

Global Optimization with a Power-Transformed Objective and Gaussian Smoothing

Chen Xu¹

Abstract

We propose a novel method, namely Gaussian Smoothing with a Power-Transformed Objective (GS-PowerOpt), that solves global optimization problems in two steps: (1) perform a (exponential) power- N transformation to the not necessarily differentiable objective $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and get f_N , and (2) optimize the Gaussian-smoothed f_N with stochastic approximations. Under mild conditions on f , for any $\delta > 0$, we prove that with a sufficiently large power N_δ , this method converges to a solution in the δ -neighborhood of f 's global optimum point, at the iteration complexity of $O(d^4\epsilon^{-2})$. If we require that f is differentiable and further assume the Lipschitz condition on f and its gradient, the iteration complexity reduces to $O(d^2\epsilon^{-2})$, which is significantly faster than the standard homotopy method. In most of the experiments performed, our method produces better solutions than other algorithms that also apply the smoothing technique.

1. Introduction

In this work, we consider the global optimization problem of

$$\max_{\mathbf{x} \in \mathcal{S} \subset \mathbb{R}^d} f(\mathbf{x}), \quad (1)$$

where \mathcal{S} is a compact set and $f : \mathcal{S} \rightarrow \mathbb{R}$ is a continuous and possibly non-concave function with a *unique* global maximum point $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x})$. The minimization version of this problem is often encountered in machine learning, such as model trainings and adversarial attacks in computer vision. The gradient-based algorithms are commonly used, such as the (stochastic) gradient descent. In general, these methods only guarantee to approximate a lo-

cally optimal solution (Mertikopoulos et al., 2020; Lei et al., 2019; Choromanska et al., 2015).

Gaussian smoothing (e.g., Section 3 in (Nesterov & Spokoiny, 2017)) refers to convolving f with a Gaussian density to obtain a surrogate objective

$$\hat{f}_\sigma(\boldsymbol{\mu}) := \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)}[f(\mathbf{x})], \quad (2)$$

where \mathcal{N} denotes a multivariate Gaussian distribution, $\sigma > 0$ is called the scaling parameter, and I_d denotes a $d \times d$ identity matrix. The smoothing effect possibly eliminates certain local extremes of the objective, and this technique is applied by the widely used homotopy methods for global optimizations.

Homotopy (e.g., (Mobahi & Fisher, 2015)), also called the graduated continuation, is a class of methods that aim to find a global solution to (1), with many applications in machine learning (Xu et al., 2016; Iwakiri et al., 2022). It converts the original problem to

$$\max_{\boldsymbol{\mu} \in \mathbb{R}^d, \sigma \geq 0} \mathbb{E}_\xi[f(\boldsymbol{\mu} + \sigma \boldsymbol{\xi})], \quad (3)$$

where $\boldsymbol{\xi}$ is a random vector with a pre-selected distribution, such as a standard multivariate Gaussian distribution (Gaussian Homotopy, GH) or a uniform distribution in a unit sphere. Based on the observation that $\boldsymbol{\mu}_\sigma^* := \arg \max_{\boldsymbol{\mu}} \mathbb{E}[f(\boldsymbol{\mu} + \sigma \boldsymbol{\xi})]$ approaches¹ \mathbf{x}^* as σ decreases to 0, the standard homotopy method adopts a double-loop mechanism: the outer loop iteratively decreases σ , and for each fixed value of σ , the inner loop solves $\max_{\boldsymbol{\mu}} \mathbb{E}[f(\boldsymbol{\mu} + \sigma \boldsymbol{\xi})]$, with the solution found in the current inner loop as the starting search point in the next inner loop.

The double-loop mechanism is costly in time. To tackle this issue, (Iwakiri et al., 2022) proposes a single-loop Gaussian homotopy (SLGH) method that updates both $\boldsymbol{\mu}$ and σ in each iteration. Although SLGH aims at the global optimum, in theory, it only guarantees to approximate a local optimum², which is not necessarily a global one. A time-efficient algorithm that aims at the global maximum is still to be found.

¹Note that $\mathbb{E}[f(\boldsymbol{\mu} + \sigma \boldsymbol{\xi})] = f(\boldsymbol{\mu})$ if $\sigma = 0$.

²Theorem 4.1 in (Iwakiri et al., 2022) shows that SLGH approximates a solution $\hat{\mathbf{x}}$ such that $\mathbb{E}[\nabla f(\hat{\mathbf{x}})] = 0$.

¹Department of Engineering, Shenzhen MSU-BIT University, Shenzhen, P. R. China. Correspondence to: Chen Xu <xuchen@smbu.edu.cn>.

Therefore, in this work, we propose a new method, namely the Gaussian Smoothing with a Power-Transformed Objective (GS-PowerOpt), which adopts a single-loop mechanism and aims at the global optimum. Specifically, with $\sigma > 0$ and a sufficiently large $N > 0$, it solves the surrogate problem of $\max_{\mu} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \sigma^2 I_d)} [f_N(\mathbf{x})]$, where $f_N(\mathbf{x})$ modifies f by putting more weight to its global maximum point \mathbf{x}^* . The functional form of $f_N(\mathbf{x})$ can be either $f^N(\mathbf{x})$ or $e^{Nf(\mathbf{x})}$. According to our theory (see Subsection 3.1 for a preview), GS-PowerOpt converges to a neighborhood of the solution \mathbf{x}^* to (1) at the iteration complexity of $O(d^4 \varepsilon^{-2})$. This complexity can be improved to $O(d^2 \varepsilon^{-2})$ if the commonly assumed Lipschitz conditions hold. The majority of experiments in Section 5 show that the GS-PowerOpt-based algorithm (e.g., EPGS, introduced later) outperforms other algorithms that also apply the smoothing technique.

RELATED WORK

The homotopy methods, firstly proposed in (Blake & Zisserman, 1987), are intensively studied in the field of machine learning for global optimization problems. (Mobahi & Fisher, 2015) derives a bound for the worst scenario of the GH algorithm in a deterministic setting (i.e., the expectation \mathbb{E} is computed accurately), while (Hazan et al., 2016) provides a convergence analysis in a stochastic setting (i.e., \mathbb{E} is estimated with samples). (Gao & Sener, 2022) changes the distribution of the perturbation ξ from the commonly used Gaussian or uniform to the distribution that minimizes the estimation error of the gradient $\nabla_{\mu} \mathbb{E}_{\xi} [f(\mu + \xi)]$. (Lin et al., 2023) proposes an algorithm for learning the whole solution path produced by the homotopy. Specifically, their algorithm learns a model $x_{\phi}(\sigma)$ that predicts (for any $\sigma > 0$) the solution to $\min_{\mu} \nabla_{\mu} \mathbb{E}_{\xi \sim \mathcal{N}(0, I_d)} [f(\mu + \sigma \xi)]$, where ϕ is the set of model parameters to be trained.

The smoothing technique and the homotopy method have a large number of successful applications in machine learning, such as neural network training (Hazan et al., 2016), adversarial attack on image classification models (Iwakiri et al., 2022), solving L_1 -regularized least-square problems ((Xiao & Zhang, 2012)), neural combinatorial optimization (Gao & Sener, 2022), improving the optimization algorithms of stochastic gradient descent and Adam (Starnes & Webster, 2024), and so on.

There are a few existing studies (Dvijotham et al., 2014; Roulet et al., 2020; Chen et al., 2024) that replace the original f with a surrogate objective that also involves the exponential transformation $e^{Nf(\mu + \xi)}$ before smoothing. But their works are different from ours. (Dvijotham et al., 2014) proposes to minimize the surrogate objective of $G(\mu) := \frac{1}{N} \log (\mathbb{E}_{\xi \sim \mathcal{N}(0, \Sigma)} [e^{Nf(\mu + \xi)}]) + \frac{1}{2} \mu^T R \mu$. The theory³ that justifies this surrogate objective requires that

³According to Theorem 3.1 in (Dvijotham et al., 2014),

N, R , and Σ be selected so that $NR - \Sigma^{-1}$ is positive semi-definite. This indicates that our EPGS (see Section 2), for which $R = \mathbf{0}$ and $\Sigma = \sigma I_d$, is not a special case of theirs, since $-\sigma^{-1} I_d$ is negative definite and violates their requirement. Moreover, their theory on the distance between the optimum point of the new surrogate and \mathbf{x}^* is incomplete (see Section 3.2 in (Dvijotham et al., 2014)). For optimal-control problems, (Roulet et al., 2020) study the surrogate objective that is similar to $G(\mu)$, and provide a theoretical analysis on the corresponding algorithm’s convergence to a stationary point. However, the relation between this stationary point and the global optimum point \mathbf{x}^* is not revealed. The proposed surrogate objective in (Chen et al., 2024) is $(1 - N)f(\mu) + \log (\mathbb{E}_{\xi \sim \mathcal{N}(0, \Sigma)} [e^{Nf(\mu + \xi)}])$, where $N \in [0, 1]$. This is very different from our requirement that N is sufficiently large (see Theorem 2.1). Also, their theory (i.e., Theorem 10 in (Chen et al., 2024)) bounds $|f(\mathbf{x}^*) - f(\mu^*)|$ with $O(N\sigma^2) + G(N, \sigma)$ where $G(N, \sigma)$ is in general nonlinear in N , which does not imply an improvement for increasing the value of N .

As shown later in Eq. (5), GS-PowerOpt iteratively generates solution candidates and update its internal state (i.e., μ_t) using candidates’ evaluations. This pattern is shared by many evolutionary algorithms (EA), such as simulated annealing (Van Laarhoven et al., 1987), particle swarm optimization (PSO, e.g., (Miranda, 2018) and Section 3.1.5 in (Locatelli & Schoen, 2013)), the cross-entropy method for optimization (Boer et al., 2005), and the covariance matrix adaptation evolution strategy (CMA-ES, (Hansen & Ostermeier, 2001)). However, GS-PowerOpt is distinct from the EA algorithms in its construction of the surrogate objective $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \sigma^2 I_d)} [f_N(\mathbf{x})]$ and its use of the stochastic gradient method for optimization.

CONTRIBUTION

This paper introduces a novel zeroth-order method, GS-PowerOpt, for solving global optimization problems, with the contributions summarized as follows.

1. To our knowledge, this is the first work that proposes the idea⁴ of putting sufficiently large weight on the global maximum values of the objective, to decrease the distance between the optimum point before and after Gaussian smoothing (i.e., $\|\mathbf{x}^* - \mu^*\|$). It provides a theoretical foundation (Theorem 2.1 and Lemma 3.4), as well as motivations, for future studies to find better ways of weight shifting than power transforms.

2. In theory, compared to the iteration complexity of

$\min_{\mu} G(\mu, \Sigma)$ is a convex problem given that $NR - \Sigma^{-1}$ is positive definite and \mathcal{C} is convex.

⁴(Dvijotham et al., 2014; Roulet et al., 2020; Chen et al., 2024), which involve power transforms, have not mentioned this idea.

$O(d^2/\varepsilon^4)$ of the standard zeroth-order (ZO) homotopy method analyzed in Theorem 5.1 of (Hazan et al., 2016), GS-PowerOpt is much faster at the iteration complexity of $O(d^2\varepsilon^{-2})$ under the Lipschitz conditions (see Section 3.5) and does not require the strong assumption of “ σ -nice”. Compared to the ZO methods in (Ghadimi & Lan, 2013; Chen et al., 2019) and ZO-SLGH, the main advantage of GS-PowerOpt lies in its ability to approximate the global solution. The majority of our experiments show that it ranks highest among all these algorithms that apply smoothing (see Section 5.4).

3. Our convergence analysis in Corollary 3.9 does not require the Lipschitz condition on the original objective f , which is assumed in the theoretical analysis of homotopy methods in other studies (Hazan et al., 2016; Iwakiri et al., 2022). Therefore, our analysis applies to more situations.
4. The theory derived in this work is on the distance between the found solution and the optimal one, while the convergence analysis in other studies on homotopy (Hazan et al., 2016; Iwakiri et al., 2022) is on the objective value of the found solution. Therefore, our theory has a wider range of applications (e.g., for problems that concern the distance between the found solution and the true one, such as inverse problems (Arridge et al., 2019) and adversarial attacks in image recognitions).

ROAD MAP

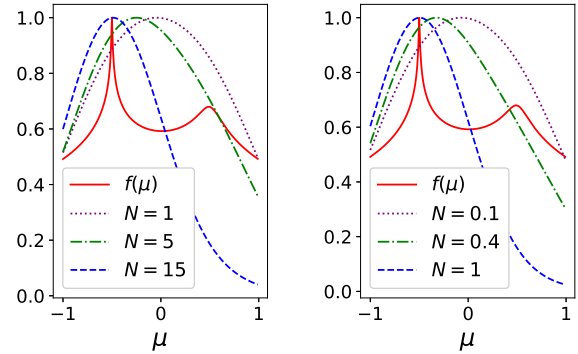
Section 2 describes GS-PowerOpt and the intuition behind it. Its theoretical convergence analysis is performed in Section 3 and 4 (see Subsection 3.1 for a preview). Numerical experiments are included in Section 5. Section 6 provides a guidance on selecting values of the hyper-parameters, and Section 7 concludes. The appendix includes a table of notations, *proofs to all our theoretical results (lemma, theorem, corollary and proposition)*, details of certain compared algorithms, and hyper-parameter values used in our experiments. All our codes can be found at <http://github.com/chen-research/GS-PowerTransform>.

2. GS-PowerOpt: The Proposed Method

2.1. Motivation

Intuitively, if we modify $f(x)$ to put sufficiently large weight on its global maximum point x^* , then $\mu^* := \arg \max_{\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2 I_d)}[f(x)]$ should get close enough to x^* . One way of such modification is by taking (exponential) powers of f , if $f(x^*) > 1$. The difference $f^N(x^*) - f^N(x)$ is positively related with the power N , which indicates that more weight is put on x^* as N increases. Figure 1 verifies

this intuition with an example. As shown in these two toy examples, μ^* approaches x^* as N increases.



(a) Gauss. Smooth of f^N . (b) Gauss. Smooth of e^{Nf} .

Figure 1: Graph of approximated $F_{N,\sigma}(\mu)$, where $F_{N,\sigma}(\mu)$ is defined as (a) $\mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2)}[f^N(x)]$ or (b) $\mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2)}[e^{Nf(x)}]$, $\sigma = 0.5$, and $f(\mu) = -\log((\mu + 0.5)^2 + 10^{-5}) - \log((\mu - 0.5)^2 + 10^{-2}) + 10$ for $|\mu| \leq 1$ and $f(\mu) = 0$ for $|\mu| > 1$. All function graphs are scaled to have a maximum value of 1 for easier comparisons.

From the above intuition, we propose GS-PowerOpt for solving the global optimization problem (1). It is a new method that places more weight on the objective f ’s maximum value (by increasing the gap between $f(x^*)$ and f -values at other points) before performing Gaussian smoothing. Based on GS-PowerOpt, we design two algorithms⁵, Power Gaussian Smoothing (PGS) and Exponential Power Gaussian Smoothing (EPGS), which are featured with replacing the original objective $f(x)$ with a (exponential) power transformation. Specifically, with σ and N as two hyper-parameters, PGS solves $\max_{\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2 I_d)}[f^N(X)]$ and EPGS solves $\max_{\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2 I_d)}[e^{Nf(x)}]$, both using a stochastic gradient ascent algorithm derived in this paper, which does not require the differentiability of f .

2.2. Theoretical Justification of the Motivation

In Theorem 2.1, we lay the ground for justifying the intuition that motivates GS-PowerOpt: Given $\sigma > 0$, for any $\delta > 0$, there exists a threshold such that whenever N exceeds this threshold, the global maximum point $\arg \max_{\mu} F_N(\mu, \sigma)$ lies within a δ -neighborhood of x^* , where $F_N(\mu, \sigma) := \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma I_d)}[f_N(x)]$ and

$$f_N(x) := \begin{cases} f^N(x), & x \in \mathcal{S}; \\ 0, & \text{otherwise,} \end{cases} \quad (\text{PGS setting});$$

$$f_N(x) := \begin{cases} e^{Nf(x)}, & x \in \mathcal{S}; \\ 0, & \text{otherwise.} \end{cases} \quad (\text{EPGS setting}). \quad (4)$$

⁵See Algorithm 1 for the two algorithms.

Theorem 2.1. Let $f : \mathcal{S} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function that is possibly non-concave (and non-negative only for the case of PGS), where \mathcal{S} is compact. Assume that f has a global maximum \mathbf{x}^* such that $\sup_{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \geq \delta} f(\mathbf{x}) < f(\mathbf{x}^*)$ for any $\delta > 0$. For any $\sigma > 0$ and any $N > 0$, define

$$F_{N,\sigma}(\boldsymbol{\mu}) := \frac{1}{(\sqrt{2\pi}\sigma)^d} \int_{\mathbf{x} \in \mathbb{R}^d} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}, \boldsymbol{\mu} \in \mathbb{R}^d$$

where f_N is defined in (4) for either PGS or EPGS. Then, for any $M > 0$ and $\delta > 0$ such that $\bar{B}(\mathbf{x}^*; \delta) := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{x}^*\| \leq \delta\} \subset \mathcal{S}$, there exists $N_{\delta,\sigma,M} > 0$, such that whenever $N > N_{\delta,\sigma,M}$, for any $\|\boldsymbol{\mu}\| \leq M$ and any $i \in \{1, 2, \dots, d\}$ we have that: $\frac{\partial F_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} > 0$ if $\mu_i < x_i^* - \delta$, and $\frac{\partial F_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} < 0$ if $\mu_i > x_i^* + \delta$. Here, μ_i and x_i^* denote the i^{th} entry of $\boldsymbol{\mu}$ and \mathbf{x}^* , respectively.

Remark 2.2. The inequalities in Theorem 2.1 imply that any stationary point $\boldsymbol{\mu}'$ (if any) of $F_{N,\sigma}(\boldsymbol{\mu})$ within the region $\|\boldsymbol{\mu}\| \leq M$ satisfies $|\mu'_i - x_i^*| \leq \delta$ for all $i \in \{1, 2, \dots, d\}$.

Proposition 2.3. Given $\sigma > 0$ and $N > 0$, $F_{N,\sigma}(\boldsymbol{\mu})$ in Theorem 2.1 attains a global max at some $\boldsymbol{\mu}^* \in \mathbb{R}^d$.

2.3. Solution Updating Rule of GS-PowerOpt

For the optimization problem (1), based on Theorem 2.1, with the pre-selected hyper-parameters N and $\sigma > 0$, GS-PowerOpt follows a stochastic gradient ascent scheme to solve $\max_{\boldsymbol{\mu}} F_{N,\sigma}(\boldsymbol{\mu})$. Specifically, the rule for updating the solution candidate is

$$\text{GS-PowerOpt : } \quad \boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t), \quad (5)$$

where $\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t) := \frac{1}{K} \sum_{k=1}^K (\mathbf{x}_k - \boldsymbol{\mu}_t) f_N(\mathbf{x}_k)$, $\{\mathbf{x}_k\}_{k=1}^K$ are independently sampled from the multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_t, \sigma^2 I_d)$, and $f_N(\mathbf{x}_k)$ is defined in (4). Note that $\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t)$ is a sample estimate of $\sigma^2 \nabla F_{N,\sigma}(\boldsymbol{\mu})$, since

$$\begin{aligned} \nabla F_{N,\sigma}(\boldsymbol{\mu}) &= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)} [f_N(\mathbf{x})] \\ &= (\sqrt{2\pi})^{-d} \sigma^{-(d+2)} \int_{\mathbf{x} \in \mathbb{R}^d} (\mathbf{x} - \boldsymbol{\mu}) f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x} \\ &= \sigma^{-2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)} [(\mathbf{x} - \boldsymbol{\mu}) f_N(\mathbf{x})], \end{aligned} \quad (6)$$

where the interchange of differentiation and integral in the second line can be justified by Lebesgue's dominated convergence theorem. A flowchart of the algorithm is given in Figure 2.

Based on GS-PowerOpt, PGS and EPGS are designed in Algorithm 1. They normalize the gradient before updating the solution, which is a common practice to stabilize results. An effective method to avoid computation overflows caused by a large N -value can be found in Appendix B.



Figure 2: A Flowchart of GS-PowerOpt

Algorithm 1 PGS/EPGS for Solving (1)

Input: The power $N > 0$, the scaling parameter $\sigma > 0$, the objective f , the initial value $\boldsymbol{\mu}_0$, the number K of sampled points for gradient approximation, the total number T of $\boldsymbol{\mu}$ -updates, and the learning rate schedule $\{\alpha_t\}_{t=1}^T$.

for t from 0 to $T - 1$ **do**

Independently sample from the multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_t, \sigma^2 I_d)$ and obtain $\{\mathbf{x}_k\}_{k=1}^K$.

$\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t) / \|\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t)\|$, where $\hat{\nabla} F(\boldsymbol{\mu}_t)$ is defined in (5).

end for

Return $\{\boldsymbol{\mu}_t\}_{t=1}^N$, from which $\boldsymbol{\mu}^*$ is selected to approximate \mathbf{x}^* (e.g., $\boldsymbol{\mu}^* := \arg \max_{t \in \{1, 2, \dots, T\}} f(\boldsymbol{\mu}_t)$).

3. Convergence Analysis

3.1. Preview of the Main Theoretical Results

In this section, we prove our main theoretical results, Lemma 3.4 and Corollary 3.9, which indicate that GS-PowerOpt converges to an arbitrarily small neighborhood of $\mathbf{x}^* := \arg \max_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x})$. Specifically, we show that, for any fixed $\delta > 0$ and any fixed $\sigma > 0$, given N larger than some (δ, σ) -dependent threshold, all the stationary points of $F_{N,\sigma}(\boldsymbol{\mu})$ lie in a δ -neighborhood of \mathbf{x}^* (Lemma 3.4), and the stochastic gradient ascent in Eq. (5) converges to one of these stationary points with the iteration complexity of $O_N(d^4 \varepsilon^{-2})$ (Corollary 3.9 and Remark 3.10). This complexity is derived from the bound given in Theorem 3.7, whose proof requires Lemma 3.5 and 3.6.

Furthermore, with additional conditions on f , we prove that the complexity is improved to $O_N(d^2 \varepsilon^{-2})$ in Subsection 3.5, and that the dependence of the O -factor on N can be removed (Corollary 4.2).

3.2. Notations in Section 3

In this section, let $\delta > 0$ and $\sigma > 0$ be fixed, and δ satisfies

$$\mathcal{S}_{\mathbf{x}^*, \delta} := \{\boldsymbol{\mu} \in \mathbb{R}^d : |\mu_i - x_i^*| \leq \delta, i \in \{1, 2, \dots, d\}\} \subset \mathcal{S}, \quad (7)$$

where i denotes the i^{th} entry, $\mathbf{x}^* := \arg \max_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x})$, and \mathcal{S} is f 's compact domain as specified in Assumption 3.1. Let $N > 0$ be such that all the stationary points (if any) of $F_{N,\sigma}(\boldsymbol{\mu})$ within the region $\|\boldsymbol{\mu}\| < \sqrt{d}M$ lie in $\mathcal{S}_{\mathbf{x}^*, \delta}$, where M is specified in Assumption 3.1. Such an N exists because of Theorem 2.1 (see Remark 2.2), and its value depends on δ, σ , and $\sqrt{d}M$. Unless needed for clarity, we omit the

subscripts in $F_{N,\sigma}(\mu)$ as they remain fixed throughout the rest of Section 3, and write $F(\mu)$ instead. $\nabla F(\mu)$ refers to the gradient of F with respect to μ .

3.3. Assumptions and Lemmas

Assumption 3.1. Assume that $f(x) : \mathcal{S} \rightarrow \mathbb{R}$ is a function satisfying the conditions specified in Theorem 2.1, where $\mathcal{S} \subset \mathcal{S}_M := \{x \in \mathbb{R}^d : |x_i| \leq M, i \in \{1, 2, \dots, d\}\}$, for some $M > 0$.

Assumption 3.2. Assume that the learning rate α_t satisfies

$$\alpha_t > 0, \sum_{t=0}^{\infty} \alpha_t = +\infty, \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < +\infty.$$

Lemma 3.3. Under Assumption 3.1, $F_{N,\sigma}(\mu)$ and $\nabla F_{N,\sigma}(\mu)$ are well-defined for all $\mu \in \mathbb{R}^d$.

Lemma 3.4. Under Assumption 3.1, any stationary point of $F_{N,\sigma}(\mu)$ in \mathbb{R}^d belongs to $\mathcal{S}_{x^*,\delta}$ (defined in (7)).

Lemma 3.5. Under Assumption 3.1, for any $\sigma > 0$, the objective function $F_{N,\sigma}(\mu)$ is Lipschitz Smooth. That is, for any $\mu_1, \mu_2 \in \mathbb{R}^d$,

$$\|\nabla F_{N,\sigma}(\mu_1) - \nabla F_{N,\sigma}(\mu_2)\| \leq L\|\mu_1 - \mu_2\|,$$

where $L = 2d\sigma^{-2}f_N(x^*)$, and $f_N(x^*) = f^N(x^*)$ for PGS and $f_N(x^*) = e^{Nf(x^*)}$ for EPGS.

Lemma 3.6. Under Assumption 3.1, for any $\sigma > 0$, let $\tilde{\nabla} F_{N,\sigma}(\mu)$ be defined in (5). Then, $\mathbb{E}[\|\tilde{\nabla} F_{N,\sigma}(\mu)\|^2] < G$, where

$$G = \begin{cases} d\sigma^2 f^{2N}(x^*), & \text{for PGS;} \\ d\sigma^2 e^{2Nf(x^*)}, & \text{for EPGS.} \end{cases}$$

3.4. Iteration Complexity

Theorem 3.7. Let $\{\mu_t\}_{t=0}^T \subset \mathbb{R}^d$ be produced by following the iteration rule of (5), with a pre-selected and deterministic μ_0 and all the involved terms defined as in Section 3.2. Then, under Assumption 3.1 and 3.2, we have that

$$\sum_{t=0}^{T-1} \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\mu_t)\|^2] \leq f_N(x^*) - F(\mu_0) + LG \sum_{t=0}^{\infty} \alpha_t^2,$$

where L and G are as defined in Lemma 3.5 and Lemma 3.6, respectively.

Remark 3.8. The inequality in Theorem 3.7 implies that $\min_{t \in \{0, 1, \dots, T-1\}} \mathbb{E}[\|\nabla F(\mu_t)\|^2]$ approaches 0 as $T \rightarrow \infty$, since the right-hand side is finite and $\sum_{t=0}^{\infty} \alpha_t = \infty$. Furthermore, the iteration complexity of this convergence is given in Corollary 3.9 when $\alpha_t := (t+1)^{-(1/2+\gamma)}$.

Corollary 3.9. Suppose Assumption 3.1 and 3.2 hold. Let $\{\mu_t\}$ be produced by the stochastic gradient ascent rule (5)

of GS-PowerOpt, with a pre-selected and deterministic μ_0 . Then, for any $\varepsilon \in (0, 1)$, after

$$T > (C_1 C_2 d^2 \varepsilon^{-1})^{2/(1-2\gamma)} = O_N((d^2 \varepsilon^{-1})^{2/(1-2\gamma)})$$

times of μ_t -updating by (5), we have that $\min_{t \in \{0, 1, \dots, T\}} \mathbb{E}[\|\nabla F(\mu_t)\|^2] < \varepsilon$. Here, $\gamma \in (0, 1/2)$ is a parameter in the learning rate $\alpha_t := (t+1)^{-(1/2+\gamma)}$, $C_0 = f_N(x^*) - F(\mu_0) + 2f_N^3(x^*) \sum_{t=1}^{\infty} t^{-(1+2\gamma)}$, $C_1 = C_0(1-2\gamma)\sigma^{-2}$, and $C_2 = \max\{1, 2/C_1\}$.

Remark 3.10. The iteration complexity is approximately $O_N(d^4 \varepsilon^{-2})$ if γ is close to 0. Also, the dependence of the big-O factor on N can be removed with an additional assumption. See Corollary 4.2 for this point.

3.5. Improved Iteration Complexity under Lipschitz Conditions

In Proposition 3.11, we re-derive the coefficient L in Lemma 3.5 under the Lipschitz condition on the objective f and its first derivative ∇f . The new value of L does not depend on the dimension d . With this result, the iteration complexity in Corollary 3.9 becomes $O_N(d^2 \varepsilon^{-2})$ when γ is close to 0.

Proposition 3.11. Assume Assumption 3.1, $f(x) = 0$ on the boundary of f 's domain, the Lipschitz condition on f : $|f(x) - f(y)| \leq L_0 \|x - y\|$, and the Lipschitz condition on ∇f : $\|\nabla f(x) - \nabla f(y)\| \leq L_1 \|x - y\|$. Then, the following two statements are true.

1. $\|\nabla F_N(\mu_1, \sigma) - \nabla F_N(\mu_2, \sigma)\| \leq L\|\mu_1 - \mu_2\|$ holds for any $\mu_1, \mu_2 \in \mathbb{R}^d$, where $L = (N f^{N-1}(x^*) L_1 + N(N-1) L_0^2 f^{N-2}(x^*))$ under the PGS setting, and $L = (N L_1 e^{Nf(x^*)} + N^2 L_0^2 e^{Nf(x^*)})$ under the EPGS setting.
2. The iteration complexity in Corollary 3.9 reduces to $O_N((d\varepsilon^{-1})^{2/(1-2\gamma)})$.

Remark 3.12. “ $f(x) = 0$ on the boundary of f 's domain” is assumed without loss of generality, since we can continuously extend f to a larger domain to make it true.

The Lipschitz condition is common in literature. For example, the convergence analysis for the standard homotopy method requires the Lipschitz assumption on f (see Theorem 5.1 in (Hazan et al., 2016)), and the convergence analysis for ZOSLGH requires the Lipschitz condition on both f and ∇f (see Assumption A1 in (Iwakiri et al., 2022)).

4. Removing the Iteration Complexity's Dependence on N

Lemma 3.4 and Corollary 3.9 imply that GS-PowerOpt converges to a δ -neighborhood of x^* at the iteration complexity of $O_N(d^4 \varepsilon^{-2})$. To ensure the convergence to a smaller

neighborhood of \mathbf{x}^* , we need to decrease δ . However, this possibly requires a larger N for Lemma 3.4 and Corollary 3.9 to hold (recall N 's definition in Subsection 3.2), and in turn increases the iteration complexity of $O_N(d^4\varepsilon^{-2})$. The following additional assumption removes the complexity's dependence on N , and guarantees an iteration complexity of $O(d^4\varepsilon^{-2})$ regardless of how close the positive δ is to 0.

Assumption 4.1. $f(\mathbf{x}) \in [0, 1)$ for the PGS case, and $f(\mathbf{x}) < 0$ for the EPGS case.

In practice, this assumption can be realized if we know an upper bound of the objective $f(\cdot)$. For example, if it is known that $f(\cdot) < B$, then we can proceed EPGS with the new objective function $f_1 := f - B$ (f and f_1 share the same global maximum point \mathbf{x}^*). Also, at least shown by our experiments (Table 1 - 4), GS-PowerOpt works well for objectives that do not necessarily satisfy this assumption.

Corollary 4.2. Under Assumption 3.1 and 4.1, the iteration complexity in Corollary 3.9 becomes $O((d^2\varepsilon^{-1})^{2/(1-2\gamma)})$, and the iteration complexity in Point 2 of Proposition 3.11 becomes $O((d\varepsilon^{-1})^{2/(1-2\gamma)})$, both of which are independent from N .

Remark 4.3. Note that the two complexities approach $O(d^4\varepsilon^{-2})$ and $O(d^2\varepsilon^{-2})$, respectively, as $\gamma \rightarrow 0$.

5. Experiments

5.1. Effects of Increasing Powers

We illustrate the improvements made by increasing N for PGS/EPGS through an example problem of $\max_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$, where

$$f(\mathbf{x}) = -\log(\|\mathbf{x} - \mathbf{m}_1\|^2 + 10^{-5}) - \log(\|\mathbf{x} - \mathbf{m}_2\|^2 + 10^{-2}), \quad (8)$$

the global maximum point $\mathbf{m}_1 \in \mathbb{R}^d$ has all its entries equal to -0.5 , and the local maximum point $\mathbf{m}_2 \in \mathbb{R}^d$ has all its entries equal to 0.5 . The graph of its 2D-version is plotted in Figure 3 (a).

With each value of N , we perform both the PGS and EPGS in Algorithm 1 to solve this problem. The experiments are done in two settings, one is two-dimensional ($d = 2$) and the other is five-dimensional ($d = 5$). More details can be found in Appendix E.

The results, plotted in Figure 4, show that, as N increases, the distance between the produced solution $\boldsymbol{\mu}^*$ and the global maximum point \mathbf{x}^* approaches zero (see the decreasing MSE curve in the plot), which is consistent with Theorem 2.1 and the idea that $F(\boldsymbol{\mu})$'s maximum $\boldsymbol{\mu}^*$ approaches the global maximum point \mathbf{x}^* of f as we put more weight on $f(\mathbf{x}^*)$.

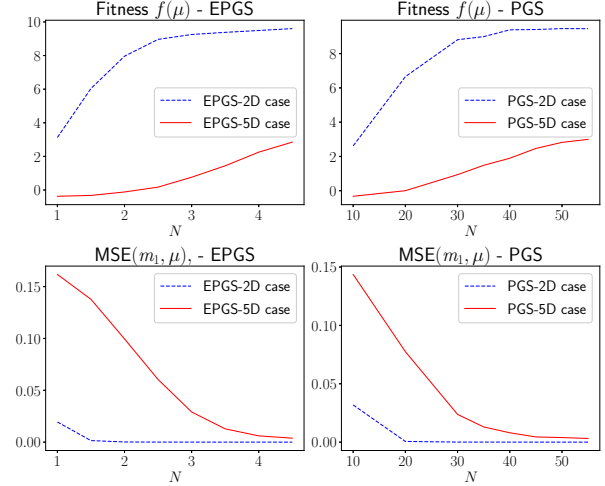


Figure 4: Effects of Increasing N . For each N , we perform the algorithm 100 times to stabilize results, and obtain $\{\boldsymbol{\mu}_k\}_{k=1}^{100}$. The average fitness $\sum_{k=1}^{100} f(\boldsymbol{\mu}_k)/100$ and $\sum_{k=1}^{100} \text{MSE}(\mathbf{m}_1, \boldsymbol{\mu}_k)/100$ are plotted, where $\text{MSE}(\mathbf{m}_1, \boldsymbol{\mu}_k) := \sum_{i=1}^d (\mu_{ki} + 0.5)^2/d$, $\sigma = 1.0$, and f is defined in (8). Note that $\mathbf{x}^* = \mathbf{m}_1$ has all its entries equal to -0.5 .

5.2. Performance on Benchmark Objective Functions

In this subsection, we test the performance of PGS and EPGS on solving (1), with the objective f being either of the two popular benchmark objective functions, the Ackley and the Rosenbrock (maximize-version).

The performances of other popular global algorithms (maximize-version) are reported for comparison, including (1) a standard homotopy method STD-Homotopy; (2) two zeroth-order single-loop Gaussian homotopy algorithms, ZOSLGHd and ZOSLGHr (see the deterministic version⁶ of Algorithm 3 in (Iwakiri et al., 2022)); (3) the gradient-ascent version of the zeroth-order algorithm of ZOSGD (see Equation (1) in (Chen et al., 2019) and Section 2.1 in (Ghadimi & Lan, 2013)) and ZOAdaMM (Algorithm 1 in (Chen et al., 2019))⁷, which were also used for comparisons in (Iwakiri et al., 2022); (4) as well as CMA-ES (Hansen & Ostermeier, 2001; Hansen et al., 2019), one of the state-of-the-art evolutionary algorithms.

More details and discussions on the compared algorithms can be found in Appendix D. The selected hyper-parameter values of all algorithms can be found in Appendix F.

⁶A deterministic version of SLGH is for solving the optimization problem of $\max_{\mathbf{x}} f(\mathbf{x})$, with f being a deterministic function.

⁷ZOSGD and ZOAdaMM were originally designed to solve (3). We take their solutions to (3) as solutions to (1), and treat the scaling parameter in (3) as a hyper-parameter.

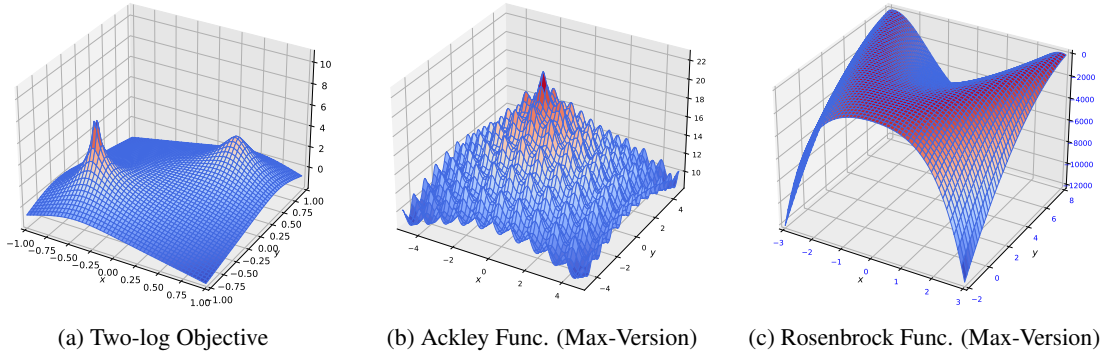


Figure 3: Graph of the Benchmark Objective Functions.

5.2.1. ACKLEY

The Ackley objective function is characterized by numerous local optima and a single global optimum. We solve the max-version of the corresponding problem, which is $\max_{(x,y) \in \mathbb{R}^2} f(x, y)$ and $f(x, y) = 20e^{-0.2\sqrt{0.5(x^2+y^2)}} + e^{0.5(\cos(2\pi x) + \cos(2\pi y))}$. The graph of this function is plotted in Figure 3(b). From both the functional form and the graph, it is not difficult to see that $f(x, y)$ attains its maximum at $\mathbf{x}^* = (0, 0)$.

The solutions and their f -values found by each of the compared algorithms are reported in Table 1. For each algorithm, the total number T of iterations is set as 200, and $K = 100$ samples are used for each solution update. The hyper-parameters are selected by trials. For each set of hyper-parameter candidates, we perform 100 experiments and *take the average* to stabilize the results (so all the reported numbers are averages). The initial solution value for each experiment is drawn from a multivariate Gaussian with mean $\boldsymbol{\mu}_0 = [5.0, 5.0]$ and covariance $0.01I_2$.

The results in the table show that all the algorithms perform well and are close to each other, except STD-Homotopy.

5.2.2. ROSENBRACK

The Rosenbrock objective is known to be difficult to optimize, since its global optimum point $\mathbf{x}^* = (1.0, 1.0)$ is surrounded by a flat curved plane (see Figure 3(c)). The problem to be solved is $\max_{(x,y) \in \mathbb{R}^2} f(x, y)$, where $f(x, y) = -100(y - x^2)^2 - (1 - x)^2$.

The performance of each algorithm is recorded in Table 2, which shows that the performances of CMA-ES and EPGS are close, and are significantly better than those of other algorithms.

Table 1: Performances on Maximizing Ackley. “Iter. Taken” refers to the number of iterations taken to reach the best found solution. All numbers are rounded to keep at most 3 decimal places.

Algorithm	Iter. Taken	Best Solution Found ($\boldsymbol{\mu}^*$)	$f(\boldsymbol{\mu}^*)$
CMA-ES	116	(0.0, 0.0)	22.718
EPGS ($N = 1$)	143	(0.001, 0.0)	22.683
PGS ($N = 20$)	141	(0.001, 0.002)	22.678
ZOSLGHr	131	(0.001, -0.002)	22.621
ZOAdaMM	158	(0.005, 0.001)	22.613
ZOSLGHd	123	(-0.003, -0.001)	22.61
ZOSGD	174	(0.005, 0.007)	22.596
STD-Homotopy	194	(0.962, 0.946)	17.58

5.3. Performance on the Black-box Targeted Adversarial Attack

Let \mathcal{C} be a black-box⁸ image classifier. The targeted adversarial attack on \mathcal{C} refers to the task of modifying the pixels of a given image \mathbf{a} so that $\mathcal{C}(\mathbf{a} + \mathbf{x})$ is equal to a pre-specified target label \mathcal{T} . We perform attacks in the most difficult case - the target label \mathcal{T} is pre-selected to be the one with the smallest predicted probability. Another goal of this task is to minimize the perturbation size $\|\mathbf{x}\|$. Hence, we set the loss as

$$L(\mathbf{x}) := \max_{i \neq \mathcal{T}} (\max_i \mathcal{C}(\mathbf{a} + \mathbf{x})_i - \mathcal{C}(\mathbf{a} + \mathbf{x})_{\mathcal{T}}, \kappa) + \lambda \|\mathbf{x}\|,$$

where $\mathcal{C}(\mathbf{a} + \mathbf{x})_i$ denotes the predicted logit (i.e., log probability) for the i^{th} class, κ is a hyper-parameter that controls the certainty level of the attack, λ is a regularization coefficient (we set $\lambda = 1$ in our experiments), $\mathcal{T} := \arg \min_i \mathcal{C}(\mathbf{a})_i$, and $\|\mathbf{x}\|$ denotes the L_2 norm of

⁸A black-box classifier refers to a classification model whose parameters are not accessible.

Table 2: Performances on Maximizing Rosenbrock. The differences between this experiment and that in Table 1 are: The number T of iterations is set as 1000, and the initial solution candidate is taken as $\mu_0 := [-3.0, 2.0]$. For PGS, the Rosenbrock is added by 20,000 to ensure the search agent only encounter positive values. The global maximum point of the Rosenbrock function is (1,1).

Algorithm	Iter. Taken	Best Solution Found (μ^*)	$f(\mu^*)$
CMA-ES	72	(1.0, 1.0).	0.0
EPGS ($N = 3$)	487.	(0.999, 1.000)	-0.017
STD-Homotopy	624	(0.903, 0.885)	-2.401
PGS ($N = 1$)	513	(0.773, 1.025)	-22.84
ZOAdaMM	852	(0.004, 0.618).	-39.206
ZOSLGHr	148	(0.105, 0.938).	-88.477
ZOSGD	45	(0.272, 1.173).	-121.14
ZOSLGHd	471	(-0.447, 1.991).	-137.016

x (i.e., the square root of the sum of squares of entries in x). This loss function resembles the popular one designed in (Carlini & Wagner, 2017).

With EPGS and other compared algorithms, we perform adversarial attacks on 100 randomly selected images from each of two image datasets, the set of MNIST handwritten digits (LeCun et al., 1998) and the CIFAR-10 set (Krizhevsky & Hinton, 2009). Specifically, the goal is to solve:

$$\max_{x \in \mathbb{R}^d} f(x) := -L(x),$$

where d is the total number of pixels in each image. For Figure-MNIST images, $d = 28 \times 28$, and for CIFAR-10 images, $d = 32 \times 32 \times 3$. We call f as the fitness function.

The classifier is a robust convolutional neural network (CNN) trained using the technique of defensive distillation⁹. The distillation temperature is set at 100, which leads to a high level of robustness (Carlini & Wagner, 2017). In the MNIST attacks, our trained classifier \mathcal{C} has a classification accuracy of 97.4% on the testing images. In the CIFAR-10 attacks, the trained \mathcal{C} has a test accuracy of 86.2%.

The hyper-parameters of all the tested algorithms are selected by trials and can be found in Appendix F. For ZO-SLGHd, ZO-SLGHr, and ZO-AdaMM, the hyper-parameters reported in Appendix E.1 of (Iwakiri et al., 2022) are included in our candidate set, since they performed similar tasks.

⁹We borrow the code by (Carlini & Wagner, 2017) for training the classifier from https://github.com/carlini/nn_robust_attacks, which applies TensorFlow (Abadi et al., 2015) for model training. We changed the layer structure in the neural network to either increase accuracy or to decrease computation complexity while maintaining the accuracy.

We choose EPGS over PGS for this task since the fitness function $f(x)$ can be negative, which makes EPGS more convenient to apply. But note that we can modify $f(x)$ by adding a large positive constant to facilitate PGS.

5.3.1. MNIST

For each image a_m that is randomly drawn from the dataset, where $m \in \{1, 2, \dots, 100\}$, and each algorithm, we perform an attack (i.e., experiment) of T_{total} iterations. Let $\{\mu_{m,t}\}_{t=0}^{T_{total}-1}$ denote all the perturbations (solutions) produced in these T_{total} iterations. We say that a perturbation μ is successful if the predicted log probability of the target label is at least $\kappa = 0.001$ greater than that of other classes (i.e., $\mathcal{C}(a + \mu)_{\mathcal{T}} - \max_{i \neq \mathcal{T}} \mathcal{C}(a + \mu)_i > \kappa$). We say that an attack is successful if the produced $\{\mu_{m,t}\}_{t=0}^{T_{total}-1}$ contain at least one successful perturbation. If the attack is successful, let μ_m^* denote the successful perturbation with the largest R^2 -value among $\{\mu_{m,t}\}_{t=0}^{T_{total}-1}$, and let T_m denote the number of iterations taken by the algorithm to produce μ_m^* . We use $R^2(a, a + \mu) := 1 - \frac{\sum_{i=1}^d \mu_i^2}{\sum_{i=1}^d (a_i - \bar{a})^2}$ to measure the similarity between a and the perturbed image $a + \mu$, where a_i and μ_i denote the i^{th} pixel (entry) of a and μ , respectively, and \bar{a} denotes the mean of $\{a_i\}$.

With the above notations, we construct four measures on the performances of an algorithm. One is the success rate, which refers to the ratio of successful image attacks out of the total number of attacks (100). The second measure is the average R^2 , which equals $\bar{R}^2 := \sum_{m \in \mathbb{S}} R^2(a_m, a_m + \mu_m^*) / |\mathbb{S}|$, where \mathbb{S} denotes the set of indices of the successful attacks and μ_m^* denotes the optimal perturbation for the m^{th} figure. The third one is the average $\|\mu^*\|$ of $\{\|\mu_m^*\|\}_{m \in \mathbb{S}}$, where $\|\cdot\|$ denotes the L_2 -norm. The last measure is the average \bar{T} of $\{T_m\}_{m \in \mathbb{S}}$.

The results are reported in Table 3, from which we see that EPGS has a high \bar{R}^2 -score of 87% (ranks among the top 2), indicating that the perturbed image is closest to the original one ($R^2 = 100\%$ implies that a and $a + \mu$ are identical). Also, the average number of iterations taken by EPGS to reach the optimal perturbation is 397, which is among the two fastest algorithms.

5.3.2. CIFAR-10

With 100 randomly drawn images from the CIFAR-10 test set, we repeat the per-image targeted adversarial attacks in Section 5.3.1. Their results are reported in Table 4. These results are in general better than those in the MNIST attacks, which we believe is because of the lower accuracy of the trained CIFAR-10 CNN (86.0% vs 97.8%).

The results show that EPGS produces a success rate of 98% and scores in the top three with respect to \bar{R}^2 . Specifically, the \bar{R}^2 -value of 98% indicates that the perturbed image

Table 3: Targeted Adversarial Attack on 100 MNIST images (per-image). For each image attack, we set the initial perturbation $\mu_0 = \mathbf{0}$, and $T_{total} = 1,500$. The success rate (SR) is the portion of successful attacks out of the 100 attacks. \bar{R}^2 , $\|\mu^*\|$, and \bar{T} are defined in Section 5.3.1. The numbers in the parentheses are sample standard deviations. STD-Htp is STD-Homotopy for short. For EP GS, N is selected to be 0.05. See Table 10 in the Appendix for the produced adversarial images for four randomly selected images.

Algorithm	SR	\bar{R}^2	$\ \mu^*\ $	\bar{T}
CMA-ES	100%	89%(4%)	2.81(0.61)	1489(12)
EP GS	100%	87%(5%)	3.01(0.60)	397(101)
ZOSGD	100%	85%(5%)	3.14(0.61)	1427(242)
ZOSLGHd	100%	74%(9%)	4.21(0.71)	1490(24)
ZOSLGHr	100%	65%(13%)	4.86(0.81)	476(658)
ZOAdaMM	100%	29%(27%)	6.88(1.15)	45(15)
STD-Htp	97%	-4%(37%)	8.25(1.09)	530(264)

Table 4: Targeted Adversarial Attack on 100 CIFAR-10 images (per-image). For each image attack, we set $\mu_0 = \mathbf{0}$ and $T_{total} = 1,500$. For EP GS, N is selected as 0.03. See Table 11 in the Appendix for the produced adversarial images for four randomly selected images.

Algorithm	SR	\bar{R}^2	$\ \mu^*\ $	\bar{T}
ZOSLGHd	98%	99%(1%)	1.72(0.32)	1290(411)
ZOSLGHr	98%	98%(3%)	2.66(0.68)	456(345)
EP GS	98%	98%(2%)	3.05(0.57)	748(248)
CMA-ES	99%	75%(25%)	10.06(2.35)	158(399)
ZOAdaMM	100%	58%(39%)	13.13(2.71)	58(31)
ZOSGD	62%	99%(1%)	1.19(0.19)	764(349)
STD-Htp	52%	87%(13%)	7.54(1.57)	566(396)

produced by EP GS is very close to the original image, which is consistent with our goal that the perturbation size $\|\mu^*\|$ should be small.

5.4. Summary on Experiment Results

Comparing to algorithms that use the smoothing technique (which **excludes** CMA-ES), PGS and EP GS ranked among the tops in all the tasks we performed. Specifically, in the tasks of Ackley, Rosenbrock, and MNIST-attack (Table 1, 2 and 3), EP GS has the highest fitness value $f(\mu^*)$ than other algorithms except CMA-ES. For CIFAR-10 attacks (Table 4), EP GS’s \bar{R}^2 score equals that of the best (ZOSLGHd). In sum, EP GS outperforms other algorithms that also apply smoothing techniques in the majority of the performed experiments.

While EP GS underperforms CMA-ES in experiments of

Ackley, Rosenbrock, and MNIST, it beats CMA-ES in the CIFAR-10 task with a comparable success rate and a significantly smaller perturbation norm (see Table 4). Moreover, to the best of our knowledge, the theoretical convergence guarantee for CMA-ES on general non-convex objectives is less developed, providing a competitive edge to GS-PowerOpt.

5.5. The Case of Multiple Global Maxima

Although the convergence theory (e.g., Theorem 2.1) requires the condition that the objective f has a unique global maximum, it is a sufficient rather than a necessary condition for GS-PowerOpt to work. To illustrate this point, we perform an additional experiment with EP GS on an objective with two global maxima: $f(\mathbf{x}) = -\log(\|\mathbf{x} - \mathbf{m}_1\|^2 + 10^{-5}) - \log(\|\mathbf{x} - \mathbf{m}_2\|^2 + 10^{-5})$, where $\mathbf{m}_1 = [-.5, -.5]$ and $\mathbf{m}_2 = [.5, .5]$. A hundred trials are performed, where $N = 1$ and each μ_0 is randomly sampled around the origin. The mean square error between the found solution and one of the two global maximum point is close to 0, with its average (over the 100 trials) equal to 1.4×10^{-5} and its sample standard deviation equal to 1.6×10^{-5} . This example indicates that GS-PowerOpt is capable of locating at least one of the multiple global optima. We leave the corresponding theoretical analysis in our future works.

6. Guidance on Choosing Hyper-parameters

Although a larger N guarantees a convergence of GS-PowerOpt to a smaller neighborhood of \mathbf{x}^* in theory (see Subsection 3.1), it also leads to a higher variance of the update term $\hat{\nabla} F_{N,\sigma}(\mu_t)$ in Eq. (5), which in turn decreases the efficiency of the algorithm in practice. Therefore, we recommend to start from a moderate N and incrementally increase its value during tuning. Although the proper starting value of N may vary for different problems, based on our experience, 5 for PGS and 0.1 for EP GS are good choices.

While a large σ increases the exploration range for each μ update, it also increases the sample complexity of the algorithm. Our experiments show that a σ -value of 10% of the search radius for \mathbf{x} is a good starting value for tuning.

7. Conclusion and Future Work

The convergence analysis and numerical results show that the easily implemented optimization method of GS-PowerOpt stands out among its peers that also apply smoothing techniques. Our work provides a foundation for future studies to explore more efficient ways to increase the gap between $f(\mathbf{x}^*)$ and f values at other points before smoothing.

Acknowledgements

We thank the anonymous reviewers for their valuable comments and suggestions for improving this paper.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., H., A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. Tensorflow: Large-scale machine learning on heterogeneous distributed systems, 2015. URL <https://www.tensorflow.org/>.
- Arridge, S., Maass, P., Öktem, O., and Schönlieb, C.-B. Solving inverse problems using data-driven models. *Acta Numerica*, 28:1–174, 2019.
- Blake, A. and Zisserman, A. *Visual Reconstruction*. MIT press, 1987.
- Boer, P. D., Kroese, D. P., Mannor, S., and Rubinstein, R. Y. A tutorial on the cross-entropy method. *Ann. Oper. Res.*, 134(1):19–67, 2005. doi: 10.1007/s10479-005-5724-z.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, 2017. doi: 10.1109/SP.2017.49.
- Chen, J., Guo, Z., Li, H., and Chen, C. P. Regularizing scale-adaptive central moment sharpness for neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 35(5), 2024.
- Chen, X., Liu, S., Xu, K., Li, X., Lin, X., Hong, M., and Cox, D. Zo-adamm: Zeroth-order adaptive momentum method for black-box optimization. *Advances in neural information processing systems*, 32, 2019.
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B., and LeCun, Y. The loss surfaces of multilayer networks. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pp. 192–204. PMLR, 2015.
- Dvijotham, K., Fazel, M., and Todorov, E. Universal convexification via risk-aversion. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, 2014.
- Gao, K. and Sener, O. Generalizing Gaussian smoothing for random search. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 7077–7101. PMLR, 2022.
- Ghadimi, S. and Lan, G. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.
- Hansen, N. and Ostermeier, A. Completely derandomized self-adaptation in evolution strategies. *Evolutionary computation*, 9(2):159–195, 2001.
- Hansen, N., Akimoto, Y., and Baudis, P. CMA-ES/pycma on Github, 2019. URL <https://doi.org/10.5281/zenodo.2559634>.
- Hazan, E., Levy, K. Y., and Shalev-Shwartz, S. On graduated optimization for stochastic non-convex problems. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pp. 1833–1841, 2016.
- Iwakiri, H., Wang, Y., Ito, S., and Takeda, A. Single loop gaussian homotopy method for non-convex optimization. In *Advances in Neural Information Processing Systems*, volume 35, pp. 7065–7076, 2022.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical Report TR-2009, University of Toronto, 2009.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Lei, Y., Hu, T., Li, G., and Tang, K. Stochastic gradient descent for nonconvex learning without bounded gradient assumptions. *IEEE transactions on neural networks and learning systems*, 31(10):4394–4400, 2019.
- Lin, X., Yang, Z., Zhang, X., and Zhang, Q. Continuation path learning for homotopy optimization. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 21288–21311. PMLR, 2023.
- Locatelli, M. and Schoen, F. *Global Optimization: Theory, Algorithms, and Applications*. SIAM, Philadelphia, PA, 2013. doi: 10.1137/1.9781611972672.

- Magnus, J. R. and Neudecker, H. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley & Sons Ltd, 2019.
- Mertikopoulos, P., Hallak, N., Kavis, A., and Cevher, V. On the almost sure convergence of stochastic gradient descent in non-convex problems. *Advances in Neural Information Processing Systems*, 33:1117–1128, 2020.
- Miranda, L. J. V. PySwarms, a research-toolkit for particle swarm optimization in python. *Journal of Open Source Software*, 3, 2018. doi: 10.21105/joss.00433.
- Mobahi, H. and Fisher, J. W. A theoretical analysis of optimization by gaussian continuation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, pp. 1205–1211, 2015. doi: 10.1609/aaai.v29i1.9356.
- Nesterov, Y. and Spokoiny, V. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- Roulet, V., Fazel, M., Srinivasa, S., and Harchaoui, Z. On the convergence of the iterative linear exponential quadratic gaussian algorithm to stationary points. In *2020 American Control Conference (ACC)*, pp. 132–137. IEEE, 2020.
- Starnes, A. and Webster, C. Improved performance of stochastic gradients with gaussian smoothing, 2024. URL <https://arxiv.org/abs/2311.00531>.
- Van Laarhoven, P. J., Aarts, E. H., van Laarhoven, P. J., and Aarts, E. H. *Simulated annealing*. Springer, 1987.
- Xiao, L. and Zhang, T. A proximal-gradient homotopy method for the l_1 -regularized least-squares problem. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pp. 839–846, 2012.
- Xu, Y., Yan, Y., Lin, Q., and Yang, T. Homotopy smoothing for non-smooth problems with lower complexity than $o(1/\epsilon)$. In *Advances in Neural Information Processing Systems*, volume 29, 2016.

A. Table of Notations

Notation	Description	Reference
\mathcal{S}	A compact set in \mathbb{R}^d .	Theorem 2.1
f	The real-valued objective function, which is further assumed to be non-negative for PGS.	Eq. (1) or Theorem 2.1
\mathbf{x}^*	The global maximum point of f .	
x_i^*	The i^{th} entry of \mathbf{x}^* , where $i \in \{1, 2, \dots, d\}$.	(7)
μ_i	The i^{th} entry of $\boldsymbol{\mu} \in \mathbb{R}^d$, where $i \in \{1, 2, \dots, d\}$.	(7)
$\mathcal{N}(\boldsymbol{\mu}, \sigma I_d)$	A multivariate standard Gaussian distribution, where $\boldsymbol{\mu} \in \mathbb{R}^d$, $\sigma > 0$, and I_d denotes the $d \times d$ identity matrix.	Eq. (2)
$\hat{f}_\sigma(\boldsymbol{\mu})$	$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)}[f(\mathbf{x})]$	(2)
N	$N > 0$ and $N \in \mathbb{R}$. In Section 3, it is assumed to be greater than a threshold (see Section 3.2 for details).	
$f_N(\mathbf{x})$	The (exponential) power of f , extended to \mathbb{R}^d .	Eq. (4)
$F_{N,\sigma}(\boldsymbol{\mu})$	$F_{N,\sigma}(\boldsymbol{\mu}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)}[f_N(\mathbf{x})]$.	Theorem 2.1
$F(\boldsymbol{\mu})$	The abbreviation of $F_{N,\sigma}(\boldsymbol{\mu})$, where N and σ are fixed.	Section 3.2.
$\nabla F(\boldsymbol{\mu})$,	The gradient of $F_{N,\sigma}(\boldsymbol{\mu})$ with respect to $\boldsymbol{\mu}$.	Eq. (6).
$\nabla F_N(\boldsymbol{\mu}, \sigma)$		
$\hat{F}_{N,\sigma}(\boldsymbol{\mu})$	A sample estimate of $\sigma^2 \nabla F_{N,\sigma}(\boldsymbol{\mu})$.	Eq. (5)
$\mathcal{S}_{\mathbf{x}^*, \delta}$	A δ -neighborhood of \mathbf{x}^* in \mathbb{R}^d , where $\delta > 0$.	Eq. (7)
L	The Lipschitz coefficient of $\nabla F_{N,\sigma}(\boldsymbol{\mu})$.	Lemma 3.5.
G	The upper bound of $\mathbb{E}[\ \hat{\nabla} F(\boldsymbol{\mu})\ ^2]$.	Lemma 3.6.
$\ \mathbf{x}\ $	$\ \mathbf{x}\ = \sqrt{\sum_{i=1}^d x_i^2}$, where $\mathbf{x} = [x_1, x_2, \dots, x_d] \in \mathbb{R}^d$.	

B. PGS/EGPS with A Baseline to Avoid Computation Overflows

An effective method to avoid computation overflows caused by a large N -value in Algorithm 1 is to modify the gradient estimate as $\hat{\nabla}_{\boldsymbol{\mu}} F(\boldsymbol{\mu}_t) = \frac{1}{K} \sum_{k=1}^K (\mathbf{x}_k - \boldsymbol{\mu}_t) f_N^{(b)}(\mathbf{x}_k)$, where $f_N^{(b)}(\mathbf{x}_k) = 0$ if $\mathbf{x}_k \notin \mathcal{S}$, and

$$f_N^{(b)}(\mathbf{x}_k) = \begin{cases} f_N^N(\mathbf{x}_k) / f_N^N(\boldsymbol{\mu}_t), & \text{for PGS;} \\ e^{N(f(\mathbf{x}_k) - f(\boldsymbol{\mu}_t))}, & \text{for EPGS,} \end{cases}$$

if $\mathbf{x}_k \in \mathcal{S}$. This modification produces the same $\boldsymbol{\mu}$ -updates as in Algorithm 1 because of the gradient-normalization step. We call this algorithm as *PGS (EPGS) with a baseline*.

C. Proofs to Theoretical Results

C.1. Proof to Theorem 2.1 for EPGS

Theorem 2.1 *Let $f : \mathcal{S} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuous function that is possibly non-concave (and non-negative only for the case of PGS), where \mathcal{S} is compact. Assume that f has a global maximum \mathbf{x}^* such that $\sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}^*\| \geq \delta} f(\mathbf{x}) < f(\mathbf{x}^*)$ for any $\delta > 0$. For $\sigma > 0$ and any $N > 0$, define*

$$F_{N,\sigma}(\boldsymbol{\mu}) := (\sqrt{2\pi}\sigma)^{-d} \int_{\mathbf{x} \in \mathbb{R}^d} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}, \quad \boldsymbol{\mu} \in \mathbb{R}^d,$$

where f_N is defined in (4) for either PGS or EPGS. Then, for any $M > 0$ and $\delta > 0$ such that $\bar{B}(\mathbf{x}^*; \delta) := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{x}^*\| \leq \delta\} \subset \mathcal{S}$, there exists $N_{\delta,\sigma,M} > 0$, such that whenever $N > N_{\delta,\sigma,M}$, for any $\|\boldsymbol{\mu}\| < M$ and any $i \in \{1, 2, \dots, d\}$ we have that: $\frac{\partial F_N(\boldsymbol{\mu}, \sigma)}{\partial \mu_i} > 0$ if $\mu_i < x_i^* - \delta$, and $\frac{\partial F_N(\boldsymbol{\mu}, \sigma)}{\partial \mu_i} < 0$ if $\mu_i > x_i^* + \delta$. Here, μ_i and x_i^* denote the i^{th} entry of $\boldsymbol{\mu}$ and \mathbf{x}^* , respectively.

Here, we provide the proof for the EPGS setting, and the proof for the PGS setting is similar.

Proof. Recall that for EPGS, $f_N(\mathbf{x}_k) := \begin{cases} e^{Nf(\mathbf{x}_k)}, & \mathbf{x} \in \mathcal{S}; \\ 0, & \text{otherwise.} \end{cases}$ For any given $\delta > 0$, define $V_\delta := \sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}^*\| \geq \delta} f(\mathbf{x})$ and $D_\delta := (V_\delta + f(\mathbf{x}^*))/2$. Using this symbol, we re-write $F_{N,\sigma}(\boldsymbol{\mu})$ as

$$F_{N,\sigma}(\boldsymbol{\mu}) = e^{D_\delta N} G_{N,\sigma}(\boldsymbol{\mu}) = e^{D_\delta N} (H_{N,\sigma}(\boldsymbol{\mu}) + R_{N,\sigma}(\boldsymbol{\mu})), \quad (9)$$

where

$$\begin{aligned} G_{N,\sigma}(\boldsymbol{\mu}) &:= (\sqrt{2\pi}\sigma)^{-d} \int_{\mathbf{x} \in \mathbb{R}^d} e^{-ND_\delta} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}, \\ H_{N,\sigma}(\boldsymbol{\mu}) &:= (\sqrt{2\pi}\sigma)^{-d} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta)} e^{-ND_\delta} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}, \\ R_{N,\sigma}(\boldsymbol{\mu}) &:= (\sqrt{2\pi}\sigma)^{-d} \int_{\mathbf{x} \notin B(\mathbf{x}^*; \delta)} e^{-ND_\delta} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}, \end{aligned}$$

where $B(\mathbf{x}^*; \delta) := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{x}^*\| < \delta\}$. The main idea of the proof is to first show that $\left| \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right|$ dominates $\left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right|$ if N is sufficiently large, and then that the sign of $\frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i}$ satisfies the declared property in the theorem.

We derive an upper bound for $\left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right|$. For any $\boldsymbol{\mu} \in \mathbb{R}^d$,

$$\begin{aligned} \left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| &\leq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \notin B(\mathbf{x}^*; \delta)} |x_i - \mu_i| e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} e^{-ND_\delta} f_N(\mathbf{x}) d\mathbf{x} \\ &\leq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \notin B(\mathbf{x}^*; \delta)} |x_i - \mu_i| e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} e^{N(f(\mathbf{x}) - D_\delta)} d\mathbf{x} \\ &\leq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \notin B(\mathbf{x}^*; \delta)} |x_i - \mu_i| e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} e^{N(V_\delta - D_\delta)} d\mathbf{x} \\ &\leq \frac{e^{N(V_\delta - D_\delta)}}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in \mathbb{R}^d} |x_i - \mu_i| e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x} \\ &\leq e^{N(V_\delta - D_\delta)} \left(\prod_{j \neq i} \frac{1}{\sqrt{2\pi}\sigma} \int_{x_j \in \mathbb{R}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma^2}} dx_j \right) \\ &\quad \cdot \frac{1}{\sqrt{2\pi}\sigma^3} \int_{x_i \in \mathbb{R}} |x_i - \mu_i| e^{-\frac{(x_i - \mu_i)^2}{2\sigma^2}} dx_i \\ &= e^{N(V_\delta - D_\delta)} \frac{1}{\sqrt{2\pi}\sigma^3} \int_{y \in \mathbb{R}} \sqrt{2}\sigma |y| e^{-y^2} d(\sqrt{2}\sigma y), \quad y := \frac{x_i - \mu_i}{\sqrt{2}\sigma}, \\ &= e^{N(V_\delta - D_\delta)} \frac{\sqrt{2}}{\sqrt{\pi}\sigma} \cdot 2 \int_0^\infty y e^{-y^2} dy, \\ &= e^{N(V_\delta - D_\delta)} \frac{\sqrt{2}}{\sqrt{\pi}\sigma} \cdot \int_0^\infty e^{-y^2} dy^2, \\ &= e^{N(V_\delta - D_\delta)} \frac{\sqrt{2}}{\sqrt{\pi}\sigma} \cdot \int_0^\infty e^{-z} dz, \\ &= \frac{\sqrt{2} e^{N(V_\delta - D_\delta)}}{\sqrt{\pi}\sigma} \end{aligned} \quad (10)$$

where the third inequality sign is because $\|\mathbf{x} - \mathbf{x}^*\| \geq \delta \Rightarrow f(\mathbf{x}) \leq V_\delta$, and the fifth inequality sign is from the separability of a multivariate integral.

Since f is continuous, for $\epsilon_\delta := f(\mathbf{x}^*) - D_\delta > 0$ (because $V_\delta < f(\mathbf{x}^*)$), there exists $\delta' \in (0, \delta)$ such that whenever $\|\mathbf{x} - \mathbf{x}^*\| \leq \delta'$,

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) - \epsilon_\delta = D_\delta > V_\delta. \quad (11)$$

Using this result, we derive a lower bound for $\left| \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right|$ when $\|\boldsymbol{\mu}\| \leq M$ and $|\mu_i - x_i^*| > \delta$.

$$\begin{aligned}
 \left| \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| &= \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta)} |x_i - \mu_i| e^{-ND\delta} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x} \\
 &= \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta)} |x_i - \mu_i| e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}} e^{N(f(\mathbf{x})-D\delta)} d\mathbf{x}, \text{ since } B(\mathbf{x}^*; \delta) \subset \mathcal{S}, \\
 &\geq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta')} |x_i - \mu_i| e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}} e^{N(f(\mathbf{x})-D\delta)} d\mathbf{x} \\
 &\geq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta')} (\delta - \delta') e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x} \\
 &\geq \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta')} (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} e^{-\frac{\|\mathbf{x}\|^2}{\sigma^2}} d\mathbf{x}, \quad \|\mathbf{x} - \boldsymbol{\mu}\|^2 \leq 2(\|\mathbf{x}\|^2 + \|\boldsymbol{\mu}\|^2), \\
 &\geq (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} V(\delta', d, \sigma)
 \end{aligned} \tag{12}$$

where the fourth line is implied by

- $e^{N(f(\mathbf{x})-D\delta)} \geq 1$ because of (11);
- $|x_i - \mu_i| = |(x_i - x_i^*) + (x_i^* - \mu_i)| \geq |x_i^* - \mu_i| - |x_i - x_i^*| > \delta - \delta'$.

In the last line of (12),

$$V(\delta', d, \sigma) := \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta')} e^{-\frac{\|\mathbf{x}\|^2}{\sigma^2}} d\mathbf{x}.$$

The positive number $N_{\delta,\sigma,M}$ is constructed by solving the following inequality for N , which involves the two bounds in (10) and (12).

$$\frac{\sqrt{2}e^{N(V_\delta-D_\delta)}}{\sqrt{\pi}\sigma} < (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} V(\delta', d, \sigma).$$

The solution of this inequality is

$$N > \frac{\ln \left(\frac{\sqrt{\pi}\sigma}{\sqrt{2}} (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} V(\delta', d, \sigma) \right)}{V_\delta - D_\delta},$$

where $V_\delta - D_\delta < 0$ and the numerator is negative for sufficiently large $M > 0$. Therefore, whenever $\|\boldsymbol{\mu}\| \leq M$, $|\mu_i - x_i^*| > \delta$, and

$$N > N_{\delta,\sigma,M} := \max \left\{ 0, \frac{\ln \left(\frac{\sqrt{\pi}\sigma}{\sqrt{2}} (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} V(\delta', d, \sigma) \right)}{V_\delta - D_\delta} \right\},$$

we have

$$\left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| \leq \frac{\sqrt{2}e^{N(V_\delta-D_\delta)}}{\sqrt{\pi}\sigma} < (\delta - \delta') e^{-\frac{M^2}{\sigma^2}} V(\delta', d, \sigma) \leq \left| \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right|. \tag{13}$$

When $N > N_{\delta,\sigma,M}$, $\|\boldsymbol{\mu}\| \leq M$, and $\mu_i > x_i^* + \delta$,

$$\begin{aligned}
 \frac{\partial G_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} &= \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} + \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \\
 &= \frac{1}{(\sqrt{2\pi})^k \sigma^{k+2}} \int_{\mathbf{x} \in B(\mathbf{x}^*; \delta)} (x_i - \mu_i) e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}} e^{N(f(\mathbf{x})-D\delta)} d\mathbf{x} + \frac{\partial R_{N,\sigma}(\boldsymbol{\mu}, \sigma)}{\partial \mu_i} \\
 &= - \left| \frac{\partial H_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| + \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \\
 &\stackrel{\text{by (13)}}{<} - \left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| + \left| \frac{\partial R_{N,\sigma}(\boldsymbol{\mu})}{\partial \mu_i} \right| \\
 &= 0,
 \end{aligned} \tag{14}$$

where the third line is because the integrand of the first term is always negative in the integration region.

On the other hand, when $N > N_{\delta,\sigma,M}$, $\|\mu\| \leq M$, and $\mu_i < x_i^* - \delta$,

$$\begin{aligned} \frac{\partial G_{N,\sigma}(\mu)}{\partial \mu_i} &= \frac{\partial H_{N,\sigma}(\mu)}{\partial \mu_i} + \frac{\partial R_{N,\sigma}(\mu)}{\partial \mu_i} \\ &= \left| \frac{\partial H_{N,\sigma}(\mu)}{\partial \mu_i} \right| + \frac{\partial R_{N,\sigma}(\mu)}{\partial \mu_i} \\ &\stackrel{\text{by (13)}}{>} \left| \frac{\partial R_{N,\sigma}(\mu)}{\partial \mu_i} \right| - \left| \frac{\partial R_{N,\sigma}(\mu)}{\partial \mu_i} \right| \\ &= 0. \end{aligned} \tag{15}$$

Then, (14) and (15) imply the result in the theorem since $\frac{\partial G_{N,\sigma}(\mu)}{\partial \mu_i}$ and $\frac{\partial F_{N,\sigma}(\mu)}{\partial \mu_i}$ share the same sign (see Eq. (9)). \square

C.2. Proof to Proposition 2.3

Proposition 2.3 Given $\sigma > 0$ and $N > 0$, $F_{N,\sigma}(\mu)$ in Theorem 2.1 has a global max point μ^* .

Proof. Since $\sigma > 0$ and $N > 0$ are fixed, we denote $F_{N,\sigma}(\mu)$ by $F(\mu)$ for convenience. Recall that $F(\mu) \geq 0$ because of the definition of f_N in (4). Our proof applies the following two results, whose proofs will be given later.

1. $F(\mu)$ is Lipschitz, which indicates the continuity.
2. $\lim_{\|\mu\| \rightarrow \infty} F(\mu) = 0$.

From F 's definition, it is trivial to see that there exists some $\mu' \in \mathbb{R}^d$ such that $F(\mu') > 0$. Define $\epsilon' := F(\mu')$. From point 2, there exists $M > 0$ such that, whenever $\|\mu\| > M$, $F(\mu) < \epsilon'$.

Now, consider the closed ball $B(\mathbf{0}; M') := \{\mu \in \mathbb{R}^d : \|\mu\| \leq M'\}$, where $M' > M$ and $\mu' \in B(\mathbf{0}; M')$. From the extreme value theorem, the continuity of F and the compactness of $B(\mathbf{0}; M')$ indicate that F attains a maximum μ^* on $B(\mathbf{0}; M')$. That is, $\mu^* \in \mathbb{R}^d$, and $F(\mu^*) \geq F(\mu)$ for any $\mu \in B(\mathbf{0}; M')$. Therefore,

$$F(\mu^*) \geq F(\mu') = \epsilon' > F(\mu), \text{ for any } \mu \notin B(\mathbf{0}; M').$$

This indicates that μ^* is the global maximum point of F in \mathbb{R}^d .

Finally, it remains to prove Point 1 and 2.

Proof to Point 1. From Lemma 3.3, the gradient ∇F exists. Hence, from the mean value theorem (e.g., Theorem 5.10 in (Magnus & Neudecker, 2019)), for any $\mu_1, \mu_2 \in \mathbb{R}^d$, there exists v that lies on the line segment connecting μ_1 and μ_2 that

$$\begin{aligned} |F(\mu_1) - F(\mu_2)| &= |\nabla F(v)(\mu_1 - \mu_2)| \\ &\leq \|\nabla F(v)\| \|\mu_1 - \mu_2\|, \text{ Cauchy Schwarz Inequality,} \\ &= \sigma^{-2} \|\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(v, \sigma^2 I_d)}[(\mathbf{x} - \mu)f_N(\mathbf{x})]\| \|\mu_1 - \mu_2\|, \text{ proof same as (6),} \\ &\leq \frac{\sqrt{2d}}{\sigma\sqrt{\pi}} f_N(\mathbf{x}^*) \|\mu_1 - \mu_2\|, \end{aligned}$$

where the last inequality is derived using the last inequality derived in the proof for Lemma 3.3. This proves Point 1.

Proof to Point 2. Since $S \in \mathbb{R}^d$ is compact, it is bounded. Assume that $S \subset B(\mathbf{0}; M)$ for some $M > 0$. Then, $f_N(\mathbf{x}) = 0$ if $\|\mathbf{x}\| > M$. Hence,

$$\begin{aligned} F(\mu) &= \frac{1}{(\sqrt{2\pi}\sigma)^d} \int_{\|\mathbf{x}\| \leq M} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \mu\|^2}{2\sigma^2}} d\mathbf{x} \\ &\leq \frac{f_N(\mathbf{x}^*)}{(\sqrt{2\pi}\sigma)^d} \int_{\|\mathbf{x}\| \leq M} e^{-\frac{(\|\mu\| - M)^2}{2\sigma^2}} d\mathbf{x}, \text{ whenever } \|\mu\| > M, \\ &\rightarrow 0, \quad \text{as } \|\mu\| \rightarrow \infty, \end{aligned}$$

where the second line is because $\|\mathbf{x} - \boldsymbol{\mu}\| \geq \|\boldsymbol{\mu}\| - \|\mathbf{x}\| \geq \|\boldsymbol{\mu}\| - M$, which further implies $-\|\mathbf{x} - \boldsymbol{\mu}\|^2 \leq -(\|\boldsymbol{\mu}\| - M)^2$ if $\|\boldsymbol{\mu}\| > M$. \square

C.3. Proof to Lemma 3.3

Lemma 3.3 *Under Assumption 3.1, $F_{N,\sigma}(\boldsymbol{\mu})$ and $\nabla F_{N,\sigma}(\boldsymbol{\mu})$ are well-defined for all $\boldsymbol{\mu} \in \mathbb{R}^d$.*

Proof. By the definition of $F_{N,\sigma}(\boldsymbol{\mu})$ in Theorem 2.1, $F_{N,\sigma}(\boldsymbol{\mu}) = \mathbb{E}[f_N(\mathbf{x})]$ and $\nabla F_{N,\sigma}(\boldsymbol{\mu}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})f_N(\mathbf{x})]$, where the expectation is taken with respect to $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I_d)$. To prove that $F_{N,\sigma}(\boldsymbol{\mu})$ and $\nabla F_{N,\sigma}(\boldsymbol{\mu})$ are well-defined, it suffices to prove that $\mathbb{E}|f_N(\mathbf{x})| < +\infty$ and $\mathbb{E}|(x_i - \mu_i)f_N(\mathbf{x})| < +\infty$, for each $i \in \{1, 2, \dots, d\}$. The former inequality holds since $|f_N(\mathbf{x})|$ is bounded by Assumption 3.1 (the condition that f is continuous on its compact domain implies the boundedness of $|f_N(\mathbf{x})|$). The latter inequality holds because

$$\begin{aligned} & \mathbb{E}|(x_i - \mu_i)f_N(\mathbf{x})| \\ & \leq f_N(\mathbf{x}^*)\mathbb{E}|x_i - \mu_i| \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} |x_i - \mu_i| e^{-\frac{(x_i - \mu_i)^2}{2\sigma^2}} dx_i \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} |y| e^{-\frac{y^2}{2\sigma^2}} dy, \quad y := x_i - \mu_i, \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_0^{\infty} 2ye^{-\frac{y^2}{2\sigma^2}} dy, \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_0^{\infty} e^{-\frac{y^2}{2\sigma^2}} dy^2, \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_0^{\infty} e^{-\frac{t}{2\sigma^2}} dt, \quad t := y^2, \\ & = \frac{f_N(\mathbf{x}^*)}{\sqrt{2\pi}\sigma} \int_0^{\infty} (-2\sigma^2) de^{-\frac{t}{2\sigma^2}} dt \\ & = \sqrt{2}f_N(\mathbf{x}^*)\sigma/\sqrt{\pi} < +\infty. \end{aligned}$$

This finishes the proof for Lemma 3.3. \square

C.4. Proof to Lemma 3.4

Lemma 3.4 *Under Assumption 3.1, any stationary point of $F(\boldsymbol{\mu})$ in \mathbb{R}^d belongs to $\mathcal{S}_{\mathbf{x}^*,\delta}$, which is defined in (7).*

Remark C.1. Proposition 2.3 shows that $F(\boldsymbol{\mu})$ has at least one stationary point.

Proof. For any point $\boldsymbol{\mu} \notin \mathcal{S}_{\mathbf{x}^*,\delta}$, we show that $\nabla F(\boldsymbol{\mu}) \neq 0$.

On one hand, $\boldsymbol{\mu} \in \mathcal{S}_M - \mathcal{S}_{\mathbf{x}^*,\delta}$, then $\|\boldsymbol{\mu}\| \leq \sqrt{d}M$ and $\nabla F(\boldsymbol{\mu}) \neq 0$ because of the definition of N in Section 3.2.

On the other hand, if $\boldsymbol{\mu} \notin \mathcal{S}_M$, there is at least one j such that $|\mu_j| > M$. Then,

$$\begin{aligned} \frac{\partial F(\boldsymbol{\mu})}{\partial \mu_j} &= \frac{\int_{\mathbf{x} \in \mathbb{R}^d} (x_j - \mu_j) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} f_N(\mathbf{x}) d\mathbf{x}}{(\sqrt{2\pi})^d \sigma^{d+2}} \\ &= \frac{\int_{\mathbf{x} \in \mathcal{S}} (x_j - \mu_j) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} f_N(\mathbf{x}) d\mathbf{x}}{(\sqrt{2\pi})^d \sigma^{d+2}}, \text{ by (4),} \\ &= \begin{cases} \text{negative,} & \text{if } \mu_j > M; \\ \text{positive,} & \text{if } \mu_j < -M, \end{cases} \end{aligned}$$

where the last equality is because $x_j \in \mathcal{S} \Rightarrow |x_j| \leq M$ by Assumption 3.1. In sum, for any point $\boldsymbol{\mu} \notin \mathcal{S}_{\mathbf{x}^*,\delta}$, $\nabla F(\boldsymbol{\mu}) \neq 0$, which further implies that any stationary point of $F(\boldsymbol{\mu})$ belongs to $\mathcal{S}_{\mathbf{x}^*,\delta}$. \square

C.5. Proof to Lemma 3.5

Lemma 3.5 Under Assumption 3.1, for any $\sigma > 0$, the objective function $F_{N,\sigma}(\boldsymbol{\mu})$ is Lipschitz Smooth. That is, for any $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^d$,

$$\|\nabla F_{N,\sigma}(\boldsymbol{\mu}_1) - \nabla F_{N,\sigma}(\boldsymbol{\mu}_2)\| \leq L\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|,$$

where $L = 2d\sigma^{-2}f_N(\mathbf{x}^*)$, and $f_N(\mathbf{x}^*) = f^N(\mathbf{x}^*)$ for the case of PGS and $f_N(\mathbf{x}^*) = e^{Nf(\mathbf{x}^*)}$ for the case of EPGS.

Proof. The idea is to bound each entry of the difference, a vector in \mathbb{R}^d , on the left-hand side (LHS). For each $i \in \{1, 2, \dots, d\}$, define $g_i(\boldsymbol{\mu})$ as the i^{th} entry of $\nabla F_{N,\sigma}(\boldsymbol{\mu}) \in \mathbb{R}^d$. Then, for the i^{th} entry of the difference on LHS,

$$\begin{aligned} |g_i(\boldsymbol{\mu}_1) - g_i(\boldsymbol{\mu}_2)| &= |\nabla g_i(\mathbf{v})(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)|, \text{ by Mean Value Theorem (e.g., Theorem 5.10 in (Magnus \& Neudecker, 2019))}, \\ &\leq \|\nabla g_i(\mathbf{v})\| \|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|, \text{ by Cauchy Schwarz Inequality,} \\ &\leq 2\sqrt{d}\sigma^{-2}f_N(\mathbf{x}^*)\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|, \text{ proof given later,} \end{aligned}$$

where \mathbf{v} is some point on the line segment connecting $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$, and $\nabla g_i(\mathbf{v})(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$ denotes the inner product of the two vectors in \mathbb{R}^d . This further implies

$$\|\nabla F_{N,\sigma}(\boldsymbol{\mu}_1) - \nabla F_{N,\sigma}(\boldsymbol{\mu}_2)\| \leq 2d\sigma^{-2}f_N(\mathbf{x}^*)\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|,$$

which is the desired result.

Now, it remains to prove for each $i \in \{1, 2, \dots, d\}$ and each $\boldsymbol{\mu} \in \mathbb{R}^d$ that $\|\nabla g_i(\boldsymbol{\mu})\| \leq 2\sqrt{d}\sigma^{-2}f_N(\mathbf{x}^*)$. It suffices to derive a bound for each entry of $\nabla g_i(\boldsymbol{\mu}) \in \mathbb{R}^d$. Specifically, for its j^{th} entry:

$$\begin{aligned} \left| \frac{\partial g_i(\mathbf{u})}{\partial \mu_j} \right| &=_{\text{by (6)}} \left| \frac{1}{(\sqrt{2\pi})^d \sigma^{d+2}} \frac{\partial \int_{\mathbf{x} \in \mathbb{R}^d} (x_i - \mu_i) f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}}{\partial \mu_j} \right| \\ &= \left| \frac{-\delta_{ij} \int_{\mathbf{x} \in \mathbb{R}^d} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x} + \sigma^{-2} \int_{\mathbf{x} \in \mathbb{R}^d} (x_i - \mu_i)(x_j - \mu_j) f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}}{(\sqrt{2\pi})^d \sigma^{d+2}} \right| \\ &\leq \sigma^{-2} \delta_{ij} \mathbb{E}[f_N(\mathbf{x})] + \frac{\sigma^{-2} f_N(\mathbf{x}^*) \int_{\mathbf{x} \in \mathbb{R}^d} |x_i - \mu_i| |x_j - \mu_j| e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}}{(\sqrt{2\pi})^d \sigma^{d+2}}, \quad f_N(\mathbf{x}) \geq 0 \text{ by (4)}, \\ &\leq \sigma^{-2} \delta_{ij} f_N(\mathbf{x}^*) + \frac{\sigma^{-2} f_N(\mathbf{x}^*) \int_{\mathbf{x} \in \mathbb{R}^d} ((x_i - \mu_i)^2 + (x_j - \mu_j)^2) e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}} d\mathbf{x}}{2(\sqrt{2\pi})^d \sigma^{d+2}} \\ &= \sigma^{-2} \delta_{ij} f_N(\mathbf{x}^*) + \frac{\sigma^{-2} f_N(\mathbf{x}^*) 2 \int_{x_j \in \mathbb{R}} (x_j - \mu_j)^2 e^{-\frac{(x_j - \mu_j)^2}{2\sigma^2}} dx_j}{2\sqrt{2\pi}\sigma^3} \\ &= \sigma^{-2} \delta_{ij} f_N(\mathbf{x}^*) + f_N(\mathbf{x}^*) \sigma^{-4} \text{Var}(x_j) \\ &= \sigma^{-2} (\delta_{ij} + 1) f_N(\mathbf{x}^*) \\ &\leq 2\sigma^{-2} f_N(\mathbf{x}^*), \end{aligned}$$

where the interchange of differentiation and integral in the second line is by Lebesgue's dominated convergence theorem, and the fifth line is because of the separability of the multivariate integral. Here, $\delta_{ij} = 1$ if $i = j$, and $\delta_{ij} = 0$ if $i \neq j$. The above bound for $\left| \frac{\partial g_i(\mathbf{u})}{\partial \mu_j} \right|$ implies $\|\nabla g_i(\boldsymbol{\mu})\| \leq 2\sqrt{d}\sigma^{-2}f_N(\mathbf{x}^*)$. This finishes the proof for Lemma 3.5. \square

C.6. Proof to Lemma 3.6

Lemma 3.6 Under Assumption 3.1, for any $\sigma > 0$, let $\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu})$ be defined in (5). Then, $\mathbb{E}[\|\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu})\|^2] < G$, where

$$G = \begin{cases} d\sigma^2 f^{2N}(\mathbf{x}^*), & \text{for PGS;} \\ d\sigma^2 e^{2Nf(\mathbf{x}^*)}, & \text{for EPGS.} \end{cases}$$

Proof.

$$\begin{aligned}
 \mathbb{E} \left[\|\hat{\nabla} F_{N,\sigma}(\boldsymbol{\mu}_t)\|^2 \right] &= \frac{1}{K^2} \sum_{k=1}^K \sum_{l=1}^K \mathbb{E}[(\mathbf{x}_k - \boldsymbol{\mu}_t)(\mathbf{x}_l - \boldsymbol{\mu}_t) f_N(\mathbf{x}_k) f_N(\mathbf{x}_l)] \\
 &\leq f_N^2(\mathbf{x}^*) \frac{1}{K^2} \sum_{k=1}^K \sum_{l=1}^K \mathbb{E}[(\mathbf{x}_k - \boldsymbol{\mu}_t)(\mathbf{x}_l - \boldsymbol{\mu}_t)] \\
 &\leq f_N^2(\mathbf{x}^*) \frac{1}{K^2} \sum_{k=1}^K \sum_{l=1}^K \mathbb{E}[\|(\mathbf{x}_k - \boldsymbol{\mu}_t)\| \|(\mathbf{x}_l - \boldsymbol{\mu}_t)\|], \\
 &\leq f_N^2(\mathbf{x}^*) \frac{1}{K^2} \sum_{k=1}^K \sum_{l=1}^K \sqrt{\mathbb{E}[\|\mathbf{x}_k - \boldsymbol{\mu}_t\|^2] \mathbb{E}[\|\mathbf{x}_l - \boldsymbol{\mu}_t\|^2]} \\
 &= f_N^2(\mathbf{x}^*) \sigma^2 d, \quad d \text{ denotes the number of dimensions,}
 \end{aligned}$$

where $(\mathbf{x}_k - \boldsymbol{\mu}_t)(\mathbf{x}_l - \boldsymbol{\mu}_t)$ denotes the inner product of the two vectors, the expectation \mathbb{E} is over the random vectors $\mathbf{x}_k, \mathbf{x}_l \sim \mathcal{N}(\boldsymbol{\mu}_t, \sigma^2 I_d)$, and the fourth line is from the application of the Cauchy-Schwarz inequality to expectations. Replacing $f_N^2(\mathbf{x}^*) = e^{2Nf(\mathbf{x}^*)}$ for EPGS and $f_N^2(\mathbf{x}^*) = f^{2N}(\mathbf{x}^*)$ for PGS gives the desired result. \square

C.7. Proof to Theorem 3.7

Theorem 3.7 *Let $\{\boldsymbol{\mu}_t\}_{t=0}^T \subset \mathbb{R}^d$ be produced by following the iteration rule of (5), with a pre-selected and deterministic $\boldsymbol{\mu}_0$ and all the involved terms defined as in Section 3.2. Then, under Assumption 3.1 and 3.2, we have that*

$$\sum_{t=0}^{T-1} \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] \leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + LG \sum_{t=0}^{\infty} \alpha_t^2,$$

where L and G are as defined in Lemma 3.5 and Lemma 3.6, respectively.

Proof. By the Gradient Mean Value Theorem, there exists $\boldsymbol{\nu}_t \in \mathbb{R}^d$ such that $\nu_{t,i}$ lies between $\mu_{t+1,i}$ and $\mu_{t,i}$ for each of the i th entry, and

$$\begin{aligned}
 F(\boldsymbol{\mu}_{t+1}) &= F(\boldsymbol{\mu}_t) + (\nabla F(\boldsymbol{\nu}_t))'(\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t), \\
 &= F(\boldsymbol{\mu}_t) + (\nabla F(\boldsymbol{\mu}_t))'(\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t) + (\nabla F(\boldsymbol{\nu}_t) - \nabla F(\boldsymbol{\mu}_t))'(\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t) \\
 &= F(\boldsymbol{\mu}_t) + \alpha_t (\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t)) - (\nabla F(\boldsymbol{\mu}_t) - \nabla F(\boldsymbol{\nu}_t))'(\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t) \\
 &\geq F(\boldsymbol{\mu}_t) + \alpha_t (\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t)) - L \|\boldsymbol{\nu}_t - \boldsymbol{\mu}_t\| \cdot \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t\|, \text{ by Lemma 3.5,} \\
 &\geq F(\boldsymbol{\mu}_t) + \alpha_t (\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t)) - L \|\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t\|^2, \\
 &= F(\boldsymbol{\mu}_t) + \alpha_t (\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t)) - \alpha_t^2 L \|\hat{\nabla} F_N(\boldsymbol{\mu}_t)\|^2.
 \end{aligned}$$

where $'$ denotes the vector transpose, and the second inequality is because $\nu_{t,i}$ is between $\mu_{t+1,i}$ and $\mu_{t,i}$. Taking the expectation of the left-end and right-end of the above derived inequality gives

$$\begin{aligned}
 \mathbb{E}[F(\boldsymbol{\mu}_{t+1})] &\geq \mathbb{E}[F(\boldsymbol{\mu}_t)] + \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] - \alpha_t^2 L \mathbb{E}[\|\hat{\nabla} F_N(\boldsymbol{\mu}_t)\|^2] \\
 &\geq \mathbb{E}[F(\boldsymbol{\mu}_t)] + \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] - \alpha_t^2 LG,
 \end{aligned} \tag{16}$$

where the second inequality is because of Lemma 3.6, and for the first inequality, note that

$$\mathbb{E}[(\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t))] = \mathbb{E} \left[\mathbb{E}[(\nabla F(\boldsymbol{\mu}_t))'(\hat{\nabla} F(\boldsymbol{\mu}_t)) | \boldsymbol{\mu}_t] \right] \stackrel{\text{by (6)}}{=} \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2].$$

Taking the sum from $t = 0$ to $t = T - 1$ on both sides of (16) gives

$$\mathbb{E}[F(\boldsymbol{\mu}_T)] \geq \mathbb{E}[F(\boldsymbol{\mu}_0)] + \sum_{t=0}^{T-1} \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] - LG \sum_{t=0}^{T-1} \alpha_t^2.$$

Re-organizing the terms gives

$$\begin{aligned} \sum_{t=0}^{T-1} \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] &\leq \mathbb{E}[F(\boldsymbol{\mu}_T)] - \mathbb{E}[F(\boldsymbol{\mu}_0)] + LG \sum_{t=0}^{T-1} \alpha_t^2 \\ &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + LG \sum_{t=0}^{\infty} \alpha_t^2 \end{aligned}$$

This finishes the proof for Theorem 3.7. \square

C.8. Proof to Corollary 3.9

Corollary 3.9 Suppose Assumption 3.1 and 3.2 hold. Let $\{\boldsymbol{\mu}_t\}$ be produced by the stochastic gradient ascent rule (5) of GS-PowerOpt, with a pre-selected and deterministic $\boldsymbol{\mu}_0$. Then, for any $\varepsilon \in (0, 1)$, after

$$T > (C_1 C_2 d^2 \varepsilon^{-1})^{2/(1-2\gamma)} = O_N((d^2 \varepsilon^{-1})^{2/(1-2\gamma)})$$

times of $\boldsymbol{\mu}_t$ -updating by (5), we have that $\min_{t \in \{0, 1, \dots, T\}} \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] < \varepsilon$. Here, $\gamma \in (0, 1/2)$ is a parameter in the learning rate $\alpha_t := (t+1)^{-(1/2+\gamma)}$, $C_0 = f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + 2f_N^3(\mathbf{x}^*) \sum_{t=1}^{\infty} t^{-(1+2\gamma)}$, $C_1 = C_0(1-2\gamma)\sigma^{-2}$, and $C_2 = \max\{1, 2/C_1\}$.

Proof. For any non-negative integer t , define

$$\nu_t := \min_{\tau \in \{0, 1, \dots, t\}} \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_\tau)\|^2]. \quad (17)$$

Then,

$$\begin{aligned} \sum_{t=0}^{T-1} \alpha_t \sigma^2 \nu_T &\leq \sum_{t=0}^{T-1} \alpha_t \sigma^2 \nu_t, \quad \text{since } \nu_{t+1} \leq \nu_t \text{ for each } t \geq 0, \\ &\leq \sum_{t=0}^{T-1} \alpha_t \sigma^2 \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2], \quad \text{since } \nu_t \leq \mathbb{E}[\|\nabla F(\boldsymbol{\mu}_t)\|^2] \text{ from (17),} \\ &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + LG \sum_{t=0}^{\infty} \alpha_t^2, \quad \text{from Theorem 3.7,} \\ &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + 2d^2 f_N^3(\mathbf{x}^*) \sum_{t=1}^{\infty} t^{-(1+2\gamma)}, \quad d \text{ comes from } L \text{ in Lemma 3.5 and } G \text{ in Lemma 3.6,} \\ &\leq \left(f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + 2f_N^3(\mathbf{x}^*) \sum_{t=1}^{\infty} t^{-(1+2\gamma)} \right) d^2, \quad \text{since } d > 1 \text{ and } f_N(\mathbf{x}^*) \geq F(\boldsymbol{\mu}_0), \\ &= C_0 d^2, \quad \text{where } C_0 := f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + 2f_N^3(\mathbf{x}^*) \sum_{t=1}^{\infty} t^{-(1+2\gamma)} < \infty. \end{aligned} \quad (18)$$

Therefore, we have

$$\sum_{t=0}^{T-1} \alpha_t \sigma^2 \nu_T \leq C_0 d^2.$$

Since $\alpha_t = (t+1)^{-(1/2+\gamma)}$, dividing both the left-end and the right-end of the above inequality by $\sum_{t=0}^{T-1} \alpha_t \sigma^2$ gives

$$\begin{aligned}
 \nu_T &\leq \frac{C_0 d^2}{\sigma^2 \sum_{t=1}^T t^{-(1/2+\gamma)}} \\
 &< \frac{C_0 d^2}{\sigma^2 \int_1^T t^{-(1/2+\gamma)} dt} \\
 &= \frac{C_0 d^2}{\sigma^2 (T^{\frac{1}{2}-\gamma} - 1) / (\frac{1}{2} - \gamma)} \\
 &< \frac{C_0 d^2}{\sigma^2 (T^{\frac{1}{2}-\gamma} / 2) / (\frac{1}{2} - \gamma)}, \quad \text{when } T > 2^{2/(1-2\gamma)}, \\
 &= C_1 \frac{d^2}{T^{\frac{1}{2}-\gamma}}, \quad \text{where } C_1 := \frac{C_0}{\sigma^2 (1/2) / (\frac{1}{2} - \gamma)}.
 \end{aligned}$$

In sum, we have

$$\nu_T \leq C_1 \frac{d^2}{T^{\frac{1}{2}-\gamma}}, \quad \text{whenever } T > 2^{2/(1-2\gamma)}. \quad (19)$$

Define $C_2 := \max\{1, 2/C_1\}$. Given any $\varepsilon \in (0, 1)$, whenever $T > (C_2 C_1 d^2 \varepsilon^{-1})^{2/(1-2\gamma)} = O((d^2 \varepsilon^{-1})^{2/(1-2\gamma)})$, we have

$$T > (C_2 C_1 d^2 \varepsilon^{-1})^{2/(1-2\gamma)} > (C_1 C_2)^{2/(1-2\gamma)} \geq 2^{2/(1-2\gamma)},$$

and

$$\nu_T \stackrel{\text{from (19)}}{\leq} C_1 \frac{d^2}{T^{\frac{1}{2}-\gamma}} < C_1 \frac{d^2}{(C_2 C_1 d^2 \varepsilon^{-1})^{\frac{2}{1-2\gamma}(\frac{1}{2}-\gamma)}} = \frac{\varepsilon}{C_2} \leq \varepsilon.$$

This implies that after $T = O((d^2 \varepsilon^{-1})^{2/(1-2\gamma)})$ times of updating μ_t by GS-PowerOpt, $v_T = \min_{t \in \{0, 1, 2, \dots, T\}} \mathbb{E}[\|\nabla F(\mu_t)\|^2] < \varepsilon$. \square

C.9. Proof to Proposition 3.11

Proposition 3.11 Assume Assumption 3.1, $f(\mathbf{x}) = 0$ on the boundary of f 's domain, the Lipschitz condition on f : $|f(\mathbf{x}) - f(\mathbf{y})| \leq L_0 \|\mathbf{x} - \mathbf{y}\|$, and the Lipschitz condition on ∇f : $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L_1 \|\mathbf{x} - \mathbf{y}\|$. Then, the following two statements are true.

1. $\|\nabla F_{N,\sigma}(\mu_1) - \nabla F_{N,\sigma}(\mu_2)\| \leq L \|\mu_1 - \mu_2\|$ holds for any $\mu_1, \mu_2 \in \mathbb{R}^d$, where $L = (N f^{N-1}(\mathbf{x}^*) L_1 + N(N-1) L_0^2 f^{N-2}(\mathbf{x}^*))$ under the PGS setting, and $L = (N L_1 e^{N f(\mathbf{x}^*)} + N^2 L_0^2 e^{N f(\mathbf{x}^*)})$ under the EPGS setting.
2. The iteration complexity in Corollary 3.9 reduces to $O_N((d \varepsilon^{-1})^{2/(1-2\gamma)})$.

Proof to Point 1. We extend the domain of f from \mathcal{S} to \mathbb{R}^d such that $f(\mathbf{x}) = 0$ for any $\mathbf{x} \notin \mathcal{S}$. Clearly, after the domain extension, the two Lipschitz conditions hold in \mathbb{R}^d . From (4), with the extended f ,

$$f_N(\mathbf{x}) = \begin{cases} f^N(\mathbf{x}), & \text{(PGS setting);} \\ e^{N f(\mathbf{x})}, & \text{(EPGS setting).} \end{cases}$$

By definition,

$$\begin{aligned}
 F_{N,\sigma}(\mu) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \sigma I_d)}[f_N(\mathbf{x})] \\
 &= (\sqrt{2\pi}\sigma)^{-d} \int_{\mathbf{x} \in \mathbb{R}^d} f_N(\mathbf{x}) e^{-\frac{\|\mathbf{x} - \mu\|^2}{2\sigma^2}} d\mathbf{x} \\
 &= (\sqrt{2\pi})^{-d} \int_{\xi \in \mathbb{R}^d} f_N(\mu + \sigma \xi) e^{-\frac{\|\xi\|^2}{2}} d\xi, \quad \mathbf{x} := \mu + \sigma \xi, \\
 &= \mathbb{E}_{\xi \sim \mathcal{N}(\mathbf{0}, I_d)}[f_N(\mu + \sigma \xi)].
 \end{aligned}$$

Let ∇ denotes the first derivative with respect to μ . Under the PGS setting, $\nabla f_N(\mu + \sigma\xi) = Nf^{N-1}(\mu + \sigma\xi)\nabla f(\mu + \sigma\xi)$. Then,

$$\begin{aligned}
 & \|\nabla F_{N,\sigma}(\mu_1) - \nabla F_{N,\sigma}(\mu_2)\| \\
 &= \|\mathbb{E}[\nabla f_N(\mu_1 + \sigma\xi) - \nabla f_N(\mu_2 + \sigma\xi)]\| \\
 &= \|\mathbb{E}[Nf^{N-1}(\mu_{1,\xi})\nabla f_N(\mu_{1,\xi}) - Nf^{N-1}(\mu_{2,\xi})\nabla f(\mu_{2,\xi})]\| \\
 &= N\|\mathbb{E}[f^{N-1}(\mu_{1,\xi})(\nabla f_N(\mu_{1,\xi}) - \nabla f_N(\mu_{2,\xi})) + (f^{N-1}(\mu_{1,\xi}) - f^{N-1}(\mu_{2,\xi}))\nabla f(\mu_{2,\xi})]\| \\
 &\leq Nf^{N-1}(\mathbf{x}^*)\mathbb{E}\|\nabla f_N(\mu_{1,\xi}) - \nabla f_N(\mu_{2,\xi})\| + NL_0\mathbb{E}|f^{N-1}(\mu_{1,\xi}) - f^{N-1}(\mu_{2,\xi})| \\
 &\leq Nf^{N-1}(\mathbf{x}^*)L_1\|\mu_1 - \mu_2\| + N(N-1)f^{N-2}(\mathbf{x}^*)L_0\mathbb{E}|f(\mu_{1,\xi}) - f(\mu_{2,\xi})| \\
 &\leq Nf^{N-1}(\mathbf{x}^*)L_1\|\mu_1 - \mu_2\| + N(N-1)f^{N-2}(\mathbf{x}^*)L_0^2\|\mu_1 - \mu_2\| \\
 &\leq (Nf^{N-1}(\mathbf{x}^*)L_1 + N(N-1)L_0^2f^{N-2}(\mathbf{x}^*))\|\mu_1 - \mu_2\|,
 \end{aligned}$$

where $\mu_{1,\xi} := \mu_1 + \sigma\xi$, $\mu_{2,\xi} := \mu_2 + \sigma\xi$, the first inequality sign is derived using the formula¹⁰ of $\|\mathbb{E}x\| \leq \mathbb{E}\|x\|$ and the result that $\|\nabla f\| \leq L_0$ from the Lipschitz condition on f , and the second inequality sign is derived using the formula $a^{N-1} - b^{N-1} = (a-b)\sum_{n=0}^{N-2} a^n b^{N-2-n}$.

Under the EPGS setting, $\nabla f_N(\mu + \sigma\xi) = Ne^{Nf(\mu + \sigma\xi)}\nabla f(\mu + \sigma\xi)$. Then,

$$\begin{aligned}
 & \|\nabla F_{N,\sigma}(\mu_1) - \nabla F_{N,\sigma}(\mu_2)\| \\
 &= \|\mathbb{E}[\nabla f_N(\mu_1 + \sigma\xi) - \nabla f_N(\mu_2 + \sigma\xi)]\| \\
 &= \|\mathbb{E}[Ne^{Nf(\mu_{1,\xi})}\nabla f_N(\mu_{1,\xi}) - Ne^{Nf(\mu_{2,\xi})}\nabla f(\mu_{2,\xi})]\| \\
 &= N\|\mathbb{E}[e^{Nf(\mu_{1,\xi})}(\nabla f_N(\mu_{1,\xi}) - \nabla f_N(\mu_{2,\xi})) + (e^{Nf(\mu_{1,\xi})} - e^{Nf(\mu_{2,\xi})})\nabla f(\mu_{2,\xi})]\| \\
 &\leq Ne^{Nf(\mathbf{x}^*)}\mathbb{E}\|\nabla f_N(\mu_{1,\xi}) - \nabla f_N(\mu_{2,\xi})\| + NL_0\mathbb{E}|e^{Nf(\mu_{1,\xi})} - e^{Nf(\mu_{2,\xi})}| \\
 &\leq Ne^{Nf(\mathbf{x}^*)}L_1\|\mu_1 - \mu_2\| + NL_0Ne^{Nf(\mathbf{x}^*)}L_0\|\mu_1 - \mu_2\|, \text{ by mean-value theorem,} \\
 &\leq (NL_1e^{Nf(\mathbf{x}^*)} + N^2L_0^2e^{Nf(\mathbf{x}^*)})\|\mu_1 - \mu_2\|.
 \end{aligned}$$

□

Proof to Point 2. Point 2 can be proved from a modified version of the proof to Corollary 3.9. Specifically, when deriving (18), replacing Lemma 3.5 by Point 1 of Proposition 3.11 removes the dependence of L on d , which decreases the power of d by 1 on the right-end of (18). This in turn replaces d^2 by d in the rest of the proof to Corollary 3.9, and finally leads to the desired iteration complexity in Point 2 of Proposition 3.11. Note that C_0 needs to be redefined accordingly in (18). □

C.10. Proof to Corollary 4.2

Corollary 4.2 Under Assumption 3.1 and 4.1, the iteration complexity in Corollary 3.9 becomes $O((d^2\varepsilon^{-1})^{2/(1-2\gamma)})$, and the iteration complexity in Point 2 of Proposition 3.11 becomes $O((d\varepsilon^{-1})^{2/(1-2\gamma)})$, both of which are independent from N .

Proof. The dependence of the iteration complexity $O_N((d^2\varepsilon^{-1})^{2/(1-2\gamma)})$ in Corollary 3.9 on N comes from $C_0 = f_N(\mathbf{x}^*) - F(\mu_0) + 2f_N^3(\mathbf{x}^*)\sum_{t=1}^{\infty} t^{-(1+2\gamma)}$. Assumption 4.1 guarantees that Corollary 3.9 still holds true if C_0 is redefined as

$$C_0 = 1 + 2\sum_{t=1}^{\infty} t^{-(1+2\gamma)}.$$

This makes the iteration complexity in Corollary 3.9 become $O((d^2\varepsilon^{-1})^{2/(1-2\gamma)})$, which is independent from N .

To prove that Corollary 3.9 still holds true with the above new definition of C_0 , we only need to modify the last line of (18) in the proof to Corollary 3.9:

$$\left(f_N(\mathbf{x}^*) - F(\mu_0) + 2f_N^3(\mathbf{x}^*)\sum_{t=1}^{\infty} t^{-(1+2\gamma)}\right)d^2 \leq \left(1 + 2\sum_{t=1}^{\infty} t^{-(1+2\gamma)}\right)d^2 = C_0d^2,$$

¹⁰This inequality can be derived from Jensen's inequality and the fact that the Euclidean norm $\|\cdot\|$ is convex.

where the inequality can be justified by $f_N(\mathbf{x}^*) \in [0, 1)$ (a result immediately implied by Assumption 4.1).

Next, we prove that the iteration complexity in Point 2 of Proposition 3.11 becomes $O((d\varepsilon^{-1})^{2/(1-2\gamma)})$ under Assumption 4.1. Let ν_t be defined as in (17).

$$\begin{aligned}
 \sum_{t=0}^{T-1} \alpha_t \sigma^2 \nu_T &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + LG \sum_{t=0}^{\infty} \alpha_t^2, \quad \text{from the third line in (18),} \\
 &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + Ld\sigma^2 f_N^2(\mathbf{x}^*) \sum_{t=0}^{\infty} \alpha_t^2, \quad \text{from Lemma 3.6,} \\
 &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + Ld\sigma^2 \sum_{t=0}^{\infty} \alpha_t^2, \quad \text{since } f_N(\mathbf{x}^*) \in [0, 1), \\
 &\leq f_N(\mathbf{x}^*) - F(\boldsymbol{\mu}_0) + B(L_1 + L_0^2)d\sigma^2 \sum_{t=0}^{\infty} \alpha_t^2, \quad \text{from (20) given below,} \\
 &\leq (1 + B(L_1 + L_0^2)\sigma^2 \sum_{t=0}^{\infty} \alpha_t^2)d, \\
 &= C_0 d, \quad \text{where } C_0 := 1 + B(L_1 + L_0^2)\sigma^2 \sum_{t=0}^{\infty} \alpha_t^2.
 \end{aligned}$$

In sum, we have $\sum_{t=0}^{T-1} \alpha_t \sigma^2 \nu_T \leq C_0 d$, where $C_0 = 1 + B(L_1 + L_0^2)\sigma^2 \sum_{t=0}^{\infty} \alpha_t^2$ is independent from N . This result, and the rest of the proof (where the term d^2 is replaced by d) to Corollary 3.9 after Eq. (18), gives the desired iteration complexity of $O((d\varepsilon^{-1})^{2/(1-2\gamma)})$.

Finally, it remains to prove for some constant $B > 0$,

$$L \leq B(L_1 + L_0^2). \quad (20)$$

Under the PGS setting, from Point 1 of Proposition 3.11, $L = (Nf^{N-1}(\mathbf{x}^*)L_1 + N(N-1)L_0^2 f^{N-2}(\mathbf{x}^*))$. Since $f(\mathbf{x}^*) \in [0, 1)$,

$$\lim_{N \rightarrow +\infty} Nf^{N-1}(\mathbf{x}^*) = \lim_{N \rightarrow +\infty} N(N-1)f^{N-2}(\mathbf{x}^*) = 0,$$

which implies both $Nf^{N-1}(\mathbf{x}^*)$ and $N(N-1)f^{N-2}(\mathbf{x}^*)$ are upper bounded by a constant B for all $N \geq 0$. This implies $L \leq B(L_1 + L_0^2)$.

Under the EPGS setting, from Point 1 of Proposition 3.11, $L = (NL_1 e^{Nf(\mathbf{x}^*)} + N^2 L_0^2 e^{Nf(\mathbf{x}^*)})$. Since $f(\mathbf{x}^*) < 0$,

$$\lim_{N \rightarrow +\infty} N e^{Nf(\mathbf{x}^*)} = \lim_{N \rightarrow +\infty} N^2 e^{Nf(\mathbf{x}^*)} = 0,$$

which implies both $N e^{Nf(\mathbf{x}^*)}$ and $N^2 e^{Nf(\mathbf{x}^*)}$ are upper bounded by a constant B for all $N \geq 0$. This implies $L \leq B(L_1 + L_0^2)$. □

D. More Details and Discussions on Compared Algorithms

D.1. Zeroth-order Gradient Algorithms

The two zeroth-order gradient algorithms, ZO-SGD (Ghadimi & Lan, 2013) and ZO-AdaMM (Chen et al., 2019), aim to solve $\max_{\mathbf{x}} \mathbb{E}[\mathbf{x} + \sigma \boldsymbol{\xi}]$. Here, $\boldsymbol{\xi}$ denotes a random noise and we select it to be a standard multi-variate Gaussian vector. For non-concave maximization problems, they are theoretically guaranteed to converge to stationary points of f . In our experiments, we take the scaling parameter σ as a hyper-parameter.

D.2. STD-Homotopy

STD-Homotopy is a standard homotopy algorithm for optimization. It has a double-loop mechanism. The inner loop updates the solution μ_t with a fixed scaling parameter σ until no improvements of $f(\mu_t)$ are made, and the outer loop decays σ iteratively. In this algorithm, the term $\widehat{\nabla F}(\mu, \sigma)$ is an estimate of $\mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma I_d)}[(x - \mu)f(x)] = \sigma^2 \nabla_{\mu} \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma I_d)}[f(x)]$, which is used to update μ .

Algorithm 2 STD-Homotopy

Input: The maximum of iteration number $T_{total} > 0$, the initial scaling parameter $\sigma > 0$, the objective f , the initial value μ_0 , the number K of sampled points for gradient approximation, the maximum number N_{σ} of times σ gets updated, the maximum T_{μ} of the number of times μ gets updated for each value of σ , the tolerance number τ for no improvements of $f(\mu_t)$ for any fixed σ , the decay factor $\gamma \in (0, 1)$, and the learning rate $\alpha_t > 0$.
 Initialize $t_1 = 0, n_{\sigma} = 0$.
while $t_1 \leq T_{total}$ and $n_{\sigma} < N_{\sigma}$ **do**
 $t = 0, I = True$.
 while $t < T_{\mu}$ and $I == True$ and $t_1 \leq T_{total}$ **do**
 Independently and uniformly sample K points $\{v_k\}_{k=1}^K$ from the uniform sphere in \mathbb{R}^d . Compute $x_k := \mu_t + \sigma v_k$, for each k .
 Compute the gradient estimate $\widehat{\nabla F}(\mu_t, \sigma) = \frac{1}{K} \sum_{k=1}^K (x_k - \mu_t)f(x_k)$.
 $\mu_{t+1} = \mu_t + \alpha_t \widehat{\nabla F}_N(\mu_t, \sigma) / \|\widehat{\nabla F}_N(\mu_t, \sigma)\|$.
 if $\max\{f(\mu_{t+1}), f(\mu_t), \dots, f(\mu_{t-\tau+1})\} \leq f(\mu_{t-\tau})$ **then**
 $I = False$.
 end if
 $t_1 = t_1 + 1, t = t + 1$.
 end while
 $\sigma = \gamma \sigma, n_{\sigma} = n_{\sigma} + 1$.
end while
 Return(μ_{t_1}).

D.3. ZO-SGD

The zeroth-order stochastic gradient ascent (ZO-SGD) is a maximize-version of Equation (1) in (Chen et al., 2019), whose gradient estimate method is from (Nesterov & Spokoiny, 2017). The gradient-ascent version of this algorithm is in Algorithm 3.

Algorithm 3 ZO-SGD

Input:: The scaling parameter $\sigma > 0$, the objective f , the initial value $\mu_0 \in \mathbb{R}^d$, the number K of sampled points for gradient approximation, the total number T of μ -updates, and the learning rate schedule $\{\alpha_t\}_{t=1}^T$.
for t from 0 to $T - 1$ **do**
 Independently and uniformly sample K points $\{v_k\}_{k=1}^K$ from the uniform sphere in \mathbb{R}^d .
 Compute the gradient estimate

$$\widehat{\nabla F}(\mu_t, \sigma) = \frac{1}{K} \sum_{k=1}^K \frac{(f(\mu_t + \sigma v_k) - f(\mu_t))v_k d}{\sigma}.$$

$\mu_{t+1} = \mu_t + \alpha_t \widehat{\nabla F}_N(\mu_t, \sigma)$.
end for
 Return(μ_T).

E. Details on the Experiments for Figure 4

The set of N -values is $\{10, 20, \dots, 65\}$ for PGS and $\{1.0, 1.5, 2.0, \dots, 4.5\}$ for EGS. For each N value and each algorithm, we do 100 trials to stabilize the result. In each trial, the initial solution candidate μ_0 is uniformly sampled from $C := \{\mathbf{x} \in \mathbb{R}^d | x_i \in [-1.0, 1.0], i \in \{1, 2, \dots, d\}\}$, where x_i represents the i^{th} entry of \mathbf{x} . We set the initial learning rate as 0.1, the scaling parameter σ as 0.5, and the total number of solution updates as 1000. The objective for PGS is modified to be $f_1(\mathbf{x}) := f(\mathbf{x}) + 10$ to ensure that the PGS agent will not encounter negative fitness values during the 1000 updates.

F. Hyper-parameters

For experiments on the benchmark test functions (Table 1 and 2), the set of hyper-parameter values with the smallest mean square error (averaged over the 100 experiments) between the true and estimated solutions are selected. The set of candidate values, as well as the selected values, are listed in Table 6 and 7.

For the image attacks (Table 3 and 4), for each set of the hyper-parameter candidate values, we randomly choose 10 images to attack. The set with the highest average fitness value (average over the 10 image attacks) will be selected. Table 8 and 9 reports the candidate set and the selected values.

Table 6: Hyper-parameters for Optimizing Ackley. The candidate set for learning rates is $\mathcal{L} := \{.1, .001\}$. The candidate set for smoothing parameters is $\mathcal{S} := \{.1, 1.0, 2.0\}$. t_1 in ZO-SLGHd and ZO-SLGHr is the initial scaling parameter. μ in ZO-AdaMM is the scaling parameter. The set of candidate values that lead to the minimum mean square error between the true solution and the found solution will be selected. For CMA-ES, σ_0 denotes the initial standard deviation.

	Selected Values	Candidates
CMA-ES	$\sigma_0 = 1.5$	$\sigma_0 \in \{.5, 1.0, 1.5\}$.
EPGS	$\alpha_0 = .1, N = 1, \sigma = 1.0$.	$N \in \{1, 2, 3\}, \alpha_0 \in \mathcal{L}, \sigma \in \mathcal{S}$.
PGS	$\alpha_0 = .1, N = 20, \sigma = 1.0$.	$N \in \{5, 10, 20, 30\}, \alpha_0 \in \mathcal{L}, \sigma \in \mathcal{S}$.
STD-Homotopy	$\alpha = .1, \gamma = .5, \sigma = 2.0,$ $T_\mu = 500, \tau = 100, N_\sigma = 10$.	$\gamma \in \{.2, .5, .8\}, \alpha \in \mathcal{L}, \sigma \in \mathcal{S}$.
ZO-SLGHd	$\beta = .001, \eta = .001, t_1 = 2.0, \gamma = .99$	$\beta \in \{.1, .001\}, \eta \in \{.1, .01, .001\}, t_1 \in \mathcal{S}, \gamma \in \{.99, .95\}$.
ZO-SLGHr	$\beta = .001, t_1 = 0.1, \gamma = .995$	$\beta \in \mathcal{L}, \gamma \in \{.999, .995\}, t_1 \in \mathcal{S}$.
ZO-AdaMM	$\beta_1 = .5, \beta_2 = .5, \alpha = 0.1, \mu = 1.0$.	$\alpha \in \mathcal{L}, \beta_1 \in \{.5, .7, .9\}, \beta_2 \in \{.1, .3, .5\}, \mu \in \mathcal{S}$
ZO-SGD	$\alpha = .1, \sigma = 1.0$.	$\alpha \in \mathcal{L}, \mu \in \mathcal{S}$.

G. Adversarial Images

We randomly select four images $\{\mathbf{a}_m\}_{m=1}^4$ from the test image set and perform a per-image adversarial attack on each of them, using each of the compared algorithms. Table 10 and 11 report the adversarial image $\mathbf{a}_m + \mu_m^*$ with the highest R^2 score, given that the attack is successful. That is, $\mu_m^* = \arg \max_{\mu_{m,t} \in \{\mu_{m,t}\}_{t=0}^{T-1}} R^2(\mathbf{a}_m, \mathbf{a}_m + \mu_{m,t})$, where $\{\mu_{m,t}\}_{t=0}^{T-1}$ denote the sequence of perturbations generated during the m^{th} attack, and T is set at the same value (1,500) as in Table 3 and 4. According to the results in Table 10 and 11, the performances of the compared algorithms are consistent with the statistics reported in Table 3 and 4.

Table 7: Hyper-parameters for Optimizing Rosenbrock. The candidate set for learning rates is $\mathcal{L} := \{.2, .1, .01, .001, .0001\}$. The candidate set for smoothing parameters is $\mathcal{S} := \{1.0, 2.0\}$. t_1 in ZO-SLGHD and ZO-SLGHR is the initial scaling parameter. μ in ZO-AdaMM is the scaling parameter. For EPGS and PGS, $\alpha_t = \frac{1,000\alpha_0}{1,000+t}$. The set of candidate values that lead to the minimum mean square error between the true solution and the found solution will be selected.

	Selected Values	Candidates
CMA-ES	$\sigma_0 = .1$.	$\sigma_0 \in \{.1, .5, 1.0\}$.
EPGS	$\alpha_0 = .1, N = 1, \sigma = 1.0$.	$N \in \{1, 2, 3\}, \alpha_0 \in \mathcal{L}, \sigma \in \mathcal{S}$.
PGS	$\alpha_0 = .1, N = 1, \sigma = 1.0$.	$N \in \{1, 3, 5\}, \alpha_0 \in \mathcal{L}, \sigma \in \mathcal{S}$.
STD-Homotopy	$\alpha = .2, \gamma = .2, \sigma = 2.0,$ $T_\mu = 500, \tau = 100, N_\sigma = 10$.	$\gamma \in \{.2, .5, .8\}, \alpha \in \mathcal{L}, \sigma \in \mathcal{S}$.
ZO-SLGHD	$\beta = .0001, \eta = .001, t_1 = 2.0, \gamma = .999$	$\beta \in \mathcal{L}, \eta \in \{.1, .01, .001\}, t_1 \in \mathcal{S}, \gamma \in \{.99, .995, .999\}$.
ZO-SLGHR	$\beta = .0001, t_1 = 2.0, \gamma = .999$.	$\beta \in \mathcal{L}, \gamma \in \{.999, .995\}, t_1 \in \mathcal{S}$.
ZO-AdaMM	$\beta_1 = .5, \beta_2 = .5, \alpha = .2, \mu = 2.0$.	$\alpha \in \mathcal{L}, \beta_1 \in \{.5, .7, .9\}, \beta_2 \in \{.1, .3, .5\}, \mu \in \mathcal{S}$
ZO-SGD	$\alpha = .001, \sigma = 2.0$.	$\alpha \in \mathcal{L}, \mu \in \mathcal{S}$.

Table 8: Hyper-parameters for MNIST Attack. The candidate set for (initial) smoothing parameters is $\mathcal{S} := \{1.0, .1\}$. The hyper-parameter symbols for each algorithm are the same as their source publications. For example, t_1 in ZO-SLGHD and ZO-SLGHR denotes the initial scaling parameter, μ in ZO-AdaMM is the scaling parameter, and α denotes a constant learning rate. The number 784 equals the dimensional number d . It appears in the candidate set since it is taken from (Iwakiri et al., 2022), who performed similar experiments. The set of candidate values that lead to the highest fitness (averaged over the 10 image attacks) are be selected. For CMA-ES, σ_0 denotes the initial standard deviation.

	Selected Values	Candidates (μ^*)
CMA-ES	$\sigma_0 = .05$.	$\sigma_0 \in \{.05, .1, .5\}$.
EPGS	$\alpha = .1, N = .05, \sigma = .1$.	$N \in \{.02, .03, .04, .05\}, \alpha_t \in \{.1, .05\}, \sigma \in \mathcal{S}$.
STD-Homotopy	$\alpha = .5, \gamma = .5, \sigma = 1.0,$ $T_\mu = 500, \tau = 100, N_\sigma = 10$.	$\gamma \in \{.5, .8\}, \alpha \in \{.5, .1\}, \sigma \in \mathcal{S}$.
ZO-SLGHD	$\beta = 10^{-4}, \eta = .1/784, t_1 = .1, \gamma = .995$	$\beta \in \{1/784, 10^{-4}, .1\}, \eta \in \{.1/784, 10^{-3}\}, t_1 \in \mathcal{S}, \gamma \in \{.999, .995\}$.
ZO-SLGHR	$\beta = 10^{-4}, t_1 = .1, \gamma = .995$.	$\beta \in \{10^{-4}, 1/784, .1\}, \gamma \in \{.999, .995\}, t_1 \in \mathcal{S}$.
ZO-AdaMM	$\beta_1 = .9, \beta_2 = .1, \alpha = .1, \mu = .1$.	$\alpha \in \{100/784, .001, .1\}, \beta_1 \in \{.5, .9\}, \beta_2 \in \{.1, .3\}, \mu \in \mathcal{S}$
ZO-SGD	$\alpha = 10^{-4}, \mu = .1$.	$\alpha \in \{10^{-4}, .1, 1/784\}, \mu \in \mathcal{S}$.

Table 9: Hyper-parameters for CIFAR-10 Attack. The candidate set for (initial) smoothing parameters is $\mathcal{S} := \{1.0, .1\}$. The symbols are the same as those in Table 8. The set of candidate values that lead to the highest fitness (averaged over the 10 image attacks) are be selected.

	Selected Values	Candidates (μ^*)
CMA-ES	$\sigma_0 = .05$.	$\sigma_0 \in \{.05, .1, .5\}$.
EPGS	$\alpha = .1, N = .03, \sigma = .1$.	$N \in \{.02, .03, .04\}, \alpha_t \in \{.1, .05\}, \sigma \in \mathcal{S}$.
STD-Homotopy	$\alpha = .5, \gamma = .8, \sigma = .1, T_\mu = 300, \tau = 100, N_\sigma = 10$.	$\gamma \in \{.5, .8\}, \alpha \in \{.1, .5\}, \sigma \in \mathcal{S}$.
ZO-SLGHD	$\beta = .01/3072, \eta = 10^{-5}, t_1 = 1.0, \gamma = .999$	$\beta \in \{.01/3072, .1\}, \eta \in \{10^{-4}/784, 10^{-5}\}, t_1 \in \mathcal{S}, \gamma \in \{.995, .999\}$.
ZO-SLGHR	$\beta = .01/3072, t_1 = .1, \gamma = .995$.	$\beta \in \{.01/3072, .1\}, \gamma \in \{.999, .995\}, t_1 \in \mathcal{S}$.
ZO-AdaMM	$\beta_1 = .9, \beta_2 = .1, \alpha = .1, \mu = .1$.	$\alpha \in \{.5/3072, .001, .1\}, \beta_1 \in \{.5, .9\}, \beta_2 \in \{.1, .3\}, \mu \in \mathcal{S}$
ZO-SGD	$\alpha = .01/3072, \sigma = .1$.	$\alpha \in \{.01/3072, 10^{-4}, .05\}, \sigma \in \mathcal{S}$.

Table 10: Adversarial images for MNIST, produced by algorithms that apply the smoothing techniques. The adversarial target label $\mathcal{T} := \arg \min_i (\mathcal{C}(\mathbf{a}))_i$ is the label with the minimum predicted probability. All attacks are successful.

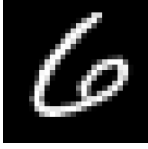


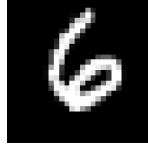






















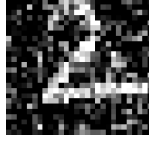

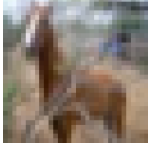


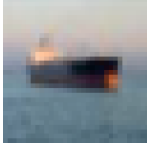
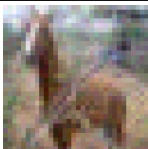

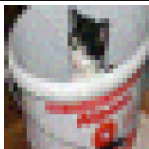

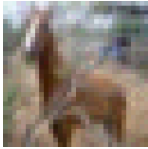



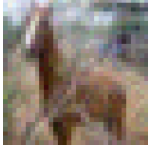



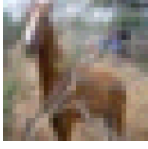


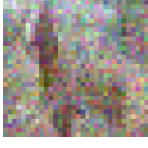


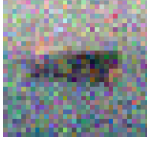


Test Image ID	9953	3850	4962	3886
Original Image Adversarial Target Label \mathcal{T}	 9	 0	 9	 3
EPGS $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 90%	 88%	 78%	 93%
ZO-SLGHD $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 70%	 61%	 41%	 76%
ZO-SLGHR $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 68%	 55%	 14%	 74%
ZO-SGD $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 85%	 76%	 61%	 87%
ZO-AdaMM $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 50%	 36%	 15%	 60%
STD-Homotopy $R^2(\mathbf{a}, \mathbf{a} + \mu^*)$	 21%	 -15%	 -19%	 61%

Table 11: Adversarial images for CIFAR-10, produced by algorithms that apply the smoothing techniques. The adversarial target label $\mathcal{T} := \arg \min_i (\mathcal{C}(a))_i$ is the label with the minimum predicted probability. Unsuccessful attacks (i.e., predicted label is different from the adversarial target) are marked with ‘Unsuccessful’.

Test Image ID	9953	3850	4962	3886
Original Image True Label \mathcal{T} Adversarial Target Label \mathcal{T}	 7 (i.e., Horse) 1(i.e., Automobile)	 9 (i.e., Truck) 4 (i.e., Deer)	 3 (i.e., Cat) 4 (i.e., Deer)	 8 (i.e., Ship) 7 (i.e., Horse)
EPGS $R^2(a, a + \mu^*)$	 96.0%	 99.3%	 99.0%	 96.1%
ZO-SLGHd $R^2(a, a + \mu^*)$	 99.2%	 99.6%	 99.6%	 98.3%
ZO-SLGHr $R^2(a, a + \mu^*)$	 97.7%	 99.5%	 98.8%	 94.3%
ZO-SGD $R^2(a, a + \mu^*)$	 99.6%	 99.8%	Unsuccessful.	 99.1%
ZO-AdaMM $R^2(a, a + \mu^*)$	 39.6%	 92.0%	 84.9%	 39.7%
STD-Homotopy $R^2(a, a + \mu^*)$	Unsuccessful.	 93.2%	 92.8%	Unsuccessful.