

# STOCHASTIC GAUSSIAN ZERO-ORDER OPTIMIZATION: IMPROVED CONVERGENCE ANALYSIS UNDER SKEWED HESSIAN SPECTRA

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper addresses large-scale finite-sum optimization problems, which are particularly prevalent in the big data era. In the field of zeroth-order optimization, stochastic methods have become essential tools. Natural zeroth-order stochastic methods primarily rely on stochastic gradient descent (SGD). Preprocessing the stochastic gradient using a Gaussian vector defines the method ZO-SGD-Gauss (ZSG), whereas estimating coordinate-wise partial derivatives defines ZO-SGD-Coordinate (ZSC). Compared to ZSC, ZSG often demonstrates superior performance in practice. However, the underlying mechanisms behind this phenomenon remain unclear in the academic community. To the best of our knowledge, our work is the first to theoretically analyze the potential advantages of ZSG compared to ZSC. To facilitate convergence analysis, the quadratic regularity assumption is introduced to generalize the smoothness and strong convexity to the Hessian matrix. This assumption makes it possible to integrate Hessian information into the complexity analysis. We provide a theoretical analysis proving the significant convergence improvement of ZSG. Finally, experiments on both synthetic and real-world datasets validate the effectiveness of our theoretical analysis.

## 1 INTRODUCTION

Modern machine learning presents significant challenges for optimization due to the large scale of the problems involved. Contemporary datasets are both enormous and high-dimensional, often with millions of samples and features. Because evaluating the full objective or gradient even once is too slow to be useful, stochastic optimization methods have emerged in response.

Throughout the paper, we aim to solve finite-sum minimization problems of the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}). \quad (1)$$

An optimization method that solves the problem 1 with function value access only is known as zeroth-order optimization or black-box optimization (Ghadimi & Lan, 2013; Nesterov & Spokoiny, 2017). In recent years, zeroth-order optimization has attracted widespread attention from both the machine learning community and the optimization community (Nesterov & Spokoiny, 2017; Ilyas et al., 2018). One important application of the zeroth-order optimization is the black-box adversarial attack on deep neural networks (Chen et al., 2017; Zhao et al., 2020; Zhang et al., 2020; Bai et al., 2023). In the black-box adversarial attack, only the inputs and outputs of the neural network are available and back propagation is often prohibited (Papernot et al., 2017). In the above situation, the evaluation of gradient  $\nabla f(\mathbf{x})$  is infeasible. So, applying zeroth-order optimization methods becomes a natural choice. Additional application scenarios in the field of artificial intelligence where zeroth-order optimization algorithms demonstrate significant effectiveness are deep reinforcement learning (Salimans et al., 2017; Mania et al., 2018; Zhang & Zavlanos, 2023; Jing et al., 2024), hyper-parameter tuning (Snoek et al., 2012; Rapin & Teytaud, 2018), the problem of optimizing functions with only ranking feedback (Tang et al., 2023), learning linear quadratic regulators (Malik et al., 2020; Mohammadi et al., 2020), and so on. Zeroth-order optimization has even played a

054 significant role in fine-tuning LLMs. Malladi et al. (2023) and Zhao et al. (2024) use the zeroth-  
055 order optimization methods for fine-tuning, in addressing the significant memory overhead of first-  
056 order optimizers. Zeroth-order optimization achieves a substantial memory reduction and makes it  
057 possible to train and store LLMs on low-cost hardware.

058 The ZO-SGD-Gauss (ZSG) algorithm is based on the Gaussian version of SPSA (Spall, 1992).  
059 The ZO-SGD-Coordinate (ZSC) algorithm is based on the finite-difference stochastic approxima-  
060 tion (Kiefer & Wolfowitz, 1952). Although the ZSG algorithm and the ZSC algorithm share the  
061 same theoretical convergence rate and both have sample complexity that is linear in the dimension  
062 (Ghadimi & Lan, 2013), ZSG has a wider range of applications and performs better than ZSC in  
063 practice. For example, ZSG has been widely used in fine-tuning LLMs (Malladi et al., 2023; Zhao  
064 et al., 2024) and black-box attacks (Ilyas et al., 2018). The academic community is still unclear  
065 about the underlying mechanism why ZSG outperforms ZSC. For the gradient descent method,  
066 recent works by Yue et al. (2023) and Wang et al. (2024) show that zeroth-order Gaussian gra-  
067 dient descent can outperform coordinate descent under skewed Hessian spectra, which indicates  
068 ill-conditioning and anisotropy in the loss landscape. An intriguing question is whether the zeroth-  
069 order SGD algorithm possesses a similar property to the zeroth-order gradient descent algorithm.  
070 Motivated by these works, we try to investigate whether ZSG can theoretically achieve better com-  
071 plexity than ZSC. We obtain a surprising result: compared to ZSC, ZSG exhibits weak dimensional  
072 dependence—meaning that the dimension  $d$  does not explicitly appear in the complexity bounds.  
073 Our work fills a theoretical gap in the field of zeroth-order optimization.

## 074 1.1 LITERATURE REVIEW

075 Here, we present a concise overview of stochastic optimization methods.

076  
077 An optimization method that solves the problem 1 by accessing gradient information from a sub-  
078 set of samples is called SGD. SGD and its variance reduction variants, which operate on only a  
079 small mini-batch of data at each iteration, have become the preferred methods (Robbins & Monro,  
080 1951; Moulines & Bach, 2011; Johnson & Zhang, 2013; Allen-Zhu, 2018). However, stochastic  
081 optimizers sacrifice stability in favor of speed. Parameters such as the learning rate are challenging  
082 to choose (Nemirovski et al., 2009), and for ill-conditioned large-scale machine learning problems,  
083 even finding the optimal learning rate can lead to very slow convergence. Second-order optimizers  
084 based on the Hessian, such as Newton’s method (Battiti, 1992) and quasi-Newton methods (Dennis  
085 & Moré, 1977; Jin & Mokhtari, 2023), are the classic remedy for solving above challenges. Some  
086 researchers have proposed using stochastic Hessian approximations while still utilizing the full gra-  
087 dient (Lacotte et al., 2021; Tong et al., 2021). Then, Frangella et al. (2022) propose the SketchySGD  
088 algorithm whose excellent performance suggests it could potentially replace SGD.

089  
090 When the gradient is difficult to calculate or cannot be obtained, researchers shift their attention from  
091 the study of SGD to stochastic zeroth-order optimization algorithms, estimating the gradient using  
092 function value differences (Ghadimi & Lan, 2013; Duchi et al., 2015; Nesterov & Spokoiny, 2017).  
093 Malladi et al. (2023) directly use zeroth-order optimizer (ZOO) for fine-tuning LLMs. However, the  
094 zeroth-order optimization algorithms mentioned above overlook the use of higher-order information  
095 about the objective, leading to less competitive convergence in practice. Similar to the development  
096 of SGD, researchers have begun to introduce second-order Hessian information into zeroth-order  
097 optimization algorithms. This idea holds promise for the design of efficient and competitive algo-  
098 rithms. Chen et al. (2017) utilize the second-order Hessian information in a relatively coarsened  
099 manner. Ye et al. (2018) take a first step to efficiently incorporate second-order Hessian information  
100 of the objective function and propose a novel class of algorithms called the ZOHA algorithm. Zhao  
101 et al. (2024) propose HiZOO, which is the first work to leverage the diagonal Hessian to enhance  
ZOO for fine-tuning LLMs.

102 It is worth noting that Nesterov & Spokoiny (2017) conduct a theoretical analysis of the com-  
103 plexity bounds for three random gradient-free oracles. However, the potential advantages of us-  
104 ing Gaussian preconditioning vectors remain unexplored. The essential reason is that they do  
105 not effectively utilize the information from the Hessian matrix in their theoretical analysis pro-  
106 cess. Therefore, in essence, our work is different from that of (Nesterov & Spokoiny, 2017).  
107 In addition, we would like to highlight some differences between previous works and ours. Al-  
though Malladi et al. (2023) propose the descent theorem for ZO-SGD, the ultimately proven con-

vergence rate  $t = \mathcal{O}\left((r+1) \cdot \left(\frac{L}{\mu} + \frac{L\alpha}{\mu^2|\mathcal{S}|}\right) \log \frac{f(\mathbf{x}^0) - f^*}{\epsilon}\right)$  is essentially that of gradient descent rather than stochastic gradient descent, where  $r$  is effective rank and  $\alpha$  is used to control the covariance of the gradient estimation. So, Malladi et al. (2023) do not reveal the true convergence rate of ZSG, and determining it remains a challenging open problem. Yue et al. (2023) enhance the convergence rate  $\mathcal{O}\left(\frac{\sum_{i=1}^d \lambda_i(\nabla^2 f(\mathbf{x}))}{\mu} \log \frac{1}{\epsilon}\right)$  for standard zeroth-order optimization algorithm (Nesterov & Spokoiny, 2017) on quadratic objectives, where  $\lambda_i$  represents the  $i$ -th eigenvalue of  $\nabla^2 f(\mathbf{x})$ . Similarly, their work is based on gradient descent rather than SGD. The iterative algorithm  $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C\left(\frac{f_i(\mathbf{x}^t + \alpha \mathbf{s}_t) - f_i(\mathbf{x}^t - \alpha \mathbf{s}_t)}{2\alpha} \mathbf{s}_t\right) + \mathbf{u}_t\right)$  proposed by Zhang et al. (2023) still relies on full gradient information rather than stochastic gradient information. In summary, our theoretical analysis is different from above works. The theoretical result we obtain is unique.

## 1.2 CONTRIBUTIONS

The main contributions are summarized as follows:

- Compared to ZSC, we establish an accelerated convergence rate for ZSG. We successfully reveal that ZSG also exhibits weak dimensional dependence, which explains the fundamental reason behind its superior empirical performance. The advantage of ZSG becomes more pronounced under skewed Hessian spectra. It is worth noting that, in practice, this condition often holds because the singular values of Hessian matrices tend to decrease rapidly (Yue et al., 2023). To the best of our knowledge, our work is the first to theoretically analyze the potential advantages of ZSG compared to ZSC and our conclusion is novel.
- Our theoretical analysis is based on the quadratic regularity assumptions. This assumption helps leverage Hessian structure and broadens the applicability of our complexity results beyond the quadratic case.
- Our research indicates that ZSG also exhibits weak dimensional dependence, similar to zeroth-order gradient descent. This fills a theoretical gap in the field of zeroth-order optimization, and our analytical results provide significant theoretical insights.
- Extensive experiments confirm the reliability of our theoretical analysis. On both synthetic and real-world datasets, the performance of ZSG outperforms that of ZSC. This observation is in line with established practices in the optimization community.

## 2 NOTATION AND ASSUMPTIONS

Let us define the weighted Euclidean norm and weighted inner product associated with a positive definite weight matrix  $\mathbf{M} \succ 0$

$$\|\mathbf{x}\|_{\mathbf{M}} \stackrel{def}{=} \langle \mathbf{x}, \mathbf{x} \rangle_{\mathbf{M}}^{\frac{1}{2}}, \quad \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{M}} \stackrel{def}{=} \langle \mathbf{M}\mathbf{x}, \mathbf{y} \rangle.$$

We define the stochastic gradient  $\nabla f(\mathbf{x}, \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \nabla f_j(\mathbf{x})$ , where  $\mathcal{S}$  represents the sample set and  $|\mathcal{S}|$  represents the sample size.

A widely accepted notion is that the assumptions of  $f$  being  $L$ -smooth and  $\mu$ -strongly convex are standard in the analysis of stochastic gradient methods for solving the problem 1. As research on stochastic algorithms deepens, many researchers have proposed more generalized assumptions. Hanzely et al. (2018) introduce the  $\mathbf{M}$ -smoothness assumption, which is a common assumption in modern analyses of stochastic methods. Gower et al. (2019) present the relative smoothness assumption and relative convexity assumption to exploit information from the Hessian matrix. Frangella et al. (2023) utilize the quadratic regularity assumption to overcome the dilemma of infrequent preconditioner updates. Frangella et al. (2022) propose the relative quadratic regularity assumption, which replaces the Hessian matrix with any positive definite matrix.

Based on these developments, we present the following assumptions on the objective function  $f$ . First, we introduce the quadratic regularity assumption (Frangella et al., 2023), which generalizes classical notions of smoothness and strong convexity to the Hessian norm, thereby enabling the incorporation of rich Hessian information into the complexity analysis.

**Algorithm 1** ZSG: ZO-SGD-Gauss Method

---

**Input and Initialize:** parameters  $\mathbf{x} \in \mathbb{R}^d$ , loss function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , step budget  $t$ , step size  $\eta_t > 0$ , perturbation scale  $\alpha$ , sample distribution  $\mathcal{D}$ , initial point  $\mathbf{x}^0 \in \mathbb{R}^d$

**for**  $t = 0, 1, \dots$  **do**

Sample  $\mathcal{S}_t \sim \mathcal{D}$  and  $\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$

Query the zeroth-order oracle  $f_+^t = f(\mathbf{x}^t + \alpha \mathbf{u}_t, \mathcal{S}_t)$

Query the zeroth-order oracle  $f_-^t = f(\mathbf{x}^t - \alpha \mathbf{u}_t, \mathcal{S}_t)$

Estimating the gradient  $\hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t) = \frac{(f_+^t - f_-^t)}{2\alpha} \cdot \mathbf{u}_t$

$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)$

**end for**

---

**Assumption 2.1.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a twice differentiable function, and let  $\mathbf{M}$  denote the Hessian matrix of  $f$ . The function  $f$  is said to be upper quadratically regular with respect to  $\mathbf{M}$ , if there exists a global constant  $0 \leq \gamma_u < \infty$ , such that for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^d$ ,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\gamma_u}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{M}(\mathbf{z})}^2. \quad (2)$$

Similarly,  $f$  is said to be lower quadratically regular with respect to  $\mathbf{M}$ , if there exists a global constant  $\gamma_l > 0$ , such that for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^d$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\gamma_l}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{M}(\mathbf{z})}^2. \quad (3)$$

Next, we introduce the standard variance assumption.

**Assumption 2.2.** The variance of the stochastic gradient can be bounded by  $\sigma^2$ , which means

$$\mathbb{E} \left[ \|\nabla f(\mathbf{x}, \mathcal{S}) - \nabla f(\mathbf{x})\|^2 \right] \leq \sigma^2. \quad (4)$$

### 3 ALGORITHM DESCRIPTION

This section commences with a detailed description of ZSG the algorithm. According to the formulation in Nesterov & Spokoiny (2017), the zeroth-order gradient estimator can be expressed as  $\hat{\nabla} f(\mathbf{x}, \mathcal{S}) = \frac{[f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S}) - f(\mathbf{x} - \alpha \mathbf{u}, \mathcal{S})]}{2\alpha} \cdot \mathbf{u}$ , where  $\mathbf{u} \in \mathbb{R}^d$  is sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\alpha$  is a very small perturbation scale. In order to help us prove complexity, we need to find the connection between the zeroth-order oracles and the gradient.

**Lemma 3.1.** *We access to the  $f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S})$  and  $f(\mathbf{x} - \alpha \mathbf{u}, \mathcal{S})$ . Through the upper quadratically regular assumption, we yield the following equality*

$$\hat{\nabla} f(\mathbf{x}, \mathcal{S}) = \mathbf{u} \mathbf{u}^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}, \alpha, \mathbf{x}), \quad (5)$$

with

$$\|\phi(\mathbf{u}, \alpha, \mathbf{x})\| \leq \frac{\gamma_u \alpha}{2} \|\mathbf{u}\|_{\mathbf{M}(\mathbf{z})}^2 \cdot \|\mathbf{u}\|, \quad (6)$$

where  $\mathbf{z}_1 \in (\mathbf{x}, \mathbf{x} + \alpha \mathbf{u})$ ,  $\mathbf{z}_2 \in (\mathbf{x} - \alpha \mathbf{u}, \mathbf{x})$  and  $\mathbf{M}(\mathbf{z}) = \begin{cases} \mathbf{M}(\mathbf{z}_1) & \text{if } \mathbf{M}(\mathbf{z}_1) \succeq \mathbf{M}(\mathbf{z}_2) \\ \mathbf{M}(\mathbf{z}_2) & \text{otherwise} \end{cases}$ .

The detailed proof is presented in C.1. The aforementioned relationships can help us conduct convergence analysis. This paper focuses on analyzing the convergence properties of the following update rule:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t). \quad (7)$$

The main algorithmic procedure of ZSG is provided in Algorithm 1. We can also obtain the ZSC algorithm by simply substituting  $\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  with  $\mathbf{e}_t \sim \mathcal{U}^d$ , where  $\mathcal{U}^d$  denotes the uniform distribution over the standard basis vectors in  $\mathbb{R}^d$ . The main algorithmic procedure of ZSC is provided in Algorithm 2.

## 4 MAIN THEORETICAL RESULTS

This section provides an in-depth examination of the iterative complexity of ZSG under the assumptions we introduced. First, we study the convergence properties of quadratic functions. To explain the superiority of ZSG conveniently, we assume that  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{M}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ . If the objective function  $f$  in Assumption 2.1 is quadratic function, we need to point that  $\gamma_l = \gamma_u = 1$  and  $\mathbf{M}(\mathbf{z}) \equiv \mathbf{M}$ , meaning the Hessian matrix is independent of the iteration points.

We begin by presenting several essential lemmas that support the derivation of the main theorems in this section. The detailed proofs of Lemma 4.1 and Lemma 4.2 are provided in Section C. In addition, several other lemmas along with their proofs are given in Section B. The complete proofs of the main theorems and corollaries are deferred to Section D and Section E.

**Lemma 4.1.** *Consider an arbitrary point  $\mathbf{x} \in \mathbb{R}^d$  and a Gaussian vector  $\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . For any  $t > 0$ , the zeroth-order approximation of the gradient at  $\mathbf{x}$  admits the following upper bound:*

$$\mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}}^2 \right] \leq 3\text{tr}(\mathbf{M}) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2. \quad (8)$$

**Lemma 4.2.** *Let  $f^*$  denote the optimum of the objective function. For all  $t > 0$ , if  $\mathbf{z} \in (\mathbf{x}^t, \mathbf{x}^*)$ , the difference between the function value at  $\mathbf{x}^t$  and the optimum  $f^*$  can be bounded as follows:*

$$f(\mathbf{x}^t) - f^* \leq \frac{1}{2\gamma_l \lambda_{\min}(\mathbf{M}(\mathbf{z}))} \|\nabla f(\mathbf{x}^t)\|^2. \quad (9)$$

**Theorem 4.3.** *Let  $f$  be a quadratic function, and assume that  $f$  is both upper and lower quadratically regular with respect to  $\mathbf{M}$ . That is, Assumption 2.1 holds. In addition, the variance of stochastic gradient is bounded, i.e., Assumption 2.2 holds. Suppose the update rule of  $\mathbf{x}^{t+1}$  follows Eq. 7. We define  $P_1(\alpha^2) = \frac{[1+2\lambda_{\max}(\mathbf{M})\eta]\lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})}$  and choose  $\eta_t \equiv \eta \leq \frac{1}{12\text{tr}(\mathbf{M})}$ , then, we obtain*

$$\mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] \leq \frac{6\eta\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}(\mathbf{M})} + P_1(\alpha^2) + \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right]^t [f(\mathbf{x}^0) - f^*].$$

We can observe that ZSG converges to a ball around the optimum from Theorem 4.3. This phenomenon is analogous to the classic SGD which employs a fixed step size (Moulines & Bach, 2011).

**Corollary 4.4.** *Theorem 4.3 suggests that with a fixed step size, the algorithm may fail to converge in the presence of noise. Assume that  $f$  and the parameters satisfy the conditions specified in Theorem 4.3. Since we can choose a sufficiently small  $\alpha$  in practice, we can omit it. If  $\sigma^2 = 0$ , to find an  $\varepsilon$ -suboptimal solution, the iteration complexity is*

$$t = \mathcal{O} \left( \frac{\text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon} \right). \quad (10)$$

When  $\sigma^2 = 0$ , the update of  $\mathbf{x}$  depends on the full gradient, reducing to the deterministic setting. The result in Wang et al. (2024) can be viewed as an intermediate product of our analysis. Their purpose is to compare it with the coordinate-sketched SEGAs (Hanzely et al., 2018), which achieves an iteration complexity of  $\mathcal{O} \left( \frac{d\lambda_{\max}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon} \right)$  without importance sampling. However, our work focus on the analysis of zeroth-order stochastic optimization. The following theorem and corollary will indicate that ZSG outperforms ZSC.

**Theorem 4.5.** *Let  $f$  be a quadratic function, and suppose that Assumption 2.1 and Assumption 2.2 hold. Suppose the update rule of  $\mathbf{x}^{t+1}$  follows Eq. 7. Assume the step size follows the decreasing form  $\eta_t = \frac{l}{\gamma+t}$ , where  $\gamma > 0$  and the intermediate parameter  $l = \frac{3}{\lambda_{\min}(\mathbf{M})}$ . Let  $t_{\max} = T$  and  $\alpha \leq \sqrt{\frac{\alpha_0}{T+1}}$ , where  $\alpha_0$  is a tunable perturbation scale. We define  $Q_1(\alpha_0^2) = \frac{3[6\lambda_{\max}(\mathbf{M})+36\text{tr}(\mathbf{M})]\lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha_0^2}{4\lambda_{\min}^2(\mathbf{M})}$ . The initial step size satisfies  $\eta_0 = \frac{l}{\gamma} \leq \frac{1}{12\text{tr}(\mathbf{M})}$ , which implies  $\gamma \geq \frac{36\text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})}$ . Next, define the auxiliary parameter  $v = \max \left\{ \gamma(f(\mathbf{x}^0) - f^*), \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + Q_1(\alpha_0^2) \right\}$ . Then, for every integer  $t$  with  $0 \leq t \leq T$ , we can obtain the following result*

$$\mathbb{E} [f(\mathbf{x}^t) - f^*] \leq \frac{v}{\gamma + t}.$$

**Corollary 4.6.** *Theorem 4.5 implies that with a decreasing step size, ZSG converges in the presence of noise. Assume that the function  $f$  and the parameters satisfy the conditions specified in Theorem 4.5. Then, to obtain an  $\varepsilon$ -suboptimal solution, the iteration complexity satisfies*

$$t = \mathcal{O} \left( \left[ \frac{\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + Q_1(\alpha_0^2) \right] \frac{1}{\varepsilon} \right). \quad (11)$$

When  $\sigma^2 > 0$  and a sufficiently small  $\alpha_0$  is chosen in practice, the iteration complexity of ZSG is given by  $\mathcal{O} \left( \frac{\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon} \right)$ . In this case, the zeroth-order oracle is queried twice per iteration. Therefore, the query complexity of ZSG is also  $\mathcal{O} \left( \frac{\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon} \right)$ . Notably, the iteration complexity of SGD is  $\mathcal{O} \left( \frac{\lambda_{\max}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon} \right)$  (Rakhlin et al., 2011). The prevailing view in the optimization community (Ghadimi & Lan, 2013; Nesterov & Spokoiny, 2017) is that zeroth-order optimization methods typically have  $\mathcal{O}(d)$  times the complexity of their first-order counterparts. Thus, the query complexity of ZSG becomes  $\mathcal{O} \left( \frac{d\lambda_{\max}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon} \right)$ . To better elucidate our theoretical contribution, we rigorously establish that the ZSG algorithm achieves an iteration complexity of  $\mathcal{O} \left( \left[ \frac{d\lambda_{\max}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + C_1(\alpha_0^2) \right] \frac{1}{\varepsilon} \right)$ , where  $C_1(\alpha_0^2) = \frac{3[6\lambda_{\max}(\mathbf{M})+12d\lambda_{\max}(\mathbf{M})]\lambda_{\max}^2(\mathbf{M})d^2\alpha_0^2}{4\lambda_{\min}^2(\mathbf{M})}$ . The core of the proof lies in substituting  $\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  with  $\mathbf{e}_t \sim \mathcal{U}^d$ , where  $\mathcal{U}^d$  denotes the uniform distribution over the standard basis vectors in  $\mathbb{R}^d$ . Accordingly, we prove that ZSG enjoys a superior convergence rate over ZSC, and this advantage becomes more pronounced under skewed Hessian spectra.

Next, we generalize our results to other function classes satisfying Assumption 2.1, where  $\gamma_l \neq 1$  or  $\gamma_u \neq 1$  may hold.

**Theorem 4.7.** *Suppose  $f$  is in the general form described in problem 1 and assumption 2.1,2.2 hold. Suppose the update rule of  $\mathbf{x}^{t+1}$  follows Eq. 7, using a fixed step size*

$$\eta_t \equiv \eta \leq \frac{1}{12\gamma_u \text{tr}(\mathbf{M})}. \quad (12)$$

We define  $\text{tr}(\mathbf{M}) = \max_{\mathbf{z}^t} \text{tr}(\mathbf{M}(\mathbf{z}^t))$ , with similar definitions for  $\lambda_{\min}(\mathbf{M})$  and  $\lambda_{\max}(\mathbf{M})$ . We also define  $P_2(\alpha^2) = \frac{[1+2\gamma_u\lambda_{\max}(\mathbf{M})\eta]\lambda_{\max}^2(\mathbf{M})(6+d)^3\gamma_u^2\alpha^2}{4\gamma_l\lambda_{\min}(\mathbf{M})}$ . Then, it holds that

$$\mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] \leq \frac{6\eta\gamma_u \text{tr}(\mathbf{M})\sigma^2}{\gamma_l\lambda_{\min}(\mathbf{M})} + P_2(\alpha^2) + \left[ 1 - \frac{1}{2}\eta\gamma_l\lambda_{\min}(\mathbf{M}) \right]^t [f(\mathbf{x}^0) - f^*].$$

If  $\sigma^2 = 0$  and  $\alpha$  is sufficiently small, the complexity to achieve an  $\varepsilon$ -suboptimal solution satisfies

$$t = \mathcal{O} \left( \frac{\gamma_u \text{tr}(\mathbf{M})}{\gamma_l\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon} \right). \quad (13)$$

As shown in Theorem 4.7, ZSG under the full-gradient setting achieves a faster convergence rate than the coordinate-sketched SEGA when  $\frac{\gamma_u}{\gamma_l} = \mathcal{O}(1)$ . This intermediate result extends the result in Wang et al. (2024) from quadratic functions to a broader class of objective functions.

**Theorem 4.8.** *Suppose  $f$  is in the general form described in problem 1 and assumption 2.1,2.2 hold. Suppose the update rule of  $\mathbf{x}^{t+1}$  follows Eq. 7. Assume the step size follows the decreasing form  $\eta_t = \frac{l}{\gamma+t}$ , where  $\gamma > 0$  and the intermediate parameter  $l = \frac{3}{\gamma_l\lambda_{\min}(\mathbf{M})}$ . Let  $t_{\max} = T$  and  $\alpha \leq \sqrt{\frac{\alpha_0}{T+1}}$ , where  $\alpha_0$  is a tunable perturbation scale. We define  $Q_2(\alpha_0^2) = \frac{3[6\gamma_u\lambda_{\max}(\mathbf{M})+36\gamma_u\gamma_l\text{tr}(\mathbf{M})]\lambda_{\max}^2(\mathbf{M})(6+d)^3\gamma_u^2\alpha_0^2}{4\gamma_l^2\lambda_{\min}^2(\mathbf{M})}$ . The initial step size satisfies  $\eta_0 = \frac{l}{\gamma} \leq \frac{1}{12\gamma_u \text{tr}(\mathbf{M})}$ , which implies  $\gamma \geq \frac{36\gamma_u \text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})}$ . Next, define the auxiliary parameter  $v = \max \{ \gamma(f(\mathbf{x}^0) - f^*), \frac{54\gamma_u \text{tr}(\mathbf{M})\sigma^2}{\gamma_l^2\lambda_{\min}^2(\mathbf{M})} + Q_2(\alpha_0^2) \}$ . We define  $\text{tr}(\mathbf{M}) = \max_{\mathbf{z}^t} \text{tr}(\mathbf{M}(\mathbf{z}^t))$ , with similar definitions for  $\lambda_{\min}(\mathbf{M})$  and  $\lambda_{\max}(\mathbf{M})$ . Then, for every integer  $t$  with  $0 \leq t \leq T$ , we can obtain*

$$\mathbb{E} [f(\mathbf{x}^t) - f^*] \leq \frac{v}{\gamma + t}.$$

324 Finally, to obtain an  $\varepsilon$ -suboptimal solution, the iteration complexity satisfies

$$325 \quad t = \mathcal{O} \left( \left[ \frac{\gamma_u \text{tr}(\mathbf{M}) \sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} + Q_2(\alpha_0^2) \right] \frac{1}{\varepsilon} \right). \quad (14)$$

326 As shown in Theorem 4.8, when  $\sigma^2 > 0$  and a sufficiently small  $\alpha_0$  is chosen, the query  
 327 complexity of ZSG is given by  $\mathcal{O} \left( \frac{\gamma_u \text{tr}(\mathbf{M}) \sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon} \right)$ . We rigorously establish that the ZSC algo-  
 328 rithm achieves an iteration complexity of  $\mathcal{O} \left( \left[ \frac{\gamma_u d \lambda_{\max}(\mathbf{M}) \sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} + C_2(\alpha_0^2) \right] \frac{1}{\varepsilon} \right)$ , where  $C_2(\alpha_0^2) =$   
 329  $\frac{3[6\gamma_u \lambda_{\max}(\mathbf{M}) + 12\gamma_u \gamma_l d \lambda_{\max}(\mathbf{M})] \lambda_{\max}^2(\mathbf{M}) d^2 \gamma_u^2 \alpha_0^2}{4\gamma_l^2 \lambda_{\min}^2(\mathbf{M})}$ .

330 The proof of ZSC proceeds along the same lines as the previous analysis. For general functions,  
 331 we accordingly prove that ZSG enjoys a superior convergence rate over ZSC, and this advantage  
 332 also becomes more pronounced under skewed Hessian spectra. For quadratic functions, it follows  
 333 that  $\gamma_u = \gamma_l = 1$ , which is consistent with our previous analysis in Theorem 4.5. That is, based on  
 334 assumption 2.1, we establish that ZSG exhibits weak dimensional dependence for both quadratic and  
 335 general non-quadratic objectives, thereby providing a novel theoretical explanation for the empirical  
 336 observation that ZSG tends to outperform ZSC in practice.

## 342 5 EXPERIMENTS

343 In the preceding sections, we have conducted a comprehensive theoretical analysis of ZSG and  
 344 ZSC, highlighting their respective convergence behaviors. This section is dedicated to the empirical  
 345 validation of ZSG’s effectiveness and superiority. We choose  $\alpha = 10^{-6}$  for all experiments.

### 348 5.1 QUADRATIC FUNCTIONS

349 In this part, our experiments will focus on the quadratic minimization problem, whose objective  
 350 function adheres to the form delineated in the problem 1, characterized by

$$351 \quad \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{2n} \mathbf{x}^\top \mathbf{A} \mathbf{A}^\top \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

352 where  $\mathbf{M} = \frac{1}{n} \mathbf{A} \mathbf{A}^\top$ . The parameters of the quadratic function which we construct as follows.  
 353 The dimension of feature vector  $\mathbf{x}$  is  $d$ . We set  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top$ , where  $\mathbf{U}$  is obtained from QR  
 354 decomposition of a random matrix with entries sampled independently from  $\mathcal{N}(0, 1)$ , and  $\mathbf{\Sigma}$  is set  
 355 as in Table 1. The vector  $\mathbf{b}$  is generated with entries independently drawn from  $\mathcal{N}(0, 1)$ . For each  
 356 problem instance, the initial point is randomly initialized from  $\mathcal{N}(0, 1)$  as well.

357 The decreasing step sizes for both algorithms are set appropriately. According to the theoretical  
 358 results of ZSG and ZSC, the step sizes for them are set proportional to  $\mathcal{O}(1/(\text{tr}(\mathbf{M}) + \lambda_{\min}(\mathbf{M})t))$   
 359 and  $\mathcal{O}(1/(d\lambda_{\max}(\mathbf{M}) + \lambda_{\min}(\mathbf{M})t))$ , respectively. We report the experimental results in Figure 1.

360 When  $d = 100$ , we observe that in the first two experiments, ZSG outperforms ZSC. As  $\text{tr}(\mathbf{M})$   
 361 increases, the performance of ZSG deteriorates. Nevertheless, ZSG remains superior to ZSC under  
 362 Skewed Hessian Spectra.

363 When  $d = 500$ , we find that ZSG significantly outperforms ZSC in the remaining two experiments.  
 364 Similar performance trends are observed as  $\text{tr}(\mathbf{M})$  increases. Notably, as the problem dimension  
 365 grows, the eigenvalues of the Hessian matrix become more diverse, making ZSG increasingly ad-  
 366 vantageous over ZSC. All results are consistent with our theoretical analysis.

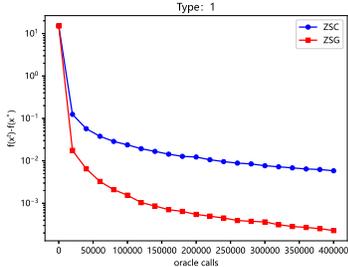
### 371 5.2 LOGISTIC REGRESSION FOR BINARY CLASSIFICATION

372 In this part, we use real datasets to compare the convergence behavior of ZSG and ZSC on strongly  
 373 convex problems. We consider logistic regression the following loss function

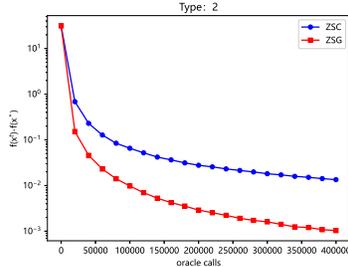
$$374 \quad f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{a}_i, \mathbf{x} \rangle)] + \frac{\beta}{2} \|\mathbf{x}\|^2,$$

Table 1: Setting of diagonal matrix  $\Sigma$  used to construct  $\mathbf{A}$ .

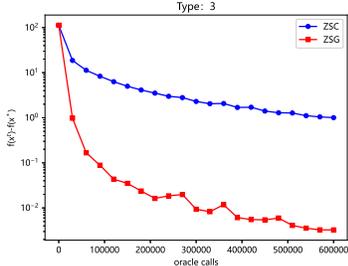
Type	$\Sigma$
1	$d = 100$ Matrix with first 99 components equal to 10 and the remaining one equal to $10\sqrt{10}$
2	$d = 100$ Matrix with first 80 components equal to 10 and the rest equal to $10\sqrt{10}$
3	$d = 500$ Matrix with first 499 components equal to $10\sqrt{5}$ and the remaining one equal to $100\sqrt{5}$
4	$d = 500$ Matrix with first 480 components equal to $10\sqrt{5}$ and the rest equal to $100\sqrt{5}$



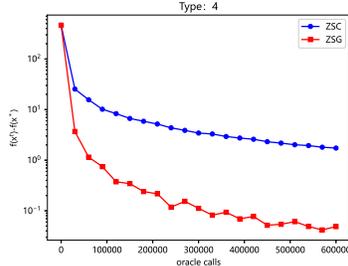
(a) The comparison on the first type diagonal matrix of Table 1



(b) The comparison on the second type diagonal matrix of Table 1



(c) The comparison on the third type diagonal matrix of Table 1



(d) The comparison on the fourth type diagonal matrix of Table 1

Figure 1: Comparison of running results of ZSG and ZSC on quadratic functions.

where  $\mathbf{a}_i \in \mathbb{R}^d$  denotes the  $i$ -th input vector,  $y_i \in \{-1, 1\}$  is the corresponding label and  $\beta$  is the regularization parameter. We conduct experiments on the ‘mushrooms’, ‘phishing’ and ‘a8a’ datasets, with  $d = 112, 68$  and  $123$ , respectively. All three datasets can be downloaded from libsvm datasets. Through the analysis in Section F, we find that all of these datasets fall into the skewed-Hessian setting. In our experiments on ‘mushrooms’ and ‘phishing’, we divide the training set and test set in a ratio of 4:1 and set  $\beta = 0.001$ . We properly choose the decreasing step sizes of them. We report the experimental results in Figure 2.

The first row of subfigures presents the training loss across all experiments. It can be observed that ZSG outperforms ZSC. The second row displays the corresponding test accuracy. The test accuracy achieved by ZSG is more competitive across all experiments. Therefore, the results consistently demonstrate the superiority of ZSG. Intuitively, ZSG benefits from simultaneously incorporating all coordinates in each oracle call, whereas ZSC estimates gradients along individual coordinates. In addition, we conduct a sensitivity analysis of  $\alpha$  in Section G. We find that a large  $\alpha$  significantly hinders the convergence process, while  $\alpha = 10^{-2}$  is already sufficient for most cases.

## 6 CONCLUSION AND FUTURE WORK

In this paper, we are the first to theoretically analyze the potential advantages of ZSG compared to ZSC and obtain the best result from quadratic functions to a broader class of objective func-

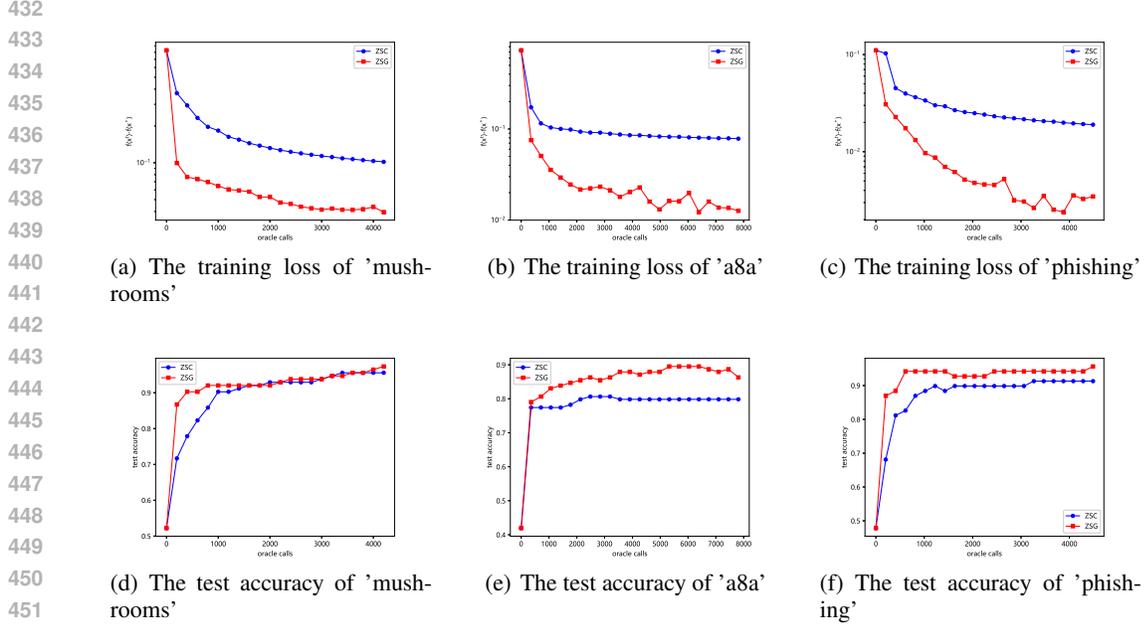


Figure 2: Comparison of running results of ZSG and ZSC on binary classification problem.

tions. When  $\sigma^2 = 0$  and the objective function is quadratic, we recover the main conclusion proposed by Wang et al. (2024): the complexity of ZSG, given by  $\mathcal{O}\left(\frac{\text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon}\right)$ , outperforms the coordinate-sketched SEGA algorithm, whose complexity is  $\mathcal{O}\left(\frac{d\lambda_{\max}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon}\right)$  in the field of zeroth-order optimization. Notably, leveraging the quadratic regularity assumption, we extend the result of Wang et al. (2024) beyond quadratic functions to a broader class of objectives, and rigorously establish an  $\mathcal{O}\left(\frac{\gamma_u \text{tr}(\mathbf{M})}{\gamma_l \lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon}\right)$  complexity guarantee for ZSG. When  $\sigma^2 > 0$  and the objective function is quadratic, we establish the main conclusion of our paper: the query complexity of ZSG is  $\mathcal{O}\left(\frac{\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon}\right)$ , which outperforms ZSC algorithm, whose query complexity is  $\mathcal{O}\left(\frac{d\lambda_{\max}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon}\right)$ . We further establish that, even for non-quadratic objectives, ZSG achieves a complexity of  $\mathcal{O}\left(\frac{\gamma_u \text{tr}(\mathbf{M})\sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon}\right)$ , outperforming the  $\mathcal{O}\left(\frac{\gamma_u d\lambda_{\max}(\mathbf{M})\sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} \frac{1}{\varepsilon}\right)$  complexity of ZSC. ZSG also exhibits weak dimensional dependence and demonstrates a notable advantage in practice, primarily skewed Hessian spectra are commonly observed in real-world problems (Yue et al., 2023).

The upper and lower quadratic regularity constants enable us to generalize the results from quadratic functions to a broader class of objective functions, although  $\gamma_u$  and  $\gamma_l$  are indeed difficult to control. Additional assumptions or a more refined analysis may be required to verify whether  $\frac{\gamma_u}{\gamma_l} = \mathcal{O}(1)$  or  $\frac{\gamma_u}{\gamma_l^2} = \mathcal{O}(1)$  holds, thereby guiding better parameter tuning strategies in future settings. Frangella et al. (2023) point out that for any objective with a Lipschitz Hessian,  $\frac{\gamma_u}{\gamma_l}$  approaches 1 as the optimal objective value is approached. This insight helps explain the empirical advantage of ZSG in the fine-tuning of large language models. Moreover, our experimental evaluation includes not only quadratic objectives but also logistic regression tasks. These results highlight the practical significance of extending our convergence analysis beyond the quadratic setting.

In addition, a promising research direction is to incorporate Hessian information into the gradient estimator  $\hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)$ . The motivation stems from the observation that a significant difference in the curvature of the loss function can lead to training instability or slow convergence. Hessian information can be leveraged to adaptively scale parameter updates, thereby mitigating this issue. We expect that integrating such techniques into our analytical framework could lead to improved practical performance in terms of query complexity.

## REFERENCES

- 486  
487  
488 Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. *Journal*  
489 *of Machine Learning Research*, 18(221):1–51, 2018.
- 490 Yang Bai, Yisen Wang, Yuyuan Zeng, Yong Jiang, and Shu-Tao Xia. Query efficient black-box  
491 adversarial attack on deep neural networks. *Pattern Recognition*, 133:109037, 2023.
- 492  
493 Roberto Battiti. First-and second-order methods for learning: between steepest descent and newton’s  
494 method. *Neural computation*, 4(2):141–166, 1992.
- 495 Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order opti-  
496 mization based black-box attacks to deep neural networks without training substitute models. In  
497 *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 15–26, 2017.
- 498 John E Dennis, Jr and Jorge J Moré. Quasi-newton methods, motivation and theory. *SIAM review*,  
499 19(1):46–89, 1977.
- 500  
501 John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono. Optimal rates for  
502 zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on*  
503 *Information Theory*, 61(5):2788–2806, 2015.
- 504 Zachary Frangella, Pratik Rathore, Shipu Zhao, and Madeleine Udell. Sketchysgd: reliable stochas-  
505 tic optimization via randomized curvature estimates. *arXiv preprint arXiv:2211.08597*, 2022.
- 506  
507 Zachary Frangella, Pratik Rathore, Shipu Zhao, and Madeleine Udell. Promise: Preconditioned  
508 stochastic optimization methods by incorporating scalable curvature estimates. *arXiv preprint*  
509 *arXiv:2309.02014*, 2023.
- 510 Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochas-  
511 tic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.
- 512  
513 Robert Gower, Dmitry Kovalev, Felix Lieder, and Peter Richtárik. Rsn: randomized subspace new-  
514 ton. *Advances in Neural Information Processing Systems*, 32, 2019.
- 515 Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. Sega: Variance reduction via gradient  
516 sketching. *Advances in Neural Information Processing Systems*, 31, 2018.
- 517  
518 Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with  
519 limited queries and information. In *International conference on machine learning*, pp. 2137–  
520 2146. PMLR, 2018.
- 521 Qiujiang Jin and Aryan Mokhtari. Non-asymptotic superlinear convergence of standard quasi-  
522 newton methods. *Mathematical Programming*, 200(1):425–473, 2023.
- 523  
524 Gangshan Jing, He Bai, Jemin George, Aranya Chakraborty, and Piyush K Sharma. Asynchronous  
525 distributed reinforcement learning for lqr control via zeroth-order block coordinate descent. *IEEE*  
526 *Transactions on Automatic Control*, 2024.
- 527 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance  
528 reduction. *Advances in neural information processing systems*, 26, 2013.
- 529  
530 Jack Kiefer and Jacob Wolfowitz. Stochastic estimation of the maximum of a regression function.  
531 *The Annals of Mathematical Statistics*, pp. 462–466, 1952.
- 532 Jonathan Lacotte, Yifei Wang, and Mert Pilanci. Adaptive newton sketch: Linear-time optimization  
533 with quadratic convergence and effective hessian dimensionality. In *International Conference on*  
534 *Machine Learning*, pp. 5926–5936. PMLR, 2021.
- 535  
536 Jan R Magnus et al. *The moments of products of quadratic forms in normal variables*. Univ.,  
537 Instituut voor Actuarieat en Econometrie, 1978.
- 538 Dhruv Malik, Ashwin Pananjady, Kush Bhatia, Koulik Khamaru, Peter L Bartlett, and Martin J  
539 Wainwright. Derivative-free methods for policy optimization: Guarantees for linear quadratic  
systems. *Journal of Machine Learning Research*, 21(21):1–51, 2020.

- 540 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev  
541 Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information*  
542 *Processing Systems*, 36:53038–53075, 2023.
- 543
- 544 Horia Mania, Aurelia Guy, and Benjamin Recht. Simple random search of static linear policies is  
545 competitive for reinforcement learning. *Advances in neural information processing systems*, 31,  
546 2018.
- 547 Hesameddin Mohammadi, Mahdi Soltanolkotabi, and Mihailo R Jovanović. On the linear conver-  
548 gence of random search for discrete-time lqr. *IEEE Control Systems Letters*, 5(3):989–994, 2020.
- 549
- 550 Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms  
551 for machine learning. *Advances in neural information processing systems*, 24, 2011.
- 552 Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic  
553 approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–  
554 1609, 2009.
- 555
- 556 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.  
557 *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- 558 Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram  
559 Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM*  
560 *on Asia conference on computer and communications security*, pp. 506–519, 2017.
- 561
- 562 Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for  
563 strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
- 564
- 565 Jérémy Rapin and Olivier Teytaud. Nevergrad-a gradient-free optimization platform, 2018.
- 566
- 566 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathemati-*  
567 *cal statistics*, pp. 400–407, 1951.
- 568
- 568 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a  
569 scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- 570
- 571 Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine  
572 learning algorithms. *Advances in neural information processing systems*, 25, 2012.
- 573
- 573 James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient  
574 approximation. *IEEE transactions on automatic control*, 37(3):332–341, 1992.
- 575
- 576 Zhiwei Tang, Dmitry Rybin, and Tsung-Hui Chang. Zeroth-order optimization meets human feed-  
577 back: Provable learning via ranking oracles. *arXiv preprint arXiv:2303.03751*, 2023.
- 578
- 578 Tian Tong, Cong Ma, and Yuejie Chi. Accelerating ill-conditioned low-rank matrix estimation via  
579 scaled gradient descent. *Journal of Machine Learning Research*, 22(150):1–63, 2021.
- 580
- 581 Yilong Wang, Haishan Ye, Guang Dai, and Ivor Tsang. Can gaussian sketching converge faster on  
582 a preconditioned landscape? In *Forty-first International Conference on Machine Learning*, 2024.
- 583
- 583 Haishan Ye, Zhichao Huang, Cong Fang, Chris Junchi Li, and Tong Zhang. Hessian-aware zeroth-  
584 order optimization for black-box adversarial attack. *arXiv preprint arXiv:1812.11377*, 2018.
- 585
- 586 Pengyun Yue, Long Yang, Cong Fang, and Zhouchen Lin. Zeroth-order optimization with weak  
587 dimension dependency. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 4429–  
588 4472. PMLR, 2023.
- 589
- 589 Liang Zhang, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. Dpzero: dimension-  
590 independent and differentially private zeroth-order optimization. In *International Workshop on*  
591 *Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023*, 2023.
- 592
- 593 Yan Zhang and Michael M Zavlanos. Cooperative multiagent reinforcement learning with partial  
observations. *IEEE Transactions on Automatic Control*, 69(2):968–981, 2023.

594 Yonggang Zhang, Ya Li, Tongliang Liu, and Xinmei Tian. Dual-path distillation: A unified frame-  
595 work to improve black-box attacks. In *International Conference on Machine Learning*, pp.  
596 11163–11172. PMLR, 2020.

597 Pu Zhao, Pin-Yu Chen, Siyue Wang, and Xue Lin. Towards query-efficient black-box adversary  
598 with zeroth-order natural gradient descent. In *Proceedings of the AAAI Conference on Artificial  
599 Intelligence*, volume 34, pp. 6909–6916, 2020.

601 Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor W Tsang. Second-order  
602 fine-tuning without pain for llms: A hessian informed zeroth-order optimizer. *arXiv preprint  
603 arXiv:2402.15173*, 2024.

## 605 A THE USE OF LARGE LANGUAGE MODELS (LLMs)

606 This paper uses large language models only to polish the language and adjust the paragraph structure.  
607

## 608 B SEVERAL USEFUL LEMMAS

609 In this section, we introduce several useful lemmas. The following lemma shows that the expectation  
610 of the product of two quadratic forms of the random Gaussian vector is related to the trace of the  
611 corresponding matrix.

612 **Lemma B.1** (Magnus et al. (1978)). *Let  $\mathbf{A}$  and  $\mathbf{B}$  be two symmetric matrices, and  $\mathbf{u}$  obeys the  
613 Gaussian distribution, that is,  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Define  $z = \mathbf{u}^\top \mathbf{A} \mathbf{u} \cdot \mathbf{u}^\top \mathbf{B} \mathbf{u}$ . The expectation of  $z$  is*

$$614 \mathbb{E}_{\mathbf{u}}[z] = (\text{tr}\mathbf{A})(\text{tr}\mathbf{B}) + 2(\text{tr}\mathbf{A}\mathbf{B}). \quad (15)$$

615 **Lemma B.2** (Nesterov & Spokoiny (2017)). *Let  $\mathbf{u}$  obeys the Gaussian distribution, that is,  
616  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . We define normalization constant  $\kappa = \int e^{-\frac{1}{2}\|\mathbf{u}\|^2} d\mathbf{u}$  and define moments  
617  $\mathbf{M}_p = \frac{1}{\kappa} \int \|\mathbf{u}\|^p e^{-\frac{1}{2}\|\mathbf{u}\|^2} d\mathbf{u}$ . For  $p \geq 2$ , we can obtain upper bounds*

$$618 n^{p/2} \leq \mathbf{M}_p \leq (p+d)^{p/2}. \quad (16)$$

619 **Lemma B.3.** *If we have a positive definite matrix  $\mathbf{M}$  defined as weighted inner product, for all  
620  $\mathbf{x} \in \mathbb{R}^d$ , we can obtain the following inequalities*

$$621 \|\mathbf{x}\|_{\mathbf{M}}^2 \leq \text{tr}(\mathbf{M}) \|\mathbf{x}\|^2, \quad (17)$$

$$622 \lambda_{\min}(\mathbf{M}) \|\mathbf{x}\|^2 \leq \|\mathbf{x}\|_{\mathbf{M}}^2 \leq \lambda_{\max}(\mathbf{M}) \|\mathbf{x}\|^2. \quad (18)$$

623 *Proof.* For a positive definite matrix  $\mathbf{M}$ , there must exist an orthogonal matrix  $\mathbf{T}$  such that  $\mathbf{M}$  is  
624 similar to a diagonal matrix whose elements are eigenvalues of matrix  $\mathbf{M}$ . We denote  $\lambda_i$  be the  $i$ -th  
625 eigenvalue of matrix  $\mathbf{M}$ , then, we can obtain an equation as follows

$$626 \mathbf{M} = \mathbf{T} \text{diag} \{\lambda_1, \lambda_2, \dots, \lambda_d\} \mathbf{T}^{-1}. \quad (19)$$

627 Let  $\mathbf{y} = \mathbf{T}^\top \mathbf{x}$ , then, we can easily prove this Lemma. We first prove Eq. (17)

$$\begin{aligned} 628 \|\mathbf{x}\|_{\mathbf{M}}^2 &= \langle \mathbf{M}\mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^\top \mathbf{M}\mathbf{x} \stackrel{(19)}{=} \mathbf{x}^\top \mathbf{T} \text{diag} \{\lambda_1, \lambda_2, \dots, \lambda_d\} \mathbf{T}^{-1} \mathbf{x} \\ 629 &= \mathbf{x}^\top \mathbf{T} \text{diag} \{\lambda_1, \lambda_2, \dots, \lambda_d\} \mathbf{T}^\top \mathbf{x} \\ 630 &= \mathbf{y}^\top \text{diag} \{\lambda_1, \lambda_2, \dots, \lambda_d\} \mathbf{y} \\ 631 &\leq \text{tr}(\mathbf{M}) \mathbf{x}^\top \mathbf{T} \mathbf{T}^\top \mathbf{x} \\ 632 &= \text{tr}(\mathbf{M}) \|\mathbf{x}\|^2. \end{aligned}$$

633 Similarly, we can prove the Eq. (18). □

634 **Lemma B.4.** *For the sake of simplicity in the subsequent proof, we first derive the upper bound of  
635  $\hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)$ . The upper bound is related to  $\nabla f(\mathbf{x}^t, \mathcal{S}_t)$  and  $\alpha$ :*

$$636 \mathbb{E}_{\mathbf{u}_t} \left[ \|\hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] \leq 6\text{tr}(\mathbf{M}(\mathbf{z}^t)) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2 + \frac{(6+d)^3 \gamma_u^2 \alpha^2}{2}. \quad (20)$$

648 *Proof.* This part of the proof involves the basic properties of the norm and some important lemmas.

$$\begin{aligned}
649 & \mathbb{E}_{\mathbf{u}_t} \left[ \|\hat{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] \stackrel{(5)}{\leq} \mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] \\
650 & \leq 2\mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] + 2\mathbb{E}_{\mathbf{u}_t} \left[ \|\phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] \\
651 & \stackrel{(6)}{\leq} 2\mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}(\mathbf{z}^t)}^2 \right] + \frac{\gamma_u^2 \alpha^2}{2} \mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t\|_{\mathbf{M}(\mathbf{z}^t)}^6 \right] \\
652 & \stackrel{(8)+(16)+(18)}{\leq} 6\text{tr}(\mathbf{M}(\mathbf{z}^t)) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2 + \frac{\lambda_{\max}^3(\mathbf{M}(\mathbf{z}^t))(6+d)^3 \gamma_u^2 \alpha^2}{2}.
\end{aligned}$$

□

653 **Lemma B.5.** For the sake of simplicity in the subsequent proof, we will derive the upper bound of an important inner product  $\langle \nabla f(\mathbf{x}^t), \phi(\mathbf{u}_t, \alpha) \rangle$ . The upper bound is related to real gradient  $\nabla f(\mathbf{x}^t)$  and  $\alpha$ :

$$654 \quad -\mathbb{E}_{\mathbf{u}_t} [\langle \nabla f(\mathbf{x}^t), \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \rangle] \leq \frac{1}{2} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{\lambda_{\max}^2(\mathbf{M}(\mathbf{z}^t))(6+d)^3 \gamma_u^2 \alpha^2}{8}. \quad (21)$$

655 *Proof.* The techniques involved in this part are similar to those in Lemma B.4.

$$\begin{aligned}
656 & -\mathbb{E}_{\mathbf{u}_t} [\langle \nabla f(\mathbf{x}^t), \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \rangle] \leq \mathbb{E}_{\mathbf{u}_t} [\|\nabla f(\mathbf{x}^t)\| \|\phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|] \\
657 & \leq \frac{1}{2} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{1}{2} \mathbb{E}_{\mathbf{u}_t} [\|\phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|^2] \\
658 & \stackrel{(6)}{\leq} \frac{1}{2} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{\gamma_u^2 \alpha^2}{8} \mathbb{E}_{\mathbf{u}_t} [\|\mathbf{u}_t\|_{\mathbf{M}(\mathbf{z}^t)}^4 \cdot \|\mathbf{u}_t\|^2] \\
659 & \stackrel{(18)}{\leq} \frac{1}{2} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{\lambda_{\max}^2(\mathbf{M}(\mathbf{z}^t)) \gamma_u^2 \alpha^2}{8} \mathbb{E}_{\mathbf{u}_t} [\|\mathbf{u}_t\|^6] \\
660 & \stackrel{(16)}{\leq} \frac{1}{2} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{\lambda_{\max}^2(\mathbf{M}(\mathbf{z}^t))(6+d)^3 \gamma_u^2 \alpha^2}{8}.
\end{aligned}$$

□

## 661 C PROOF OF IMPORTANT LEMMAS

662 In this section, we give some details of proof about some important Lemmas.

### 663 C.1 PROOF OF LEMMA 3.1

664 *Proof.* By the Taylor's expansion, we can obtain that

$$665 \quad f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S}) = f(\mathbf{x}) + \alpha \langle \nabla f(\mathbf{x}, \mathcal{S}), \mathbf{u} \rangle + \phi'(\mathbf{u}, \alpha, \mathbf{x})$$

666 where  $\phi'(\mathbf{u}, \alpha, \mathbf{x}) = f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S}) - f(\mathbf{x}) - \alpha \langle \nabla f(\mathbf{x}, \mathcal{S}), \mathbf{u} \rangle$ . Similarly, we can obtain

$$667 \quad f(\mathbf{x} - \alpha \mathbf{u}, \mathcal{S}) = f(\mathbf{x}) - \alpha \langle \nabla f(\mathbf{x}, \mathcal{S}), \mathbf{u} \rangle + \phi'(\mathbf{u}, -\alpha, \mathbf{x}).$$

$$668 \quad \hat{\nabla} f(\mathbf{x}, \mathcal{S}) = \frac{[f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S}) - f(\mathbf{x} - \alpha \mathbf{u}, \mathcal{S})]}{2\alpha} \cdot \mathbf{u} = \mathbf{u} \mathbf{u}^\top \nabla f(\mathbf{x}, \mathcal{S}) + \frac{\phi'(\mathbf{u}, \alpha, \mathbf{x}) - \phi'(\mathbf{u}, -\alpha, \mathbf{x})}{2\alpha} \cdot \mathbf{u}.$$

669 By the upper quadratically regular assumption, we can obtain that

$$670 \quad |\phi'(\mathbf{u}, \alpha, \mathbf{x})| = |f(\mathbf{x} + \alpha \mathbf{u}, \mathcal{S}) - f(\mathbf{x}) - \alpha \langle \nabla f(\mathbf{x}, \mathcal{S}), \mathbf{u} \rangle| \leq \frac{\gamma_u \alpha^2}{2} \|\mathbf{u}\|_{\mathbf{M}(\mathbf{z}_1)}^2,$$

$$671 \quad |\phi'(\mathbf{u}, -\alpha, \mathbf{x})| = |f(\mathbf{x} - \alpha \mathbf{u}, \mathcal{S}) - f(\mathbf{x}) + \alpha \langle \nabla f(\mathbf{x}, \mathcal{S}), \mathbf{u} \rangle| \leq \frac{\gamma_u \alpha^2}{2} \|\mathbf{u}\|_{\mathbf{M}(\mathbf{z}_2)}^2.$$

672 Then, we can finally obtain that

$$673 \quad \left\| \frac{\phi'(\mathbf{u}, \alpha, \mathbf{x}) - \phi'(\mathbf{u}, -\alpha, \mathbf{x})}{2\alpha} \cdot \mathbf{u} \right\| \leq \frac{|\phi'(\mathbf{u}, \alpha, \mathbf{x})| + |\phi'(\mathbf{u}, -\alpha, \mathbf{x})|}{2\alpha} \|\mathbf{u}\| \leq \frac{\gamma_u \alpha}{2} \|\mathbf{u}\|_{\mathbf{M}(\mathbf{z})}^2 \cdot \|\mathbf{u}\|.$$

□

## C.2 PROOF OF LEMMA 4.1

*Proof.* This part of the proof mainly relies on the properties of the matrix trace.

$$\begin{aligned}
\mathbb{E}_{\mathbf{u}_t} \left[ \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}}^2 \right] &= \mathbb{E}_{\mathbf{u}_t} \left[ \nabla f(\mathbf{x}^t, \mathcal{S}_t)^\top \mathbf{u}_t \mathbf{u}_t^\top \mathbf{M}^\top \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t) \right] \\
&= \mathbb{E}_{\mathbf{u}_t} \left[ \text{tr}(\nabla f(\mathbf{x}^t, \mathcal{S}_t)^\top \mathbf{u}_t \mathbf{u}_t^\top \mathbf{M}^\top \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)) \right] \\
&= \mathbb{E}_{\mathbf{u}_t} \left[ \text{tr}(\mathbf{u}_t^\top \mathbf{M}^\top \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t) \nabla f(\mathbf{x}^t, \mathcal{S}_t)^\top \mathbf{u}_t) \right] \\
&\stackrel{(15)}{=} \text{tr}(\mathbf{M}) \text{tr}(\nabla f(\mathbf{x}^t, \mathcal{S}_t) \nabla f(\mathbf{x}^t, \mathcal{S}_t)^\top) \\
&\quad + 2 \text{tr}(\nabla f(\mathbf{x}^t, \mathcal{S}_t)^\top \mathbf{M}^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t)) \\
&= \text{tr}(\mathbf{M}) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2 + 2 \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|_{\mathbf{M}}^2 \\
&\stackrel{(17)}{\leq} 3 \text{tr}(\mathbf{M}) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2.
\end{aligned}$$

□

## C.3 PROOF OF LEMMA 4.2

*Proof.* We use the lower quadratically regular introduced in Assumption 2.1,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\gamma_l}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{M}(\mathbf{z})}^2.$$

Then, we construct an auxiliary function,

$$F(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\gamma_l}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{M}(\mathbf{z})}^2.$$

To obtain the minimum of the auxiliary function, we need to make

$$\nabla F(\mathbf{y}^*) = \nabla f(\mathbf{x}) + 2\gamma_l \mathbf{M}(\mathbf{z})(\mathbf{y}^* - \mathbf{x}) = 0.$$

So, we can find that

$$\mathbf{y}^* = \mathbf{x} - \frac{1}{\gamma_l} \mathbf{M}(\mathbf{z})^{-1} \nabla f(\mathbf{x}). \quad (22)$$

Using the above information, we can continue to deduce that

$$\begin{aligned}
f(\mathbf{y}) &\geq F(\mathbf{y}) \\
&\geq F(\mathbf{y}^*) \\
&\stackrel{(22)}{=} f(\mathbf{x}) - \left\langle \nabla f(\mathbf{x}), \frac{1}{\gamma_l} \mathbf{M}(\mathbf{z})^{-1} \nabla f(\mathbf{x}) \right\rangle + \frac{\gamma_l}{2} \left\| \frac{1}{\gamma_l} \mathbf{M}(\mathbf{z})^{-1} \nabla f(\mathbf{x}) \right\|_{\mathbf{M}(\mathbf{z})}^2 \\
&= f(\mathbf{x}) - \frac{1}{\gamma_l} \|\nabla f(\mathbf{x})\|_{\mathbf{M}(\mathbf{z})^{-1}}^2 + \frac{1}{2\gamma_l} \nabla f(\mathbf{x})^\top (\mathbf{M}(\mathbf{z})^{-1})^\top \mathbf{M}(\mathbf{z})^\top \mathbf{M}(\mathbf{z})^{-1} \nabla f(\mathbf{x}) \\
&= f(\mathbf{x}) - \frac{1}{\gamma_l} \|\nabla f(\mathbf{x})\|_{\mathbf{M}(\mathbf{z})^{-1}}^2 + \frac{1}{2\gamma_l} \|\nabla f(\mathbf{x})\|_{\mathbf{M}(\mathbf{z})^{-1}}^2 \\
&= f(\mathbf{x}) - \frac{1}{2\gamma_l} \|\nabla f(\mathbf{x})\|_{\mathbf{M}(\mathbf{z})^{-1}}^2.
\end{aligned}$$

Let  $\mathbf{x} = \mathbf{x}^t$ ,  $\mathbf{y} = \mathbf{x}^*$ , and rearrange the above formula, we can obtain

$$\begin{aligned}
f(\mathbf{x}^t) - f^* &\leq \frac{1}{2\gamma_l} \|\nabla f(\mathbf{x}^t)\|_{\mathbf{M}(\mathbf{z})^{-1}}^2 \\
&\stackrel{(18)}{\leq} \frac{\lambda_{\max}(\mathbf{M}(\mathbf{z})^{-1})}{2\gamma_l} \|\nabla f(\mathbf{x}^t)\| \\
&= \frac{1}{2\gamma_l \lambda_{\min}(\mathbf{M}(\mathbf{z}))} \|\nabla f(\mathbf{x}^t)\|.
\end{aligned}$$

□

## D PROOF OF MAIN THEOREMS

In this section, we give some details of proof about some important Theorem.

### D.1 PROOF OF THEOREM 4.3

*Proof.* Firstly, we can deduce the expectation of  $f(\mathbf{x}^{t+1})$ ,

$$\begin{aligned} f(\mathbf{x}^{t+1}) &\stackrel{(2)}{\leq} f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x}^{t+1} - \mathbf{x}^t \rangle + \frac{1}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_{\mathbf{M}}^2 \\ &\stackrel{(7)+(5)}{=} f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \rangle \\ &\quad + \frac{\eta^2}{2} \|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|_{\mathbf{M}}^2. \end{aligned} \quad (23)$$

Let us deduce the expectation of  $f(\mathbf{x}^{t+1})$  for  $u$ ,

$$\begin{aligned} \mathbb{E}_{\mathbf{u}_t} [f(\mathbf{x}^{t+1})] &= f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \mathbb{E}_{\mathbf{u}_t} [\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)] \rangle \\ &\quad + \frac{\eta^2}{2} \mathbb{E}_{\mathbf{u}_t} [\|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|_{\mathbf{M}}^2] \\ &\leq f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \nabla f(\mathbf{x}^t, \mathcal{S}_t) \rangle - \eta \mathbb{E}_{\mathbf{u}_t} [\langle \nabla f(\mathbf{x}^t), \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \rangle] \\ &\quad + \frac{\eta^2}{2} \mathbb{E}_{\mathbf{u}_t} [\|\mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}, \mathcal{S}) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t)\|_{\mathbf{M}}^2] \\ &\stackrel{(20)+(21)}{\leq} f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \nabla f(\mathbf{x}^t, \mathcal{S}_t) \rangle + \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|^2 + 3\eta^2 \text{tr}(\mathbf{M}) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2 \\ &\quad + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8}. \end{aligned}$$

Then, let us deduce the expectation of  $\mathbb{E}_u [f(\mathbf{x}^{t+1})]$ ,

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1})] &\leq f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \mathbb{E} [\nabla f(\mathbf{x}^t, \mathcal{S}_t)] \rangle + \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|^2 \\ &\quad + 3\eta^2 \text{tr}(\mathbf{M}) \mathbb{E} [\|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2] + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &= f(\mathbf{x}^t) - \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|^2 + 3\eta^2 \text{tr}(\mathbf{M}) \mathbb{E} [\|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2] \\ &\quad + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &\stackrel{(4)}{\leq} f(\mathbf{x}^t) - \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|^2 + 3\eta^2 \text{tr}(\mathbf{M})(\sigma^2 + \|\nabla f(\mathbf{x}^t)\|^2) \\ &\quad + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &= f(\mathbf{x}^t) + \left[ 3\eta^2 \text{tr}(\mathbf{M}) - \frac{\eta}{2} \right] \|\nabla f(\mathbf{x}^t)\|^2 + 3\eta^2 \sigma^2 \text{tr}(\mathbf{M}) \\ &\quad + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &= f(\mathbf{x}^t) + 3\eta^2 \sigma^2 \text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &\quad - \frac{\eta}{2} [1 - 6\eta \text{tr}(\mathbf{M})] \|\nabla f(\mathbf{x}^t)\|^2 \\ &\stackrel{(9)}{\leq} f(\mathbf{x}^t) + 3\eta^2 \sigma^2 \text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &\quad - \eta \lambda_{\min}(\mathbf{M}) [1 - 6\eta \text{tr}(\mathbf{M})] (f(\mathbf{x}^t) - f^*) \\ &\leq f(\mathbf{x}^t) + 3\eta^2 \sigma^2 \text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta}{8} \\ &\quad - \frac{1}{2} \eta \lambda_{\min}(\mathbf{M}) (f(\mathbf{x}^t) - f^*). \end{aligned} \quad (24)$$

And then, let us use the optimal value  $f^*$  to transform the inequality,

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] + f^* - f(\mathbf{x}^t) &\leq 3\eta^2\sigma^2\text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2\eta}{8} \\ &\quad - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\mathbb{E} [f(\mathbf{x}^t) - f^*]. \end{aligned}$$

Rearranging the above formula, we can obtain,

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] &\leq 3\eta^2\sigma^2\text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2\eta}{8} \\ &\quad + \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right] \mathbb{E} [f(\mathbf{x}^t) - f^*]. \end{aligned}$$

We need to construct a recursive relation with the following structure,

$$\mathbb{E} [f(\mathbf{x}^{t+1}) - f^* - \beta] \leq \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right] \mathbb{E} [f(\mathbf{x}^t) - f^* - \beta].$$

If  $\beta = \frac{24\eta\text{tr}(\mathbf{M})\sigma^2 + [1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})}$ , the above formula can be derived as

$$\begin{aligned} &\mathbb{E} \left[ f(\mathbf{x}^{t+1}) - f^* - \frac{24\eta\text{tr}(\mathbf{M})\sigma^2 + [1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})} \right] \\ &\leq \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right] \mathbb{E} \left[ f(\mathbf{x}^t) - f^* - \frac{24\eta\text{tr}(\mathbf{M})\sigma^2 + [1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})} \right] \\ &\leq \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right]^t \left[ f(\mathbf{x}^0) - f^* - \frac{24\eta\text{tr}(\mathbf{M})\sigma^2 + [1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})} \right] \\ &\leq \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right]^t [f(\mathbf{x}^0) - f^*]. \end{aligned}$$

Thus, we can obtain that

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] &\leq \frac{24\eta\text{tr}(\mathbf{M})\sigma^2 + [1 + 2\lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2}{4\lambda_{\min}(\mathbf{M})} \\ &\quad + \left[1 - \frac{1}{2}\eta\lambda_{\min}(\mathbf{M})\right]^t [f(\mathbf{x}^0) - f^*]. \end{aligned}$$

□

## D.2 PROOF OF THEOREM 4.5

*Proof.* Firstly, if we choose decreasing step size  $\eta_t$ , based on D.1, we can obtain the following formula

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f(\mathbf{x}^t)] &\leq 3\eta_t^2\sigma^2\text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2\eta_t}{8} \\ &\quad - \eta_t\lambda_{\min}(\mathbf{M}) [1 - 6\eta_t\text{tr}(\mathbf{M})] \mathbb{E} [f(\mathbf{x}^t) - f^*] \\ &\leq 3\eta_t^2\sigma^2\text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2\eta_t}{8} \\ &\quad - \eta_t\lambda_{\min}(\mathbf{M}) [1 - 6\eta_0\text{tr}(\mathbf{M})] \mathbb{E} [f(\mathbf{x}^t) - f^*] \\ &\leq 3\eta_t^2\sigma^2\text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3\alpha^2\eta_t}{8} \\ &\quad - \frac{1}{2}\eta_t\lambda_{\min}(\mathbf{M})\mathbb{E} [f(\mathbf{x}^t) - f^*]. \end{aligned}$$

Let us prove the final result by induction, for  $t = 0$

$$\mathbb{E} [f(\mathbf{x}^0) - f^*] = f(\mathbf{x}^0) - f^* = \frac{\gamma}{\gamma + 0} [f(\mathbf{x}^0) - f^*] \leq \frac{v}{\gamma + 0},$$

864 by the definition of  $v$ .

865 Suppose that holds for  $t > 0$ , then

$$\begin{aligned}
866 \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] &\leq 3\eta_t^2 \sigma^2 \text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta_t}{8} \\
867 &\quad + \left[1 - \frac{1}{2}\eta_t \lambda_{\min}(\mathbf{M})\right] \mathbb{E} [f(\mathbf{x}^t) - f^*] \\
868 &\leq 3\eta_t^2 \sigma^2 \text{tr}(\mathbf{M}) + \frac{[1 + 2\lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 \eta_t}{8} \\
869 &\quad + \left[1 - \frac{1}{2}\eta_t \lambda_{\min}(\mathbf{M})\right] \frac{v}{\gamma + t} \\
870 &= \frac{3\sigma^2 l^2 \text{tr}(\mathbf{M})}{(\gamma + t)^2} + \frac{\lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 l}{8(\gamma + t)} + \frac{\lambda_{\max}^3(\mathbf{M})(6+d)^3 \alpha^2 l^2}{4(\gamma + t)^2} \\
871 &\quad + \left[1 - \frac{l\lambda_{\min}(\mathbf{M})}{2(\gamma + t)}\right] \frac{v}{\gamma + t} \\
872 &= \frac{(\gamma + t - 1)v}{(\gamma + t)^2} + \frac{3\sigma^2 l^2 \text{tr}(\mathbf{M})}{(\gamma + t)^2} + \frac{\lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 l}{8(\gamma + t)} \\
873 &\quad + \frac{\lambda_{\max}^3(\mathbf{M})(6+d)^3 \alpha^2 l^2}{4(\gamma + t)^2} - \frac{(l\lambda_{\min}(\mathbf{M}) - 2)v}{2(\gamma + t)^2}.
\end{aligned}$$

874 We let  $\frac{3\sigma^2 l^2 \text{tr}(\mathbf{M})}{(\gamma + t)^2} + \frac{\lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 l}{8(\gamma + t)} + \frac{\lambda_{\max}^3(\mathbf{M})(6+d)^3 \alpha^2 l^2}{4(\gamma + t)^2} - \frac{(l\lambda_{\min}(\mathbf{M}) - 2)v}{2(\gamma + t)^2} \leq 0$ .

875 This is equivalent to

$$\begin{aligned}
876 6\sigma^2 l^2 \text{tr}(\mathbf{M}) + \frac{\lambda_{\max}^3(\mathbf{M})(6+d)^3 \alpha^2 l^2}{2} + \frac{\lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 l(\gamma + t)}{4} &\leq (l\lambda_{\min}(\mathbf{M}) - 2)v. \\
877 &\Rightarrow v \geq \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + \frac{9\lambda_{\max}^3(\mathbf{M})(6+d)^3 \alpha^2}{2\lambda_{\min}^2(\mathbf{M})} + \frac{3\lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2 (\gamma + t)}{4\lambda_{\min}(\mathbf{M})}. \\
878 &\Rightarrow v \geq \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + \frac{3[6\lambda_{\max}(\mathbf{M}) + 36\text{tr}(\mathbf{M}) + \lambda_{\min}(\mathbf{M})T] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha^2}{4\lambda_{\min}^2(\mathbf{M})}. \\
879 &\Rightarrow v \geq \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + \frac{3[6\lambda_{\max}(\mathbf{M}) + 36\text{tr}(\mathbf{M}) + \lambda_{\min}(\mathbf{M})T] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha_0^2}{4(T+1)\lambda_{\min}^2(\mathbf{M})}. \\
880 &\Rightarrow v \geq \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + \frac{3[6\lambda_{\max}(\mathbf{M}) + 36\text{tr}(\mathbf{M})] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \alpha_0^2}{4\lambda_{\min}^2(\mathbf{M})} \\
881 &= \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + Q_1(\alpha_0^2).
\end{aligned}$$

882 So, we can finally obtain  $v \geq \frac{54\text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}^2(\mathbf{M})} + Q_1(\alpha_0^2)$ .

883 Due to the facts

$$884 (\gamma + t)^2 \geq (\gamma + t + 1)(\gamma + t - 1) = (\gamma + t)^2 - 1,$$

885 then

$$886 \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] \leq \frac{v}{\gamma + t + 1}.$$

887  $\square$

918 D.3 PROOF OF THEOREM 4.7  
919

920 If the objective function is not quadratic function, we notice that  $\gamma_u \neq 1$  and  $\gamma_l \neq 1$ . So, we can  
921 transform inequality (23) into

$$922 \quad f(\mathbf{x}^{t+1}) \leq f(\mathbf{x}^t) - \eta \langle \nabla f(\mathbf{x}^t), \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \rangle \\ 923 \quad + \frac{\gamma_u \eta^2}{2} \left\| \mathbf{u}_t \mathbf{u}_t^\top \nabla f(\mathbf{x}^t, \mathcal{S}_t) + \phi(\mathbf{u}_t, \alpha, \mathbf{x}^t) \right\|_{\mathbf{M}(\mathbf{z}^t)}^2.$$

926 Let us deduce the expectation of  $f(\mathbf{x}^{t+1})$  for  $\mathbf{u}$ ,

$$927 \quad \mathbb{E}_{\mathbf{u}_t} [f(\mathbf{x}^{t+1})] \stackrel{(20)+(21)}{\leq} f(\mathbf{x}^t) + 3\eta^2 \gamma_u \text{tr}(\mathbf{M}(\mathbf{z}^t)) \|\nabla f(\mathbf{x}^t, \mathcal{S}_t)\|^2 + \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|^2 \\ 928 \quad - \eta \langle \nabla f(\mathbf{x}^t), \nabla f(\mathbf{x}^t, \mathcal{S}_t) \rangle \\ 929 \quad + \frac{[1 + 2\gamma_u \lambda_{\max}(\mathbf{M}(\mathbf{z}^t))\eta] \lambda_{\max}^2(\mathbf{M}(\mathbf{z}^t)) (6+d)^3 \gamma_u^2 \alpha^2 \eta}{8}.$$

933 And we can transform inequality (24) into

$$934 \quad \mathbb{E} [f(\mathbf{x}^{t+1})] \leq f(\mathbf{x}^t) + 3\eta^2 \sigma^2 \gamma_u \text{tr}(\mathbf{M}(\mathbf{z}^t)) \\ 935 \quad + \frac{[1 + 2\gamma_u \lambda_{\max}(\mathbf{M}(\mathbf{z}^t))\eta] \lambda_{\max}^2(\mathbf{M}(\mathbf{z}^t)) (6+d)^3 \gamma_u^2 \alpha^2 \eta}{8} \\ 936 \quad - \frac{1}{2} \gamma_l \eta \lambda_{\min}(\mathbf{M}(\mathbf{z}^t)) (f(\mathbf{x}^t) - f^*).$$

940 If we let  $\text{tr}(\mathbf{M}) = \max_{\mathbf{z}^t} \text{tr}(\mathbf{M}(\mathbf{z}^t))$ ,  $\lambda_{\min}(\mathbf{M}) = \min_{\mathbf{z}^t} \lambda_{\min}(\mathbf{M}(\mathbf{z}^t))$  and  $\lambda_{\max}(\mathbf{M}) =$   
941  $\max_{\mathbf{z}^t} \lambda_{\max}(\mathbf{M}(\mathbf{z}^t))$  in the subsequent analysis. Then, we can obtain

$$942 \quad \mathbb{E} \left[ f(\mathbf{x}^{t+1}) - f^* - \frac{24\eta \gamma_u \text{tr}(\mathbf{M}) \sigma^2 + [1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M}) (6+d)^3 \gamma_u^2 \alpha^2}{4\gamma_l \lambda_{\min}(\mathbf{M})} \right] \\ 943 \leq \left[ 1 - \frac{1}{2} \eta \gamma_l \lambda_{\min}(\mathbf{M}) \right] \mathbb{E} \left[ f(\mathbf{x}^t) - f^* - \frac{24\eta \gamma_u \text{tr}(\mathbf{M}) \sigma^2 + [1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M}) (6+d)^3 \gamma_u^2 \alpha^2}{4\gamma_l \lambda_{\min}(\mathbf{M})} \right] \\ 944 \leq \left[ 1 - \frac{1}{2} \eta \gamma_l \lambda_{\min}(\mathbf{M}) \right]^t \left[ f(\mathbf{x}^0) - f^* - \frac{24\eta \gamma_u \text{tr}(\mathbf{M}) \sigma^2 + [1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M}) (6+d)^3 \gamma_u^2 \alpha^2}{4\gamma_l \lambda_{\min}(\mathbf{M})} \right] \\ 945 \leq \left[ 1 - \frac{1}{2} \eta \gamma_l \lambda_{\min}(\mathbf{M}) \right]^t [f(\mathbf{x}^0) - f^*].$$

952 Thus, we can obtain that

$$953 \quad \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] \leq \frac{24\eta \gamma_u \text{tr}(\mathbf{M}) \sigma^2 + [1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta] \lambda_{\max}^2(\mathbf{M}) (6+d)^3 \gamma_u^2 \alpha^2}{4\gamma_l \lambda_{\min}(\mathbf{M})} \\ 954 \quad + \left[ 1 - \frac{1}{2} \eta \gamma_l \lambda_{\min}(\mathbf{M}) \right]^t [f(\mathbf{x}^0) - f^*].$$

959 Let  $\sigma = 0$  and a sufficiently small  $\alpha$  is chosen, similar to the proof process of E.1, we can obtain the  
960 iteration complexity

$$961 \quad t = \mathcal{O} \left( \frac{\gamma_u \text{tr}(\mathbf{M})}{\gamma_l \lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon} \right). \quad (25)$$

964 D.4 PROOF OF THEOREM 4.8  
965

966 Firstly, if we choose decreasing step size  $\eta_t$ , we can obtain the following formula

$$967 \quad \mathbb{E} [f(\mathbf{x}^{t+1})] \leq f(\mathbf{x}^t) + 3\eta_t^2 \sigma^2 \gamma_u \text{tr}(\mathbf{M}(\mathbf{z}^t)) \\ 968 \quad + \frac{[1 + 2\gamma_u \lambda_{\max}(\mathbf{M}(\mathbf{z}^t))\eta_t] \lambda_{\max}^3(\mathbf{M}(\mathbf{z}^t)) (6+d)^3 \gamma_u^2 \alpha^2 \eta_t}{8} \\ 969 \quad - \frac{1}{2} \gamma_l \eta_t \lambda_{\min}(\mathbf{M}(\mathbf{z}^t)) (f(\mathbf{x}^t) - f^*).$$

We let  $\text{tr}(\mathbf{M}) = \max_{\mathbf{z}^t} \text{tr}(\mathbf{M}(\mathbf{z}^t))$ ,  $\lambda_{\min}(\mathbf{M}) = \min_{\mathbf{z}^t} \lambda_{\min}(\mathbf{M}(\mathbf{z}^t))$  and  $\lambda_{\max}(\mathbf{M}) = \max_{\mathbf{z}^t} \lambda_{\max}(\mathbf{M}(\mathbf{z}^t))$  in the subsequent analysis. Then, we need to add  $\gamma_u$  and  $\gamma_l$  to the appropriate position in the proof process of D.2 like the similar ways we operated in D.3. Suppose that holds for  $t > 0$ , then

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] &\leq 3\eta_t^2 \sigma^2 \gamma_u \text{tr}(\mathbf{M}) + \frac{[1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \gamma_u^2 \alpha^2 \eta_t}{8} \\ &\quad + \left[1 - \frac{1}{2} \gamma_l \eta_t \lambda_{\min}(\mathbf{M})\right] \mathbb{E} [f(\mathbf{x}^t) - f^*] \\ &\leq 3\eta_t^2 \sigma^2 \gamma_u \text{tr}(\mathbf{M}) + \frac{[1 + 2\gamma_u \lambda_{\max}(\mathbf{M})\eta_t] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \gamma_u^2 \alpha^2 \eta_t}{8} \\ &\quad + \left[1 - \frac{1}{2} \gamma_l \eta_t \lambda_{\min}(\mathbf{M})\right] \frac{v}{\gamma + t} \\ &= \frac{(\gamma + t - 1)v}{(\gamma + t)^2} + \frac{3\sigma^2 l^2 \gamma_u \text{tr}(\mathbf{M})}{(\gamma + t)^2} + \frac{\lambda_{\max}^2(\mathbf{M})(6+d)^3 \gamma_u^2 \alpha^2 l}{8(\gamma + t)} \\ &\quad + \frac{\lambda_{\max}^3(\mathbf{M})(6+d)^3 \gamma_u^3 \alpha^2 l^2}{4(\gamma + t)^2} - \frac{(\gamma_l \lambda_{\min}(\mathbf{M}) - 2)v}{2(\gamma + t)^2}. \end{aligned}$$

$$\text{We define } Q_2(\alpha_0^2) = \frac{3[6\gamma_u \lambda_{\max}(\mathbf{M}) + 36\gamma_u \gamma_l \text{tr}(\mathbf{M})] \lambda_{\max}^2(\mathbf{M})(6+d)^3 \gamma_u^2 \alpha_0^2}{4\gamma_l^2 \lambda_{\min}^2(\mathbf{M})}.$$

$$\text{Then, we can obtain } v \geq \frac{54\gamma_u \text{tr}(\mathbf{M})\sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} + Q_2(\alpha_0^2).$$

Finally, we obtain the iteration complexity

$$t = \mathcal{O} \left( \left[ \frac{\gamma_u \text{tr}(\mathbf{M})\sigma^2}{\gamma_l^2 \lambda_{\min}^2(\mathbf{M})} + Q_2(\alpha_0^2) \right] \frac{1}{\varepsilon} \right). \quad (26)$$

## E PROOF OF MAIN COROLLARIES

In this section, we give some details of proof about Corollary.

### E.1 PROOF OF COROLLARY 4.4

*Proof.* From the proof of Theorem 4.3, if we choose a sufficiently small  $\alpha$  in practice, we can find that

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{t+1}) - f^*] &\leq \frac{6\eta \text{tr}(\mathbf{M})\sigma^2}{\lambda_{\min}(\mathbf{M})} + \left[1 - \frac{1}{2} \eta \lambda_{\min}(\mathbf{M})\right]^t [f(\mathbf{x}^0) - f^*] \\ &\stackrel{\sigma=0}{=} \left[1 - \frac{1}{2} \eta \lambda_{\min}(\mathbf{M})\right]^t [f(\mathbf{x}^0) - f^*] \\ &\leq \exp \left( -\frac{1}{2} \eta \lambda_{\min}(\mathbf{M}) t \right) [f(\mathbf{x}^0) - f^*] \\ &\leq \exp \left( -\frac{\lambda_{\min}(\mathbf{M})}{24 \text{tr}(\mathbf{M})} t \right) [f(\mathbf{x}^0) - f^*]. \end{aligned}$$

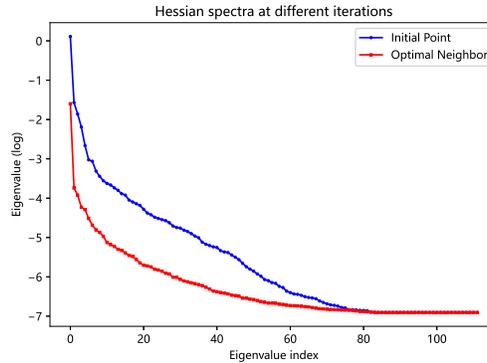
Thus, in order to achieve  $\varepsilon$ -suboptimal solution,  $t$  is required to be

$$\begin{aligned} t &= \frac{24 \text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \left( \log \frac{1}{\varepsilon} + \log (f(\mathbf{x}^0) - f^*) \right) \\ &= \mathcal{O} \left( \frac{\text{tr}(\mathbf{M})}{\lambda_{\min}(\mathbf{M})} \log \frac{1}{\varepsilon} \right). \end{aligned}$$

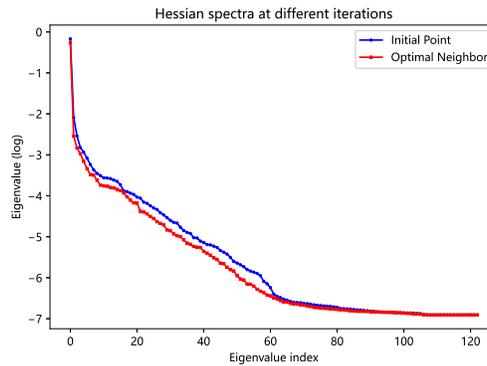
□

## F EIGENVALUE SKEWNESS ANALYSIS ON REAL-WORLD DATASETS

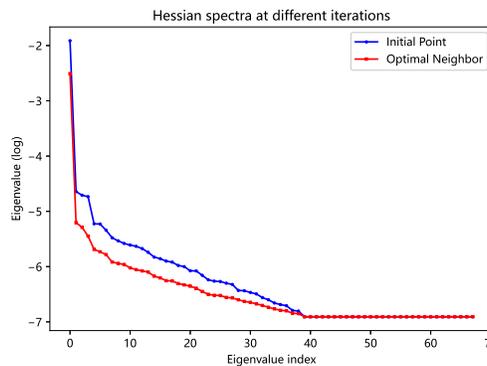
We analyze the spectral properties of the Hessian matrices at both the initial points and near the optimal solutions across all three datasets. We observe that the eigenvalue spectra of the Hessians exhibit a rapid decay in all cases. Consistent with the findings of Yue et al. (2023), this indicates that many real-world problems naturally possess rapidly decaying Hessian spectra. This phenomenon fundamentally explains why ZSG-type algorithms—with their inherently weak dimensional dependence—tend to perform well and are widely adopted in practice.



(a) Eigenvalue distribution of the Hessian of ‘mushrooms’



(b) Eigenvalue distribution of the Hessian of ‘a8a’

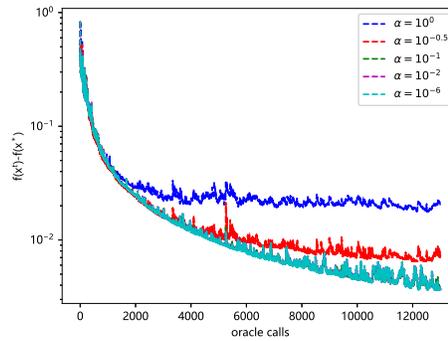


(c) Eigenvalue distribution of the Hessian of ‘phishing’

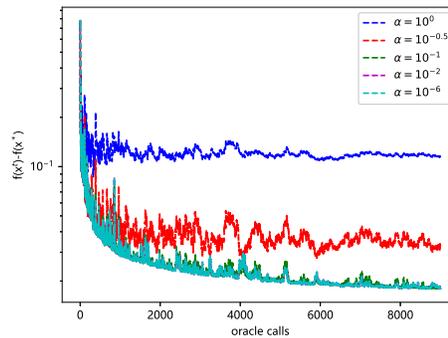
Figure 3: Eigenvalue distributions of the Hessian matrices at the initial points and near the optimal solutions for three logistic regression datasets.

## G SENSITIVITY ANALYSIS OF $\alpha$ IN LOGISTIC REGRESSION

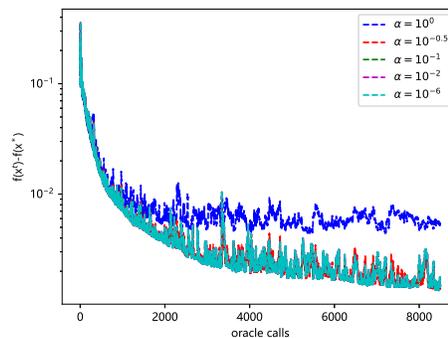
In the quadratic objective setting, the  $\alpha$ -related terms do not introduce additional bias into the zeroth-order gradient estimator. Therefore, in this section we analyze the sensitivity of  $\alpha$  in the Logistic Regression tasks. For each of the three datasets, we conduct ZSG experiments with  $\alpha = 10^0$ ,  $10^{-0.5}$ ,  $10^{-1}$ ,  $10^{-2}$ , and  $10^{-6}$ . The results show that when  $\alpha = 10^0$ , the large noise markedly impedes convergence across all datasets. When  $\alpha = 10^{-0.5}$ , the noise has a mild negative effect on convergence for 'mushrooms' and 'a8a'. In addition, we find that  $\alpha = 10^{-2}$  is already sufficiently small for most practical problems, indicating that the noise induced by  $\alpha$  is easily controlled.



(a) Sensitivity analysis on 'mushrooms'



(b) Sensitivity analysis on 'a8a'



(c) Sensitivity analysis on 'phishing'

Figure 4: The effect of different values of  $\alpha$  on model training across three logistic-regression datasets.

## H ALGORITHM DESCRIPTION OF ZSC

---

**Algorithm 2** ZSC: ZO-SGD-Coordinate Method

---

**Input and Initialize:** parameters  $\mathbf{x} \in \mathbb{R}^d$ , loss function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , step budget  $t$ , step size  $\eta_t > 0$ , perturbation scale  $\alpha$ , sample distribution  $\mathcal{D}$ , initial point  $\mathbf{x}^0 \in \mathbb{R}^d$

**for**  $t = 0, 1, \dots$  **do**

Sample  $\mathcal{S}_t \sim \mathcal{D}$  and  $\mathbf{e}_t \sim \mathcal{U}^d$

Query the zeroth-order oracle  $f_+^t = f(\mathbf{x}^t + \alpha \mathbf{e}_t, \mathcal{S}_t)$

Query the zeroth-order oracle  $f_-^t = f(\mathbf{x}^t - \alpha \mathbf{e}_t, \mathcal{S}_t)$

Estimating the gradient  $\tilde{\nabla} f(\mathbf{x}^t, \mathcal{S}_t) = \frac{(f_+^t - f_-^t)}{2\alpha} \cdot \mathbf{e}_t$

$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \tilde{\nabla} f(\mathbf{x}^t, \mathcal{S}_t)$

**end for**

---