

# Model-Based Relative Entropy Stochastic Search

Abbas Abdolmaleki  
IEETA, University of Aveiro  
DSI, University of Minho  
LIACC, University of Porto  
Portugal  
abbas.a@ua.pt

Luis Paulo Reis  
DSI, University of Minho  
LIACC, University of Porto  
Portugal  
lpreis@dsi.uminho.pt

Rudolf Lioutikov  
IAS, TU Darmstadt, Darmstadt  
Lioutikov@ias.tu-  
darmstadt.de

Jan Peters  
IAS, TU Darmstadt, Darmstadt  
MPI for Intelligent Systems  
Stuttgart, Germany  
peters@ias.tu-  
darmstadt.de

Nuno Lau  
DETI, IEETA, University of  
Aveiro, Portugal  
nunolau@ua.pt

Gerhard Neumann  
CLAS, TU Darmstadt,  
Darmstadt, Germany  
neumann@ias.tu-  
darmstadt.de

## 1. INTRODUCTION

Stochastic search algorithms [1, 2] are black box optimizers of an objective function  $R(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$ . The goal is to find one or more parameter vectors  $\theta \in \mathbb{R}^n$  which have the highest possible objective value. The only accessible information on  $R(\theta)$  are (possibly noisy) evaluations  $\{R^{[k]}\}_{k=1\dots N}$  of parameter vectors  $\{\theta^{[k]}\}_{k=1\dots N}$ , where  $k$  is the index of the sample and  $N$  is number of samples. Stochastic search algorithms typically maintain a search distribution  $\pi(\theta)$  over the parameter space  $\theta$  of the objective function  $R(\theta)$ . The search distribution  $\pi(\theta)$  is implemented as a multivariate Gaussian distribution, i.e.,  $\pi(\theta) = \mathcal{N}(\theta|\mu, \Sigma)$ . In each iteration, the search distribution  $\pi(\theta)$  is used to create samples  $\theta^{[k]}$  of the parameter vector  $\theta$ . Subsequently, the (possibly noisy) evaluation  $R^{[k]}$  of  $\theta^{[k]}$  is obtained by querying the objective function. Subsequently using the samples  $\{\theta^{[k]}, R^{[k]}\}_{k=1\dots N}$ , a new stochastic search distribution is computed. Information-theoretic search distribution updates [3] bound the Kullback Leibler divergence between two subsequent search distributions. Using a KL-bound for the update of the search distribution is a common approach in stochastic search. However, such information theoretic bounds could so far only be approximately applied either by using Taylor-expansions of the KL-divergence resulting in natural evolutionary strategies (NES) [2], or sample-based approximations, resulting in the relative entropy policy search (REPS) [3] algorithm. In this paper, we present a novel stochastic search algorithm which is called MOdel-based Relative-Entropy stochastic search (MORE). Our algorithm bounds the KL divergence of the new and old search distribution in closed form. In order to do so, we locally learn a simple, quadratic surrogate of the objective function. The quadratic surrogate allows us to compute the new search distribution analytically where the KL divergence of the new and old distribution is bounded. Therefore, we only exploit the surrogate

model locally which prevents the algorithm to be misled by inaccurate optima introduced by an inaccurate surrogate model. In addition to solving the search distribution update in closed form, we also upper-bound the entropy of the new search distribution to ensure that exploration is sustained in the search distribution throughout the learning progress, and, hence, premature convergence is avoided. We provide a comparison of stochastic search algorithms.

## 2. MODEL-BASED RELATIVE ENTROPY STOCHASTIC SEARCH

Similar as in [3], we can formulate an optimization problem to obtain a new search distribution that maximizes the expected objective value while upper-bounding the KL-divergence and lower-bounding the entropy of the distribution, i.e.,

$$\begin{aligned} \max_{\pi} \int \pi(\theta) \mathcal{R}_{\theta} d\theta & \quad (1) \\ \text{s.t. } \text{KL}(\pi(\theta)||q(\theta)) \leq \epsilon, \quad H(\pi) \geq \beta, \quad 1 = \int \pi(\theta) d\theta. \end{aligned}$$

where  $\mathcal{R}_{\theta}$  denotes the expected objective when evaluating parameter vector  $\theta$ . The term  $H(\pi) = -\int \pi(\theta) \log \pi(\theta) d\theta$  denotes the entropy of the distribution  $\pi$  and  $q(\theta)$  is the old distribution. The parameters  $\epsilon$  and  $\beta$  are user-specified parameters to control the exploration-exploitation trade-off of the algorithm. We can obtain a closed form solution for  $\pi(\theta)$  by optimizing the Lagrangian for the optimization problem given above. This solution is given as

$$\pi(\theta) \propto q(\theta)^{\eta/(\eta+\omega)} \exp\left(\frac{\mathcal{R}_{\theta}}{\eta+\omega}\right), \quad (2)$$

where  $\eta$  and  $\omega$  are the Lagrangian multipliers. The optimal value for  $\eta$  and  $\omega$  can be obtained by minimizing the convex dual function  $g(\eta, \omega)$  such that  $\eta > 0$  and  $\omega > 0$ . The dual function  $g(\eta, \omega)$  is given by

$$g(\eta, \omega) = \eta\epsilon - \omega\beta + (\eta + \omega) \log \left( \int q(\theta)^{\frac{\eta}{\eta+\omega}} \exp\left(\frac{\mathcal{R}_{\theta}}{\eta+\omega}\right) d\theta \right). \quad (3)$$

As we are dealing with continuous distributions, the entropy can also be negative. We specify  $\beta$  such that the relative difference of  $H(\pi)$  to a minimum exploration policy  $H(\pi_0)$  is decreased for a certain percentage, i.e., we change the entropy constraint to

$$H(\pi) - H(\pi_0) \geq \gamma(H(q) - H(\pi_0)) \Rightarrow \beta = \gamma(H(q) - H(\pi_0)) + H(\pi_0).$$

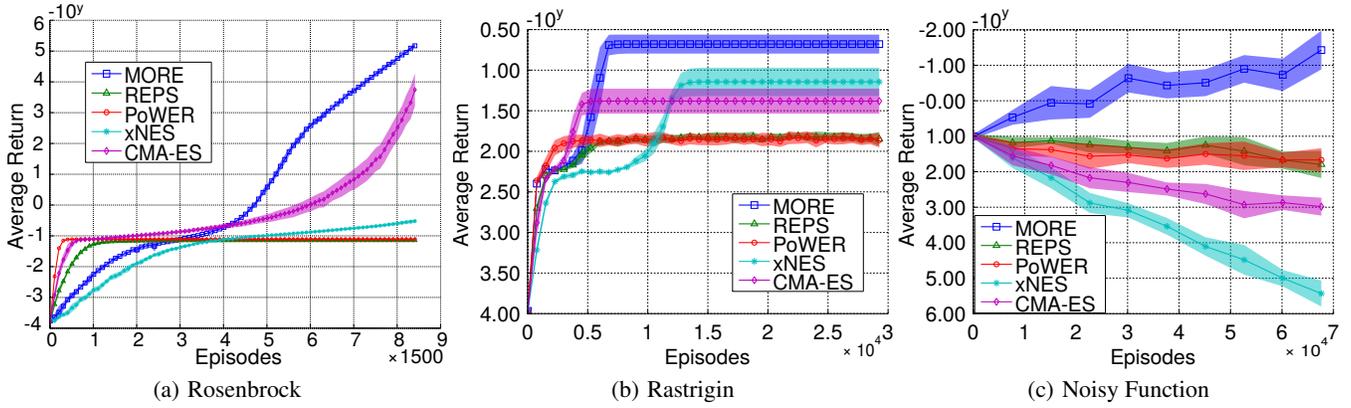
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

GECCO'16 Companion July 20-24, 2016, Denver, CO, USA

© 2016 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-4323-7/16/07...\$15.00

DOI: <http://dx.doi.org/10.1145/2908961.2930952>



**Figure 1:** Comparison of stochastic search methods for optimizing (a) Rosenbrock function and (b) Rastrigin function. (c) Comparison for a noisy objective function. All results show that MORE performs favourably.

We set minimum entropy  $H(\pi_0)$  of search distribution to a small enough value like  $-75$ .

### 2.1 Analytic Solution of the Dual-Function and the Search Distribution

Using a quadratic surrogate model of the objective function, we can compute the integrals in the dual function analytically, and, hence, we can satisfy the introduced bounds exactly in the MORE framework. At the same time, we take advantage of surrogate models such as a smoothed estimate in the case of noisy objective functions. We will assume that we are given a quadratic surrogate model  $\mathcal{R}_\theta \approx \theta^T \mathbf{R} \theta + \theta^T \mathbf{r} + r_0$  of the objective function  $\mathcal{R}_\theta$  which we will learn from data<sup>1</sup>. Moreover, the search distribution is Gaussian, i.e.,  $q(\theta) = \mathcal{N}(\theta | \mathbf{b}, \mathbf{Q})$ . In this case the integrals in the dual function given in Equation 3 can be solved in closed form. The integral inside the log-term in Equation (3) now represents an integral over an un-normalized Gaussian distribution. Hence, the integral evaluates to the inverse of the normalization factor of the corresponding Gaussian. The dual can be written as

$$g(\eta, \omega) = \eta\epsilon - \beta\omega + \frac{1}{2} \left( \mathbf{f}^T \mathbf{F} \mathbf{f} - \eta \mathbf{b}^T \mathbf{Q}^{-1} \mathbf{b} - \eta \log |2\pi \mathbf{Q}| \right) + (\eta + \omega) \log |2\pi(\eta + \omega) \mathbf{F}| \quad (4)$$

with  $\mathbf{F} = (\eta \mathbf{Q}^{-1} - 2\mathbf{R})^{-1}$  and  $\mathbf{f} = \eta \mathbf{Q}^{-1} \mathbf{b} + \mathbf{r}$ . Hence, the dual function  $g(\eta, \omega)$  can be efficiently evaluated by matrix inversions and matrix products. Note that, for a large enough value of  $\eta$ , the matrix  $\mathbf{F}$  will be positive definite and hence invertible even if  $\mathbf{R}$  is not. In our optimization, we always restrict the  $\eta$  values such that  $\mathbf{F}$  stays positive definite. Nevertheless, we could always find the  $\eta$  value with the correct KL-divergence. We can also obtain the update rule for the new policy  $\pi(\theta)$ . From Equation (2), we know that the new policy is the geometric average of the Gaussian sampling distribution  $q(\theta)$  and a squared exponential given by the exponentially transformed surrogate. After rearranging terms and completing the square, the new policy can be written as  $\pi(\theta) = \mathcal{N}(\theta | \mathbf{F} \mathbf{f}, \mathbf{F}(\eta + \omega))$ , where  $\mathbf{F}$ ,  $\mathbf{f}$  are given in the previous section. objective function.

<sup>1</sup>In order to learn the local quadratic surrogate, we can use linear regression to fit a function of the form  $f(\theta) = \phi(\theta)\beta$ , where  $\phi(\theta)$  is a feature function that returns a bias term, all linear and all quadratic terms of  $\theta$ .

## 3. EXPERIMENTS

We compare MORE with state of the art methods in stochastic search and policy search such as CMA-ES [1], NES [2], PoWER [4] and episodic REPS [3]. We use standard optimization test functions. We chose Rosenbrock function  $f(\mathbf{x}) = \sum_{i=1}^{n-1} [100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2]$ , and a multi-modal function which is known as the Rastrigin function  $f(\mathbf{x}) = 10n + \sum_{i=1}^n [x_i^2 - 10 \cos(2\pi x_i)]$ . All these functions have a global minimum equal  $f(\mathbf{x}) = 0$ . We use a 15 dimensional version of these functions. In our experiments, the mean of the initial distributions has been chosen randomly.

**Algorithmic Comparison.** We compared our algorithm against CMA-ES, NES, PoWER and REPS. In each iteration, we generated 15 new samples<sup>2</sup>. For MORE, REPS and PoWER, we always keep the last  $L = 150$  samples, while for NES and CMA-ES only the 15 current samples are kept<sup>3</sup>. As we can see in the Figure 1, MORE performs favourably in terms of both the learning speed and the final performance. However, in terms of the computation time, MORE was slower than the other algorithms. Yet, MORE was sufficiently fast as one policy update took less than 1s.

**Performance on a Noisy Function.** We also conducted an experiment on optimizing the Sphere function where we add multiplicative noise to the reward samples, i.e.,  $y = f(\mathbf{x}) + \epsilon |f(\mathbf{x})|$ , where  $\epsilon \sim \mathcal{N}(0, 1.0)$  and  $f(\mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{x}$  with a randomly chosen  $\mathbf{M}$  matrix. Figure 1(c) shows that MORE successfully smooths out the noise and converges, while other methods diverge.

## 4. REFERENCES

- [1] N. Hansen, S.D. Muller, and P. Koumoutsakos. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 2003.
- [2] Y. Sun, D. Wierstra, T. Schaul, and J. Schmidhuber. Efficient Natural Evolution Strategies. In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation (GECCO)*, 2009.
- [3] A. Kupcsik, M. P. Deisenroth, J. Peters, and G. Neumann. Data-Efficient Contextual Policy Search for Robot Movement Skills. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, 2013.
- [4] J. Kober and J. Peters. Policy Search for Motor Primitives in Robotics. *Machine Learning*, pages 1–33, 2010.

<sup>2</sup>We use the heuristics introduced in [1, 2] for CMA-ES and NES.

<sup>3</sup>NES and CMA-ES algorithms typically discard the old samples.