
Efficient Privacy-Preserving Stochastic Nonconvex Optimization (Supplementary material)

Lingxiao Wang¹

Bargav Jayaraman²

David Evans²

Quanquan Gu³

¹Toyota Technological Institute at Chicago

²Department of Computer Science, University of Virginia

³Department of Computer Science, University of California, Los Angeles

1 ADDITIONAL EXPERIMENTS

In this section, we present additional experiment results on nonconvex logistic regression and convolutional neural networks.

1.1 RESULTS ON *IJCNN1* DATASET

In this subsection, we present the additional experiment of our method on *ijcnn1* dataset. In this dataset, we follow the same settings as before: we set the clipping thresholds $C_1 = 1, C_2 = 0.01$, and set the momentum parameter $\gamma = C_2$. Figure 1 illustrates the objective function value and the gradient norm of different algorithms under various privacy budgets $\epsilon \in \{0.2, 0.5\}$. We can see that our proposed algorithm (DP-SRM) outperforms the other three baseline algorithms (RRPSGD, DP-GD, and DP-AGD) in terms of the objective loss, gradient norm, and convergence rate by a large margin. Table 1 shows the test error of different algorithms as well as the CPU time (in seconds) of the training process on *ijcnn1* dataset. It demonstrates that our algorithm converges faster and can achieve a better test error on the test set than other baselines.

Table 1: Comparison of different algorithms on *ijcnn1* dataset under different privacy budgets $\epsilon \in \{0.2, 0.5\}$ and $\delta = 10^{-5}$. Note that the non-private baseline denotes the test error of the non-private STORM algorithm [Cutkosky and Orabona, 2019].

Privacy Budget	Non-private Baseline	Method	Test Error	Data Passes	CPU time	Gradient Norm
$\epsilon = 0.2$	0.2096 (0.002)	DP-GD	0.3160 (0.0120)	20	0.5180	0.0184 (0.0024)
		DP-AGD	0.2645 (0.0044)	346	90.05	0.0133 (0.0018)
		RRPSGD	0.3110 (0.0106)	8	47.64	0.0175 (0.0023)
		DP-SRM	0.2503 (0.0090)	4	0.4748	0.0117 (0.0008)
$\epsilon = 0.5$	0.2096 (0.002)	DP-GD	0.2717 (0.0081)	20	0.4990	0.0171 (0.0024)
		DP-AGD	0.2416 (0.0029)	365	94.28	0.0397 (0.0025)
		RRPSGD	0.3033 (0.0110)	10	59.06	0.0160 (0.0018)
		DP-SRM	0.2341 (0.0042)	5	0.4368	0.0082 (0.0005)

1.2 ADDITIONAL EXPERIMENTS ON CONVOLUTIONAL NEURAL NETWORKS

In this subsection, we present additional experiment results on training convolutional neural networks. Figure 2 shows the average test error (over 30 trials) and the corresponding 95% confidence interval of different methods versus the number of iterations as well as the training time under different privacy budgets on MNIST and CIFAR-10 datasets.

Results on MNIST dataset. We can see from Figure 2(a) and Figure 2(b) that our proposed method can achieve 2.91%

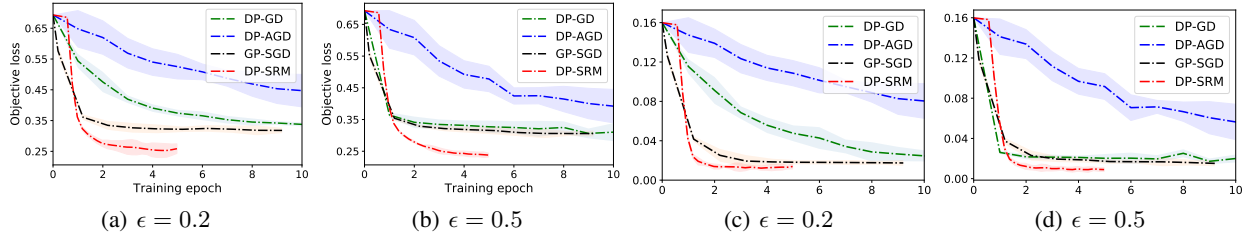


Figure 1: Results for nonconvex logistic regression on *ijcnn1* dataset. (a), (b) show the objective loss versus the number of epochs. (c), (d) illustrate the gradient norm versus the number of epochs.

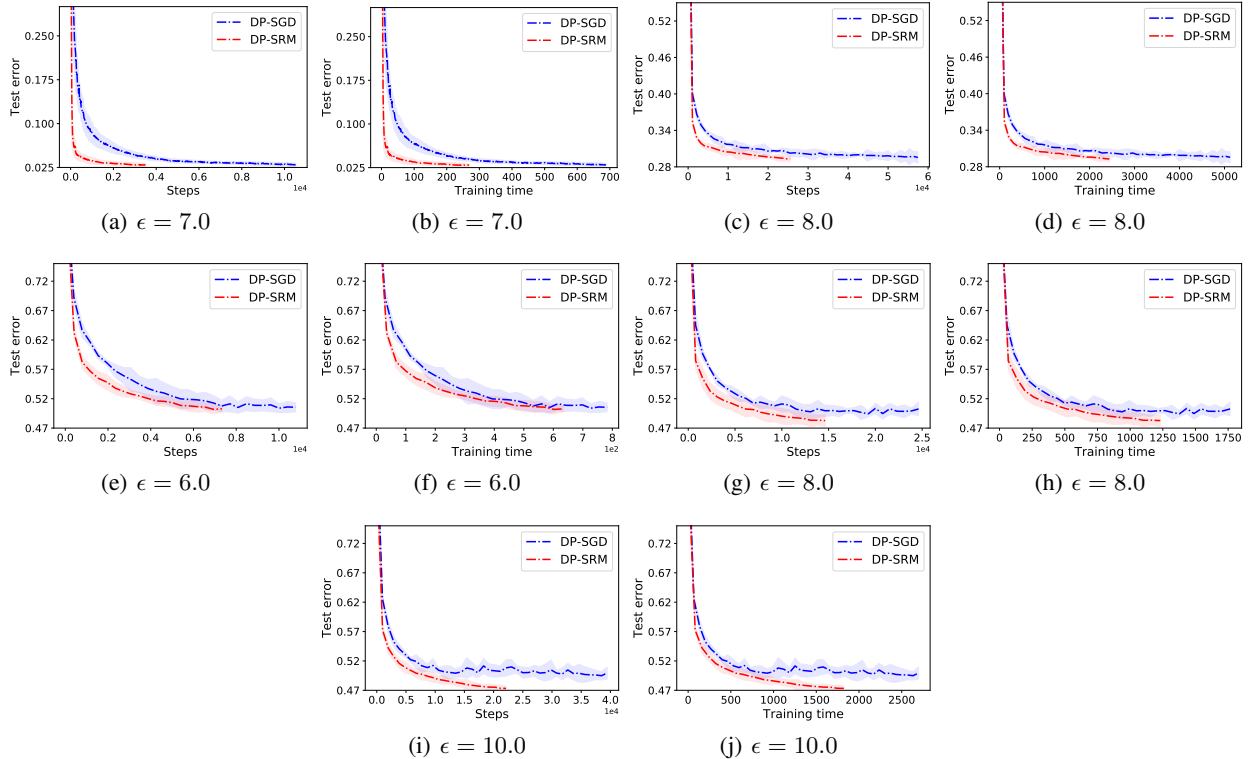


Figure 2: Results for CNN on MNIST and CIFAR-10 datasets. (a), (b) illustrate the results on MNIST dataset. (c), (d) demonstrate the results for CNN6 on CIFAR-10 dataset. (e)-(j) show the results for CNN5 on CIFAR-10 dataset.

test error when $\epsilon = 7.0$, which is comparable to the 2.93% test errors achieved by DP-SGD. Furthermore, the results show that our method is more efficient than DP-SGD in terms of iteration numbers and the training time. More specifically, our method is more than $2\times$ faster than DP-SGD to achieve the desired test error.

Parameters for CNN5. We choose three different privacy budgets $\epsilon \in \{6.0, 8.0, 10.0\}$, and set $\delta = 10^{-5}$. We set the clipping parameter $C_1 = 2$ for the term $\|\nabla f_i(\theta^t)\|_2$. For the term $\|\nabla f_i(\theta^t) - \nabla f_i(\theta^{t-1})\|_2$, we choose the clipping parameter C_2 by searching the grid $\{0.01, 0.1, 0.3, 0.5, 0.7, 0.9, 0.99\}$. For DP-SGD, we tune the batch size by searching the grid $\{32, 64, 128\}$ and the step size by $\{0.01, 0.02, 0.05, 0.1, 0.2\}$. For DP-SRM, we tune the batch size b by searching the grid $\{32, 64, 128\}$, step size by $\{0.01, 0.02, 0.05, 0.1, 0.2\}$, and b_0 by $\{b, 2b, 4b\}$. In addition, we set the momentum parameter $\gamma = C_2$.

Results for CNN5 on CIFAR-10 dataset. Figures 2(e)-2(j) present the average test error of different methods versus the number of iterations as well as the training time under different privacy budgets for CNN5 on CIFAR-10 dataset. The CNN5 trained by the non-private SGD will have 39.5% test error after 100 epochs. The results show that that our proposed method has 50.3%, 48.2% and 47.1% test errors when $\epsilon = 6.0$, $\epsilon = 8.0$ and $\epsilon = 10.0$. Nevertheless, DP-SGD has 51.0%,

50.2% and 49.3% test errors under the privacy budgets $\epsilon = 6.0$, $\epsilon = 8.0$ and $\epsilon = 10.0$, which are worse than our method. Furthermore, we can see from the plots that compared with DP-SGD, our method can reduce both the iteration numbers and the training time.

Results for CNN6 on CIFAR-10 dataset. Figure 2(c) and Figure 2(d) illustrate the average test error of different methods versus the number of iterations and the training time for CNN6 on CIFAR-10 dataset. We can see from the results that that our proposed method can achieve 29.3% test errors given the privacy budget $\epsilon = 8.0$, which are comparable to the results of DP-SGD with 29.4% under the same privacy budget. However, we can see from the plots that our method can significantly reduce the iteration numbers and the training time. When $\epsilon = 8$, DP-SGD takes 5.8×10^4 iterations and 5176 seconds to achieve 29.4% test error. In sharp contrast, our method only takes 2.6×10^4 iterations and 2589 seconds to achieve 29.3% test error.

2 PROOF OF MAIN RESULTS

In this section, we present the proofs of our main results.

2.1 PROOF OF THEOREM 5.1

We will provide the privacy guarantee of Algorithm 1 in this subsection. To this end, we need the following composition rule for RDP.

Lemma 2.1 (Mironov [2017]).] If k randomized mechanisms $\mathcal{M}_i : \mathcal{S}^n \rightarrow \mathcal{R}$ for $i \in [k]$, satisfy (α, ρ_i) -RDP, then their composition $(\mathcal{M}_1(S), \dots, \mathcal{M}_k(S))$ satisfies $(\alpha, \sum_{i=1}^k \rho_i)$ -RDP. Moreover, the input of the i -th mechanism can base on the outputs of previous $(i - 1)$ mechanisms.

We will first show that our proposed algorithm satisfies RDP using Lemma 3.7 and Lemma 2.1. Then we will transform it into (ϵ, δ) -DP based on Lemma 3.9. For the given dataset S , we use S' to denote its neighboring dataset with one different example indexed by i' in the following discussion. According to Algorithm 1, we use the following \mathcal{M}_t to denote the mechanism at t -th iteration

$$\mathcal{M}_t = \begin{cases} \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) + (1 - \gamma)(\mathbf{v}_p^{t-1} - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1})) + \mathbf{u}^t, & t > 0, \\ \mathbf{v}^0 + \mathbf{u}^0, & t = 0. \end{cases} \quad (2.1)$$

Therefore, our goal is to show the privacy guarantees of \mathcal{M}_t for $t = 0, 1, \dots, T$.

Case 1: If $t = 0$, we have $\mathbf{v}^0 = \nabla F_{\mathcal{B}_0}(\boldsymbol{\theta}^0)$ and \mathcal{M}_0 is equivalent to the following Gaussian mechanism

$$\mathcal{G}_0 = \nabla F_{\mathcal{B}_0}(\boldsymbol{\theta}^0) + \mathbf{u}^0,$$

where $\mathbf{u}^0 \sim N(0, \sigma_0^2 \mathbf{I}_d)$. Note that the mechanism \mathcal{G}_0 is based on the subsampling, thus we will use the results of privacy-amplification by subsampling, i.e., Lemma 3.7, to show that \mathcal{G}_0 satisfies RDP given appropriate \mathbf{u}^0 . To this end, we first consider the following Gaussian mechanism without subsampling

$$\tilde{\mathcal{G}}_0 = \frac{1}{b_0} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^0) + \mathbf{u}^0.$$

Sensitivity. Consider the query on the dataset S as follows $\tilde{\mathbf{q}}_0(S) = \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^0)/b_0$, where $\tilde{\mathbf{q}}_0(S)$ denotes that the query is based on the dataset S . Thus, we have

$$\tilde{\mathbf{q}}_0(S) - \tilde{\mathbf{q}}_0(S') = \frac{1}{b_0} (\nabla f_i(\boldsymbol{\theta}^0) - \nabla f_{i'}(\boldsymbol{\theta}^0)).$$

Since each component function is G -Lipschitz, we can obtain the ℓ_2 -sensitivity of this query as follows

$$\tilde{\Delta}_0 = \frac{1}{b_0} \|\nabla f_i(\boldsymbol{\theta}^0) - \nabla f_{i'}(\boldsymbol{\theta}^0)\|_2 \leq \frac{2G}{b_0}. \quad (2.2)$$

Privacy guarantee of \mathcal{G}_0 . By Lemma 3.7, if the Gaussian noise \mathbf{u}^0 in $\tilde{\mathcal{G}}_0$ has the following variance

$$\sigma_0^2 = \frac{14T\alpha G^2}{\beta n^2 \epsilon}, \quad (2.3)$$

the mechanism $\tilde{\mathcal{G}}_0$ satisfies $(\alpha, \beta\epsilon n^2/(7b_0^2 T))$ -RDP. Therefore, according to the privacy-amplification by subsampling result in Lemma 3.7, we have that the mechanism \mathcal{G}_0 satisfies (α, ρ_0) -RDP, where $\rho_0 = \beta\epsilon/T$. Furthermore, the variance σ_0^2 should satisfy the following condition

$$\frac{\sigma_0^2}{\tilde{\Delta}_0^2} = \frac{\sigma_0^2 b_0^2}{4G^2} = \frac{7b_0^2 T \alpha}{\beta n^2 \epsilon} \geq 0.7.$$

And the parameter α should satisfy $\alpha \leq 1 + 2(\sigma_0/\tilde{\Delta}_0)^2 \log(1/\tau\alpha(1 + (\sigma_0/\tilde{\Delta}_0)^2))/3$.

Case 2: If $t > 0$, according to the definition of \mathcal{M}_t in (2.1), we consider the following Gaussian mechanism

$$\mathcal{G}_t = \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - (1 - \gamma)\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1}) + \mathbf{u}^t.$$

Now, we are going to show that \mathcal{G}_t satisfies RDP given appropriate \mathbf{u}^t . Since the mechanism \mathcal{G}_t is based on the subsampling, we will use the similar proof procedure as in **Case 1** to show that \mathcal{G}_t satisfies RDP. Thus we consider the following Gaussian mechanism without subsampling

$$\tilde{\mathcal{G}}_t = \frac{1}{b} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^t) - (1 - \gamma) \frac{1}{b} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^{t-1}) + \mathbf{u}^t.$$

Sensitivity. We consider the following query without subsampling

$$\tilde{\mathbf{q}}_t(S) = \frac{1}{b} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^t) - (1 - \gamma) \frac{1}{b} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}^{t-1}).$$

Thus we have

$$\tilde{\mathbf{q}}_t(S) - \tilde{\mathbf{q}}_t(S') = \frac{1}{b} (\nabla f_i(\boldsymbol{\theta}^t) - (1 - \gamma)\nabla f_i(\boldsymbol{\theta}^{t-1}) - \nabla f_{i'}(\boldsymbol{\theta}^t) + (1 - \gamma)\nabla f_{i'}(\boldsymbol{\theta}^{t-1})).$$

As a result, we can obtain the ℓ_2 -sensitivity of the query $\tilde{\mathbf{q}}_t$ as follows

$$\begin{aligned} \tilde{\Delta}_t &= \frac{1}{b} \left\| (1 - \gamma)(\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}^{t-1}) - \nabla f_{i'}(\boldsymbol{\theta}^t) + \nabla f_{i'}(\boldsymbol{\theta}^{t-1})) \right. \\ &\quad \left. + \gamma(\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_{i'}(\boldsymbol{\theta}^t)) \right\|_2 \\ &\leq \frac{2L(1 - \gamma)}{b} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2 + \frac{2\gamma G}{b}, \end{aligned}$$

where the inequality is due to L -Lipschitz continuous gradient and G -Lipschitz of each component function. Furthermore, according to the update rule of Algorithm 1 and the definition of η_{t-1} , we have

$$\|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2 \leq \eta_{t-1} \|\mathbf{v}_p^{t-1}\|_2 \leq \min \left\{ \frac{\zeta}{n_0 L \|\mathbf{v}_p^{t-1}\|_2}, \frac{1}{2n_0 L} \right\} \cdot \|\mathbf{v}_p^{t-1}\|_2 \leq \frac{\zeta}{n_0 L},$$

which implies that

$$\tilde{\Delta}_t \leq \frac{2L(1 - \gamma)}{b} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2 + \frac{2\gamma G}{b} \leq \frac{2((1 - \gamma)\zeta/n_0 + \gamma G)}{b}. \quad (2.4)$$

Privacy guarantee of \mathcal{G}_t . By Lemma 3.7, if we the Gaussian noise \mathbf{u}^t in $\tilde{\mathcal{G}}_t$ has the variance as follows

$$\sigma_t^2 = \frac{14T\alpha((1 - \gamma)\zeta/n_0 + \gamma G)^2}{\beta n^2 \epsilon}, \quad (2.5)$$

the mechanism $\tilde{\mathcal{G}}_t$ satisfies $(\alpha, \beta\epsilon n^2/(7b^2T))$ -RDP. Thus based on the privacy-amplification by subsampling result (Lemma 3.7), we can get that the mechanism \mathcal{G}_t satisfies (α, ρ) -RDP, where $\rho = \beta\epsilon/T$. In addition, the variance σ_t^2 should satisfy the following condition

$$\frac{\sigma_t^2}{\tilde{\Delta}_t^2} = \frac{\sigma_t^2 b^2}{4((1-\gamma)\zeta/n_0 + \gamma G)^2} = \frac{7b^2 T \alpha}{\beta n^2 \epsilon} \geq 0.7.$$

And the parameter α should satisfy $\alpha \leq 1 + 2(\sigma_t/\tilde{\Delta}_t)^2 \log(1/\tau\alpha(1 + (\sigma_t/\tilde{\Delta}_t)^2))/3$. As a result, we show that \mathcal{G}_t satisfies (α, ρ) -RDP.

Privacy guarantee of \mathcal{M}_t . By the definition of the mechanism \mathcal{M}_t in (2.1), \mathcal{M}_t is a composition of $\mathcal{G}_0, \dots, \mathcal{G}_t$, i.e., $\mathcal{M}_t = (\mathcal{G}_0, \dots, \mathcal{G}_t)$. According to the composition property of RDP, i.e., Lemma 2.1, we have \mathcal{M}_t satisfies $(\alpha, \rho_0 + (t-1)\rho)$ -RDP. Since $\rho_0 = \rho = \beta\epsilon/T$, we have that after T' iterations of Algorithm 1, it satisfies $(\alpha, \beta T' \epsilon/T)$ -RDP. According to Lemma 3.9 and $\alpha = \log(1/\delta)/((1-\beta)\epsilon) + 1$, we have that after T' iterations, Algorithm 1 satisfies $(T' \epsilon/T, \delta)$ -DP. As a result, we have that for each θ^t , where $t = 1, \dots, T$, it satisfies (ϵ, δ) -DP. Finally, by the definition of $\tilde{\theta}$, we have $\tilde{\theta}$ satisfies (ϵ, δ) -DP.

2.2 PROOF OF COROLLARY 5.3

In this subsection, we show that by choosing a larger mini-batch size, we can get rid of the constraints in Theorem 5.1. More specifically, let $b_0^2 = b^2 = n^2 \epsilon/T$ and $\beta = 1/2$, we have $\sigma'^2 = 7T\alpha b^2/(\beta n^2 \epsilon) = 14\alpha$. Furthermore, we have

$$\tau\alpha(1 + \sigma'^2) \stackrel{(i)}{\leq} 15\tau\alpha^2 \stackrel{(ii)}{=} 15(2\log(1/\delta)/\epsilon + 1)^2 \sqrt{\epsilon/T},$$

where (i) uses $\sigma'^2 = 14\alpha$, (ii) uses $\tau = b/n = \sqrt{\epsilon/T}$ and $\epsilon = 2\log(1/\delta)/(\alpha - 1)$. If $\epsilon \leq 2\log(1/\delta)$, we can obtain $\tau\alpha(1 + \sigma'^2) \leq 1/3$ if T is larger than $O(\log^4(1/\delta)/\epsilon^3)$. If $\epsilon > 2\log(1/\delta)$, we can obtain $\tau\alpha(1 + \sigma'^2) \leq 1/3$ if T is larger than $O(\epsilon)$. Therefore, we can get $\log(1/\tau\alpha(1 + \sigma'^2)) \geq 1$. As a result, we have $2(\sigma'^2 \log(1/\tau\alpha(1 + \sigma'^2)))/3 \geq 28\alpha/3 > \alpha - 1$.

2.3 PROOF OF THEOREM 5.4

In this subsection, we provide the utility guarantee of our method. According to the assumption that each component function has L -Lipschitz continuous gradient, we can obtain that

$$\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|_2 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2,$$

which implies that $F(\mathbf{x})$ has L -Lipschitz continuous gradient. Thus we have

$$\begin{aligned} F(\theta^{t+1}) &\leq F(\theta^t) + \langle \nabla F(\theta^t), \theta^{t+1} - \theta^t \rangle + \frac{L}{2} \|\theta^{t+1} - \theta^t\|_2^2 \\ &= F(\theta^t) - \eta_t \langle \nabla F(\theta^t), \mathbf{v}_p^t \rangle + \frac{\eta_t^2 L}{2} \|\mathbf{v}_p^t\|_2^2 \\ &= F(\theta^t) + \frac{\eta_t}{2} \|\nabla F(\theta^t) - \mathbf{v}_p^t\|_2^2 - \frac{\eta_t}{2} \|\nabla F(\theta^t)\|_2^2 - \eta_t \left(\frac{1}{2} - \frac{\eta_t L}{2} \right) \|\mathbf{v}_p^t\|_2^2, \end{aligned}$$

where the last equality is due to the fact that $2\langle \nabla F(\theta^t), \mathbf{v}_p^t \rangle = \|\nabla F(\theta^t)\|_2^2 + \|\mathbf{v}_p^t\|_2^2 - \|\nabla F(\theta^t) - \mathbf{v}_p^t\|_2^2$. Since $\eta_t \leq 1/(2n_0L)$, we can obtain that

$$F(\theta^{t+1}) \leq F(\theta^t) + \frac{1}{4n_0L} \|\nabla F(\theta^t) - \mathbf{v}_p^t\|_2^2 - \frac{\eta_t}{4} \|\mathbf{v}_p^t\|_2^2.$$

In addition, we have

$$\frac{\eta_t}{4} \|\mathbf{v}_p^t\|_2^2 = \frac{\zeta^2}{8n_0L} \min \{2\|\mathbf{v}_p^t/\zeta\|_2, \|\mathbf{v}_p^t/\zeta\|_2^2\} \geq \frac{\zeta \|\mathbf{v}_p^t\|_2 - 2\zeta^2}{4n_0L}.$$

Thus we have

$$F(\boldsymbol{\theta}^{t+1}) \leq F(\boldsymbol{\theta}^t) + \frac{1}{4n_0L} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2^2 - \frac{\zeta \|\mathbf{v}_p^t\|_2}{4n_0L} + \frac{\zeta^2}{2n_0L}. \quad (2.6)$$

Summing over $t = 0, \dots, T-1$ and taking expectation in (2.6), we can get

$$\begin{aligned} \frac{\zeta}{4n_0L} \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{v}_p^t\|_2 &\leq F(\boldsymbol{\theta}^0) - \mathbb{E} F(\boldsymbol{\theta}^T) + \frac{1}{4n_0L} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2^2 + \frac{T\zeta^2}{2n_0L} \\ &\leq F(\boldsymbol{\theta}^0) - F(\boldsymbol{\theta}^*) + \frac{1}{4n_0L} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2^2 + \frac{T\zeta^2}{2n_0L}. \end{aligned} \quad (2.7)$$

For the term $\mathbb{E} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2^2$, we can bound it as follows: we first consider the conditional expectation

$$\begin{aligned} \mathbb{E}_t \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 &= \mathbb{E}_t \|(1-\gamma)(\mathbf{v}_p^{t-1} - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1})) + \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^t) + \mathbf{u}^t\|_2^2 \\ &= \mathbb{E}_t \|(1-\gamma)(\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})) + (1-\gamma)\nabla F(\boldsymbol{\theta}^{t-1}) \\ &\quad - (1-\gamma)\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1}) + \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^t)\|_2^2 + \mathbb{E}_t \|\mathbf{u}^t\|_2^2 \\ &= \mathbb{E}_t \|(1-\gamma)(\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})) + (1-\gamma)(\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1})) \\ &\quad + \nabla F(\boldsymbol{\theta}^{t-1}) - \nabla F(\boldsymbol{\theta}^t)\|_2^2 + \mathbb{E}_t \|\mathbf{u}^t\|_2^2, \end{aligned} \quad (2.8)$$

where \mathbb{E}_t is taken over the randomness at the t -th iteration given the observations after $(t-1)$ -th iteration, the first equation comes from the definition of \mathbf{v}_p^t , the second one is due to the independence of the random variables. Therefore, we can obtain that

$$\begin{aligned} \mathbb{E}_t \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 &= (1-\gamma)^2 \mathbb{E}_t \|\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})\|_2^2 \\ &\quad + 2\gamma^2 \mathbb{E}_t \|\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^t)\|_2^2 + \mathbb{E}_t \|\mathbf{u}^t\|_2^2 \\ &\quad + 2(1-\gamma)^2 \mathbb{E}_t \|\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1}) + \nabla F(\boldsymbol{\theta}^{t-1}) - \nabla F(\boldsymbol{\theta}^t)\|_2^2, \end{aligned} \quad (2.9)$$

where the equality is due to the expansion of (2.8) and Cauchy-Schwartz inequality. In addition, we have

$$\begin{aligned} &\mathbb{E}_t \|\nabla F(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^{t-1}) - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) + \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1})\|_2^2 \\ &\leq \frac{1}{b} \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla F(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^{t-1}) - \nabla f_i(\boldsymbol{\theta}^t) + \nabla f_i(\boldsymbol{\theta}^{t-1})\|_2^2 \\ &\leq \frac{1}{b} \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}^{t-1})\|_2^2 \\ &\leq \frac{L^2}{b} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2^2, \end{aligned}$$

where the first inequality is due to Lemma 4.1, the second one comes from the fact that $\mathbb{E} \|\mathbf{X} - \mathbb{E}\mathbf{X}\|_2^2 \leq \mathbb{E} \|\mathbf{X}\|_2^2$ for any random variable \mathbf{X} , and the last one is due to the gradient Lipschitz property of each component function. According to the update rule, we have

$$\|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2 \leq \eta_{t-1} \|\mathbf{v}_p^{t-1}\|_2 \leq \min \left\{ \frac{\zeta}{n_0L \|\mathbf{v}_p^{t-1}\|_2}, \frac{1}{2n_0L} \right\} \cdot \|\mathbf{v}_p^{t-1}\|_2 \leq \frac{\zeta}{n_0L},$$

which implies

$$\mathbb{E}_t \|\nabla F(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^{t-1}) - \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) + \nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^{t-1})\|_2^2 \leq \frac{\zeta^2}{n_0^2 b}. \quad (2.10)$$

Thus plugging (2.10) into (2.9), we can obtain that

$$\begin{aligned}\mathbb{E}_t \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 &\leq (1-\gamma)^2 \|\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})\|_2^2 + \frac{2(1-\gamma)^2 L^2}{b} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t-1}\|_2^2 \\ &\quad + 2\gamma^2 \mathbb{E}_t \|\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^t)\|_2^2 + \mathbb{E}_t \|\mathbf{u}^t\|_2^2 \\ &\leq (1-\gamma)^2 \|\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})\|_2^2 + \frac{2(1-\gamma)^2 \zeta^2}{n_0^2 b} + \frac{2\gamma^2 G^2}{b} + \mathbb{E}_t \|\mathbf{u}^t\|_2^2,\end{aligned}\quad (2.11)$$

where the second inequality follows the following inequality (using Lemma 4.1, $\mathbb{E} \|\mathbf{X} - \mathbb{E}\mathbf{X}\|_2^2 \leq \mathbb{E} \|\mathbf{X}\|_2^2$, and the G -Lipschitz of each component function)

$$\mathbb{E}_t \|\nabla F_{\mathcal{B}_t}(\boldsymbol{\theta}^t) - \nabla F(\boldsymbol{\theta}^t)\|_2^2 \leq \frac{1}{b} \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\boldsymbol{\theta}^t)\|_2^2 \leq \frac{G^2}{b}.\quad (2.12)$$

Therefore, taking expectations over all iterations in (2.11), we can get

$$\mathbb{E} \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 \leq (1-\gamma)^2 \mathbb{E} \|\mathbf{v}_p^{t-1} - \nabla F(\boldsymbol{\theta}^{t-1})\|_2^2 + \frac{2(1-\gamma)^2 \zeta^2}{n_0^2 b} + \frac{2\gamma^2 G^2}{b} + d\sigma^2.\quad (2.13)$$

Following the proof of Lemma 9 in Yuan et al. [2020], we have

$$\begin{aligned}\gamma \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 &\leq \frac{2T(1-\gamma)^2 \zeta^2}{n_0^2 b} + \frac{2T\gamma^2 G^2}{b} + Td\sigma^2 + \mathbb{E} \|\mathbf{v}_p^0 - \nabla F(\boldsymbol{\theta}^0)\|_2^2 \\ &\leq \frac{2T(1-\gamma)^2 \zeta^2}{n_0^2 b} + \frac{2T\gamma^2 G^2}{b} + Td\sigma^2 + \frac{G^2}{b_0} + d\sigma_0^2,\end{aligned}$$

where the last line comes from the definition of $\mathbf{v}_p^0 = \nabla F_{\mathcal{B}_0}(\boldsymbol{\theta}^0) + \mathbf{u}^0$ and the inequality $\mathbb{E} \|\nabla F_{\mathcal{B}_0}(\boldsymbol{\theta}^0) - \nabla F(\boldsymbol{\theta}^0)\|_2^2 \leq G^2/b_0$ (see equation (2.12)). Therefore, we can obtain that

$$\sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{v}_p^t - \nabla F(\boldsymbol{\theta}^t)\|_2^2 \leq \frac{2T(1-\gamma)^2 \zeta^2}{n_0^2 \gamma b} + \frac{2T\gamma G^2}{b} + \frac{Td\sigma^2 + d\sigma_0^2}{\gamma} + \frac{G^2}{\gamma b_0}.\quad (2.14)$$

Combining (2.7) and (2.14), we can get

$$\begin{aligned}\frac{\zeta}{4n_0 L} \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{v}_p^t\|_2 &\leq F(\boldsymbol{\theta}^0) - F(\boldsymbol{\theta}^*) + \frac{1}{4n_0 L} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2 + \frac{T\zeta^2}{2n_0 L} \\ &\leq F(\boldsymbol{\theta}^0) - F(\boldsymbol{\theta}^*) + \frac{T(1-\gamma)^2 \zeta^2}{2n_0^3 L \gamma b} + \frac{T\gamma G^2}{4Ln_0 b} \\ &\quad + \frac{Td\sigma^2 + d\sigma_0^2}{4n_0 L \gamma} + \frac{G^2}{4L\gamma n_0 b_0} + \frac{T\zeta^2}{2n_0 L}.\end{aligned}$$

Hence we have

$$\begin{aligned}\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{v}_p^t\|_2 &\leq \frac{4n_0 L}{T\zeta} (F(\boldsymbol{\theta}^0) - F(\boldsymbol{\theta}^*)) + \frac{2\zeta}{n_0^2 \gamma b} + \frac{\gamma G^2}{\zeta b} + \frac{d\sigma^2 + d\sigma_0^2/T}{\zeta \gamma} + \frac{G^2}{T\zeta \gamma b_0} + 2\zeta \\ &\leq 6\zeta + \frac{2\zeta}{n_0^2 \gamma b} + \frac{\gamma G^2}{\zeta b} + \frac{d\sigma^2 + d\sigma_0^2/T}{\zeta \gamma} + \frac{G^2}{T\zeta \gamma b_0},\end{aligned}\quad (2.15)$$

where the first inequality is due to $T = \lfloor 4n_0 L (F(\boldsymbol{\theta}^0) - F(\boldsymbol{\theta}^*)) / \zeta^2 \rfloor + 1$. In addition, according to (2.14) and Jensen's inequality, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla F(\boldsymbol{\theta}^t) - \mathbf{v}_p^t\|_2 \leq \frac{\sqrt{2}\zeta}{n_0 \sqrt{\gamma b}} + \frac{\sqrt{2}\gamma G}{\sqrt{b}} + \frac{\sqrt{d}\sigma + \sqrt{d}\sigma_0/\sqrt{T}}{\sqrt{\gamma}} + \frac{G}{\sqrt{T\gamma b_0}}.\quad (2.16)$$

Thus by the definition of $\tilde{\theta}$, we have

$$\begin{aligned}
\mathbb{E}\|\nabla F(\tilde{\theta})\|_2 &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla F(\theta^t)\|_2 \\
&\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\mathbf{v}_p^t\|_2 + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla F(\theta^t) - \mathbf{v}_p^t\|_2 \\
&\leq 6\zeta + \frac{2\zeta}{n_0^2\gamma b} + \frac{\gamma G^2}{\zeta b} + \frac{d\sigma^2}{\zeta\gamma} + \frac{d\sigma_0^2}{T\zeta\gamma} + \frac{G^2}{T\zeta\gamma b_0} + \frac{\sqrt{2}\zeta}{n_0\sqrt{\gamma b}} + \frac{\sqrt{2\gamma}G}{\sqrt{b}} \\
&\quad + \frac{\sqrt{d}\sigma}{\sqrt{\gamma}} + \frac{\sqrt{d}\sigma_0}{\sqrt{T\gamma}} + \frac{G}{\sqrt{T\gamma b_0}}, \tag{2.17}
\end{aligned}$$

where the second inequality comes from (2.15) and (2.16). Let $\gamma^2 = 2\zeta^2/(n_0^2G^2)$, $b = G/(n_0\zeta)$, $b_0 = G^3/(\zeta LD_F)$, where $D_F = F(\theta^0) - F(\theta^*)$ and $F(\theta^*)$ is a global minimum of F , by the definition of T , we can get

$$\mathbb{E}\|\nabla F(\tilde{\theta})\|_2 \leq 15\zeta + \frac{d\sigma^2}{\zeta\gamma} + \frac{\sqrt{d}\sigma}{\sqrt{\gamma}} + \frac{d\sigma_0^2}{T\zeta\gamma} + \frac{\sqrt{d}\sigma_0}{\sqrt{T\gamma}}. \tag{2.18}$$

Furthermore, we have

$$\sigma^2 = \frac{14T((1-\gamma)\zeta/n_0 + \gamma G)^2 \log(1/\delta)}{n^2\epsilon^2}, \quad \sigma_0^2 = \frac{14TG^2 \log(1/\delta)}{n^2\epsilon^2}. \tag{2.19}$$

Plugging (2.19) into (2.18), we can obtain

$$\begin{aligned}
\mathbb{E}\|\nabla F(\tilde{\theta})\|_2 &\leq 15\zeta + \frac{C_1TdG \log(1/\delta)}{n_0n^2\epsilon^2} + \frac{\sqrt{C_1T\zeta dG \log(1/\delta)}}{n\epsilon\sqrt{n_0}} + \frac{C_2dn_0G^3 \log(1/\delta)}{n^2\epsilon^2\zeta^2} \\
&\quad + \frac{\sqrt{C_2n_0dG^3 \log(1/\delta)}}{n\epsilon\sqrt{\zeta}} \\
&\leq 15\zeta + \frac{C_3LD_FGd \log(1/\delta)}{n^2\epsilon^2\zeta^2} + \frac{\sqrt{C_4GLD_Fd \log(1/\delta)}}{n\epsilon\sqrt{\zeta}} \\
&\quad + \frac{C_5n_0dG^3 \log(1/\delta)}{n^2\epsilon^2\zeta^2} + \frac{\sqrt{C_6n_0dG^3 \log(1/\delta)}}{n\epsilon\sqrt{\zeta}}, \tag{2.20}
\end{aligned}$$

where the second inequality is due to the fact that $T = \lceil 4n_0LD_F/\zeta^2 \rceil + 1$. Without loss of generality, we can assume $G \geq 1$ and $\zeta \leq 1$. Therefore, let $n_0 = LD_F/G^2 \cdot (G/\zeta)^\kappa$ with $\kappa \in [0, 1]$, and plugging n_0 into (2.20), we can obtain

$$\mathbb{E}\|\nabla F(\tilde{\theta})\|_2 \leq 15\zeta + \frac{C_7LD_FGd \log(1/\delta)G^\kappa}{n^2\epsilon^2\zeta^{2+\kappa}} + \frac{C_8\sqrt{GLD_Fd \log(1/\delta)}G^{\frac{\kappa}{2}}}{n\epsilon\zeta^{\frac{1+\kappa}{2}}}. \tag{2.21}$$

Thus, choosing

$$\zeta = C_9 \left(\frac{G^{\frac{\kappa}{2}} \sqrt{GLD_Fd \log(1/\delta)}}{n\epsilon} \right)^{\frac{2}{3+\kappa}}, \tag{2.22}$$

we can get

$$\mathbb{E}\|\nabla F(\tilde{\theta})\|_2 \leq C_{10} \left(\frac{G^{\frac{\kappa}{2}} \sqrt{GLD_Fd \log(1/\delta)}}{n\epsilon} \right)^{\frac{2}{3+\kappa}}. \tag{2.23}$$

Note that we require $\gamma \leq 1$, which gives us $n\epsilon \geq O(G^2(d \log(1/\delta))^{1/2}/(LD_F))$.

Furthermore, according to Theorem 5.1 to achieve the desired privacy guarantee, we require $\sigma'^2 = \min\{b^2\sigma^2/(4((1-\gamma)\zeta/n_0 + \gamma G)^2), b_0^2\sigma_0^2/(4G^2)\} \geq 0.7$. Note that $b = G/(n_0\zeta)$, $b_0 = G^3/(\zeta LD_F)$, $n_0 = LD_F/G^2 \cdot (G/\zeta)^\kappa$, we have

$b = b_0 \cdot (\zeta/G)^\kappa$. Thus, the aforementioned requirement reduces to

$$\begin{aligned} \frac{14b_0^2 T \log(1/\delta)}{4n^2 \epsilon^2} \cdot \frac{\zeta^{2\kappa}}{G^{2\kappa}} &= \frac{14b_0^2 n_0 L D_F \log(1/\delta)}{\zeta^4 n^2 \epsilon^2} \cdot \frac{\zeta^{2\kappa}}{G^{2\kappa}} \\ &\geq \frac{14b_0 n_0 L D_F \log(1/\delta)}{\zeta^4 n^2 \epsilon^2} \cdot \frac{\zeta^{2\kappa}}{G^{2\kappa}} \\ &= \frac{14G L D_F \log(1/\delta)}{\zeta^3 n^2 \epsilon^2} \cdot \frac{\zeta^\kappa}{G^\kappa} \\ &\geq 0.7, \end{aligned}$$

where the first equality comes from the definition of T and the first inequality is due to $b_0 \geq 1$. Therefore, we need

$$\zeta \leq \left(4 \frac{G^{-\frac{\kappa}{2}} \sqrt{G L D_F d \log(1/\delta)}}{n \epsilon} \right)^{\frac{2}{3-\kappa}}. \quad (2.24)$$

Combining (2.22) and (2.24), we need to choose $\kappa = 0$ in n_0 , which gives us

$$\mathbb{E} \|\nabla F(\tilde{\theta})\|_2 \leq C_{10} \left(\frac{\sqrt{G L D_F d \log(1/\delta)}}{n \epsilon} \right)^{\frac{2}{3}}, \quad (2.25)$$

where $\{C_i\}_{i=1}^{10}$ are absolute constants. Furthermore, the requirement $\alpha - 1 = \log(1/\delta)/((1-\beta)\epsilon) \leq 2\sigma'^2 \log(1/(\tau\alpha(1+\sigma'^2)))/3$ in Theorem 5.1 can be satisfied under our choice of parameters given large enough n . Since we have $\sigma'^2 \geq 0.7$, we have $2\sigma'^2 \log(1/(\tau\alpha(1+\sigma'^2)))/3 \geq 0.4 \log(1/(\tau\alpha(1+\sigma'^2))) \geq 0.4 \log(1/(3\tau\alpha\sigma'^2))$. Furthermore, we have

$$\begin{aligned} \tau\alpha\sigma'^2 &= \frac{G^3}{n\zeta L D_F} \cdot \frac{\log(1/\delta) + (1-\beta)\epsilon}{(1-\beta)\epsilon} \cdot \frac{14G L D_F \log(1/\delta)}{\zeta^3 n^2 \epsilon^2} \\ &\leq \frac{28G^4 \log^2(1/\delta)}{(1-\beta)n^3 \epsilon^3 \zeta^4} \\ &\leq C_{11} \frac{G^4 \log^2(1/\delta)}{(n\epsilon)^3} \cdot \frac{(n\epsilon)^{8/3}}{(G L D_F d \log(1/\delta))^{4/3}} \\ &= C_{11} \frac{G^{8/3} \log^{2/3}(1/\delta)}{(n\epsilon)^{1/3} (L D_F d)^{4/3}}, \end{aligned}$$

where the first inequality comes from assuming $\epsilon \leq \log(1/\delta)$ without loss of generality, and the second inequality is due to the definition of ζ . Thus we have

$$\log(1/(3\tau\alpha\sigma'^2)) \geq \log\left(3C_{11} \frac{(n\epsilon)^{1/3} (L D_F d)^{4/3}}{G^{8/3} \log^{2/3}(1/\delta)}\right).$$

As a result, the requirement reduces to

$$0.4 \log\left(3C_{11} \frac{(n\epsilon)^{1/3} (L D_F d)^{4/3}}{G^{8/3} \log^{2/3}(1/\delta)}\right) \geq \frac{\log(1/\delta)}{(1-\beta)\epsilon},$$

which can be satisfied if we have

$$n \geq C_{12} \frac{G^8 \log^2(1/\delta)}{(L D_F d)^4 \epsilon},$$

where C_{11}, C_{12} are some large constants.

Gradient Complexity. Since we have $b = b_0 = G^3/(\zeta L D_F)$, the total gradient complexity is

$$2(T-1)b + b_0 \leq \frac{8L D_F n_0}{\zeta^2} \cdot \frac{G^3}{L D_F \zeta} + \frac{G^3}{L D_F \zeta}.$$

According to the definition of ζ and n_0 , we have the total gradient complexity is $O(n^2 \epsilon^2 / (d \log(1/\delta)))$.

3 PROOF OF LEMMA 3.7

Without loss of generality, we assume $\Delta(q) = 1$. According to Theorem 9 in Wang et al. [2019b], we have

$$\rho'(\alpha) \leq \frac{1}{\alpha-1} \log \left(1 + \tau^2 \binom{\alpha}{2} \min \left\{ 4(e^{\rho(2)} - 1), 2e^{\rho(2)} \right\} + \sum_{j=3}^{\alpha} \tau^j \binom{\alpha}{j} 2e^{(j-1)\rho(j)} \right), \quad (3.1)$$

where τ is the subsample rate, $\rho(j) = j/(2\sigma^2)$. Next, we will show that the summation term in the right hand side of the above inequality is dominated by the second term under certain conditions. First of all, when σ^2 is large, i.e., $\sigma^2 \geq 0.7$, we have

$$\min \left\{ 4(e^{\rho(2)} - 1), 2e^{\rho(2)} \right\} \leq 6/\sigma^2,$$

which implies that

$$\tau^2 \binom{\alpha}{2} \min \left\{ 4(e^{\rho(2)} - 1), 2e^{\rho(2)} \right\} \leq \tau^2 \binom{\alpha}{2} 6/\sigma^2.$$

Next, we consider the summation term in (3.1), and we have

$$\begin{aligned} \sum_{j=3}^{\alpha} \tau^j \binom{\alpha}{j} 2e^{(j-1)\rho(j)} &\leq \tau^2 \binom{\alpha}{2} \left(\sum_{j=3}^{\alpha} \tau^{j-2} \alpha^{j-2} e^{\frac{(\alpha-1)j}{2\sigma^2}} \right) \\ &\leq \tau^2 \binom{\alpha}{2} \frac{\tau \alpha e^{\frac{3(\alpha-1)}{2\sigma^2}}}{1 - \tau \alpha e^{\frac{\alpha-1}{2\sigma^2}}}, \end{aligned}$$

where the first inequality is due to the fact that

$$e^{(j-1)\rho(j)} = e^{\frac{(j-1)j}{2\sigma^2}} \leq e^{\frac{(\alpha-1)j}{2\sigma^2}} \quad \text{and} \quad \binom{\alpha}{j} = \frac{\alpha!}{j!(\alpha-j)!} \leq \frac{\alpha^2 \alpha^{j-2}}{3!}.$$

In addition, the last inequality comes from the condition that $\tau \alpha \exp((\alpha-1)/(2\sigma^2)) < 1$ and the sum of the geometric sequence. Therefore, as long as

$$\alpha - 1 \leq \frac{2}{3} \sigma^2 \log \frac{1}{\tau \alpha (1 + \sigma^2)}, \quad (3.2)$$

we have

$$\sum_{j=3}^{\alpha} \tau^j \binom{\alpha}{j} 2e^{(j-1)\rho(j)} \leq \tau^2 \binom{\alpha}{2} \frac{1}{\sigma^2}.$$

In addition, we require that $\tau \alpha \exp((\alpha-1)/(2\sigma^2)) < 1$. By plugging the condition of α into the above requirement, we can obtain that this condition can hold if $\tau < 1$.

As a result, under the conditions that $\sigma^2 \geq 0.7$, $\alpha \leq \log(1/\tau(1 + \sigma^2))$, we can obtain that

$$\rho'(\alpha) \leq \frac{1}{\alpha-1} \log \left(1 + \tau^2 \binom{\alpha}{2} \frac{10}{\sigma^2} \right) \leq \frac{1}{\alpha-1} \tau^2 \binom{\alpha}{2} \frac{7}{\sigma^2} \leq 3.5\alpha\tau^2/\sigma^2.$$

4 AUXILIARY LEMMAS

Lemma 4.1 (Lei et al., 2017). Consider vectors \mathbf{a}_i satisfying $\sum_{i=1}^n \mathbf{a}_i = 0$. Let \mathcal{B} be a uniform random subset of $\{1, 2, \dots, n\}$ with size m , we have

$$\mathbb{E} \left\| \frac{1}{m} \sum_{i \in \mathcal{B}} \mathbf{a}_i \right\|_2^2 \leq \frac{\mathbb{1}\{|\mathcal{B}| < n\}}{mn} \sum_{i=1}^n \|\mathbf{a}_i\|_2^2.$$