) \right).$$

246 Moreover, if we run Algorithm 2 with inputs $\{x_i\}_{i \in [T]}, \mathcal{D}, B, M, G, \rho, \alpha_1$, with probability at least
 247 $1 - \omega$, we can get an α_2 -SOSP of $F_{\mathcal{D}}$ with $\alpha_2 = O \left(\alpha_1 + \frac{G \log(n/G\omega)}{n\varepsilon} + \frac{M \log(ndBGM/\rho\omega)}{n\varepsilon\sqrt{\rho}} \sqrt{\alpha_1} \right)$.
 248 The whole procedure is (ε, δ) -DP.

249 **Remark 4.3.** It's worth noting that the cost of gradient computation can be reduced by utilizing
 250 smaller batch sizes. Additionally, the application of Rényi Differential Privacy techniques may
 251 enhance results by some logarithmic terms. However, our work does not focus on optimizing these
 252 specific aspects.

253 4.2 Convergence to the SOSP of the Population Risk

254 This subsection aims at getting an α -SOSP for $F_{\mathcal{P}}$ (the population function). Differing from the
 255 stochastic oracles used for empirical function $F_{\mathcal{D}}$, we do not use full batch in the oracle. As an
 256 alternative, we draw fresh samples from \mathcal{D} without replacement with a smaller batch size:

$$\mathcal{O}_1(x) := \frac{1}{b_1} \sum_{z \in S_1} \nabla f(x; z) + g_1, \text{ and } \mathcal{O}_2(x, y) := \frac{1}{b_2} \sum_{z \in S_2} (\nabla f(x; z) - \nabla f(y; z)) + g_2, \quad (3)$$

257 where S_1 and S_2 are sets of size of b_1 and b_2 respectively drawn from \mathcal{D} without replacement,
 258 $g_1 \sim \mathcal{N}(0, \sigma_1^2 I_d)$ and $g_2 \sim \mathcal{N}(0, \sigma_2^2 \|x - y\|_2^2 \cdot I_d)$ are added for privacy guarantee. These gradient
 259 oracles satisfy the following.

260 **Claim 4.4.** The gradient oracles \mathcal{O}_1 and \mathcal{O}_2 constructed in Equation (3) are a first kind of
 261 $O(\frac{L\sqrt{\log d}}{\sqrt{b_1}} + \sqrt{d}\sigma_1)$ norm-subGaussian gradient oracle and second kind of $O(\frac{M\sqrt{\log d}}{\sqrt{b_2}} + \sqrt{d}\sigma_2)$
 262 norm-subGaussian gradient oracle respectively.

263 Recall that in the empirical case, we use Algorithm 2 to choose the SOSP for $F_{\mathcal{D}}$. But in the
 264 population case, we need to find SOSP for $F_{\mathcal{P}}$, and what we have are samples from \mathcal{P} . We need
 265 the following technical results to help us find the SOSP from the set, which follows from Hoeffding
 266 inequality for norm-subGaussians (Lemma 2.4) and Matrix Bernstein inequality (in the Appendix).

267 **Lemma 4.5.** Fix a point $x \in \mathbb{R}^d$. Given a set S of m samples drawn i.i.d. from the distribution \mathcal{P} ,
 268 then we know with probability at least $1 - \omega$, we have

$$\|\nabla F_S(x) - \nabla F_{\mathcal{P}}(x)\|_2 \leq O\left(\frac{G \log(d/\omega)}{\sqrt{m}}\right) \wedge \|\nabla^2 F_S(x) - \nabla^2 F_{\mathcal{P}}(x)\|_{op} \leq O\left(\frac{M \log(d/\omega)}{\sqrt{m}}\right).$$

269 By choosing the appropriate noise scales σ_1 and σ_2 to ensure the privacy guarantee, we can bound
 270 the population bound similar to the empirical bound with these tools.

271 **Theorem 4.6** (Population). *Divide the dataset \mathcal{D} into two disjoint datasets \mathcal{D}_1 and \mathcal{D}_2 of size*
272 *$\lfloor n/2 \rfloor$ and $\lfloor n/2 \rfloor$ respectively. Set $b_1 = \frac{n\kappa}{B\eta}$, $b_2 = \frac{n\alpha_1^2}{BM}$, $\sigma_1 = \frac{3G\sqrt{\log(1/\delta)}}{b_1\varepsilon}$, $\sigma_2 = \frac{3M\sqrt{\log(1/\delta)}}{b_2\varepsilon}$*
273 *and $\kappa = \max(\frac{G^{4/3}B^{1/3}\log^{1/3}d}{M^{5/3}}n^{-1/3}, (\frac{GB^{2/3}}{M^{5/3}})^{6/7}(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon})^{4/7})$ in Equation (3) and use them as*
274 *gradient oracles. Running Algorithm 1 with \mathcal{D}_1 , and outputting the set $\{x_i\}_{i \in [T]}$ if the total time to*
275 *query \mathcal{O}_1 is bounded by $O(B\eta \log^4(\frac{dMB}{\rho\gamma\omega})/\kappa)$, otherwise outputting a set of T arbitrary points, is*
276 *$(\varepsilon/2, \delta)$ -DP. is $(\varepsilon/2, \delta)$ -DP, and with probability at least $1 - \omega$, at least one point in the output is*
277 *α_1 -SOSP of $F_{\mathcal{P}}$ with*

$$\alpha_1 = O\left(\left((BGM \cdot \log d)^{1/3} \frac{1}{n^{1/3}} + (G^{1/7} B^{3/7} M^{3/7}) \left(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon}\right)^{3/7}\right) \log^3(nBMd/\rho\omega)\right).$$

278 *Moreover, if we run Algorithm 2 with inputs $\{x_i\}_{i \in [T]}$, \mathcal{D}_2 , B , M , G , ρ , α_1 , with prob-*
279 *ability at least $1 - \omega$, Algorithm 2 can output an α_2 -SOSP of $F_{\mathcal{P}}$ with $\alpha_2 =$*
280 *$O\left(\alpha_1 + \frac{M \log(ndBGM/\rho\omega)}{\sqrt{\rho} \min(n\varepsilon, n^{1/2})} \sqrt{\alpha_1} + G\left(\frac{\log(n/G\omega)}{n\varepsilon} + \frac{\log(d/\omega)}{\sqrt{n}}\right)\right)$. The whole procedure is (ε, δ) -DP.*

281 5 Bounding the Private Excess Risk

282 In this section, we consider the private risk bounds.

283 5.1 Polynomial Time Approach

284 If we want the algorithm to be efficient and implementable in polynomial time, to our knowledge
285 the only known bound is $O(\frac{d\log(1/\delta)}{\varepsilon^2 \log n})$ in [45] for smooth functions. [45] used Gradient Langevin
286 Dynamics, a popular variant of SGD to solve this problem, and prove the privacy by advanced
287 composition. We generalize the exponential mechanism to the non-convex case and implement it
288 without a smoothness assumption.

289 First recall the Log-Sobolev inequality: We say a probability distribution π satisfies LSI with constant
290 C_{LSI} if for all $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $\mathbb{E}_\pi[f^2 \log f^2] - \mathbb{E}_\pi[f^2] \log \mathbb{E}_\pi[f^2] \leq 2C_{\text{LSI}} \mathbb{E}_\pi \|\nabla f\|_2^2$. A well-known
291 result ([39]) says if f is μ -strongly convex, then the distribution proportional to $\exp(-f)$ satisfies
292 LSI with constant $1/\mu$. Recall the results from previous results [38] about LSI and DP:

293 **Theorem 5.1** ([38]). *Sampling from $\exp(-\beta F(x; \mathcal{D}) - r(x))$ for some public regularizer r is (ε, δ) -*
294 *DP, where $\varepsilon \leq 2 \frac{G\beta}{n} \sqrt{C_{\text{LSI}}} \sqrt{1 + 2\log(1/\delta)}$, and C_{LSI} is the worst LSI constant.*

295 We can apply the classic perturbation lemma to get the new LSI constant in the non-convex case.
296 Suppose we add a regularizer $\frac{\mu}{2}\|x\|^2$, and try to sample from $\exp(-\beta(F(x; \mathcal{D}) + \frac{\mu}{2}\|x\|^2))$.

297 **Lemma 5.2** (Stroock perturbation). *Suppose π satisfies LSI with constant $C_{\text{LSI}}(\pi)$. If $0 < c \leq$*
298 *$\frac{d\pi'}{d\pi} \leq C$, then $C_{\text{LSI}}(\pi') \leq \frac{C}{c} C_{\text{LSI}}(\pi)$.*

299 Lemma 5.3 is a more general version of Theorem 3.4 in [22] and can be used to bound the empirical
300 risk.

301 **Lemma 5.3.** *Let $\pi(x) \propto \exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$. Then for $\beta GD > d$, we know*

$$\mathbb{E}_{x \sim \pi} (F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2) - \min_{x^* \in \mathcal{K}} (F_{\mathcal{D}}(x^*) + \frac{\mu}{2}\|x^*\|_2^2) \leq \frac{d}{\beta} \log(\beta GD/d)$$

302 We now turn to bound the generalization error, and use the notion of uniform stability:

303 **Lemma 5.4** (Stability and Generalization [10]). *Given a dataset $\mathcal{D} = \{s_i\}_{i \in [n]}$ drawn i.i.d. from*
304 *some underlying distribution \mathcal{P} , and given any algorithm \mathcal{A} , suppose we randomly replace a*
305 *sample s in \mathcal{D} by an independent fresh one s' from \mathcal{P} and get the neighboring dataset \mathcal{D}' , then*
306 *$\mathbb{E}_{\mathcal{D}, \mathcal{A}}[F_{\mathcal{P}}(\mathcal{A}(\mathcal{D})) - F_{\mathcal{D}}(\mathcal{A}(\mathcal{D}))] = \mathbb{E}_{\mathcal{D}, s', \mathcal{A}}[f(\mathcal{A}(\mathcal{D}); s') - f(\mathcal{A}(\mathcal{D}'); s')]$, where $\mathcal{A}(\mathcal{D})$ is the out-*
307 *put of \mathcal{A} with input \mathcal{D} .*

308 As each function $f(\cdot; s')$ is G -Lipschitz, it suffices to bound the W_2 distance of $\mathcal{A}(\mathcal{D})$ and $\mathcal{A}(\mathcal{D}')$.
309 If \mathcal{A} is sampling from the exponential mechanism, letting $\pi_{\mathcal{D}} \propto \exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|^2))$ and
310 $\pi_{\mathcal{D}'} \propto \exp(-\beta(F_{\mathcal{D}'}(x) + \frac{\mu}{2}\|x\|^2))$, it suffices to bound the W_2 distance between $\pi_{\mathcal{D}}$ and $\pi_{\mathcal{D}'}$. The
311 following lemma can bound the generalization risk of the exponential mechanism under LSI:

312 **Lemma 5.5** (Generalization error bound). *Let $\pi_{\mathcal{D}} \propto \exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$. Then we have*
 313 $\mathbb{E}_{\mathcal{D}, x \sim \pi_{\mathcal{D}}} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)] \leq O\left(\frac{G^2 \exp(\beta GD)}{n\mu}\right)$.

314 We get the following results:

315 **Theorem 5.6** (Risk bound). *We are given $\varepsilon, \delta \in (0, 1/2)$. Sampling from $\exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$*
 316 *with $\beta = O\left(\frac{\varepsilon \log(nd)}{GD\sqrt{\log(1/\delta)}}\right)$, $\mu = \frac{d}{D^2\beta}$ is (ε, δ) -DP. The empirical risk and population risk are*
 317 *bounded by $O\left(GD \frac{d \cdot \log \log(n) \sqrt{\log(1/\delta)}}{\varepsilon \log(nd)}\right)$.*

318 **Implementation** There are multiple existing algorithms that can sample efficiently from density
 319 with LSI, under mild assumptions. For example, when the functions are smooth or weakly smooth,
 320 one can turn to the Langevin Monte Carlo [15], and [35]. The algorithm in [45] also requires mild
 321 smoothness assumptions. We discuss the implementation of non-smooth functions in bit more details,
 322 which is more challenging.

323 We can adopt the rejection sampler in [25], which is based on the alternating sampling algorithm
 324 in [34]. Both [34] and [25] are written in the language of log-concave and strongly log-concave
 325 densities, but their results hold as long as LSI holds. By combining them together, we can get the
 326 following risk bounds. The details of the implementation can be found in Appendix D.3.

327 **Theorem 5.7** (Implementation, risk bound). *For $\varepsilon, \delta \in (0, 1/2)$, there is an $(\varepsilon, 2\delta)$ -DP efficient*
 328 *sampler that can achieve the empirical and population risks $O\left(GD \frac{d \cdot \log \log(n) \sqrt{\log(1/\delta)}}{\varepsilon \log(nd)}\right)$. Moreover,*
 329 *in expectation, the sampler takes $\tilde{O}\left(n\varepsilon^3 \log^3(d) \sqrt{\log(1/\delta)} / (GD)\right)$ function values query and some*
 330 *Gaussian random variables restricted to the convex set \mathcal{K} in total.*

331 5.2 Exponential Time Approach

332 In [22], it is shown that sampling from $\exp(-\frac{\varepsilon n}{GD} F_{\mathcal{D}}(x))$ is ε -DP, and a nearly tight empirical risk
 333 bound of $\tilde{O}\left(\frac{DGd}{n\varepsilon}\right)$ is achieved for convex functions. It is open what is the bound we can get for
 334 non-convex DP-SO.

335 **Upper Bound** Given exponential time we can use a discrete exponential mechanism as considered
 336 in [9]. We recap the argument and extend it to DP-SO. The proof is based on a simple packing
 337 argument, and can be found in the Appendix.

338 **Theorem 5.8.** *There exists an ε -DP differentially private algorithm that achieves a population risk*
 339 *of $O\left(GD \left(d \log(\varepsilon n/d) / (\varepsilon n) + \sqrt{d \log(\varepsilon n/d)} / (\sqrt{n})\right)\right)$.*

340 **Lower Bound** Results in [22] imply that the first term of $\tilde{O}(GDd/\varepsilon n)$ is tight, even if we relax
 341 to approximate DP with $\delta > 0$. A reduction from private selection problem shows the $\tilde{O}(\sqrt{d/n})$
 342 generalization term is also nearly-tight (Theorem 5.11). In the selection problem, we have k coins,
 343 each with an unknown probability p_i . Each coin is flipped n times such that $\{x_{i,j}\}_{j \in [n]}$, each $x_{i,j}$
 344 i.i.d. sampled from $\text{Bern}(p_i)$, and we want to choose a coin i with the smallest p_i . The risk of
 345 choosing i is $p_i - \min_{i^*} p_{i^*}$.

346 **Theorem 5.9.** *Any algorithm for the selection problem has excess population risk $\tilde{\Omega}\left(\sqrt{\frac{\log k}{n}}\right)$.*

347 This follows from a folklore result on the selection problem (see e.g. [5]). We can combine this with
 348 the following reduction from selection to non-convex optimization:

349 **Theorem 5.10** (Restatement of results in [22]). *If any (ε, δ) -DP algorithm for selection has risk $R(k)$,*
 350 *then any (ε, δ) -DP algorithm for minimizing 1-Lipschitz losses over $B_d(0, 1)$ (the d -dimensional unit*
 351 *ball) has risk $R(2^{\Theta(d)})$.*

352 From this and the aforementioned lower bounds in empirical non-convex optimization we get the
 353 following:

354 **Theorem 5.11.** *For $\varepsilon \leq 1, \delta \in [2^{-\Omega(n)}, 1/n^{1+\Omega(1)}]$, any (ε, δ) -DP algorithm for minimizing 1-*
 355 *Lipschitz losses over $B_d(0, 1)$ has excess population risk $\max\{\Omega(d \log(1/\delta) / (\varepsilon n)), \tilde{\Omega}(\sqrt{d/n})\}$.*

References

- 356
- 357 [1] Naman Agarwal, Zeyuan Allen-Zhu, Brian Bullins, Elad Hazan, and Tengyu Ma. Finding
358 approximate local minima faster than gradient descent. In *Proceedings of the 49th Annual ACM*
359 *SIGACT Symposium on Theory of Computing*, pages 1195–1199, 2017.
- 360 [2] Raman Arora, Raef Bassily, Tomás González, Cristóbal Guzmán, Michael Menart, and Enayat
361 Ullah. Faster rates of convergence to stationary points in differentially private optimization.
362 *arXiv preprint arXiv:2206.00846*, 2022.
- 363 [3] Hilal Asi, Vitaly Feldman, Tomer Koren, and Kunal Talwar. Private stochastic convex optimiza-
364 tion: Optimal rates in ℓ_1 geometry. In *International Conference on Machine Learning*, pages
365 393–403. PMLR, 2021.
- 366 [4] Hilal Asi, Daniel Asher Nathan Levy, and John Duchi. Adapting to function difficulty and
367 growth conditions in private optimization. In *Advances in Neural Information Processing*
368 *Systems*, 2021.
- 369 [5] Mitali Bafna and Jonathan Ullman. The price of selection in differential privacy. In Satyen Kale
370 and Ohad Shamir, editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65
371 of *Proceedings of Machine Learning Research*, pages 151–168. PMLR, 07–10 Jul 2017.
- 372 [6] Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic
373 gradient descent on nonsmooth convex losses. *arXiv preprint arXiv:2006.06914*, 2020.
- 374 [7] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Thakurta. Private stochastic convex
375 optimization with optimal rates. In *Advances in Neural Information Processing Systems*, pages
376 11279–11288, 2019.
- 377 [8] Raef Bassily, Cristóbal Guzmán, and Michael Menart. Differentially private stochastic opti-
378 mization: New results in convex and non-convex settings. *Advances in Neural Information*
379 *Processing Systems*, 34:9317–9329, 2021.
- 380 [9] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization:
381 Efficient algorithms and tight error bounds. In *Proc. of the 2014 IEEE 55th Annual Symp. on*
382 *Foundations of Computer Science (FOCS)*, pages 464–473, 2014.
- 383 [10] Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine*
384 *Learning Research*, 2:499–526, 2002.
- 385 [11] Yair Carmon, Arun Jambulapati, Yujia Jin, Yin Tat Lee, Daogao Liu, Aaron Sidford, and
386 Kevin Tian. Resqueing parallel and private stochastic convex optimization. *arXiv preprint*
387 *arXiv:2301.00457*, 2023.
- 388 [12] Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. *Advances*
389 *in neural information processing systems*, 21, 2008.
- 390 [13] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical
391 risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.
- 392 [14] Yongxin Chen, Sinho Chewi, Adil Salim, and Andre Wibisono. Improved analysis for a
393 proximal algorithm for sampling. In *Conference on Learning Theory*, pages 2984–3014. PMLR,
394 2022.
- 395 [15] Sinho Chewi, Murat A Erdogdu, Mufan Li, Ruoqi Shen, and Shunshi Zhang. Analysis of
396 langevin monte carlo from poincare to log-sobolev. In *Conference on Learning Theory*, pages
397 1–2. PMLR, 2022.
- 398 [16] Rishav Chourasia, Jiayuan Ye, and Reza Shokri. Differential privacy dynamics of langevin
399 diffusion and noisy gradient descent. In *Advances in Neural Information Processing Systems*,
400 2021.
- 401 [17] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex
402 sgd. *Advances in neural information processing systems*, 32, 2019.

- 403 [18] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to
404 sensitivity in private data analysis. In *Proc. of the Third Conf. on Theory of Cryptography*
405 (*TCC*), pages 265–284, 2006.
- 406 [19] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Founda-*
407 *tions and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- 408 [20] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-
409 convex optimization via stochastic path-integrated differential estimator. *Advances in Neural*
410 *Information Processing Systems*, 31, 2018.
- 411 [21] Vitaly Feldman, Tomer Koren, and Kunal Talwar. Private stochastic convex optimization:
412 Optimal rates in linear time. In *Proc. of the Fifty-Second ACM Symp. on Theory of Computing*
413 (*STOC’20*), 2020.
- 414 [22] Arun Ganesh, Abhradeep Thakurta, and Jalaj Upadhyay. Langevin diffusion: An almost univer-
415 sal algorithm for private euclidean (convex) optimization. *arXiv preprint arXiv:2204.01585*,
416 2022.
- 417 [23] Changyu Gao and Stephen J Wright. Differentially private optimization for smooth nonconvex
418 erm. *arXiv preprint arXiv:2302.04972*, 2023.
- 419 [24] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online
420 stochastic gradient for tensor decomposition. In *Conference on learning theory*, pages 797–842.
421 PMLR, 2015.
- 422 [25] Sivakanth Gopi, Yin Tat Lee, and Daogao Liu. Private convex optimization via exponential
423 mechanism. In *Conference on Learning Theory*, pages 1948–1989. PMLR, 2022.
- 424 [26] Sivakanth Gopi, Yin Tat Lee, Daogao Liu, Ruoqi Shen, and Kevin Tian. Private convex
425 optimization in general norms. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on*
426 *Discrete Algorithms (SODA)*, pages 5068–5089. SIAM, 2023.
- 427 [27] Roger Iyengar, Joseph P Near, Dawn Song, Om Thakkar, Abhradeep Thakurta, and Lun Wang.
428 Towards practical differentially private convex optimization. In *2019 IEEE Symposium on*
429 *Security and Privacy (SP)*, 2019.
- 430 [28] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape
431 saddle points efficiently. In *International conference on machine learning*, pages 1724–1732.
432 PMLR, 2017.
- 433 [29] Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. A short
434 note on concentration inequalities for random vectors with subgaussian norm. *arXiv preprint*
435 *arXiv:1902.03736*, 2019.
- 436 [30] Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The composition theorem for differential
437 privacy. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International*
438 *Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*,
439 pages 1376–1385, Lille, France, 07–09 Jul 2015. PMLR.
- 440 [31] Gautam Kamath, Xingtu Liu, and Huanyu Zhang. Improved rates for differentially private
441 stochastic convex optimization with heavy-tailed data. In *International Conference on Machine*
442 *Learning*, pages 10633–10660. PMLR, 2022.
- 443 [32] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. Private convex empirical risk minimization
444 and high-dimensional regression. In *Conference on Learning Theory*, pages 25–1, 2012.
- 445 [33] Janardhan Kulkarni, Yin Tat Lee, and Daogao Liu. Private non-smooth erm and sco in sub-
446 quadratic steps. *Advances in Neural Information Processing Systems*, 34, 2021.
- 447 [34] Yin Tat Lee, Ruoqi Shen, and Kevin Tian. Structured logconcave sampling with a restricted
448 gaussian oracle. In *Conference on Learning Theory*, pages 2993–3050. PMLR, 2021.
- 449 [35] Jiaming Liang and Yongxin Chen. A proximal algorithm for sampling from non-smooth
450 potentials. In *2022 Winter Simulation Conference (WSC)*, pages 3229–3240. IEEE, 2022.

- 451 [36] Songtao Lu, Meisam Razaviyayn, Bo Yang, Kejun Huang, and Mingyi Hong. Finding second-
452 order stationary points efficiently in smooth nonconvex linearly constrained optimization prob-
453 lems. *Advances in Neural Information Processing Systems*, 33:2811–2822, 2020.
- 454 [37] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual*
455 *IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.
- 456 [38] Kentaro Minami, Hitomi Arai, Issei Sato, and Hiroshi Nakagawa. Differential privacy without
457 sensitivity. *Advances in Neural Information Processing Systems*, 29, 2016.
- 458 [39] Felix Otto and Cédric Villani. Generalization of an inequality by talagrand and links with the
459 logarithmic sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000.
- 460 [40] Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with
461 differentially private updates. In *2013 IEEE Global Conference on Signal and Information*
462 *Processing*, pages 245–248. IEEE, 2013.
- 463 [41] Shuang Song, Thomas Steinke, Om Thakkar, and Abhradeep Thakurta. Evading the curse
464 of dimensionality in unconstrained private glms. In *International Conference on Artificial*
465 *Intelligence and Statistics*, pages 2638–2646. PMLR, 2021.
- 466 [42] Kunal Talwar, Abhradeep Guha Thakurta, and Li Zhang. Nearly optimal private lasso. *Advances*
467 *in Neural Information Processing Systems*, 28, 2015.
- 468 [43] Hoang Tran and Ashok Cutkosky. Momentum aggregation for private non-convex erm. In
469 *Advances in Neural Information Processing Systems*, 2022.
- 470 [44] Joel A Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends®*
471 *in Machine Learning*, 8(1-2):1–230, 2015.
- 472 [45] Di Wang, Changyou Chen, and Jinhui Xu. Differentially private empirical risk minimization
473 with non-convex loss functions. In *International Conference on Machine Learning*, pages
474 6526–6535. PMLR, 2019.
- 475 [46] Di Wang and Jinhui Xu. Differentially private empirical risk minimization with smooth non-
476 convex loss functions: A non-stationary view. In *Proceedings of the AAAI Conference on*
477 *Artificial Intelligence*, volume 33, pages 1182–1189, 2019.
- 478 [47] Di Wang and Jinhui Xu. Escaping saddle points of empirical risk privately and scalably via
479 dp-trust region method. In *Joint European Conference on Machine Learning and Knowledge*
480 *Discovery in Databases*, pages 90–106. Springer, 2020.
- 481 [48] Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited:
482 Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017.
- 483 [49] Lingxiao Wang, Bargav Jayaraman, David Evans, and Quanquan Gu. Efficient privacy-
484 preserving stochastic nonconvex optimization. *arXiv preprint arXiv:1910.13659*, 2019.
- 485 [50] Yongqiang Wang and Tamer Başar. Decentralized nonconvex optimization with guaranteed
486 privacy and accuracy. *Automatica*, 150:110858, 2023.
- 487 [51] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. Spiderboost and momentum:
488 Faster variance reduction algorithms. *Advances in Neural Information Processing Systems*, 32,
489 2019.
- 490 [52] Yi Xu, Rong Jin, and Tianbao Yang. First-order stochastic algorithms for escaping from saddle
491 points in almost linear time. *Advances in neural information processing systems*, 31, 2018.
- 492 [53] Xiaodong Yang, Huishuai Zhang, Wei Chen, and Tie-Yan Liu. Normalized/clipped sgd with per-
493 turbation for differentially private non-convex optimization. *arXiv preprint arXiv:2206.13033*,
494 2022.
- 495 [54] Qiuchen Zhang, Jing Ma, Jian Lou, and Li Xiong. Private stochastic non-convex optimization
496 with improved utility rates. In *Proceedings of the Thirtieth International Joint Conference on*
497 *Artificial Intelligence*, 2021.

- 498 [55] Yingxue Zhou, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Arindam Banerjee. Private
499 stochastic non-convex optimization: Adaptive algorithms and tighter generalization bounds.
500 *arXiv preprint arXiv:2006.13501*, 2020.

501 A Other Preliminary

502 **Definition A.1** (Laplace distribution). We say $X \sim \text{Lap}(b)$ if X has density $f(X = x) =$
 503 $\frac{1}{2b} \exp(\frac{-|x|}{b})$.

504 **Theorem A.2** (Basic composition, [19]). If \mathcal{A}_1 is $(\varepsilon_1, \delta_1)$ -DP and \mathcal{A}_2 is $(\varepsilon_2, \delta_2)$ -DP, then their
 505 combination is $(\varepsilon_1 + \varepsilon_2, \delta_1 + \delta_2)$ -DP.

506 **Theorem A.3** (Advanced composition, [30]). For $\varepsilon \leq 0.9$, an end-to-end guarantee of (ε, δ) -
 507 differential privacy is satisfied if a database is accessed at most k times, where each time with a
 508 $(\varepsilon/(2\sqrt{2k \log(2/\delta)}), \delta/(2k))$ -differentially private mechanism.

509 Due to space limit, some other preliminaries and proofs are left in the Appendix.

510 **Theorem A.4** (Gaussian Mechanism, [19]). Given a randomized algorithm $\mathcal{A} : P^* \rightarrow \mathbb{R}^d$, let
 511 $\Delta_2 f = \max_{\text{neighboring } \mathcal{D}, \mathcal{D}'} \|\mathcal{A}(\mathcal{D}) - \mathcal{A}(\mathcal{D}')\|_2$, then adding noise scaled to $\mathcal{N}(0, \sigma^2)$ with $\sigma \geq$
 512 $\frac{\sqrt{2 \log(1.25/\delta)} \Delta_2 f}{\varepsilon}$ is (ε, δ) -DP.

513 **Theorem A.5** (Matrix Bernstein inequality, [44]). Consider a sequence $\{X_i\}_{i \in [m]}$ of independent,
 514 mean-zero, symmetric $d \times d$ random matrices. If for each matrix X_i , we know $\|X_i\|_{op} \leq M$,
 515 then for all $t \geq 0$, we have $\Pr \left[\left\| \sum_{i \in [m]} X_i \right\|_{op} \geq t \right] \leq d \exp \left(\frac{-t^2}{2(\sigma^2 + Mt/3)} \right)$, where $\sigma^2 =$
 516 $\left\| \sum_{i \in [m]} \mathbb{E} X_i^2 \right\|_{op}$.

517 B Omitted Proof of Section 3

518 B.1 Proof of Lemma 3.3

519 **Lemma 3.3.** For any $0 \leq t \leq T$ and letting $\tau_t \leq t$ be the largest integer such that drift $_{\tau_t}$ is set to
 520 be 0, with probability at least $1 - \omega/T$, for some universal constant $C > 0$, we have

$$\|\nabla_t - \nabla F(x_t)\|^2 \leq (\zeta_2^2 \cdot \sum_{i=\tau_t+1}^t \|x_i - x_{i-1}\|^2 + 4\zeta_1^2) \cdot C \cdot \log(Td/\omega). \quad (1)$$

521 Hence with probability at least $1 - \omega$, we know for each $t \leq T$, $\|\nabla_t - \nabla F(x_t)\|^2 \leq \gamma^2/16$, where
 522 $\gamma^2 := 16C(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(Td/\omega)$ and κ is a parameter we can choose in the algorithm.

523 *Proof.* If drift $_{\tau_t} = 0$ happens, we use the first kind oracle to query the gradient, and hence $\nabla_{\tau_t} -$
 524 $\nabla F(x_{\tau_t})$ is zero-mean and nSG($2\zeta_1$). If $t = \tau_t$, Equation (1) holds by the property of norm-
 525 subGaussian.

526 For each $\tau_t + 1 \leq i \leq t$, conditional on ∇_{i-1} , we know $\Delta_i - (\nabla F(x_i) - \nabla F(x_{i-1}))$ is zero-mean
 527 and nSG($\zeta_2 \|x_i - x_{i-1}\|$). Note that

$$\nabla_t - \nabla F(x_t) = \nabla_{\tau_t} - \nabla F(x_{\tau_t}) + \sum_{i=\tau_t+1}^t [\Delta_i - (\nabla F(x_i) - \nabla F(x_{i-1}))].$$

528 Equation (1) follows from Lemma 2.4.

529 We know drift $_{t-1} = \sum_{i=\tau_t+1}^t \|x_i - x_{i-1}\|^2 \leq \kappa$ almost surely by the design of the algorithm. By
 530 union bound, we know with probability at least $1 - \omega$, for each $t \in [T]$,

$$\|\nabla_t - \nabla F(x_t)\|^2 \leq C(\zeta_2^2 \kappa + 4\zeta_1^2) \cdot \log(Td/\omega) = \gamma^2/16.$$

531 □

532 B.2 Discussion of Lemma 3.4

533 **Lemma 3.4** (Essentially from [45]). Under Assumption 3.1, run SGD iterations $x_{t+1} = x_t -$
 534 $\eta \nabla_t$, with step size $\eta = 1/M$. Suppose x_0 is a saddle point satisfying $\|\nabla F(x_0)\| \leq \alpha$ and
 535 $\text{smin}(\nabla^2 F(x_0)) \leq -\sqrt{\rho\alpha}$, $\alpha = \gamma \log^3(dBM/\rho\omega)$. If $\nabla_0 = \nabla F(x_0) + \zeta_1 + \zeta_2$ where $\|\zeta_1\| \leq \gamma$,

536 $\zeta_2 \sim \mathcal{N}(0, \frac{\gamma^2}{d \log(d/\omega)} I_d)$, and $\|\nabla_t - \nabla F(x_t)\| \leq \gamma$ for all $t \in [\Gamma]$, with probability at least
 537 $1 - \omega \cdot \log(1/\omega)$, one has $F(x_\Gamma) - F(x_0) \leq -\Omega(\frac{\gamma^{3/2}}{\sqrt{\rho} \log^3(\frac{dMB}{\rho\gamma\omega})})$, where $\Gamma = \frac{M \log(\frac{dMB}{\rho\gamma\omega})}{\sqrt{\rho\gamma}}$.

538 We briefly recap the proof of Lemma 3.4 in [45]. One observation between the decreased function
 539 value, and the distance solutions moved is stated below:

540 **Lemma B.1** (Lemma 11, [45]). *For each $t \in [\Gamma]$, we know*

$$\|x_{t+1} - x_0\|_2^2 \leq 8\eta(\Gamma(F(x_0) - F(x_\Gamma))) + 50\eta^2\Gamma \sum_{i \in [\Gamma]} \|\nabla_i - \nabla F(x_t)\|_2^2.$$

541 The difference between our algorithm and the DP-GD in [45] is the noise on the gradient. Note that
 542 with high probability, $\sum_{i \in [\Gamma]} \|\nabla_i - \nabla F(x_t)\|_2^2$ in our algorithm is controlled and small, and hence
 543 does not change the other proofs in [45]. Hence if $F(x_0) - F(x_\Gamma)$ is small, i.e., the function value
 544 does not decrease significantly, we know x_t is close to x_0 .

545 Let $B_x(r)$ be the unit ball of radius r around point x . Denote the $(x)_\Gamma$ the point x_Γ after running
 546 SGD mentioned in Lemma 3.4 for Γ steps beginning at x . With this observation, denote $B^\gamma(x_0) :=$
 547 $\{x \mid x \in B_{x_0}(\eta\alpha), \Pr[F((x)_\Gamma) - F(x) \geq -\Phi] \geq \omega\}$. [45] demonstrates the following lemma:

548 **Lemma B.2.** *If $\|\nabla F(x_0)\| \leq \alpha$ and $\text{smin}(\nabla^2 F(x_0)) \leq -\sqrt{\rho\gamma}$, then the width of $B^\gamma(x_0)$ along the
 549 along the minimum eigenvector of $\nabla^2 F(x_0)$ is at most $\frac{\omega\eta\gamma}{\log(1/\omega)} \sqrt{\frac{2\pi}{d}}$.*

550 The intuition is that if two different points $x^1, x^2 \in B_{x_0}(\eta\alpha)$, and $x^1 - x^2$ is large along the minimum
 551 eigenvector, then with high probability, the distance between $\|(x^1)_\Gamma - (x^2)_\Gamma\|$ will be large, and either
 552 $\|(x^1)_\Gamma - x^1\|$ or $\|(x^2)_\Gamma - x^2\|$ is large, and hence either $F(x^1) - F((x^1)_\Gamma)$ or $F(x^2) - F((x^2)_\Gamma)$
 553 is large. The Lemma 3.4 follows from Lemma B.2 by using the Gaussian ζ_2 to kick off the point.

554 B.3 Proof of Lemma 3.5

555 **Lemma 3.5.** *By setting $\eta = 1/M$, we have $F(x_{t+1}) \leq F(x_t) + \eta\|\nabla_t\| \cdot \|\nabla F(x_t) - \nabla_t\| - \frac{\eta}{2}\|\nabla_t\|^2$.
 556 Moreover, with probability at least $1 - \omega$, for each $t \leq T$ such that $\|\nabla F(x_t)\| \geq \gamma$, we have*

$$F(x_{t+1}) - F(x_t) \leq -\eta\|\nabla_t\|^2/6 \leq -\eta\gamma^2/6.$$

557 *Proof.* By the assumption on smoothness, we know

$$\begin{aligned} F(x_{t+1}) &\leq F(x_t) + \langle \nabla F(x_t), x_{t+1} - x_t \rangle + \frac{M}{2} \|x_{t+1} - x_t\|^2 \\ &= F(x_t) - \eta/2 \|\nabla_t\|^2 - \langle \nabla F(x_t) - \nabla_t, \eta \nabla_t \rangle \\ &\leq F(x_t) + \eta \|\nabla F(x_t) - \nabla_t\| \cdot \|\nabla_t\| - \frac{\eta}{2} \|\nabla_t\|^2. \end{aligned}$$

558 By Lemma 3.3, with probability at least $1 - \omega$, for each $t \in [T]$ we have $\|\nabla F(x_t) - \nabla_t\|_2 \leq \gamma/4$.
 559 Hence we know if $\|\nabla F(x_t)\| \geq \gamma$, we have

$$F(x_{t+1}) - F(x_t) \leq -\eta\|\nabla_t\|^2/6 \leq -\eta\gamma^2/6.$$

560 □

561 B.4 Proof of Lemma 3.6

562 **Lemma 3.6.** *Suppose \mathcal{O}_1 and \mathcal{O}_2 are ζ_1 and ζ_2 norm-subGaussian respectively. If one sets $\gamma =$
 563 $O(1)\sqrt{(\zeta_2^2\kappa + 4\zeta_1^2) \cdot \log(Td/\omega)}$, with probability at least $1 - \omega$, at least one point in the output set
 564 $\{x_1, \dots, x_T\}$ of Algorithm 1 is α -SOSP, where*

$$\alpha = \gamma \log^3(BMd/\rho\omega\gamma) = \sqrt{(\zeta_2^2\kappa + 4\zeta_1^2) \cdot \log(\frac{d/\omega}{\zeta_2^2\kappa + \zeta_1^2})} \cdot \log^3\left(\frac{BMd}{\rho\omega(\zeta_2^2\kappa + \zeta_1^2)}\right).$$

565 *Proof.* By Lemma 3.5, we know if the gradient $\|\nabla F(x_t)\| \geq \gamma$, then with high probability that
 566 $F(x_{t+1}) - F(x_t) \leq -\eta\gamma^2/6$. By Lemma 3.4, if x_t is a saddle point (with small gradient norm
 567 but the Hessian has a small eigenvalue), then with high probability that $F(x_{\Gamma+t}) - F(x_t) \leq$
 568 $-\Omega\left(\frac{\gamma^{3/2}}{\sqrt{\rho} \log^3\left(\frac{dMB}{\rho\gamma\omega}\right)}\right)$, and the function values decrease $\Omega\left(\frac{\gamma^2}{M \log^4\left(\frac{dMB}{\rho\gamma\omega}\right)}\right)$ on average for each step.

569 Recall the assumption that the risk is upper bounded by B , by our setting $T = \Omega\left(\frac{BM}{\gamma^2} \log^4\left(\frac{dMB}{\rho\gamma\omega}\right)\right)$,
 570 the statement is proved. \square

571 B.5 Proof of Lemma 3.7

572 **Lemma 3.7.** *Under the event that $\|\nabla_t - \nabla F(x_t)\| \leq \gamma/4$ for all $t \in [T]$ and our parameter settings,*
 573 *letting $K = \{t \in [T] : \text{drift}_t \geq \kappa\}$ be the set of iterations where the drift is large, we know*
 574 $|K| \leq O\left(\frac{B\eta}{\kappa} + T\gamma^2\eta^2/\kappa\right) = O\left(B\eta \log^4\left(\frac{dMB}{\rho\gamma\omega}\right)/\kappa\right)$.

575 *Proof.* By Lemma 3.5, if $\|F(x_t)\| \geq \gamma$, we know $F(x_{t+1}) - F(x_t) \leq -\eta\|\nabla_t\|^2/6$, and $F(x_{t+1}) -$
 576 $F(x_t) \leq \eta\gamma^2$ otherwise. Index the items in $K = \{t_1, t_2, \dots, t_{|K|}\}$ such that $t_i < t_{i+1}$. We know

$$F(x_{t_{i+1}}) - F(x_{t_i}) \leq -\frac{1}{6\eta} \text{drift}_{t_{i+1}} + (t_{i+1} - t_i)\gamma^2\eta \leq -\frac{1}{6\eta}\kappa + (t_{i+1} - t_i)\gamma^2\eta.$$

577 Recall by the assumption that $\max_y F(y) - \min_x F(x) \leq B$. And hence $-B \leq F(x_{t_{|K|}}) - F(x_{t_1}) \leq$
 578 $-\frac{|K|}{6\eta}\kappa + T\gamma^2\eta$, and we know

$$|K| \leq O\left(\frac{B\eta}{\kappa} + T\gamma^2\eta^2/\kappa\right) = O\left(B\eta \log^4\left(\frac{dMB}{\rho\gamma\omega}\right)/\kappa\right).$$

579 \square

580 C Appendix for Section 4

The pseudocode of Algorithm 2 is stated below:

Algorithm 2 AboveThreshold

- 1: **Input:** A set of points $\{x_i\}_{i=1}^T$, dataset S , parameters of objective function B, M, G, ρ , objective error α
 - 2: Set $\widehat{T}_1 = \alpha + \text{Lap}\left(\frac{4G}{n\varepsilon}\right) + \frac{16 \log(2T/\omega)G}{n\varepsilon}$, $\widehat{T}_2 = -\sqrt{\rho\alpha} + \text{Lap}\left(\frac{4M}{n\varepsilon}\right) - \frac{16 \log(2T/\omega)M}{n\varepsilon}$
 - 3: **for** $i = 1, \dots, T$ **do**
 - 4: **if** $\|\nabla F_S(x_i)\| + \text{Lap}\left(\frac{8G}{n\varepsilon}\right) \leq \widehat{T}_1 \wedge \text{smin}(\nabla^2 F_S(x_i)) + \text{Lap}\left(\frac{8M}{n\varepsilon}\right) \geq \widehat{T}_2$ **then**
 - 5: **Output:** x_i
 - 6: **Halt**
 - 7: **end if**
 - 8: **end for**
-

581

582 C.1 Proof of Theorem 4.2

583 **Theorem 4.2** (Empirical). *For $\varepsilon \leq 10, \delta \in (0, 1/2)$, use Equation (2) as gradient or-*
 584 *acles with $\kappa = \frac{G^{4/3} B^{1/3}}{M^{5/3}} \left(\frac{\sqrt{d \log(1/\delta)}}{n\varepsilon}\right)^{2/3}$, $\sigma_1 = \frac{G\sqrt{B\eta \log^2(n/\delta)/\kappa \log^2(ndMB/\omega)}}{n\varepsilon}$, $\sigma_2 =$*
 585 $\frac{M\sqrt{\log^2(n/\delta)BM/\alpha_1^2 \log^5(ndMB/\omega)}}{n\varepsilon}$. *Running Algorithm 1, outputting the set $\{x_i\}_{i \in [T]}$ if the total*
 586 *time to query \mathcal{O}_1 is bounded by $O\left(B\eta \log^4\left(\frac{dMB}{\rho\gamma\omega}\right)/\kappa\right)$, otherwise outputting a set of T arbitrary*
 587 *points is $(\varepsilon/2, \delta)$ -DP. With probability at least $1 - \omega$, at least one point in the output set is α_1 -SOSP*
 588 *of $F_{\mathcal{D}}$ with*

$$\alpha_1 = O\left(\left(\frac{\left(\sqrt{dBGM \log^2(1/\delta)}\right)^{2/3}}{n\varepsilon}\right) \cdot \log^6\left(\frac{nBMd}{\rho\omega}\right)\right).$$

589 Moreover, if we run Algorithm 2 with inputs $\{x_i\}_{i \in [T]}$, \mathcal{D} , B , M , G , ρ , α_1 , with probability at least
590 $1 - \omega$, we can get an α_2 -SOSP of $F_{\mathcal{D}}$ with $\alpha_2 = O\left(\alpha_1 + \frac{G \log(n/G\omega)}{n\varepsilon} + \frac{M \log(ndBGM/\rho\omega)}{n\varepsilon\sqrt{\rho}}\sqrt{\alpha_1}\right)$.
591 The whole procedure is (ε, δ) -DP.

592 *Proof.* The privacy guarantee can be proved by composition theorems (Theorem A.2 and Theo-
593 rem A.3), Gaussian Mechanism (Theorem A.4) and Lemma 3.7. Specifically, by Assumption 3.1 and
594 our settings of parameters, we know the sensitivity of \mathcal{O}_1 and \mathcal{O}_2 are bounded by $\frac{G}{n}$ and $\frac{M\|x-y\|}{n}$
595 respectively, and querying \mathcal{O}_1 and \mathcal{O}_2 each time are $(\frac{\varepsilon}{\sqrt{B\eta \log(n/\delta) \log^2(ndMB/\omega)}}, \delta/n^2)$ -DP and
596 $(\frac{\varepsilon}{\sqrt{\log(n/\delta)BM/\alpha_1^2 \log^5(ndMB/\omega)}}, \delta/n^2)$ -DP respectively. We can apply the advanced composition to
597 prove the privacy guarantee of the whole algorithm. As the total number of iterations T is determined,
598 and the privacy cost to query \mathcal{O}_2 for T times is controlled. It suffices to bound the total time to
599 query \mathcal{O}_1 , which is guaranteed in the statement. That is if the total time to query \mathcal{O}_1 is bounded by
600 $O(B\eta \log^4(\frac{dMB}{\rho\gamma\omega})/\kappa)$, the privacy guarantee follows from the advanced composition. If the time
601 exceeds $O(B\eta \log^4(\frac{dMB}{\rho\gamma\omega})/\kappa)$, then we will output a set of arbitrary points which does not occur
602 additional privacy cost.

603 As for the utility, we know the \mathcal{O}_1 and \mathcal{O}_2 constructed in Equation (2) are first kind of $\sigma_1\sqrt{d}$ and
604 second kind of $\sigma_2\sqrt{d}$ norm-subGaussian gradient oracle by Fact 2.3. Hence by Lemma 3.6, the utility
605 α_1 satisfies that

$$\begin{aligned} \alpha_1 &= O(\sigma_1\sqrt{d} + \sigma_2\sqrt{d\kappa}) \cdot \log^3(BMd/\rho\omega) \\ &= O\left(\frac{L\sqrt{dB\eta \log^2(1/\delta)/\kappa}}{n\varepsilon} + \frac{M \log^3(ndMB/\omega)\sqrt{\log^2(1/\delta)BM}}{n\varepsilon\alpha_1}\sqrt{d\kappa}\right) \cdot \log^5(nBMd/\rho\omega). \end{aligned}$$

606 By Lemma 3.7, with probability at least $1 - \omega$, the total time to query \mathcal{O}_1 is controlled and the final
607 output will not be arbitrary points. Choosing the best κ demonstrates the bound on α_1 . The bound
608 for α_2 follows from the value of α_1 and Lemma 4.1. Combining the two items in Lemma 4.1, we
609 know with probability at least $1 - \omega$, the output point x of Algorithm 2 satisfies that

$$\|\nabla F_{\mathcal{D}}(x)\| \leq \alpha_1 + \frac{32 \log(2T/\omega)G}{n\varepsilon}, \text{ and } \text{smin}(\nabla^2 F_{\mathcal{D}}(x)) \geq -\sqrt{\rho\alpha_1} - \frac{32 \log(2T/\omega)M}{n\varepsilon}.$$

610 Hence we know x is an α_2 -SOSP for α_2 stated in the statement. \square

611 C.2 Proof of Claim 4.4

612 **Claim 4.4.** The gradient oracles \mathcal{O}_1 and \mathcal{O}_2 constructed in Equation (3) are a first kind of
613 $O(\frac{L\sqrt{\log d}}{\sqrt{b_1}} + \sqrt{d}\sigma_1)$ norm-subGaussian gradient oracle and second kind of $O(\frac{M\sqrt{\log d}}{\sqrt{b_2}} + \sqrt{d}\sigma_2)$
614 norm-subGaussian gradient oracle respectively.

615 *Proof.* For the oracle \mathcal{O}_1 , we know for each $z \in S_1$, $\mathbb{E}_{z \sim \mathcal{P}}[\nabla f(x, z)] = \nabla F_{\mathcal{P}}(x)$ and $\nabla f(x, z) -$
616 $\nabla F_{\mathcal{P}}(x)$ is nSG(L) due to the Lipschitzness assumption. The statement follows from Fact 2.3 and
617 Lemma 2.4. As for the \mathcal{O}_2 , the statement follows similarly with the smoothness assumption. \square

618 C.3 Proof of Lemma 4.5

619 **Lemma 4.5.** Fix a point $x \in \mathbb{R}^d$. Given a set S of m samples drawn i.i.d. from the distribution \mathcal{P} ,
620 then we know with probability at least $1 - \omega$, we have

$$\|\nabla F_S(x) - \nabla F_{\mathcal{P}}(x)\|_2 \leq O\left(\frac{G \log(d/\omega)}{\sqrt{m}}\right) \wedge \|\nabla^2 F_S(x) - \nabla^2 F_{\mathcal{P}}(x)\|_{op} \leq O\left(\frac{M \log(d/\omega)}{\sqrt{m}}\right).$$

621 *Proof.* As for any $s \in S$, $\nabla f(x; s) - \nabla F_{\mathcal{P}}(x)$ is zero-mean nSG(G). Then the Hoeffding inequality
622 for norm-subGaussians (Lemma 2.4) demonstrates with probability at least $1 - \omega/2$, we have
623 $\|\nabla F_S(x) - \nabla F_{\mathcal{P}}(x)\|_2 \leq O\left(\frac{G \log(d/\omega)}{\sqrt{m}}\right)$.

624 As for the other term, we know for any $s \in S$, $\mathbb{E}[\nabla^2 f(x; s) - \nabla^2 F_{\mathcal{P}}(x)] = 0$, and $\|\nabla^2 f(x; s) -$
625 $\nabla^2 F_{\mathcal{P}}(x)\|_{op} \leq 2M$ almost surely. Hence applying Matrix Bernstein inequality (Theorem A.5) with
626 $\sigma^2 = 4M^2m$, $t = O(\sqrt{m}M \log(d/\omega))$, we know with probability at least $1 - \omega/2$, $\|\nabla^2 F_S(x) -$
627 $\nabla^2 F_{\mathcal{P}}(x)\|_{op} \leq t/m$.

628 Applying the Union bound completes the proof. \square

629 C.4 Proof of Theorem 4.6

630 **Theorem 4.6** (Population). *Divide the dataset \mathcal{D} into two disjoint datasets \mathcal{D}_1 and \mathcal{D}_2 of size*
631 *$\lfloor n/2 \rfloor$ and $\lfloor n/2 \rfloor$ respectively. Set $b_1 = \frac{n\kappa}{B\eta}$, $b_2 = \frac{n\alpha_1^2}{BM}$, $\sigma_1 = \frac{3G\sqrt{\log(1/\delta)}}{b_1\varepsilon}$, $\sigma_2 = \frac{3M\sqrt{\log(1/\delta)}}{b_2\varepsilon}$*
632 *and $\kappa = \max(\frac{G^{4/3}B^{1/3}\log^{1/3}d}{M^{5/3}}n^{-1/3}, (\frac{GB^{2/3}}{M^{5/3}})^{6/7}(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon})^{4/7})$ in Equation (3) and use them as*
633 *gradient oracles. Running Algorithm 1 with \mathcal{D}_1 , and outputting the set $\{x_i\}_{i \in [T]}$ if the total time to*
634 *query \mathcal{O}_1 is bounded by $O(B\eta \log^4(\frac{dMB}{\rho\omega})/\kappa)$, otherwise outputting a set of T arbitrary points, is*
635 *$(\varepsilon/2, \delta)$ -DP. is $(\varepsilon/2, \delta)$ -DP, and with probability at least $1 - \omega$, at least one point in the output is*
636 *α_1 -SOSP of $F_{\mathcal{P}}$ with*

$$\alpha_1 = O\left(\left((BGM \cdot \log d)^{1/3} \frac{1}{n^{1/3}} + (G^{1/7}B^{3/7}M^{3/7})\left(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon}\right)^{3/7}\right) \log^3(nBMd/\rho\omega)\right).$$

637 *Moreover, if we run Algorithm 2 with inputs $\{x_i\}_{i \in [T]}$, \mathcal{D}_2 , B , M , G , ρ , α_1 , with prob-*
638 *ability at least $1 - \omega$, Algorithm 2 can output an α_2 -SOSP of $F_{\mathcal{P}}$ with $\alpha_2 =$*
639 *$O\left(\alpha_1 + \frac{M \log(ndBGM/\rho\omega)}{\sqrt{\rho} \min(n\varepsilon, n^{1/2})} \sqrt{\alpha_1} + G\left(\frac{\log(n/G\omega)}{n\varepsilon} + \frac{\log(d/\omega)}{\sqrt{n}}\right)\right)$. The whole procedure is (ε, δ) -DP.*

640 *Proof.* Recall that we draw the samples to construct the gradient oracles (Equation 3) without
641 replacement, and we should have all samples to be fresh to avoid dependency, and hence we need

$$b_1 \cdot |K| + b_2 \cdot T \leq n/2,$$

642 which is satisfied by the procedure in the statement, as if the total time to query the \mathcal{O}_1 exceeds the
643 threshold, the algorithm fails and outputs a set of arbitrary points. As we never reuse a sample, the
644 privacy guarantee follows directly from the Gaussian Mechanism [19]. Specifically, the sensitivity of
645 querying \mathcal{O}_1 and \mathcal{O}_2 are bounded by G/b_1 and $M\|x - y\|/b_2$ respectively, and querying \mathcal{O}_1 and \mathcal{O}_2
646 are $(\varepsilon/3, \delta)$ -DP by Theorem A.4.

647 The Norm-subGaussian parameters of the oracles follow from Claim 4.4. By lemma 3.6, we have

$$\begin{aligned} & \frac{\alpha_1}{\log^3(nBMd/\rho\omega)} \\ &= O\left(\sigma_1 \sqrt{d} + \frac{G\sqrt{\log d}}{\sqrt{b_1}} + \sigma_2 \sqrt{d\kappa} + \frac{M\sqrt{\kappa \log d}}{\sqrt{b_2}}\right). \\ &= O\left(\frac{GB\eta\sqrt{d\log(1/\delta)}}{n\varepsilon\kappa} + \frac{BM^2\sqrt{\log(1/\delta)}}{n\varepsilon\alpha_1^2} \sqrt{d\kappa} + \frac{G\sqrt{B\eta \log d}}{\sqrt{n\kappa}} + M\sqrt{\kappa} \frac{\sqrt{BM \log d}}{\sqrt{n\alpha_1}}\right). \end{aligned}$$

648 Setting $\kappa = \max(\frac{G^{4/3}B^{1/3}\log^{1/3}d}{M^{5/3}}(n)^{-1/3}, (\frac{GB^{2/3}}{M^{5/3}})^{6/7}(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon})^{4/7})$, we get

$$\alpha_1 = O\left(\left((BGM \log d)^{1/3} \frac{1}{n^{1/3}} + (G^{1/7}B^{3/7}M^{3/7})\left(\frac{\sqrt{d\log(1/\delta)}}{n\varepsilon}\right)^{3/7}\right) \log^3(nBMd/\rho\omega)\right).$$

649 Then we use the other half fresh samples \mathcal{D}_2 to find the point in the set by Algorithm 2. By Lemma 4.1
650 and Lemma 4.5, we know with probability at least $1 - \omega$, for some large enough constant $C > 0$, the
651 output point x of Algorithm 2 satisfies that

$$\begin{aligned} \|\nabla F_{\mathcal{P}}(x)\|_2 &\leq \alpha_1 + G\left(\frac{32 \log(2T/\omega)}{n\varepsilon} + \frac{C \log(dn/\omega)}{\sqrt{n}}\right), \\ \text{smin}(\nabla^2 F_{\mathcal{P}}(x)) &\geq -\sqrt{\rho\alpha_1} - M\left(\frac{32 \log(2T/\omega)}{n\varepsilon} + \frac{C \log(dn/\omega)}{\sqrt{n}}\right) \end{aligned}$$

652 Hence we know x is an α_2 -SOSP for α_2 stated in the statement. The privacy guarantee follows from
653 Basic composition and Lemma 4.1. \square

654 **D Omitted proof of Section 5**

655 **D.1 Proof of Lemma 5.5**

656 **Lemma 5.5** (Generalization error bound). *Let $\pi_{\mathcal{D}} \propto \exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$. Then we have*

657 $\mathbb{E}_{\mathcal{D}, x \sim \pi_{\mathcal{D}}} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)] \leq O\left(\frac{G^2 \exp(\beta GD)}{n\mu}\right).$

658 *Proof.* We know how to bound the KL divergence by LSI:

$$\begin{aligned} KL(\pi_{\mathcal{D}}, \pi_{\mathcal{D}'}) &:= \int \log \frac{d\pi_{\mathcal{D}}}{d\pi_{\mathcal{D}'}} d\pi_{\mathcal{D}} \\ &\leq \frac{C_{\text{LSI}}}{2} \mathbb{E}_{\pi_{\mathcal{D}}} \left\| \nabla \log \frac{d\pi_{\mathcal{D}}}{d\pi_{\mathcal{D}'}} \right\|_2^2 \\ &\leq 2C_{\text{LSI}} G^2 \beta^2 / n^2. \end{aligned}$$

659 LSI can lead to Talagrand transportation inequality [Theorem 1 in [39]], i.e.,

$$W_2(\pi_{\mathcal{D}}, \pi_{\mathcal{D}'}) \lesssim \sqrt{C_{\text{LSI}} \cdot KL(\pi_{\mathcal{D}}, \pi_{\mathcal{D}'})} = C_{\text{LSI}} G \beta / n.$$

660 The generalization error is bounded by $O(C_{\text{LSI}} G^2 \beta / n)$. Using Holley-Stroock perturbation, we

661 know $C_{\text{LSI}}(\pi_{\mathcal{D}}) \leq \frac{\exp(\beta GD)}{\beta \mu}$ and hence the W_2 distance between $\pi_{\mathcal{D}}$ and $\pi_{\mathcal{D}'}$ can be bounded by

662 $O\left(\frac{G \exp(\beta GD)}{n\mu}\right)$. The statement follows the Lipschitzness constant and Lemma 5.4. \square

663 **D.2 Proof of Theorem 5.6**

664 **Theorem 5.6** (Risk bound). *We are given $\varepsilon, \delta \in (0, 1/2)$. Sampling from $\exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$*

665 *with $\beta = O\left(\frac{\varepsilon \log(nd)}{GD \sqrt{\log(1/\delta)}}\right)$, $\mu = \frac{d}{D^2 \beta}$ is (ε, δ) -DP. The empirical risk and population risk are*

666 *bounded by $O\left(GD \frac{d \cdot \log \log(n) \sqrt{\log(1/\delta)}}{\varepsilon \log(nd)}\right)$.*

667 *Proof.* Denote $\pi(x) \propto \exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$. By Lemma 5.2, we know $C_{\text{LSI}}(\pi) \leq \frac{1}{\beta \mu} \cdot$

668 $\exp(\beta GD)$. Plugging in the parameters and applying Theorem 5.1, we get

$$\frac{2G\beta}{n} \cdot \sqrt{\frac{\exp(\beta GD)}{\beta \mu}} \sqrt{3 \log(1/\delta)} = O(1) \frac{GD\beta}{n\sqrt{d}} \sqrt{\exp(\beta GD) \log(1/\delta)} \leq 1$$

669 and hence prove the privacy guarantee.

670 As for the empirical risk bound, by Lemma 5.3, we know

$$\mathbb{E}_{x \sim \pi} (F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2) - \min_{x^* \in \mathcal{K}} (F_{\mathcal{D}}(x^*) + \frac{\mu}{2}\|x^*\|_2^2) \lesssim \frac{d \log(\beta GD/d)}{\beta},$$

671 and we know

$$\mathbb{E}_{x \sim \pi} F_{\mathcal{D}}(x) - \min_{x^* \in \mathcal{K}} F_{\mathcal{D}}(x^*) \lesssim \frac{d \log(\beta GD/d)}{\beta} + \mu D^2.$$

672 Replacing the value of β achieves the empirical risk bound.

673 As for the population risk, we have

$$\begin{aligned} &\mathbb{E}_{x \sim \pi} F_{\mathcal{P}}(x) - \min_{y^* \in \mathcal{K}} F_{\mathcal{P}}(y^*) \\ &= \mathbb{E}_{x \sim \pi} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)] + \mathbb{E}[F_{\mathcal{D}}(x) - \min_{x^* \in \mathcal{K}} F_{\mathcal{D}}(x^*)] + \mathbb{E}[\min_{x^* \in \mathcal{D}} F_{\mathcal{D}}(x^*) - \min_{y^* \in \mathcal{K}} F_{\mathcal{P}}(y^*)] \\ &\leq \mathbb{E}_{x \sim \pi} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)] + \mathbb{E}[F_{\mathcal{D}}(x) - \min_{x^* \in \mathcal{K}} F_{\mathcal{D}}(x^*)]. \end{aligned}$$

674 We can bound $\mathbb{E}_{x \sim \pi} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)] \leq O\left(\frac{G^2 \exp(\beta GD)}{n\mu}\right) \leq O\left(\frac{GD\varepsilon \log(n)}{n^{1-c} d \sqrt{\log(1/\delta)}}\right)$ by Lemma 5.5

675 for an arbitrarily small constant $c > 0$. Hence the empirical risk is dominated term compared to

676 $\mathbb{E}_{x \sim \pi} [F_{\mathcal{P}}(x) - F_{\mathcal{D}}(x)]$, and we complete the proof. \square

677 **D.3 Implementation**

678 We rewrite them below: Let $\widehat{F}(x) := F(x) + r(x)$ where $r(x)$ is some regularizer, and $F = \mathbb{E}_{i \in I} f_i$ is the expectation of a family of G -Lipschitz functions.

Algorithm 3 AlternateSample, [34]

- 1: **Input:** Function \widehat{F} , initial point $x_0 \sim \pi_0$, step size η
 - 2: **for** $t \in [T]$ **do**
 - 3: $y_t \leftarrow x_{t-1} + \sqrt{\eta}\zeta$ where $\zeta \sim \mathcal{N}(0, I_d)$
 - 4: Sample $x_t \leftarrow \exp(-\widehat{F}(x) - \frac{1}{2\eta}\|x - y_t\|_2^2)$
 - 5: **end for**
 - 6: **Output:** x_T
-

679

680 **Theorem D.1** (Guarantee of Algorithm 3, [14]). *Let $\mathcal{K} \subset \mathbb{R}^d$ be a convex set of diameter D , and*
 681 *$\widehat{F} : \mathcal{K} \rightarrow \mathbb{R}$, and $\pi \propto \exp(-\widehat{F})$ satisfies LSI with constant C_{LSI} . Then set $\eta \geq 0$, we have*

$$R_q(\pi_t, \pi) \leq \frac{R_q(\pi_0, \pi)}{(1 + \eta/C_{\text{LSI}})^{2t/q}},$$

682 where $R_q(\pi', \pi)$ is the q -th order of Renyi divergence between π' and π .

683 To get a sample from $\exp(-\widehat{F}(x) - \frac{1}{2\eta}\|x - y_t\|_2^2)$, we use the rejection sampler from [25], whose
 684 guarantee is stated below:

685 **Lemma D.2** (Rejection Sampler, [25]). *If the step size $\eta \lesssim G^{-2} \log^{-1}(1/\delta_{\text{inner}})$ and the inner*
 686 *accuracy $\delta_{\text{inner}} \in (0, 1/2)$, there is an algorithm that can return a random point x that has δ_{inner}*
 687 *total variation distance to the distribution proportional to $\exp(-\widehat{F}(x) - \frac{1}{2\eta}\|x - y\|_2^2)$. Moreover, the*
 688 *algorithm accesses $O(1)$ different f_i function values and $O(1)$ samples from the density proportional*
 689 *to $\exp(-r(x) - \frac{1}{2\eta}\|x - y\|_2^2)$.*

690 Combining Theorem 5.6, Theorem D.1 and Lemma D.2, we can get the following implementation of
 691 the exponential mechanism for non-smooth functions.

692 **Theorem 5.7** (Implementation, risk bound). *For $\varepsilon, \delta \in (0, 1/2)$, there is an $(\varepsilon, 2\delta)$ -DP efficient*
 693 *sampler that can achieve the empirical and population risks $O(GD \frac{d \cdot \log \log(n) \sqrt{\log(1/\delta)}}{\varepsilon \log(nd)})$. Moreover,*
 694 *in expectation, the sampler takes $\tilde{O}\left(n\varepsilon^3 \log^3(d) \sqrt{\log(1/\delta)} / (GD)\right)$ function values query and some*
 695 *Gaussian random variables restricted to the convex set \mathcal{K} in total.*

696 *Proof.* By Theorem 5.6, it suffices to get a good sample from π with density proportional to
 697 $\exp(-\beta(F_{\mathcal{D}}(x) + \frac{\mu}{2}\|x\|_2^2))$ where $\beta = O(\frac{\varepsilon \log(nd)}{GD \sqrt{\log(1/\delta)}})$, $\mu = \frac{d}{D^2 \beta}$. Set $q = 1$, which gives
 698 that $R_q(\cdot, \cdot)$ is the KL-divergence. Suppose we let x_0 is drawn from density proportional to
 699 $\exp(-\frac{\beta}{2}\mu\|x\|_2^2)$, then the KL divergence between π_0 and π is bounded by $\exp(q\beta GD)$.

Now let $\pi_T^{(i)}$ be the distribution we get over x_T from Algorithm 3 if we use an exact sampler for
 i iterations, then the sampler of Lemma D.2 for the remaining $T - i$ iterations. The output of
 Algorithm 3 that we actually get is $\pi_T^{(0)}$. Note that $C_{\text{LSI}} \leq D^2 n$, and $\eta \lesssim \beta^{-2} G^{-2} \log^{-1}(2T/\delta)$.
 Setting

$$T = O\left(\frac{C_{\text{LSI}}}{\eta} \log(\exp(q\beta GD)/\delta^2)\right) = \tilde{O}\left(\frac{n\varepsilon^3 \log^3(d) \sqrt{\log(1/\delta)}}{GD}\right)$$

700 we get $\delta_{\text{inner}} = \delta/2T$ in Lemma D.2 and that $R_1(\pi_T^{(T)}, \pi) \leq \delta^2/8$. This implies the total variation
 701 distance between $\pi_T^{(T)}$ and π is at most $\delta/2$ by Pinsker's inequality. Furthermore, by the post-
 702 processing inequality, the total variation distance between $\pi_T^{(i)}$ and $\pi_T^{(i+1)}$ is at most $\delta/2T$ for all i .
 703 Then by triangle inequality the total variation distance between $\pi_T^{(0)}$ and π is at most δ . \square

704 **D.4 Proof of Theorem 5.8**

705 **Theorem 5.8.** *There exists an ε -DP differentially private algorithm that achieves a population risk*
 706 *of $O\left(GD\left(\frac{d\log(\varepsilon n/d)}{\varepsilon n} + \sqrt{\frac{d\log(\varepsilon n/d)}{n}}\right)\right)$.*

707 *Proof.* We pick a maximal packing P of $O((D/r)^d)$ points, such that every point in \mathcal{K} is distance at
 708 most r from some point in P . By G -Lipschitzness, the risk of any point in P for the DP-ERM/SCO
 709 problems over \mathcal{K} are at most Gr plus the risk of the same point for DP-ERM/SCO over P . The
 710 exponential mechanism over P gives a DP-ERM risk bound of $O\left(\frac{GD}{\varepsilon n} \log |P|\right)$. Next, note that
 711 the empirical loss of each point in P is the average of n random variables in $[0, GD]$ wlog. So,
 712 the expected maximum difference between the empirical and population loss of any point in P is
 713 $O\left(\frac{GD\sqrt{\log |P|}}{\sqrt{n}}\right)$. Putting it all together we get a DP-SCO expected risk bound of:

$$O\left(Gr + GD\left(\frac{d\log(D/r)}{\varepsilon n} + \sqrt{\frac{d\log(D/r)}{n}}\right)\right).$$

714 This is approximately minimized by setting $r = Dd/\varepsilon n$. This gives a bound of:

$$O\left(GD\left(\frac{d\log(\varepsilon n/d)}{\varepsilon n} + \sqrt{\frac{d\log(\varepsilon n/d)}{n}}\right)\right).$$

715 □

716 **E Conclusion**

717 We present a novel framework that can improve upon the state-of-the-art rates for locating second-
 718 order stationary points for both empirical and population risks. We also examine the utilization of the
 719 exponential mechanism to attain favorable excess risk bounds for both a polynomial time sampling
 720 approach and an exponential time sampling approach. Despite the progress made, several interesting
 721 questions remain. There is still a gap between the upper and lower bounds for finding stationary
 722 points. As noted in [2], it is quite challenging to beat the current $(\frac{\sqrt{d}}{n})^{2/3}$ empirical upper bound, and
 723 overcoming this challenge may require the development of new techniques. A potential avenue for
 724 improving the population rate for SOSp could be combining our drift-controlled framework with the
 725 tree-based private SpiderBoost algorithm in [2]. Additionally, it is worth exploring if it is possible to
 726 achieve better excess risk bounds within polynomial time, and what the optimal risk bound could be.

727 **F Extended related work**

728 In the convex setting, it is feasible to achieve efficient algorithms with good risk guarantees. In turn,
 729 differentially private empirical risk minimization (DP-ERM) [12, 13, 16, 27, 32, 9, 42, 40, 41] and
 730 differentially private stochastic optimization [4, 7, 6, 21, 33, 3, 31, 25, 22, 11, 26] have been two of
 731 the most extensively studied problems in the DP literature. Most common approaches are variants of
 732 DP-SGD [13] or the exponential mechanism [37].

733 As for the non-convex optimization, due to the intrinsic challenges in minimizing general non-convex
 734 functions, most of the previous works [48, 49, 46, 45, 55, 41, 43, 53, 2, 50, 23] adopted the gradient
 735 norm as the accuracy metric rather than risk. Instead of minimizing the gradient norm discussed
 736 before, [8] tried to minimize the stationarity gap of the population function privately, which is defined
 737 as $\text{Gap}_{F_{\mathcal{P}}}(x) := \max_{y \in \mathcal{K}} \langle \nabla F_{\mathcal{P}}(x), x - y \rangle$, which requires \mathcal{K} to be a bounded domain. There are
 738 also some different definitions of the second order stationary point. We refer the readers to [36] for
 739 more details.

740 The risk bound achieved by algorithms with polynomial running time is weak and requires $n \gg d$
 741 to be meaningful. Many previous works consider minimizing risks of non-convex functions under
 742 stronger assumptions, such as, Polyak-Lojasiewicz condition [48, 54], Generalized linear model
 743 (GLM) [45] and weakly convex functions [8].

744 In the (non-private) classic stochastic optimization, there is a long line of influential works on finding
 745 the first and second-order stationary points for non-convex functions, [1, 28, 20, 52, 17].

746 **First order stationary points.** Progress towards privately finding a first-order stationary point is
 747 measured in (i) the norm of the empirical gradient at the solution x , i.e., $\|\nabla F_{\mathcal{D}}(x)\|$, and (ii) the
 748 norm of the population gradient, i.e., $\|\nabla F_{\mathcal{P}}(x)\|$. We summarize compare these first-order guarantees
 749 achieved by Algorithm 1 with previous algorithms in Table 2:

References	Empirical	Population
[48]	$\frac{d^{1/4}}{\sqrt{n}}$	N/A
[46]	$\frac{d^{1/4}}{\sqrt{n}}$	$\frac{\sqrt{d}}{\sqrt{n}}$
[49]	$(\frac{\sqrt{d}}{n})^{2/3}$	N/A
[55]	$\frac{d^{1/4}}{\sqrt{n}}$	$\frac{d^{1/4}}{\sqrt{n}}$
[43]	$\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n})^{2/3}$	N/A
[2]	$(\frac{\sqrt{d}}{n})^{2/3}$	$\frac{1}{n^{1/3}} + (\frac{\sqrt{d}}{n})^{1/2}$

Table 2: Previous work in finding first-order stationary points. We omit logarithmic terms and dependencies on other parameters such as Lipschitz constant. “N/A” means we do not find an explicit result in the work.

750 **Second order stationary points.** We say a point x is a Second-Order Stationary Point (SOSP),
 751 or a local minimum of a twice differentiable function g if $\|\nabla g(x)\|_2 = 0$ and $\text{smin}(\nabla^2 g(x)) \geq 0$.
 752 Exact second-order stationary points can be extremely challenging to find [24]. Instead, it is common
 753 to measure the progress in terms of how well the solution approximates an SOSP.

754 **Definition F.1** (approximate-SOSP, [1]). We say $x \in \mathbb{R}^d$ is an α -second order stationary point
 755 (α -SOSP) for ρ -Hessian Lipschitz function g , if

$$\|\nabla g(x)\|_2 \leq \alpha \bigwedge \text{smin}(\nabla^2 g(x)) \geq -\sqrt{\rho\alpha}.$$

References	Empirical	Population
[45]	$\frac{d^{1/4}}{\sqrt{n}}$	N/A
[47]	$(\frac{d}{n})^{4/7}$	N/A
[23]	$(\frac{d}{n})^{1/2}$	N/A
Ours	$(\frac{\sqrt{d}}{n})^{2/3}$	$\frac{1}{n^{1/3}} + (\frac{\sqrt{d}}{n})^{3/7}$

Table 3: Summary of previous results in finding α -SOSP, where α is demonstrated in the Table. Omit the logarithmic terms and the dependencies on other parameters like Lipschitz constant. “N/A” means we do not find an explicit result in the work.

756 Existing works in finding approximate SOSP privately give guarantees for the empirical function $F_{\mathcal{D}}$.
 757 We improve upon the state-of-the-art result and give the first guarantee for the population function
 758 $F_{\mathcal{P}}$, which is summarized in Table 3.