

# SHARPNESS-AWARE BLACK-BOX OPTIMIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Black-box optimization algorithms have been widely used in various machine learning problems, including reinforcement learning and prompt fine-tuning. However, directly optimizing the training loss value, as commonly done in existing black-box optimization methods, could lead to suboptimal model quality and generalization performance. To address those problems in black-box optimization, we propose a novel Sharpness-Aware Black-box Optimization (SABO) algorithm, which applies a sharpness-aware minimization strategy to improve the model generalization. Specifically, the proposed SABO method first reparameterizes the objective function by its expectation over a Gaussian distribution. Then it iteratively updates the parameterized distribution by approximated stochastic gradients of the maximum objective value within a small neighborhood around the current solution in the Gaussian distribution space. Theoretically, we prove the convergence rate and generalization bound of the proposed SABO algorithm. Empirically, extensive experiments on the black-box prompt fine-tuning tasks demonstrate the effectiveness of the proposed SABO method in improving model generalization performance.

## 1 INTRODUCTION

Black-box optimization involves optimizing one objective function by using function queries only. In this work, we study the black-box optimization problem (Jones et al., 1998), which is formulated as

$$\min_{\mathbf{x}} F(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathcal{X}, \quad (1)$$

where  $\mathcal{X} \subseteq \mathbb{R}^d$ , and  $d$  represents the parameter dimension. The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , which satisfies  $F(\mathbf{x}) \geq -\infty$  (lower bounded), can only be queried to obtain function values and we cannot get the gradient of  $F$  w.r.t.  $\mathbf{x}$ . In this work, we focus on the online setting for black-box optimization, where different from the offline setting (Chen et al., 2022; Qi et al., 2022), we do not have a prior dataset containing the variable  $\mathbf{x}$  and its corresponding objective value.

Black-box optimization has drawn intensive attention in a wide range of applications, such as deep reinforcement learning (Salimans et al., 2017; Conti et al., 2018), black-box adversarial attacks of deep neural networks (Ilyas et al., 2018; Dong et al., 2019), etc. Recently, black-box optimization has shown increasing power on real-world natural language processing tasks, especially with the emergence of large language models (LLMs). Since a common practice is to release LLMs as a service and allow users to access it through their inference APIs. In such a scenario, called Languaged-Model-as-a-Service (LMaaS) (Sun et al., 2022b;a; 2023), users cannot access or tune model parameters but can only tune their prompts without model backpropagation to accomplish language tasks of interest, which directly increase the demand of black-box optimization methods.

Although black-box optimization algorithms have been successfully applied to various learning tasks, most existing works directly optimize the training loss value (Sun et al., 2022b), which may lead to suboptimal model quality and generalization performance. Since the training loss landscape is complex and has many local minima with different generalization abilities (Zhang et al., 2021), the learned model may suffer from the overfitting problem, causing poor generalization performance (Foret et al., 2021). Hence, reducing the performance gap between training and testing is an important research topic in deep learning (Neyshabur et al., 2017). Recently, there have been many works exploring the close relationship between loss geometry and generalization performance, and it has been observed that flat minima often imply better generalization (Dziugaite & Roy, 2017; Chatterji et al., 2019; Petzka et al., 2021). This inspires us to design a black-box optimization algorithm to improve the model generalization by finding the flat minima.

054 Sharpness-aware minimization (SAM) (Foret et al., 2021) is a state-of-the-art method to seek flat  
 055 minima in white-box cases by solving a min-max optimization problem. SAM minimizes the  
 056 maximum objective value within a small neighborhood around current solution. Since SAM considers  
 057 the geometry of the Euclidean parameter space, it uses the Euclidean ball to define the neighborhood.  
 058 In SAM, each update consists of two forward-backward computations: one for computing the  
 059 perturbation and the other for computing the actual update direction. SAM has been proven to  
 060 perform better than SGD and its variants (Kwon et al., 2021; Zhuang et al., 2022; Zhao et al., 2022;  
 061 Kim et al., 2022; Jiang et al., 2023; Liu et al., 2022) yield significant performance gains in various  
 062 fields such as computer vision and natural language processing (Bahri et al., 2022; Foret et al., 2021).  
 063 However, SAM and its variants rely on the availability of true gradients or stochastic gradients  
 064 w.r.t. the variable  $x$ , so they are inapplicable to black-box optimization.

065 To take advantage of SAM to improve the generalization performance of black-box optimization,  
 066 we propose a **Sharpness-Aware Black-box Optimization (SABO)** algorithm. Specifically, SABO  
 067 first reparameterizes the objective function via its expectation over a Gaussian distribution, which  
 068 can help to optimize the objective by only accessing the function value (Wierstra et al., 2014; Lyu  
 069 & Tsang, 2021). Then the SABO method seeks to identify the robust minimum region over the  
 070 space of Gaussian distributions, which is different from SAM that finds the flat minimum over  
 071 the parameter space. To achieve that, the SABO method iteratively updates the parameterized  
 072 distribution via a search direction obtained by approximated stochastic gradients for the maximum  
 073 objective value within a small neighborhood around the current solution in the space of Gaussian  
 074 distributions. Theoretically, we analyze the convergence rate and provide a generalization error bound  
 075 for the proposed SABO algorithm. Empirically, we verify the convergence result of the proposed  
 076 algorithm on the synthetic problems, and extensive experimental results on a black-box prompt  
 077 fine-tuning problem demonstrate the effectiveness of the proposed SABO method. Our contributions  
 078 are summarized as follows.

- 079 • We propose the SABO algorithm for black-box optimization. To the best of our knowledge,  
 080 we are the first to design a stochastic gradient approximation algorithm to improve the model  
 081 generalization in black-box optimization by using the sharpness-aware minimization strategy.
- 082 • Theoretically, we prove that the proposed SABO algorithm possesses a convergence rate  $\mathcal{O}(\frac{\log T}{T})$   
 083 in a full-batch function query setting and  $\mathcal{O}(\frac{1}{\sqrt{T}})$  in a mini-batch function query setting, respec-  
 084 tively. Moreover, we provide a generalization error analysis for the proposed SABO method.
- 085 • Empirically, we verify the convergence result of the SABO algorithm on the synthetic numerical  
 086 problems. Moreover, extensive experiments on black-box prompt fine-tuning tasks demon-  
 087 strate the effectiveness of the proposed SABO method in improving the model generalization  
 088 performance.

090 **Notation and Symbols.** We denote by  $\|\cdot\|_2$  and  $\|\cdot\|_\infty$  the  $l_2$  norm and the  $l_\infty$  norm for vectors,  
 091 respectively.  $\|\cdot\|_F$  denotes the Frobenius norm for matrices.  $\mathcal{S}^+$  denotes the set of positive semi-  
 092 definite matrices. For a square matrix  $\mathbf{X}$ ,  $\text{diag}(\mathbf{X})$  represents a vector with diagonal entries in  $\mathbf{X}$ ,  
 093 and if  $\mathbf{x}$  is a vector,  $\text{diag}(\mathbf{x})$  represents a diagonal matrix with  $\mathbf{x}$  as its diagonal entries. We define  
 094  $\|X\|_Y := \sqrt{\langle X, YX \rangle}$  for a positive semi-definite matrix  $Y \in \mathcal{S}^+$  or a non-negative vector  $Y$ , where  
 095  $\langle \cdot, \cdot \rangle$  denotes the inner product under the Frobenius norm for matrices and inner product under the  $l_2$   
 096 norm for vectors.  $\frac{X}{Y}$  denotes the elementwise division operation when  $X$  and  $Y$  are vectors (for the  
 097 diagonal matrix), and the elementwise division operation for diagonal elements in  $X$  and  $Y$  when  
 098 they are diagonal matrices.

## 100 2 BACKGROUND

102 **Stochastic Gradient Approximation** The stochastic gradient approximation method (Wierstra  
 103 et al., 2014; Lyu & Tsang, 2021; Ye et al., 2024) is a representative strategy for solving black-box  
 104 optimization problems, which instead of maintaining a population of searching points, iteratively  
 105 updates a search distribution by stochastic gradient approximation. The general procedure of stochas-  
 106 tic gradient approximation methods is to first generate a batch of sample points by a parameterized  
 107 search distribution. Then the sample points allow the algorithm to capture the local structure of the  
 fitness function and appropriately estimate the stochastic gradient to update the distribution.

Specifically, the stochastic gradient approximation method reparameterizes  $F(\mathbf{x})$  as

$$J(\boldsymbol{\theta}) = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})}[F(\mathbf{x})] = \int F(\mathbf{x})p(\mathbf{x}; \boldsymbol{\theta})d\mathbf{x}, \quad (2)$$

where  $\boldsymbol{\theta}$  denotes the parameters of density  $p(\mathbf{x}; \boldsymbol{\theta})$  or  $p_{\boldsymbol{\theta}}$  and  $F(\mathbf{x})$  is also referred to as the fitness function for  $\mathbf{x}$ . Based on this definition, we can obtain the Monte Carlo estimation of the search gradient as

$$\bar{\nabla}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{j=1}^N F(\mathbf{x}_j) \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}_j), \quad (3)$$

where  $\mathbf{x}_j$  denotes the  $j$ -th sample and  $N$  denotes the number of samples. Therefore, the stochastic gradient  $\bar{\nabla}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  provides a search direction in the space of search distributions.

**Sharpness-Aware Minimization** SAM (Foret et al., 2021) attempts to improve generalization by finding flat minima. This is achieved by minimizing the worst-case loss within some perturbation radius. Mathematically, it is formulated as the following minimax optimization problem:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho^2} F(\mathbf{x} + \boldsymbol{\epsilon}), \quad (4)$$

where  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  is the objective function,  $\mathbf{x}$  denotes variables that can represent model parameters,  $\rho > 0$  is a positive constant, and  $\boldsymbol{\epsilon}$  is the perturbation whose magnitude is bounded by  $\rho^2$ . By taking the first-order approximation of  $F(\mathbf{x} + \boldsymbol{\epsilon})$  over  $F(\mathbf{x})$ , a solution of  $\boldsymbol{\epsilon}$  for the maximization subproblem can be obtained as

$$\boldsymbol{\epsilon}(\mathbf{x}) = \frac{\rho^2 \nabla F(\mathbf{x})}{\|\nabla F(\mathbf{x})\|_2}. \quad (5)$$

Then problem (4) can be solved by performing the gradient descent method for the minimization subproblem as  $\mathbf{x}_{t+1} = \mathbf{x}_t - \beta_t \nabla_{\mathbf{x}} F(\mathbf{x}_t + \boldsymbol{\epsilon}(\mathbf{x}_t))$ , where  $\beta_t$  represents the step size in the  $t$ -th iteration. Note that this gradient implicitly depends on the Hessian of  $F(\mathbf{x})$  because  $\boldsymbol{\epsilon}(\mathbf{x})$  is a function of  $\mathbf{x}$ . To accelerate the computation, a common approach in SAM-based methods (Foret et al., 2021; Kim et al., 2022; Jiang et al., 2023) is to apply a first-order gradient approximation, so we obtain the update rule for SAM as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \beta_t \nabla_{\mathbf{x}} F(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_t+\boldsymbol{\epsilon}_t}, \quad (6)$$

where  $\boldsymbol{\epsilon}_t = \boldsymbol{\epsilon}(\mathbf{x}_t)$  is viewed as a constant w.r.t.  $\mathbf{x}$ .

### 3 METHODOLOGY

In this section, we introduce the proposed SABO algorithm. Firstly, we formulate the sharpness-aware black-box optimization as a min-max optimization problem and solve it in Section 3.1, and in Section 3.2, we derive the update formula of parameters in the search distribution. The detailed derivations in this section are put in Appendix A.

#### 3.1 SHARPNESS-AWARE BLACK-BOX OPTIMIZATION

Suppose we are given a training set  $\mathcal{D}$  with i.i.d. samples  $\{(X_i, y_i)\}$ . The main objective is defined as

$$F(\mathbf{x}; \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(X_i, y_i) \in \mathcal{D}} l(\mathbf{x}; (X_i, y_i)), \quad (7)$$

where  $\mathbf{x}$  denotes model parameters,  $|\mathcal{D}|$  denotes the number of data in the dataset  $\mathcal{D}$ , and  $l(\mathbf{x}; (X, y))$  denotes the loss function (e.g., the cross-entropy loss for classification). To simplify the notation, we define  $F(\mathbf{x}) := F(\mathbf{x}; \mathcal{D})$ .

In black-box optimization, we aim at minimizing the objective function  $F(\mathbf{x})$ , with only function queries. Due to the lack of gradient information, we first apply the stochastic gradient approximation method (Wierstra et al., 2014; Lyu & Tsang, 2021). We denote by  $\boldsymbol{\theta}$  the parameters of the search distribution  $p_{\boldsymbol{\theta}}$  and define the expected fitness of  $F(\mathbf{x})$  under the parametric search distribution

162  $p_{\theta}(\mathbf{x})$  as  $J(\theta) = \mathbb{E}_{p_{\theta}(\mathbf{x})}[F(\mathbf{x})]$ . Then the optimal parameter  $\theta$  can be found by minimizing the  
 163 reparameterized objective  $J(\theta)$ .  
 164

165 Inspired by SAM (Foret et al., 2021), we attempt to improve generalization by finding flat minima of  $\theta$ .  
 166 However, for the reparameterized objective  $J(\theta)$ , the geometry of the corresponding distribution space  
 167 is not Euclidean but a statistical manifold, where the distance between two probability distributions is  
 168 defined by some statistical distance, e.g., Kullback-Leibler (KL) divergence. Therefore, instead of  
 169 restricting the perturbation in an Euclidean ball, we restrict the perturbation distribution to be inside  
 170 a small neighborhood of the unperturbed distribution w.r.t. the KL divergence (Amari, 2016). The  
 171 proposed optimization problem for black-box optimization is formulated as

$$172 \min_{\theta} \max_{\delta \in \mathcal{C}(\theta)} J(\theta + \delta), \quad (8)$$

173 where  $\mathcal{C}(\theta) = \{\delta \mid \text{KL}(p_{\theta+\delta} \| p_{\theta}) \leq \rho^2\}$ ,  $\rho$  is a positive constant, and  $J(\theta + \delta) = \mathbb{E}_{x \sim p_{\theta+\delta}}[F(x)]$ .  
 174 Note that  $\mathcal{C}(\theta)$  defines the neighborhood around a given distribution  $p_{\theta}$  in the distribution space,  
 175 which is different from the neighborhood of SAM that is defined in the parameter space. Here  $\rho$   
 176 defines the size of the neighborhood.

177 Generally, problem (8) is applicable to any family of distribution  $p_{\theta}$ . For computational consideration,  
 178 the search distribution is assumed to be a Gaussian distribution, i.e.,  $p_{\theta}(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mu, \Sigma)$  where  $\mu$   
 179 denotes the mean and  $\Sigma$  denotes the covariance matrix, and correspondingly  $\theta$  includes  $\mu$  and  $\Sigma$ ,  
 180 i.e.,  $\theta = \{\mu, \Sigma\}$ . In this work, we assume that the covariance matrix  $\Sigma$  is a diagonal matrix. For a  
 181 perturbation  $\delta = \{\delta_{\mu}, \delta_{\Sigma}\}$  where  $\delta_{\Sigma}$  is a diagonal matrix, the perturbed distribution is a Gaussian  
 182 distribution  $p_{\theta+\delta}(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mu + \delta_{\mu}, \Sigma + \delta_{\Sigma})$ .

183 We need to solve problem (8) to derive the update formulation for  $\theta$ . By using the first-order Taylor  
 184 expansion, the maximization subproblem can be approximated as a quadratically constrained linear  
 185 programming problem:

$$186 \max_{\delta \in \mathcal{C}(\theta)} J(\theta + \delta) \approx \max_{\delta \in \mathcal{C}(\theta)} \langle \nabla_{\theta} J(\theta), \delta \rangle, \quad (9)$$

187 The corresponding Lagrangian of problem (9) is

$$188 \mathcal{L}(\delta, \lambda) = -\langle \nabla_{\theta} J(\theta), \delta \rangle + \lambda(\text{KL}(p_{\theta+\delta} \| p_{\theta}) - \rho^2), \quad (10)$$

189 where  $\lambda$  is the Lagrange multiplier. We can see that problem (10) is convex with respect to  $\delta$ .  
 190 Therefore, by setting the derivatives w.r.t.  $\delta_{\mu}$  and  $\delta_{\Sigma}$  to zero, we can obtain  $\delta$  as

$$192 \delta_{\mu}(\theta) = \frac{1}{\lambda} \Sigma \nabla_{\mu} J(\theta), \quad \delta_{\Sigma}(\theta) = \frac{2 \Sigma \nabla_{\Sigma} J(\theta)}{\lambda \Sigma^{-1} - 2 \nabla_{\Sigma} J(\theta)}. \quad (11)$$

193 As shown in Eq. (11), to calculate the perturbation, we need to calculate the inverse covariance  
 194 matrix, i.e.,  $\Sigma^{-1}$ , which is computationally expensive for high-dimensional problems. Therefore,  
 195 assuming that  $\Sigma$  is a diagonal matrix can significantly reduce the computation cost. Then by plugging  
 196  $\delta_{\mu}(\theta)$  and  $\delta_{\Sigma}(\theta)$  into the neighborhood constraint, i.e.,  $\delta \in \mathcal{C}(\theta)$ , we can determine the optimal  $\lambda$  as

$$198 \lambda = \frac{1}{\rho} \sqrt{\|\Sigma \nabla_{\Sigma} J(\theta)\|_{\text{F}}^2 + 0.5 \|\Sigma^{\frac{1}{2}} \nabla_{\mu} J(\theta)\|_2^2}. \quad (12)$$

200 Based on Eqs. (11) and (12), we obtain the approximated closed-form solution  $\delta(\theta)$  for a given  $\theta$ .  
 201 With  $\delta(\theta)$ , the minimization subproblem of problem (8) can be reformulated as

$$202 \min_{\theta} J(\theta + \delta(\theta)). \quad (13)$$

203 To solve problem (13), in the  $t$ -th iteration, we add a regularization term  $\frac{1}{\beta_t} \text{KL}(p_{\theta} \| p_{\theta_t})$  to problem  
 204 (13) to enforce  $\theta$  to be close to  $\theta_t$ , and obtain  $\theta_{t+1}$  by solving the following problem as

$$206 \theta_{t+1} = \arg \min_{\theta} J(\theta + \delta(\theta_t)) - J(\theta_t) + \frac{1}{\beta_t} \text{KL}(p_{\theta} \| p_{\theta_t}). \quad (14)$$

207 Following standard SAM-based methods (Foret et al., 2021), we treat  $\delta(\theta_t)$  as a constant  $\delta_t$  instead  
 208 of a function of  $\theta_t$  to accelerate the computation. Then by using the first-order Taylor expansion,  
 209 problem (14) can be approximated as

$$211 \theta_{t+1} = \arg \min_{\theta} \langle \theta - \theta_t, \nabla_{\theta} J(\theta_t + \delta_t) \rangle + \frac{1}{\beta_t} \text{KL}(p_{\theta} \| p_{\theta_t}). \quad (15)$$

212 By solving problem (15), we can obtain the update formulations for  $\mu$  and  $\Sigma$  in the  $t$ -th iteration as

$$214 \mu_{t+1} = \mu_t - \beta_t \Sigma_t \nabla_{\mu} J(\theta_t + \delta_t), \quad \Sigma_{t+1}^{-1} = \Sigma_t^{-1} + 2\beta_t \nabla_{\Sigma} J(\theta_t + \delta_t), \quad (16)$$

215 where  $\nabla_{\mu} J(\theta_t + \delta_t)$  and  $\nabla_{\Sigma} J(\theta_t + \delta_t)$  denote the derivative of  $J(\theta)$  w.r.t.  $\mu$  and  $\Sigma$  at  $\mu = \theta_t + \delta_{\mu_t}$   
 and  $\Sigma = \Sigma_t + \delta_{\Sigma_t}$ , respectively.

### 3.2 UPDATE FORMULATIONS FOR SABO

The gradients of the reparameterized objective  $J(\theta)$  w.r.t.  $\mu$  and  $\Sigma$  rely on the expectations of the black-box function and can be obtained with only function queries (Wierstra et al., 2014) (see Theorem A.1 in Appendix A.4). Hence, we estimate them by Monte Carlo sampling. Specifically, the stochastic approximation of the gradients  $\nabla_{\mu}J(\theta_t)$  and  $\nabla_{\Sigma}J(\theta_t)$  are given as

$$\mathbf{g}'_t = \frac{1}{N} \sum_{j=1}^N \Sigma_t^{-1} (\mathbf{x}'_j - \mu_t) (F(\mathbf{x}'_j) - F(\mu_t)), \quad (17)$$

$$\mathbf{G}'_t = \frac{1}{2N} \sum_{j=1}^N \text{diag} \left[ \Sigma_t^{-1} \left[ \text{diag}((\mathbf{x}'_j - \mu_t)(\mathbf{x}'_j - \mu_t)^\top \Sigma_t^{-1} - \mathbf{I}) (F(\mathbf{x}'_j) - F(\mu_t)) \right] \right], \quad (18)$$

where  $\mathbf{x}'_j$  denotes the  $j$ -th sample sampled from the distribution  $\mathcal{N}(\mathbf{x} \mid \mu_t, \Sigma_t)$ . Note that  $\mathbf{g}'_t$  is an unbiased estimator for the gradient  $\nabla_{\mu}J(\theta)$  as proved in Lemma C.5, and inspired by Lyu & Tsang (2021), we subtract  $F(\mu_t)$  to improve the computational stability while keeping them as unbiased estimations.

Then according to Eq. (11), we obtain the perturbation  $\delta_t$  in the  $t$ -th iteration as

$$\delta_t = \left\{ \frac{1}{\lambda} \Sigma_t \mathbf{g}'_t, \frac{2 \Sigma_t \mathbf{G}'_t}{\lambda \Sigma_t^{-1} - 2 \mathbf{G}'_t} \right\}, \quad (19)$$

where  $\lambda$  is approximated by  $\frac{1}{\rho} \sqrt{\|\Sigma_t \mathbf{G}'_t\|_F^2 + 0.5 \|\Sigma_t^{\frac{1}{2}} \mathbf{g}'_t\|_2^2}$ . Similarly, the gradients  $\nabla_{\mu}J(\theta_t + \delta_t)$  and  $\nabla_{\Sigma}J(\theta_t + \delta_t)$  can be approximated as follows:

$$\mathbf{g}_t = \frac{1}{N} \sum_{j=1}^N \widehat{\Sigma}_t^{-1} (\mathbf{x}_j - \widehat{\mu}_t) (F(\mathbf{x}_j) - F(\widehat{\mu}_t)), \quad (20)$$

$$\mathbf{G}_t = \frac{1}{2N} \sum_{j=1}^N \text{diag} \left[ \widehat{\Sigma}_t^{-1} \left[ \text{diag}((\mathbf{x}_j - \widehat{\mu}_t)(\mathbf{x}_j - \widehat{\mu}_t)^\top \widehat{\Sigma}_t^{-1} - \mathbf{I}) (F(\mathbf{x}_j) - F(\widehat{\mu}_t)) \right] \right], \quad (21)$$

where  $\widehat{\Sigma}_t = \Sigma_t + \delta_{\Sigma_t}$ ,  $\widehat{\mu}_t = \mu_t + \delta_{\mu_t}$ , and  $\mathbf{x}_j$  denotes the  $j$ -th sample sampled from the distribution  $\mathcal{N}(\mathbf{x} \mid \widehat{\mu}_t, \widehat{\Sigma}_t)$ . Then the updated formulations for  $\mu$  and  $\Sigma$  are rewritten as

$$\mu_{t+1} = \mu_t - \beta_t \Sigma_t \mathbf{g}_t, \quad \Sigma_{t+1}^{-1} = \Sigma_t^{-1} + 2\beta_t \mathbf{G}_t. \quad (22)$$

The entire algorithm is shown in Algorithm 1. For the proposed SABO method, the computation cost per iteration is of order  $\mathcal{O}(Nd)$ .

**Mini-batch SABO.** In Algorithm 1, we assume full access to the objective function  $F(\mathbf{x}; \mathcal{D})$ , while in practice, a full-batch function query might be costly. Therefore, we can perform a mini-batch function query. Specifically, in each iteration, we query the expected fitness by a mini-batch of data and approximate  $F(\mathbf{x}; \mathcal{D})$  by

$$F(\mathbf{x}; \mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{(X,y) \in \mathcal{B}} l(\mathbf{x}; (X, y)), \quad (23)$$

where  $|\mathcal{B}|$  denotes the number of data in the mini-batch  $\mathcal{B}$ . The corresponding expected fitness of  $F(\mathbf{x}; \mathcal{B})$  under the distribution  $p_{\theta}(\mathbf{x})$  is formulated as  $J(\theta; \mathcal{B}) = \mathbb{E}_{p_{\theta}(\mathbf{x})} [F(\mathbf{x}; \mathcal{B})]$ . Then similar to the full-batch function query setting, we can approximate the gradients  $\nabla_{\mu}J(\theta_t, \mathcal{B})$  and  $\nabla_{\Sigma}J(\theta_t, \mathcal{B})$ , where the detailed formulations can be found in Appendix B. The entire SABO algorithm with mini-batch function queries is shown in Algorithm 2 in Appendix B.

## 4 ANALYSIS

In this section, we provide comprehensive theoretical analyses for the proposed SABO method with all the detailed proofs in Appendix D.

### 4.1 CONVERGENCE ANALYSIS OF SABO

Firstly, we make an assumption for the reparameterized objective function.

**Algorithm 1** SABO**Require:** Neighborhood size  $\rho$ , learning rate  $\beta_t$ 

- 1: Initialized  $\theta_0 = (\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ ;
- 2: **for**  $t = 0$  to  $T - 1$  **do**
- 3:   Take i.i.d. samples  $\mathbf{z}'_j \sim \mathcal{N}(0, I)$  and set  $\mathbf{x}'_j = \boldsymbol{\mu}_t + \boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{z}'_j$  for  $j \in \{1, \dots, N\}$ ;
- 4:   Query the batch observations  $\{F(\mathbf{x}'_1), \dots, F(\mathbf{x}'_N)\}$ ;
- 5:   Compute the gradient  $\mathbf{g}'_t$  via Eq. (17) and compute the gradient  $\mathbf{G}'_t$  via Eq. (18);
- 6:   Compute  $\lambda = \frac{1}{\rho} \sqrt{\|\boldsymbol{\Sigma}_t \mathbf{G}'_t\|_F^2 + 0.5 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}'_t\|_2^2}$ ;
- 7:   Compute  $\boldsymbol{\delta}_{\boldsymbol{\mu}_t}$  and  $\boldsymbol{\delta}_{\boldsymbol{\Sigma}_t}$  via Eq. (19);
- 8:   Take i.i.d. samples  $\mathbf{z}_j \sim \mathcal{N}(0, I)$  for  $j \in \{1, \dots, N\}$ ;
- 9:   Set  $\mathbf{x}_j = \boldsymbol{\mu}_t + \boldsymbol{\delta}_{\boldsymbol{\mu}_t} + (\boldsymbol{\Sigma}_t + \boldsymbol{\delta}_{\boldsymbol{\Sigma}_t})^{\frac{1}{2}} \mathbf{z}_j$  for  $j \in \{1, \dots, N\}$ ;
- 10:   Query the batch observations  $\{F(\mathbf{x}_1), \dots, F(\mathbf{x}_N)\}$ ;
- 11:   Compute the gradient  $\mathbf{g}_t$  via Eq. (20) and compute the gradient  $\mathbf{G}_t$  via Eq. (21);
- 12:   Set  $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \mathbf{g}_t$  and set  $\boldsymbol{\Sigma}_{t+1}^{-1} = \boldsymbol{\Sigma}_t^{-1} + 2\beta_t \mathbf{G}_t$ ;
- 13: **end for**
- 14: **return**  $\theta_T = (\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_T)$ .

**Assumption 4.1** *The function  $J(\boldsymbol{\theta})$  satisfies that  $\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta})$  is  $L$ -Lipschitz w.r.t.  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\} \in \Theta$ , where  $\Theta := \{\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \in \mathcal{S}^+\}$ .*

Note that the proposed SABO algorithm approximates the gradients of the reparameterized objective function. It is necessary to study the relation between the optimal solutions of the original objective functions  $F(\mathbf{x})$  and  $J(\boldsymbol{\theta})$ , and we put the results in the following proposition.

**Proposition 4.2** (Lyu & Tsang, 2021) *Suppose  $p_{\boldsymbol{\theta}}(\mathbf{x})$  is a Gaussian distribution with  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$  and  $F(\mathbf{x})$  is a convex function. Let  $J(\boldsymbol{\theta}) = \mathbb{E}_{p_{\boldsymbol{\theta}}}[F(\mathbf{x})]$ , and  $J(\boldsymbol{\mu}^*, \mathbf{0}) := F(\boldsymbol{\mu}^*)$ . Then we have*

$$F(\boldsymbol{\mu}) - F(\boldsymbol{\mu}^*) \leq J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) - J(\boldsymbol{\mu}^*, \mathbf{0}), \quad (24)$$

where  $\mathbf{0}$  denotes a zero matrix with appropriate size.

The convexity assumption of the objective function in Theorem 4.3 has been widely adopted in the area of stochastic gradient approximation black-box optimization (Beyer, 2014; Wierstra et al., 2014; Lyu & Tsang, 2021; Ye, 2023). Since the Gaussian-smooth approximation function is always an upper bound of the true target function in convex cases, i.e.,  $F(\boldsymbol{\mu}) \leq \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})}[F(\mathbf{x})]$ . When  $\boldsymbol{\mu}^*$  is an optimal solution of minimization problem  $\min_{\mathbf{x}} F(\mathbf{x})$ , Proposition 4.2 implies that the difference between the objective value at  $\boldsymbol{\mu}$  and the optimal objective value of the original problem is upper-bounded by that of the expected objective function. Then for Algorithm 1, the following theorem captures the convergence of  $\boldsymbol{\mu}$  for a convex objective function.

**Theorem 4.3** *Suppose that  $F(\mathbf{x})$  is a convex function,  $J(\boldsymbol{\theta})$  is  $c$ -strongly convex w.r.t.  $\boldsymbol{\mu}$ , the gradient estimator  $\mathbf{G}_t$  (w.r.t. the covariance matrix) is positive semi-definite matrix such that  $\xi \mathbf{I} \preceq \mathbf{G}_t \preceq \frac{c\mathbf{I}}{4}$  with  $\xi \geq 0$ ,  $\boldsymbol{\Sigma}_0 \in \mathcal{S}^+$ , and  $\boldsymbol{\Sigma}_0 \preceq R\mathbf{I}$  where  $R > 0$ . Suppose the sequence  $\{\boldsymbol{\mu}_t\}$  generated by Algorithm 1 satisfies that the distance between the sequence  $\{\boldsymbol{\mu}_t\}$  and the optimal solution of  $F(\mathbf{x})$  is bounded, i.e.,  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$ ,  $\|\nabla_{\boldsymbol{\Sigma}=\boldsymbol{\Sigma}_t} J(\boldsymbol{\theta})\|_F \leq H$ ,  $\beta_t = \mathcal{O}(1)$ , and  $\rho < \frac{\sqrt{d}}{2}$  satisfies  $\rho = \mathcal{O}(\frac{1}{\sqrt{T}})$ , then with Assumption 4.1, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, \mathbf{0})] = \mathcal{O}\left(\frac{\log T}{T}\right). \quad (25)$$

Based on Theorem 4.3 and Proposition 4.2, when  $\beta_t = \mathcal{O}(1)$  and  $\rho = \mathcal{O}(\frac{1}{\sqrt{T}})$ , we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[F(\boldsymbol{\mu}_{t+1}) - F(\boldsymbol{\mu}^*)] = \mathcal{O}\left(\frac{\log T}{T}\right). \quad (26)$$

Therefore, the proposed SABO algorithm with full-batch function query possesses a convergence rate  $\mathcal{O}(\frac{\log T}{T})$  for the convex objective function. Additionally, in Theorem 4.3, when  $\beta = \mathcal{O}(1)$  and

$\rho = \mathcal{O}(1)$ , the proposed SABO algorithm still maintains a convergence rate  $\mathcal{O}(T^{-\frac{1}{2}})$ . The detailed discussion is provided in Remark D.1.

For the mini-batch setting, we make an additional assumption for the objective function  $F(\mathbf{x}; \mathcal{D})$ .

**Assumption 4.4** *It is assumed that the datasets  $\mathcal{D}$  and  $\mathcal{B}$  are i.i.d. sampled from a data distribution  $P(X, y)$ , and the variance of the mini-batch estimation of the objective function is bounded, i.e.,  $\|F(\mathbf{x}; \mathcal{B}) - \mathbb{E}F(\mathbf{x}; \mathcal{B})\|_2^2 \leq \varepsilon_{\mathcal{B}}^2$  and  $\|F(\mathbf{x}; \mathcal{D}) - \mathbb{E}F(\mathbf{x}; \mathcal{D})\|_2^2 \leq \varepsilon_{\mathcal{D}}^2$ .*

Note that for the standard stochastic gradient descent (SGD) method (Shamir & Zhang, 2013), the unbiased estimation and bounded variance assumptions were made for the approximated gradient. However, in black-box optimization, the gradient of the objective function  $F(\mathbf{x})$  w.r.t.  $\mathbf{x}$  is unavailable. Hence we can only make assumptions for the batch estimations  $F(\mathbf{x}; \mathcal{B})$  and  $F(\mathbf{x}; \mathcal{D})$ .

Then we have the following result for the mini-batch estimation.

**Proposition 4.5** *Suppose Assumption 4.4 holds, then we have  $\mathbb{E}F(\mathbf{x}; \mathcal{B}) = \mathbb{E}F(\mathbf{x}; \mathcal{D})$ , and  $\|F(\mathbf{x}; \mathcal{B}) - F(\mathbf{x}; \mathcal{D})\|_2^2 \leq \varepsilon^2$ , where  $\varepsilon^2 = 2(\varepsilon_{\mathcal{B}}^2 + \varepsilon_{\mathcal{D}}^2)$ .*

Then with Assumption 4.4, the following theorem shows the convergence of  $\boldsymbol{\mu}$  for Algorithm 2.

**Theorem 4.6** *Suppose that  $F(\mathbf{x})$  is a convex function,  $J(\boldsymbol{\theta})$  is  $c$ -strongly convex w.r.t.  $\boldsymbol{\mu}$ , the gradient estimator  $\mathbf{G}_t$  (w.r.t. the covariance matrix) is positive semi-definite matrix such that  $\xi \mathbf{I} \preceq \mathbf{G}_t \preceq \frac{c}{4} \mathbf{I}$  with  $\xi \geq 0$ ,  $\boldsymbol{\Sigma}_0 \in \mathcal{S}^+$ , and  $\boldsymbol{\Sigma}_0 \preceq R\mathbf{I}$  where  $R > 0$ . Suppose the sequence  $\{\boldsymbol{\mu}_t\}$  generated by Algorithm 2 satisfies that the distance between the sequence  $\{\boldsymbol{\mu}_t\}$  and the optimal solution of  $F(\mathbf{x})$  is bounded, i.e.,  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$ ,  $\|\nabla_{\boldsymbol{\Sigma}=\boldsymbol{\Sigma}_t} J(\boldsymbol{\theta})\|_F \leq H$ ,  $\beta_t = \mathcal{O}(1)$ , and  $\rho < \frac{\sqrt{d}}{2}$  satisfies  $\rho = \mathcal{O}(1)$ . Then with Assumptions 4.1 and 4.4, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [J(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, 0)] = \mathcal{O} \left( \frac{1}{\sqrt{T}} \right). \quad (27)$$

Based on Theorem 4.6 and Proposition 4.2, the proposed SABO algorithm with mini-batch function query possesses a convergence rate  $\mathcal{O}(T^{-\frac{1}{2}})$ .

**Remark 4.7** *Note that in Theorem 4.3 and Theorem 4.6, the convergence results do not require the objective function  $F(\mathbf{x})$  to be strongly convex or differentiable. Hence, the convergence holds for a non-smooth convex function  $F(\mathbf{x})$  (as long as  $J(\boldsymbol{\theta})$  being a  $c$ -strongly convex function w.r.t.  $\boldsymbol{\mu}$ ).*

## 4.2 GENERALIZATION ERROR ANALYSIS

In this subsection, we analyze the generalization bound of the proposed SABO algorithm. Specifically, we bound the expectation of the objective function over the Gaussian perturbation.

We denote by  $(X, y)$  a data pair drawn from a data distribution  $P(X, y)$  and by  $F(\mathbf{x}; (X, y))$  the corresponding loss of parameter  $\mathbf{x}$  on  $(X, y)$ . So we have  $F(\mathbf{x}; (X, y)) = l(\mathbf{x}; (X, y))$ . We define the population loss over the data distribution  $P(X, y)$  as  $\mathbb{E}_{P(X, y)}[F(\mathbf{x}; (X, y))]$ , and the empirical loss over a dataset  $\mathcal{S}$ , which consists of  $M$  i.i.d. samples drawn from  $P(X, y)$ , as  $F(\mathbf{x}; \mathcal{S}) = \frac{1}{M} \sum_{i=1}^M l(\mathbf{x}; (X_i, y_i))$ . Then we have following result.

**Theorem 4.8** *Let the loss function  $F(\mathbf{x}; (X, y))$  be a convex function w.r.t.  $\mathbf{x}$ , then for any  $\boldsymbol{\mu} \in \mathbb{R}^d$ , with probability at least  $1 - \kappa$ , we have*

$$\mathbb{E}_{P(X, y)}[F(\boldsymbol{\mu}; (X, y))] \leq \max_{\boldsymbol{\delta} \in \mathcal{C}(\boldsymbol{\theta})} \mathbb{E}_{p_{\boldsymbol{\theta}+\boldsymbol{\delta}}} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\rho^2 + \log(\frac{M}{\kappa})}{2(M-1)}}, \quad (28)$$

where  $p_{\boldsymbol{\theta}} := \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathcal{C}(\boldsymbol{\theta}) = \{\boldsymbol{\delta} \mid \text{KL}(p_{\boldsymbol{\theta}+\boldsymbol{\delta}} \| p_{\boldsymbol{\theta}}) \leq \rho^2\}$ , and  $\mathcal{S}$  denotes the training set that consists of  $M$  i.i.d. samples drawn from data distribution  $P(X, y)$ .

Theorem 4.8 provides a generalization bound for the proposed SABO algorithm. Compared with the generalization bound of SAM presented in Appendix A.1 of Foret et al. (2021), Theorem 4.8 has an

asymptotically identical order in the complexity term. However, the expected generalization loss on the right-hand side of Eq. (28) is different in that we reparameterize the objective function and have perturbation of  $\theta$  in its neighborhood of a statistical manifold, i.e.,  $\delta \in \mathcal{C}(\theta)$ , while SAM bounds the generalization loss averaged over a spherical Gaussian perturbation on parameters.

## 5 RELATED WORKS

**Black-Box Optimization.** Many methods have been proposed for black-box optimization, including Bayesian optimization (BO) methods (Srinivas et al., 2010; Gardner et al., 2017; Nayebi et al., 2019), stochastic optimization methods such as evolution strategies (ES) (Hansen, 2006; Wierstra et al., 2014; Lyu & Tsang, 2021), and genetic algorithms (GA) (Srinivas & Patnaik, 1994; Mirjalili, 2019). Among those methods, BO achieves good performance for low-dimensional problems, but it often fails to handle high-dimensional problems through a global surrogate model, as shown in (Eriksson et al., 2019) and (Nguyen et al., 2022). As a result, TuRBO (Eriksson et al., 2019) and GIBO (Nguyen et al., 2022) try to address this problem with the local BO approach. (Ziomek & Ammar, 2023) further showed that decomposition is important for alleviating the high-dimensional problems in BO. Although BO is not our main focus, we further compare our method with these local Bayesian optimization methods on the black-box prompt fine-tuning problem, and the corresponding exploration is shown in Appendix F.2. GA method is computationally expensive for machine learning problems and usually lacks convergence analysis. The stochastic optimization methods such as CMA-ES (Hansen, 2006) and INGO (Lyu & Tsang, 2021) can scale up to higher-dimensional problems compared with BO. Hence we mainly consider stochastic optimization methods as baseline methods in our experiments.

**Sharpness-Aware Minimization.** SAM has been widely studied for improving the model generalization. Among previous works on SAM (Kwon et al., 2021; Zhuang et al., 2022; Zhao et al., 2022; Jiang et al., 2023), the most relevant method to our approach is the FSAM method (Kim et al., 2022) which also finds the worst-case objective function via a statistical manifold instead of the Euclidean space. However, the loss function of the model studied in FSAM is a predictive distribution conditional on both model parameters and data, while in our case, we consider the parameter as a Gaussian distribution. The bSAM method (Möllenhoff & Khan, 2022) builds a connection between the SAM objective and Bayes objective by Fenchel biconjugate of the loss function. Möllenhoff & Khan (2022) shares a similarity with our work in developing SAM w.r.t. the expected loss. However, The bSAM method relies on the derivation of a convex lower bound of the expected loss by the Fenchel biconjugate and the perturbation is still w.r.t. each point inside the expected loss as standard SAM. Hence FSAM and bSAM are different from the proposed SABO method. Additionally, like other variants of SAM, FSAM and bSAM are both inapplicable to black-box optimization. The STABLEOPT method (Bogunovic et al., 2018) proposes a SAM-like optimization formulation in the Bayesian optimization area. They aim to improve the robustness w.r.t. the adversarial perturbation of the return point by a GP-based optimization. Their method relies on a GP surrogate model that is expensive for training and inference. In addition, it is challenging for the proposed adversarial robust GP-based optimization to handle high-dimensional problems. In contrast, our work aims to improve the generalization property in high-dimensional black-box optimization.

## 6 EMPIRICAL STUDY

In this section, we empirically evaluate the proposed SABO method, and compare it with four representative black-box methods, i.e., CMA-ES (Hansen, 2006), MMES (He et al., 2020), BES (Gao & Sener, 2022), and INGO (Lyu & Tsang, 2021). All the experiments are conducted on a single NVIDIA GeForce RTX 3090 GPU.

### 6.1 SYNTHETIC PROBLEMS

To verify the convergent results of the proposed SABO method in Section 4. We compare the proposed SABO method with baseline methods on minimizing four  $d$ -dimensional synthetic benchmark test functions, i.e., ellipsoid function,  $l_{\frac{1}{2}}$ -ellipsoid function, different powers function, and Levy function. All the test functions are listed in Appendix F.1.



The results are evaluated by calculating the Euclidean distance between the solution  $\mathbf{x}$  and the optimal solution  $\mathbf{x}^*$ , i.e.,  $\mathcal{E} = \|\mathbf{x} - \mathbf{x}^*\|_2$ . We then assess the baseline methods using varying dimensions, i.e.,  $d \in \{200, 500, 1000\}$ . Due to the page limitation, the implementation details and more detailed experimental results are put in Appendix F.1.

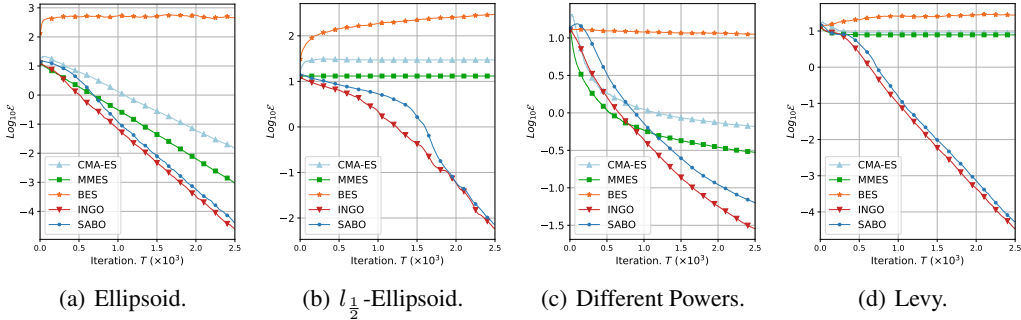


Figure 1: Results on the four test functions with problem dimension  $d = 500$  and  $N = 50$ .

**Result.** Figure 1 shows the results on four  $d$ -dimensional synthetic problems with  $d = 500$  and population size  $N = 50$ . The proposed SABO method approximately achieves a linear convergence rate similar to the INGO method. This is reasonable as these two methods have the same theoretical convergent rate, i.e.,  $\mathcal{O}(\frac{\log T}{T})$ , according to Theorem 4.3 and Theorem 5 in Lyu & Tsang (2021). Since the SABO method perturbs the main objective in each iteration, its practical convergence speed is slightly slower than INGO. The CMA-ES method and MMES methods can converge on the ellipsoid problem and the different powers problem, but they do not converge as fast as SABO and INGO methods. Moreover, CMA-ES and MMES methods fail on the  $l_{\frac{1}{2}}$ -ellipsoid problem and Levy problem. The BES method fails on all test problems. This shows that it could be challenging for BES to optimize non-smooth or high-dimensional test functions without adaptively updating mean and covariance. These results demonstrate the superiority of the SABO method in optimizing high-dimensional problems, and verify our theoretical convergence results.

## 6.2 BLACK-BOX PROMPT FINE-TUNING

Black-box prompt fine-tuning of large language models (Ding et al., 2023; Sun et al., 2022b;a; 2023) is a promising direction to achieve expertise models efficiently for downstream tasks. In such an LMaaS setting, we cannot access the model parameter and can only tune their prompts without backpropagation. We evaluate the proposed SABO method in improving generalization performance on the black-box prompt fine-tuning task.

**Datasets.** We conduct experiments on six language understanding benchmark datasets: *SST-2* (Socher et al., 2013) and *Yelp polarity* (Zhang et al., 2015) for sentiment analysis, *AG’s News* (Zhang et al., 2015) for topic classification, *MRPC* (Dolan & Brockett, 2005) for paraphrase, *RTE* (Wang et al., 2018) and *SNLI* (Bowman et al., 2015) for natural language inference. Each dataset contains a classification task. The statistics of six datasets are summarized in Table 1 of (Sun et al., 2022b). By following (Sun et al., 2022b), the testing accuracy is used to measure the performance of all the methods on the *SST-2*, *AG’s News*, *RTE*, *SNLI*, and *Yelp P*. datasets, and the F1 score is used to measure the performance on the *MRPC* datasets.

**Implementation Details.** Following Sun et al. (2022b), we employ a fixed randomly initialized matrix  $\mathbf{A} \in \mathbb{R}^{d \times D}$  to project a vector  $\mathbf{v} \in \mathbb{R}^d$  onto the token embedding space  $\mathbb{R}^D$ . Then we optimize the vector  $\mathbf{v} \in \mathbb{R}^d$  instead of directly optimizing the prompt  $\mathbf{p} \in \mathbb{R}^D$ . The pre-trained RoBERTa<sub>LARGE</sub> model (Liu et al., 2019) is used as the backbone model. The matrix  $\mathbf{A}$  is sampled from the normal distribution as described in Sun et al. (2022a), i.e.,  $\mathcal{N}(0, \frac{\sigma_e}{\sqrt{d}})$ , where  $\sigma_e$  is the standard deviation of word embeddings in RoBERTa<sub>LARGE</sub>. The templates and label words in Table 1 of Sun et al. (2022b) are used to conduct the zero-shot baseline.

Table 1: Performance (%) on *SST-2*, *AG’s News*, *MRPC*, *RTE*, *SNLI* and *Yelp P*. datasets. We report the mean and standard deviation over 3 random seeds. The best result across all groups is highlighted in **bold** and the best result in each group is marked with underlined.

| Methods                   | <i>SST-2</i>            | <i>AG’s News</i>        | <i>MRPC</i>             | <i>RTE</i>              | <i>SNLI</i>             | <i>Yelp P</i>           |
|---------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| Zero-shot                 | 79.82                   | 76.96                   | 67.40                   | 51.62                   | 38.82                   | 89.64                   |
| <i>Dimension d = 200</i>  |                         |                         |                         |                         |                         |                         |
| CMA-ES                    | 85.74 $\pm$ 0.35        | 82.09 $\pm$ 0.56        | 74.98 $\pm$ 2.16        | 51.02 $\pm$ 2.14        | 34.27 $\pm$ 1.18        | 90.57 $\pm$ 0.05        |
| MMES                      | 83.98 $\pm$ 0.78        | 80.52 $\pm$ 0.99        | 76.54 $\pm$ 4.34        | 48.50 $\pm$ 0.45        | 40.39 $\pm$ 1.83        | 90.94 $\pm$ 0.36        |
| BES                       | 83.52 $\pm$ 0.11        | 75.44 $\pm$ 0.31        | 79.23 $\pm$ 0.20        | 53.07 $\pm$ 0.29        | 38.73 $\pm$ 0.17        | 89.65 $\pm$ 0.01        |
| INGO                      | 83.57 $\pm$ 0.11        | 76.47 $\pm$ 0.03        | 78.87 $\pm$ 0.20        | 53.07 $\pm$ 0.00        | 38.86 $\pm$ 0.06        | 89.84 $\pm$ 0.04        |
| <b>SABO</b>               | <u>87.88</u> $\pm$ 0.53 | <u>82.22</u> $\pm$ 0.41 | <u>79.35</u> $\pm$ 0.12 | <b>53.67</b> $\pm$ 0.17 | <u>40.72</u> $\pm$ 0.15 | <u>91.50</u> $\pm$ 0.13 |
| <i>Dimension d = 500</i>  |                         |                         |                         |                         |                         |                         |
| CMA-ES                    | 86.12 $\pm$ 0.59        | 82.50 $\pm$ 0.23        | 77.10 $\pm$ 1.90        | 52.71 $\pm$ 0.51        | 41.34 $\pm$ 1.49        | 91.19 $\pm$ 0.44        |
| MMES                      | 85.28 $\pm$ 0.94        | 81.67 $\pm$ 0.80        | 77.31 $\pm$ 1.24        | 48.74 $\pm$ 0.59        | 42.07 $\pm$ 2.62        | 91.39 $\pm$ 0.24        |
| BES                       | 83.56 $\pm$ 0.05        | 75.93 $\pm$ 0.17        | 79.21 $\pm$ 0.09        | 52.95 $\pm$ 0.17        | 38.64 $\pm$ 0.28        | 89.62 $\pm$ 0.07        |
| INGO                      | 84.29 $\pm$ 0.34        | 76.54 $\pm$ 0.20        | 79.09 $\pm$ 0.15        | 53.19 $\pm$ 0.17        | 38.91 $\pm$ 0.10        | 89.90 $\pm$ 0.13        |
| <b>SABO</b>               | <u>87.31</u> $\pm$ 0.38 | <u>82.65</u> $\pm$ 0.59 | <u>79.62</u> $\pm$ 0.07 | <u>53.55</u> $\pm$ 0.17 | <b>42.29</b> $\pm$ 2.48 | <u>91.83</u> $\pm$ 0.16 |
| <i>Dimension d = 1000</i> |                         |                         |                         |                         |                         |                         |
| CMA-ES                    | 86.85 $\pm$ 0.57        | 82.21 $\pm$ 0.36        | 78.98 $\pm$ 0.17        | 52.35 $\pm$ 0.17        | 38.40 $\pm$ 1.83        | 90.46 $\pm$ 0.62        |
| MMES                      | 84.98 $\pm$ 0.52        | 80.86 $\pm$ 1.95        | 76.43 $\pm$ 0.82        | 49.22 $\pm$ 1.23        | 39.82 $\pm$ 3.43        | 91.63 $\pm$ 0.20        |
| BES                       | 83.11 $\pm$ 0.11        | 75.66 $\pm$ 0.09        | 79.09 $\pm$ 0.08        | 53.19 $\pm$ 0.17        | 38.57 $\pm$ 0.13        | 89.61 $\pm$ 0.04        |
| INGO                      | 84.36 $\pm$ 0.23        | 76.35 $\pm$ 0.14        | 78.97 $\pm$ 0.08        | 53.07 $\pm$ 0.29        | 39.05 $\pm$ 0.06        | 89.95 $\pm$ 0.08        |
| <b>SABO</b>               | <b>87.96</b> $\pm$ 0.83 | <b>82.77</b> $\pm$ 0.41 | <b>79.68</b> $\pm$ 0.23 | <u>53.31</u> $\pm$ 0.17 | <u>40.32</u> $\pm$ 0.27 | <b>91.96</b> $\pm$ 0.41 |

For CMA-ES, MMES, BES, INGO, and SABO methods, we employ the cross-entropy loss of training data as the black-box objective for six datasets and optimize the vector  $v$  with 100 iterations. The Gaussian distributions are initialized as  $\mu_0 = \mathbf{0}$  and  $\Sigma_0 = \mathbf{I}$ , and the population size  $N$  is set to 100. We perform a grid search for hyperparameters of INGO, SABO, and BES methods. Specifically, we search the learning rate  $\beta$  over  $\{0.1, 0.5, 1, 5\}$  for INGO, SABO, and BES, the neighborhood size  $\rho$  over  $\{10, 50, 100, 500\}$  for SABO, and the spacing  $c$  over  $\{0.1, 1, 10\}$  for BES. Additionally, we evaluate the performance of all methods on different dimensions of  $v$ , specifically  $d \in \{200, 500, 1000\}$ . All the experiments are performed in three independent runs, and the experimental results of mean objective  $\pm$  std are reported.

**Results.** Table 1 presents experimental results on these six benchmark datasets for three different dimensions of the vector  $v$ . We can see that the SABO method consistently outperforms all baselines in terms of testing classification accuracy or testing F1 scores across different settings, highlighting its effectiveness in improving generalization performance. Notably, even in the high-dimensional setting (i.e.,  $d = 1000$ ), our method maintains good performance. Moreover, we can see that when  $d = 1000$ , SABO achieves the best performance on *SST-2*, *AG’s News*, *MRPC*, and *Yelp P*. datasets. The SABO method also achieves the best performance on *RTE* and *SNLI* datasets with  $d = 200$  and  $d = 500$ , respectively.

## 7 CONCLUSION

In this work, we have introduced SABO, a novel black-box optimization algorithm that improves generalization by utilizing a sharpness-aware minimization strategy. Theoretically, we provide a convergence guarantee for the proposed SABO algorithm in both full-batch function query and mini-batch function query settings. Additionally, we prove the generalization bound for the proposed method. Empirical studies on synthetic numerical problems verify the convergence properties of the proposed method. Moreover, extensive experimental results on black-box prompt fine-tuning problems demonstrate the effectiveness of the proposed SABO method in improving the generalization performance.

## REFERENCES

- 540  
541  
542 Shun-ichi Amari. *Information geometry and its applications*, volume 194. Springer, 2016.
- 543  
544 Dara Bahri, Hossein Mobahi, and Yi Tay. Sharpness-aware minimization improves language model  
545 generalization. In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- 546  
547 Hans-Georg Beyer. Convergence analysis of evolutionary algorithms that are based on the paradigm  
548 of information geometry. *Evolutionary Computation*, 22(4):679–709, 2014.
- 549  
550 Ilija Bogunovic, Jonathan Scarlett, Stefanie Jegelka, and Volkan Cevher. Adversarially robust  
551 optimization with gaussian processes. *Advances in neural information processing systems*, 31,  
2018.
- 552  
553 Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated  
554 corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- 555  
556 Niladri S Chatterji, Behnam Neyshabur, and Hanie Sedghi. The intriguing role of module criticality  
557 in the generalization of deep networks. *arXiv preprint arXiv:1912.00528*, 2019.
- 558  
559 Can Chen, Yingxueff Zhang, Jie Fu, Xue Steve Liu, and Mark Coates. Bidirectional learning for  
560 offline infinite-width model-based optimization. *Advances in Neural Information Processing  
561 Systems*, 35:29454–29467, 2022.
- 562  
563 Edoardo Conti, Vashisht Madhavan, Felipe Petroski Such, Joel Lehman, Kenneth Stanley, and  
564 Jeff Clune. Improving exploration in evolution strategies for deep reinforcement learning via a  
565 population of novelty-seeking agents. In *Neural Information Processing Systems*, 2018.
- 566  
567 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
568 Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained  
569 language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- 570  
571 Bill Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In  
572 *Third International Workshop on Paraphrasing*, 2005.
- 573  
574 Yinpeng Dong, Hang Su, Baoyuan Wu, Zhifeng Li, Wei Liu, Tong Zhang, and Jun Zhu. Efficient  
575 decision-based black-box adversarial attacks on face recognition. In *IEEE/CVF Conference on  
576 Computer Vision and Pattern Recognition*, 2019.
- 577  
578 Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for  
579 deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint  
580 arXiv:1703.11008*, 2017.
- 581  
582 David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. Scalable  
583 global optimization via local bayesian optimization. *Advances in neural information processing  
584 systems*, 32, 2019.
- 585  
586 Clara Fannjiang and Jennifer Listgarten. Autofocused oracles for model-based design. *Advances in  
587 Neural Information Processing Systems*, 33:12945–12956, 2020.
- 588  
589 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization  
590 for efficiently improving generalization. In *International Conference on Learning Representations*,  
591 2021.
- 592  
593 Katelyn Gao and Ozan Sener. Generalizing gaussian smoothing for random search. In *International  
594 Conference on Machine Learning*, pp. 7077–7101. PMLR, 2022.
- 595  
596 Jacob Gardner, Chuan Guo, Kilian Weinberger, Roman Garnett, and Roger Grosse. Discovering and  
597 exploiting additive structure for bayesian optimization. In *Artificial Intelligence and Statistics*,  
598 2017.
- 599  
600 Nikolaus Hansen. The cma evolution strategy: a comparing review. *Towards a new evolutionary  
601 computation: Advances in the estimation of distribution algorithms*, pp. 75–102, 2006.

- 594 Xiaoyu He, Zibin Zheng, and Yuren Zhou. Mmes: Mixture model-based evolution strategy for  
595 large-scale optimization. *IEEE Transactions on Evolutionary Computation*, 25(2):320–333, 2020.  
596
- 597 Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with  
598 limited queries and information. In *International Conference on Machine Learning*, 2018.
- 599 Weisen Jiang, Hansi Yang, Yu Zhang, and James Kwok. An adaptive policy to employ sharpness-  
600 aware minimization. In *International Conference on Learning Representations*, 2023.  
601
- 602 Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive  
603 black-box functions. *Journal of Global optimization*, 13:455–492, 1998.
- 604 Minyoung Kim, Da Li, Shell X Hu, and Timothy Hospedales. Fisher SAM: Information geometry  
605 and sharpness aware minimisation. In *International Conference on Machine Learning*, 2022.  
606
- 607 Aviral Kumar and Sergey Levine. Model inversion networks for model-based optimization. *Advances  
608 in neural information processing systems*, 33:5126–5137, 2020.
- 609 Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-aware  
610 minimization for scale-invariant learning of deep neural networks. In *International Conference on  
611 Machine Learning*, 2021.
- 612 Tianyi Lin, Zeyu Zheng, and Michael Jordan. Gradient-free methods for deterministic and stochastic  
613 nonsmooth nonconvex optimization. *Advances in Neural Information Processing Systems*, 35:  
614 26160–26175, 2022.  
615
- 616 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
617 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
618 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 619 Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scal-  
620 able sharpness-aware minimization. In *IEEE/CVF Conference on Computer Vision and Pattern  
621 Recognition*, 2022.  
622
- 623 Yueming Lyu and Ivor W Tsang. Black-box optimizer with stochastic implicit natural gradient. In  
624 *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery  
625 in Databases*, 2021.
- 626 David A McAllester. Pac-bayesian model averaging. In *Annual Conference on Computational  
627 Learning Theory*, 1999.
- 628 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and  
629 projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.  
630
- 631 Seyedali Mirjalili. Evolutionary algorithms and neural networks. In *Studies in computational  
632 intelligence*, volume 780. Springer, 2019.
- 633 Thomas Möllenhoff and Mohammad Emtiyaz Khan. Sam as an optimal relaxation of bayes. In *The  
634 Eleventh International Conference on Learning Representations*, 2022.  
635
- 636 Amin Nayebi, Alexander Munteanu, and Matthias Poloczek. A framework for bayesian optimization  
637 in embedded subspaces. In *International Conference on Machine Learning*, 2019.
- 638 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.  
639 *Foundations of Computational Mathematics*, 17(2):527–566, 2017.  
640
- 641 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generaliza-  
642 tion in deep learning. In *Neural Information Processing Systems*, 2017.
- 643 Quan Nguyen, Kaiwen Wu, Jacob Gardner, and Roman Garnett. Local bayesian optimization  
644 via maximizing probability of descent. *Advances in neural information processing systems*, 35:  
645 13190–13202, 2022.  
646
- 647 Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, and Mario Boley. Relative  
flatness and generalization. In *Neural Information Processing Systems*, 2021.

- 648 Han Qi, Yi Su, Aviral Kumar, and Sergey Levine. Data-driven offline decision-making via invariant  
649 representation learning. *Advances in Neural Information Processing Systems*, 35:13226–13237,  
650 2022.
- 651 Shogo Sagawa and Hideitsu Hino. Gradual domain adaptation via normalizing flows. *arXiv preprint*  
652 *arXiv:2206.11492*, 2022.
- 654 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a  
655 scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- 656 Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence  
657 results and optimal averaging schemes. In *International Conference on Machine Learning*, 2013.
- 659 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng,  
660 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment  
661 treebank. In *Conference on Empirical Methods in Natural Language Processing*, 2013.
- 662 Mandavilli Srinivas and Lalit M Patnaik. Genetic algorithms: A survey. *Computer*, 27(6):17–26,  
663 1994.
- 665 Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimiza-  
666 tion in the bandit setting: no regret and experimental design. In *International Conference on*  
667 *Machine Learning*, 2010.
- 668 Tianxiang Sun, Zhengfu He, Hong Qian, Yunhua Zhou, Xuan-Jing Huang, and Xipeng Qiu. Bbtv2:  
669 towards a gradient-free future with large language models. In *Conference on Empirical Methods*  
670 *in Natural Language Processing*, 2022a.
- 672 Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for  
673 language-model-as-a-service. In *International Conference on Machine Learning*, 2022b.
- 674 Tianxiang Sun, Zhengfu He, Qin Zhu, Xipeng Qiu, and Xuan-Jing Huang. Multitask pre-training  
675 of modular prompt for chinese few-shot learning. In *Annual Meeting of the Association for*  
676 *Computational Linguistics*, 2023.
- 677 Brandon Trabucco, Aviral Kumar, Xinyang Geng, and Sergey Levine. Conservative objective models  
678 for effective offline model-based optimization. In *International Conference on Machine Learning*,  
679 pp. 10358–10368. PMLR, 2021.
- 681 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue:  
682 A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint*  
683 *arXiv:1804.07461*, 2018.
- 684 Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber.  
685 Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014.
- 687 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking  
688 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 689 Feiyang Ye, Yueming Lyu, Xuehao Wang, Yu Zhang, and Ivor Tsang. Adaptive stochastic gradient  
690 algorithm for black-box multi-objective learning. In *The Twelfth International Conference on*  
691 *Learning Representations*, 2024.
- 692 Haishan Ye. Mirror natural evolution strategies. *arXiv preprint arXiv:2308.00469*, 2023.
- 694 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep  
695 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,  
696 2021.
- 698 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text  
699 classification. In *Neural Information Processing Systems*, 2015.
- 700 Yang Zhao, Hao Zhang, and Xiuyuan Hu. Penalizing gradient norm for efficiently improving  
701 generalization in deep learning. In *International Conference on Machine Learning*, 2022.

702 Juntang Zhuang, Boqing Gong, Liangzhe Yuan, Yin Cui, Hartwig Adam, Nicha Dvornek, Sekhar  
703 Tatikonda, James Duncan, and Ting Liu. Surrogate gap minimization improves sharpness-aware  
704 training. In *International Conference on Learning Representations*, 2022.  
705  
706 Juliusz Krzysztof Ziomek and Haitham Bou Ammar. Are random decompositions all we need in  
707 high dimensional bayesian optimisation? In *International Conference on Machine Learning*, pp.  
708 43347–43368. PMLR, 2023.  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A ADDITIONAL MATERIAL FOR SECTION 3

### A.1 DETERMINE THE PERTURBATION $\delta$

The Lagrangian of problem (9) is

$$\mathcal{L}(\delta, \lambda) = -\langle \nabla_{\theta} J(\theta), \delta \rangle + \lambda(\text{KL}(p_{\theta+\delta} \| p_{\theta}) - \rho^2) \quad (29)$$

$$\begin{aligned} &= -\delta_{\mu}^{\top} \nabla_{\mu} J(\theta_t) - \text{tr}(\delta_{\Sigma} \nabla_{\Sigma} J(\theta_t)) \\ &\quad + \frac{\lambda}{2} \left[ \text{tr}(\Sigma^{-1}(\Sigma + \delta_{\Sigma})) + \delta_{\mu}^{\top} \Sigma^{-1} \delta_{\mu} + \log \frac{|\Sigma|}{|(\Sigma + \delta_{\Sigma})|} - d \right] - \lambda \rho^2. \end{aligned} \quad (30)$$

Taking the derivative  $\delta_{\mu}$  and  $\delta_{\Sigma}$  and setting them to zero, we can obtain that

$$\nabla_{\mu} J(\theta) - \lambda \Sigma^{-1} \delta_{\mu} = 0, \quad (31)$$

$$\nabla_{\Sigma} J(\theta) - \frac{\lambda}{2} [(\Sigma + \delta_{\Sigma})^{-1} - \Sigma^{-1}] = 0. \quad (32)$$

Note that  $\Sigma$  is a diagonal matrix. Therefore, we can achieve that

$$\delta_{\mu}(\theta) = \frac{1}{\lambda} \Sigma \nabla_{\mu} J(\theta), \quad (33)$$

$$\delta_{\Sigma}(\theta) = \frac{2 \Sigma \nabla_{\Sigma} J(\theta)}{\lambda \Sigma^{-1} - 2 \nabla_{\Sigma} J(\theta)}. \quad (34)$$

### A.2 DETERMINE THE OPTIMAL $\lambda$

Note that we have

$$\text{KL}(p_{\theta+\delta} \| p_{\theta}) = \frac{1}{2} \left[ \text{tr}(\Sigma^{-1}(\Sigma + \delta_{\Sigma})) + \delta_{\mu}^{\top} \Sigma^{-1} \delta_{\mu} + \log \frac{|\Sigma|}{|(\Sigma + \delta_{\Sigma})|} - d \right] \quad (35)$$

$$= \frac{1}{2} \left[ \text{tr}(I + \frac{2}{\lambda} \nabla_{\Sigma} J(\theta)(\Sigma + \delta_{\Sigma})) + \delta_{\mu}^{\top} \Sigma^{-1} \delta_{\mu} + \log \frac{|\Sigma|}{|(\Sigma + \delta_{\Sigma})|} - d \right] \quad (36)$$

$$= \frac{1}{2} \left[ \text{tr}(\frac{2}{\lambda} \nabla_{\Sigma} J(\theta) \delta_{\Sigma}) + \delta_{\mu}^{\top} \Sigma^{-1} \delta_{\mu} + Q \right], \quad (37)$$

where

$$Q = \left| \frac{2}{\lambda} \Sigma \nabla_{\Sigma} J(\theta) \right| + \log(I - \frac{2}{\lambda} \Sigma \nabla_{\Sigma} J(\theta)). \quad (38)$$

Since  $\Sigma$  and  $\nabla_{\Sigma} J(\theta)$  are both diagonal matrix, we denote  $\text{diag}(\Sigma \nabla_{\Sigma} J(\theta)) = (v^1, \dots, v^d)$ , then we have

$$Q = \log\left(\prod_{i=1}^d (1 - \frac{2}{\lambda} v^i)\right) + \sum_{i=1}^d \frac{2}{\lambda} v^i = \sum_{i=1}^d \left( \log(1 - \frac{2}{\lambda} v^i) + \frac{2}{\lambda} v^i \right) = \sum_{i=1}^d -\frac{2}{\lambda^2} (v^i)^2 - \mathcal{O}\left(\frac{1}{\lambda^3} (v^i)^3\right). \quad (39)$$

We denote  $\text{diag}(\nabla_{\Sigma} J(\theta)) = (\hat{v}^1, \dots, \hat{v}^d)$ , then we have  $\hat{v}^i \sigma^i = v^i$  and

$$\text{tr}\left(\frac{2}{\lambda} \nabla_{\Sigma} J(\theta) \delta_{\Sigma}\right) = \sum_{i=1}^d \frac{2}{\lambda} \hat{v}^i \left( \frac{1}{(\sigma^i)^{-1} - \frac{2}{\lambda} \hat{v}^i} - \sigma^i \right) = \sum_{i=1}^d \frac{2}{\lambda} \hat{v}^i \sigma^i \left( \frac{2}{\lambda} \sigma^i \hat{v}^i \right) + \mathcal{O}\left(\frac{4}{\lambda^2} (\sigma^i \hat{v}^i)^2\right). \quad (40)$$

Then substituting  $\delta_{\mu}(\theta)$  and  $\delta_{\Sigma}(\theta)$  into the inequality  $\text{KL}(p_{\theta+\delta} \| p_{\theta_t}) \leq \rho^2$ , we can obtain that

$$\text{KL}(p_{\theta+\delta} \| p_{\theta}) = \frac{1}{2} \left[ \frac{2}{\lambda^2} \|\Sigma \nabla_{\Sigma} J(\theta)\|_{\text{F}}^2 + \frac{1}{\lambda^2} \|\Sigma^{\frac{1}{2}} \nabla_{\mu} J(\theta)\|_2^2 \right] + \epsilon \leq \rho^2, \quad (41)$$

where  $\epsilon = \mathcal{O}\left(\frac{4(\sigma^i \hat{v}^i)^2}{\lambda^2}\right) = \mathcal{O}\left(\frac{1}{\lambda^2}\right)$ . Let the equality holds and solve Eq. (41), we have

$$\lambda \approx \frac{1}{\rho} \sqrt{\|\Sigma \nabla_{\Sigma} J(\theta)\|_{\text{F}}^2 + 0.5 \|\Sigma^{\frac{1}{2}} \nabla_{\mu} J(\theta)\|_2^2}. \quad (42)$$

### 810 A.3 UPDATE RULE FOR $\theta_t$

811 Note that we have

$$\begin{aligned}
& \langle \theta - \theta_t, \nabla_{\theta} J(\theta_t + \delta_t) \rangle + \frac{1}{\beta_t} \text{KL}(p_{\theta} \| p_{\theta_t}) \\
&= (\boldsymbol{\mu} - \boldsymbol{\mu}_t)^{\top} \nabla_{\boldsymbol{\mu}} J(\theta_t + \delta_t) + \text{tr}((\boldsymbol{\Sigma} - \boldsymbol{\Sigma}_t) \nabla_{\boldsymbol{\Sigma}} J(\theta_t + \delta_t)) \\
&\quad + \frac{1}{2\beta_t} \left[ \text{tr}(\boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \boldsymbol{\mu}_t)^{\top} \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_t) + \log \frac{|\boldsymbol{\Sigma}_t|}{|\boldsymbol{\Sigma}|} - d \right],
\end{aligned}$$

812 where  $\nabla_{\boldsymbol{\mu}} J(\theta_t + \delta_t)$  and  $\nabla_{\boldsymbol{\Sigma}} J(\theta_t + \delta_t)$  denotes the derivative w.r.t.  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  taking at  $\boldsymbol{\mu} = \boldsymbol{\mu}_t + \delta_{\boldsymbol{\mu}_t}$   
813 and  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_t + \delta_{\boldsymbol{\Sigma}_t}$ , respectively. We can see the above problem is convex with respect to  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ .  
814 Taking the derivative w.r.t.  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  and setting them to zero, we can obtain that

$$815 \nabla_{\boldsymbol{\mu}} J(\theta_t + \delta_t) + \frac{1}{\beta_t} \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_t) = 0, \quad (43)$$

$$816 \nabla_{\boldsymbol{\Sigma}} J(\theta_t + \delta_t) + \frac{1}{2\beta_t} [\boldsymbol{\Sigma}_t^{-1} - \boldsymbol{\Sigma}^{-1}] = 0. \quad (44)$$

817 Therefore, we obtain the update rule for  $\theta_t$  as

$$818 \boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \nabla_{\boldsymbol{\mu}} J(\theta_t + \delta_t), \quad (45)$$

$$819 \boldsymbol{\Sigma}_{t+1}^{-1} = \boldsymbol{\Sigma}_t^{-1} + 2\beta_t \nabla_{\boldsymbol{\Sigma}} J(\theta_t + \delta_t). \quad (46)$$

### 820 A.4 THE GRADIENTS OF $J(\theta)$

821 To obtain the gradients of the reparameterized objective  $J(\theta)$  w.r.t.  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , we use the following  
822 theorem to show that only function queries are needed.

823 **Theorem A.1** (Wierstra et al., 2014) *The gradient of the expectation of an integrable function  $F(\mathbf{x})$   
824 under a Gaussian distribution  $p_{\theta} := \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the mean  $\boldsymbol{\mu}$  and the covariance  $\boldsymbol{\Sigma}$  can  
825 be expressed as*

$$826 \nabla_{\boldsymbol{\mu}} J(\theta) = \mathbb{E}_{p_{\theta}} [\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) F(\mathbf{x})], \quad (47)$$

$$827 \nabla_{\boldsymbol{\Sigma}} J(\theta) = \frac{1}{2} \mathbb{E}_{p_{\theta}} [(\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1}) F(\mathbf{x})]. \quad (48)$$

## 828 B MINI-BATCH SABO

829 For the proposed SABO with a mini-batch function query, the stochastic approximation of the  
830 gradients  $\nabla_{\boldsymbol{\mu}} J(\theta_t)$  and  $\nabla_{\boldsymbol{\Sigma}} J(\theta_t)$  using Monte Carlo sampling are given as

$$831 \mathbf{g}'_t = \frac{1}{NM} \sum_{j=1}^N \sum_{i=1}^M \boldsymbol{\Sigma}_t^{-1} (\mathbf{x}'_j - \boldsymbol{\mu}_t) (F(\mathbf{x}'_j; (X'_i, y'_i)) - F(\boldsymbol{\mu}_t; (X'_i, y'_i))), \quad (49)$$

$$\begin{aligned}
\mathbf{G}'_t &= \frac{1}{2NM} \sum_{j=1}^N \sum_{i=1}^M \text{diag} \left[ \boldsymbol{\Sigma}_t^{-1} \left[ \text{diag}((\mathbf{x}'_j - \boldsymbol{\mu}_t)(\mathbf{x}'_j - \boldsymbol{\mu}_t)^{\top} \boldsymbol{\Sigma}_t^{-1} - \mathbf{I}) \right. \right. \\
&\quad \left. \left. \times (F(\mathbf{x}'_j; (X'_i, y'_i)) - F(\boldsymbol{\mu}_t; (X'_i, y'_i))) \right] \right], \quad (50)
\end{aligned}$$

832 where  $\mathbf{x}'_j$  denotes the  $j$ -th sample sampled from the distribution  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ , and  $(X'_i, y'_i)$  denotes  
833 the  $i$ -th data in the mini-batch dataset  $\mathcal{B}'$ . Then by setting  $\delta_t = \left\{ \frac{1}{\lambda} \boldsymbol{\Sigma}_t \mathbf{g}'_t, \frac{2\boldsymbol{\Sigma}_t \mathbf{G}'_t}{\lambda \boldsymbol{\Sigma}_t^{-1} - 2\mathbf{G}'_t} \right\}$ , the gradients



**Algorithm 2** Mini-batch SABO**Require:** Neighborhood size  $\rho$ , learning rate  $\beta_t$ , batch size  $M$ 

- 1: Initialized  $\theta_0 = (\mu_0, \Sigma_0)$ ;
- 2: **for**  $t = 0$  to  $T - 1$  **do**
- 3:   Take i.i.d. samples  $\mathbf{z}'_j \sim \mathcal{N}(0, I)$  and set  $\mathbf{x}'_j = \mu_t + \Sigma_t^{\frac{1}{2}} \mathbf{z}'_j$  for  $j \in \{1, \dots, N\}$ ;
- 4:   Sample a batch of data  $\mathcal{B}'$ ;
- 5:   Query the batch observations  $\{F(\mathbf{x}'_1; \mathcal{B}'), \dots, F(\mathbf{x}'_N; \mathcal{B}')\}$ ;
- 6:   Set the gradient  $\mathbf{g}'_t$  via Eq. (49) and set the gradient  $\mathbf{G}'_t$  via Eq. (50);
- 7:   Compute  $\lambda = \frac{1}{\rho} \sqrt{\|\Sigma_t \mathbf{G}'_t\|_F^2 + 0.5 \|\Sigma_t^{\frac{1}{2}} \mathbf{g}'_t\|_2^2}$ ;
- 8:   Compute  $\delta_{\mu_t} = \frac{1}{\lambda} \Sigma_t \mathbf{g}'_t$  and  $\delta_{\Sigma_t} = \frac{2 \Sigma_t \mathbf{G}'_t}{\lambda \Sigma_t^{-1} - 2 \mathbf{G}'_t}$ ;
- 9:   Take i.i.d. samples  $\mathbf{z}_j \sim \mathcal{N}(0, I)$  for  $j \in \{1, \dots, N\}$ ;
- 10:   Set  $\mathbf{x}_j = \mu_t + \delta_{\mu_t} + (\Sigma_t + \delta_{\Sigma_t})^{\frac{1}{2}} \mathbf{z}_j$  for  $j \in \{1, \dots, N\}$ ;
- 11:   Sample a batch of data  $\mathcal{B}$ ;
- 12:   Query the batch observations  $\{F(\mathbf{x}_1; \mathcal{B}), \dots, F(\mathbf{x}_N; \mathcal{B})\}$ ;
- 13:   Set the gradient  $\mathbf{g}_t$  via Eq. (51) and set the gradient  $\mathbf{G}_t$  via Eq. (52);
- 14:   Set  $\mu_{t+1} = \mu_t - \beta_t \Sigma_t \mathbf{g}_t$ ;
- 15:   Set  $\Sigma_{t+1}^{-1} = \Sigma_t^{-1} + 2\beta_t \mathbf{G}_t$ ;
- 16: **end for**
- 17: **return**  $\theta_T = (\mu_T, \Sigma_T)$ .

$\nabla_{\mu} J(\theta_t + \delta_t)$  and  $\nabla_{\Sigma} J(\theta_t + \delta_t)$  can also be approximated as

$$\mathbf{g}_t = \frac{1}{NM} \sum_{j=1}^N \sum_{i=1}^M \widehat{\Sigma}_t^{-1} (\mathbf{x}_j - \widehat{\mu}_t) (F(\mathbf{x}_j; (X_i, y_i)) - F(\widehat{\mu}_t; (X_i, y_i))), \quad (51)$$

$$\mathbf{G}_t = \frac{1}{2NM} \sum_{j=1}^N \sum_{i=1}^M \text{diag} \left[ \widehat{\Sigma}_t^{-1} \left[ \text{diag}((\mathbf{x}_j - \widehat{\mu}_t)(\mathbf{x}_j - \widehat{\mu}_t)^\top \widehat{\Sigma}_t^{-1} - \mathbf{I}) \right. \right. \\ \left. \left. \times (F(\mathbf{x}_j; (X_i, y_i)) - F(\widehat{\mu}_t; (X_i, y_i))) \right] \right], \quad (52)$$

where  $\widehat{\Sigma}_t = \Sigma_t + \delta_{\Sigma_t}$ ,  $\widehat{\mu}_t = \mu_t + \delta_{\mu_t}$ ,  $\mathbf{x}_j$  denotes the  $j$ -th sample sampled from the distribution  $\mathcal{N}(\mathbf{x} \mid \widehat{\mu}_t, \widehat{\Sigma}_t)$ , and  $(X_i, y_i)$  denotes the  $i$ -th data in the mini-batch dataset  $\mathcal{B}$ . Then we can update  $\theta_t$  by Eq. (22). The entire algorithm of SABO with a mini-batch function query is shown in Algorithm 2.

## C TECHNICAL LEMMAS

In this section, we introduce the following technical lemmas for analysis. The proof of all technical lemmas is put in Appendix E.

**Lemma C.1** Suppose  $\Sigma$  and  $\widehat{\Sigma}$  are two  $d$ -dimensional diagonal matrix and  $\mathbf{z}$  is a  $d$ -dimensional vector, then we have  $\|\Sigma \mathbf{z}\| \leq \|\Sigma\|_F \|\mathbf{z}\|$  and  $\|\Sigma \widehat{\Sigma}\|_F \leq \|\Sigma\|_F \|\widehat{\Sigma}\|_F$ .

**Lemma C.2** Given a convex function  $F(\mathbf{x})$ , for Gaussian distribution with parameters  $\theta := \{\mu, \Sigma^{\frac{1}{2}}\}$ , let  $J(\theta) := \mathbb{E}_{p(\mathbf{x}; \theta)}[F(\mathbf{x})]$ . Then  $J(\theta)$  is a convex function with respect to  $\theta$ .

**Lemma C.3** Suppose that the gradient  $\mathbf{G}_t$  are positive semi-definite matrix and satisfies  $\xi \mathbf{I} \preceq \mathbf{G}_t \preceq b \mathbf{I}$ . Then for algorithm 1 and 2, we have the following results.

- (a) The (diagonal) covariance matrix  $\Sigma_T$  satisfies  $\frac{1}{2b \sum_{t=1}^T \beta_t \mathbf{I} + \Sigma_0^{-1}} \preceq \Sigma_T \preceq \frac{1}{2\xi \sum_{t=1}^T \beta_t \mathbf{I} + \Sigma_0^{-1}}$ .
- (b) The Frobenius norm for the covariance matrix  $\Sigma_t$  satisfies  $\|\Sigma_t\|_F \leq \frac{\sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}$ .

**Lemma C.4** For given  $\theta_t$ , denote the approximated gradients of  $\nabla_{\boldsymbol{\mu}} J(\theta_t)$  and  $\nabla_{\boldsymbol{\Sigma}} J(\theta_t)$  by  $\mathbf{g}'_t$  and  $\mathbf{G}'_t$ , respectively. Then for Algorithm 1 and Algorithm 2, if  $\rho \leq \frac{\sqrt{d}}{2}$ , then the perturbation satisfies,  $\|\delta_{\boldsymbol{\mu}_t}\| \leq \rho\sqrt{2}\|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_F$  and  $\|\delta_{\boldsymbol{\Sigma}_t}\|_F \leq 2\rho r\|\boldsymbol{\Sigma}_t\|_F$ , where  $r$  is a positive constant.

**Lemma C.5** In Algorithm 1, suppose the gradient estimator  $\mathbf{g}'_t$  in  $t$ -th iteration as

$$\mathbf{g}'_t = \boldsymbol{\Sigma}_t^{-\frac{1}{2}} \mathbf{z} (F(\boldsymbol{\mu}_t + \boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{z}) - F(\boldsymbol{\mu}_t)), \quad (53)$$

where  $\mathbf{z} \sim \mathcal{N}(0, I)$ . Then  $\mathbf{g}'_t$  is an unbiased estimator of the gradient  $\nabla_{\boldsymbol{\mu}} \mathbb{E}_{p_{\theta_t}}[F(\mathbf{x})]$ .

**Lemma C.6** In Algorithm 1, suppose the gradient estimator  $\mathbf{g}_t$  in  $t$ -th iteration as

$$\mathbf{g}_t = \widehat{\boldsymbol{\Sigma}}_t^{-\frac{1}{2}} \mathbf{z} (F(\widehat{\boldsymbol{\mu}}_t + \widehat{\boldsymbol{\Sigma}}_t^{\frac{1}{2}} \mathbf{z}) - F(\widehat{\boldsymbol{\mu}}_t)), \quad (54)$$

where  $\mathbf{z} \sim \mathcal{N}(0, I)$ . Suppose that Assumption 4.1 holds,  $\rho < \frac{\sqrt{d}}{2}$ , the gradient  $\mathbf{G}_t$  is positive semi-definite matrix and satisfies  $\xi \mathbf{I} \preceq \mathbf{G}_t \preceq b\mathbf{I}$  and  $\boldsymbol{\Sigma}_0 \preceq R\mathbf{I}$ , where  $\xi, b, R \geq 0$ . Then we have

$$\mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2] \leq \frac{L_F^2 (1 + 2\rho r') (d + 4)^2}{2\xi (\sum_{k=1}^t \beta_k)}, \quad (55)$$

where  $r'$  is a positive constant.

**Lemma C.7** In Algorithm 2, suppose assumption 4.4 holds and the gradient estimator  $\mathbf{g}'_t$  in  $t$ -th iteration as

$$\mathbf{g}'_t = \boldsymbol{\Sigma}_t^{-\frac{1}{2}} \mathbf{z} (F(\boldsymbol{\mu}_t + \boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{z}; \mathcal{B}) - F(\boldsymbol{\mu}_t; \mathcal{B})), \quad (56)$$

where  $\mathbf{z} \sim \mathcal{N}(0, I)$ . Then  $\mathbf{g}'_t$  is an unbiased estimator of the gradient  $\nabla_{\boldsymbol{\mu}} \mathbb{E}_{p_{\theta_t}}[F(\mathbf{x})]$ .

**Lemma C.8** In Algorithm 2, suppose the gradient estimator  $\mathbf{g}_t$  in  $t$ -th iteration as

$$\mathbf{g}_t = \widehat{\boldsymbol{\Sigma}}_t^{-\frac{1}{2}} \mathbf{z} (F(\widehat{\boldsymbol{\mu}}_t + \widehat{\boldsymbol{\Sigma}}_t^{\frac{1}{2}} \mathbf{z}; \mathcal{B}) - F(\widehat{\boldsymbol{\mu}}_t; \mathcal{B})), \quad (57)$$

where  $\mathbf{z} \sim \mathcal{N}(0, I)$ . Suppose that Assumption 4.1 and Assumption 4.4 hold,  $\rho < \frac{\sqrt{d}}{2}$ , the gradient  $\mathbf{G}_t$  is positive semi-definite matrix and satisfies  $\xi \mathbf{I} \preceq \mathbf{G}_t \preceq b\mathbf{I}$  and  $\boldsymbol{\Sigma}_0 \preceq R\mathbf{I}$ , where  $\xi, b, R \geq 0$ . Then we have

$$\mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2] \leq \frac{L_F^2 (1 + 2\rho r') (d + 4)^2}{6\xi (\sum_{k=1}^t \beta_k)} + \frac{2(d + 4)\varepsilon^2}{3}, \quad (58)$$

where  $r'$  is a positive constant,  $\varepsilon^2 = 2\varepsilon_B^2 + 2\varepsilon_D^2$ .

## D PROOF OF THE RESULT IN SECTION 4

In this section, we provide the proof of the result in Section 4.

### D.1 PROOF OF THE PROPOSITION 4.2

Note that  $F(\mathbf{x})$  is a convex function, we have

$$F(\boldsymbol{\mu}) = F(\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})}[\mathbf{x}]) \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})}[F(\mathbf{x})] = J(\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (59)$$

Since  $F(\boldsymbol{\mu}^*) = J(\boldsymbol{\mu}^*, \mathbf{0})$ , it follows that

$$F(\boldsymbol{\mu}) - F(\boldsymbol{\mu}^*) \leq J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) - J(\boldsymbol{\mu}^*, \mathbf{0}), \quad (60)$$

where we reach the conclusion.

## D.2 PROOF OF THEOREM 4.3

The update rule of  $\boldsymbol{\mu}$  can be represented as  $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \mathbf{g}_t$ . Since  $F(\mathbf{x})$  is convex function, we have  $J(\boldsymbol{\theta})$  is convex w.r.t.  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}^{\frac{1}{2}}\}$  by Lemma C.2, together with  $J(\boldsymbol{\theta})$  is  $c$ -strongly convex w.r.t.  $\boldsymbol{\mu}$  we obtain

$$J(\boldsymbol{\theta}_t) \leq J(\boldsymbol{\mu}^*, 0) + \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) + \nabla_{\boldsymbol{\Sigma}^{\frac{1}{2}}} J(\boldsymbol{\theta}_t)^\top \boldsymbol{\Sigma}_t^{\frac{1}{2}} - \frac{c}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2. \quad (61)$$

Note that  $\nabla_{\boldsymbol{\Sigma}^{\frac{1}{2}}} J(\boldsymbol{\theta}_t) = \boldsymbol{\Sigma}_t^{\frac{1}{2}} \nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t) + \nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t) \boldsymbol{\Sigma}_t^{\frac{1}{2}}$ , we have

$$J(\boldsymbol{\theta}_t) \leq J(\boldsymbol{\mu}^*, 0) + \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) + 2 \nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t) \boldsymbol{\Sigma}_t - \frac{c}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2. \quad (62)$$

Let  $A_t = J(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, 0)$ , we have

$$\beta_t \mathbb{E}[A_t] \leq \beta_t \mathbb{E}_{\mathbf{z}}[\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] + 2\beta_t \mathbb{E}_{\mathbf{z}}[\nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t)^\top \boldsymbol{\Sigma}_t] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2. \quad (63)$$

Note that

$$\begin{aligned} & \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 \\ &= \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - \|\boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \mathbf{g}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 \end{aligned} \quad (64)$$

$$= \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - (\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - 2\beta_t \langle \boldsymbol{\mu}_t - \boldsymbol{\mu}^*, \mathbf{g}_t \rangle + \beta_t^2 \langle \boldsymbol{\Sigma}_t \mathbf{g}_t, \mathbf{g}_t \rangle) \quad (65)$$

$$= 2\beta_t \mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) - \beta_t^2 (\boldsymbol{\Sigma}_t \mathbf{g}_t)^\top \mathbf{g}_t. \quad (66)$$

Therefore we have

$$\mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) = \frac{1}{2\beta_t} (\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2) + \frac{\beta_t}{2} (\boldsymbol{\Sigma}_t \mathbf{g}_t)^\top \mathbf{g}_t. \quad (67)$$

According to Lemma C.5, we have  $\mathbb{E} \mathbf{g}_t = \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t)$ . Then we have

$$\begin{aligned} & \mathbb{E}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \\ &= \mathbb{E}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t))^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] + \mathbb{E}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \end{aligned} \quad (68)$$

$$\leq \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \|\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t)\| \quad (69)$$

$$\leq DL \|\boldsymbol{\delta}_t\| \quad (70)$$

$$\leq DL \rho_t U_t, \quad (71)$$

where the first inequality is due to the Cauchy-Schwarz inequality, the second inequality is due to  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$  and Assumption 4.1, the last inequality is due to Lemma C.4, and  $U_t = \max(2\sqrt{2} \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}}, 4r \|\boldsymbol{\Sigma}_t\|_{\text{F}})$ . Then we have

$$\mathbb{E}_{\mathbf{z}} [\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \quad (72)$$

$$= \mathbb{E}_{\mathbf{z}} [\mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) + (\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \quad (73)$$

$$\leq \frac{1}{2\beta_t} \mathbb{E}_{\mathbf{z}} [\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2] + \beta_t (\boldsymbol{\Sigma}_t \mathbf{g}_t)^\top \mathbf{g}_t + DL \rho_t U_t, \quad (74)$$

where the inequality is due to Eq. (67), Eq. (71), and  $\beta_t \geq 0$ . Note that

$$\mathbb{E}[\nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t)^\top \boldsymbol{\Sigma}_t] \leq \mathbb{E}[\|\nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t)\|] \|\boldsymbol{\Sigma}_t\|_{\text{F}} \leq H \|\boldsymbol{\Sigma}_t\|_{\text{F}}, \quad (75)$$

where the first inequality is due to Lemma C.1 and the second inequality is due to the Lipschitz continuous assumption of the function  $J(\boldsymbol{\theta})$ . Then substituting Eq. (74) and Eq. (75) into Eq. (63) and multiplying  $\beta_t$  on both sides of the inequality, we have

$$\begin{aligned} \beta_t \mathbb{E}[A_t] &\leq \frac{1}{2} \mathbb{E}[\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t^{-1}}^2] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 + 2\beta_t^2 H \|\boldsymbol{\Sigma}_t\|_{\text{F}} \\ &\quad + DL \rho_t U_t + \beta_t^2 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2], \end{aligned} \quad (76)$$

We further obtain that

$$\sum_{t=0}^{T-1} \left[ \frac{1}{2} \mathbb{E} [\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\Sigma_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\Sigma_{t+1}^{-1}}^2] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] \quad (77)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\Sigma_t^{-1}}^2 - \|\boldsymbol{\mu}_{t-1} - \boldsymbol{\mu}^*\|_{\Sigma_{t-1}^{-1}}^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] \quad (78)$$

$$+ \frac{1}{2} \left[ \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}^*\|_{\Sigma_0^{-1}}^2 - \|\boldsymbol{\mu}_T - \boldsymbol{\mu}^*\|_{\Sigma_{T-1}^{-1}}^2 \right] \quad (79)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{2\beta_t \mathbf{G}_t}^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] + \|\Sigma_0^{-1}\|_{\text{F}} D^2 \quad (80)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] + \|\Sigma_0^{-1}\|_{\text{F}} D^2 \quad (81)$$

$$= \|\Sigma_0^{-1}\|_{\text{F}} D^2, \quad (82)$$

where the second inequality is due to the update rule of  $\Sigma_t$  and  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$ , and the third inequality is due to Cauchy-Schwarz inequality and  $\mathbf{G}_t \preceq \frac{c}{4} \mathbf{I}$ . Then we have

$$\sum_{t=0}^{T-1} \beta_t \mathbb{E}[A_t] \leq \|\Sigma_0^{-1}\|_{\text{F}} D^2 + \sum_{t=0}^{T-1} \left( 2H\beta_t^2 \|\Sigma_t\|_{\text{F}} + DL\rho_t U_t + \beta_t \|\Sigma_t^{\frac{1}{2}}\|_{\text{F}} \mathbb{E} \|\Sigma_t^{\frac{1}{2}} \mathbf{g}_t\|^2 \right) \quad (83)$$

$$\leq \|\Sigma_0^{-1}\|_{\text{F}} D^2 + \sum_{t=0}^{T-1} \left( \frac{2H\beta_t^2 \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k} + DL\rho_t U_t + \frac{L_F^2 R^{\frac{1}{2}} (1 + 2\rho r')(d+4)^2 \beta_t^2}{2\xi (\sum_{k=1}^t \beta_k)} \right), \quad (84)$$

where the first inequality is due to Eq. (82) and Eq. (76), and the second inequality is due to Lemma C.6 and  $\|\Sigma_t^{\frac{1}{2}}\|_{\text{F}} \leq R^{\frac{1}{2}}$ . According to C.3 (b), we have  $U_t \leq \max\left(\frac{2\sqrt{2}d^{\frac{1}{4}}}{(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}}, \frac{4r\sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}\right)$ .

Therefore, there exists a constant  $t^*$ , when  $t > t^* - 1$ ,  $U_t \leq \frac{2\sqrt{2}d^{\frac{1}{4}}}{(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}}$ . Denote  $\sum_{t=0}^{t^*-1} \frac{4r\rho_t \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}$  by a constant  $\Gamma$ . Then if  $T > t^*$ , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[A_t] &\leq \frac{\|\Sigma_0^{-1}\|_{\text{F}} D^2}{T\beta_t} + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{2H\beta_t \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k} + \frac{L_F^2 R^{\frac{1}{2}} (1 + 2\rho r')(d+4)^2 \beta_t}{2\xi \sum_{k=1}^t \beta_k} \right) \\ &\quad + \frac{1}{T} \sum_{t=t^*}^{T-1} \left( \frac{2\sqrt{2}DL\rho_t d^{\frac{1}{4}}}{\beta_t (2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}} \right) + \frac{DL\Gamma}{T\beta_t}, \end{aligned} \quad (85)$$

Let  $\beta_t = \beta$  and  $\rho_t = \frac{\rho_0}{\sqrt{t}}$ , we can obtain that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[A_t] \leq \frac{\|\Sigma_0^{-1}\|_{\text{F}} D^2}{T\beta} + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{C}{2\xi t} \right) + \frac{1}{T} \sum_{t=t^*}^{T-1} \left( \frac{2\sqrt{2}DL\rho_0 d^{\frac{1}{4}}}{\sqrt{2\xi} \beta^{\frac{3}{2}} t} \right) + \frac{DL\Gamma}{T\beta}, \quad (86)$$

where  $C = 2H\sqrt{d} + L_F^2 R^{\frac{1}{2}} (1 + 2\rho r')(d+4)^2$ . Since we have  $\sum_{t=1}^T \frac{1}{t} \leq 1 + \log(T)$ , we obtain

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\boldsymbol{\mu}_{t+1}, \Sigma_t) - J(\boldsymbol{\mu}^*, 0)] = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[A_t] = \mathcal{O} \left( \frac{1}{T} + \frac{\log T}{T} \right), \quad (87)$$

where we reach the conclusion.

**Remark D.1** In Theorem 4.3, if we set  $\beta = \mathcal{O}(1)$  and  $\rho = \mathcal{O}(\frac{1}{\sqrt{T}})$ , then the third term in right-hand side of Eq. (85) has a convergence rate of  $\mathcal{O}(\frac{\log T}{T})$ . If we set  $\beta = \mathcal{O}(1)$  and  $\rho = \mathcal{O}(1)$ , then

1080 the third term in right-hand side of Eq. (85) has a convergence rate of  $\mathcal{O}(\frac{1}{\sqrt{T}})$ , since we have  
 1081  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$ . Then we obtain

$$1082 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, 0)] = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[A_t] = \mathcal{O}\left(\frac{1}{T} + \frac{\log T}{T} + \frac{1}{\sqrt{T}}\right). \quad (88)$$

1083 This implies  $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, 0)] = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ .

### 1084 D.3 PROOF OF THE PROPOSITION 4.5

1085 Since the datasets  $\mathcal{D}$  and  $\mathcal{B}$  are i.i.d. sampled from one data distribution  $P(X, y)$ . We have  
 1086  $\mathbb{E}_{(X, y) \sim P}[F(\mathbf{x}; (X, y))] = \mathbb{E}[F(\mathbf{x}; \mathcal{B})] = \mathbb{E}[F(\mathbf{x}; \mathcal{D})]$ . We obtain that

$$1087 \|\mathbb{E}[F(\mathbf{x}; \mathcal{B})] - \mathbb{E}[F(\mathbf{x}; \mathcal{D})]\|_2^2 \leq 2\|F(\mathbf{x}; \mathcal{B}) - \mathbb{E}[F(\mathbf{x}; \mathcal{B})]\|_2^2 + 2\|F(\mathbf{x}; \mathcal{D}) - \mathbb{E}[F(\mathbf{x}; \mathcal{D})]\|_2^2 = 2(\varepsilon_{\mathcal{B}}^2 + \varepsilon_{\mathcal{D}}^2) \quad (89)$$

### 1088 D.4 PROOF OF THEOREM 4.6

1089 Note that for SABO with a mini-batch function query, the update rule of  $\boldsymbol{\mu}$  can be represented as  
 1090  $\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \mathbf{g}_t$ . Let  $B_t = J(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) - J(\boldsymbol{\mu}^*, 0)$ , then by using Eq. (62), we have

$$1091 \beta_t \mathbb{E}[B_t] \leq \beta_t \mathbb{E}[\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] + 2\beta_t \mathbb{E}[\nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\theta}_t)^\top \boldsymbol{\Sigma}_t] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \quad (90)$$

$$1092 \leq \beta_t \mathbb{E}[\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] + 2\beta_t H \|\boldsymbol{\Sigma}_t\|_{\text{F}} - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2, \quad (91)$$

1093 where the second inequality is due to Lemma C.1 and Lipschitz continuous assumption of the function  
 1094  $J(\boldsymbol{\theta})$ . Note that

$$1095 \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2$$

$$1096 = \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - \|\boldsymbol{\mu}_t - \beta_t \boldsymbol{\Sigma}_t \mathbf{g}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 \quad (92)$$

$$1097 = \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - (\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - 2\beta_t \langle \boldsymbol{\mu}_t - \boldsymbol{\mu}^*, \mathbf{g}_t \rangle + \beta_t^2 \langle \boldsymbol{\Sigma}_t \mathbf{g}_t, \mathbf{g}_t \rangle) \quad (93)$$

$$1098 = 2\beta_t \mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) - \beta_t^2 (\boldsymbol{\Sigma}_t \mathbf{g}_t)^\top \mathbf{g}_t. \quad (94)$$

1099 It follows that

$$1100 \mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) = \frac{1}{2\beta_t} (\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2) + \frac{\beta_t}{2} (\boldsymbol{\Sigma}_t \mathbf{g}_t)^\top \mathbf{g}_t \quad (95)$$

$$1101 \leq \frac{1}{2\beta_t} (\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2) + \beta_t \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2, \quad (96)$$

1102 where the inequality is due to  $\beta_t \geq 0$  and Lemma C.1. According to Lemma C.7, we have  $\mathbb{E} \mathbf{g}_t =$   
 1103  $\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t)$ . Therefore

$$1104 \mathbb{E}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)]$$

$$1105 = \mathbb{E}_{\mathbf{z}}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t))^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] + \mathbb{E}_{\mathbf{z}}[(\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \quad (97)$$

$$1106 \leq \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \|\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t + \boldsymbol{\delta}_t)\| \quad (98)$$

$$1107 \leq DL \|\boldsymbol{\delta}_t\| \quad (99)$$

$$1108 \leq DL \rho_t U_t, \quad (100)$$

1109 where the first inequality is due to the Cauchy-Schwarz inequality, the second inequality is due to  
 1110  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$  and smoothness assumption of the function  $J(\boldsymbol{\theta})$ , the last inequality is due to Lemma  
 1111 C.4, and  $U_t = \max(2\sqrt{2} \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}}, 4r \|\boldsymbol{\Sigma}_t\|_{\text{F}})$ . Then we have

$$1112 \mathbb{E}_{\mathbf{z}}[\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \quad (101)$$

$$1113 = \mathbb{E}_{\mathbf{z}}[\mathbf{g}_t^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*) + (\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\theta}_t) - \mathbf{g}_t)^\top (\boldsymbol{\mu}_t - \boldsymbol{\mu}^*)] \quad (102)$$

$$1114 \leq \frac{1}{2\beta_t} \mathbb{E}_{\mathbf{z}}[\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\boldsymbol{\Sigma}_t}^2] + \beta_t \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2] + DL \rho_t U_t, \quad (103)$$

where the inequality is due to Eq. (96) and Eq. (100). Then substituting Eq. (103) into Eq. (91) and multiplying  $\beta_t$  on both sides of the inequality, we have

$$\begin{aligned} \beta_t \mathbb{E}[B_t] &\leq \frac{1}{2} \mathbb{E}[\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\Sigma_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\Sigma_{t+1}^{-1}}^2] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 + 2\beta_t^2 H \|\boldsymbol{\Sigma}_t\|_{\text{F}} \\ &\quad + DL\rho_t U_t + \beta_t^2 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2], \end{aligned} \quad (104)$$

We further obtain that

$$\sum_{t=0}^{T-1} \left[ \frac{1}{2} \mathbb{E}[\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\Sigma_t^{-1}}^2 - \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}^*\|_{\Sigma_{t+1}^{-1}}^2] - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] \quad (105)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\Sigma_t^{-1}}^2 - \|\boldsymbol{\mu}_{t-1} - \boldsymbol{\mu}^*\|_{\Sigma_{t-1}^{-1}}^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] \quad (106)$$

$$+ \frac{1}{2} \left[ \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}^*\|_{\Sigma_0^{-1}}^2 - \|\boldsymbol{\mu}_T - \boldsymbol{\mu}^*\|_{\Sigma_{T-1}^{-1}}^2 \right] \quad (107)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{2\beta_t \mathbf{G}_t}^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] + \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2 \quad (108)$$

$$\leq \frac{1}{2} \sum_{t=0}^{T-1} \left[ \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 - \frac{c\beta_t}{2} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|^2 \right] + \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2 \quad (109)$$

$$= \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2, \quad (110)$$

where the second inequality is due to the update rule of  $\boldsymbol{\Sigma}_t$  and  $\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\| \leq D$ , and the third inequality is due to Cauchy-Schwarz inequality and  $\mathbf{G}_t \preceq \frac{c}{4} \mathbf{I}$ . Then we have

$$\sum_{t=0}^{T-1} \beta_t \mathbb{E}[B_t] \leq \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2 + \sum_{t=0}^{T-1} \left( 2H\beta_t^2 \|\boldsymbol{\Sigma}_t\|_{\text{F}} + DL\rho_t U_t + \beta_t^2 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \mathbb{E}[\|\boldsymbol{\Sigma}_t^{\frac{1}{2}} \mathbf{g}_t\|^2] \right) \quad (111)$$

$$\begin{aligned} &\leq \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2 + \sum_{t=0}^{T-1} \left( \frac{2H\beta_t^2 \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k} + DL\rho_t U_t + \frac{2\beta_t^2 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} (d+4)\varepsilon^2}{3} \right. \\ &\quad \left. + \frac{L_F^2 (1+2\rho r')(d+4)^2 \beta_t^2 \|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}}}{6\xi (\sum_{k=1}^t \beta_k)} \right) \end{aligned} \quad (112)$$

$$\begin{aligned} &\leq \|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2 + \sum_{t=0}^{T-1} \left( \frac{2H\beta_t^2 \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k} + DL\rho_t U_t + \frac{2\beta_t^2 (d+4)d^{\frac{1}{4}}\varepsilon^2}{3(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}} \right. \\ &\quad \left. + \frac{L_F^2 (1+2\rho r')(d+4)^2 \beta_t^2 R^{\frac{1}{2}}}{6\xi (\sum_{k=1}^t \beta_k)} \right), \end{aligned} \quad (113)$$

where the first inequality is due to Eq. (104) and Eq. (110), the second inequality is due to Lemma C.3 (b) and Lemma C.8, and the third inequality is due to  $\|\boldsymbol{\Sigma}_t^{\frac{1}{2}}\|_{\text{F}} \leq R^{\frac{1}{2}}$  and Lemma C.3 (b). According

to C.3 (b), we have  $U_t \leq \max\left(\frac{2\sqrt{2}d^{\frac{1}{4}}}{(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}}, \frac{4r\sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}\right)$ . Therefore, there exists a constant  $t^*$ ,

when  $t > t^* - 1$ ,  $U_t \leq \frac{2\sqrt{2}d^{\frac{1}{4}}}{(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}}$ . Denote  $\sum_{t=0}^{t^*-1} \frac{4r\rho_t \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}$  by a constant  $\Gamma$ . Then if  $T > t^*$ ,

we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[B_t] &\leq \frac{\|\boldsymbol{\Sigma}_0^{-1}\|_{\text{F}} D^2}{T\beta_t} + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{2H\beta_t \sqrt{d}}{2\xi \sum_{k=1}^t \beta_k} + \frac{L_F^2 R^{\frac{1}{2}} (1+2\rho r')(d+4)^2 \beta_t}{6\xi \sum_{k=1}^t \beta_k} \right) \\ &\quad + \frac{1}{T} \sum_{t=t^*}^{T-1} \left( \frac{2\sqrt{2}DL\rho_t d^{\frac{1}{4}}}{\beta_t (2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}} \right) + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{2\beta_t (d+4)d^{\frac{1}{4}}\varepsilon^2}{3(2\xi \sum_{k=1}^t \beta_k)^{\frac{1}{2}}} \right) + \frac{DL\Gamma}{T\beta_t}, \end{aligned} \quad (114)$$

Let  $\beta_t = \beta$  and  $\rho_t = \rho$ , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[B_t] &\leq \frac{\|\Sigma_0^{-1}\|_F D^2}{T\beta} + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{2H\sqrt{d}}{2\xi t} + \frac{L_F^2 R^{\frac{1}{2}} (1 + 2\rho r')(d+4)^2}{6\xi t} \right) \\ &\quad + \frac{1}{T} \sum_{t=t^*}^{T-1} \left( \frac{2\sqrt{2}DL\rho d^{\frac{1}{4}}}{\sqrt{2\xi}\beta^{\frac{3}{2}}\sqrt{t}} \right) + \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{2(d+4)d^{\frac{1}{4}}\varepsilon^2}{3\sqrt{2\xi}\sqrt{t}} \right) + \frac{DL\Gamma}{T\beta}. \end{aligned} \quad (115)$$

Since we have  $\sum_{t=1}^T \frac{1}{t} \leq 1 + \log(T)$  and  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$ , we obtain

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\boldsymbol{\mu}_{t+1}, \Sigma_t) - J(\boldsymbol{\mu}^*, 0)] = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[B_t] = \mathcal{O}\left(\frac{1}{T} + \frac{\log T}{T} + \frac{1}{\sqrt{T}}\right), \quad (116)$$

where we reach the conclusion.

## D.5 PROOF OF THEOREM 4.8

Using PAC-Bayesian bound (McAllester, 1999) and following (Dziugaite & Roy, 2017), for any prior distribution  $p$  and any posterior distribution  $q$  that may be dependent on finite dataset  $\mathcal{S}$ , where data set  $\mathcal{S}$  with  $M$  i.i.d. samples drawn from data distribution  $P(X, y)$ , we have

$$\forall q, \quad \mathbb{E}_{q(\mathbf{x})} [\mathbb{E}_{P(X,y)} [F(\mathbf{x}; (X, y))]] \leq \mathbb{E}_{q(\mathbf{x})} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\text{KL}(q||p) + \log(\frac{M}{\kappa})}{2(M-1)}}. \quad (117)$$

Set the posterior distribution as  $q = p_{\boldsymbol{\theta}} := \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ . Denote the set  $\mathcal{M}(\boldsymbol{\theta}) = \{\boldsymbol{\delta} \mid \text{KL}(p_{\boldsymbol{\theta}+\boldsymbol{\delta}}||p_{\boldsymbol{\theta}}) + \text{KL}(p_{\boldsymbol{\delta}}||p_{\boldsymbol{\theta}+\boldsymbol{\delta}}) \leq \rho^2\}$ . We can choose  $p \in \mathcal{M}(\boldsymbol{\theta})$ . Then, we know that for the prior distribution  $p$ , the following inequality holds with a probability at least  $1 - \kappa$ ,

$$\mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [\mathbb{E}_{P(X,y)} [F(\mathbf{x}; (X, y))]] \leq \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\text{KL}(p_{\boldsymbol{\theta}}||p) + \log(\frac{M}{\kappa})}{2(M-1)}}. \quad (118)$$

Note that for any density  $p, q$ , we have  $\text{KL}(p||q) \geq 0$ . Thus, we know  $\mathcal{M}(\boldsymbol{\theta}) \subset \mathcal{C}(\boldsymbol{\theta})$ . It follows that

$$\mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [\mathbb{E}_{P(X,y)} [F(\mathbf{x}; (X, y))]] \leq \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x}) \in \mathcal{C}(\boldsymbol{\theta})} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\text{KL}(p_{\boldsymbol{\theta}}||p) + \log(\frac{M}{\kappa})}{2(M-1)}} \quad (119)$$

$$\leq \max_{\boldsymbol{\delta} \in \mathcal{C}(\boldsymbol{\theta})} \mathbb{E}_{p_{\boldsymbol{\theta}+\boldsymbol{\delta}}} [F(\mathbf{x}; \mathcal{S})] + \max_{\boldsymbol{\delta} \in \mathcal{M}(\boldsymbol{\theta})} \sqrt{\frac{\text{KL}(p_{\boldsymbol{\theta}}||p_{\boldsymbol{\theta}+\boldsymbol{\delta}}) + \log(\frac{M}{\kappa})}{2(M-1)}} \quad (120)$$

$$\leq \max_{\boldsymbol{\delta} \in \mathcal{C}(\boldsymbol{\theta})} \mathbb{E}_{p_{\boldsymbol{\theta}+\boldsymbol{\delta}}} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\rho^2 + \log(\frac{M}{\kappa})}{2(M-1)}}. \quad (121)$$

Note that  $F(\mathbf{x}; (X, y))$  is convex function w.r.t.  $\mathbf{x}$ , we know that  $\mathbb{E}_{P(X,y)} [F(\mathbf{x}; (X, y))]$  is a convex function w.r.t.  $\mathbf{x}$ . It follows that

$$\mathbb{E}_{P(X,y)} [F(\boldsymbol{\mu}; (X, y))] = \mathbb{E}_{P(X,y)} [F(\mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})}[\mathbf{x}]; (X, y))] \quad (122)$$

$$\leq \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} [\mathbb{E}_{P(X,y)} [F(\mathbf{x}; (X, y))]]. \quad (123)$$

Finally, we know that with a probability of at least  $1 - \kappa$ , the following inequality holds.

$$\mathbb{E}_{P(X,y)} [F(\boldsymbol{\mu}; \mathbf{z}, y)] \leq \max_{\boldsymbol{\delta} \in \mathcal{C}(\boldsymbol{\theta})} \mathbb{E}_{p_{\boldsymbol{\theta}+\boldsymbol{\delta}}} [F(\mathbf{x}; \mathcal{S})] + \sqrt{\frac{\rho^2 + \log(\frac{M}{\kappa})}{2(M-1)}}. \quad (124)$$

## E PROOF OF TECHNICAL LEMMAS

In this section, we provide the proof of lemmas in Appendix C. Note that Lemma C.1 and Lemma C.2 can be directly obtained by Lemma B.1. and Lemma B.2. in (Ye et al., 2024), respectively.

## E.1 PROOF OF LEMMA C.3

(a): Since we have  $\Sigma_{t+1}^{-1} = \Sigma_t^{-1} + 2\beta_t \mathbf{G}_t$ . We can obtain that

$$\Sigma_t^{-1} + 2b\beta_t \mathbf{I} \succeq \Sigma_{t+1}^{-1} \succeq \Sigma_t^{-1} + 2\xi\beta_t \mathbf{I}. \quad (125)$$

Summing up it over  $t = 0, \dots, T-1$ , we have

$$\Sigma_0^{-1} + 2b \sum_{t=1}^T \beta_t \mathbf{I} \succeq \Sigma_T^{-1} \succeq \Sigma_0^{-1} + 2\xi \sum_{t=1}^T \beta_t \mathbf{I}. \quad (126)$$

Therefore, we have

$$\frac{1}{2b \sum_{t=1}^T \beta_t \mathbf{I} + \Sigma_0^{-1}} \preceq \Sigma_T \preceq \frac{1}{2\xi \sum_{t=1}^T \beta_t \mathbf{I} + \Sigma_0^{-1}}. \quad (127)$$

(b): We have

$$\|\Sigma_t\|_{\text{F}} \leq \left\| \frac{1}{2\xi \sum_{k=1}^t \beta_k \mathbf{I} + \Sigma_0^{-1}} \right\|_{\text{F}} \leq \left\| \frac{1}{2\xi \sum_{k=1}^t \beta_k \mathbf{I}} \right\|_{\text{F}} = \frac{\sqrt{d}}{2\xi \sum_{k=1}^t \beta_k}. \quad (128)$$

## E.2 PROOF OF LEMMA C.4

In Algorithm 1 and Algorithm 2, for given  $\theta_t$ , the perturbation  $\delta_t$  satisfies

$$\delta_{\mu_t} = \frac{1}{\lambda} \Sigma_t \mathbf{g}'_t = \frac{\rho \Sigma_t \mathbf{g}'_t}{\sqrt{\|\Sigma_t \mathbf{G}'_t\|_{\text{F}}^2 + 0.5 \|\Sigma_t^{\frac{1}{2}} \mathbf{g}'_t\|_2^2}} \leq \frac{\rho \sqrt{2} \Sigma_t^{\frac{1}{2}} \Sigma_t^{\frac{1}{2}} \mathbf{g}'_t}{\sqrt{\|\Sigma_t^{\frac{1}{2}} \mathbf{g}'_t\|_2^2}}. \quad (129)$$

Therefore, we have  $\|\delta_{\mu_t}\| \leq \rho \sqrt{2} \|\Sigma_t^{\frac{1}{2}}\|_{\text{F}}$ . For  $\delta_{\Sigma_t}$ , we have

$$\delta_{\Sigma_t} = \frac{2\Sigma_t \mathbf{G}'_t}{\lambda \Sigma_t^{-1} - 2\mathbf{G}'_t} = \Sigma_t \frac{\frac{2}{\lambda} \Sigma_t \mathbf{G}'_t}{\mathbf{I} - \frac{2}{\lambda} \Sigma_t \mathbf{G}'_t}. \quad (130)$$

Note that  $\frac{2}{\lambda} \leq \frac{2\rho}{\|\Sigma_t \mathbf{G}'_t\|_{\text{F}}}$ . Therefore, we can obtain that

$$\|\delta_{\Sigma_t}\|_{\text{F}} \leq \|\Sigma_t\|_{\text{F}} \frac{2\rho}{\|\mathbf{I} - \frac{2}{\lambda} \Sigma_t \mathbf{G}'_t\|_{\text{F}}}. \quad (131)$$

If  $\rho < \frac{\sqrt{d}}{2}$ , there exist a constant  $r > 0$ ,  $\|\mathbf{I} - \frac{2}{\lambda} \Sigma_t \mathbf{G}'_t\|_{\text{F}} > \frac{1}{r}$  holds. Therefore we have  $\|\delta_{\Sigma_t}\|_{\text{F}} \leq 2\rho r \|\Sigma_t\|_{\text{F}}$ .

## E.3 PROOF OF LEMMA C.5

We have

$$\mathbb{E}_{\mathbf{z}}[\mathbf{g}'_t] = \mathbb{E}_{\mathbf{z}}[\Sigma_t^{-\frac{1}{2}} \mathbf{z} F(\boldsymbol{\mu}_t + \Sigma_t^{\frac{1}{2}} \mathbf{z})] - \mathbb{E}_{\mathbf{z}}[\Sigma_t^{-\frac{1}{2}} \mathbf{z} F(\boldsymbol{\mu}_t)] \quad (132)$$

$$= \mathbb{E}_{\mathbf{z}}[\Sigma_t^{-\frac{1}{2}} \mathbf{z} F(\boldsymbol{\mu}_t + \Sigma_t^{\frac{1}{2}} \mathbf{z})] \quad (133)$$

$$= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)}[\Sigma_t^{-1} (\mathbf{x} - \boldsymbol{\mu}_t) F(\mathbf{x})] \quad (134)$$

$$= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{p_{\theta_t}}[F(\mathbf{x})], \quad (135)$$

where we reach the conclusion.

## E.4 PROOF OF LEMMA C.6

Denote the diagonal elements of  $\Sigma$  and  $\hat{\Sigma}$  by  $\boldsymbol{\sigma}$  and  $\hat{\boldsymbol{\sigma}}$ , respectively. Then we have

$$\|\Sigma_t^{\frac{1}{2}} \mathbf{g}_{t,j}\|^2 = \|\boldsymbol{\sigma}_t^{\frac{1}{2}} \odot \hat{\boldsymbol{\sigma}}_t^{-\frac{1}{2}} \odot \mathbf{z}_j (F(\hat{\boldsymbol{\mu}}_t + \hat{\boldsymbol{\sigma}}_t^{\frac{1}{2}} \odot \mathbf{z}_j) - F(\hat{\boldsymbol{\mu}}_t))\|^2. \quad (136)$$



Note that  $\widehat{\Sigma}_{t+1}^{-1} = \Sigma_t^{-1} - 2\lambda G'_t$ , we have  $\sigma_t \leq \widehat{\sigma}_t$ . Since we have

$$\mathbb{E}\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_t\|_2^2 = \mathbb{E}\left\|\frac{1}{N}\sum_{j=1}^N\Sigma_t^{\frac{1}{2}}\mathbf{g}_{t,j}\right\|^2 \leq \frac{1}{N}\sum_{j=1}^N\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_{t,j}\|^2. \quad (137)$$

It follows that

$$\mathbb{E}\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_t\|_2^2 \leq \frac{1}{N}\sum_{j=1}^N\mathbb{E}\|\widehat{\sigma}_t^{\frac{1}{2}}\odot\widehat{\sigma}_t^{-\frac{1}{2}}\odot\mathbf{z}_j(F(\widehat{\boldsymbol{\mu}}_t+\widehat{\sigma}_t^{\frac{1}{2}}\odot\mathbf{z}_j)-F(\widehat{\boldsymbol{\mu}}_t))\|^2 \quad (138)$$

$$\leq \frac{1}{N}\sum_{j=1}^N\mathbb{E}\left[\|\mathbf{z}_j\|^2(F(\widehat{\boldsymbol{\mu}}_t+\widehat{\sigma}_t^{\frac{1}{2}}\odot\mathbf{z}_j)-F(\widehat{\boldsymbol{\mu}}_t))^2\right] \quad (139)$$

$$\leq \frac{1}{N}\sum_{j=1}^N\mathbb{E}\left[L_F^2\|\mathbf{z}_j\|^4\|\widehat{\sigma}_t\|_\infty\right]. \quad (140)$$

Note that

$$\widehat{\Sigma}_t = \Sigma_t + \delta_{\Sigma_t} = \Sigma_t + \Sigma_t \frac{\frac{2}{\lambda}\Sigma_t G'_t}{\mathbf{I} - \frac{2}{\lambda}\Sigma_t G'_t}. \quad (141)$$

Note that  $\frac{2}{\lambda} \leq \frac{2\rho}{\|\Sigma_t G'_t\|_F}$ . Then if  $\rho < \frac{\sqrt{d}}{2}$ , there exist a constant  $r' > 0$ , the inequality  $\widehat{\sigma}_t \leq \sigma_t + 2\rho r' \sigma_t$  holds. Therefore, we have

$$\|\widehat{\sigma}_t\|_\infty \leq (1 + 2\rho r')\|\sigma_t\|_\infty \leq \frac{1 + 2\rho r'}{\|\sigma_0^{-1}\|_{min} + 2(\sum_{k=1}^t \beta_k)\xi}, \quad (142)$$

where  $\|\cdot\|_{min}$  denotes the minimum element in the input. Noticed that  $\mathbb{E}_z[\|\mathbf{z}\|^2] \leq d + 4$  as shown in (Nesterov & Spokoiny, 2017), we obtain

$$\mathbb{E}\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_t\|_2^2 \leq \frac{L_F^2(1 + 2\rho r')(d + 4)^2}{2\xi(\sum_{k=1}^t \beta_k)}, \quad (143)$$

where we reach the conclusion.

## E.5 PROOF OF LEMMA C.7

We have

$$\mathbb{E}[\mathbf{g}'_t] = \mathbb{E}[\Sigma_t^{-\frac{1}{2}}\mathbf{z}F(\boldsymbol{\mu}_t + \Sigma_t^{\frac{1}{2}}\mathbf{z}; \mathcal{B})] - \mathbb{E}[\Sigma_t^{-\frac{1}{2}}\mathbf{z}F(\boldsymbol{\mu}_t; \mathcal{B})] \quad (144)$$

$$= \mathbb{E}_z[\Sigma_t^{-\frac{1}{2}}\mathbf{z}F(\boldsymbol{\mu}_t + \Sigma_t^{\frac{1}{2}}\mathbf{z}; \mathcal{D})] - \mathbb{E}_z[\Sigma_t^{-\frac{1}{2}}\mathbf{z}F(\boldsymbol{\mu}_t; \mathcal{D})] \quad (145)$$

$$= \mathbb{E}_z[\Sigma_t^{-\frac{1}{2}}\mathbf{z}F(\boldsymbol{\mu}_t + \Sigma_t^{\frac{1}{2}}\mathbf{z})] \quad (146)$$

$$= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)}[\Sigma_t^{-1}(\mathbf{x} - \boldsymbol{\mu}_t)F(\mathbf{x})] \quad (147)$$

$$= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{p_{\theta_t}}[F(\mathbf{x})], \quad (148)$$

where the second equality is due to Proposition 4.5.

## E.6 PROOF OF LEMMA C.8

Denote the diagonal elements of  $\Sigma$  and  $\widehat{\Sigma}$  by  $\sigma$  and  $\widehat{\sigma}$ , respectively. Then we have

$$\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_{t,j}\|^2 = \|\sigma_t^{\frac{1}{2}}\odot\widehat{\sigma}_t^{-\frac{1}{2}}\odot\mathbf{z}_j(F(\widehat{\boldsymbol{\mu}}_t+\widehat{\sigma}_t^{\frac{1}{2}}\odot\mathbf{z}_j; \mathcal{B})-F(\widehat{\boldsymbol{\mu}}_t; \mathcal{B}))\|^2 \quad (149)$$

Note that  $\widehat{\Sigma}_{t+1}^{-1} = \Sigma_t^{-1} - 2\lambda G'_t$ , we have  $\sigma_t \leq \widehat{\sigma}_t$ . Since we have

$$\mathbb{E}\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_t\|_2^2 = \mathbb{E}\left\|\frac{1}{N}\sum_{j=1}^N\Sigma_t^{\frac{1}{2}}\mathbf{g}_{t,j}\right\|^2 \leq \frac{1}{N}\sum_{j=1}^N\|\Sigma_t^{\frac{1}{2}}\mathbf{g}_{t,j}\|^2. \quad (150)$$

1350 It follows that

$$1351 \mathbb{E} \|\Sigma_t^{\frac{1}{2}} \mathbf{g}_t\|_2^2 \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E} \|\hat{\sigma}_t^{\frac{1}{2}} \odot \hat{\sigma}_t^{-\frac{1}{2}} \odot \mathbf{z}_j (F(\hat{\boldsymbol{\mu}}_t + \hat{\sigma}_t^{\frac{1}{2}} \odot \mathbf{z}_j; \mathcal{B}) - F(\hat{\boldsymbol{\mu}}_t; \mathcal{B}))\|^2 \quad (151)$$

$$1352 \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E} \left[ \|\mathbf{z}_j\|^2 (F(\hat{\boldsymbol{\mu}}_t + \hat{\sigma}_t^{\frac{1}{2}} \odot \mathbf{z}_j; \mathcal{B}) - F(\hat{\boldsymbol{\mu}}_t; \mathcal{B}))^2 \right] \quad (152)$$

$$1353 \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E} \left[ \|\mathbf{z}_j\|^2 \left( \frac{1}{3} L_F^2 \|\hat{\sigma}_t^{\frac{1}{2}} \odot \mathbf{z}_j\|^2 + \frac{2}{3} \varepsilon^2 \right) \right], \quad (153)$$

1354 where the third inequality is due to Proposition 4.5. Since  $\hat{\Sigma}_t = \Sigma_t + \Sigma_t \frac{2}{\lambda} \frac{\Sigma_t \mathbf{G}'_t}{I - \frac{2}{\lambda} \Sigma_t \mathbf{G}'_t}$ , and  $\frac{2}{\lambda} \leq \frac{2\rho}{\|\Sigma_t \mathbf{G}'_t\|_F}$ .

1355 Then if  $\rho < \frac{\sqrt{d}}{2}$ , there exist a constant  $r' > 0$ , the inequality  $\hat{\sigma}_t \leq \sigma_t + 2\rho r' \sigma_t$  holds. Therefore, we

$$1356 \text{ have} \quad \|\hat{\sigma}_t\|_\infty \leq (1 + 2\rho r') \|\sigma_t\|_\infty \leq \frac{1 + 2\rho r'}{\|\sigma_0^{-1}\|_{\min} + 2(\sum_{k=1}^t \beta_k) \xi}, \quad (154)$$

1357 where  $\|\cdot\|_{\min}$  denotes the minimum element in the input. Noticed that  $\mathbb{E}_z[\|z\|^2] \leq d + 4$  as shown

$$1358 \text{ in (Nesterov \& Spokoiny, 2017), we obtain} \quad \mathbb{E} \|\Sigma_t^{\frac{1}{2}} \mathbf{g}_t\|_2^2 \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E} \left[ \frac{(1 + 2\rho r') L_F^2}{3} \|\sigma_t\|_\infty \times \|\mathbf{z}_j\|^4 + \frac{2\varepsilon^2}{3} \|\mathbf{z}_j\|^2 \right] \quad (155)$$

$$1359 \leq \frac{L_F^2 (1 + 2\rho r') (d + 4)^2}{6\xi (\sum_{k=1}^t \beta_k)} + \frac{2(d + 4)\varepsilon^2}{3}, \quad (156)$$

1360 where we reach the conclusion.

## 1361 F ADDITIONAL MATERIALS FOR SECTION 6

1362 In this section, we provide additional experiments and more implementation details of Section 6.

### 1363 F.1 SYNTHETIC PROBLEMS

1364 The four numerical benchmark test functions employed in Section 6.1 are listed as follows:

$$1365 F(\mathbf{x}) = \sum_{i=1}^d 10^{\frac{2(i-1)}{d-1}} \mathbf{x}_i^2, \quad (157)$$

$$1366 F(\mathbf{x}) = \sum_{i=1}^d 10^{\frac{2(i-1)}{d-1}} |\mathbf{x}_i|^{\frac{1}{2}}, \quad (158)$$

$$1367 F(\mathbf{x}) = \sqrt{\sum_{i=1}^d |\mathbf{x}_i|^{2+4\frac{i-1}{d-1}}}, \quad (159)$$

$$1368 F(\mathbf{x}) = \sin^2(\pi\omega_1) + \sum_{i=1}^{d-1} (\omega_i - 1)^2 (1 + 10 \sin^2(\omega_i \pi + 1)) + (\omega_d - 1)^2 (1 + \sin^2(2\pi\omega_d)), \quad (160)$$

$$1369 \text{ where } \omega_i = 1 + \frac{x_i - 1}{4}, i \in \{1, \dots, d\}.$$

1370 Test functions (157)-(160) are called the ellipsoid function,  $l_{\frac{1}{2}}$ -ellipsoid function, different powers

1371 function and Levy function, respectively.

1372 **Implementation Details.** For all the methods, we initialize  $\boldsymbol{\mu}_0$  from the uniform distribution

1373 Uni[0, 1], and set  $\Sigma_0 = \mathbf{I}$ . For the INGO, BES, and SABO methods, we use a fixed step size of

1374  $\beta = 0.1$ . According to our assumption in Theorem 4.3, we set  $\rho = 100/\sqrt{T+1}$  for the proposed

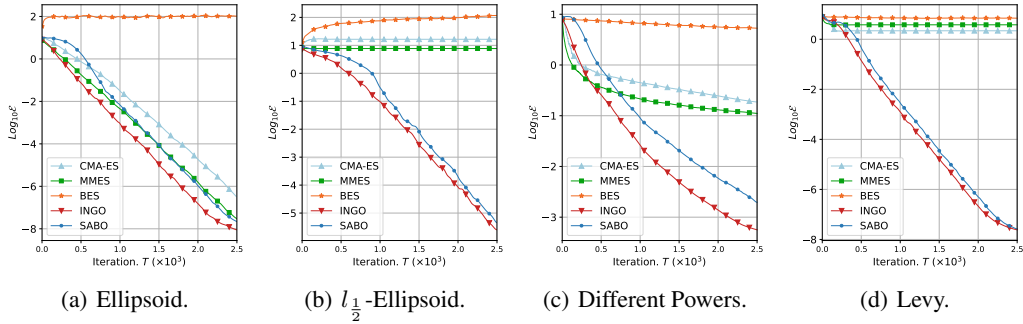
1375 SABO method. We set the spacing  $c = 1$  for the BES method and employ the default hyperparameter

1376 setting from He et al. (2020) for the MMES method. We then assess these methods using varying

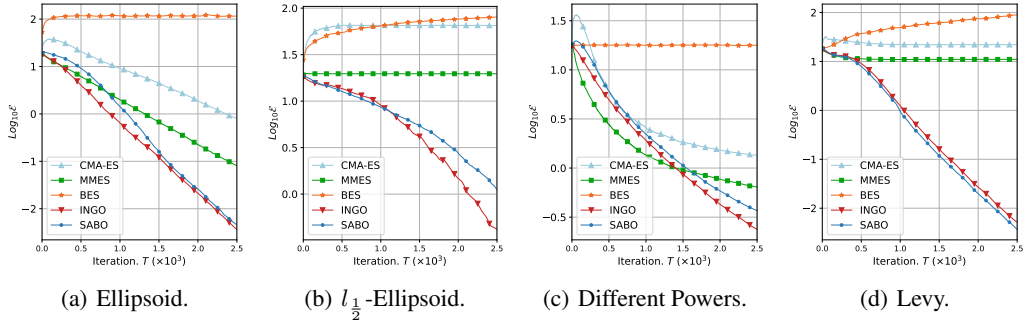
1377 dimensions, i.e.,  $d \in \{200, 500, 1000\}$ . For  $d = 200$ , we assess these methods using varying sample

1378 sizes, i.e.,  $N \in \{10, 50, 100\}$ . The mean value of  $\mathcal{E}$  over 3 independent runs is reported.

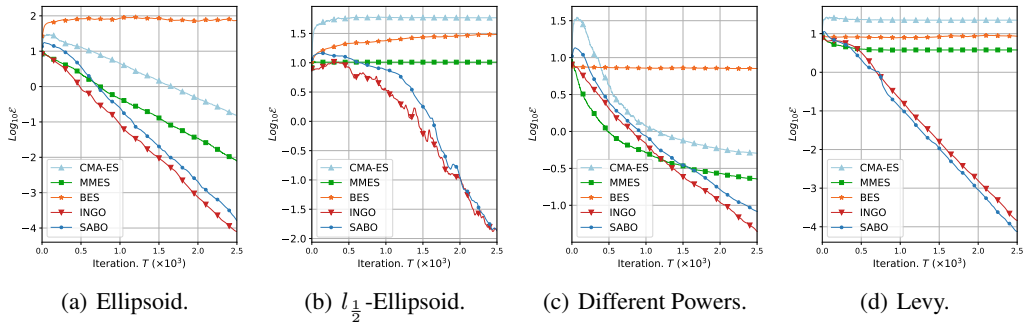
1404 **Results.** Figure 2 and 3 show the results on four test functions with problem dimensions  $d = 200$   
 1405 and  $d = 1000$ , respectively. Figure 4 and 5 show the results on 200-dimensional test functions  
 1406 with sample size  $N = 10$  and  $N = 100$ , respectively. Combining these results with the result from  
 1407 Figure 1, we observe consistent performance for the proposed ASBO method. It achieves a similar  
 1408 convergence result to the INGO method, as they have the same theoretical convergent rate. In some  
 1409 cases, i.e., Figure 3 (d) and Figure 4 (d), it converges slightly faster than INGO. The CMA-ES method  
 1410 and MMES method both work for ellipsoid and different powers functions, but fail in  $l_{\frac{1}{2}}$ -ellipsoid and  
 1411 Levy functions. In most cases, they converge slower than INGO and SABO. With a large sample size,  
 1412 i.e.,  $N = 100$ , the CMA-ES method can maintain a fast converge rate according to Figure 4 (a). The  
 1413 BES method fails to achieve high precision in all test functions. It diverges in ellipsoid,  $l_{\frac{1}{2}}$ -ellipsoid,  
 1414 Levy functions, and only achieves a low precision in the different power functions. These results  
 1415 demonstrate the superiority of the SABO method in optimizing high-dimensional problems, and  
 1416 verify our theoretical convergence results in Section 4.



1429 Figure 2: Results on the four test functions with problem dimension  $d = 200$  and  $N = 50$ .



1443 Figure 3: Results on the four test functions with problem dimension  $d = 1000$  and  $N = 50$ .



1457 Figure 4: Results on the four test functions with problem dimension  $d = 200$  and  $N = 10$ .

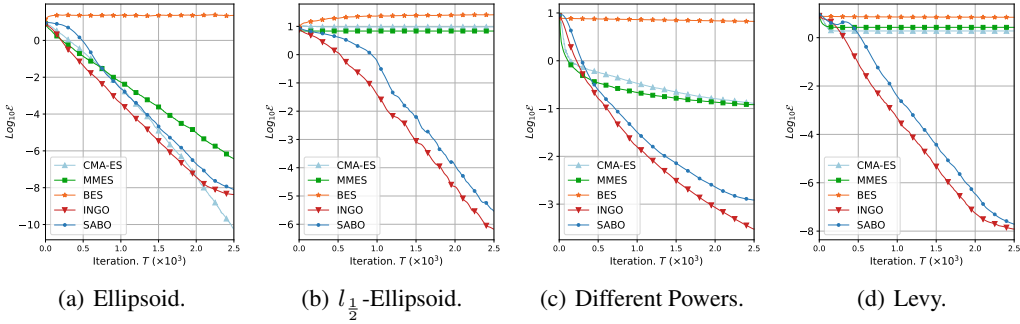


Figure 5: Results on the four test functions with problem dimension  $d = 200$  and  $N = 100$ .

Table 2: Performance (%) on *SST-2*, *AG’s News*, *MRPC*, *RTE*, *SNLI*, and *Yelp P*. datasets. We report the mean and standard deviation over 3 random seeds. The best result across all groups is highlighted in **bold** and the best result in each group is marked with underlined.

| Methods                   | <i>SST-2</i>            | <i>AG’s News</i>        | <i>MRPC</i>             | <i>RTE</i>              | <i>SNLI</i>             | <i>Yelp P</i> .         |
|---------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| Zero-shot                 | 79.82                   | 76.96                   | 67.40                   | 51.62                   | 38.82                   | 89.64                   |
| <i>Dimension d = 200</i>  |                         |                         |                         |                         |                         |                         |
| GIBO                      | 83.53 $\pm$ 0.15        | 75.79 $\pm$ 0.08        | 79.21 $\pm$ 0.09        | 53.07 $\pm$ 0.29        | 38.73 $\pm$ 0.09        | 89.63 $\pm$ 0.03        |
| TurBo                     | 83.30 $\pm$ 0.19        | 79.01 $\pm$ 1.68        | 69.59 $\pm$ 8.51        | 46.57 $\pm$ 2.30        | 40.27 $\pm$ 0.69        | 90.16 $\pm$ 0.19        |
| <b>SABO</b>               | <u>87.88</u> $\pm$ 0.53 | <u>82.22</u> $\pm$ 0.41 | <u>79.35</u> $\pm$ 0.12 | <b>53.67</b> $\pm$ 0.17 | <u>40.72</u> $\pm$ 0.15 | <u>91.50</u> $\pm$ 0.13 |
| <i>Dimension d = 500</i>  |                         |                         |                         |                         |                         |                         |
| GIBO                      | 83.49 $\pm$ 0.09        | 75.70 $\pm$ 0.05        | 79.03 $\pm$ 0.08        | 52.95 $\pm$ 0.17        | 38.71 $\pm$ 0.16        | 89.65 $\pm$ 0.02        |
| TurBo                     | 84.52 $\pm$ 0.65        | 80.03 $\pm$ 1.97        | 75.30 $\pm$ 2.34        | 48.01 $\pm$ 0.59        | 38.82 $\pm$ 0.34        | 90.20 $\pm$ 0.45        |
| <b>SABO</b>               | <u>87.31</u> $\pm$ 0.38 | <u>82.65</u> $\pm$ 0.59 | <u>79.62</u> $\pm$ 0.07 | <u>53.55</u> $\pm$ 0.17 | <b>42.29</b> $\pm$ 2.48 | <u>91.83</u> $\pm$ 0.16 |
| <i>Dimension d = 1000</i> |                         |                         |                         |                         |                         |                         |
| GIBO                      | 83.45 $\pm$ 0.11        | 75.67 $\pm$ 0.10        | 79.15 $\pm$ 0.0         | 52.95 $\pm$ 0.17        | 38.87 $\pm$ 0.18        | 89.65 $\pm$ 0.04        |
| TurBo                     | 85.90 $\pm$ 0.95        | 82.36 $\pm$ 0.21        | 77.30 $\pm$ 0.86        | 50.30 $\pm$ 1.11        | 39.87 $\pm$ 1.07        | 90.14 $\pm$ 0.20        |
| <b>SABO</b>               | <b>87.96</b> $\pm$ 0.83 | <u>82.77</u> $\pm$ 0.41 | <u>79.68</u> $\pm$ 0.23 | <u>53.31</u> $\pm$ 0.17 | <u>40.32</u> $\pm$ 0.27 | <b>91.96</b> $\pm$ 0.41 |

## F.2 ADDITIONAL RESULTS ON BLACK-BOX PROMPT FINE-TUNING

We conduct an additional experiment on the black-box prompt fine-tuning task to compare the proposed SABO method with high-dimensional BO methods, i.e., TuRBO (Eriksson et al., 2019) and GIBO (Nguyen et al., 2022) methods discussed in Section 5. We employ the default setting of TuRBO-1 from He et al. (2020) for the TuRBO method, and the default setting from Nguyen et al. (2022) for the GIBO method. The results on six benchmark datasets are reported in Table 2. According to the results, the SABO method consistently outperforms these two baselines, highlighting its effectiveness.

## F.3 SYNTHETIC IMAGE CLASSIFICATION PROBLEM

We conduct experiments on a synthetic image classification problem to empirically evaluate the performance of the proposed mini-batch SABO method. Specifically, we apply black-box optimization methods to train a model to classify the images accurately. Moreover, following Foret et al. (2021), we conduct additional experiments on the noisy label setting. Particularly, we train the model on a corrupted version of the *Fashion-MNIST* dataset, where some of its training labels are randomly flipped according to different noise rates, while the testing set is clean. To construct problems with different dimensions, we first adopt UMAP (McInnes et al., 2018) to reduce the dimension of *Fashion-MNIST* to 8 while preserving the class discriminability (Sagawa & Hino, 2022), and then employ a fixed randomly initialized matrix  $P \in \mathbb{R}^{8 \times \tilde{d}}$  to project the extracted features to a

Table 3: Performance (%) on *Fashion-MNIST* dataset with noise labels. The best result across all groups is highlighted in **bold** and the best result in each group is marked with underlined.

| Methods                   | noise=0%           | noise=20%          | noise=40%          | noise=60%          | noise=80%          |
|---------------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| <i>Dimension d = 100</i>  |                    |                    |                    |                    |                    |
| CMA-ES                    | 76.68±0.44         | 77.25±1.26         | 75.71±1.15         | 70.81±1.01         | 57.86±3.65         |
| MMES                      | 76.20±1.45         | 74.68±0.33         | 74.46±0.32         | 71.12±2.24         | 65.92±1.30         |
| BES                       | 76.26±0.21         | 73.60±0.50         | 67.60±0.72         | 55.14±6.26         | 45.65±4.35         |
| INGO                      | 76.68±0.90         | 75.16±0.55         | 71.26±0.74         | 61.70±3.89         | 46.75±1.35         |
| SABO                      | <b>78.10</b> ±0.90 | <u>77.41</u> ±0.68 | <b>77.10</b> ±1.69 | <u>73.41</u> ±0.33 | <b>66.45</b> ±0.43 |
| <i>Dimension d = 1000</i> |                    |                    |                    |                    |                    |
| CMA-ES                    | 76.75±0.54         | 74.85±0.06         | 73.21±0.52         | 67.32±0.84         | 50.00±0.98         |
| MMES                      | 75.85±0.43         | 72.39±3.59         | 73.73±0.50         | 72.51±1.47         | 56.60±2.72         |
| BES                       | 76.18±0.63         | 72.45±0.44         | 66.64±2.09         | 58.74±1.50         | 46.99±4.30         |
| INGO                      | 74.73±0.62         | 74.87±1.72         | 69.75±0.73         | 63.18±4.04         | 41.80±1.36         |
| SABO                      | <u>76.87</u> ±1.41 | <b>77.73</b> ±0.85 | <u>74.97</u> ±1.12 | <b>73.81</b> ±1.26 | <u>61.96</u> ±0.08 |

$\tilde{d}$ -dimensional space. After preprocessing, a linear layer is used as a classifier. Therefore, the total number of the trainable parameters is  $d = 10 \times \tilde{d}$ .

**Datasets.** *Fashion-MNIST* (Xiao et al., 2017) is a image classifications dataset. It contains 60,000 training samples and 10,000 test samples, each representing a  $28 \times 28$ -pixel grayscale image of fashion items from 10 different classes.

**Implementation Details.** For a fair comparison, we set the same population size, number of batch samples, and the initialization for CMA-ES (Hansen, 2006), INGO (Lyu & Tsang, 2021), BES (Gao & Sener, 2022), and SABO. The population size and the number of batch samples are set to  $N = 100$  and  $M = 2048$ , respectively. The Gaussian distributions are initialized with  $\mu_0 = \mathbf{0}$  and  $\Sigma_0 = 0.5\mathbf{I}$ . For BES, INGO and SABO methods, we search the learning rate  $\beta$  over  $\{0.1, 0.5, 1, 5\}$ . Moreover, we perform grid-search on  $\rho$  over  $\{100, 500, 1000, 5000\}$  for SABO. We employ the default hyperparameter setting from He et al. (2020) for the MMES method. For the BES method, we perform grid-search on the spacing  $c$  over  $\{0.1, 1, 10\}$ . The cross-entropy loss is used as the training objective. All experiments are repeatedly run with three independent seeds and the mean and standard deviation are reported.

**Results.** Table 3 shows experimental results on *Fashion-MNIST* dataset with different  $\tilde{d}$  and different noise rates. We can see that the SABO method consistently outperforms all baselines across different noise rates and dimensions, highlighting its effectiveness. These results show that the proposed SABO method can achieve good robustness performance and demonstrate the effectiveness of the proposed SABO method in improving model generalization performance.

## G DISCUSSION WITH OFFLINE BLACK-BOX OPTIMIZATION

The typical black-box optimization problem we studied in this work can also be called the online black-box optimization problem. Since we have access to the objective function  $F$ , the problem can be solved in an **online** iterative manner, where in each iteration the solver proposes new  $\mathbf{x}$  and queries the objective function for feedback in order to inform better solution proposals at the next iteration.

The offline black-box optimization (Chen et al., 2022; Qi et al., 2022) is different from the research line of standard online black-box optimization. In offline black-box optimization, access to the true objective  $F$  is not available. Instead, the offline black-box optimization algorithm is provided access to a static dataset  $\mathcal{D} = \{\mathbf{x}_i, F(\mathbf{x}_i)\}$  of the variable  $\mathbf{x}_i$  and its corresponding objective value  $F(\mathbf{x}_i)$ . Therefore, the basic settings of offline black-box optimization and online black-box optimization are different.

1566 Moreover, the challenges of offline black-box optimization and online black-box optimization are  
 1567 also different. Offline black-box optimization focuses on producing query candidates by a surrogate  
 1568 model trained with a prior static dataset (Trabucco et al., 2021; Kumar & Levine, 2020). The goal is  
 1569 to produce a good query set based on the fixed model at one time without considering query feedback  
 1570 update, exploration and exploitation balance in the long term. Along this line, Kumar & Levine  
 1571 (2020) proposed a model-based offline optimization by training a conditional generative model that  
 1572 conditions the objective value. In Fannjiang & Listgarten (2020), the authors formulated the problem  
 1573 as a non-zero-sum game and proposed an alternating ascent-descent algorithm for model-based offline  
 1574 optimization. Trabucco et al. (2021) proposed the conservative objective models, which presents a  
 1575 technique similar to adversarial training that avoids overestimation of out-of-distribution inputs. In  
 1576 contrast to offline black-box optimization, standard online black-box optimization needs to balance  
 1577 exploration and exploitation, which focuses on long-term convergence performance. As a result,  
 1578 offline black-box optimization is not suitable for our online prompt fine-tuning tasks.

1579 **H DISCUSSIONS ABOUT THE SOCIETAL IMPACT AND LIMITATIONS**

1580  
 1581 This work only focuses on black-box optimization in deep learning, so it has no negative societal  
 1582 impact. The main theoretical analysis in this work focuses on convex black-box functions. It is  
 1583 technically challenging to analyze non-convex cases considering both the black-box nature and the  
 1584 sharpness-aware properties, and we leave this as one of our future work. Additionally, the SABO  
 1585 method employs a standard Monte Carlo sampling for gradient approximation. Other sampling  
 1586 techniques might be more efficient, but those are out of the scope of this work. We will study it in the  
 1587 future.  
 1588

1589 **I ADDITIONAL RESULTS ON SYNTHETIC PROBLEMS**

1590  
 1591 We conduct additional experiments to compare the proposed SABO method and GFM method Lin  
 1592 et al. (2022). For GFM, we employ its default hyperparameter setting from Lin et al. (2022). We set its  
 1593 smoothing parameter  $\delta = 0.2$  and conduct experiments on the step-size over  $\{0.001, 0.0005, 0.0001\}$ .  
 1594 The initial point is set the same as SABO.  
 1595

1596 The results are shown in Figure 6 with problem dimension  $d = 200$  and  $N = 50$ . The GFM method  
 1597 can converge slowly for the Ellipsoid problem. It fails in the  $l_{\frac{1}{2}}$ -Ellipsoid problem and cannot achieve  
 1598 high precision for Different Powers and Levy problem. This shows that it could be challenging for  
 1599 GFM to optimize non-smooth or high-dimensional test functions without adaptively updating mean  
 1600 and covariance. The proposed SABO method takes advantage of the second-order information, which  
 1601 gains great acceleration for solving these problems.  
 1602

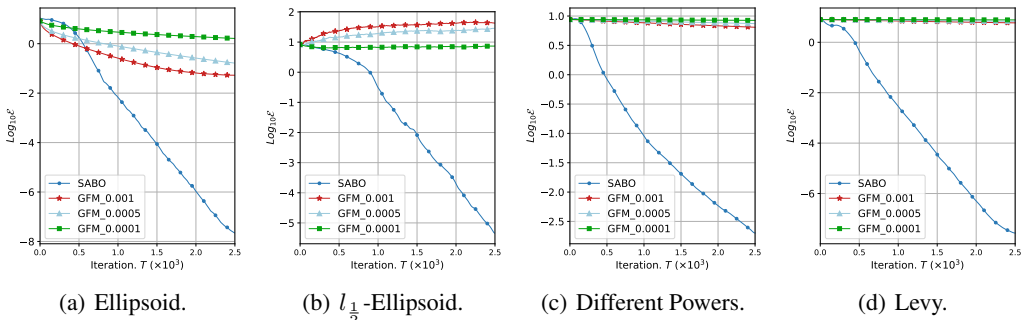


Figure 6: Results on the four test functions with problem dimension  $d = 200$  and  $N = 50$ .