

Dynamic Privacy Budget Allocation Improves Data Efficiency of Differentially Private Gradient Descent

Junyuan Hong
hongju12@msu.edu
Michigan State University
East Lansing, MI, USA

Zhangyang Wang
atlaswang@utexas.edu
University of Texas at Austin
Austin, TX, USA

Jiayu Zhou
jiayuz@msu.edu
Michigan State University
East Lansing, MI, USA

ABSTRACT

Protecting privacy in learning while maintaining the model performance has become increasingly critical in many applications that involve sensitive data. A popular private learning framework is differentially private learning composed of many privatized gradient iterations by noising and clipping. Under the privacy constraint, it has been shown that the dynamic policies could improve the final iterate loss, namely the quality of published models. In this talk, we will introduce these dynamic techniques for learning rate, batch size, noise magnitude and gradient clipping. Also, we discuss how the dynamic policy could change the convergence bounds which further provides insight of the impact of dynamic methods.

CCS CONCEPTS

• Security and privacy → Privacy protections.

KEYWORDS

machine learning, privacy

ACM Reference Format:

Junyuan Hong, Zhangyang Wang, and Jiayu Zhou. 2022. Dynamic Privacy Budget Allocation Improves Data Efficiency of Differentially Private Gradient Descent. In *2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*, June 21–24, 2022, Seoul, Republic of Korea. ACM, New York, NY, USA, 25 pages. <https://doi.org/10.1145/3531146.3533070>

1 INTRODUCTION

In the era of big data, privacy protection in machine learning systems is becoming a crucial topic as increasing personal data involved in training models [9] and the presence of malicious attackers [12, 31]. In response to the growing demand, differential-private (DP) machine learning [10] provides a computational framework for privacy protection and has been widely studied in various settings, including both convex and non-convex optimization [15, 33, 34].

One widely used procedure for privacy-preserving learning is the (Differentially) Private Gradient Descent (PGD) [1, 3]. A typical gradient descent procedure updates its model by gradients of the loss evaluated on a training dataset. When the data is sensitive, the gradients should be *privatized* to prevent excess privacy

leakage. The PGD privatizes a gradient by adding controlled noise. As such, the models from PGD is expected to have a lower utility as compared to those from unprotected algorithms. In the cases where strict privacy control is exercised, or equivalently, a tight *privacy budget*, accumulating effects from highly-noised gradients may lead to unacceptable model performance. It is thus critical to design effective privatization procedures for PGD to maintain a great balance between utility and privacy.

Recent years witnessed a promising direction of privatization that *dynamically allocate a privacy budget* for each iteration to boost utility, under the constraint of a specified total privacy budget. One example is [19], which reduces the budget-bonded noise magnitude when the loss does not decrease, due to the observation that gradients become very small when approaching convergence, and a static noise scale will overwhelm these gradients. Another example is [40], which periodically decreases the magnitude following a pre-defined strategy, e.g., exponential decaying or step decaying. Both approaches confirmed the empirically advantages of decreasing noise magnitudes. Intuitively, the dynamic mechanism may coordinate with certain properties of the learning task, e.g., training data and loss surface. Following the work, improved allocation policies are proposed, e.g., policies transferred from auxiliary datasets [13], policies with distributed budgets [6], and a combination with adaptive learning rate [38]. Yet there is little theoretical analysis available and two important questions remain unanswered: 1) *What is the form of utility-preferred budget (or noise equivalently) schedules?* 2) *When and to what extent such an allocation policy improves utility?*

Though there are theoretical studies of static-allocation policies, e.g., [34], the data efficiency is not the focus as discussions usually assume an unlimited amount of data is available. However, we argue that the data efficiency with limited data size is critical in practice, especially when DP makes the learning more data-hungry [21]. One example is federated learning [20, 21], a distributed learning framework that aggregates many local models to form a stronger global model, where each model is privately trained on a local client, typically with very limited private data. Another example is biomedical applications, where collecting samples involves expansive clinical trials or cohort studies, resulting the scarcity of training set. To study biomarkers of Alzheimer's, NIH has funded Alzheimer's Disease Neuroimaging Initiative for \$40 million, which collected imaging and genetic biomarkers from only 800 patients after 5 years' extensive and collaborative efforts [35]. Therefore, we believe data efficiency needs to be taken into account in developing a private learning algorithm.

To answer these questions, in this paper we develop a principled approach to construct dynamic schedules and quantify their utility bounds in different learning algorithms. Our contributions are

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

FAccT '22, June 21–24, 2022, Seoul, Republic of Korea

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9352-2/22/06...\$15.00

<https://doi.org/10.1145/3531146.3533070>

summarized as follows. 1) For the class of loss functions satisfying the Polyak-Lojasiewicz condition [24], we show that dynamic schedules, that improve the utility upper bound with high data-efficiency, are shaped by the changing influence of per-iteration noise on the final loss. As the influence is tightly connected to the loss curvature, the advantage of using dynamic schedules therefore depends on the loss function. 2) Beyond vanilla gradient descent, our results show the gradient methods with momentum implicitly introduce a dynamic schedule and result in a non-monotonous influence trend. 3) We also show that our results are generalizable to population bounds in high probability or based on uniform stability theorems. Though our major focus is on the theoretic study, we empirically validate the results on a non-convex loss function formulated by a neural network. The empirical results suggest that a preferred dynamic schedule admits the exponentially decaying form, and works better when learning with high-curvature loss functions. Moreover, dynamic schedules give higher utility under stricter privacy conditions (e.g., smaller sample size and less privacy budget).

2 RELATED WORK

Differentially Private Learning. Differential privacy (DP) characterizes the chance that an algorithm output (e.g., a learned model) leaks private information of its training data when the output distribution is known. Since outputs of many learning algorithms have undetermined distributions, the probabilistic risk is hard to measure. A common approach to tackle this issue is to inject randomness with known distribution to *privatize* the learning procedures. Classical methods include output perturbation [5], objective perturbation [5] and gradient perturbation [1, 3, 36]. Among these approaches, the Private Gradient Descent (PGD) has attracted extensive attention in recent years because it can be flexibly integrated with variants of gradient-based iteration methods, e.g., stochastic gradient descent, momentum methods [26], and Adam [17], for (non-)convex problems.

Dynamic Policies for Privacy Protection. Wang et al. [34] studied the empirical risk minimization using dynamic variation reduction of perturbed gradients. They showed that the utility upper bound can be achieved by gradient methods under uniform noise parameters. Instead of enhancing the gradients, Lee and Kifer [19], Yu et al. [40] showed the benefits of using a dynamic schedule of privacy parameters or equivalently noise scales. Following [19], a series of work [6, 14, 37, 41] adaptively allocate privacy budget towards better privacy-utility trade-off. Moreover, adaptive sensitivity control [23, 32] and dynamic batch sizes [11] are also shown to improve convergence.

Utility Upper Bounds. Utility upper bounds are a critical metric for privacy schedules, which characterizes the maximum utility that a schedule can deliver in theory. Wang et al. [34] is the first to prove the utility bound under the PL condition. Recently, Zhou et al. proved the utility bound by using the momentum of gradients [17, 25]. In this paper, we improve the upper bound by a more accurate estimation of the dynamic influence of step noise. We show that introducing a dynamic schedule further boosts the sample-efficiency of the upper bound. Table 1 summarizes the upper bounds of a selection of state-of-the-art algorithms based on

private gradients (up block, see Appendix B for the full list), and methods studied in this paper (down block), showing the benefits of dynamic influence.

Especially, a closely-related work by Feldman *et al.* achieved a convergence rate similar to ours in terms of generalization error bounds (c.f. SSGD in Table 2), by dynamically adjusting batch sizes [11]. However, the approach requires controllable batch sizes, which may not be feasible in many applications. In federated learning, for example, where users update models locally and then pass the parameters to server for aggregation, the server has no control over batch sizes, and coordinating users to use varying batch sizes may not be realistic. On the other hand, our proposed method can still be applied for enhancing utility, as the server can dynamically allocate privacy budget for each round when the presence of a user in the global aggregation is privatized [21].

3 PRIVATE GRADIENT DESCENT

Notations. We consider a learning task by empirical risk minimization (ERM) $f(\theta) = \frac{1}{N} \sum_{n=1}^N f(\theta; x_n)$ on a private dataset $\{x_n\}_{n=1}^N$ and $\theta \in \mathbb{R}^D$. The gradient methods are defined as $\theta_{t+1} = \theta_t - \eta_t \nabla_t$, where $\nabla_t = \nabla f(\theta_t) = \frac{1}{N} \sum_n \nabla f(\theta_t; x_n)$ denotes the non-private gradient at iteration t , η_t is the step learning rate. $\nabla_t^{(n)} = \nabla f(\theta_t; x_n)$ denotes the gradient on a sample x_n . $\mathbb{1}_c$ denotes the indicator function that returns 1 if the condition c holds, otherwise 0.

Assumptions. (1) In this paper, we assume $f(\theta)$ is continuous and differentiable. Many commonly used loss functions satisfy this assumption, e.g., the logistic function. (2) For a learning task, only finite amount of privacy cost is allowed where the maximum cost is called *privacy budget* and denoted as R . (3) Generally, we assume that loss functions $f(\theta; x)$ (sample-wise loss) are G -Lipschitz continuous and $f(\theta)$ (the empirical loss) is M -smooth.

DEFINITION 3.1 (G-LIPSCHITZ CONTINUITY). A function $f(\cdot)$ is G -Lipschitz continuous if, for $G > 0$ and all x, y in the domain of $f(\cdot)$, $f(\cdot)$ satisfies $\|f(y) - f(x)\| \leq G \|y - x\|$.

DEFINITION 3.2 (m-STRONGLY CONVEXITY). A function $f(\cdot)$ is m -strongly convex if $f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|^2$, for some $m > 0$ and all x, y in the domain of $f(\cdot)$.

DEFINITION 3.3 (M-SMOOTHNESS). A function is M -smooth w.r.t. l_2 norm if $f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{M}{2} \|y - x\|^2$, for some constant $M > 0$ and all x, y in the domain of $f(\cdot)$.

For a private algorithm $\mathcal{M}(d)$ which maps a dataset d to some output, the privacy cost is measured by the bound of the output difference on the adjacent datasets. *Adjacent datasets* are defined to be datasets that only differ in one sample. In this paper, we use the zero-Concentrated Differential Privacy (zCDP, see Definition 3.4) as the privacy measurement, because it provides the simplicity and possibility of adaptively composing privacy costs at each iteration. Various privacy metrics are discussed or reviewed in [8]. A notable example is Moment Accountant (MA) [1], which adopts similar principle for composing privacy costs while is less tight for a smaller privacy budget. We note that alternative metrics can be adapted to our study without major impacts to the analysis.

Table 1: Comparison of utility upper bound using different privacy schedules. The algorithms are T -iteration $\frac{1}{2}R$ -zCDP under the PL condition (unless marked with * for convexity). The O notation in this table drops other ln terms. Unless otherwise specified, all algorithms (including non-private GD) terminate at step $T = O(\ln \frac{N^2 R}{D})$. Assume loss functions are 1-smooth and 1-Lipschitz continuous, and all parameters satisfy their numeric assumptions. Key notations: O_p – bound occurs in probability p ; D – feature dimension; N – sample size; R – privacy budget where $R_{\epsilon, \delta}$ is the equivalent budget accounted by (ϵ, δ) -DP; c_i – constant; other notations can be found in Section 4. We provide citation after algorithm names for concise reference here and defer detailed explanations of baseline algorithms to Appendix B. Extensions to generalization error bounds are given in Table 2.

Algorithm	Schedule (σ_t^2)	Utility Upper Bound
*GD+Adv [3]	$O\left(\frac{\ln(N/\delta)}{R_{\epsilon, \delta}}\right)$	$O\left(\frac{D \ln^3 N}{NR_{\epsilon, \delta}}\right)$
GD+MA [34]	$O\left(\frac{T}{R_{\epsilon, \delta}}\right)$	$O\left(\frac{D \ln^2 N}{N^2 R_{\epsilon, \delta}}\right)$
GD+MA (adjusted utility) [39]	$O\left(\frac{T}{R_{\epsilon, \delta}}\right)$	$O\left(\min \frac{\sqrt{D}}{NR_{\epsilon, \delta}}, \frac{D \ln N}{N^2 R_{\epsilon, \delta}^2}\right)$
*GD+Adv+BBImp [7]	$O\left(\frac{n^2 \ln(n/\delta)}{R_{\epsilon, \delta}}\right)$	$O_p\left(\frac{D^2 \ln^2(1/p)}{R_{\epsilon, \delta} N^{1-c}}\right)$
Adam+MA [42]	$O\left(\frac{T}{R_{\epsilon, \delta}}\right)$	$O_p\left(\frac{\sqrt{D} \ln(ND\epsilon/(1-p))}{NR_{\epsilon, \delta}}\right)$
GD, Non-Private	0	$O\left(\frac{D}{N^2 R}\right)$
GD+zCDP, Static Schedule	$\frac{T}{R}$	$O\left(\frac{D \ln N}{N^2 R}\right)$
GD+zCDP, Dynamic Schedule	$O\left(\frac{\gamma^{(t-T)/2}}{R}\right)$	$O\left(\frac{D}{N^2 R}\right)$
Momentum+zCDP, Static Schedule	$\frac{T}{R}$	$O\left(\frac{D}{N^2 R}(c + \ln N \mathbb{I}_{T > \hat{T}})\right)$
Momentum+zCDP, Dynamic Schedule	$O\left(\frac{c_1 \gamma^{T+t} + c_2 \gamma^{(T-t)/2}}{R}\right)$	$O\left(\frac{D}{N^2 R}(1 + \frac{cD}{N^2 R} \mathbb{I}_{T > \hat{T}})\right)$

DEFINITION 3.4 (ρ -zCDP [4]). *Let $\rho > 0$. A randomized algorithm $M : \mathcal{D}^n \rightarrow \mathbb{R}$ satisfies ρ -zCDP if, for all adjacent datasets $d, d' \in \mathcal{D}^n$, $D_\alpha(M(d) \| M(d')) \leq \rho\alpha$, $\forall \alpha \in (1, \infty)$ where $D_\alpha(\cdot \| \cdot)$ denotes the Rényi divergence [28] of order α .*

The zCDP provides a linear composition of privacy costs of sub-route algorithms. When the input vector is privatized by injecting Gaussian noise of $\mathcal{N}(0, \sigma_t^2 I)$ for the t -th iteration, the composed privacy cost is proportional to $\sum_t \rho_t$ where the step cost is $\rho_t = \frac{1}{\sigma_t^2}$. For simplicity, we absorb the constant coefficient into the (residual) *privacy budget* R . The formal theorems for the privacy cost computation of composition and Gaussian noising is included in Lemmas C.1 and C.2.

Algorithm 1 Privatizing Gradients

Input: Raw gradients $[\nabla_t^{(1)}, \dots, \nabla_t^{(n)}]$ ($n = N$ by default), v_t , residual privacy budget R_t assuming the full budget is R and $R_1 = R$.

- 1: $\rho_t \leftarrow 1/\sigma_t^2$, $\nabla_t \leftarrow \frac{1}{n} \sum_{i=1}^n \nabla_t^{(i)}$ (Budget request)
- 2: **if** $\rho_t < R_t$ **then**
- 3: $R_{t+1} \leftarrow R_t - \rho_t$
- 4: $g_t \leftarrow \nabla_t + G\sigma_t v_t/N$, $v_t \sim \mathcal{N}(0, I)$ (Privacy noise)
- 5: $m_{t+1} \leftarrow \phi(m_t, g_t)$ or g_t if $t = 1$
- 6: **return** $\eta_t m_{t+1}$, R_{t+1} (Utility projection)
- 7: **else**
- 8: Terminate

Generally, we define the Private Gradient Descent (PGD) method as iterations for $t = 1 \dots T$:

$$\theta_{t+1} = \theta_t - \eta_t \phi_t = \theta_t - \eta_t (\nabla_t + \sigma_t G v_t / N), \quad (1)$$

where $\phi_t = g_t$ is the gradient privatized from ∇_t as shown in Algorithm 1, G/N is the bound of sensitivity of the gathered gradient excluding one sample gradient, and $v_t \sim \mathcal{N}(0, I)$ is a vector elementwisely subject to Gaussian distribution. We use σ_t to denote the noise scale at step t and use σ to collectively represents the schedule $(\sigma_1, \dots, \sigma_T)$ if not confusing. When the Lipschitz constant is unknown, we can control the upper bound by scaling the gradient if it is over some constant. The scaling operation is often called *clipping* in literatures since it clips the gradient norm at a threshold. After the gradient is noised, we apply a modification, $\phi(\cdot)$, to enhance its utility. In this paper, we consider two types of $\phi(\cdot)$:

$$\phi(m_t, g_t) = g_t \text{ (GD)}, \quad (2)$$

$$\phi(m_t, g_t) = [\beta(1 - \beta^{t-1})m_t + (1 - \beta)g_t] / (1 - \beta^t) \text{ (Momentum)} \quad (3)$$

We now show that the PGD using Algorithm 1 guarantees a privacy cost less than R :

THEOREM 3.1. *Suppose $f(\theta; x)$ is G -Lipschitz continuous and the PGD algorithm with privatized gradients defined by Algorithm 1, stops at step T . The PGD algorithm outputs θ_T and satisfies ρ -zCDP where $\rho \leq \frac{1}{2}R$.*

Note that Theorem 3.1 allows σ_t to be different throughout iterations. Next we present a principled approach for deriving dynamic schedules optimized for the final loss $f(\theta_T)$.

4 DYNAMIC POLICIES BY MINIMIZING UTILITY UPPER BOUNDS

To characterize the utility of the PGD, we adopt the Expected Excess Risk (EER), which notion is widely used for analyzing the convergence of random algorithms, e.g., [3, 34]. Due to the presence of the noise and the limitation of learning iterations, optimization using private gradients is expected to reach a point with a higher loss (i.e., excess risk) as compared to the optimal solution without private protection. Define $\theta^* = \arg \min_{\theta} f(\theta)$, after [Algorithm 1](#) is iterated for T times in total, the EER gives the expected utility degradation:

$$\text{EER} = \mathbb{E}_v[f(\theta_{T+1})] - f(\theta^*).$$

Due to the variety of loss function and complexity of recursive iterations, an exact EER with noise is intractable for most functions. Instead, we study the worst case scenario, i.e., the upper bound of the EER, and our goal is to minimize the upper bound. For consistency, we call the upper bound of EER divided by the initial error as ERUB. Since the analytical form of EER is either intractable or complicated due to the recursive iterations of noise, studying the ERUB is a convenient and tractable alternative. The upper bound often has convenient functional forms which are (1) sufficiently simple, such that we can directly minimize it, and (2) closely related to the landscape of the objective depending on both the training dataset and the loss function. As a consequence, it is also used in previous PGD literature [23, 34] for choosing proper parameters. Moreover, we let ERUB_{\min} be the achievable optimal upper bound by a specific choice of parameters, e.g., the σ and T .

As the EER is iteratively determined by [Eq. \(1\)](#), we define the influence of the dynamics in noise magnitude σ_t as the derivative: $q_t^* = \frac{\partial \text{EER}}{\partial \sigma_t}$. Accordingly, we can approximate the EER shift as $q_t^* \Delta \sigma_t$ when σ_t increases by $\Delta \sigma_t$. However, because the EER is strongly data-dependent, the derived q_t^* on a given dataset may not generalize to another dataset. Instead, we consider a more general term based on ERUB, i.e., $q_t = \frac{\partial \text{ERUB}}{\partial \sigma_t}$.

In this paper, we consider the class of loss functions satisfying the Polyak-Lojasiewicz (PL) condition which bounds losses by corresponding gradient norms. It is more general than the m -strongly convexity. If f is differentiable and M -smooth, then m -strongly convexity implies the PL condition.

DEFINITION 4.1 (POLYAK-LOJASIEWICZ CONDITION [24]). For $f(\theta)$, there exists $\mu > 0$ and for every θ , $\|\nabla f(\theta)\|^2 \geq 2\mu(f(\theta) - f(\theta^*))$.

The PL condition helps us to reveal how the influence of step noise propagates to the final excess error, i.e., EER. Though the assumption was also used previously in Wang et al. [34], Zhou et al. [42], neither did they discuss the propagated influence of noise. In the following sections, we will show how the influence can tighten the upper bound in gradient descent and its momentum variant.

4.1 Gradient Descent Methods and Noise Influences

For the brevity of variables, we first define the following summarized constants:

$$\text{non-private ERUB} : \alpha \triangleq \frac{DG^2}{2RMN^2(f(\theta_1) - f(\theta^*))} \leq O\left(\frac{DG^2}{RMN^2}\right), \quad (4)$$

$$\text{curvature} : \kappa \triangleq \frac{M}{\mu}, \quad (5)$$

$$\text{non-private convergence rate} : \gamma \triangleq 1 - \frac{1}{\kappa}, \quad (6)$$

which satisfy $\kappa \geq 1$ and $\gamma \in [0, 1)$. Here, α is upper bounded by non-private ERUB within $T = \left\lceil O\left(\ln \frac{N^2 R(M-\mu)}{DG^2}\right) \right\rceil$ iterations. Therefore, α provide a simple reference of an ideal convergence bound, reaching which indicates a superior performance with privacy guarantee. κ characterizes the curvature of $f(\cdot)$ which is the condition number of $f(\cdot)$ if $f(\cdot)$ is strongly convex, and γ is the convergence rate for non-private SGD (c.f. [Theorem 4.1](#) with $\sigma_t = 0$). κ tends to be large if the function is sensitive to small differences in inputs, and $1/\alpha$ tends to be large if more samples are provided and with a less strict privacy budget. The convergence of PGD under the PL condition has been studied for private [34] and non-private [16, 22, 27] ERM. Below we extend the bound in [34] by considering dynamic influence of noise and relax σ_t to be dynamic:

THEOREM 4.1. Let α , κ and γ be defined in [Eq. \(6\)](#), and $\eta_t = \frac{1}{M}$. Suppose $f(\theta; x_t)$ is G -Lipschitz and $f(\theta)$ is M -smooth satisfying the Polyak-Lojasiewicz condition. For PGD, the following holds:

$$\text{ERUB} = \gamma^T + R \sum_{t=1}^T q_t \sigma_t^2, \text{ where } q_t \triangleq \gamma^{T-t} \alpha. \quad (7)$$

[Theorem 4.1](#) degenerates to a non-private variant as no noise is applied, i.e., $\sigma_t = 0$ for all t . In [Eq. \(7\)](#), the step noise magnitude σ_t^2 has an exponential influence, q_t , on the EER. Note we ignore the constant factor R in the influence. The [Eq. \(7\)](#) implies that the influence of noise at step t increase quickly by an exponential rate. Importantly, the increasing rate is the same as the convergence rate, i.e., the first term in [Eq. \(7\)](#). The dynamic characteristic of the influence is the key to prove a tighter bound. Plus, on the presence of the dynamic influence, it is natural to choose a dynamic σ_t^2 . When relaxing q_t to a static 1, a static σ_t^2 was studied by Wang et al. They proved a bound which is nearly optimal except a $\ln^2 N$ factor. To get the optimal bound, in the following sections, we look for the σ and T that minimize the upper bound.

4.1.1 Uniform Schedule. The uniform setting of σ_t has been previously studied in Wang et al. [34]. Here, we show that the bound can be further tightened by considering the dynamic influence of iterations and a proper T .

THEOREM 4.2. Suppose conditions in [Theorem 4.1](#) are satisfied. When $\sigma_t^2 = T/R$, let α , γ and κ be defined in [Eq. \(6\)](#) and let T be: $T = \left\lceil O\left(\kappa \ln\left(1 + \frac{1}{\kappa\alpha}\right)\right) \right\rceil$. Meanwhile, if $\kappa \geq \frac{1}{1-c} > 1$, $1/\alpha > 1/\alpha_0$ for some constant $c \in (0, 1)$ and $\alpha_0 > 0$, the corresponding bound is:

$$\text{ERUB}_{\min}^{\text{uniform}} = \Theta\left(\frac{\kappa^2}{\kappa + 1/\alpha} \ln\left(1 + \frac{1}{\kappa\alpha}\right)\right). \quad (8)$$

SKETCH OF PROOF. The key of proof is to find a proper T to minimize

$$\begin{aligned} \text{ERUB} = E &= \gamma^T + \sum_{t=1}^T \gamma^{T-t} \alpha R \sigma^2 \\ &= \gamma^T + \alpha T \frac{1 - \gamma^T}{1 - \gamma} = \gamma^T + \alpha \kappa (1 - \gamma^T) T \end{aligned} \quad (9)$$

where we use $\sigma_t = \sqrt{T/R}$. Vanishing its gradient is to solve $\gamma^T \ln \gamma + \alpha \kappa (1 - \gamma^T) - \alpha \kappa T \gamma^T \ln \gamma = 0$, which however is intractable. In [34], T is chosen to be $\mathcal{O}(\ln(1/\alpha))$ and ERUB is relaxed as $\gamma^T + \alpha \kappa T^2$. The approximation results in a less tight bound as $\mathcal{O}(\alpha(1 + \kappa \ln^2(1/\alpha)))$ which explodes as $\kappa \rightarrow \infty$.

We observe that for a super sharp loss function, i.e., a large κ , any minor perturbation may result in tremendously fluctuating loss values. In this case, not-stepping-forward will be a good choice. Thus, we choose $T = \frac{1}{\ln(1/\gamma)} \ln\left(1 + \frac{\ln(1/\gamma)}{\alpha}\right) \leq \mathcal{O}\left(\kappa \ln\left(1 + \frac{1}{\kappa \alpha}\right)\right)$ which converges to 0 as $\kappa \rightarrow +\infty$. The full proof is deferred to the appendix. \square

4.1.2 Dynamic Schedule. A dynamic schedule can improve the upper bound delivered by the uniform schedule. First, we observe that the excess risk in Eq. (7) is upper bounded by two terms: the first term characterizes the error due to the finite iterations of gradient descents; the second term, a weighted sum, comes from error propagated from noise at each iteration. Now we show for any $\{q_t | q_t > 0, t = 1, \dots, T\}$ (not limited to the q_t defined in Eq. (7)), there is a unique σ_t minimizing the weighted sum:

LEMMA 4.3 (DYNAMIC SCHEDULE). *Suppose σ_t satisfy $\sum_{t=1}^T \sigma_t^{-2} = R$. Given a positive sequence $\{q_t\}$, the following equation holds:*

$$\min_{\sigma} R \sum_{t=1}^T q_t \sigma_t^2 = \left(\sum_{t=1}^T \sqrt{q_t} \right)^2, \text{ when } \sigma_t^2 = \frac{1}{R} \sum_{i=1}^T \sqrt{\frac{q_i}{q_t}}. \quad (10)$$

Remarkably, the difference between the minimum and $T \sum_{t=1}^T q_t$ (uniform σ_t) monotonically increases by the variance of $\sqrt{q_t}$ w.r.t. t .

We see that the dynamics in σ_t come from the non-uniform nature of the weight q_t . Since q_t presents the impact of the σ_t on the final error, we denote it as *influence*. Given the dynamic schedule in Eq. (10), it is of our interest to which extent the ERUB can be improved. First, we present Theorem 4.4 to show the optimal T and ERUB.

THEOREM 4.4. *Suppose conditions in Theorem 4.1 are satisfied. Let α, κ and γ be defined in Eq. (6). When $\eta_t = \frac{1}{M}$, σ_t (based on Eqs. (7) and (10)) and the T minimizing ERUB are, i.e., $\sigma_t^2 = \frac{1}{R} \frac{\sqrt{(1/\gamma)^T - 1}}{1 - \sqrt{\gamma}} \sqrt{\gamma^t}$, $T = \left\lceil \left(2\kappa \ln\left(1 + \frac{1}{\kappa \alpha}\right) \right) \right\rceil$. Meanwhile, when $\kappa \geq 1$ and $1/\alpha \geq 1/\alpha_0$ for some positive constant α_0 , the minimal bound is:*

$$\text{ERUB}_{\min}^{\text{dynamic}} = \Theta\left(\frac{\kappa^2}{\kappa^2 + 1/\alpha}\right). \quad (11)$$

4.1.3 Discussion. In Theorems 4.2 and 4.4, we present the tightest bounds for functions satisfying the PL condition, to our best knowledge. We further analyze the advantages of our bounds from two aspects: sample efficiency and robustness to sharp losses.

Sample efficiency. Since dataset cannot be infinitely large, it is critical to know how accurate the model can be trained privately

with a limited number of samples. Formally, it is of interest to study when κ is fixed and N is large enough such that $\alpha \gg 1$. Then we have the upper bound in Eq. (8) as

$$\text{ERUB}_{\min}^{\text{uniform}} \leq \mathcal{O}\left(\kappa^2 \alpha \ln\left(\frac{1}{\kappa \alpha}\right)\right) \leq \tilde{\mathcal{O}}\left(\frac{DG^2 \ln(N)}{MN^2 R}\right), \quad (12)$$

where we ignore κ and other logarithmic constants with $\tilde{\mathcal{O}}$ as done in Wang et al. [34]. As a result, we get a bound very similar to [34], except that R is replaced by $R_{MA} = \epsilon^2 / \ln(1/\delta)$ using Moment Accountant. In comparison, based on Lemma C.3, $R = 2\rho = 2\epsilon + 4\ln(1/\delta) + 4\sqrt{\ln(1/\delta)(\epsilon + \ln(1/\delta))}$ if θ_T satisfies ρ -zCDP. Because $\ln(1/\delta) > 1$, it is easy to see $R = R_{zCDP} > R_{MA}$ when $\epsilon \leq 2\ln(1/\delta)$. As compared to the one reported in [34], our bound saved a factor of $\ln N$ and thus require less sample to achieve the same accuracy. Remarkably, the saving is due to the maintaining of the influence terms as shown in the proof of Theorem 4.2.

Using the dynamic schedule, we have $\text{ERUB}_{\min}^{\text{dynamic}} \leq \mathcal{O}(\alpha) = \mathcal{O}\left(\frac{DG^2}{MN^2 R}\right)$, which saved another $\ln N$ factor in comparison to the one using the uniform schedule Eq. (12). As shown in Table 1, such advantage maintains when comparing with other baselines and reaches the ideal non-private case, recalling the meaning of α .

Stability on ill-conditioned loss. Besides sample efficiency, we are also interested in robustness of the convergence under the presence of privacy noise. Because of the privacy noise, the private gradient descent will be unable to converge to where a non-private algorithm can reach. Specifically, when the samples are noisy or have noisy labels, the loss curvature may be sharp. The sharpness also implies lower smoothness, i.e., a small M or has a very small PL parameter. Thus, gradients may change tremendously at some steps especially in the presence of privacy noise. Such changes have more critical impact when only a less number of iterations can be executed due to the privacy constraint. Assume α is some constant while $\kappa \gg 1/\alpha$, we immediately get:

$$\text{ERUB}_{\min}^{\text{uniform}} = \Theta\left(\kappa \ln\left(1 + \frac{1}{\kappa \alpha}\right)\right) = \Theta\left(\frac{1}{\alpha}\right) \leq \mathcal{O}\left(\frac{MN^2 R}{DG^2}\right),$$

$$\text{ERUB}_{\min}^{\text{dynamic}} = \Theta(1).$$

Both are robust, but the dynamic schedule has a smaller factor since $1/\alpha$ could be a large number. In addition, the factor implies that when more samples are used, the dynamic schedule is more robust.

4.2 Gradient Descent Methods with Momentum

Section 4.1 shows that the step noise has an exponentially increasing influence on the final loss, and therefore a decreasing noise magnitude improves the utility upper bound by a $\ln N$ factor. However, the proper schedule can be hard to find when the curvature information, e.g., κ , is absent. A parameterized method that less depends on the curvature information is preferred. On the other hand, long-term iterations will result in forgetting of the initial iterations, since accumulated noise overwhelmed the propagated information from the beginning. This effect will reduce the efficiency of the recursive learning frameworks.

Alternative to GD, the momentum method can mitigate the two issues. It was originally proposed to stabilize the gradient estimation

[25]. In this section, we show that momentum (agnostic about the curvature) can flatten the dynamic influence and improve the utility upper bound. Previously, [Pichapati et al.](#) used the momentum as an estimation of gradient mean, without discussions of convergence improvements. [Zhou et al.](#) gave a bound for the Adam with DP. However, the derivation is based on gradient norm, which results in a looser bound (see [Table 1](#)).

The momentum method stabilizes gradients by moving average history coordinate values and thus greatly reduces the variance. The $\phi(m_t, g_t)$ can be rewritten as:

$$m_{t+1} = \phi(m_t, g_t) = \frac{v_{t+1}}{1 - \beta^t},$$

$$v_{t+1} = \beta v_t + (1 - \beta)g_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i, \quad v_1 = 0, \quad (13)$$

where $\beta \in [0, 1]$. Note v_{t+1} is a biased estimation of the gradient expectation while m_{t+1} is unbiased.

THEOREM 4.5 (CONVERGENCE UNDER PL CONDITION). *Suppose $f(\theta; x_i)$ is G -Lipschitz, and $f(\theta)$ is M -smooth and satisfies the Polyak-Lojasiewicz condition. Assume $\beta \neq \gamma$ and $\beta \in (0, 1)$. Let $\eta_t = \frac{\eta_0}{2M}$ and $\eta_0 \leq 8 \left(\sqrt{1 + 64\beta\gamma(\gamma - \beta)^{-2}(1 - \beta)^{-3}} + 1 \right)^{-1}$. Then the following holds:*

$$\text{EER} \leq (\gamma^T + 2R\eta_0\alpha \underbrace{U_3(\sigma, T)}_{\text{noise variance}}) (f(\theta_1) - f(\theta^*))$$

$$- \zeta \underbrace{\frac{\eta_0}{2M} \sum_{t=1}^T \gamma^{T-t} \mathbb{E} \|v_{t+1}\|^2}_{\text{momentum effect}} \quad (14)$$

where $\gamma = 1 - \frac{\eta_0}{\kappa}$, $\zeta = 1 - \frac{1}{\beta(1-\beta)^3} \eta_0^2 - \frac{1}{4} \eta_0 \geq 0$, and $U_3 = \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2$.

The upper bound includes three parts that influence the bound differently: (1) *Convergence*. The convergence term is mainly determined by η_0 and κ . η_0 should be in $(0, \kappa)$ such that the upper bound can converge. A large η_0 will be preferred to speed up convergence if it does not make the other two terms worse. (2) *Noise Variance*. The second term compressed in U_3 is the effect of the averaged noise, $\sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2$. One difference introduced by the momentum is the factor $(1 - \beta)/(1 - \beta^t)$ which is less than γ^t at the beginning and converges to a non-zero constant $1 - \beta$. Therefore, in U_3 , $\gamma^{T-t}(1 - \beta)/(1 - \beta^t)$ will be constantly less than γ^T meanwhile. Furthermore, when $t > \hat{T}$, the moving average $\sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2$ smooths the influence of each σ_t . (3) *Momentum Effect*. The momentum effect term can improve the upper bound when η_0 is small. For example, when $\beta = 0.9$ and $\gamma = 0.99$, then $\eta_0 \leq 0.98/M$ which is a rational value. Following the analysis, when M is large which means the gradient norms will significantly fluctuate, the momentum term may take the lead. Adjusting the noise scale in this case may be less useful for improving utility.

To give an insight on the effect of dynamic schedule, we provide the following utility bounds.

THEOREM 4.6 (UNIFORM SCHEDULE). *Suppose the assumptions in [Theorem 4.5](#) are satisfied. Let $\sigma_t^2 = T/R$, and let: $\hat{T} = \max t$ s.t. $\gamma^{t-1} \geq \frac{1-\beta}{1-\beta^t}$, $T = \left\lceil \mathcal{O} \left(\frac{\kappa}{\eta_0} \ln \left(1 + \frac{\eta_0}{\kappa\alpha} \right) \right) \right\rceil$. Given some positive constant c and*

$\alpha_0 > 0$ with $1/\alpha > 1/\alpha_0$, the following inequality holds:

$$\text{ERUB}_{\min} \leq \mathcal{O} \left(\frac{\kappa^2}{\kappa + \eta_0/\alpha} \left[\mathbb{I}_{T \leq \hat{T}} + \gamma^{\hat{T}-1} \ln \left(1 + \frac{\eta_0}{\kappa\alpha} \right) \mathbb{I}_{T > \hat{T}} \right] \right).$$

THEOREM 4.7 (DYNAMIC SCHEDULE). *Suppose the assumptions in [Theorem 4.5](#) are satisfied. Let $\alpha' = \frac{2\eta_0\alpha}{\gamma(1-\gamma\beta^2)}$, $\beta < \gamma$ and $\hat{T} = \max t$ s.t. $\gamma^{t-1} \geq \frac{1-\beta}{1-\beta^t}$. Use the following schedule: $\sigma_t^2 = \frac{1}{R} \sum_{i=1}^T \sqrt{\frac{q_i}{q_t}}$, $T^{\text{dyn}} = \left\lceil \mathcal{O} \left(\frac{2\kappa}{\eta_0} \ln \left(1 + \frac{\eta_0}{\kappa\alpha} \right) \right) \right\rceil$, where $q_t = c_1 \gamma^{T+t} \mathbb{I}_{T \leq \hat{T}} + \gamma^{\hat{T}-1} c_2 \gamma^{T-t} \mathbb{I}_{T > \hat{T}}$ for some positive constants c_1 and c_2 . The following inequality holds:*

$$\text{ERUB} \leq \gamma^T + 2\eta_0\alpha \sum_{t=1}^T Rq_t \sigma_t^2,$$

$$\text{ERUB}_{\min} \leq \mathcal{O} \left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \mathbb{I}_{T \leq \hat{T}} + \mathbb{I}_{T > \hat{T}} \right) \right).$$

Discussion. Theoretically, the dynamic schedule is more influential in vanilla gradient descent methods than the momentum variant. The result is mainly attributed to the averaging operation. The moving averaging, $(1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i / (1 - \beta^t)$, increase the influence of the under-presented initial steps and decrease the one of the over-sensitive last steps. Counterintuitively, the preferred dynamic schedule should be increasing since q_t decreases when $t \leq \hat{T}$.

4.3 Extension to Private Stochastic Gradient Descent

Though PGD provides a guarantee both for utility and privacy, computing gradients of the whole dataset is impractical for large-scale problems. For this sake, studying the convergence of Private Stochastic Gradient Descent (PSGD) is meaningful. The [Algorithm 1](#) can be easily extended to PSGD by subsampling n gradients where the batch size $n \ll N$. According to [\[40\]](#), when privacy is measured by zCDP, there are two ways to account for the privacy cost of PSGD depending on the batch-sampling method: sub-sampling with or without replacement. In this paper, we focus on the random subsampling with replacement since it is widely used in deep learning in literature, e.g., [\[1, 11\]](#). Accordingly, we replace N in the definition of α by n because the term is from the sensitivity of batch data (see [Eq. \(1\)](#)). For clarity, we assume that T is the number of iterations rather than epochs and that $\tilde{\nabla}_t$ is mean stochastic gradient.

When a batch of data are randomly sampled, the privacy cost of one iteration is cp^2/σ_t^2 where c is some constant, $p = n/N$ is the sample rate, and $1/\sigma_t^2$ is the full-batch privacy cost. Details of the sub-sampling theorems are referred to the [Theorem 3](#) of [\[40\]](#) and their empirical setting. Therefore, we can replace the privacy constraint $\sum_t p^2/\sigma_t^2 = R$ by $\sum_t 1/\sigma_t^2 = R'$ where $R' = R/p^2 = \frac{N^2}{n^2}R$. Remarkably, we omit the constant c because it will not affect the results regarding uniform or dynamic schedules. Notice N^2R in the α is replaced by $n^2R' = N^2R$. Thus, the form of α is not changed which provides convenience for the following derivations.

Now we study the utility bound of PSGD. To quantify the randomness of batch sampling, we define a random vector ξ_t with $\mathbb{E}[\xi_t] = 0$ and $\mathbb{E} \|\xi_t\|^2 \leq D$ such that $\tilde{\nabla}_t \leq \nabla_t + \sigma_g \xi_t / n$ for some positive constant σ_g . Because ξ_t has similar property to the privacy

noise v_t , we can easily extend the PGD bounds to PSGD bounds by following theories.

THEOREM 4.8 (UTILITY BOUNDS OF PSGD). *Let α , κ and γ be defined in Eq. (6), and $\eta_t = \frac{1}{M}$. Suppose $f(\theta; x_i)$ is G -Lipschitz and $f(\theta)$ is M -smooth satisfying the Polyak-Lojasiewicz condition. For PSGD, when batch size satisfies $n = \max\{N\sqrt{R}, 1\}$, the following holds: $\text{ERUB} = \gamma^T + \alpha_g \sigma_g^2 + R' \sum_{t=1}^T q_t \sigma_t^2$, where $q_t \triangleq \gamma^{T-t} \alpha$, $\sum_t 1/\sigma_t^2 = R'$. where $\alpha_g = \frac{D}{2\mu N^2 R(f(\theta_1) - f(\theta^*))}$.*

THEOREM 4.9 (PSGD WITH MOMENTUM). *Let $\alpha_g = \frac{D}{2\mu N^2 R(f(\theta_1) - f(\theta^*))}$. Suppose assumptions in Theorem 4.5 holds. When batch size satisfies $n = \max\{N\sqrt{R}, 1\}$, the $U_3(\sigma, T)$ has to be replaced by $\tilde{U}_3 = U_3^g + U_3$, with $\alpha R' U_3^g \leq \alpha_g \sigma_g^2$ when PSGD is used.*

As shown above, the utility bound of PSGD differs from the PGD merely by $\alpha_g \sigma_g^2$. Note $\alpha_g = O(\frac{D}{N^2 R})$ which fits the order of dynamic-schedule bounds. In addition, α and other variables are not changed. Hence, the conclusions w.r.t. the dynamic/uniform schedules maintain the same.

4.4 Comparison of generalization bounds

In addition to the empirical risk bounds in Table 1, in this section we study the *true risk bounds*, or generalization error bounds. True risk bounds characterize how well the learnt model can generalize to unseen samples subject to the inherent data distribution. By leveraging the generic learning-theory tools, we extend our results to the *True Excess Risk* (TER) for strongly convex functions as follows. For a model θ , its TER is defined as follows:

$$\text{TER} \triangleq \mathbb{E}_{x \sim \mathcal{X}} [\mathbb{E}[f(\theta; x)]] - \min_{\hat{\theta}} \mathbb{E}_{x \sim \mathcal{X}} [f(\hat{\theta}; x)],$$

where the second expectation is over the randomness of generating θ (e.g., the noise and stochastic batches). Assume a dataset d consist of N samples drawn i.i.d. from the distribution \mathcal{X} . Two approaches could be used to extend the empirical bounds to the true excess risk: One is proposed by [30] where the true excess risk of PGD can be bounded in high probability. For example, [3] achieved a $\frac{\ln^2 N}{N}$ bound with N^2 iterations. Alternatively, instead of relying on the probabilistic bound, Bassily et al. [2] used the uniform stability to give a tighter bound. Later, Feldman et al. [11] improve the efficiency of gradient computation to achieve a similar bound. Both approaches introduce an additive term to the empirical bounds. In this section, we adopt both approaches to investigate the two types of resulting true risk bounds.

(1) True Risk in High Probability. First, we consider the high-probability true risk bound. Based on Section 5.4 from [30] (restated in Theorem 4.10), we can relate the EER to the TER.

THEOREM 4.10. *Let $f(\theta; x)$ be G -Lipschitz, and $f(\theta)$ be μ -strong convex loss function given any $x \in \mathcal{X}$. With probability at least $1 - p$ over the randomness of sampling the data set d , the following inequality holds:*

$$\text{TER}(\theta) \leq \sqrt{\frac{2G^2}{\mu N} \sqrt{f(\theta) - f(\theta^*)}} + \frac{4G^2}{p\mu N}, \quad (15)$$

where $\theta^* = \arg \min_{\theta} f(\theta)$.

To apply the Eq. (15), we need to extend EER, the expectation bound, to a high-probability bound. Following [3] (Section D), we repeat the PGD with privacy budget R/k for k times. Note, the output of all repetitions is still of R budget. When $k = 1$, let the EER of the algorithm be denoted as $F(R)$. Then the EER of one execution of the k repetitions is $F(R/k)$ where privacy is accounted by zCDP. When $k = \log_2(1/p)$ for $p \in [0, 1]$, by Markov's inequality, there exists one repetition whose EER is $F(R/\log_2(1/p))$ with probability at least $1 - 1/2^k = 1 - p$. Combined with Eq. (15), we use the bounds of uniform schedule and dynamic schedules in Section 4.1.3 to obtain:

$$\text{TER}^{\text{uniform}} \leq \tilde{O} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(N) \ln(1/p)}{NR}} + \frac{4}{p} \right) \right), \quad (16)$$

$$\text{TER}^{\text{dynamic}} \leq \tilde{O} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(1/p)}{NR}} + \frac{4}{p} \right) \right), \quad (17)$$

where we again ignore the κ and other constants. Similarly, we can extend the momentum methods.

(2) True Risk by Uniform Stability. Following Bassily et al. [2], we use the uniform stability (defined in Definition 4.2) to extend the empirical bounds. We restate the related definition and theorems as follows.

DEFINITION 4.2 (UNIFORM STABILITY). *Let $s > 0$. A randomized algorithm $\mathcal{M} : \mathcal{D}^N \rightarrow \Theta$ is s -uniformly stable w.r.t. the loss function f if for any neighbor datasets d and d' , we have:*

$$\sup_{x \in \mathcal{X}} \mathbb{E}[f(\mathcal{M}(d); x) - f(\mathcal{M}(d'); x)] \leq s,$$

where the expectation is over the internal randomness of \mathcal{M} .

THEOREM 4.11 (SEE, E.G., [29]). *Suppose $\mathcal{M} : \mathcal{D}^N \rightarrow \Theta$ is a s -uniformly stable algorithm w.r.t. the loss function f . Let \mathcal{D} be any distribution over data space and let $d \sim \mathcal{D}^N$. The following holds true.*

$$\mathbb{E}_{d \sim \mathcal{D}^N} [\mathbb{E}[f(\mathcal{M}(d); \mathcal{D}) - f(\mathcal{M}(d); d)]] \leq s,$$

where the second expectation is over the internal randomness of \mathcal{M} . $f(\mathcal{M}(d); \mathcal{D})$ and $f(\mathcal{M}(d); d)$ represent the true loss and the empirical loss, respectively.

THEOREM 4.12 (UNIFORM STABILITY OF PGD FROM [2]). *Suppose $\eta < 2/M$ for M smooth, G -Lipschitz $f(\theta; x)$. Then PGD is s -uniformly stable with $s = G^2 T \eta / N$.*

Combining Theorems 4.11 and 4.12, we obtain the following:

$$\text{TER} \leq \text{EER} + G^2 \frac{\eta T}{N}.$$

Because EER in this paper compresses a γ^T or similar exponential terms, unlike [2], we cannot directly minimize the TER upper bound w.r.t. T and η in the presence of a polynomial form of γ^T and T . Therefore, we still use $T = O(\ln \frac{N^2 R}{D})$ and η for minimizing EER. Note that

$$G^2 \frac{\eta T}{N} \leq O\left(\frac{G^2}{MN} \ln \frac{N^2 R}{D}\right) \leq O\left(\frac{G^2}{M}\right)$$

where we assume $N \gg D$ and use $\ln N \leq N$. Because the term $O(G^2/M)$ is constant and independent from dimension, we follow [2] to drop the term when comparing the bounds. After dropping

Table 2: Comparison of true excess risk bounds. The algorithms are T -iteration $\frac{1}{2}R$ -zCDP or equivalently (ϵ, δ) -DP under the μ -strongly-convex condition. The \mathcal{O} notation in this table drops other \ln terms. Assume loss functions are 1-smooth and 1-Lipschitz continuous, and all parameters satisfy their numeric assumptions. * marks the method with convex assumption.

Algorithm	Utility Upper Bd.	T
GD+Adv [3]	$\mathcal{O}_{1-p} \left(\frac{\sqrt{D} \ln^2 N \ln(1/p)}{p \mu N R \epsilon, \delta} \right)$	$\mathcal{O}(N^2)$
SVRG+MA [34]	$\mathcal{O} \left(\frac{D \ln N}{\mu N^2 R \epsilon, \delta} \right)$	$\mathcal{O}(\ln \frac{N^2 R \epsilon, \delta}{D})$
SSGD+zCDP [11]	$\mathcal{O} \left(\left(\frac{1}{\sqrt{N}} + \frac{2\sqrt{D}}{\sqrt{RN}} \right) \ln N \right)$	$\mathcal{O}(\frac{N^2}{16D/R^2+4N})$
* SGD+MA [2]	$\mathcal{O} \left(\max \left\{ \frac{\sqrt{D}}{N \sqrt{R \epsilon, \delta}}, \frac{1}{\sqrt{N}} \right\} \right)$	$\mathcal{O}(\min \{ \frac{N}{8}, \frac{N^2 R \epsilon, \delta}{32D} \})$
True risk in high probability $(1-p)$		
GD+zCDP, Static Schedule	$\mathcal{O}_{1-p} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(N) \ln(1/p)}{NR}} + \frac{4}{p} \right) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
GD+zCDP, Dynamic Schedule	$\mathcal{O}_{1-p} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(1/p)}{NR}} + \frac{4}{p} \right) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
Momentum+zCDP, Static Sch.	$\mathcal{O}_{1-p} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(1/p)}{NR}} (c + \ln N \mathbb{I}_{T > \hat{T}}) + \frac{4}{p} \right) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
Momentum+zCDP, Dynamic Sch.	$\mathcal{O}_{1-p} \left(\frac{G^2}{\mu N} \left(\sqrt{\frac{D \ln(1/p)}{NR}} (1 + \frac{cD}{N^2 R} \mathbb{I}_{T > \hat{T}}) + \frac{4}{p} \right) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
True risk by uniform stability		
GD, Non-Private	$\mathcal{O} \left(\frac{D}{N^2 R} \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
GD+zCDP, Static Schedule	$\mathcal{O} \left(\frac{D \ln N}{N^2 R} \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
GD+zCDP, Dynamic Schedule	$\mathcal{O} \left(\frac{D}{N^2 R} \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
Momentum+zCDP, Static Sch.	$\mathcal{O} \left(\frac{D}{N^2 R} (c + \ln N \mathbb{I}_{T > \hat{T}}) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$
Momentum+zCDP, Dynamic Sch.	$\mathcal{O} \left(\frac{D}{N^2 R} (1 + \frac{cD}{N^2 R} \mathbb{I}_{T > \hat{T}}) \right)$	$\mathcal{O}(\ln \frac{N^2 R}{D})$

the additive term, it is obvious to see that the advantage of dynamic schedules still maintains since $\text{TER} \leq \text{EER}$. A similar extension can be derived for [34].

We summarize the results and compare them to prior works in Table 2 where we include an additional method: Snowball Stochastic Gradient Descent (SSGD). SSGD dynamically schedule the batch size to achieve an optimal convergence rate in linear time.

Discussion. By using uniform stability, we successfully transfer the advantage of our dynamic schedules from empirical bounds to true risk bounds. The inherent reason is that our bounds only need $\ln N$ iterations to reach the preferred final loss. With uniform stability, the logarithmic T reduce the gap caused by transferring. Compared to the [2, 11], our method has remarkably improved efficiency in T from N or N^2 to $\ln(N)$. That implies fewer iterations are required for converging to the same generalization error.

5 EXPERIMENTS

We empirically validate the properties of privacy schedules and their connections to learning algorithms. In this section, we briefly review the schedule behavior on quadratic losses under varying data sensitivity.

Dataset. We create a subset of the MNIST dataset [18] including 1000 handwritten images of 10 digits (MNIST). We also construct a subset of the MNIST dataset with digit 3 and 5 only, denoted as

MNIST35. Compared to the original dataset (70,000 samples), the small set will be more vulnerable to attack and the private learning will require larger noise (see the $1/N$ factor in Eq. (1)). Following the preprocessing in [1], we project the vectorized images into a 60-dimensional subspace extracted by PCA.

Setup. The samples are first normalized so that $\sum_{n=1}^N x_n = 0$ and the standard deviation is 1. Then the sample norms are scaled such that $\max_n \|x_n\| = 10$ (i.e., *data scales*). Upon the scaled data, we train a 2-layer Deep Neural Network (DNN) with 1000 hidden units by logistic regression. We fix the learning rate to 0.1 based on the corresponding experiments of non-private training (same setting without noise). The total privacy budget is $(4, 10^{-8})$ -DP, equal to 0.1963-zCDP, which implies $R = 0.3927$. To control the sensitivity of the gradients, we clip gradients by a clipping norm fixed at 4. Formally, we scale down the sample gradients to length 4 if its norm is larger than 4. Because the schedule highly depends on the iteration number T , we grid search the best T in range [50, 150] for compared methods. Therefore, we ignore the privacy cost of such tuning in our experiments which protocol is also used in previous work [1, 36]. All the experiments are repeated 100 times and metrics are averaged afterwards.

We first show the estimated influence of step noise q_t (by re-training the private learning algorithms) in Fig. 1 Left. We see the trends of influence are approximately in an exponential form of

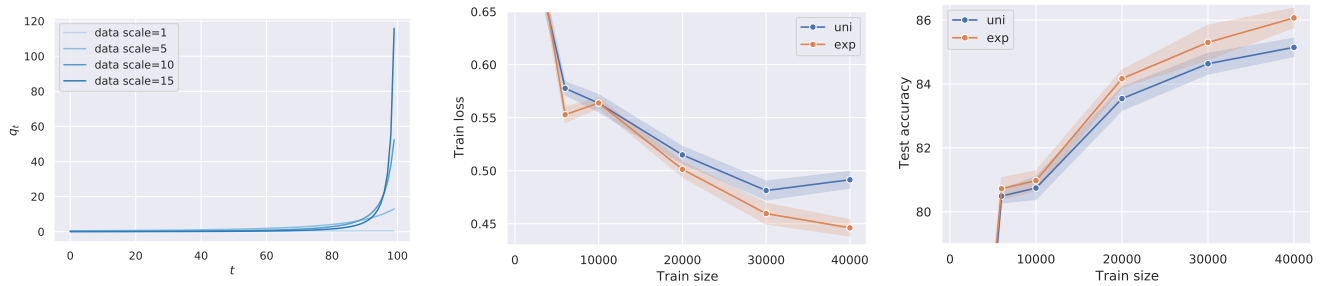


Figure 1: Comparison of dynamic schedule and uniform schedule on different data scale. Left pane is the influence by iteration estimated by retraining. The rest two panes are performance of DNN trained on the MNIST35 dataset with a varying total number of training samples, when the exponential influence is estimated on a randomly-generated auxiliary dataset.

t . By Eq. (10), the resultant schedule on noise scale σ_t will be a normalized exponential decay. This observation motivates the use of exponential decay schedule in practice.

To estimate the influence without extra privacy costs, we use an auxiliary set, which is randomly sampled from Gaussian distribution, to pick the proper influence curvature parameterized by an exponential function. We use auxiliary synthesized datasets of the same size as the corresponding private datasets to tune the parameters. We vary the size of training data to examine the data efficiency of the dynamic schedule denoted as exp. For a fair comparison, we also choose the hyper-parameters of uniform schedule (uni) on the same auxiliary dataset. We show that as the training size increases, exp outperforms uni both on the training loss and the test accuracy. The result verifies our theoretic conclusion: dynamic schedule is more data efficient than the static schedule.

6 CONCLUSION

When a privacy budget is provided for a certain learning task, one has to carefully schedule the privacy usage through the learning process. Uniformly scheduling the budget has been widely used in literature whereas increasing evidence suggests that dynamically schedules could empirically outperform the uniform one. This paper provided a principled analysis on the problem of optimal budget allocation and connected the advantages of dynamic schedules to both the loss structure and the learning behavior. We further validated our results through empirical studies.

ACKNOWLEDGMENTS

This material is based in part upon work supported by National Institute of Aging (1RF1AG072449), Office of Naval Research (N00014-20-1-2382), National Science Foundation (IIS-1749940). Z. Wang is in part supported by Good Systems, a UT Austin Grand Challenge to develop responsible AI technologies

REFERENCES

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep Learning with Differential Privacy. In *CCS: Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS '16)*. ACM, New York, NY, USA, 308–318.
- [2] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Guha Thakurta. 2019. Private Stochastic Convex Optimization with Optimal Rates. In *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 11282–11291.
- [3] R. Bassily, A. Smith, and A. Thakurta. 2014. Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*. 464–473.
- [4] Mark Bun and Thomas Steinke. 2016. Concentrated Differential Privacy: Simplifications, Extensions, and Lower Bounds. In *Theory of Cryptography*. Vol. 9985. Springer Berlin Heidelberg, Berlin, Heidelberg, 635–658.
- [5] Kamalika Chaudhuri, Claire Monteleoni, and Anand D. Sarwate. 2011. Differentially Private Empirical Risk Minimization. *Journal of Machine Learning Research* 12, Mar (2011), 1069–1109.
- [6] Junhong Cheng, Wenyang Liu, Xiaoling Wang, Xingjian Lu, Jing Feng, Yi Li, and Chaofan Duan. 2020. Adaptive Distributed Differential Privacy with SGD. *Workshop on Privacy-Preserving Artificial Intelligence* (2020), 6.
- [7] Rachel Cummings, Sara Krehbiel, Kevin A Lai, and Uthaipon Tantipongpipat. 2018. Differential Privacy for Growing Databases. In *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., 8864–8873.
- [8] Damien Desfontaines and Balázs Pejő. 2019. SoK: Differential Privacies. *arXiv:1906.01337 [cs]* (June 2019).
- [9] Cynthia Dwork, Alan Karr, Kobbi Nissim, and Lars Vilhuber. 2020. On Privacy in the Age of COVID-19. *Journal of Privacy and Confidentiality* 10, 2 (June 2020).
- [10] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. In *Theory of Cryptography (Lecture Notes in Computer Science)*. Springer Berlin Heidelberg, 265–284.
- [11] Vitaly Feldman, Tomer Koren, and Kunal Talwar. 2020. Private stochastic convex optimization: optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing (STOC 2020)*. Association for Computing Machinery, New York, NY, USA, 439–449.
- [12] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model Inversion Attacks That Exploit Confidence Information and Basic Countermeasures. In *CCS: Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS '15)*. ACM, New York, NY, USA, 1322–1333.
- [13] Junyuan Hong, Haotao Wang, Zhangyang Wang, and Jiayu Zhou. 2021. Learning Model-Based Privacy Protection under Budget Constraints. In *AAAI* 9.
- [14] Xixi Huang, Jian Guan, Bin Zhang, Shuhan Qi, Xuan Wang, and Qing Liao. 2019. Differentially Private Convolutional Neural Networks with Adaptive Gradient Descent. In *2019 IEEE Fourth International Conference on Data Science in CyberSpace (DSC)*. 642–648.
- [15] Prateek Jain, Dheeraj Nagaraj, and Praneeth Netrapalli. 2019. Making the Last Iterate of SGD Information Theoretically Optimal. In *Conference on Learning Theory*. 1752–1755.
- [16] Hamed Karimi, Julie Nutini, and Mark Schmidt. 2016. Linear Convergence of Gradient and Proximal-Gradient Methods Under the Polyak-Lojasiewicz Condition. In *Machine Learning and Knowledge Discovery in Databases (Lecture Notes in Computer Science)*. Springer International Publishing, Cham, 795–811.
- [17] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *the 3rd International Conference for Learning Representations*. San Diego, CA.
- [18] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (Nov. 1998), 2278–2324.
- [19] Jaewoo Lee and Daniel Kifer. 2018. Concentrated Differentially Private Gradient Descent with Adaptive per-Iteration Privacy Budget. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '18)*. ACM, New York, NY, USA, 1656–1665.
- [20] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Artificial Intelligence and Statistics*. 1273–1282.

- [21] H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. 2018. Learning Differentially Private Recurrent Language Models. In *International Conference on Learning Representations*.
- [22] Yurii Nesterov and B.T. Polyak. 2006. Cubic regularization of Newton method and its global performance. *Mathematical Programming* 108, 1 (Aug. 2006), 177–205.
- [23] Venkatadheeraj Pichapati, Ananda Theertha Suresh, Felix X. Yu, Sashank J. Reddi, and Sanjiv Kumar. 2019. AdaClip: Adaptive Clipping for Private SGD. *arXiv:1908.07643 [cs, stat]* (Oct. 2019).
- [24] B. T. Polyak. 1963. Gradient methods for the minimisation of functionals. *U. S. S. R. Comput. Math. and Math. Phys.* 3, 4 (Jan. 1963), 864–878.
- [25] B. T. Polyak. 1964. Some methods of speeding up the convergence of iteration methods. *U. S. S. R. Comput. Math. and Math. Phys.* 4, 5 (Jan. 1964), 1–17.
- [26] Ning Qian. 1999. On the momentum term in gradient descent learning algorithms. *Neural Networks* 12, 1 (Jan. 1999), 145–151.
- [27] Sashank J. Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. 2016. Stochastic Variance Reduction for Nonconvex Optimization. In *International Conference on Machine Learning*, 314–323.
- [28] Alfréd Rényi. 1961. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California.
- [29] Shai Shalev-Shwartz and Shai Ben-David. 2014. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press.
- [30] Shai Shalev-Shwartz, Nathan Srebro, and Karthik Sridharan. 2009. Stochastic Convex Optimization. In *Proceedings of the 22nd Annual Conference on Learning Theory, COLT '09*, 11.
- [31] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. 2017. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy (SP)*, 3–18.
- [32] Om Thakkar, Galen Andrew, and H. Brendan McMahan. 2019. Differentially Private Learning with Adaptive Clipping. *arXiv:1905.03871 [cs, stat]* (May 2019).
- [33] Di Wang, Changyou Chen, and Jinhui Xu. 2019. Differentially Private Empirical Risk Minimization with Non-convex Loss Functions. In *International Conference on Machine Learning*, 6526–6535.
- [34] Di Wang, Minwei Ye, and Jinhui Xu. 2017. Differentially Private Empirical Risk Minimization Revisited: Faster and More General. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2722–2731.
- [35] Michael W. Weiner, Dallas P. Veitch, Paul S. Aisen, Laurel A. Beckett, Nigel J. Cairns, Robert C. Green, Danielle Harvey, Clifford R. Jack, William Jagust, Enchi Liu, John C. Morris, Ronald C. Petersen, Andrew J. Saykin, Mark E. Schmidt, Leslie Shaw, Li Shen, Judith A. Siuciak, Holly Soares, Arthur W. Toga, and John Q. Trojanowski. 2013. The Alzheimer’s Disease Neuroimaging Initiative: A review of papers published since its inception. *Alzheimer’s & Dementia* 9, 5 (Sept. 2013), e111–e194.
- [36] Xi Wu, Fengang Li, Arun Kumar, Kamalika Chaudhuri, Somesh Jha, and Jeffrey Naughton. 2017. Bolt-on Differential Privacy for Scalable Stochastic Gradient Descent-based Analytics. In *Proceedings of the 2017 ACM International Conference on Management of Data (SIGMOD '17)*. ACM, New York, NY, USA, 1307–1322.
- [37] Yun Xie, Peng Li, Chao Wu, and Qiuling Wu. 2021. Differential Privacy Stochastic Gradient Descent with Adaptive Privacy Budget Allocation. In *2021 IEEE International Conference on Consumer Electronics and Computer Engineering (ICCECE)*, 227–231.
- [38] Zhiying Xu, Shuyu Shi, Alex X. Liu, Jun Zhao, and Lin Chen. 2020. An Adaptive and Fast Convergent Approach to Differentially Private Deep Learning. *the Proceedings of IEEE International Conference on Computer Communications (2020)*.
- [39] Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. 2020. Gradient Perturbation is Underrated for Differentially Private Convex Optimization. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, Yokohama, Japan, 3117–3123.
- [40] Lei Yu, Ling Liu, Calton Pu, Mehmet Emre Gursoy, and Stacey Truex. 2019. Differentially Private Model Publishing for Deep Learning. *proceedings of 40th IEEE Symposium on Security and Privacy* (April 2019).
- [41] Xinyue Zhang, Jiahao Ding, Maoqiang Wu, Stephen T. C. Wong, Hien Van Nguyen, and Miao Pan. 2021. Adaptive Privacy Preserving Deep Learning Algorithms for Medical Data. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1169–1178.
- [42] Yingxue Zhou, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Arindam Banerjee. 2020. Private Stochastic Non-Convex Optimization: Adaptive Algorithms and Tighter Generalization Bounds. *arXiv:2006.13501 [cs, stat]* (Aug. 2020).

A SOCIAL IMPACT

The wide usage of personal data in training machine learning has led to huge successes in many application domains but is also accompanied by rising concerns on privacy protection due to the sensitive information in the data. The development of privacy-preserving algorithms has become one of critical research areas of machine learning, in which the key challenge is to train high performance models under the constraint of a given privacy budget, or how much sensitive information can be accessed during the training phase. Differential privacy provided a principled framework to quantify the privacy budget, under which researchers proposed various schemes to schedule the budget usage during a learning process, yet there is a lack of systematical studies on when and why some schedules are better than other. Our efforts in this paper are among the first to study and compare the effectiveness of these schedules from a rigorous optimization perspective. Our theoretical results can benefit any privacy-preserving machine learning practitioners to efficiently and effectively choose proper privacy schedules tailored to their learning tasks.

B COMPARISON OF ALGORITHMS

Here we elaborate the meaning of algorithm names in [Table 1](#). Asymptotic upper bounds are achieved when sample size N approaches infinity. Both R and $R_{\epsilon, \delta}$ with $R_{\epsilon, \delta} < R$ are the privacy budgets of corresponding algorithms. Specifically, $R_{\epsilon, \delta} = \epsilon^2 / \ln(1/\delta) < R$ when the private algorithm is (ϵ, δ) -DP with $\epsilon \leq 2 \ln(1/\delta)$.

PGD+Adv. Adv denotes the Advanced Composition method [3]. The method assumes that loss function is 1-strongly convex which implies the PL condition and optimized variable is in a convex set of diameter 1 w.r.t. l_2 norm.

PGD+MA and the adjusted-utility version. MA denotes the Moment Accountant [1] which improve the composed privacy bound versus the Advanced Composition. The improvement on privacy bound lead to a enhanced utility bound, as a result.

PGD+Adv+BBImp. The dynamic method assumes that the loss is 1-strongly convex and data comes in stream with $n \leq N$ samples at each round. Their utility upper bound is achieved at some probability p with any positive c .

Adam+MA. The authors prove a convergence bound for the gradient norms which is extended to loss bound by using PL condition. They also presents the results for AdaGrad and GD which are basically of the same upper bound. Our theorems improve their bound by using the recursive derivation based on the PL condition, while their bound is a simple application of the condition on the gradient norm bound.

GD, Non-Private. This method does not inject noise into gradients but limit the number of iterations. With the bound, we can see that our utility bound are optimal with dynamic schedule.

GD+zCDP. We discussed the static and dynamic schedule for the gradient descent method where the dynamic noise influence is the key to tighten the bound.

Momentum+zCDP. Different from the GD+zCDP, momentum methods will have two phase of utility upper bound. When T is small than some positive constant \hat{T} , the bound is as tight as the non-private one. Afterwards, the momentum has a bound degraded as the GD bound.

C PRELIMINARIES

C.1 Privacy

LEMMA C.1 (COMPOSITION & POST-PROCESSING). *Let two mechanisms be $M : \mathcal{D}^n \rightarrow \mathcal{Y}$ and $M' : \mathcal{D}^n \times \mathcal{Y} \rightarrow \mathcal{Z}$. Suppose M satisfies (ρ_1, a) -zCDP and $M'(\cdot, y)$ satisfies (ρ_2, a) -zCDP for $\forall y \in \mathcal{Y}$. Then, mechanism $M'' : \mathcal{D}^n \rightarrow \mathcal{Z}$ (defined by $M''(x) = M'(x, M(x))$) satisfies $(\rho_1 + \rho_2)$ -zCDP.*

DEFINITION C.1 (SENSITIVITY). *The sensitivity of a gradient query ∇_t to the dataset $\{x_i\}_{i=1}^N$ is*

$$\begin{aligned} \Delta_2(\nabla_t) &= \max_n \left\| \frac{1}{N} \sum_{j=1, j \neq n}^N \nabla_t^{(j)} - \frac{1}{N} \sum_{j=1}^N \nabla_t^{(j)} \right\|_2 \\ &= \frac{1}{N} \max_n \left\| \nabla_t^{(n)} \right\|_2 \end{aligned} \quad (18)$$

where $\nabla_t^{(n)}$ denotes the gradient of the n -th sample.

LEMMA C.2 (GAUSSIAN MECHANISM [4]). *Let $f : \mathcal{D}^n \rightarrow \mathcal{Z}$ have sensitivity Δ . Define a randomized algorithm $M : \mathcal{D}^n \rightarrow \mathcal{Z}$ by $M(x) \leftarrow f(x) + \mathcal{N}(0, \Delta^2 \sigma^2 I)$. Then M satisfies $\frac{1}{2\sigma^2}$ -zCDP.*

LEMMA C.3 ([4]). *If M is a mechanism satisfying ρ -zCDP, then M is $(\rho + 2\sqrt{\rho \ln(1/\delta)}, \delta)$ -DP for any $\delta > 0$.*

By solving $\rho + 2\sqrt{\rho \ln(1/\delta)} = \epsilon$, we can get $\rho = \epsilon + 2 \ln(1/\delta) + 2\sqrt{\ln(1/\delta)(\epsilon + \ln(1/\delta))}$.

C.2 Auxiliary lemmas

LEMMA C.4. *If $\max_n \|x_n\|_2 = 1$ and $\frac{1}{N} \sum_n x_n = 0$, then the gradient sensitivity of the squared loss will be*

$$\Delta_2(\nabla) = \max_i \frac{1}{N} \sqrt{2f(\theta; x_i)} \|x_i\|_2 \leq \frac{1}{2} (DM \|\theta\|^2 + 1),$$

where $\Theta_{\mathcal{M}}$ is the set of all possible parameters θ_t generated by the learning algorithm \mathcal{M} .

PROOF. According to the definition of sensitivity in Eq. (18), we have

$$\Delta_2(\nabla) = \max_i \left\| \nabla^{(i)} \right\|_2 = \max_n \frac{1}{n} \left\| A^{(i)} \theta - x_i \right\|_2$$

where we use i denotes the index of sample in the dataset. Here, we assume it is constant 1. We may get

$$\begin{aligned} \left\| A^{(i)} \theta - x_i \right\|_2^2 &= \left\| x_i (x_i^\top \theta - 1) \right\|_2^2 \\ &= (x_i^\top \theta - 1)^2 \left\| x_i \right\|_2^2 = 2f(\theta; x_i) \left\| x_i \right\|_2^2 \end{aligned} \quad (19)$$

where $f(\theta; x_i) = \frac{1}{2}(x_i^\top \theta - 1)^2$. Thus,

$$\Delta_2(\nabla) = \max_i \frac{1}{N} \sqrt{2f(\theta; x_i)} \left\| x_i \right\|_2$$

Since $\left\| x_n \right\|_2 \leq 1$ and $\frac{1}{N} \sum_{n=1}^N x_n = 0$,

$$\begin{aligned} f(\theta) &= \frac{1}{2N} \sum_{n=1}^N [(x_n^\top \theta)^2 - 2x_n^\top \theta + 1] \\ &\leq \frac{1}{2N} \sum_{n=1}^N [(\left\| x_n \right\| \left\| \theta \right\|)^2 + 1] \\ &\leq \frac{1}{2} (DM \left\| \theta \right\|^2 + 1) \end{aligned}$$

□

LEMMA C.5. Assume assumptions in Theorem 4.5 are satisfied. Given variables defined in Theorem 4.5, the following inequality holds true:

$$\begin{aligned} &\sum_{t=1}^T \gamma^{T-t} \frac{2(1-\beta)\eta_t}{b_t} \sum_{i=1}^t \beta^{t-i} \left\| \nabla_t - \nabla_i \right\|^2 \\ &\leq \frac{\eta_0^3 \beta \gamma}{2M(1-\beta)^3(\gamma-\beta)^2} \sum_{i=1}^{T-1} \gamma^{T-i} \left\| v_{i+1} \right\|^2. \end{aligned}$$

PROOF. We first handle the inner summation. By smoothness, the inequality $\left\| \nabla f(x) - \nabla f(y) \right\| \leq M \left\| x - y \right\|$ holds true. Thus,

$$\begin{aligned} &\sum_{i=1}^t \beta^{t-i} \left\| \nabla_t - \nabla_i \right\|^2 \\ &\leq M^2 \sum_{i=1}^t \beta^{t-i} \left\| \theta_t - \theta_i \right\|^2 \\ &= M^2 \sum_{k=0}^{t-1} \beta^k \left\| \theta_t - \theta_{t-k} \right\|^2 \\ &= M^2 \sum_{k=0}^{t-1} \beta^k \left\| \sum_{i=t-k}^{t-1} \eta_i v_{i+1} / b_i \right\|^2 \\ &\leq M^2 \sum_{k=0}^{t-1} \beta^k \left(\sum_{j=t-k}^{t-1} \eta_j^2 / b_j^2 \right) \left(\sum_{i=t-k}^{t-1} \left\| v_{i+1} \right\|^2 \right) \end{aligned}$$

where the last inequality is by Cauchy-Schwartz inequality. Because $\frac{1}{b_t} = \frac{1}{1-\beta^t} \leq \frac{1}{1-\beta}$ and $\eta_t = \frac{\eta_0}{2M}$,

$$\sum_{i=1}^t \beta^{t-i} \left\| \nabla_t - \nabla_i \right\|^2 \quad (20)$$

$$\begin{aligned} &\leq \frac{\eta_0^2}{4(1-\beta)^2} \sum_{k=0}^{t-1} \beta^k k \sum_{i=t-k}^{t-1} \left\| v_{i+1} \right\|^2 \\ &= \frac{\eta_0^2}{4(1-\beta)^2} \sum_{k=0}^{t-1} \beta^k k \sum_{i=1}^{t-1} \left\| v_{i+1} \right\|^2 \mathbb{I}(i \geq t-k) \\ &= \frac{\eta_0^2}{4(1-\beta)^2} \sum_{i=1}^{t-1} \left\| v_{i+1} \right\|^2 \sum_{k=0}^{t-1} \beta^k k \mathbb{I}(k \geq t-i) \\ &= \frac{\eta_0^2}{4(1-\beta)^2} \sum_{i=1}^{t-1} \left\| v_{i+1} \right\|^2 \sum_{k=t-i}^{t-1} \beta^k k \end{aligned} \quad (21)$$

where $\mathbb{I}(\cdot)$ is the indicating function which output 1 if the condition holds true, otherwise 0.

Denote the left-hand-side of the conclusion as LHS. We plug Eq. (21) into LHS to get

$$\begin{aligned} \text{LHS} &\leq \sum_{t=1}^T \gamma^{T-t} \frac{1}{b_t} \frac{\eta_0^3}{4M(1-\beta)} \sum_{i=1}^{t-1} \|v_{i+1}\|^2 \sum_{k=t-i}^{t-1} \beta^k k \\ &\leq \frac{\eta_0^3}{4M(1-\beta)^2} \sum_{t=1}^T \gamma^{T-t} \sum_{i=1}^{t-1} \|v_{i+1}\|^2 \sum_{k=t-i}^{t-1} \beta^k k \end{aligned}$$

where we relax the upper bound by $\frac{1}{b_t} = \frac{1}{1-\beta^t} \leq \frac{1}{1-\beta}$. Using Lemma C.6 can directly lead to the conclusion:

$$\text{LHS} \leq \frac{\eta_0^3 \beta \gamma}{2M(1-\beta)^3 (\gamma-\beta)^2} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2.$$

□

LEMMA C.6. Given variables defined in Theorem 4.5, the following inequality holds true:

$$\begin{aligned} &\sum_{t=1}^T \gamma^{T-t} \sum_{i=1}^{t-1} \|v_{i+1}\|^2 \sum_{k=t-i}^{t-1} k \beta^k \\ &\leq \frac{2\beta\gamma}{(\gamma-\beta)^2(1-\beta)} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2. \end{aligned}$$

PROOF. We first derive the summation:

$$\begin{aligned} U_1(t, i) &\triangleq \sum_{k=t-i}^{t-1} \beta^k k \\ &= \sum_{k=t-i}^{t-1} \sum_{j=1}^k \beta^k \\ &= \sum_{k=t-i}^{t-1} \sum_{j=1}^{t-1} \beta^k \mathbb{I}(j \leq k) \\ &= \sum_{j=1}^{t-1} \sum_{k=\max(t-i, j)}^{t-1} \beta^k \\ &= \sum_{j=1}^{t-1} \frac{\beta^{\max(t-i, j)} - \beta^t}{1-\beta} \\ &= \frac{1}{1-\beta} \left((t-i)\beta^{t-i} + \frac{\beta^{t-i+1} - \beta^t}{1-\beta} - \frac{\beta - \beta^t}{1-\beta} \right) \\ &= \frac{1}{1-\beta} \left((t-i)\beta^{t-i} + \frac{\beta^{t-i+1} - \beta}{1-\beta} \right) \end{aligned}$$

Now, we substitute $U_1(t, i)$ into LHS and replace $t-i$ by j , i.e., $t = j+i$, to get

$$\begin{aligned} \text{LHS} &= \sum_{t=1}^T \gamma^{T-t} \sum_{i=1}^{t-1} \|v_{i+1}\|^2 \frac{1}{1-\beta} \left((t-i)\beta^{t-i} + \frac{\beta^{t-i+1} - \beta}{1-\beta} \right) \\ &= \sum_{i=1}^{T-1} \|v_{i+1}\|^2 \sum_{t=i+1}^T \gamma^{T-t} \frac{1}{1-\beta} \left((t-i)\beta^{t-i} + \frac{\beta^{t-i+1} - \beta}{1-\beta} \right) \\ &= \sum_{i=1}^{T-1} \|v_{i+1}\|^2 \sum_{j=1}^{T-i} \gamma^{T-(j+i)} \frac{1}{1-\beta} \left(j\beta^j + \frac{\beta^{j+1} - \beta}{1-\beta} \right) \\ &= \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2 \sum_{j=1}^{T-i} \gamma^{-j} \frac{1}{1-\beta} \left(j\beta^j + \frac{\beta^{j+1} - \beta}{1-\beta} \right) \\ &\leq \frac{1}{1-\beta} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2 \sum_{j=1}^{T-i} \left(j \left(\frac{\beta}{\gamma} \right)^j + \frac{\beta}{1-\beta} \left(\frac{\beta}{\gamma} \right)^j \right) \end{aligned}$$

Let $a = \beta/\gamma$, we show

$$\begin{aligned}
\sum_{j=1}^{T-i} ja^j &= \sum_{j=1}^{T-i} \sum_{o=1}^j a^j \\
&= \sum_{o=1}^{T-i} \sum_{j=o}^{T-i} a^j \\
&= \sum_{o=1}^{T-i} \left(\frac{a^o - a^{T-i+1}}{1-a} \right) \\
&= \frac{a - a^{T-i+1}}{(1-a)^2} - (T-i) \frac{a^{T-i+1}}{1-a} \\
&\leq \frac{a}{(1-a)^2}.
\end{aligned}$$

Thus,

$$\begin{aligned}
\text{LHS} &\leq \frac{1}{1-\beta} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2 \left(\frac{a}{(1-a)^2} + \frac{\beta}{1-\beta} \sum_{j=1}^{T-i} a^j \right) \\
&\leq \frac{1}{1-\beta} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2 \left(\frac{a}{(1-a)^2} + \frac{\beta}{1-\beta} \frac{a}{1-a} \right) \\
&\leq \frac{a}{(1-a)^2(1-\beta)} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2
\end{aligned}$$

Because $\gamma < 1$, $\beta < a = \beta/\gamma$ and

$$\frac{a}{(1-a)^2} + \frac{\beta}{1-\beta} \frac{a}{1-a} \leq \frac{2a}{(1-a)^2}.$$

Therefore,

$$\begin{aligned}
\text{LHS} &\leq \frac{2a}{(1-a)^2(1-\beta)} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2 \\
&= \frac{2\beta\gamma}{(\gamma-\beta)^2(1-\beta)} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2
\end{aligned}$$

□

LEMMA C.7. Suppose $\gamma \in (0, 1)$ and $\beta \in (0, 1)$. Define

$$\hat{T} = \max t \text{ s.t. } \gamma^{t-1} \geq \frac{1-\beta}{1-\beta^t}.$$

If $t \leq \hat{T}$, $\frac{1-\beta}{1-\beta^t} \leq \gamma^{t-1}$ for $t = 1, \dots, T$. If $t > \hat{T}$, $\frac{1-\beta}{1-\beta^t} < \gamma^{\hat{T}-1}$.

PROOF. Define $h(t) = \gamma^{t-1}(1-\beta^t)$ whose derivatives are

$$\begin{aligned}
h'(t) &= \gamma^{t-1}(1-\beta^t) \ln \gamma + \gamma^{t-1}(-\beta^t) \ln \beta \\
&= \gamma^{t-1} [\ln \gamma - \beta^t (\ln \gamma + \ln \beta)] \\
&= \gamma^{t-1} [1 - \beta^t (1 + \log_\gamma \beta)] \ln \gamma.
\end{aligned}$$

Simple calculation shows $1 - \beta^t (1 + \log_\gamma \beta) \Big|_{t=0} = -\log_\gamma \beta < 0$ and $\lim_{t \rightarrow +\infty} 1 - \beta^t (1 + \log_\gamma \beta) = 1$. When $t = -\log_\beta(1 + \log_\gamma \beta)$ denoted as t_0 , $1 - \beta^t (1 + \log_\gamma \beta) = 0$. Because $1 - \beta^t (1 + \log_\gamma \beta)$ is monotonically increasing by t and $\gamma^{t-1} \ln \gamma$ is negative, $h'(t) \geq 0$ if $t \leq t_0$. Otherwise, $h'(t) < 0$. Therefore, $h(t)$ is a concave function. Because $h(1) = 1 - \beta$ and $h(\hat{T}) = \gamma^{\hat{T}-1}(1 - \beta^{\hat{T}}) \geq 1 - \beta > 0$, $h(t) \geq 1 - \beta$ for $t = 1, \dots, \hat{T}$. Thus, for all $t \in [1, \hat{T}]$, we have $\frac{1-\beta}{1-\beta^t} \leq \gamma^{t-1}$.

For $t > \hat{T}$, because $\frac{1-\beta}{1-\beta^t}$ monotonically increases by t , we have $\frac{1-\beta}{1-\beta^t} < \frac{1-\beta}{1-\beta^{\hat{T}}} \leq \gamma^{\hat{T}-1}$.

□

D PROOFS

PROOF OF THEOREM 3.1. Because all sample gradient are G -Lipschitz continuous, the sensitivity of the averaged gradient is upper bounded by G/N . Based on Lemma C.2, the privacy cost of g_t is $\frac{1}{2\sigma_t^2}$ ¹.

Here, we make the output of each iteration a tuple of (θ_{t+1}, v_{t+1}) . For the 1st iteration, because θ_1 does not embrace private information by random initialization, the mapping,

$$\begin{bmatrix} v_2 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} g_1 \\ \theta_1 - \eta_1 g_1 \end{bmatrix},$$

is $\hat{\rho}_1$ -zCDP where $\hat{\rho}_1 = \frac{1}{2\sigma_1^2}$.

Suppose the output of the t -th iteration, (θ_t, v_t) , is $\hat{\rho}_t$ -zCDP. At each iteration, we have the following mapping $(\theta_t, v_t) \rightarrow (\theta_{t+1}, v_{t+1})$ defined as

$$\begin{bmatrix} v_{t+1} \\ \theta_{t+1} \end{bmatrix} = \begin{bmatrix} \phi(v_t, g_t) \\ \theta_t - \eta_t \phi(v_t, g_t) \end{bmatrix}.$$

Thus, the output tuple (θ_{t+1}, v_{t+1}) is $(\hat{\rho}_t + \frac{1}{2\sigma_t^2})$ -zCDP by Lemma C.1.

Thus, the recursion implies that (θ_{T+1}, v_{T+1}) has privacy cost as

$$\hat{\rho}_{T+1} = \hat{\rho}_T + \frac{1}{2\sigma_T^2} = \dots = \sum_{t=1}^T \frac{1}{2\sigma_t^2} = \frac{1}{2} \sum_{t=1}^T \rho_t \leq \frac{1}{2}(R - R_T) \leq \frac{1}{2}R.$$

Let $\rho = \hat{\rho}_{T+1}$. Then we can get the conclusion. □

D.1 Gradient Descents

PROOF OF THEOREM 4.1. With the definition of smoothness in Definition 3.3 and Eq. (1), we have

$$\begin{aligned} f(\theta_{t+1}) - f(\theta_t) &\leq -\eta_t \nabla_t^\top (\nabla_t + G\sigma_t v_t/N) + \frac{1}{2} M\eta_t^2 \|\nabla_t + G\sigma_t v_t/N\|^2 \\ &= -\eta_t (1 - \frac{1}{2} M\eta_t) \|\nabla_t\|^2 - (1 - M\eta_t) \eta_t \nabla_t^\top G\sigma_t v_t/N + \frac{1}{2} M\eta_t^2 \|G\sigma_t v_t/N\|^2 \\ &\leq -2\mu\eta_t (1 - \frac{1}{2} M\eta_t) (f(\theta_t) - f(\theta^*)) - (1 - M\eta_t) \eta_t \nabla_t^\top G\sigma_t v_t/N \\ &\quad + \frac{1}{2} M\eta_t^2 \|G\sigma_t v_t/N\|^2. \end{aligned}$$

where the last inequality is due to the Polyak-Lojasiewicz condition. Taking expectation on both sides, we can obtain

$$\mathbb{E}[f(\theta_{t+1})] - \mathbb{E}[f(\theta_t)] \leq -2\mu\eta_t (1 - \frac{M}{2}\eta_t) (\mathbb{E}[f(\theta_t)] - f(\theta^*)) + \frac{M}{2} (\eta_t G\sigma_t/N)^2 \mathbb{E}\|v_t\|^2$$

which can be reformulated by subtracting $f(\theta^*)$ on both sides and re-arranged as

$$\mathbb{E}[f(\theta_{t+1})] - f(\theta^*) \leq \left(1 - 2\mu\eta_t (1 - \frac{M}{2}\eta_t)\right) (\mathbb{E}[f(\theta_t)] - f(\theta^*)) + \frac{M}{2} (\eta_t G\sigma_t/N)^2 D$$

Recursively using the inequality, we can get

$$\begin{aligned} \mathbb{E}[f(\theta_{T+1})] - f(\theta^*) &\leq \prod_{t=1}^T \left(1 - 2\mu\eta_t (1 - \frac{M}{2}\eta_t)\right) (\mathbb{E}[f(\theta_1)] - f(\theta^*)) \\ &\quad + \frac{MD}{2} \sum_{t=1}^T \prod_{i=t+1}^T \left(1 - 2\mu\eta_i (1 - \frac{M}{2}\eta_i)\right) (\eta_t G\sigma_t/N)^2. \end{aligned}$$

Let $\eta_t \equiv 1/M$. Then the above inequality can be simplified as

$$\begin{aligned} \mathbb{E}[f(\theta_{T+1})] - f(\theta^*) &\leq \gamma^T (\mathbb{E}[f(\theta_1)] - f(\theta^*)) + R \sum_{t=1}^T \gamma^{T-t} \frac{MD}{2R} \left(\frac{\eta_t G}{N}\right)^2 \sigma_t^2 \\ &= \gamma^T (\mathbb{E}[f(\theta_1)] - f(\theta^*)) + R \sum_{t=1}^T \gamma^{T-t} \alpha \sigma_t^2 (\mathbb{E}[f(\theta_1)] - f(\theta^*)) \\ &= \left(\gamma^T + R \sum_{t=1}^T \gamma^t \alpha \sigma_t^2\right) (f(\theta_1) - f(\theta^*)) \end{aligned}$$

¹For brevity, when we say the privacy cost of some value, e.g., gradient, we actually refer to the cost of mechanism that output the value.

□

PROOF OF THEOREM 4.2. The minimizer of the upper bound of Eq. (7) can be written as

$$T^* = \arg \min_T \gamma^T + \alpha\kappa(1 - \gamma^T)T \quad (22)$$

where we substitute $\sigma^2 = T/R$ in the second line. To find the convex minimization problem, we need to vanishing its gradient which involves an equation like $T\gamma^T = c$ for some real constant c . However, the solution is $W_k(c)$ for some integer k where W is Lambert W function which does not have a simple analytical form. Instead, because $\gamma^T > 0$, we can minimize a surrogate upper bound as following

$$T^* = \arg \min_T \gamma^T + \alpha\kappa T = \frac{1}{\ln(1/\gamma)} \ln \left(\frac{\ln(1/\gamma)}{\kappa\alpha} \right), \text{ if } \kappa\alpha + \ln \gamma < 0 \quad (23)$$

where we use the surrogate upper bound in the second line and utilize $\gamma = 1 - \frac{1}{\kappa}$. However, the minimizer of the surrogate objective is not optimal for the original objective. When κ is large, the term, $-\alpha\kappa\gamma^T T$, cannot be neglected as we expect. On the other hands, T suffers from explosion if $\kappa \rightarrow \infty$ and meanwhile $1/\gamma \rightarrow 1$. The tendency is counterintuitive since a small T should be taken for sharp losses. To fix the issue, we change the form of T^* as

$$T^* = \frac{1}{\ln(1/\gamma)} \ln \left(1 + \frac{\ln(1/\gamma)}{\alpha} \right), \quad (24)$$

which gradually converges to 0 as $\kappa \rightarrow \infty$.

Now we substitute Eq. (24) into the original objective function, Eq. (22), to get

$$\text{ERUB}^{\text{uniform}} = \frac{1}{1 + \ln(1/\gamma)/\alpha} \left[1 + \kappa \ln \left(1 + \frac{\ln(1/\gamma)}{\alpha} \right) \right]. \quad (25)$$

Notice that

$$\ln(1/\gamma) = \ln(\kappa/(\kappa - 1)) = \ln(1 + 1/(\kappa - 1)) \leq \frac{1}{\kappa - 1} \leq \frac{1}{c\kappa}$$

because $\kappa \geq \frac{1}{1-c} > 1$ for some constant $c \in (0, 1)$. In addition,

$$\ln(1/\gamma) = -\ln(1 - 1/\kappa) \geq 1/\kappa.$$

Now, we can get the upper bound of Eq. (25) as

$$\begin{aligned} \text{ERUB}^{\text{uniform}} &\leq \frac{\kappa}{\kappa + 1/\alpha} \left[1 + \kappa \ln \left(1 + \frac{1}{c\kappa\alpha} \right) \right] \\ &\leq c_1 \frac{\kappa}{\kappa + 1/\alpha} \kappa \left[\ln \left(1 + \frac{1}{\kappa\alpha} \right) + \ln \left(\frac{1}{c} \right) \right] \\ &\leq c_1 c_2 \frac{\kappa^2}{\kappa + 1/\alpha} \ln \left(1 + \frac{1}{\kappa\alpha} \right) \end{aligned}$$

for some constants c_1, c_2 and large enough $\frac{1}{\alpha}$. Also, we can get the lower bound

$$\text{ERUB}^{\text{uniform}} \geq \frac{c\kappa}{c\kappa + 1/\alpha} \left[1 + \kappa \ln \left(1 + \frac{1}{\kappa\alpha} \right) \right] \geq c \frac{\kappa^2}{\kappa + 1/\alpha} \ln \left(1 + \frac{1}{\kappa\alpha} \right).$$

where we use the condition $c \in (0, 1)$. Thus, $\text{ERUB}^{\text{uniform}} = \Theta \left(\frac{\kappa^2}{\kappa + 1/\alpha} \ln \left(1 + \frac{1}{\kappa\alpha} \right) \right)$. □

PROOF OF LEMMA 4.3. By $\sum_{t=1}^T \sigma^{-2} = R$ and Cauchy-Schwarz inequality, we can derive the achievable lower bound as

$$R \sum_t q_t \sigma_t^2 = \sum_t \frac{1}{\sigma_t^2} \sum_t q_t \sigma_t^2 \geq \left(\sum_{t=1}^T \sqrt{q_t} \right)^2$$

where the inequality becomes equality if and only if $s/\sigma_t^2 = q_t \sigma_t^2$, i.e., $\sigma_t = (s/q_t)^{1/4}$, for some positive constant s . The equality $\sum_{t=1}^T \sigma_t^{-2} = R$ immediately suggests $\sqrt{s} = \frac{1}{R} \sum_{t=1}^T \sqrt{q_t}$. Thus, we get the σ_t .

Notice

$$T \sum_{t=1}^T q_t - \left(\sum_{t=1}^T \sqrt{q_t} \right)^2 = T^2 \frac{1}{T} \sum_{t=1}^T \left(\sqrt{q_t} - \frac{1}{T} \sum_{i=1}^T \sqrt{q_i} \right)^2 = T^2 \text{Var}[q_t] \quad (26)$$

where the variance is w.r.t. t . □

PROOF OF THEOREM 4.4. The upper bound of Eq. (7) can be written as

$$\begin{aligned} \text{ERUB}^{\text{dyn}} &= \gamma^T + \sum_{t=1}^T \gamma^{T-t} \alpha R \sigma_t^2 \\ &= \gamma^T + \alpha \left(\sum_{t=1}^T \sqrt{\gamma^{T-t}} \right)^2 \\ &= \gamma^T + \alpha \left(\frac{1 - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right)^2 \end{aligned}$$

where we make use of Lemma 4.3. Then, the minimizer of the ERUB is

$$\begin{aligned} T^* &= \arg \min_T \gamma^T + \alpha \left(\frac{1 - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right)^2 \\ &= 2 \log_{\gamma} \left(\frac{\alpha}{\alpha + (1 - \sqrt{\gamma})^2} \right). \end{aligned} \quad (27)$$

We can substitute Eq. (27) into ERUB^{dyn} to get

$$\begin{aligned} \text{ERUB}_{\min}^{\text{dyn}} &= \left(\frac{\alpha}{\alpha + (1 - \sqrt{\gamma})^2} \right)^2 + \alpha \left(\frac{1}{1 - \sqrt{\gamma}} \right)^2 \left(1 - \frac{\alpha}{\alpha + (1 - \sqrt{\gamma})^2} \right)^2 \\ &= \left(\frac{\alpha(1 - \sqrt{\gamma})^{-2}}{\alpha(1 - \sqrt{\gamma})^{-2} + 1} \right)^2 + \frac{\alpha(1 - \sqrt{\gamma})^{-2}}{(\alpha(1 - \sqrt{\gamma})^{-2} + 1)^2} \\ &= \frac{\alpha(1 - \sqrt{\gamma})^{-2}}{\alpha(1 - \sqrt{\gamma})^{-2} + 1} \end{aligned}$$

Notice that $(1 - \sqrt{\gamma})^{-2} = \kappa^2 + \kappa^2 - \kappa + 2\kappa\sqrt{\kappa(\kappa - 1)} = \kappa(2\kappa - 1 + 2\sqrt{\kappa(\kappa - 1)})$ and it is bounded by

$$\begin{aligned} \kappa(2\kappa - 1 + 2\sqrt{\kappa(\kappa - 1)}) &\leq 4\kappa^2, \\ \kappa(2\kappa - 1 + 2\sqrt{\kappa(\kappa - 1)}) &\geq \kappa(2\kappa - (3\kappa - 2) + 2\sqrt{(\kappa - 1)(\kappa - 1)}) = \kappa(-\kappa + 2 + 2\kappa - 2) = \kappa^2. \end{aligned}$$

Therefore, $\kappa \leq (1 - \sqrt{\gamma})^{-1} \leq 2\kappa$, with which we can derive

$$\begin{aligned} \text{ERUB}_{\min}^{\text{dyn}} &\leq 4 \frac{\kappa^2 \alpha}{\kappa^2 \alpha + 1}, \\ \text{ERUB}_{\min}^{\text{dyn}} &\geq \frac{\kappa^2 \alpha}{4\kappa^2 \alpha + 1} \geq \frac{1}{4} \frac{\kappa^2 \alpha}{\kappa^2 \alpha + 1}. \end{aligned}$$

Thus, $\text{ERUB}_{\min}^{\text{dyn}} = \Theta\left(\frac{\kappa^2 \alpha}{\kappa^2 \alpha + 1}\right)$. \square

D.2 Gradient Descents with Momentum

PROOF OF THEOREM 4.5. Without loss of generality, we absorb the $C\sigma_t/N$ into the variance of v_t such that $v_t \sim \mathcal{N}(0, \frac{C\sigma_t^2}{N}I)$ and $g_t \leftarrow \nabla_t + v_t$. Define $b_t = 1 - \beta^t$.

By smoothness and Eq. (1), we have

$$\begin{aligned} f(\theta_{t+1}) - f(\theta_t) &\leq \nabla_t^\top (\theta_{t+1} - \theta_t) + \frac{1}{2} M \|\theta_{t+1} - \theta_t\|^2 \\ &= -\frac{\eta_t}{b_t^2} b_t \nabla_t^\top v_{t+1} + \frac{1}{2} M \frac{\eta_t^2}{b_t^2} \|v_{t+1}\|^2 \\ &= \frac{\eta_t}{b_t^2} \left(\|b_t \nabla_t - v_{t+1}\|^2 - \|b_t \nabla_t\|^2 - \|v_{t+1}\|^2 \right) + \frac{1}{2} M \frac{\eta_t^2}{b_t^2} \|v_{t+1}\|^2 \\ &= \frac{\eta_t}{b_t^2} \underbrace{\|b_t \nabla_t - v_{t+1}\|^2}_{U_1(t)} - \eta_t \|\nabla_t\|^2 - \frac{\eta_t}{b_t^2} \left(1 - \frac{1}{2} M \eta_t \right) \|v_{t+1}\|^2, \end{aligned} \quad (28)$$

where only the $U_1(t)$ is non-negative. Specifically, $U_1(t)$ describes the difference between current gradient and the average. We can expand v_{t+1} to get an upper bound:

$$\begin{aligned}
U_1(t) &= \|b_t \nabla_t - v_{t+1}\|^2 \\
&= \left\| (1-\beta) \sum_{i=1}^t \beta^{t-i} \nabla_t - (1-\beta) \sum_{i=1}^t \beta^{t-i} g_i \right\|^2 \\
&= (1-\beta)^2 \left\| \sum_{i=1}^t \beta^{t-i} (\nabla_t - g_i) \right\|^2 \\
&= (1-\beta)^2 \left\| \sum_{i=1}^t \beta^{t-i} (\nabla_t - \nabla_i) + \sum_{i=1}^t \beta^{t-i} (\nabla_i - g_i) \right\|^2 \\
&\leq 2(1-\beta)^2 \left[\left\| \sum_{i=1}^t \beta^{t-i} (\nabla_t - \nabla_i) \right\|^2 + \left\| \sum_{i=1}^t \beta^{t-i} (\nabla_i - g_i) \right\|^2 \right] \\
&\leq 2(1-\beta) \left[\underbrace{b_t \sum_{i=1}^t \beta^{t-i} \|\nabla_t - \nabla_i\|^2}_{U_2(t) \text{ (gradient variance)}} + (1-\beta) \underbrace{\left\| \sum_{i=1}^t \beta^{t-i} v_i \right\|^2}_{\text{noise variance}} \right]
\end{aligned}$$

where we use $\|x + y\|^2 \leq (\|x\| + \|y\|)^2 \leq 2(\|x\|^2 + \|y\|^2)$. The last inequality can be proved by Cauchy-Schwartz inequality for each coordinate.

We plug the $U_1(t)$ into Eq. (28) and use the PL condition to get

$$\begin{aligned}
f(\theta_{t+1}) - f(\theta_t) &\leq \frac{\eta_t}{b_t^2} U_1(t) - \eta_t \|\nabla_t\|^2 - \frac{\eta_t}{b_t^2} (1 - \frac{1}{2} M \eta_t) \|v_{t+1}\|^2 \\
&\leq -\eta_t \|\nabla_t\|^2 + \frac{\eta_t}{b_t^2} 2(1-\beta) \left[b_t U_2(t) + (1-\beta) \left\| \sum_{i=1}^t \beta^{t-i} v_i \right\|^2 \right] \\
&\quad - \frac{\eta_t}{b_t^2} (1 - \frac{1}{2} M \eta_t) \|v_{t+1}\|^2 \\
&\leq -2\mu \eta_t (f(\theta_t) - f(\theta^*)) + \frac{2(1-\beta)\eta_t}{b_t} U_2(t) + \frac{2(1-\beta)^2 \eta_t}{b_t^2} \left\| \sum_{i=1}^t \beta^{t-i} v_i \right\|^2 \\
&\quad - \frac{\eta_t}{b_t^2} (1 - \frac{1}{2} M \eta_t) \|v_{t+1}\|^2.
\end{aligned}$$

Rearranging terms and taking expectation to show

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1})] - f(\theta^*) &\leq \gamma (\mathbb{E}[f(\theta_t)] - f(\theta^*)) + \frac{2(1-\beta)^2 \eta_t}{b_t^2} \sum_{i=1}^t \beta^{t-i} \mathbb{E} \|v_i\|^2 \\
&\quad + \frac{2(1-\beta)\eta_t}{b_t} \mathbb{E}[U_2(t)] - \frac{\eta_t}{b_t^2} (1 - \frac{1}{2} M \eta_t) \mathbb{E} \|v_{t+1}\|^2 \\
&= \gamma (\mathbb{E}[f(\theta_t)] - f(\theta^*)) + \frac{2(1-\beta)^2 \eta_t}{b_t^2} \sum_{i=1}^t \beta^{2(t-i)} \frac{C^2 D \sigma_t^2}{N^2} \\
&\quad + \frac{2(1-\beta)\eta_t}{b_t} \mathbb{E}[U_2(t)] - \frac{\eta_t}{b_t^2} (1 - \frac{1}{2} M \eta_t) \mathbb{E} \|v_{t+1}\|^2
\end{aligned}$$

where $\gamma = 1 - \eta_0/\kappa = 1 - 2\mu\eta_t$. The recursive inequality implies

$$\begin{aligned} \mathbb{E}[f(\theta_{T+1})] - f(\theta^*) &\leq \gamma^T (f(\theta_1) - f(\theta^*)) + \sum_{t=1}^T \gamma^{T-t} \frac{2(1-\beta)^2 \eta_t}{b_t^2} \sum_{i=1}^t \beta^{2(t-i)} \frac{C^2 D \sigma_t^2}{N^2} \\ &\quad + \sum_{t=1}^T \gamma^{T-t} \frac{2(1-\beta)\eta_t}{b_t} \mathbb{E}[U_2(t)] - \sum_{t=1}^T \gamma^{T-t} \frac{\eta_t}{b_t^2} \left(1 - \frac{1}{2} M \eta_t\right) \mathbb{E} \|v_{t+1}\|^2 \\ &= \left(\gamma^T + 2\eta_0 \alpha R \underbrace{\sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{b_t^2} \sum_{i=1}^t \beta^{2(t-i)} \sigma_t^2}_{U_3} \right) (f(\theta_1) - f(\theta^*)) \\ &\quad + \underbrace{\sum_{t=1}^T \gamma^{T-t} \frac{2(1-\beta)\eta_t}{b_t} \mathbb{E}[U_2(t)] - \sum_{t=1}^T \gamma^{T-t} \frac{\eta_t}{b_t^2} \left(1 - \frac{1}{2} M \eta_t\right) \mathbb{E} \|v_{t+1}\|^2}_{U_4(t)}. \end{aligned}$$

where we utilize $\alpha = \frac{DC^2}{2MN^2R} \frac{1}{f(\theta_1) - f(\theta^*)}$ and $\eta_t = \frac{\eta_0}{2M}$.

By Lemma C.5, we have

$$\sum_{t=1}^T \gamma^{T-t} \frac{2(1-\beta)\eta_t}{b_t} U_2(t) \leq \frac{\eta_0^3 \beta \gamma}{2M(1-\beta)^3(\gamma-\beta)^2} \sum_{i=1}^{T-1} \gamma^{T-i} \|v_{i+1}\|^2.$$

Thus, by $\frac{1}{b_t} \geq 1$,

$$\begin{aligned} U_4(t) &\leq \frac{\eta_0^3 \beta \gamma}{2M(1-\beta)^3(\gamma-\beta)^2} \sum_{i=1}^{T-1} \gamma^{T-i} \mathbb{E} \|v_{i+1}\|^2 - \frac{\eta_0}{2M} \left(1 - \frac{\eta_0}{4}\right) \sum_{t=1}^T \gamma^{T-t} \mathbb{E} \|v_{t+1}\|^2 \\ &= -\frac{\eta_0}{2M} \zeta \sum_{t=1}^T \gamma^{T-t} \mathbb{E} \|v_{t+1}\|^2 \end{aligned}$$

where

$$\zeta = 1 - \frac{1}{4} \eta_0 - \frac{\beta \gamma}{(\gamma-\beta)^2(1-\beta)^3} \eta_0^2 = 1 - \frac{1}{4} \eta_0 - \frac{\beta/\gamma}{(1-\beta/\gamma)^2(1-\beta)^3} \eta_0^2$$

When a small enough η_0 , e.g., Specifically,

$$\begin{aligned} \eta_0 &\leq \frac{(\gamma-\beta)^2(1-\beta)^3}{8\beta\gamma} \left[\sqrt{1 + \frac{64\beta\gamma}{(\gamma-\beta)^2(1-\beta)^3} - 1} \right] \\ &= \frac{8}{\sqrt{1 + 64\beta\gamma(\gamma-\beta)^{-2}(1-\beta)^{-3} + 1}} \end{aligned}$$

We can have $\zeta \geq 0$.

By the definition of $U_3(T, \sigma)$, we can get

$$\mathbb{E}[f(\theta_{T+1})] - f(\theta^*) \leq \left(\gamma^T + 2\eta_0 \alpha R U_3(T, \sigma) \right) (f(\theta_1) - f(\theta^*)) - \frac{\eta_0}{2M} \zeta \sum_{t=1}^T \gamma^{T-t} \mathbb{E} \|v_{t+1}\|^2.$$

□

PROOF OF THEOREM 4.6. Since σ_t is static, by definition of U_3 in Theorem 4.5,

$$\begin{aligned} U_3 &= \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \sigma^2 \\ &= \sigma^2 \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \\ &= \sigma^2 \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \frac{1-\beta^{2t}}{1-\beta^2} \\ &= \sigma^2 \sum_{t=1}^T \gamma^{T-t} \frac{1-\beta}{1-\beta^t} \frac{1+\beta^t}{1+\beta}. \end{aligned}$$

Because $\frac{1-\beta}{1-\beta^t} \frac{1+\beta^t}{1+\beta} \leq 1$, the U_3 will be smaller than the corresponding summation in GD with uniform schedule.

By Lemma C.7, when $T > \hat{T}$, we can rewrite U_3 as

$$\begin{aligned}
U_3 &\leq \sigma^2 \sum_{t=1}^T \gamma^{T-t} \frac{1-\beta}{1-\beta^t} \\
&= \sigma^2 \sum_{t=1}^{\hat{T}} \gamma^{T-t} \frac{1-\beta}{1-\beta^t} + \sigma^2 \sum_{t=\hat{T}+1}^T \gamma^{T-t} \frac{1-\beta}{1-\beta^t} \\
&\leq \sigma^2 \sum_{t=1}^{\hat{T}} \gamma^{T-t} \gamma^{t-1} + \sigma^2 \sum_{t=\hat{T}+1}^T \gamma^{T-t} \gamma^{\hat{T}-1} \\
&= \sigma^2 \gamma^{T-1} \hat{T} + \sigma^2 \gamma^{\hat{T}-1} \sum_{t=1}^{T-\hat{T}} \gamma^{T-\hat{T}-t} \\
&= \sigma^2 \gamma^{T-1} \hat{T} + \sigma^2 \frac{\gamma^{\hat{T}-1} - \gamma^{T-1}}{1-\gamma} \\
&= \frac{T}{\gamma R} \gamma^T \left(\hat{T} + \frac{\gamma^{\hat{T}-T} - 1}{1-\gamma} \right)
\end{aligned}$$

where we use $\sigma^2 = T/R$ in the last line. Without assuming $T > \hat{T}$, we can generally write the upper bound as

$$U_3 \leq \frac{T}{\gamma R} \gamma^T \left(\min\{\hat{T}, T\} + \max\left\{\frac{\gamma^{\hat{T}-T} - 1}{1-\gamma}, 0\right\} \right).$$

By Theorem 4.5, because $\zeta \geq 0$, we have

$$\begin{aligned}
\text{ERUB} &\leq \gamma^T + 2R\eta_0\alpha U_3 \\
&= \gamma^T \left(1 + \frac{\alpha'}{\gamma} T \left(\min\{\hat{T}, T\} + \max\left\{\frac{\gamma^{\hat{T}-T} - 1}{1-\gamma}, 0\right\} \right) \right)
\end{aligned}$$

where $\alpha' = 2\eta_0\alpha$.

First, we consider $T \leq \hat{T}$. Use $T = \frac{1}{\ln(1/\gamma)} \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right) = \left\lceil \mathcal{O}\left(\frac{\kappa}{\eta_0} \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right)\right) \right\rceil$ to get

$$\begin{aligned}
\text{ERUB} &\leq \left(\frac{\alpha}{\alpha + \eta_0/\kappa} \right) \left(1 + \alpha' \gamma^{-1} \left(\frac{2}{\ln(1/\gamma)} \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right)^2 \right) \\
&\leq \left(\frac{\alpha}{\alpha + \eta_0/\kappa} \right) \left(1 + \frac{8\kappa^2\alpha}{\eta_0\gamma} \ln^2\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right) \\
&\leq \mathcal{O}\left(\frac{\kappa}{\kappa + \eta_0/\alpha} \left(1 + \frac{8\kappa^2\alpha}{\eta_0\gamma} \ln^2\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right) \right) \\
&= \mathcal{O}\left(\frac{\kappa}{\kappa + \eta_0/\alpha} \left(1 + \frac{4\kappa}{\gamma} \right) \right) \\
&= \mathcal{O}\left(\frac{\kappa^2}{\kappa + \eta_0/\alpha} \right)
\end{aligned}$$

where we used $\ln(1/\gamma) \geq \eta_0/\kappa$ and $\ln(1+x) \leq \sqrt{x}$ for any $x > 0$.

Second, when $T > \hat{T}$,

$$\begin{aligned}
\text{ERUB} &\leq \gamma^T \left(1 + \frac{\alpha'}{\gamma} T \left(\hat{T} + \frac{\gamma^{\hat{T}-T} - 1}{1-\gamma} \right) \right) \\
&\leq \mathcal{O}\left(\gamma^T + \frac{2\alpha'}{\gamma} T \kappa (\gamma^{\hat{T}} - \gamma^T) \right).
\end{aligned}$$

Make use of $T = \left\lceil \frac{1}{\ln(1/\gamma)} \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right\rceil$ to show

$$\begin{aligned}
\text{ERUB} &\leq \mathcal{O}\left(\frac{\kappa}{\kappa + \eta_0/\alpha} + \frac{4\kappa^2\alpha}{\eta_0\gamma} (\gamma^{\hat{T}} - \frac{\kappa}{\kappa + \eta_0/\alpha}) \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right) \\
&\leq \mathcal{O}\left(\frac{\kappa^2}{\kappa + \eta_0/\alpha} \gamma^{\hat{T}-1} \ln\left(1 + \frac{\eta_0}{\kappa\alpha}\right) \right).
\end{aligned}$$

□

PROOF OF THEOREM 4.7. By Lemma C.7, we can rewrite U_3 as

$$\begin{aligned} U_3 &= \sum_{t=1}^T Y^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &\leq \sum_{t=1}^{\hat{T}} Y^{T-t} Y^{2(t-1)} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 + \sum_{t=\hat{T}+1}^T Y^{T-t} Y^{2(\hat{T}-1)} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &\leq Y^{T-\hat{T}} \underbrace{\sum_{t=1}^{\hat{T}} Y^{\hat{T}-t} Y^{2(t-1)} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2}_{V_1} + Y^{2(\hat{T}-1)} \underbrace{\sum_{t=\hat{T}+1}^T Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2}_{V_2} \end{aligned}$$

We derive V_1 and V_2 separately.

For V_1 , we can obtain the upper bound by

$$\begin{aligned} V_1 &= \sum_{t=1}^{\hat{T}} Y^{\hat{T}-t} Y^{2(t-1)} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= Y^{\hat{T}-2} \sum_{t=1}^{\hat{T}} Y^t \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= Y^{\hat{T}-2} \sum_{i=1}^{\hat{T}} \beta^{-2i} \sigma_i^2 \sum_{t=i}^{\hat{T}} (\gamma \beta^2)^t \\ &= Y^{\hat{T}-2} \sum_{i=1}^{\hat{T}} \beta^{-2i} \sigma_i^2 \frac{(\gamma \beta^2)^i - (\gamma \beta^2)^{\hat{T}+1}}{1 - \gamma \beta^2} \\ &= Y^{2\hat{T}-3} \sum_{i=1}^{\hat{T}} \frac{\gamma^{i-\hat{T}-1} - \beta^{2(\hat{T}+1-i)}}{1 - \gamma \beta^2} \sigma_i^2 \\ &= Y^{2\hat{T}-3} \sum_{i=1}^{\hat{T}} \frac{1 - (\gamma \beta^2)^{\hat{T}+1-i}}{1 - \gamma \beta^2} \gamma^{i-\hat{T}-1} \sigma_i^2 \\ &\leq \frac{Y^{\hat{T}}}{\gamma^2(1 - \gamma \beta^2)} \sum_{i=1}^{\hat{T}} \gamma^i \sigma_i^2 \\ &\leq \frac{Y^{\hat{T}}}{\gamma(\gamma - \beta^2)} \sum_{i=1}^{\hat{T}} \gamma^i \sigma_i^2 \end{aligned}$$

For V_2 , we can derive

$$\begin{aligned} V_2 &= \sum_{t=\hat{T}+1}^T Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= \sum_{t=1}^T Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 - \sum_{t=1}^{\hat{T}} Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= \sum_{t=1}^T Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 - Y^{T-\hat{T}} \sum_{t=1}^{\hat{T}} Y^{\hat{T}-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2. \end{aligned}$$

We first consider the first term

$$\begin{aligned} &\sum_{t=1}^T Y^{T-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= \sum_{i=1}^T \sigma_i^2 \sum_{t=i}^T Y^{T-t} \beta^{2(t-i)} \\ &= \sum_{i=1}^T Y^T \beta^{-2i} \sigma_i^2 \sum_{t=i}^T Y^{-t} \beta^{2t} \\ &= \sum_{i=1}^T Y^T \beta^{-2i} \sigma_i^2 \frac{(\beta^2/\gamma)^i - (\beta^2/\gamma)^{T+1}}{1 - (\beta^2/\gamma)} \\ &= \sum_{i=1}^T \frac{Y^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2. \end{aligned}$$

Similarly, we have

$$\begin{aligned} & \gamma^{T-\hat{T}} \sum_{t=1}^{\hat{T}} \gamma^{\hat{T}-t} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2 \\ &= \gamma^{T-\hat{T}} \sum_{i=1}^{\hat{T}} \frac{\gamma^{\hat{T}+1-i} - \beta^{2(\hat{T}+1-i)}}{\gamma - \beta^2} \sigma_i^2 \\ &= \sum_{i=1}^{\hat{T}} \frac{\gamma^{T+1-i} - \gamma^{T-\hat{T}} \beta^{2(\hat{T}+1-i)}}{\gamma - \beta^2} \sigma_i^2. \end{aligned}$$

Thus,

$$\begin{aligned} V_2 &= \sum_{i=1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 - \sum_{i=1}^{\hat{T}} \frac{\gamma^{T+1-i} - \gamma^{T-\hat{T}} \beta^{2(\hat{T}+1-i)}}{\gamma - \beta^2} \sigma_i^2 \\ &= \sum_{i=\hat{T}+1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 + \sum_{i=1}^{\hat{T}} \frac{\gamma^{T-\hat{T}} - \beta^{2(T-\hat{T})}}{\gamma - \beta^2} \beta^{2(\hat{T}+1-i)} \sigma_i^2 \\ &\leq \sum_{i=\hat{T}+1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 + \sum_{i=1}^{\hat{T}} \frac{\gamma^{T-\hat{T}}}{\gamma - \beta^2} \beta^{2(\hat{T}+1-i)} \sigma_i^2. \end{aligned}$$

Substitute V_1 and V_2 into U_3 to get

$$\begin{aligned} U_3 &\leq \gamma^T \frac{1}{\gamma(\gamma - \beta^2)} \sum_{i=1}^{\hat{T}} \gamma^i \sigma_i^2 + \gamma^{2\hat{T}-2} \sum_{i=\hat{T}+1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 \\ &\quad + \sum_{i=1}^{\hat{T}} \frac{\gamma^{T+\hat{T}-2}}{\gamma - \beta^2} \beta^{2(\hat{T}+1-i)} \sigma_i^2 \\ &\leq \left(\frac{\gamma^T}{\gamma(\gamma - \beta^2)} \sum_{i=1}^{\hat{T}} (\gamma^i + \gamma^{\hat{T}-1} \beta^{2(\hat{T}+1-i)}) \sigma_i^2 + \gamma^{2\hat{T}-2} \sum_{i=\hat{T}+1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 \right) \\ &\leq \left(\frac{2\gamma^T}{\gamma(\gamma - \beta^2)} \sum_{i=1}^{\hat{T}} \gamma^i \sigma_i^2 + \gamma^{2\hat{T}-2} \sum_{i=\hat{T}+1}^T \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \sigma_i^2 \right) \\ &= \sum_{t=1}^T q_t \sigma_t^2 \end{aligned}$$

where

$$\begin{aligned} q_t &= \frac{2}{\gamma(\gamma - \beta^2)} \gamma^{T+t} \mathbb{1}_{T \leq \hat{T}} + \gamma^{2(\hat{T}-1)} \frac{\gamma^{T+1-i} - \beta^{2(T+1-i)}}{\gamma - \beta^2} \gamma^{T-t} \mathbb{1}_{T > \hat{T}} \\ &\leq c_1 \gamma^{T+t} \mathbb{1}_{T \leq \hat{T}} + \gamma^{\hat{T}-1} c_2 \gamma^{T-t} \mathbb{1}_{T > \hat{T}} \end{aligned}$$

where $c_1 = \frac{2}{\gamma(\gamma - \beta^2)}$ and $c_2 = \frac{\gamma^{2\hat{T}}}{\gamma - \beta^2}$.

When $T > \hat{T}$, by [Lemma 4.3](#), the lower bound of $R \sum_{t=1}^T q_t \sigma_t^2$ is

$$\begin{aligned} \left(\sum_{t=1}^T \sqrt{q_t} \right)^2 &= \gamma^T \left(\sum_{t=1}^{\hat{T}} \sqrt{c_1} \gamma^t + \sum_{t=\hat{T}+1}^T \sqrt{\gamma^{\hat{T}-1} c_2} \gamma^{-t} \right)^2 \\ &= \gamma^T \left(\sqrt{c_1} \gamma \frac{1 - \gamma^{\hat{T}/2}}{1 - \sqrt{\gamma}} + \sqrt{c_2} \frac{1 - \gamma^{(\hat{T}-T)/2}}{\sqrt{\gamma} - 1} \right)^2 \\ &= \gamma^T \left(\sqrt{c_1} \gamma \frac{1 - \gamma^{\hat{T}/2}}{1 - \sqrt{\gamma}} + \sqrt{c_2} \frac{\gamma^{(\hat{T}-T)/2} - 1}{1 - \sqrt{\gamma}} \right)^2 \\ &\leq \mathcal{O} \left(c_2 \left\{ \frac{\gamma^{(\hat{T}-1)/2} - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right\}^2 \right) \end{aligned}$$

which is achieved when

$$\sigma_t^2 = \frac{1}{R} \sum_{i=1}^T \sqrt{\frac{q_i}{q_t}}.$$

By [Theorem 4.5](#), because $\zeta \geq 0$, we have

$$\begin{aligned} \text{ERUB} &\leq \gamma^T + 2R\eta_0\alpha U_3 \\ &= \gamma^T + 2\eta_0\alpha \sum_{t=1}^T Rq_t\sigma_t^2. \end{aligned}$$

And the minimum of the upper bound is

$$\text{ERUB}_{\min} = \gamma^T + \alpha' \mathcal{O} \left(\left\{ \frac{\gamma^{(\hat{T}-1)/2} - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right\}^2 \right)$$

where $\alpha' = 2\eta_0 c_2 \alpha$. Let $T = \frac{2}{\ln(1/\gamma)} \ln \left(1 + \frac{\eta_0}{\kappa\alpha} \right)$. Then,

$$\begin{aligned} \text{ERUB}_{\min} &= \mathcal{O} \left(\left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right)^2 + \frac{\alpha'}{(1 - \sqrt{\gamma})^2} \left\{ \frac{\gamma^{(\hat{T}-1)/2} - (1 - \gamma^{(\hat{T}-1)/2})\kappa\alpha}{\kappa\alpha + \eta_0} \right\}^2 \right) \\ &\leq \mathcal{O} \left(\left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right)^2 + \frac{2\eta_0 c_2 \alpha}{(1 - \sqrt{\gamma})^2} \left\{ \frac{\gamma^{(\hat{T}-1)/2}}{\kappa\alpha + \eta_0} \right\}^2 \right) \\ &= \mathcal{O} \left(\frac{\kappa\alpha}{(\kappa\alpha + \eta_0)^2} \left(\kappa\alpha + \frac{2\eta_0 c_2 / \kappa}{(1 - \sqrt{\gamma})^2} \gamma^{(\hat{T}-1)} \right) \right) \\ &= \mathcal{O} \left(\frac{\kappa\alpha}{(\kappa\alpha + \eta_0)^2} (\kappa\alpha + c_3 \eta_0) \right) \\ &\leq \mathcal{O} \left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right) \end{aligned}$$

where c_3 is some constant.

When $T \leq \hat{T}$,

$$\begin{aligned} U_3 &\leq \gamma^{T-T} \underbrace{\sum_{t=1}^T \gamma^{T-t} \gamma^{2(t-1)} \sum_{i=1}^t \beta^{2(t-i)} \sigma_i^2}_{V_1} \\ &\leq \frac{\gamma^{T-2}}{1 - \gamma\beta^2} \sum_{i=1}^T \gamma^i \sigma_i^2 \end{aligned}$$

with which we obtain

$$\begin{aligned} \text{ERUB} &\leq \gamma^T + 2R\eta_0\alpha U_3 \\ &\leq \gamma^T + 2\eta_0\alpha \frac{\gamma^{-2}}{1 - \gamma\beta^2} \sum_{t=1}^T Rq_t\sigma_t^2. \end{aligned}$$

where we let $q_t = \gamma^{T+t}$. By [Lemma 4.3](#),

$$\begin{aligned} \sum_{i=1}^T Rq_t\sigma_i^2 &\geq \left(\sum_{t=1}^T \sqrt{q_t} \right)^2 \\ &= \gamma^T \left(\sum_{t=1}^T \gamma^{t/2} \right)^2 \\ &= \gamma^{T+1} \left(\frac{1 - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right)^2. \end{aligned}$$

Thus,

$$\begin{aligned} \text{ERUB}_{\min} &\leq \gamma^T + 2\eta_0\alpha \frac{\gamma^{T-1}}{1 - \gamma\beta^2} \left(\frac{1 - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right)^2 \\ &= \gamma^T \left(1 + 2\eta_0\gamma c_1 \alpha \left(\frac{1 - \gamma^{T/2}}{1 - \sqrt{\gamma}} \right)^2 \right) \end{aligned}$$

Let $T = \left\lceil \frac{2}{\ln(1/\gamma)} \ln \left(1 + \frac{\eta_0}{\kappa\alpha} \right) \right\rceil$. Then,

$$\begin{aligned} \text{ERUB}_{\min} &\leq \left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right)^2 \left(1 + \frac{2\eta_0\gamma c_1\alpha}{(1-\sqrt{\gamma})^2} \left(\frac{1}{\kappa\alpha + 1} \right)^2 \right) \\ &\leq \left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right)^2 \left(1 + \mathcal{O}\left(\frac{1}{\kappa\alpha + 1} \right) \right) \\ &\leq \mathcal{O}\left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \right)^2. \end{aligned}$$

In summary,

$$\text{ERUB}_{\min} \leq \mathcal{O}\left(\frac{\kappa\alpha}{\kappa\alpha + \eta_0} \left(\mathbb{I}_{T \leq \hat{T}} \frac{\kappa\alpha}{\kappa\alpha + \eta_0} + \mathbb{I}_{T > \hat{T}} \right) \right)$$

□

D.3 Stochastic Gradient Descents

PROOF OF THEOREM 4.8. Let $\tilde{\nabla}_t$ be the stochastic gradient of the step t . By the smoothness, we have

$$\begin{aligned} f(\theta_{t+1}) - f(\theta_t) &\leq -\eta_t \nabla_t^\top (\tilde{\nabla}_t + G\sigma_t v_t/n) + \frac{1}{2} M\eta_t^2 \|\tilde{\nabla}_t + G\sigma_t v_t/n\|^2 \\ &= -\eta_t \nabla_t^\top (\nabla_t + \sigma_g \xi_t/n + G\sigma_t v_t/n) + \frac{1}{2} M\eta_t^2 \|\nabla_t + \sigma_g \xi_t/n + G\sigma_t v_t/n\|^2. \end{aligned}$$

Note that $\mathbb{E}(\sigma_g \xi_t/n + G\sigma_t v_t/n) = 0$ and $\mathbb{E}(\sigma_g \xi_t/n + G\sigma_t v_t/n)^2 = \sigma_g^2 + (G\sigma_t/n)^2$. Without loss of generality, we can write $\sigma_g \xi_t + G\sigma_t v_t$ as $\tilde{\sigma}_t \zeta_t$ where $\tilde{\sigma}_t \triangleq \sqrt{\sigma_g^2 + (G\sigma_t)^2}$ and ζ_t is a random vector with $\mathbb{E}\zeta_t = 0$ and $\mathbb{E}\|\zeta_t\|^2 \leq D$. Therefore,

$$\begin{aligned} f(\theta_{t+1}) - f(\theta_t) &\leq -\eta_t \nabla_t^\top (\nabla_t + \tilde{\sigma}_t \zeta_t/n) + \frac{1}{2} M\eta_t^2 \|\nabla_t + \tilde{\sigma}_t \zeta_t/n\|^2 \\ &= -\eta_t \left(1 - \frac{1}{2} M\eta_t \right) \|\nabla_t\|^2 - (1 - M\eta_t) \eta_t \nabla_t^\top \tilde{\sigma}_t \zeta_t/n + \frac{1}{2} M\eta_t^2 \|\tilde{\sigma}_t \zeta_t/n\|^2 \\ &\leq -2\mu\eta_t \left(1 - \frac{1}{2} M\eta_t \right) (f(\theta_t) - f(\theta^*)) - (1 - M\eta_t) \eta_t \nabla_t^\top \tilde{\sigma}_t \zeta_t/n \\ &\quad + \frac{1}{2} M\eta_t^2 \|\tilde{\sigma}_t \zeta_t/n\|^2. \end{aligned}$$

Then following the same proof of Theorem 4.1, we can get

$$\begin{aligned} \mathbb{E}[f(\theta_{T+1})] - f(\theta^*) &\leq \gamma^T (\mathbb{E}[f(\theta_1)] - f(\theta^*)) + R' \sum_{t=1}^T \gamma^{T-t} \alpha \frac{1}{G^2} \tilde{\sigma}_t^2 (\mathbb{E}[f(\theta_1)] - f(\theta^*)) \\ &= \left[\gamma^T + R' \sum_{t=1}^T \gamma^{T-t} \alpha \left(\frac{1}{G^2} \sigma_g^2 + \sigma_t^2 \right) \right] (\mathbb{E}[f(\theta_1)] - f(\theta^*)) \\ &= \left[\gamma^T + R' \alpha \frac{1}{G^2} \sigma_g^2 \frac{1-\gamma^T}{1-\gamma} + R' \sum_{t=1}^T \gamma^{T-t} \alpha \sigma_t^2 \right] (\mathbb{E}[f(\theta_1)] - f(\theta^*)) \\ &\leq \left[\gamma^T + \frac{R'\kappa\alpha}{G^2} \sigma_g^2 + R' \sum_{t=1}^T \gamma^{T-t} \alpha \sigma_t^2 \right] (\mathbb{E}[f(\theta_1)] - f(\theta^*)). \end{aligned}$$

where $\frac{R'\kappa\alpha}{G^2} = \frac{D}{2\mu(f(\theta_1)-f(\theta^*))} \frac{1}{n^2} = \frac{D}{2\mu(f(\theta_1)-f(\theta^*))} \min\left\{ \frac{1}{N^2R}, 1 \right\} \leq \frac{D}{2\mu(f(\theta_1)-f(\theta^*))} \frac{1}{N^2R}$.

□

PROOF OF THEOREM 4.9. Without loss of generality, we can write $\sigma_g \xi_t + G\sigma_t v_t$ as $\tilde{\sigma}_t \zeta_t$ where $\tilde{\sigma}_t \triangleq \sqrt{\sigma_g^2 + (G\sigma_t)^2}$ and ζ_t is a random vector with $\mathbb{E}\zeta_t = 0$ and $\mathbb{E}\|\zeta_t\|^2 \leq D$. Therefore, we replace v_t by ζ_t and σ_t^2 by $\tilde{\sigma}_t^2/G^2 = \sigma_g^2/G^2 + \sigma_t^2$. Now, we only need to update $U_3(\sigma, T)$

as

$$\begin{aligned}\tilde{U}_3 &= \frac{1}{G^2} \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \tilde{\sigma}_i^2 \\ &= \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)} \left(\frac{1}{G^2} \sigma_g^2 + \sigma_t^2 \right) \\ &= U_3^g + U_3\end{aligned}$$

where we define

$$U_3^g \triangleq \frac{1}{G^2} \sigma_g^2 \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \sum_{i=1}^t \beta^{2(t-i)}.$$

We can upper bound U_3^g by

$$\begin{aligned}U_3^g &= \frac{1}{G^2} \sigma_g^2 \sum_{t=1}^T \gamma^{T-t} \frac{(1-\beta)^2}{(1-\beta^t)^2} \frac{1-\beta^{2t}}{1-\beta^2} \\ &= \frac{1}{G^2} \sigma_g^2 \sum_{t=1}^T \gamma^{T-t} \frac{1-\beta}{1-\beta^t} \frac{1+\beta^t}{1+\beta} \\ &\leq \frac{1}{G^2} \sigma_g^2 \sum_{t=1}^T \gamma^{T-t} \\ &\leq \frac{1}{G^2} \sigma_g^2 \frac{1}{1-\gamma} \\ &= \frac{1}{G^2} \kappa \sigma_g^2.\end{aligned}$$

Combine with the factors of U_3 in the PGD bounds:

$$\alpha R' U_3^g \leq \frac{\alpha R'}{G^2} \kappa \sigma_g^2 = \frac{\alpha R'}{G^2} \kappa \sigma_g^2 = \frac{D \sigma_g^2}{2\mu n^2 (f(\theta_1) - f(\theta^*))} \leq \frac{D \sigma_g^2}{2\mu N^2 R (f(\theta_1) - f(\theta^*))}.$$

□