

# Global Optimization with a Power-Transformed Objective and Gaussian Smoothing<sup>1</sup>

Chen Xu

Department of Engineering, Shenzhen MSU-BIT University, China.



---

<sup>1</sup>Chen Xu. Global Optimization with a Power-Transformed Objective and Gaussian Smoothing. ICML, 2025 (Poster)

# Outline

1. The Problem to Solve
2. Existing Methods
3. Contribution
4. GS-PowerOpt: The Proposed Method
5. Experiments
6. Conclusion and Future Work

# The Problem to Solve

We propose a novel method, namely Gaussian Smoothing with a Power-Transformed Objective (GS-PowerOpt), that solves the following continuous optimization problems

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

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  $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x})$ .

## Gradient-Based Methods

The gradient-based methods (e.g., stochastic gradient ascent) are the most popular ones. However,

- they require the derivative of  $f$ ; and
- they are likely to be trapped in local optimum points ( $[11]$ ,  $[9]$ ,  $[4]$ ).

## Evolutionary Algorithms

- (Advantage) The evolutionary algorithms, such as symmetric annealing ([15]), particle swarm optimization ([12]), and CMA-ES ([6]), do not require the derivative of  $f$ , and may avoid local optimums.
- (Drawback) However, they suffer from the curse of dimensionality.
- (Drawback) The convergence theories are not complete for most EA algorithms.



## Standard Homotopy

The standard homotopy (e.g., [7]) creates a schedule  $\{\sigma_j\} \rightarrow 0$ , solve  $\mu_0^* := \arg \max_{\mu} \hat{f}_{\sigma_0}(\mu)$ , and performs the following double-loop mechanism:

1. Let  $\mu_j^*$  be the starting point for solving  $\max_{\mu} \hat{f}_{\sigma_{j+1}}(\mu)$  (outer loop for  $j$ ).
  2. Solve  $\mu_{j+1}^* := \arg \max_{\mu} \hat{f}_{\sigma_{j+1}}(\mu)$  (inner loop);
- $\mu_{\infty}^* = \mathbf{x}^* := \arg \max_{\mathbf{x}} f(\mathbf{x})$  (e.g., [7]).









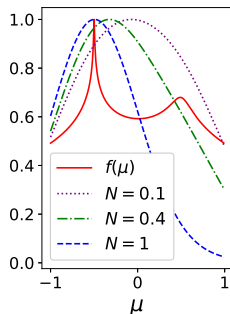
## Illustration of the Motivation

- Define  $f_N(x) := e^{Nf(x)}$ .
- The gap  $f_N(x^*) - f_N(x)$  is enlarged by increasing  $N$ .
- Gaussian smooth:  
 $F_{N,\sigma}(\mu) := \mathbb{E}_{x \sim \mathcal{N}(\mu, \sigma^2)}[f_N(x)]$ ,  
 where  $\sigma = 0.5$ .

The right figure show that

- $F_{N,\sigma}(\mu)$  has a unique maximum  $\mu^*$ .
- As we enlarge the gap  $f_N(x^*) - f_N(x)$  (i.e., increases  $N$ ),  $\mu^*$  approaches  $x^* := \arg \max_x f(x)$ .

$f(\mu) = -\log((\mu + 0.5)^2 + 10^{-5}) - \log((\mu - 0.5)^2 + 10^{-2}) + 10$  for  $|\mu| \leq 1$   
 and  $f(\mu) = 0$  for  $|\mu| > 1$ ;



**Figure:** Graphs of  $f(\mu)$  and  $F_{N,\sigma}(\mu)$ . All function graphs are scaled to have a maximum value of 1 for easier comparisons.

## The Inspired Scheme for Finding $\mathbf{x}^*$

From the previous example,  $\mu^* := \arg \max_{\mu} F_{N,\sigma}(\mu)$  approaches  $\mathbf{x}^*$  as we increase  $N$ . Therefore, the inspired method for finding  $\mathbf{x}^*$  is through finding  $\mu^*$  under a large value of  $N$ .

## Theoretical Justification of the Motivation

According to Lemma 3.4 in our paper, for any  $\sigma > 0$  and  $\delta > 0$ , as long as  $N$  is sufficiently large (depends on  $\sigma$  and  $\delta$ ), all the maximum points of  $F_{N,\sigma}(\mu)$  lie in a  $\delta$ -neighborhood of  $\mathbf{x}^*$ .

## Updating Rule of GS-PowerOpt

Since  $\mu^* := \arg \max_{\mu} F_{N,\sigma}(\mu)$  is close to  $x^*$  for sufficiently large  $N$ , GS-PowerOpt solves the surrogate objective

$$\mu^* := \arg \max_{\mu} F_{N,\sigma}(\mu),$$

where  $N$  is pre-selected (it is treated as a hyper-parameter). Stochastic gradient ascent is used to solve this objective:

$$\text{GS-PowerOpt : } \mu_{t+1} = \mu_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\mu_t), \quad (2)$$

where  $\hat{\nabla} F_{N,\sigma}(\mu_t) := \frac{1}{K} \sum_{k=1}^K (\mathbf{x}_k - \mu_t) f_N(\mathbf{x}_k)$ , and  $\{\mathbf{x}_k\}_{k=1}^K$  are independently sampled from the multivariate Gaussian distribution  $\mathcal{N}(\mu_t, \sigma^2 I_d)$ .

## A Flowchart of GS-PowerOpt



## What's New

To our knowledge, this is the first work that proposes the idea<sup>2</sup> 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.,  $\|\mathbf{x}^* - \boldsymbol{\mu}^*\|$ ).

---

<sup>2</sup>[5, 14, 3], which involve power transforms, have not mentioned this idea



## Convergence Analysis

- Let  $\{\mu_t\}$  denote the sequence produced by the GS-PowerOpt updating equation

$$\text{GS-PowerOpt : } \mu_{t+1} = \mu_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\mu_t).$$

- Let  $\nu_t$  denote the point in  $\{\mu_\tau\}_{\tau=0}^t$  that minimizes  $\mathbb{E}[\|\nabla F_{N,\sigma}(\mu_t)\|^2]$ .
- According to Corollary 3.9,  $\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla F_{N,\sigma}(\nu_t)\|^2] = 0$ , with an iteration complexity of  $O_N(d^4 \epsilon^{-2})$ . Hence,  $\nu_\infty$  is a stationary point of  $F_{N,\sigma}$ .
- From Lemma 3.4, for any  $\sigma > 0$  and arbitrarily small  $\delta > 0$ , there exists a sufficiently large  $N$  such that any stationary point of  $F_{N,\sigma}(\mu)$  lies in a  $\delta$ -neighborhood of  $\mathbf{x}^*$ .
- Hence,  $\nu_\infty$  lies in a  $\delta$ -neighborhood of  $\mathbf{x}^*$ .



## (Exponential) Power Gaussian Smoothing

- The only difference between GS-PowerOpt and the two designed algorithms, the power Gaussian smoothing (PGS) and the exponential power Gaussian smoothing (EPGS), is that PGS and EPGS update  $\mu_t$  with the normalized gradient estimate, i.e.,

$$\mu_{t+1} = \mu_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\mu_t) / \|\hat{\nabla} F_{N,\sigma}(\mu_t)\|.$$

- For PGS,  $f_N(\mathbf{x}) = f^N(\mathbf{x})$  (applied only for a non-negative objective  $f$ ).
- For EPGS,  $f_N(\mathbf{x}) = e^{Nf(\mathbf{x})}$ .

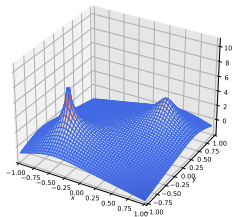
---

### Algorithm PGS/EPGS

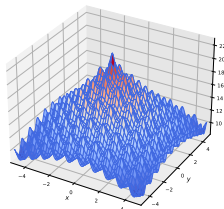
---

- 1: **Input:** The power  $N > 0$ , the scaling parameter  $\sigma > 0$ , the objective  $f$ , the initial value  $\mu_0$ , the number  $K$  of sampled points for gradient approximation, the total number  $T$  of  $\mu$ -updates, and the learning rate schedule  $\{\alpha_t\}_{t=1}^T$ .
  - 2: **for**  $t$  from 0 to  $T - 1$  **do**
  - 3:     Independently sample from  $\mathcal{N}(\mu_t, \sigma^2 I_d)$  and obtain  $\{\mathbf{x}_k\}_{k=1}^K$ .
  - 4:      $\mu_{t+1} = \mu_t + \alpha_t \hat{\nabla} F_{N,\sigma}(\mu_t) / \|\hat{\nabla} F_{N,\sigma}(\mu_t)\|$ .
  - 5: **end for**
  - 6: Return  $\{\mu_t\}_{t=1}^N$ , from which  $\mu^*$  is selected to approximate  $\mathbf{x}^*$  (e.g.,  $\mu^* := \arg \max_{t \in \{1, 2, \dots, T\}} f(\mu_t)$ ).
-

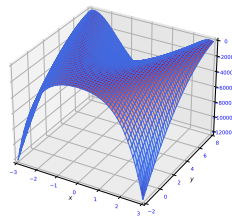
## Benchmark 2D Objectives to be Optimized



(a) Two-log Objective



(b) Ackley Func.  
(Max-Version)

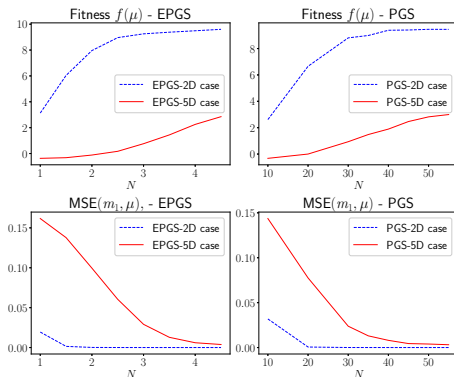


(c) Rosenbrock Func.  
(Max-Version).

Figure: Graph of the Benchmark Objective Functions.

## Experiment on Optimizing the Two-log Objective

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



**Figure:** Effects of Increasing  $N$ . We apply PGS/EPGS to solve an example problem  $\max_{\mathbf{x}} f(\mathbf{x})$ . The output from PGS/EPGS is denoted by  $\mu$ , and the global optimum of  $f$  is denoted by  $m_1$ . The graph shows that the result improves as  $N$  increases.

## Experiment on Optimizing Ackley

$$\text{Obj. } f(x, y) = 20e^{-0.2\sqrt{0.5(x^2+y^2)}} + e^{0.5(\cos(2\pi x) + \cos(2\pi y))}.$$

**Table:** Performances on Maximizing Ackley. “Iter. Taken” refers to the number of iterations taken to reach the best found solution. The true solution  $\mathbf{x}^* = (0, 0)$ .

Algorithm	Iter. Taken	Best Solution Found ( $\mu^*$ )	$f(\mu^*)$
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

## Experiment on Optimizing Rosenbrock

Obj.  $f(x, y) = -100(y - x^2)^2 - (1 - x)^2$ .

**Table:** Performances on Maximizing Rosenbrock. For PGS, the Rosenbrock is added by 20,000 to ensure the search agent only encounter positive values. The global maximum point  $\mathbf{x}^*$  of the Rosenbrock function is (1,1).

Algorithm	Iter. Taken	Best Solution Found ( $\mu^*$ )	$f(\mu^*)$
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





## Loss Function

For this task, similar to the popular loss function designed in [2], we set the loss function as

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

where

- $\mathcal{T}$  is the preselected category;
- $\mathcal{C}(\mathbf{a} + \mathbf{x})_i$  denotes the predicted probability by  $\mathcal{C}$  for  $\mathbf{a} + \mathbf{x}$  to be in category  $i$ ;
- $\kappa \geq 0$  and  $\lambda \geq 0$  are hyper-parameters.

## Explaining the Loss

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

- The smaller is  $\max_{i \neq \mathcal{T}} \mathcal{C}(\mathbf{a} + \mathbf{x})_i - \mathcal{C}(\mathbf{a} + \mathbf{x})_{\mathcal{T}}$ , the more certain for  $\mathcal{C}$  to classify  $\mathbf{a} + \mathbf{x}$  as the pre-selected category  $\mathcal{T}$ .
- When minimizing the loss, a close to zero  $\kappa$  (e.g.,  $-0.001$ ) prevents excess efforts on increasing the certainty level for  $\mathcal{C}$  to classify  $\mathbf{a} + \mathbf{x}$  as category  $\mathcal{T}$ .

## Experiment Details

- Dataset  $\mathcal{D}$ : MNIST hand-written digits or CIFAR-10 images.
- For each of the 100 randomly selected images from  $\mathcal{D}$ , we apply EPGS, and other compared algorithms, to solve

$$\max_{\mathbf{x}} \{-L(\mathbf{x})\}.$$

- If  $\mathbf{a} + \boldsymbol{\mu}^*$  is classified by  $\mathcal{C}$  as the preselected category  $\mathcal{T}$ , where  $\boldsymbol{\mu}^*$  denotes the best solution (i.e., perturbation) found by the tested algorithm, we say the attack is successful.
- Hyper-parameters are selected by trials.

## MNIST-Attack Results

**Table:** Targeted Adversarial Attack on 100 MNIST images (per-image). The success rate (SR) is the portion of successful attacks out of the 100 attacks.  $\bar{R}^2$  measures the similarity between the original image  $\mathbf{a}$  and the perturbed one  $\mathbf{a} + \mu^*$ ,  $\|\mu^*\|$  denotes the norm of the perturbation.  $\bar{T}$  denotes the number of steps taken to find the best solution.

Algorithm	SR	$\bar{R}^2$	$\ \mu^*\ $	$\bar{T}$
CMA-ES	100%	89%(4%)	2.81(0.61)	1489(12)
EPGS	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)

## Image Adversarial Attack: CIFAR-10

**Table:** Targeted Adversarial Attack on 100 CIFAR-10 images (per-image).

Algorithm	SR	$\bar{R}^2$	$\ \bar{\mu}^*\ $	$\bar{T}$
ZOSLGHd	98%	99%(1%)	1.72(0.32)	1290(411)
ZOSLGHr	98%	98%(3%)	2.66(0.68)	456(345)
EPGS	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)

## Summary on Experiment Results

Compared to algorithms using smoothing (which excludes CMA-ES):

- EPGS ranked first in the experiment of Ackley, Rosenbrock, and MNIST.
- EPGS ranked the 3rd but has a performance close to the winner in the CIFAR-10 task.

While EPGS under performs CMA-ES in the first three tasks, it beats CMA-ES in the CIFAR-10 task. Moreover, the theoretical convergence guarantee for EPGS is more developed than that of CMA-ES.

## The Case of Multiple Global Maxima

Although our convergence analysis assumes that  $f$  has a unique global maximum, this condition is not necessary for GS-PowerOpt to work, at least for the following toy example, in which we apply EPGS to solve

$$\max_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) := -\log(\|\mathbf{x} - \mathbf{m}_1\|^2 + 10^{-5}) - \log(\|\mathbf{x} - \mathbf{m}_2\|^2 + 10^{-5}),$$

where  $\mathbf{m}_1 = [-.5, -.5]$  and  $\mathbf{m}_2 = [.5, .5]$  are the two global maxima. Our experiments show that EPGS is able to locate one of the two maxima.



## Guidance on Selecting the Hyperparameters $N$ and $\sigma$

- 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 EPGS are good choices.
- Our experiments show that a  $\sigma$ -value of 10% of the search radius for  $\mathbf{x}$  is a good starting value for tuning.

## 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.
- Future works would include the convergence analysis for the maximization objective  $f$  with more than one global maxima.

- [1] A. Blake and A. Zisserman. Visual Reconstruction. MIT press, 1987.
- [2] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In 2017 IEEE Symposium on Security and Privacy (SP), pages 39–57, 2017.
- [3] J. Chen, Z. Guo, H. Li, and C. P. Chen. Regularizing scale-adaptive central moment sharpness for neural networks. IEEE Transactions on Neural Networks and Learning Systems, 35(5), 2024.
- [4] A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous, and Y. LeCun. 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, pages 192–204. PMLR, 2015.
- [5] K. Dvijotham, M. Fazel, and E. Todorov. Universal convexification via risk-aversion. In Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, 2014.
- [6] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. Evolutionary computation, 9(2):159–195, 2001.
- [7] E. Hazan, K. Y. Levy, and S. Shalev-Shwartz. On graduated optimization for stochastic non-convex problems. In Proceedings of The 33rd International Conference on Machine Learning, volume 48, pages 1833–1841, 2016.

- [8] H. Iwakiri, Y. Wang, S. Ito, and A. Takeda. Single loop gaussian homotopy method for non-convex optimization. In Advances in Neural Information Processing Systems, volume 35, pages 7065–7076, 2022.
- [9] Y. Lei, T. Hu, G. Li, and K. Tang. Stochastic gradient descent for nonconvex learning without bounded gradient assumptions. IEEE transactions on neural networks and learning systems, 31(10):4394–4400, 2019.
- [10] X. Lin, Z. Yang, X. Zhang, and Q. Zhang. Continuation path learning for homotopy optimization. In Proceedings of the 40th International Conference on Machine Learning, volume 202, pages 21288–21311. PMLR, 2023.
- [11] P. Mertikopoulos, N. Hallak, A. Kavis, and V. Cevher. On the almost sure convergence of stochastic gradient descent in non-convex problems. Advances in Neural Information Processing Systems, 33:1117–1128, 2020.
- [12] L. J. V. Miranda. PySwarms, a research-toolkit for particle swarm optimization in python. Journal of Open Source Software, 3, 2018.
- [13] H. Mobahi and J. W. Fisher. A theoretical analysis of optimization by gaussian continuation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 29, pages 1205–1211, 2015.
- [14] V. Roulet, M. Fazel, S. Srinivasa, and Z. Harchaoui. On the convergence of the iterative linear exponential quadratic gaussian algorithm to

stationary points. In 2020 American Control Conference (ACC), pages 132–137. IEEE, 2020.

- [15] P. J. Van Laarhoven, E. H. Aarts, P. J. van Laarhoven, and E. H. Aarts. Simulated annealing. Springer, 1987.

# Thank you!

- Our codes are available at  
<http://github.com/chen-research/GS-PowerTransform>.
- Welcome to visit my webpage at  
<https://orcid.org/0000-0002-7238-7254>,  
<https://www.linkedin.com/in/chen-xu-quantitative/>.
- Email: [xuchen@smbu.edu.cn](mailto:xuchen@smbu.edu.cn), [chen\\_xu\\_research@outlook.com](mailto:chen_xu_research@outlook.com)