# Bayesian Optimization

A comparative study on acquisition functions

Raffaella D'Anna, Michele Di Sabato,
Martina Garavaglia, Anna Iob, Veronica Mazzola
Supervisor: Bruno Guindani

February 15, 2022

## OUTLINE

# A brief summary

## OUR PROBLEM

We want to solve black box optimization problems such as:

$$\text{find } x^* \text{ s.t. } f(x^*) = \max_{x \in A \subset \mathbb{R}^d} f(x)$$

Bayesian Optimization requires:

- **Gaussian process regression** as the prior distribution for the objective function $f$
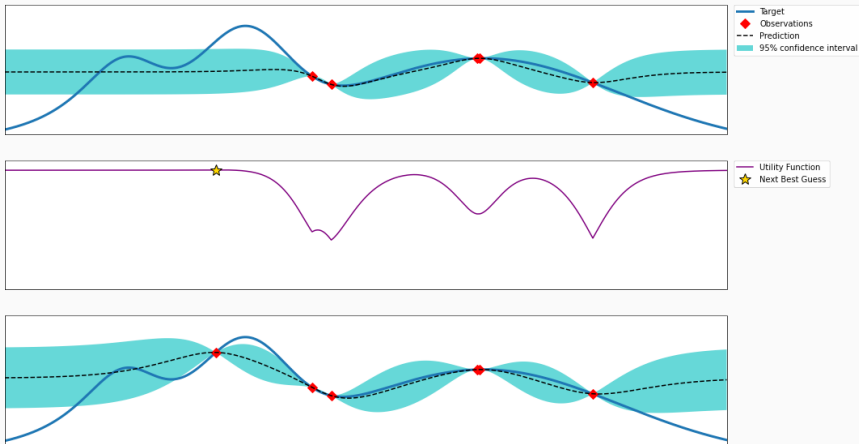- **Acquisition function** to decide where to evaluate the objective function

$$x_{n+1} = \arg\max_{x \in A} u(x|\mathcal{D}_{1:n})$$

## PSEUDO CODE

**for** $n = 1, 2, 3, ..., N$ :

1. Find the new point $x_{n+1}$ by maximizing the function $u$, called acquisition function: $x_{n+1} = \arg\max_{x \in A} u(x|\mathcal{D}_{1:n})$

2. Evaluate the objective function $f$ in the new point $x_{n+1}$ and set $f_{n+1} := f(x_{n+1})$

3. Augment the dataset: $\mathcal{D}_{1:n+1} = \{\mathcal{D}_{1:n}, (x_{n+1}, f_{n+1})\}$.

4. Compute the posterior mean $\mu(\cdot|\mathcal{D}_{1:n+1})$ and posterior variance $\sigma^2(\cdot|\mathcal{D}_{1:n+1})$

# EXAMPLE

## GOAL

The goal of our project is:

**Conducting a comparative analysis of the performance of different acquisition functions applied to a variety of objective functions**

# Library and acquisition functions

## BAYESIANOPTIMIZATION LIBRARY

- We started from the **BayesianOptimization** library by Fernando Nogueira, built in 2014

- It gives us a pure Python implementation of Bayesian global Optimization with Gaussian Processes and some known acquisition functions:
    - **PoI**: Probability of Improvement
    - **EI**: Expected Improvement
    - **UCB**: Upper Confidence Bound

## ACQUISITION FUNCTIONS

Setting $x_n^+ = \underset{i \in [1,n]}{\arg \max} f(x_i)$ :

- **PoI**:

$$PI(x) = P(f(x) \geq f_n^* + \xi) = \Phi\left(\frac{\mu(x) - f_n^* - \xi}{\sigma(x)}\right) \quad (1)$$

where $\Phi$ is the normal cumulative distribution function and $\xi$ is the trade-off parameter

- **EI**:

$$EI_n(x) := \mathbb{E}_n[\max\{0, f(x) - f(x_n^+)\}] \quad (2)$$

- **UCB**:

$$UCB(x) = \mu(x) + k\sigma(x) \quad (3)$$

where $k$ is a trade-off parameter.

## IMPROVEMENT

We have decided to improve the library adding:

- **Knowledge gradient** acquisition function
- Some **plot functions** to compare the acquisition functions and their convergence to the optimum

# Knowledge Gradient

## KNOWLEDGE GRADIENT

$$EI_n(x) := \mathbb{E}_n[\max\{0,\ f(x) - f(x_n^+)\}]$$
$$KG_n(x) := \mathbb{E}_n[\mu_{n+1}^* - \mu_n^* \mid x_{n+1} = x]$$

where

- $\mu_n^* := \max_{x'} \mu_n(x')$, with $\mu_n(x')$ being the posterior mean of the Gaussian Process given the observations up to time $n$
- $\mu_{n+1}^* := \max_{x'} \mu_{n+1}(x')$, with $\mu_{n+1}(x')$ being the posterior mean of the Gaussian Process with the new observation at time $n + 1$
- $x_n^+ = \arg\max_{x \in A} u(x|\mathcal{D}_{1:n})$

## KNOWLEDGE GRADIENT

$$EI_n(\boldsymbol{x}) = [\Delta_n(\boldsymbol{x})]^+ + \sigma_n(\boldsymbol{x})\varphi\left(\frac{\Delta_n(\boldsymbol{x})}{\sigma_n(\boldsymbol{x})}\right) + |\Delta_n(\boldsymbol{x})|\,\Phi\left(\frac{\Delta_n(\boldsymbol{x})}{\sigma_n(\boldsymbol{x})}\right)$$

$$KG_n(\boldsymbol{x}) := \mathbb{E}_n[\mu_{n+1}^* - \mu_n^* \mid \boldsymbol{x}_{n+1} = \boldsymbol{x}]$$

where $\Delta_n(\boldsymbol{x}) := \mu_n(\boldsymbol{x}) - f_n^*$ is the expected difference in quality between the proposed point $x$ and the previous best

- Initially, we tried to maximize the *KG* using a naive **grid search** approach, but it turned out to be too computationally expensive, because of the curse of dimensionality and because at each point a MC estimation of the integral needs to be performed.

- In order to speed up the iterations, we switched to the ad hoc **multi-start Stochastic Gradient Ascent** algorithm.

10

## PSEUDO CODE

### Algorithm 1: estimate $KG_n(x)$

Let $\mu_n^* = \max_{x'} \mu_n(x')$

**for** $j$ in 1 to $J$ **do**

    Generate $y_{n+1} \sim \mathcal{N}(\mu_n(x), \sigma_n(x))$

    Compute $\mu_{n+1}(x'; x, y_{n+1})$ by refitting the Gaussian Process

    $\mu_{n+1}^* = \max_{x'} \mu_{n+1}(x'; x, y_{n+1})$

    $\Delta_j = \mu_{n+1}^* - \mu_n^*$

**end for**

Estimate $KG_n(x)$ by $\frac{1}{J} \sum_{j=1}^{J} \Delta_j$

## PSEUDO CODE

### Algorithm 1: estimate $KG_n(x)$

Let $\mu_n^* = \max_{x'} \mu_n(x')$
**for** $j$ in 1 to $J$ **do**
    Generate $y_{n+1} \sim \mathcal{N}(\mu_n(x), \sigma_n(x))$
    Compute $\mu_{n+1}(x'; x, y_{n+1})$ by refitting the Gaussian Process
    $\mu_{n+1}^* = \max_{x'} \mu_{n+1}(x'; x, y_{n+1})$
    $\Delta_j = \mu_{n+1}^* - \mu_n^*$
**end for**
Estimate $KG_n(x)$ by $\frac{1}{J}\sum_{j=1}^{J} \Delta_j$

### Algorithm 2: estimate $\nabla KG_n(x)$

**for** $j$ in 1 to $J$ **do**
    Generate $y_{n+1} \sim \mathcal{N}(\mu_n(x), \sigma_n(x))$
    Compute $\mu_{n+1}(x'; x, y_{n+1})$ by refitting the Gaussian Process
    Compute $x^* := \mathrm{argmax}_{x'} \mu_{n+1}(x'; x, y_{n+1})$
    Let $G_j$ be the gradient of $\mu_{n+1}(x^*; x, y_{n+1})$ w.r.t. $x$, with $x^*$ fixed
**end for**
Estimate $\nabla KG_n(x)$ by $\frac{1}{J}\sum_{j=1}^{J} G_j$

12

## PSEUDO CODE

### Algorithm 1: estimate $KG_n(x)$

Let $\mu_n^* = \max_{x'} \mu_n(x')$
**for** $j$ in 1 to $J$ **do**
    Generate $y_{n+1} \sim \mathcal{N}(\mu_n(x), \sigma_n(x))$
    Compute $\mu_{n+1}(x'; x, y_{n+1})$ by refitting the Gaussian Process
    $\mu_{n+1}^* = \max_{x'} \mu_{n+1}(x'; x, y_{n+1})$
    $\Delta_j = \mu_{n+1}^* - \mu_n^*$
**end for**
Estimate $KG_n(x)$ by $\frac{1}{J}\sum_{j=1}^{J} \Delta_j$

### Algorithm 2: estimate $\nabla KG_n(x)$

**for** $j$ in 1 to $J$ **do**
    Generate $y_{n+1} \sim \mathcal{N}(\mu_n(x), \sigma_n(x))$
    Compute $\mu_{n+1}(x'; x, y_{n+1})$ by refitting the Gaussian Process
    Compute $x^* := \mathrm{argmax}_{x'} \mu_{n+1}(x'; x, y_{n+1})$
    Let $G_j$ be the gradient of $\mu_{n+1}(x^*; x, y_{n+1})$ w.r.t. $x$, with $x^*$ fixed
**end for**
Estimate $\nabla KG_n(x)$ by $\frac{1}{J}\sum_{j=1}^{J} G_j$

### Algorithm 3: maximize $KG_n(x)$ w.r.t $x$ (SGA)

**for** $r$ in 1 to $R$ **do**
    Choose $x_0$ at random using Latin Hypercube Design
    **for** $t$ in 1 to $T$ **do**
        Let $G$ be the estimate of $\nabla KG_n(x_{t-1}^r)$ from **Algorithm 2**
        $\alpha_t = a/(a+t)$
        $x_t^r = x_{t-1}^r + \alpha_t G$
        Check if $x_t^r$ is within bounds
    **end for**
    Estimate $KG_n(x_T^r)$ using **Algorithm 1**
**end for**
**return** $x_t^r$ with the largest estimated value of $KG_n(x_T^r)$

13

# Case studies

## A TOY EXAMPLE IN 1D - AFTER 9 ITERATIONS

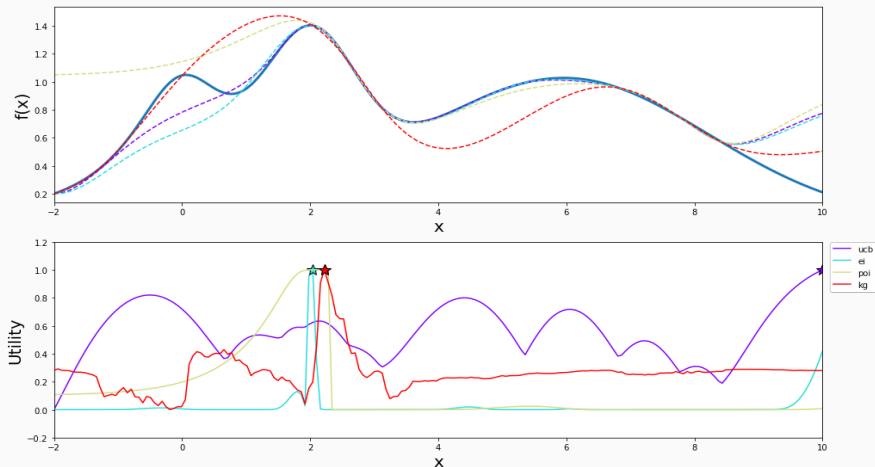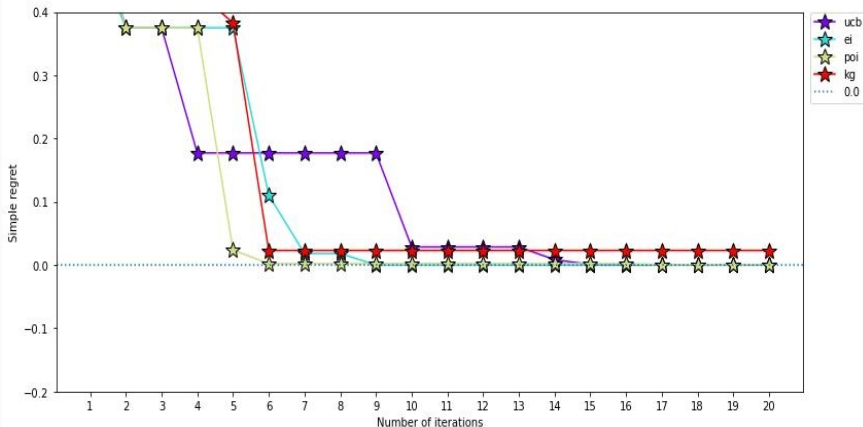$$f(x) = e^{-(x-2)^2} + e^{-\frac{(x-6)^2}{10}} + \frac{1}{x^2+1}, \quad x \in [-2, 10]$$



**Figure 1:** Argmax to be achieved: $x = 2$

14

## SIMPLE REGRET

$$r(k) = \max_{\boldsymbol{x} \in A}\{f(\boldsymbol{x})\} - \max_{t \in [1,k]}\{f(\boldsymbol{x}_t)\}$$

where $\{\boldsymbol{x}_t\}_{t=1}^{k}$ are all the points that have been suggested by the acquisition function up to iteration $k$

A brief summary
00000

Library and acquisition functions
0000

Knowledge Gradient
000000

Case studies
0000●000000000

Conclusion
000

Appendix
0000

## REGRET OF A TOY EXAMPLE IN 1D

A brief summary
○○○○○

Library and acquisition functions
○○○○

Knowledge Gradient
○○○○○○

Case studies
○○○○●○○○○○○

Conclusion
○○○

Appendix
○○○○

## 2D-ROSENBROCK FUNCTION

$$f(x,y) = 10(y - x^2)^2 + (1 - x)^2, \quad x \in [-3, 1], y \in [-2, 2]$$



**Figure 2:** The minimum is $0$ and is reached in $(0, 0)$

17

# REGRET OF ROSENBROCK FUNCTION IN 2D

## REGRET OF ROSENBROCK FUNCTION WITH NOISE IN 2D

$$f(x,y) = 10(y - x^2)^2 + (1 - x)^2 + \epsilon, \quad x \in [-3, 1], y \in [-2, 2]$$
$$\text{where } \epsilon \sim \mathcal{N}(0, 0.5^2)$$

## 2D - ACKLEY'S FUNCTION

$$f(x, y) = -20\, e^{-0.2\sqrt{\frac{x^2+y^2}{2}}} - e^{\frac{\cos(2\pi x)+\cos(2\pi y)}{2}} + 20 + e, \quad (x, y) \in [-32.7, 32.7]^2$$
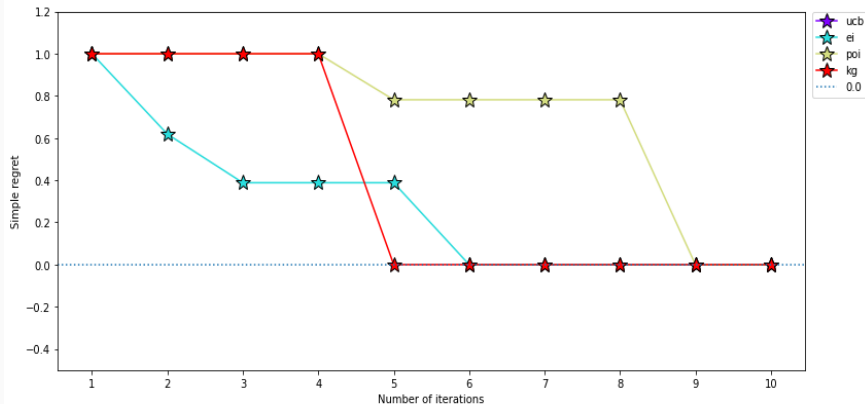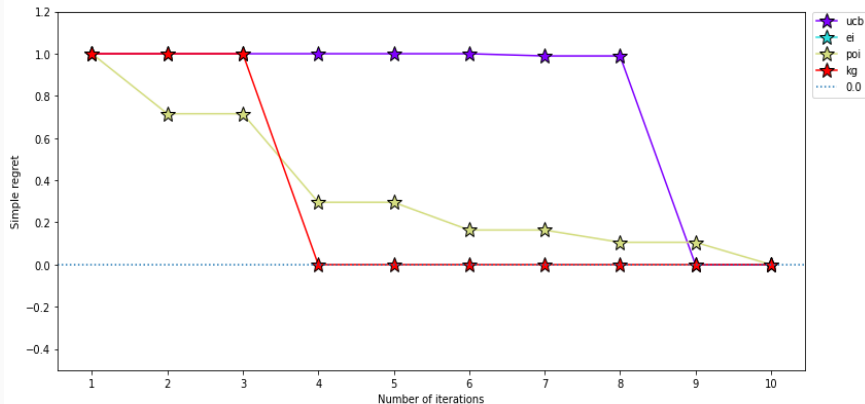


**Figure 3:** The maximum is $0$ and is reached in $(0, 0)$

20

## REGRET OF ACKLEY'S FUNCTION WITH NOISE IN 2D

$$f(x,y) = -20\,e^{-0.2\sqrt{\frac{x^2+y^2}{2}}} - e^{\frac{\cos(2\pi x)+\cos(2\pi y)}{2}} + 20 + e + \epsilon,\ (x,y) \in [-32.7, 32.7]^2$$
where $\epsilon \sim \mathcal{N}(0, 0.01^2)$



21

## REGRET OF ACKLEY'S FUNCTION WITH NOISE IN 2D

$f(x,y) = -20\, e^{-0.2\sqrt{\frac{x^2+y^2}{2}}} - e^{\frac{\cos(2\pi x)+\cos(2\pi y)}{2}} + 20 + e + \epsilon,\ (x,y) \in [-32.7, 32.7]^2$
where $\epsilon \sim \mathcal{N}(0, 0.5^2)$



22

## REGRET OF ACKLEY'S FUNCTION WITH NOISE IN 2D

$f(x,y) = -20\,e^{-0.2\sqrt{\frac{x^2+y^2}{2}}} - e^{\frac{\cos(2\pi x)+\cos(2\pi y)}{2}} + 20 + e + \epsilon,\; (x,y) \in [-32.7, 32.7]^2$
where $\epsilon \sim \mathcal{N}(0, 1.5^2)$

# Conclusion

## CONCLUSION AND FURTHER DEVELOPMENTS

- Our implementation is especially helpful when dealing with functions in higher dimensions and affected to significant noise
- However it bears computational costs that might be reduced using C++ and Parallel Programming

## REFERENCES

📄 Fernando Nogueira, *Bayesian Optimization: Open source constrained global optimization tool for Python* https://github.com/fmfn/BayesianOptimization

📄 Eric Brochu, Vlad M. Cora and Nando de Freitas, *A Tutorial on Bayesian Optimization of Expensive Cost Functions*, 2010;

📄 Peter I. Frazier, *A Tutorial on Bayesian Optimization*, 2018;

📄 Kawaguchi, Kenji, Leslie Pack Kaelbling, and Tomás Lozano-Pérez, *Bayesian Optimization with Exponential Convergence*, 2015;

📄 Candelieri Antonio, *Sequential model based optimization of partially defined functions under unknown constraints*, 2019;

📄 Marcin Molga, Czesław Smutnicki, *Test functions for optimization needs*, 2005;

# Appendix

A brief summary
00000

Library and acquisition functions
0000

Knowledge Gradient
000000

Case studies
00000000000

Conclusion
000

Appendix
0●00

## REGRET OF ROSENBROCK FUNCTION WITH NOISE IN 2D

$$f(x, y) = 10(y - x^2)^2 + (1 - x)^2 + \epsilon, \quad x \in [-3, 1], y \in [-2, 2]$$
$$\text{where } \epsilon \sim \mathcal{N}(0, 2^2)$$



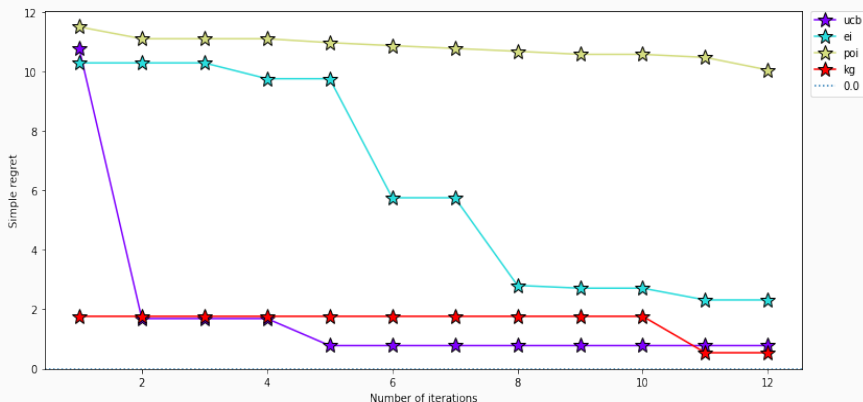Simple Regret After 12 Iterations And 5 Initial Points

**Figure 4:** standard deviation of the noise $= 2$

26

## REGRET OF ROSENBROCK FUNCTION WITH NOISE IN 2D

$f(x, y) = 10(y - x^2)^2 + (1 - x)^2 + \epsilon, \quad x \in [-3, 1], y \in [-2, 2]$
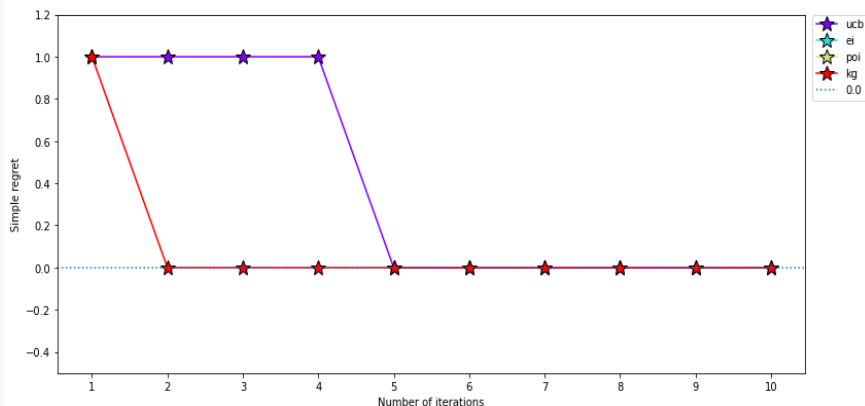where $\epsilon \sim \mathcal{N}(0, 5^2)$



**Figure 5:** standard deviation of noise $= 5$

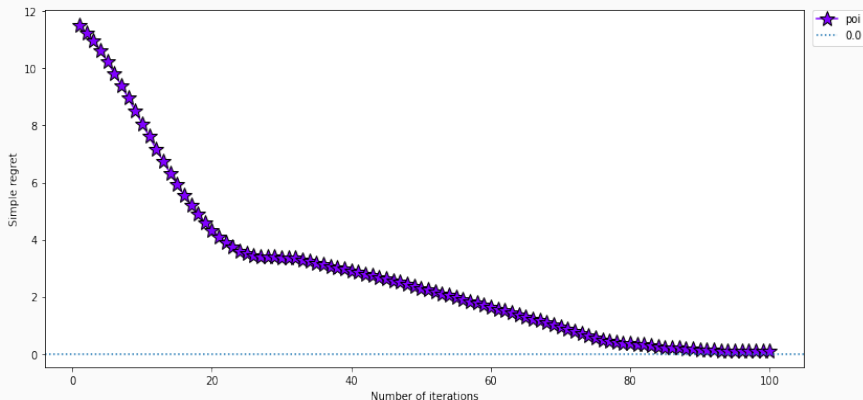27

## REGRET OF ROSENBROCK FUNCTION WITH NOISE IN 2D



**Figure 6:** simple regret of $POI$; standard deviation of noise $= 0.5$