

# Can We Do Better Than Random Start? The Power of Data Outsourcing

YI CHEN, Hong Kong University of Science and Technology, China

JING DONG, Columbia University, USA

XIN T. TONG, National University of Singapore, Singapore

Many organizations have access to abundant data but lack the computational power to process the data. While they can outsource the computational task to other facilities, there are various constraints on the amount of data that can be shared. It is natural to ask what can data outsourcing accomplish under such constraints. We address this question from a machine learning perspective. When training a model with optimization algorithms, the quality of the results often relies heavily on the points where the algorithms are initialized. Random start is one of the most popular methods to tackle this issue, but it can be computationally expensive and not feasible for organizations lacking computing resources. Based on three different scenarios, we propose simulation-based algorithms that can utilize a small amount of outsourced data to find good initial points accordingly. Under suitable regularity conditions, we provide theoretical guarantees showing the algorithms can find good initial points with high probability. We also conduct numerical experiments to demonstrate that our algorithms perform significantly better than the random start approach.

Additional Key Words and Phrases: Non-convex optimization, initialization,

## ACM Reference Format:

Yi Chen, Jing Dong, and Xin T. Tong. 2022. Can We Do Better Than Random Start? The Power of Data Outsourcing. 1, 1 (May 2022), 22 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

In this era, data is the new gold. Organizations of different sizes and sectors all realize the value of collecting data. However, it often requires substantial computational power to turn these data into valuable predictive models and not all organizations have such computational resources. One possible solution to this problem is outsourcing the data processing task to another computing facility, where computational power is substantially cheaper. However, the data organization may only be willing to share a small part of their data due to the following reasons: First, if the computing facility has access to all the available data, it can obtain an accurate predictive model which leads to potential competition risk. Second, some parts of the data may not be share-able due to privacy concerns. Third, transferring data can be expensive especially when certain encryption is required.

Given the constraint that only part of the data is “share-able”, the organization with data can only expect sub-optimal results from the computing facility, and additional learning are needed to improve these premature results. Since the data organization is assumed to have limited computational power, it is desirable if the computational cost of the additional learning can be minimized. In this context, we are interested in investigating the following two questions: 1) What type of

---

Authors’ addresses: Yi Chen, [yichen@ust.hk](mailto:yichen@ust.hk), Hong Kong University of Science and Technology, Hong Kong, China; Jing Dong, [jing.dong@gsb.columbia.edu](mailto:jing.dong@gsb.columbia.edu), Columbia University, New York, NY, USA; Xin T. Tong, [mattxin@nus.edu.sg](mailto:mattxin@nus.edu.sg), National University of Singapore, Singapore.

---

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

© 2022 Association for Computing Machinery.

XXXX-XXXX/2022/5-ART \$15.00

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

computational task should be assigned to the computing facility? 2) How much data should be outsourced? In this paper, we address these two questions from the perspective of machine learning.

Most machine learning models are trained using the risk minimization approach. That is, the unknown parameter  $\theta$  is inferred by minimizing a loss function of the form  $F(\theta) = \mathbb{E}[f(\theta, X)]$  where  $X$  is averaged over a population distribution or empirical distribution of  $N$  data points, and  $f(\theta, x)$  is the loss of using the model with parameter  $\theta$  to explain the data point  $x$ . Greedy local optimization algorithms are often applied to minimize  $F$ . If  $F$  is strongly convex, the computational cost of an algorithm  $\mathcal{T}$ ,  $c(\mathcal{T})$ , depends on the accuracy requirement  $\epsilon$  and/or the number of data points  $N$ . In this setting,  $c(\mathcal{T})$  can be large but is computationally manageable since  $\mathcal{T}$  converges to the optimal parameter regardless of the initialization [17, 29]. However, if  $F$  is non-convex, the quality of the parameter learned from  $\mathcal{T}$  can depend heavily on its initialization  $\theta^0$ . In general, greedy algorithms converge to local minimums that are close to  $\theta^0$ . Thus, in order to find the global minimum  $\theta^*$ , one needs to start  $\mathcal{T}$  in an appropriate attraction region of optimal parameter  $\theta^*$ ,  $\mathbb{B}_0^*$ . In practice, the location and shape of  $\mathbb{B}_0^*$  is unknown. A common way to deal with this issue is using "randomized initialization" where the initial points are sampled uniformly at random from the solution space. The idea is that by trying multiple, say  $m$ , random initializations, one of the initial points will be in  $\mathbb{B}_0^*$  and  $\mathcal{T}$  applied to that point will find  $\theta^*$ . Hence, the total computational cost in this case is  $m \cdot c(\mathcal{T})$ .

From the above discussion, we note that when learning a non-convex loss function, the computational cost is the product of two tasks: 1) Exploitation: running greedy algorithm starting from a given initial point and 2) Exploration: finding an initialization within the attraction region of the global minimum. To achieve high accuracy, the exploitation task given a good starting point often requires a sufficiently large amount of data and is very well understood in the literature [3]. In contrast, the exploration task is less studied. The performance can be problem dependent and the computational cost can be very high. One important insight that we will leverage in our subsequent development is that the landscape of the empirical risk based on a random sample of size  $n$ ,  $\widehat{F}_n(\theta) = \frac{1}{n} \sum_{i=1}^n f(\theta, x_i)$ , should resemble that of  $F(\theta)$  reasonably well when  $n$  is large enough. Thus, it is natural to ask that, in the data outsourcing context, if we can assign the exploration task to the computing facility. In other words, we split the computation tasks into two phases:

- *Exploration*: The computing facility is assigned to explore the energy landscape of  $\widehat{F}_n$ , where  $n$  is much smaller than the size of full dataset, and find a good initial point(s)  $\theta^0$  (or  $\theta_1, \dots, \theta_L$ ).
- *Exploitation*: The data organization can run more refined exploitation starting from  $\theta^0$  (or  $\theta_1, \dots, \theta_L$ ). In this case, the computational cost, from the data organization's perspective, can be reduced from  $m \cdot c(\mathcal{T})$  to  $c(\mathcal{T})$  (or  $L \cdot c(\mathcal{T})$ ), where  $m$  is the number of random initializations. Such a reduction can be substantial if  $m$  needs to be a large number to achieved a desired performance.

Similar computational strategy can also be applied even outside the data outsourcing context. The idea is that we can first use a less accurate loss function  $\widehat{F}_n$  with a smaller amount of data to find good initializations. We then employ greedy optimization algorithms on  $F$  starting from these carefully selected initial points.

**Our contribution.** First, we propose sampling-based algorithms to obtain good initializations for  $F(\theta)$ -optimization with outsourced data. Particularly, we design two types of procedures, sampling or optimization, depending on whether the optimization cost  $c(\mathcal{T})$  is moderate or large: If  $c(\mathcal{T})$  is moderate, multiple instances of  $\mathcal{T}$  can be implemented starting from different initial points. In this scenario, we suggest using samples from a distribution  $\pi_\beta \propto \exp(-\beta \widehat{F}_n)$  with a properly chosen  $\beta$  as initial points. If  $c(\mathcal{T})$  is large, only one instance of  $\mathcal{T}$  can be implemented. In this scenario, we

suggest starting from the global minimum of  $\widehat{F}_n$ . This minimizer can be obtained by implementing a proper selection procedure on samples from  $\pi_\beta$ .

Second, our analytical results provide rigorous justification of these procedures and guide how much data should be outsourced. In particular, we show that under appropriate conditions, when  $n \geq O(d \log(1/\rho)/\delta^2)$ , with probability  $(1 - \rho)$ , both methods can find a good initial point. Here,  $d$  is the dimension parameter and  $\delta$  is a parameter for the approximation accuracy, which may depend on the structure of the objective function. Under proper regularity conditions, when  $\mathcal{T}$  is initialized from the point(s) output from the exploration stage, with a high probability, it will find the global minimizer of  $F$ .

Noticeably, our procedures are compatible with the data outsourcing setup. In particular, the computing facility has access to only  $n$  data points. It will carry out either the sampling or the optimization procedure on  $\widehat{F}_n$  to generate good initial point(s). The data organization can then run a greedy optimization algorithm starting from these point(s) to optimize  $F$ . The data organization saves in-house computational effort in the second optimization stage. Meanwhile, it only exposes  $n$  data points to outside parties.

**Related literature.** Data outsourcing has been a problem of intensive interest in the last decade due the emergence of big data and cloud computing. Most existing works focus on data management policies and encryption [6, 12, 28]. To the best of our knowledge, this work is the first to study data outsourcing from a machine learning perspective.

Our problem can be viewed as a special non-convex stochastic optimization problem. How to efficiently solve non-convex stochastic programming is a fast developing area [1, 14, 32]. Our contribution is the development of a new initialization method. When solving non-convex optimization problems, while finding good initial points is an important problem, the related literature is rather limited. The most common approach is using crude uniform sampling, which is likely suboptimal. Our approach provides a computationally feasible refined solution to this problem. Finding good initialization in more specific problem settings has been studied in the literature. For example, [5] studies the efficacy of gradient descent with random initialization for solving systems of quadratic equations. Weight initialization for neural networks has been investigated in [2, 16, 35]. Spectral initialization has been proposed for generalized linear sensing models in the high dimensional regime [23]. The key advantage of our propose method is its general applicability and theoretical performance guarantee.

Our problem is related to but different from federated learning. Federated learning is a special form of distributed learning where the central learning agent do not have access to or control over individual agent's (distributed worker's) device and data [21]. Most existing development in federated learning try to address two main challenges: i) the communication cost, which can be extremely high (higher than the computational cost) and ii) the agents (distributed workers) are heterogeneous where the data stored with individual agent may not be representative (non i.i.d.) (see, for example, [18, 22, 36]). In contrast, our setting assume the data organization owns all the data and can decide what to distribute to outsourcing computing facilities. Thus, we can ensure that the data send to individual workers are representative. The task we assign to individual workers is also fundamentally different from federated learning. In our setting, we divide the learning into two stages. The outsourcing stage (distributed stage) where the objective is to find good initialization and the in-house stage where we try to learn the optimal solution. We also focus on the non-convex learning setting, which is not much studied in federated learning.

Our problem is also related to but different from simulated annealing or tempering-based algorithms. Simulated annealing tries to integrate the exploration of different local minimums with

the exploitation to pinpoint the global minimum [19]. We, on the other hand, separate the exploration and exploitation task to two different entities. Monte Carlo simulation can be used for the exploration task in our setting. The main advantage of our algorithm is that we use sampling-based approach to find good initial point with limited data. To the best of our knowledge, this particular setting has not been discussed in the literature.

**Notations.** We use  $\|\theta\|$  and  $\|A\|_{\text{op}}$  to denote the  $L_2$ -norm of a vector  $\theta$  and the operator norm of a matrix  $A$  respectively. For real numbers  $a, b$ , let  $a \wedge b = \min\{a, b\}$  and  $a \vee b = \max\{a, b\}$ . Lastly, given two sequences of real numbers  $\{a_n\}_{n \geq 1}$  and  $\{b_n\}_{n \geq 1}$ ,  $a_n = O(b_n)$  denotes that there exist a constant  $C > 0$ , such that  $a_n \leq Cb_n$ , and  $a_n = \Omega(b_n)$  denotes that  $a_n \geq Cb_n$ .

## 2 METHODOLOGY

We consider minimizing a smooth but non-convex function  $F(\theta)$ , which takes the form

$$F(\theta) = \mathbb{E}_{X \sim \xi} [f(\theta, X)],$$

over a  $d$ -dimensional unit ball  $\Theta = \{\theta \in \mathbb{R}^d : \|\theta\| \leq 1\}$ . It is also common to consider an empirical loss function of  $N$  data points:  $F_N(\theta) = \frac{1}{N} \sum_{i=1}^N f(\theta, x_i)$ . This can be seen as a special form of  $F(\theta)$  where  $\xi$  is the empirical distribution of the full dataset  $\{x_1, \dots, x_N\}$ . We also comment that in most applications, the solution to the optimization problem needs to be restricted to some known range. In practice, we can do whitening transformation or other rescaling so the the solution is in the unit ball. From the theoretical perspective, considering bounded domain greatly simplifies our discussion and this assumption is commonly imposed in the literature [25].

Since  $F(\cdot)$  is non-convex, the performance of any greedy deterministic optimization algorithm relies heavily on the choice of initial points. Specifically, a deterministic optimization algorithm  $\mathcal{T}$  such as gradient descent (GD) or Newton's method can be trapped in a suboptimal local minimum instead of converging to the desired global minimum if initialized inappropriately. In practice, the initial points are usually sampled uniformly at random when the structure of  $F(\theta)$  is unknown. Despite seeming simple and plausible, this approach lacks theoretical justification and can be highly inefficient. In this work, we design data outsourcing and exploration mechanisms to find good initial points for the optimization algorithm  $\mathcal{T}$ . The objective is to increase the chance that  $\mathcal{T}$  finds the global minimum successfully.

Assume the outsourced data  $\{x_1, \dots, x_n\}$  follow the same distribution as  $\xi$ . We can construct a sample approximation to  $F(\theta)$  as  $\widehat{F}_n(\theta) = \frac{1}{n} \sum_{i=1}^n f(\theta, x_i)$ . Evaluating  $\widehat{F}_n(\theta)$  or  $\nabla \widehat{F}_n(\theta)$  has a much smaller cost than evaluating  $F(\theta)$  or  $\nabla F(\theta)$  if the sample size  $n$  is not too large. This makes exploring the energy landscape of  $\widehat{F}_n(\cdot)$  more computationally friendly. Note that  $\widehat{F}_n(\cdot)$  captures certain structural information of  $F(\cdot)$ . We are interested in utilizing this information in an appropriate way. More specifically, the work of [25] has shown that the energy landscape of  $\widehat{F}_n(\theta)$  bears close similarity to that of  $F(\theta)$  when  $n$  surpasses a certain threshold. This indicates that the global minimum of  $\widehat{F}_n(\theta)$  should be closer to that of  $F(\theta)$  than a random guess. Let  $\widehat{\theta}_0^*$  denote the global minimum of  $\widehat{F}_n(\theta)$  and  $\theta_0^*$  denote the global minimum of  $F(\theta)$ . Intuitively, if we use  $\widehat{\theta}_0^*$  as the initial point to apply the optimization algorithm  $\mathcal{T}$ , we might be more likely to converge to  $\theta_0^*$ . We refer to this approach as the *optimization approach*. It is quite computational friendly to the data organization, since only one instance of  $\mathcal{T}$  is needed. However, it also comes with certain costs: 1)  $\widehat{F}_n(\theta)$  is a noisy realization of  $F(\theta)$ , especially when  $n$  is small. Using just the global minimizer of  $\widehat{F}_n$ , which is a single point, can be risky. 2)  $\widehat{F}_n$  is likely to be nonconvex as well and optimizing it can be expensive. For 2), since the task is outsourced to a computing facility, the in-house cost is reduced though.

An alternative approach we consider is to sample a Bayesian posterior distribution with the outsourced data. In Bayesian statistics, the unknown parameter  $\theta$  is usually represented using a posterior density, which is proportional to the product of a prior density and the likelihood function. Since we require  $\|\theta\| \in \Theta$ , it is natural to assume the prior distribution is the uniform distribution on  $\Theta$ . In many applications, the loss function  $f(\theta, x_i)$  is proportional to the negative log likelihood. For example, if we model the data output as a function of the input plus Gaussian noise, i.e.  $x_{\text{out}} = g(\theta, x_{\text{in}}) + \xi$  where  $\xi \sim \mathcal{N}(0, \sigma^2)$ , the likelihood function is given by  $-\log p(x|\theta) = \frac{1}{2\sigma^2} (x_{\text{out}} - g(\theta, x_{\text{in}}))^2 := \frac{1}{2\sigma^2} f(\theta, x)$ . Then, the posterior distribution is given by

$$p(\theta) \propto 1_{(\theta \in \Theta)} \prod_{i=1}^n p(x_i|\theta) = 1_{(\theta \in \Theta)} \exp\left(-\frac{n}{2\sigma^2} \widehat{F}_n(\theta)\right). \quad (1)$$

Samples from the posterior distribution learn from  $x_1, \dots, x_n$ . Thus, they are more informative than samples from the prior distribution. Comparing with the optimization approach, this *sampling approach* takes into account that  $\widehat{F}_n$  is noisy, so the candidate solution is not a single point, but a distribution which accounts for the uncertainty. In this case, the data organization needs to implement  $\mathcal{T}$  from multiple samples generated from the posterior distribution.

Given a deterministic optimization algorithm  $\mathcal{T}$ , when only partial data is available, there is in general no clear theoretical guarantee when determining whether a point is a good initial point to optimize  $F$ . Both the optimization approach and the sampling approach use criteria based on  $\widehat{F}_n$ . Our theoretical analysis shows when these criteria are sufficient. We next provide more details of these two approaches. While the optimization approach is conceptually simpler, its computation requires sampling tools. Thus, we start with the sampling approach.

*Procedures with the sampling approach.* For the exploration task, we consider sampling from a distribution

$$\pi_\beta(\theta) \propto \exp(-\beta \widehat{F}_n(\theta)) \cdot 1_{\{\theta \in \Theta\}}. \quad (2)$$

The parameter  $\beta > 0$  is often referred to as the inverse temperature [34]. The posterior distribution in (1) corresponds to  $\beta = \frac{n}{2\sigma^2}$ . We consider general  $\beta$  because in practice the observation noise  $\sigma^2$  may not be known. The parameter  $\beta$  determines how much  $\pi_\beta(\theta)$  concentrates around the global minimum of  $\widehat{F}_n(\theta)$ . A larger  $\beta$  leads to a higher concentration around  $\hat{\theta}_0^*$ . When  $\beta = \infty$ , we get  $\hat{\theta}_0^*$  with probability one. Using  $\hat{\theta}_0^*$  as a starting point is likely to be a good choice if  $n$  is large enough and  $\widehat{F}_n$  is close to  $F$ . Meanwhile, when  $\beta = 0$ ,  $\pi_\beta$  is simply the uniform distribution, which is equivalent to the standard random start. In this sense, sampling from  $\pi_\beta$  with  $\beta \in (0, \infty)$  can be viewed as an interpolation of two extreme cases.

There is a rich literature on how to sample from  $\pi_\beta(\theta)$ . When  $\pi_\beta$  is simple or close to some simple reference distributions, independent samples can be obtained through rejection sampling or importance sampling. For more complicated target distributions, Markov Chain Monte Carlo (MCMC) is usually applied. In general, these algorithms simulate stochastic processes of which  $\pi_\beta$  is the invariant distribution. Popular and simple choices include random walk Metropolis, unadjusted Langevin algorithm (ULA) [9], Metropolis-adjusted Langevin algorithm [27]. Recent studies show that these MCMC algorithms are efficient when the sampling distribution is log-concave with perturbations [10, 24]. When  $\widehat{F}_n$  is non-convex with separated local minima,  $\pi_\beta$  is a multimodal distribution, and it can be difficult to sample directly with these algorithms. This is particular the case if  $\beta$  is large, since the stochastic algorithm may stick to one mode for many iterations before visiting the other modes. This issue can often be solved using methods such as parallel tempering or simulated tempering [7, 13, 30, 33]. The papers [13, 20] show that a simulated tempering algorithm can sample a multimodal distribution with polynomial complexity.

Given the sample  $\theta_1, \dots, \theta_L$  from  $\pi_\beta$ , the data organization then implement  $\mathcal{T}$  starting from each  $\theta_i$ . Let  $\mathcal{T}(\theta)$  denote the output of the optimization algorithm  $\mathcal{T}$  starting from  $\theta$ . The actual exploration algorithm is summarized in Algorithm 1. Our theoretical analysis in the next section gives rigorous justification of this procedure assuming  $\beta$  is large enough. In practice, this approach is more efficient than the naive random start even with moderate  $\beta$  as we will demonstrate through numerical experiment in Section 4. We also emphasize that our analysis applies to most of existing sampling tools where  $\theta_1, \dots, \theta_L$  do not need to be independent.

---

**Algorithm 1** Sampling-based Initial Point Selection (SIPS)
 

---

**Input:** Outsourced data sample  $\{x_1, \dots, x_n\}$ , inverse temperature parameter  $\beta$ , sampling algorithm  $\mathcal{M}$ , exploration sample size  $L$ .

**Initialization:** Construct the empirical average  $\hat{F}_n(\theta) = \frac{1}{n} \sum f(\theta, x_i)$  and the target density  $\pi_\beta(\theta) \propto \exp(-\beta \hat{F}_n(\theta)) \cdot 1_{\{\theta \in \Theta\}}$ .

**Sampling:** Apply  $\mathcal{M}$  to draw samples  $\{\theta_1, \dots, \theta_L\}$  from distribution  $\pi_\beta$ .

**Output:** Candidate initial points  $\{\theta_1, \dots, \theta_L\}$ .

---

*Procedures with the optimization approach.* When  $\hat{F}_n$  is non-convex, there is no consensus on how to find its global minimizer. Typical choices include either using meta-heuristic algorithms or sampling-based algorithms. Here we consider using sampling-based algorithms due to their connection to the sampling approach.

One popular way to find the global minimum of  $\hat{F}_n$  involves generating samples  $\theta_1, \dots, \theta_L$  from the distribution  $\pi_\beta$  with a large  $\beta$ . This approach is investigated by [4, 26, 34] when ULA or its online version is implemented to sample from  $\pi_\beta$ . As mentioned earlier, the parameter  $\beta$  determines how much  $\pi_\beta(\theta)$  concentrates around the global minimum of  $\hat{F}_n(\theta)$  and a larger  $\beta$  leads to a higher concentration. When the samples  $\theta_1, \dots, \theta_L$  are available as candidate solutions, we can choose the one with the lowest objective value, i.e.,  $\theta_{i^*}$ , where

$$i^* = \operatorname{argmin}_{i \in \{1, \dots, L\}} \hat{F}_n(\theta_i). \quad (3)$$

This procedure is summarized as the *annealing* approach in Algorithm 2. In order for this approach to be effective at finding the global minimum of  $\hat{F}_n$ ,  $\beta$  needs to be large enough. This usually increases the difficulty of sampling from  $\pi_\beta$ . On the other hand, it is worth noticing that we are only interested in getting good starting points for optimizing  $F$ . Thus, finding the global minimum of  $\hat{F}_n$  approximately can often serve the purpose. This suggests a less extreme  $\beta$  may be sufficient.

The criterion in (3) finds the  $\theta_i$  with the lowest  $\hat{F}_n$ -value. Further refinement can be applied to improve the quality of the initial point. For example, if we apply a deterministic optimization algorithm  $\hat{\mathcal{T}}$  to  $\hat{F}_n$  initialized at  $\theta_i$ , we can achieve an even lower  $\hat{F}_n$ -value. We then pick  $\hat{\mathcal{T}}(\theta_i)$  with the lowest  $\hat{F}_n$ -value as the initial point, i.e.,  $\hat{\mathcal{T}}(\theta_{i^*})$ , where

$$i^* = \operatorname{argmin}_{i \in \{1, \dots, L\}} \hat{F}_n(\hat{\mathcal{T}}(\theta_i)). \quad (4)$$

This procedure is summarized as the *sampling-assisted-optimization (SAO) approach* in Algorithm 2. The SAO approach is similar to GDxLD developed in [8]. Comparing to the simpler criterion (3), sampling for (4) can often be done more efficiently. This is because when implementing SAO for  $\hat{F}_n$ , we separate the exploration task and the optimization task. This allows us to use a smaller  $\beta$  when sampling  $\pi_\beta$ . The cost is that invoking  $\hat{\mathcal{T}}$  to each sample in SAO can impose extra computational cost than the annealing approach. In contrast, the annealing approach combines the exploration task with the optimization task. So a larger  $\beta$  is needed in general, which increases the cost to sample from  $\pi_\beta$ .

**Algorithm 2** Optimization-based Initial Point Selection (OIPS)

**Input:** Outsourced data sample  $\{x_1, \dots, x_n\}$ , inverse temperature parameter  $\beta$ , exploration sample size  $L$ , sampling algorithm  $\mathcal{M}$ , optimization algorithm  $\hat{\mathcal{T}}$ .

**Initialization:** Construct the empirical average  $\hat{F}_n(\theta) = \frac{1}{n} \sum f(\theta, x_i)$  and the target density  $\pi_\beta(\theta) \propto \exp\{-\beta \hat{F}_n(\theta)\} \cdot 1_{\{\theta \in \Theta\}}$

**Sampling:** Apply  $\mathcal{M}$  to draw a sample  $\{\theta_1, \dots, \theta_L\}$  from distribution  $\pi_\beta$ .

**if Annealing then**

Set  $\theta^0 = \theta_{i^*}$  where  $i^* = \operatorname{argmin}_{i \in \{1, \dots, L\}} \hat{F}_n(\theta_i)$ .

**end if**

**if Sampling-assisted-optimize (SAO) then**

Set  $\theta^0 = \hat{\mathcal{T}}(\theta_{i^*})$  where  $i^* = \operatorname{argmin}_{i \in \{1, \dots, L\}} \hat{F}_n(\hat{\mathcal{T}}(\theta_i))$ .

**end if**

**Output:** Candidate initial points  $\theta^0$

### 3 THEORETICAL GUARANTEE IN FINDING THE GLOBAL MINIMUM

In this section, we analyze the performance of Algorithms 1 and 2. The key in successful implementation of the algorithms is to set the appropriate outsourcing sample size  $n$ , inverse temperature  $\beta$ , and exploration sample size  $L$ . Our performance analysis provides guidelines on choosing these parameters.

*Conditions on the energy landscapes.* We start with some assumptions on the energy landscape of  $F(\theta)$  and the randomness when evaluating  $f(\theta, x)$ . Many of them are also assumed in [25]. Since we run an optimization algorithm  $\mathcal{T}$  that converges to a stationary point in the second phase, the following assumption regularizes the configuration of the stationary points:

ASSUMPTION 1.  $F(\theta) : \Theta \rightarrow \mathbb{R}$  is  $(\sigma, \eta)$ -strongly Morse, that is,  $\|\nabla F(\theta)\| \geq \sigma$  for  $\|\theta\| = 1$ , and  $\lambda_{\min}(\nabla^2 F(\theta)) \geq \eta$  if  $\|\nabla F(\theta)\| \leq \sigma$ , where  $\lambda_{\min}(A)$  is the minimum eigenvalue of  $A$ . Moreover  $L^* := \sup_{\theta \in \Theta} \|\nabla^3 F(\theta)\|_{op} < \infty$ .

One consequence of Assumption 1 is that all the stationary points of  $F(\theta)$  in  $\Theta$  are finite and well-separated [25]. In particular, we can denote these stationary points as  $(\theta_0^*, \theta_1^*, \dots, \theta_K^*)$ . Without loss of generality, let  $\theta_0^*$  be the global minimum of  $F(\theta)$ .

For simplicity of discussion, we assume that  $\mathcal{T}$  is a deterministic optimization algorithm that is guaranteed to converge to a stationary point and the performance of  $\mathcal{T}$  is determined by the initial point. Starting from  $\theta^0$ , we denote the stationary point that  $\mathcal{T}$  converges to as  $\mathcal{T}(\theta^0)$ . Hence,  $\mathcal{T}$  can be viewed as a deterministic mapping from the parameter space  $\Theta$  to the set of stationary points  $\{\theta_0^*, \theta_1^*, \dots, \theta_K^*\}$ . Our goal is to find a  $\theta^0$  such that  $\mathcal{T}(\theta^0) = \theta_0^*$ .

Given the deterministic optimization algorithm  $\mathcal{T}$ , the attraction region of the global minimum  $\theta_0^*$  can be defined as

$$\mathbb{B}_0^* = \{\theta \in \Theta : \mathcal{T}(\theta) = \theta_0^*\}.$$

In general,  $\mathbb{B}_0^*$  cannot be characterized without  $\mathcal{T}$ . On the other hand, it is well-known that for many optimization algorithms,  $\mathcal{T}(\theta^0) = \theta_0^*$  if  $\theta^0$  is in a neighborhood of  $\theta_0^*$  in which  $F(\theta)$  is strongly convex. This indicates that a proper neighborhood of  $\theta_0^*$  can be used as a substitution of  $\mathbb{B}_0^*$ . We formalize this idea as follows.

ASSUMPTION 2. There exists a ball centered at  $\theta_0^*$  with radius  $r$ ,  $\mathcal{B}_r(\theta_0^*) = \{\theta : \|\theta - \theta_0^*\| \leq r\}$ , such that  $\mathcal{B}_r(\theta_0^*) \subseteq \mathbb{B}_0^*$  and  $F(\theta)$  is  $\mu$ -strongly convex in  $\mathcal{B}_r(\theta_0^*)$ .

Note that Assumption 2 may come as a consequence of Assumption 1. In particular,  $F(\theta)$  is  $\eta/2$ -strongly convex in  $\mathcal{B}_r(\theta_0^*)$  when  $r \leq \frac{\eta}{2L}$ .

Notably, the assumptions above enable us to derive an upper bound for the failure rate of the benchmark random start algorithm. If we draw  $m$  independent initial points from  $\Theta$  uniformly at random, the probability that none of them leads to the global minimum of is  $\mathbb{P}(\mathcal{F}_b) \leq (1 - \mathbb{P}(\mathcal{B}_r))^m$ , or

$$\log \mathbb{P}(\mathcal{F}_b) \leq -L \log(\mathbb{P}(\mathcal{B}_r)) = \Omega(Lr^d). \quad (5)$$

Then, in order for it to be lower than a threshold  $\rho$ , we need  $m = \log(\rho)/|r|^d$ , which has an exponential dependence on  $d$ .

Our next assumption concerns the uniqueness of the global minimum. When  $\theta_0^*$  is the unique global minimum, its function value needs to be strictly lower than the other stationary points.

**ASSUMPTION 3.** *There exists a constant  $\alpha > 0$ , such that for all  $\theta \notin \mathcal{B}_r(\theta_0^*)$ ,  $F(\theta) - F(\theta_0^*) \geq \alpha$ .*

In Section 3.3, we will discuss what can be achieved if this assumption does not hold.

The basic idea of our data outsourcing and exploration scheme is to approximate  $F(\theta)$  via its sample average  $\hat{F}_n(\theta)$  and then use the global minimum of  $\hat{F}_n(\theta)$  as the initial point to optimize  $F(\theta)$ . A key question is that in order for  $\hat{F}_n(\theta)$  to be a good approximation of  $F(\theta)$ , how many data points are needed. This problem has been studied [25]. We adapt some of their results into our setting. This involves the following regularity conditions on the loss function and noises (similar versions of them can be found in [25] as well).

**ASSUMPTION 4.** *The following hold for some  $\tau, c_h$*

(1) *The loss function for each data point is  $\tau^2$ -sub-Gaussian. Namely, for any  $\lambda \in \mathbb{R}^p$ , and  $\theta \in \Theta$ ,*

$$\mathbb{E} \left[ \exp \left( \langle \lambda, f(\theta; X) - \mathbb{E}_{X \sim \xi} [f(\theta; X)] \rangle \right) \right] \leq \exp \left\{ \frac{\tau^2 \|\lambda\|^2}{2} \right\}.$$

(2) *The gradient of the loss is  $\tau^2$ -sub-Gaussian. Namely, for any  $\lambda \in \mathbb{R}^p$ , and  $\theta \in \Theta$ ,*

$$\mathbb{E} \left[ \exp \left\langle \lambda, \nabla_{\theta} f(\theta; X) - \mathbb{E}_{X \sim \xi} [\nabla_{\theta} f(\theta; X)] \right\rangle \right] \leq \exp \left\{ \frac{\tau^2 \|\lambda\|^2}{2} \right\}.$$

(3) *The Hessian of the loss, evaluated on a unit vector, is  $\tau^2$ -sub-exponential. Namely, for any  $\|\lambda\| \leq 1$ , and  $\theta \in \Theta$ ,*

$$\mathbb{E} \left[ \exp \left\{ \frac{1}{\tau^2} |\mathcal{Z}_{\lambda, \theta}(X) - \mathbb{E}_{X \sim \xi} [\mathcal{Z}_{\lambda, \theta}(X)]| \right\} \right] \leq 2,$$

where  $\mathcal{Z}_{\lambda, \theta}(X) = \langle \lambda, \nabla_{\theta}^2 f(\theta; X) \lambda \rangle$ .

(4) *There exists  $J_*$  (potentially diverging polynomially in  $d$ ) such that*

$$\mathbb{E}_{X \sim \xi} [J^1(X)], \mathbb{E}_{X \sim \xi} [J^2(X)] \leq J_*$$

where

$$J^1(X) = \sup_{\theta_1, \theta_2 \in \Theta, \theta_1 \neq \theta_2} \frac{\|\nabla_{\theta} f(\theta_1; X) - \nabla_{\theta} f(\theta_2; X)\|}{\|\theta_1 - \theta_2\|},$$

$$J^2(X) = \sup_{\theta_1, \theta_2 \in \Theta, \theta_1 \neq \theta_2} \frac{\|\nabla_{\theta}^2 f(\theta_1; X) - \nabla_{\theta}^2 f(\theta_2; X)\|_{op}}{\|\theta_1 - \theta_2\|}.$$

Furthermore, there exists a constant  $c_h$  such that  $J_* \leq \tau^3 d^{c_h}$ .

(5) *There exists  $\theta^* \in \Theta$ , such that  $\|\nabla F(\theta^*)\|, \|\nabla^2 F(\theta^*)\|_{op} \leq H \leq \tau^3 d^{c_h}$ .*

Assumption 4 allows us to find a close approximation of  $F$ , which is formally defined as follows.



*Definition 3.1.* We say  $\widehat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$ , if both  $F$  and  $\widehat{F}_n$  have  $K + 1$  stationary points, denoted by  $\{\theta_i^*\}_{i=0,\dots,K}$  and  $\{\hat{\theta}_i^*\}_{i=0,\dots,K}$ , and the following inequalities hold

$$\begin{aligned} \sup_{\theta \in \Theta} |F(\theta) - \widehat{F}_n(\theta)| &\leq \delta, \quad \sup_{\theta \in \Theta} \|\nabla F(\theta) - \nabla \widehat{F}_n(\theta)\| \leq \delta, \\ \sup_{\theta \in \Theta} \|\nabla^2 F(\theta) - \nabla^2 \widehat{F}_n(\theta)\|_{\text{op}} &\leq \delta, \quad \text{and} \quad \max_{0 \leq i \leq K} \|\theta_i^* - \hat{\theta}_i^*\| \leq \delta. \end{aligned}$$

The next lemma characterizes the minimal sample size required to achieve a  $\delta$ -approximation.

**LEMMA 3.2.** *Assume that Assumptions 1 and 4 hold. Consider a given confidence level  $\rho \in (0, 1)$  and a given accuracy  $\delta$ . Let  $C = C_0 \cdot (c_h \vee 1 \vee \log(\tau/\rho)) = O(|\log \rho|)$ , where  $C_0$  is some absolute constant, and  $\eta_* = (\sigma^2/\tau^2) \wedge (\eta^2/\tau^4) \wedge (\eta^4/((L^* \tau)^2)) = \Omega(1)$ . For an arbitrary constant  $\iota > 0$ , let  $C_\iota$  be a constant such that  $\log(n) \leq C_\iota \cdot n^\iota$ . Then, when*

$$\begin{aligned} n &\geq \max \left\{ \left[ \frac{C_\iota C d}{(\delta/((2\tau/\eta) \vee \tau \vee \tau^2))^2} \right]^{\frac{1}{1-\iota}}, 4Cd (\log(d) \vee \log(n)/\eta_*^2) \right\}, \\ &:= n(\delta, \rho, d), \end{aligned}$$

with probability at least  $1 - \rho$ ,  $\widehat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$ .

The proofs of Lemma 3.2 and all subsequent results are provided in Appendix A.

Lemma 3.2 quantifies that to achieve a  $\delta$ -approximation of  $F(\theta)$  with confidence level  $1 - \rho$ , the required sample size is

$$n(\delta, \rho, d) = O\left(\frac{d \log(1/\rho)}{\delta^2}\right). \quad (6)$$

Here, we ignore the index  $\iota$  in the power since it can be made arbitrarily small.

### 3.1 Performance of the sampling approach

We denote  $\mathcal{F}_0$  as the event that using the initial point(s) constructed based on Algorithm 1, the optimization algorithm in the exploitation stage fails to find the global minimum. In this section, we establish an upper bound for  $\mathbb{P}(\mathcal{F}_0)$ .

Recall that samples are drawn from the distribution  $\pi_\beta$  defined in (2). We will justify that when  $\beta$  is large enough, a random sample  $\hat{\theta}_\beta$  from  $\pi_\beta(\theta)$  is a good starting point to optimize  $F(\theta)$ . In particular,  $\tilde{\theta}_\beta$  has a high chance to fall into an attraction basin of  $\theta_0^*$ , i.e.,  $\mathcal{B}_r(\theta_0^*)$ .

**PROPOSITION 3.3.** *Suppose Assumptions 1-4 hold and the approximation accuracy  $\delta$  satisfies  $\delta < \mu \wedge r \wedge \alpha/4$ . If  $\widehat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$  and  $\beta \geq \Omega(r^{-2})$ , then the probability that  $\tilde{\theta}_\beta$  fails to be a good starting point, i.e.  $\pi_\beta(\mathcal{B}_r^c(\theta_0^*))$ , is bounded by:*

$$\log(\pi_\beta(\mathcal{B}_r^c(\theta_0^*))) = O(-\beta\alpha/2 + d \log(\beta)).$$

Proposition 3.3 shows that as the inverse temperature parameter  $\beta$  increases, the probability that we can sample points from  $\mathcal{B}_r(\theta_0^*)$  approaches one exponentially fast. The convergence speed is determined by  $\alpha$ , the gap between the global minimum and other local minima, as well as the dimension parameter  $d$ . However, in practice, we cannot choose  $\beta$  arbitrarily large as we have to consider the computational cost in associated sampling algorithms (e.g., an MCMC algorithm). In general, when  $\beta$  increases, the difficulty of sampling from  $\pi_\beta(\theta)$  increases. In practice, we want to find a  $\beta$  that balances the estimation accuracy and the sampling difficulty.

One difficulty when applying Proposition 3.3 to the sampling approach is that in practice we may not be able to sample from  $\pi_\beta$  exactly. For example, many MCMC algorithms can only draw

samples from a distribution that is “close” to  $\pi_\beta$ . To handle this issue, To handle this issue, we impose the following assumption as a relaxation to the requirement of sampling from  $\pi_\beta$  exactly.

**ASSUMPTION 5.** *There is a sampler  $\hat{\mathcal{M}}$  such that for any fixed  $\delta_\beta \in [0, 1)$ , starting from any  $\theta_0 \in \Theta$ ,  $\hat{\mathcal{M}}$  can draw samples from a distribution  $\hat{\pi}_\beta$  which satisfies  $\|\hat{\pi}_\beta - \pi_\beta\|_{TV} \leq \delta_\beta$ .*

In addition, note that in practice, we can draw consecutive samples from the same chain of the underlying MCMC algorithm, which makes the samples correlated. The following lemma justifies the quality of the samples form  $\hat{\mathcal{M}}$  under Assumption 5.

**LEMMA 3.4.** *Given a set  $B$  and distribution  $\pi_\beta$  with  $\pi_\beta(B) > 0$ , suppose there exists a samplers  $\hat{\mathcal{M}}$  satisfying Assumption 5. If we have  $L$  samples from  $\hat{\mathcal{M}}$ , then*

$$\mathbb{P}(X_1 \notin B, \dots, X_L \notin B) \leq (\pi(B^c) + \delta_\beta)^L.$$

The following theorem then comes as a consequence of Proposition 3.3 and Lemma 3.4.

**THEOREM 3.5.** *Consider Algorithm 1. Suppose Assumptions 1-5 hold. For an arbitrary confidence level  $\rho \in (0, 1)$ , let  $\delta = \mu \wedge r \wedge \alpha/4$ . If sample size  $n \geq n(\delta, \rho, d) = O(d \log(1/\rho)/\delta^2)$  and the inverse temperature  $\beta \geq \Omega(r^{-2})$ , then there exists a constant  $C \in (0, \infty)$  such that*

$$\mathbb{P}(\mathcal{F}_0) \leq \rho + \exp\{CL \cdot \max\{-\beta\alpha/2 + d \log(\beta), \log(2\delta_\beta)\}\}. \quad (7)$$

Theorem 3.5 shows that the probability that the sampling approach fails to find the global minimum of  $F(\theta)$  decays exponentially fast as the inverse temperature  $\beta$  and sample size  $L$  increase. In particular,  $L$  only needs to surpass some dimensional independent constants, i.e., the convexity constant  $\mu$  and the separability constant  $\alpha$  of global minimum from other local minima. In contrast, by (5), the benchmark random start method would require the number of random initialization  $m$  to depend exponentially on the dimension. We comment that Algorithm 1 does require an outsourced sampling algorithm to obtain samples from  $\pi_\beta$ , which can be computationally costly, but this task is outsourced and we achieve a much smaller the in-house computational cost. Finally, it is worth mentioning that in Theorem 3.5, both  $\beta$  and  $n$  scale as  $r^{-2}$ . This is in agreement with the Bayesian setup (1), which suggests  $\beta$  should scale linearly with  $n$ .

### 3.2 Performance of the optimization approaches

We first provide an analysis of the SAO approach in Algorithm 2. Let  $\mathcal{F}_1$  denote the random event that the output of Algorithm 2-SA0 approach fails to find the global minimum of  $F(\theta)$ . The result is largely the same as Theorem 3.5, although the proof is slightly more difficult.

**THEOREM 3.6.** *Consider Algorithm 2-SA0. Suppose Assumptions 1-5 hold. For an arbitrary confidence level  $\rho \in (0, 1)$ , let  $\delta = \mu \wedge r \wedge \alpha/4$ . If the sample size  $n \geq n(\delta, \rho, d) = O(d \log(1/\rho)/\delta^2)$  and the inverse temperature  $\beta \geq \Omega(r^{-2})$ , then there exists a constant  $C \in (0, \infty)$  such that*

$$\mathbb{P}(\mathcal{F}_1) \leq \rho + \exp\{CL \cdot \max\{-\beta\alpha/2 + d \log(\beta), \log(2\delta_\beta)\}\}.$$

We next analyze the annealing approach in Algorithm 2. Let  $\mathcal{F}_2$  be the random event that Algorithm 2-annealing fails to find the global minimum of  $F(\theta)$ . The annealing approach needs more restrictions than the SAO approach. This is because: in order to generate a good starting point, one of the samples need to fall close to  $\theta_0^*$ . Moreover, its  $\hat{F}_n$ -value needs to be lower than the other samples. This can be formulated as requiring a smaller radius  $r_0$  for the attraction neighborhood:

**THEOREM 3.7.** *Consider Algorithm 2-annealing. Suppose Assumptions 1-5 hold. For an arbitrary confidence level  $\rho \in (0, 1)$ , let  $r = r_0$  where  $r_0^2 \cdot \sup_{\theta \in \Theta} \|\nabla^2 F(\theta)\|_{op} < \alpha$  and  $\delta = \mu \wedge r \wedge \alpha/4$ . If the*

sample size  $n \geq n(\delta, \rho, d) = O(d \log(1/\rho)/\delta^2)$  and the inverse temperature  $\beta \geq \Omega(r^{-2})$ , then there exists a constant  $C \in (0, \infty)$  such that

$$\mathbb{P}(\mathcal{F}_2) \leq \rho + \exp\left\{CL \cdot \max\left\{-\beta\alpha/2 + d \log(\beta), \log(2\delta\beta)\right\}\right\}.$$

### 3.3 Extension to $\epsilon$ -Global Minimum

One major constraint in our previous analysis is Assumption 3—the global minimizer is unique with a gap of  $\alpha > 0$ . In practice, there can be multiple local minima that have function values very close to the global minimum. In this setting, it can be too ambitious to find the global minimum and it may be more reasonable to find an approximately optimal solution. Given a user-specified accuracy level  $\epsilon$ , we are interested in finding a local minimum whose objective value is within  $\epsilon$ -distance from the optimal objective value, i.e.,  $\theta_i^*$  such that  $F(\theta_i^*) \leq F(\theta_0^*) + \epsilon$ . We call a such local minimum an  $\epsilon$ -global minimum of  $F(\theta)$ . In this subsection, we conduct performance analysis for our algorithms to find an  $\epsilon$ -global minimum. Let

$$\mathcal{J}_\epsilon^* = \{i : F(\theta_i^*) \leq F(\theta_0^*) + \epsilon\}$$

be the index set of the  $\epsilon$ -global minimums. To be concise, we only present the analysis for the annealing-based optimization approach (Algorithm 2-annealing). The results for the other methods are similar.

We first introduce the “attraction region” of the  $\epsilon$ -global minimums:

*Definition 3.8 (“Attraction region” of  $\epsilon$ -global minimums).* Given an optimization algorithm  $\mathcal{T}$ , we define the attraction basin of  $\epsilon$ -global minimums of  $F(\theta)$  as

$$\mathbb{B}_\epsilon^* = \{\theta \in \Theta : F(\mathcal{T}(\theta)) \leq F(\theta_0^*) + \epsilon\}.$$

By definition, the optimization algorithm  $\mathcal{T}$  converges to an  $\epsilon$ -global minimum if and only if it starts with an initial point in  $\mathbb{B}_\epsilon^*$ . However, same as before,  $\mathbb{B}_\epsilon^*$  is hard to characterize directly. So we consider the following subset as a substitution

$$\mathcal{B}_{\epsilon, r_\epsilon} := \bigcup_{i \in \mathcal{J}_\epsilon^*} \mathcal{B}_{r_\epsilon}(\theta_i^*) \subseteq \mathbb{B}_\epsilon^*.$$

Let  $\mathcal{F}_{\epsilon, 2}$  be the random event that the output of Algorithm 2-annealing fails to find the  $\epsilon$ -global minimum of  $F(\theta)$ .

**THEOREM 3.9.** *Suppose Assumptions 1, 4 and 5 hold. For any user-specified accuracy  $\epsilon > 0$ , pick  $r_\epsilon$  such that  $\sup_{\theta \in \Theta} \|\nabla^2 F(\theta)\|_{op} \cdot r_\epsilon^2 \leq \epsilon$ . In addition, assume the approximation accuracy  $\delta$  satisfies  $\delta < \mu \wedge r_\epsilon \wedge \epsilon/12$ . For an arbitrary confidence level  $\rho \in (0, 1)$ , if the sample size  $n \geq n(\delta, \rho, d) = O(d \log(1/\rho)/\delta^2)$  and the inverse temperature  $\beta \geq \Omega(r_\epsilon^{-2})$ , then there exists a constant  $C \in (0, \infty)$  such that*

$$\mathbb{P}(\mathcal{F}_{\epsilon, 2}) \leq \rho + \exp\left\{CL \cdot \max\left\{-\beta\epsilon/6 + d \log(\beta), \log(2\delta\beta)\right\}\right\}.$$

Note that Theorem 3.9 establishes a similar performance guarantee to Theorem 3.7. However, the convergence rate in Theorem 3.9 is determined by the user-specified accuracy  $\epsilon$  instead of the gap constant  $\alpha$ .

## 4 NUMERICAL EXPERIMENT

In this section, we conduct numerical experiments to demonstrate the performance of our exploration and data outsourcing mechanisms. We compare the performance of our algorithms to random start. We also run sensitivity analysis to demonstrate the robustness of our algorithms with respect to two key hyper-parameters: the outsourcing sample size  $n$  and the inverse temperature  $\beta$ .

#### 4.1 Classic Nonconvex Test Function

We first consider a classic nonconvex optimization problem – the Styblinski-Tang function (ST-function) [15]. A  $d$ -dimensional ST-function is defined as

$$F(\theta) = \frac{\sum_{i=1}^d [\theta]_i^4 - 16[\theta]_i^2 + 5[\theta]_i}{2d}, \quad -5 \leq [\theta]_i \leq 5,$$

where  $[\theta]_i$  denotes the  $i$ -th coordinate of  $\theta$ . Note that ST-function is additively separable. By the first-order optimality condition, the stationary point set of  $F(\theta)$  is  $\{\theta \in \mathbb{R}^d : 4[\theta]_i^3 - 32[\theta]_i + 5 = 0, \forall i \in [d]\}$ . Moreover, the unique global minimum of ST-function is  $\theta^* \approx (-2.903, \dots, -2.903)$  and the corresponding objective value is  $-39.165$ . In this numerical experiment, we set  $d = 5$  and use gradient descent (GD) as  $\mathcal{T}$ . We apply OIPS-annealing to generate the initial points. Since the ST-function is not defined through expectation, we do not consider data outsourcing here, i.e.,  $\hat{F}_n(\theta) = F(\theta)$ . We test different inverse temperatures  $\beta = 1, 4$ , and  $10$ . For each  $\beta$ , we use importance sampling to draw i.i.d. samples from the target distribution  $\pi_\beta(\theta) \propto \exp\{-\beta F(\theta)\}$  exactly. For GD in the optimization phase, we use a step-size  $0.05$  and run  $50$  iterations. We pick the objective value at the last iteration as the convergent value. As the benchmark, we sample the initial point uniformly at random from the cubic  $[-5, 5]^d$  (random start). Finally, for each setting, we repeat the procedure  $500$  times and record the final convergent values. Figure 1 shows the distribution of convergent function values when the initial points are drawn from OIPS-annealing algorithm with different values of  $\beta$  versus the benchmark method. Note that compared with random start, initial points obtained by OIPS-annealing typically lead to smaller objective values. Moreover, as the inverse temperature  $\beta$  increases, the performance of OIPS-annealing algorithm further improves.

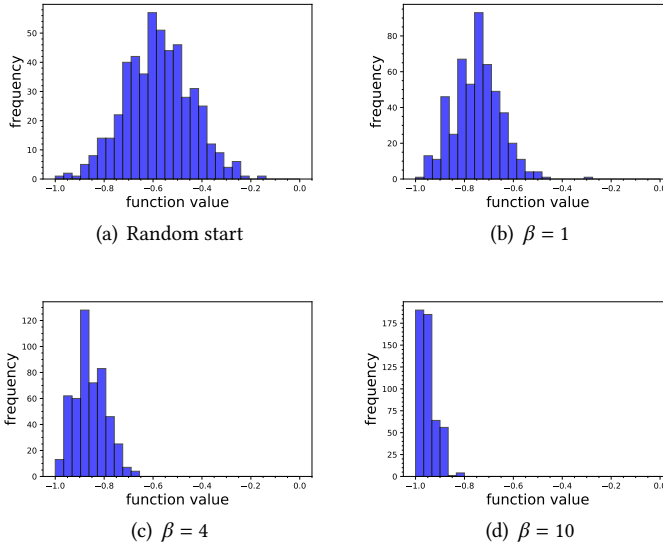


Fig. 1. Histogram of convergent function values (divided by 39.165) of ST-function

## 4.2 Gaussian Mixture Density

We study the problem of finding the largest mode of a Gaussian mixture density using kernel density estimation. In particular, the objective function

$$F(\theta) = \mathbb{E}_{X \sim \xi} \left[ (2\pi\sigma)^{-d/2} \cdot \exp \left\{ -\frac{\|\theta - X\|^2}{2\sigma^2} \right\} \right].$$

We assume  $\xi$  is a Gaussian mixture distribution, that is,  $X \sim \mathcal{N}(m_i, \sigma^2 I_d)$  with probability  $p_i$  for  $1 \leq i \leq M$ , where  $\mathcal{N}(m_i, \sigma^2 I_d)$  denotes the Gaussian distribution with mean vector  $m_i$  and covariance matrix  $\sigma^2 I_d$ ; the mixing weights  $p_i$  satisfy  $0 < p_i < 1$  and  $\sum_{i=1}^M p_i = 1$ . When  $m_i$ 's are well-separated,  $F(\theta)$  has multiple local minima located near  $m_i$ . Hence, the selection of initial point is critical to optimize  $F(\theta)$ .

We first consider a lower dimensional example with  $d = 5$  and  $M = 10$ . We implement SIPS, OIPS-annealing, and OIPS-SAO, all with  $n = 50$  and  $\beta = 10$ . ULA is used to draw  $L = 1000$  samples from  $\pi_\beta(\theta)$ . Given the initial point, GD is used to optimize  $F(\theta)$ . Moreover, to evaluate the gradient, we draw a batch  $(X_1, \dots, X_{1000})$  from  $\xi$  and approximate  $\nabla F(\theta)$  via batch means. In the optimization phase, GD is run for 20 iterations and the objective value at the last iteration is taken as the convergent value. Again, 500 independent replications of the algorithm are implemented in each setting. Figure 2 shows the distribution of convergent function values under different algorithms (number in bracket: success probability  $\mathbb{P}(\mathcal{F}_0^c)$ ). We observe that SIPS and OIPS outperform random start significantly. SIPS and OIPS-SAO perform better than OIPS-annealing with SIPS performs the best as measured by the probability of convergent function values smaller than  $-32$ . However, OIPS-annealing is the easiest and cheapest to implement in practice. Figure 3 further illustrates the success probability for different values of  $n$  in OIPS-annealing and SIPS. We observe that there is a diminishing return in the outsourcing sample size. The sample sizes that are larger than 50 in OIPS-annealing or even 30 in SIPS lead to similar performances.

We also consider a higher dimensional example with  $d = 30$  and  $M = 20$ . We focus on OIPS-annealing versus random start because of the relatively low computational cost of OIPS-annealing. We adopt the same hyper parameters as above. Figure 4 presents the results. We observe again that OIPS outperforms random start significantly.

## 4.3 Generalized multinomial logit model

We study an application of our algorithms for maximum likelihood estimation of the generalized multinomial logit (GMNL) model. Multinomial logit model is a classic model to study consumer choice. As an extension, the GMNL model accommodates the scaling heterogeneity in utility coefficients through an individual-specific scaling factor [11]. Such a generalization makes the negative log-likelihood function nonconvex. In practice, GD or BFGS with random starts are employed for the estimation [31].

Suppose that there are  $N$  customers who make a choice from  $J$  alternatives. The utility that customer  $n$  chooses alternative  $j$  is  $U_{nj} = x_j^\top \phi_n + \epsilon_{nj}$ , where  $x_j$  is a  $p$ -dimensional vector of attributes of product  $j$ ,  $\phi_n \in \mathbb{R}^p$  is the vector of utility coefficients, and  $\epsilon_{nj}$  is an idiosyncratic error term that follows standard Gumbel distribution. The customer tends to choose products with higher utilities and the probability that  $k$  is chosen is  $P_{nk} = \exp(x_k^\top \phi_n) / \sum_{j=1}^J \exp(x_j^\top \phi_n)$ . The GMNL model specifies  $\phi_n$  as  $\phi_n = \exp\{z_n^\top \psi + \xi_n\} \cdot \phi$ , where  $z_n$  is a  $q$ -dimensional vector of agent characteristics,  $\psi$  is a  $q$ -dimensional heterogeneity coefficient, and  $\xi_n$  is an independent random shock that follows the standard Gaussian distribution. Let the binary variable  $y_{nj} \in \{0, 1\}$  denote whether customer  $n$

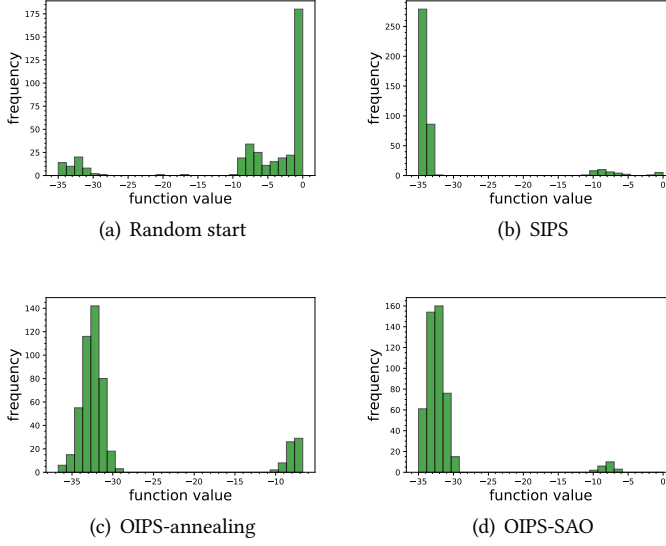


Fig. 2. Histogram of convergent function values of mixture Gaussian density ( $d = 5, M = 10$ ).

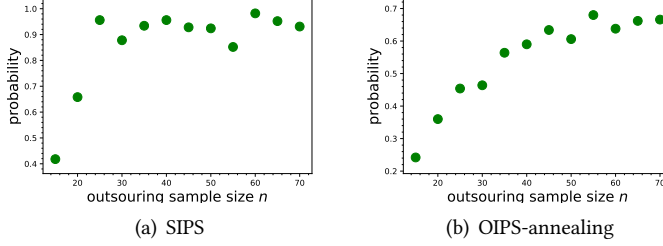


Fig. 3. Probability of finding the global minimum for different outsourcing sample sizes in OIPS-annealing and SIPS

chooses product  $j$ . Then the likelihood of customer  $n$ 's choice is

$$L_n = \mathbb{E}_{\xi_n \sim \mathcal{N}(0,1)} \left[ \prod_{k=1}^J \left( \frac{\exp(x_k^\top \phi_n)}{\sum_{j=1}^J \exp(x_j^\top \phi_n)} \right)^{y_{nk}} \right].$$

We use simulation to approximate the above expectation. The model parameter  $\theta = (\phi, \psi)$  can be estimated by maximizing the simulated negative log-likelihood function

$$F(\theta) = -\frac{1}{N} \sum_{n=1}^N \log \left( \frac{1}{R} \sum_{r=1}^R \prod_{k=1}^J \left( \frac{\exp(x_k^\top \phi_n^{[r]})}{\sum_{j=1}^J \exp(x_j^\top \phi_n^{[r]})} \right)^{y_{nk}} \right), \quad (8)$$

where  $\phi_n^{[r]} = \exp\{z_n^\top \psi + \xi_n^{[r]}\} \cdot \phi$  is the  $r$ -th draw from the distribution of  $\phi_n$  and  $R$  is the total number of draws.

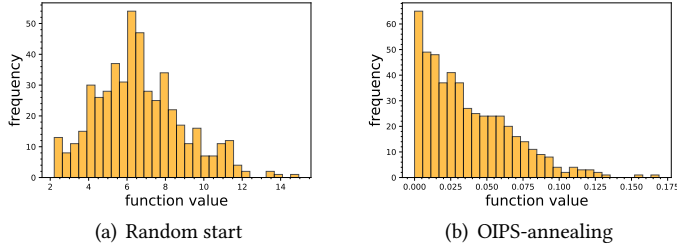


Fig. 4. Histogram of convergent function values of mixture Gaussian density ( $d = 30, M = 20$ ).

In our simulation experiment, we consider an instance with dimension parameters  $p = 10, q = 5$ , and  $J = 5$  alternatives. We generate  $N = 1000$  customers. In particular, we set the true parameter  $\phi^* = (1, \dots, 1, -1, \dots, -1)$  and  $\psi^* = (1, \dots, 1)$  and generate product attributes  $x_j$  and agent characteristics  $z_n$  from standard Gaussian distribution. Then we simulate the agents' choices following the GMNL model and obtain choice data  $y_{nj}$ . Based on the simulated dataset  $\{x_j, z_n, y_{nj}\}_{1 \leq n \leq N, 1 \leq j \leq J}$ , we use gradient descent to optimize the negative log-likelihood (8) with  $R = 100$ . We compare the performance of OIPS-annealing algorithm with random start. For OIPS-annealing, we set the outsourcing sample size  $n = 200$  and the inverse temperature  $\beta = 1$ . ULA is applied to draw  $L = 500$  samples as candidate initial points. In the optimization phase, GD is run for 100 iterations.

Figure 5 shows the distribution of convergent objective values (negative log-likelihood). We note that SIPS, OIPS-SAO, and OIPS-annealing again outperform the random start significantly. Moreover, SIPS and OIPS-SAO are performing slightly better than OIPS-annealing.

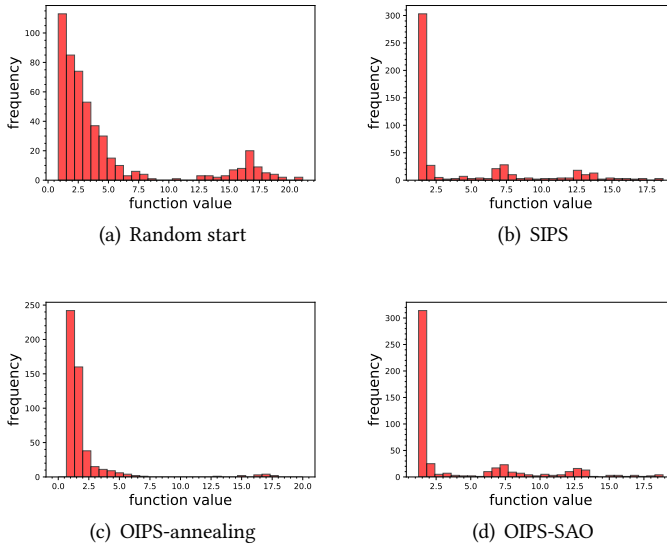


Fig. 5. Histogram of negative log-likelihood of GMNL model

## 5 CONCLUSION, LIMITATIONS, FUTURE WORKS

We have designed three algorithms using outsourced data to find good initial points. They are better than the popular random start approach. In both theoretical analysis and numerical tests, the SIPS and OIPS-SAO perform better than the OIPS-annealing, but they have computational costs in general.

Our work has the following two limitations, which can be seen as possible future directions.

1) We assume the outsourced data is drawn randomly from the true population. In practice, such data might be from a biased distribution or need additional privacy encryption. 2) Our analyses focused on the large  $\beta$  scenario. In practice, we would prefer to use a moderate  $\beta$  due to sampling complexity.

## REFERENCES

- [1] Zeyuan Allen-Zhu. 2018. Natasha 2: Faster Non-convex Optimization Than SGD. In *Advances in Neural Information Processing Systems*.
- [2] Jordan T Ash and Ryan P Adams. 2020. On Warm-Starting Neural Network Training. In *34th Conference on Neural Information Processing Systems*.
- [3] Léon Bottou, Frank E Curtis, and Jorge Nocedal. 2018. Optimization methods for large-scale machine learning. *Siam Review* 60, 2 (2018), 223–311.
- [4] Xi Chen, Simon S Du, and Xin T Tong. 2020. On Stationary-Point Hitting Time and Ergodicity of Stochastic Gradient Langevin Dynamics. *Journal of Machine Learning Research* 21, 68 (2020), 1–41.
- [5] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. 2019. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming* 176, 1 (2019), 5–37.
- [6] Sabrina De Capitani Di Vimercati, Sara Foresti, Sushil Jajodia, Stefano Paraboschi, and Pierangela Samarati. 2007. A data outsourcing architecture combining cryptography and access control. In *Proceedings of the 2007 ACM workshop on Computer security architecture*. 63–69.
- [7] Jing Dong and Xin T Tong. 2020. Spectral Gap of Replica Exchange Langevin Diffusion on Mixture Distributions. *arXiv preprint arXiv:2006.16193* (2020).
- [8] Jing Dong and Xin T Tong. 2021. Replica exchange for non-convex optimization. *Journal of Machine Learning Research* 22, 173 (2021), 1–59.
- [9] Alain Durmus, Eric Moulines, et al. 2017. Nonasymptotic Convergence Analysis for the Unadjusted Langevin Algorithm. *The Annals of Applied Probability* 27, 3 (2017), 1551–1587.
- [10] Raaz Dwivedi, Yuansi Chen, Martin J Wainwright, and Bin Yu. 2018. Log-concave sampling: Metropolis-Hastings algorithms are fast!. In *Conference on Learning Theory*. PMLR, 793–797.
- [11] Deniz G Fiebig, Michael P Keane, Jordan Louviere, and Nada Wasi. 2010. The generalized multinomial logit model: accounting for scale and coefficient heterogeneity. *Marketing Science* 29, 3 (2010), 393–421.
- [12] Sara Foresti. 2010. *Preserving privacy in data outsourcing*. Vol. 99. Springer Science & Business Media.
- [13] Rong Ge, Holden Lee, and Andrej Risteski. 2018. Simulated tempering Langevin Monte Carlo II: An improved proof using soft Markov chain decomposition. *arXiv preprint arXiv:1812.00793* (2018).
- [14] Saeed Ghadimi and Guanghui Lan. 2016. Accelerated Gradient Methods for Nonconvex Nonlinear and Stochastic Programming. *Mathematical Programming* 156, 1-2 (2016), 59–99.
- [15] Igor Grigoryev and Svetlana Mustafina. 2016. Global optimization of functions of several variables using parallel technologies. *International Journal of Pure and Applied Mathematics* 106, 1 (2016), 301–306.
- [16] Boris Hanin and David Rolnick. 2018. How to start training: The effect of initialization and architecture. In *32nd Conference on Neural Information Processing Systems*.
- [17] Rie Johnson and Tong Zhang. 2013. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems* 26 (2013), 315–323.
- [18] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. 2019. SCAFFOLD: Stochastic Controlled Averaging for On-Device Federated Learning. (2019).
- [19] Scott Kirkpatrick, C Daniel Gelatt, and Mario P Vecchi. 1983. Optimization by simulated annealing. *science* 220, 4598 (1983), 671–680.
- [20] Holden Lee, Andrej Risteski, and Rong Ge. 2018. Beyond log-concavity: Provable guarantees for sampling multi-modal distributions using simulated tempering langevin monte carlo. In *NeurIPS*.
- [21] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. 2020. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine* 37, 3 (2020), 50–60.



- [22] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. 2019. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189* (2019).
- [23] Yue M Lu and Gen Li. 2017. Spectral initialization for nonconvex estimation: high-dimensional limit and phase transitions. In *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 3015–3019.
- [24] Y. Ma, Y. Chen, C. Jin, N. Flammarin, and M. I. Jordan. 2019. Sampling can be faster than optimization. *Proceedings of the National Academy of Sciences* 116, 42 (2019), 20881–20885.
- [25] Song Mei, Yu Bai, Andrea Montanari, et al. 2018. The landscape of empirical risk for nonconvex losses. *Annals of Statistics* 46, 6A (2018), 2747–2774.
- [26] Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. 2017. Non-convex Learning via Stochastic Gradient Langevin Dynamics: A Nonasymptotic Analysis. In *Proceedings of the Conference on Learning Theory*.
- [27] Gareth O Roberts, Richard L Tweedie, et al. 1996. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli* 2, 4 (1996), 341–363.
- [28] Pierangela Samarati and Sabrina De Capitani Di Vimercati. 2010. Data protection in outsourcing scenarios: Issues and directions. In *Proceedings of the 5th ACM Symposium on Information, Computer and Communications Security*. 1–14.
- [29] Mark Schmidt, Nicolas Le Roux, and Francis Bach. 2017. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming* 162, 1-2 (2017), 83–112.
- [30] Nicholas G Tawn, Gareth O Roberts, and Jeffrey S Rosenthal. 2020. Weight-preserving simulated tempering. *Statistics and Computing* 30, 1 (2020), 27–41.
- [31] Kenneth E Train. 2009. *Discrete choice methods with simulation*. Cambridge university press.
- [32] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. 2017. Stochastic quasi-Newton methods for nonconvex stochastic optimization. *SIAM Journal on Optimization* 27, 2 (2017), 927–956.
- [33] Dawn B Woodard, Scott C Schmidler, Mark Huber, et al. 2009. Conditions for rapid mixing of parallel and simulated tempering on multimodal distributions. *The Annals of Applied Probability* 19, 2 (2009), 617–640.
- [34] Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. 2018. Global Convergence of Langevin Dynamics Based Algorithms for Nonconvex Optimization. In *Advances in Neural Information Processing Systems*.
- [35] Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. 2019. Fixup initialization: Residual learning without normalization. In *Seventh International Conference on Learning Representations*.
- [36] Xinwei Zhang, Mingyi Hong, Sairaj Dhople, Wotao Yin, and Yang Liu. 2021. FedPD: A Federated Learning Framework With Adaptivity to Non-IID Data. *IEEE Transactions on Signal Processing* 69 (2021), 6055–6070.

## A TECHNICAL VERIFICATIONS

### A.1 Approximation accuracy of $\hat{F}(\theta)$ and data complexity

PROOF OF LEMMA 3.2. To prove the results about gradient and Hessian convergence, we can apply Theorem 1 in [25] directly. Specifically, under Assumptions 1 and 4, when  $n > Cd \log(d)$ , with probability at least  $1 - \rho$ , we have

$$\begin{aligned} \sup_{\theta \in \Theta} \|\nabla F(\theta) - \nabla \hat{F}_n(\theta)\| &\leq \tau \sqrt{\frac{Cd \log(n)}{n}}, \\ \sup_{\theta \in \Theta} \|\nabla^2 F(\theta) - \nabla^2 \hat{F}_n(\theta)\|_{\text{op}} &\leq \tau^2 \sqrt{\frac{Cd \log(n)}{n}}. \end{aligned} \quad (9)$$

For the stationary points convergence, based on Theorem 2 in [25], under Assumptions 1 and 4, when  $n \geq 4Cd \log(n) \cdot ((\tau^2/\sigma^2) \vee (\tau^4/\eta^2))$ , the empirical loss function  $\hat{F}_n(\theta)$  is  $(\sigma/2, \eta/2)$ -strongly Morse and possesses  $K + 1$  stationary points with probability at least  $1 - \rho$ . Furthermore, there is a one-to-one correspondence between  $(\theta_0^*, \dots, \theta_K^*)$ , the stationary points of  $F(\theta)$ , and  $(\hat{\theta}_0^*, \dots, \hat{\theta}_K^*)$ , the stationary points of  $\hat{F}_n(\theta)$ . Moreover, when  $n \geq 4Cd \log(n)/\eta_*^2$ ,

$$\max_{0 \leq i \leq K} \|\theta_i^* - \hat{\theta}_i^*\| \leq \frac{2\tau}{\eta} \sqrt{\frac{Cd \log(n)}{n}}. \quad (10)$$

It remains to establish the uniform convergence result for  $\hat{F}_n$ . Although it is not directly available in [25], the proof follows a similar idea. For self-completeness, we provide the details here.

First of all, given the parameter space  $\Theta$ , let  $\Theta_\varepsilon := \{\theta_1, \dots, \theta_J\}$  be a  $\varepsilon$ -covering net. In other words, for arbitrary  $\theta \in \Theta$ , there exists certain  $\theta_{j(\theta)} \in \Theta_\varepsilon$  such that  $\|\theta - \theta_{j(\theta)}\| \leq \varepsilon$ . Thus, for any  $\theta \in \Theta$ , we have

$$|\hat{F}_n(\theta) - F(\theta)| \leq |\hat{F}_n(\theta) - \hat{F}_n(\theta_{j(\theta)})| + |\hat{F}_n(\theta_{j(\theta)}) - F(\theta_{j(\theta)})| + |F(\theta) - F(\theta_{j(\theta)})|. \quad (11)$$

For any  $t > 0$ , we denote by

$$A_t = \left\{ \sup_{\theta \in \Theta} |\hat{F}_n(\theta) - \hat{F}_n(\theta_{j(\theta)})| \geq t/3 \right\}, \quad B_t = \left\{ \sup_{\theta_j \in \Theta_\varepsilon} |\hat{F}_n(\theta_j) - F(\theta_j)| \geq t/3 \right\},$$

$$\text{and } C_t = \left\{ \sup_{\theta \in \Theta} |F(\theta) - F(\theta_{j(\theta)})| \geq t/3 \right\}.$$

Then we have

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |\hat{F}_n(\theta) - F(\theta)| \geq t\right) \leq \mathbb{P}(A_t) + \mathbb{P}(B_t) + \mathbb{P}(C_t).$$

In the next, we upper bound the three parts in above inequality respectively. For the last part, we have

$$|F(\theta) - F(\theta_{j(\theta)})| \leq \sup_{\theta \in \Theta} \|\nabla F(\theta)\| \cdot \|\theta - \theta_{j(\theta)}\| \leq L^* \cdot \varepsilon.$$

Hence, when  $t \geq 3\varepsilon L^*$ , the deterministic event  $C_t$  would never happen and  $P(C_t) = 0$ . For the second part, under Assumption 4, by applying the union bound and the sub-Gaussian concentration inequality, we have

$$\begin{aligned} \mathbb{P}(B_t) &\leq |\Theta_\varepsilon| \cdot P\left(|\hat{F}_n(\theta_j) - F(\theta_j)| \geq t/3\right) \\ &\leq |\Theta_\varepsilon| \cdot \exp\left\{-nt^2/(18\tau^2)\right\} \leq (2/\varepsilon)^d \cdot \exp\left\{-nt^2/(18\tau^2)\right\}. \end{aligned}$$

Thus, when

$$t > 5\tau \cdot \sqrt{\frac{\log(2/\rho) + d \log(2/\varepsilon)}{n}},$$

we have  $\mathbb{P}(B_t) \leq \rho/2$ . For the first part, by Markov inequality, we have

$$\mathbb{P}(A_t) \leq \frac{3\mathbb{E}\left[\sup_{\theta \in \Theta} |\hat{F}_n(\theta) - \hat{F}_n(\theta_{j(\theta)})|\right]}{t} \leq \frac{3\varepsilon \cdot \mathbb{E}\left[\sup_{\theta \in \Theta} \|\nabla \hat{F}_n(\theta)\|\right]}{t}.$$

By Assumption 1, we have

$$\mathbb{E}\left[\sup_{\theta \in \Theta} \|\nabla \hat{F}_n(\theta)\|\right] \leq \mathbb{E}\left[\sup_{\theta \in \Theta} \|\nabla \hat{F}_n(\theta) - \nabla \hat{F}_n(\theta^*)\|\right] + \mathbb{E}\left[\|\nabla \hat{F}_n(\theta^*)\|\right] \leq 2J^* + H,$$

which implies that

$$\mathbb{P}(A_t) \leq 3\varepsilon(2J^* + H)/t.$$

Taking  $t \geq 6\varepsilon(2J^* + H)/\rho$ , we have  $P(A_t) \leq \rho/2$ .

Finally, by taking

$$\varepsilon^* = \rho\tau/(6dn(2J^* + H)), \quad t^* = 5\tau\sqrt{(\log(2/\rho) + d \log(2/\varepsilon))/n},$$

and utilizing the fact that  $H \leq \tau^2 d^{c_h}$ ,  $J_* \leq \tau^3 d^{c_h}$ , when  $n \geq Cd \log(d)$ , we have

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |\hat{F}_n(\theta) - F(\theta)| \geq \tau\sqrt{\frac{Cd \log(n)}{n}}\right) \leq \rho.$$

Now, given an approximation accuracy  $\delta$ , we calculate the minimal required sample size. For arbitrary positive constant  $\iota$ , there exists an absolute constant  $C_\iota$  such that  $\log(n) \leq C_\iota \cdot n^\iota$ . As a result, when

$$n \geq \max \left\{ \left[ \frac{C_\iota C d}{(\delta / ((2\tau/\eta) \vee \tau \vee \tau^2))^2} \right]^{\frac{1}{1-\iota}}, 4Cd(\log(d) \vee \log(n)/\eta_*^2) \right\},$$

we have

$$\begin{aligned} \sup_{\theta \in \Theta} |F(\theta) - \hat{F}_n(\theta)| &\leq \delta, \quad \sup_{\theta \in \Theta} \|\nabla F(\theta) - \nabla \hat{F}_n(\theta)\| \leq \delta, \\ \sup_{\theta \in \Theta} \|\nabla^2 F(\theta) - \nabla^2 \hat{F}_n(\theta)\|_{\text{op}} &\leq \delta, \quad \text{and} \quad \max_{0 \leq i \leq K} \|\theta_i^* - \hat{\theta}_i^*\| \leq \delta. \end{aligned}$$

with probability at least  $1 - \rho$ . □

## A.2 Performance analysis of the sampling approach

PROOF OF PROPOSITION 3.3. First note that when  $\hat{F}_n(\theta)$  is a  $\delta$ -approximation, for any  $\theta \notin \mathcal{B}_r(\theta_0^*)$ , by Assumption 3 we have

$$\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*) \geq (F(\theta) - \delta) - (F(\theta_0^*) + \delta) \geq \alpha - 2\delta > 0.$$

Hence, by the definition of  $\pi_\beta$ , we have

$$\begin{aligned} \mathbb{P}(\tilde{\theta}_\beta \in \mathcal{B}_r(\theta_0^*)) &= \frac{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta \hat{F}_n(\theta)) d\theta}{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta \hat{F}_n(\theta)) d\theta + \int_{\Theta/\mathcal{B}_r(\theta_0^*)} \exp(-\beta \hat{F}_n(\theta)) d\theta} \\ &= \frac{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta}{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta + \int_{\Theta/\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta} \\ &\geq \frac{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta}{\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta + \exp(-\beta(\alpha - 2\delta)) \cdot \text{Vol}(\Theta/\mathcal{B}_r(\theta_0^*))}, \end{aligned}$$

where  $\text{Vol}(\Theta/\mathcal{B}_r(\theta_0^*))$  denotes the volume of set  $\Theta/\mathcal{B}_r(\theta_0^*)$ .

On the other hand, based on the regularity condition of Hessian and the definition of  $\delta$ -approximation, we have  $\|\nabla^2 \hat{F}_n(\theta)\|_{\text{op}} \leq H + L^* + \delta$ . As a result, for any  $\theta \in \mathcal{B}_r(\theta_0^*)$ ,

$$\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*) \leq 2(H + L^* + \delta) \cdot (\|\theta - \theta_0^*\|^2).$$

Hence,

$$\begin{aligned} \int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta &\geq \int_{\mathcal{B}_r(\theta_0^*)} \exp(-2\beta(H + L^* + \delta)\|\theta - \theta_0^*\|^2) d\theta \\ &= \left( \pi\beta^{-1}/(H + L^* + \delta) \cdot (\Psi(2r\sqrt{\beta(H + L^* + \delta)}) - 1/2) \right)^d, \end{aligned}$$

where  $\Psi(\cdot)$  denotes the CDF of standard normal distribution. Note that when  $\beta \geq \Omega(r^{-2})$ ,

$$\Psi(2r\sqrt{\beta(H + L^* + \delta)}) - 1/2 = O(1).$$

Then,

$$\int_{\mathcal{B}_r(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\hat{\theta}_0^*)]) d\theta = O\left((\beta^{-1}/(H + L^* + \delta))^d\right).$$

As a result,

$$1/\mathbb{P}(\tilde{\theta}_\beta \in \mathcal{B}_r(\theta_0^*)) = O\left(1 + \exp\{-\beta(\alpha - 2\delta)\} / (\beta^{-1}/(H + L^* + \delta))^d\right),$$

which further implies that

$$1 - \mathbb{P}(\tilde{\theta}_\beta \in \mathcal{B}_r(\theta_0^*)) = O\left(\exp\{-\beta(\alpha - 2\delta)\} / (\beta^{-1}/(H + L^* + \delta))^d\right).$$

Finally, by setting  $\delta = \alpha/4$ , we obtain the result.  $\square$

PROOF OF LEMMA 3.4. Let  $X_1, \dots, X_L$  be samples from  $\hat{\mathcal{M}}$ .

$$\begin{aligned} \mathbb{P}(X_1 \notin B, \dots, X_L \notin B) &= \mathbb{E}\left[\prod_{i=1}^L 1_{(X_i \notin B)}\right] \\ &= \mathbb{E}\left[\prod_{i=1}^{L-1} 1_{(X_i \notin B)} \cdot \mathbb{E}_{L-1}[1_{X_L \notin B}]\right] \\ &= \mathbb{E}\left[\prod_{i=1}^{L-1} 1_{(X_i \notin B)} \cdot \hat{\pi}_{X_{L-1}}(B^c)\right] \\ &\leq \mathbb{E}\left[\prod_{i=1}^{L-1} 1_{(X_i \notin B)} \cdot (\pi_\beta(B^c) + \delta_\beta)\right] \\ &= (\pi_\beta(B^c) + \delta_\beta) \cdot \mathbb{P}(X_1 \notin B, \dots, X_{L-1} \notin B). \end{aligned}$$

By induction, we have

$$\mathbb{P}(X_1 \notin B, \dots, X_L \notin B) \leq (\pi_\beta(B^c) + \delta_\beta)^L.$$

$\square$

PROOF OF THEOREM 3.5. We use  $\mathcal{I}_n(\delta)$  to denote the random event that  $\hat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$ . First, based on Proposition 3.3,  $\mathbb{P}(\mathcal{I}_n^c(\delta)) \leq \rho$  for  $n \geq n(\delta, \rho, d)$ . Then,

$$\mathbb{P}(\mathcal{F}_0) - \rho \leq \mathbb{P}(\mathcal{F}_0 \cap \mathcal{I}_n(\delta)) \leq \mathbb{P}(\mathcal{F}_0 | \mathcal{I}_n(\delta)).$$

By the definition of  $\delta$ -approximation, conditional on  $\mathcal{I}(\delta)$ , if at least one of  $(\theta_1, \dots, \theta_L)$  falls into  $\mathcal{B}_r(\theta_0^*)$ ,  $\mathcal{F}_0$  would not happen. Hence, by Lemma 3.4, we have

$$\log(\mathbb{P}(\mathcal{F}_0 | \mathcal{I}_n(\delta))) \leq L \cdot \log(\pi_\beta(\mathcal{B}_r^c(\theta_0^*)) + \delta_\beta) \leq L \cdot \max\{\log(2\pi_\beta(\mathcal{B}_r^c(\theta_0^*))), \log(2\delta_\beta)\}.$$

Finally, based on Lemma 3.3, when  $\beta \geq \Omega(r^{-2})$ ,

$$\log(\pi_\beta(\mathcal{B}_r^c(\theta_0^*))) = O(-\beta(\alpha - 2\delta) - d \log \beta).$$

If we set  $\delta = \alpha/4$ , the above upper bound leads to

$$\mathbb{P}(\mathcal{F}_0) \leq \rho + \exp(-CL \cdot \max\{-\beta\alpha/2 + d \log(\beta), \log(2\delta_\beta)\}).$$

for some constant  $C > 0$ , and we finish the proof.  $\square$

### A.3 Performance analysis of the optimization approaches

In this section, we establish the performance guarantee of the optimization approach, i.e., Algorithm 2. We first analyze the sample selection rule with the SAO approach, which set  $\theta_*^0 = \hat{\mathcal{T}}(\theta_{i^*}^0)$  where  $i^* = \operatorname{argmin}_{1 \leq i \leq L} \{\hat{F}_n(\hat{\mathcal{T}}(\theta_i^0))\}$ .

**PROOF OF THEOREM 3.6.** Under Assumption 2,  $F(\theta)$  is  $\mu$ -strongly convex in  $\mathcal{B}_r(\theta_0^*)$ . When  $\delta < \mu$  and  $\hat{F}_n(\theta)$  is a  $\delta$ -approximation, we have

$$\sup_{\theta \in \Theta} \|\nabla^2 F(\theta) - \nabla^2 \hat{F}_n(\theta)\|_{\text{op}} \leq \delta \text{ and } \|\hat{\theta}_0^* - \theta_0^*\| \leq \delta.$$

This implies that  $\hat{F}_n(\theta)$  is  $(\mu - \delta)$ -strongly convex in  $\mathcal{B}_r(\theta_0^*)$  and  $\hat{\theta}_0^*$  is the unique minimum of  $\hat{F}_n(\theta)$  in  $\mathcal{B}_r(\theta_0^*)$ . Hence, starting from any  $\theta \in \mathcal{B}_r(\theta_0^*)$ , the optimization algorithm  $\hat{\mathcal{T}}$  can converge to  $\hat{\theta}_0^*$ , which implies that

$$\mathbb{P}(\hat{\mathcal{T}}(\tilde{\theta}_\beta) \neq \hat{\theta}_0^*) \leq \mathbb{P}(\tilde{\theta}_\beta \notin \mathcal{B}_r(\theta_0^*)).$$

Then by Proposition 3.3, we can establish the upper bound for the probability of  $\hat{\mathcal{T}}(\tilde{\theta}_\beta) \neq \hat{\theta}_0^*$ . Finally, note that when at least one of  $(\theta_1, \theta_2, \dots, \theta_L)$  falls into  $\mathcal{B}_r(\theta_0^*)$ , the global minimum of  $\hat{F}_n(\theta)$ ,  $\hat{\theta}_0^* \in \mathcal{B}_r(\theta_0^*)$ , can be found by  $\hat{\mathcal{T}}$  and would also be selected as the initial point to optimize  $F(\theta)$ . As a result,  $\mathcal{F}_1 \subset \mathcal{F}_0$ . Thus, by Theorem 3.5, we prove the upper bound for  $\mathbb{P}(\mathcal{F}_1)$ .  $\square$

**PROOF OF THEOREM 3.7.** We first show that if at least one point of  $(\theta_1, \dots, \theta_L)$  is in  $\mathcal{B}_{r_0}(\theta_0^*)$ , the annealing approach would select a point that falls into  $\mathcal{B}_{r_0}(\theta_0^*)$ . When  $\hat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$ , if  $\theta_i \in \mathcal{B}_{r_0}(\theta_0^*)$ , we have

$$\begin{aligned} \hat{F}_n(\theta_i) &\leq F(\theta_i) + \delta \leq F(\theta_0^*) + \delta + \frac{1}{2}r_0^2 \cdot \sup_{\theta \in \Theta} \|\nabla^2 F(\theta)\|_{\text{op}} \\ &\leq F(\theta_0^*) + \delta + \frac{\alpha}{2} \\ &\leq \min_{j \neq i} F(\theta_j) - \frac{\alpha}{2} + \delta \text{ by Assumption 3} \\ &\leq \min_{j \neq i} \hat{F}_n(\theta_j) \text{ as } \delta \leq \alpha/4. \end{aligned}$$

Hence, if the algorithm selects some  $\theta_j \neq \theta_i$ , we must have  $\hat{F}_n(\theta_j) \leq \hat{F}_n(\theta_i)$ , which implies that  $\theta_j$  is in  $\mathcal{B}_{r_0}(\theta_0^*)$  as well. The remaining proof follows exactly the same line of argument as that of Theorem 3.5.  $\square$

### A.4 Performance analysis for extension to $\epsilon$ -Global Minimum

**PROOF OF THEOREM 3.9.** For Algorithm 2-annealing, note that if there is a sample  $\theta_i \in \mathcal{B}_{r_\epsilon}(\theta_i^*) \subseteq \mathcal{B}_{\epsilon, r_\epsilon}$ , then we have

$$\begin{aligned} \hat{F}_n(\theta_i) &\leq F(\theta_i) + \delta \\ &\leq F(\theta_i^*) + \sup_{\theta \in \Theta} \|\nabla^2 F(\theta)\|_{\text{op}} \cdot r_\epsilon^2 + \delta \leq F(\theta_0^*) + 2\epsilon + \delta. \end{aligned}$$

So if we select some  $\theta_j$  instead, then

$$F(\theta_j) \leq \hat{F}_n(\theta_j) + \delta \leq \hat{F}_n(\theta_i) + \delta \leq F(\theta_0^*) + 2\epsilon + 2\delta.$$

When  $\delta = \epsilon/4$ , by definition  $\theta_j$  is a  $3\epsilon$ -global minimum.

In the next, we estimate the probability that a sample drawn from  $\pi_\beta$  falls into  $\mathcal{B}_{\epsilon, r_\epsilon}$ . Note that for  $\theta \notin \mathcal{B}_{\epsilon, r_\epsilon}$ ,

$$\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*) \geq \epsilon - 2\delta.$$

Then we have

$$\begin{aligned}
\mathbb{P}(\tilde{\theta}_\beta \in \mathcal{B}_{\epsilon, r_\epsilon}) &= \frac{\int_{\mathcal{B}_{\epsilon, r_\epsilon}} \exp(-\beta \hat{F}_n(\theta)) d\theta}{\int_{\mathcal{B}_{\epsilon, r_\epsilon}} \exp(-\beta \hat{F}_n(\theta)) d\theta + \int_{\Theta \setminus \mathcal{B}_{\epsilon, r_\epsilon}} \exp(-\beta \hat{F}_n(\theta)) d\theta} \\
&\geq \frac{\int_{\mathcal{B}_{\epsilon, r_\epsilon}} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta}{\int_{\mathcal{B}_{\epsilon, r_\epsilon}} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta + \exp(-\beta(\epsilon - 2\delta)) \cdot \text{Vol}(\Theta \setminus \mathcal{B}_{\epsilon, r_\epsilon})} \\
&\geq \frac{\int_{\mathcal{B}_{r_\epsilon}(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta}{\int_{\mathcal{B}_{r_\epsilon}(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta + \exp(-\beta(\epsilon - 2\delta)) \cdot \text{Vol}(\Theta \setminus \mathcal{B}_{\epsilon, r_\epsilon})}.
\end{aligned}$$

Similar to the proof of Proposition 3.3, when  $\beta \geq \Omega(r_\epsilon^{-2})$ , we have

$$\int_{\mathcal{B}_{r_\epsilon}(\theta_0^*)} \exp(-\beta [\hat{F}_n(\theta) - \hat{F}_n(\theta_0^*)]) d\theta = O\left((\beta^{-1}/(H + L^* + \delta))^d\right).$$

As a result,

$$1/\mathbb{P}(\tilde{\theta}_\beta \in \mathcal{B}_{\epsilon, r_\epsilon}) = O\left(1 + \exp\{-\beta(\epsilon - 2\delta)\} / (\beta^{-1}/(H + L^* + \delta))^d\right),$$

which further implies that

$$\mathbb{P}(\tilde{\theta}_\beta \notin \mathcal{B}_{\epsilon, r_\epsilon}) = O\left(1 + \exp\{-\beta(\epsilon - 2\delta)\} / (\beta^{-1}/(H + L^* + \delta))^d\right).$$

When  $\delta = \epsilon/4$ , we have

$$\log(\mathbb{P}(\tilde{\theta}_\beta \notin \mathcal{B}_{\epsilon, r_\epsilon})) = O(-\beta\epsilon/2 + d \log(\beta)).$$

Hence, with probability at least  $O(\exp\{-\beta\epsilon/2 + d \log(\beta)\})$ , Algorithm 2 with subroutine 1 can find a  $3\epsilon$ -global minimum of  $F(\theta)$ .

We use  $\mathcal{I}_n(\delta)$  to denote the random event that  $\hat{F}_n(\theta)$  is a  $\delta$ -approximation of  $F(\theta)$ . Similar to the proof of Theorem 3.5, we have

$$\mathbb{P}(\mathcal{F}_{3\epsilon, 2}) - \rho \leq \mathbb{P}(\mathcal{F}_{3\epsilon, 2} \cap \mathcal{I}_n(\delta)) \leq \mathbb{P}(\mathcal{F}_{3\epsilon, 2} | \mathcal{I}_n(\delta)).$$

and

$$\log(\mathbb{P}(\mathcal{F}_{3\epsilon, 2} | \mathcal{I}_n(\delta))) \leq L \cdot \log(\pi_\beta(\mathcal{B}_{\epsilon, r_\epsilon}^c) + \delta_\beta) \leq L \cdot \max\{\log(2\pi_\beta(\mathcal{B}_{\epsilon, r_\epsilon}^c)), \log(2\delta_\beta)\}.$$

Finally, in above proof, by replacing  $\epsilon$  with  $\epsilon/3$ , we obtain the result.  $\square$