

# Lab Session 9

## CS342 Optimization Method

Junjie Qiu

Department of Statistics and Data Science  
Southern University of Science and Technology

2024-05-09

# **Task 1: Black Box Optimization**

---

We use a simple implementation of Bayesian Optimizer to handle black box optimization problems, which makes use of a surrogate model to simulate the behavior of the black box function. Here are descriptions for each API defined in the *bo.python* file which is contained in the attached *code.zip* file.

- **EI**: Calculate the Expected Improvement (EI) acquisition function value for Bayesian optimization.
- **acq\_max**: Find the maximum of the acquisition function using random search.
- **bayesian\_optimizer**: Perform Bayesian optimization.

Our bayesian optimizer that is set with  $n$  iterations will conduct  $100 + n$  evaluations on the black box function  $f$  with the default setting in Figure 1. The  $\kappa$  parameter is used to balance between exploration and exploitation, and we will show how it works later.

```
BAYESIAN_OPTIMIZER(f, n,  $\kappa$ , EI):  
  1 initialize X_train and Y_train // default: 100  
  2 for i in range(n):  
  3   model.fit(X_train, Y_train)  
  4   X_trails = random(1000) // default: 1000  
  5    $\mu, \sigma$  = model(X_trails)  
  6   while max(EI) < 0  
  7     x_best = arg maxx EI( $\kappa$ , Y_train,  $\mu$ )  
  8     y_best = f(x_best)  
  9     add x_best and f(x_best) into train data  
 10 return min(Y_train)
```

Figure 1: Pseudocode of Bayesian Optimizer.

Figure 2 shows that with the increasement of observations, the posterior function computed by Gaussian Process will approximate the Black Box Function better with lower uncertainty.

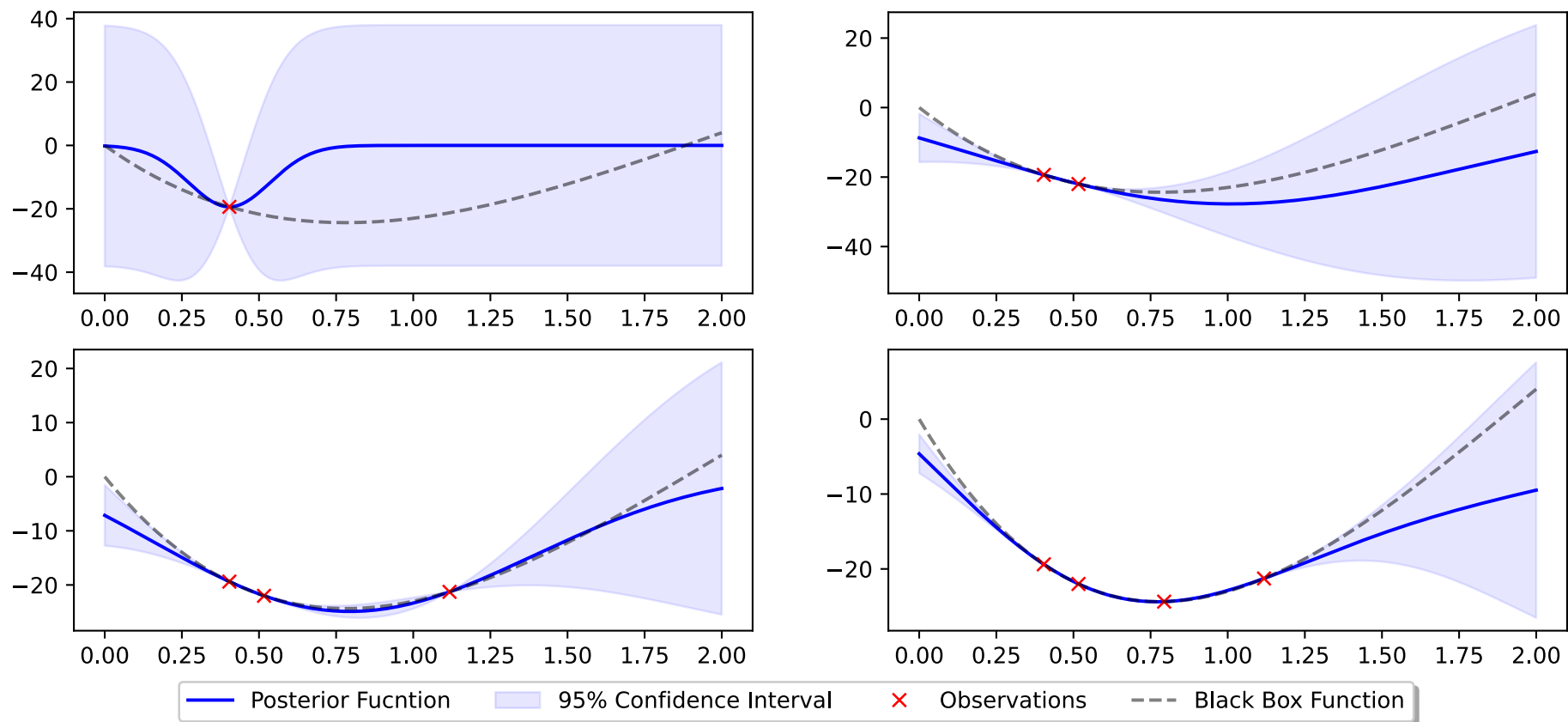


Figure 2: Simple 1D Minimization Example. Each plot shows the behavior of posterior function after adding one observation.

In addition, the next point to pick is a random but good evaluated by the Expected Increasement (EI) function. Next we will talk about the function of  $\kappa$  in EI.

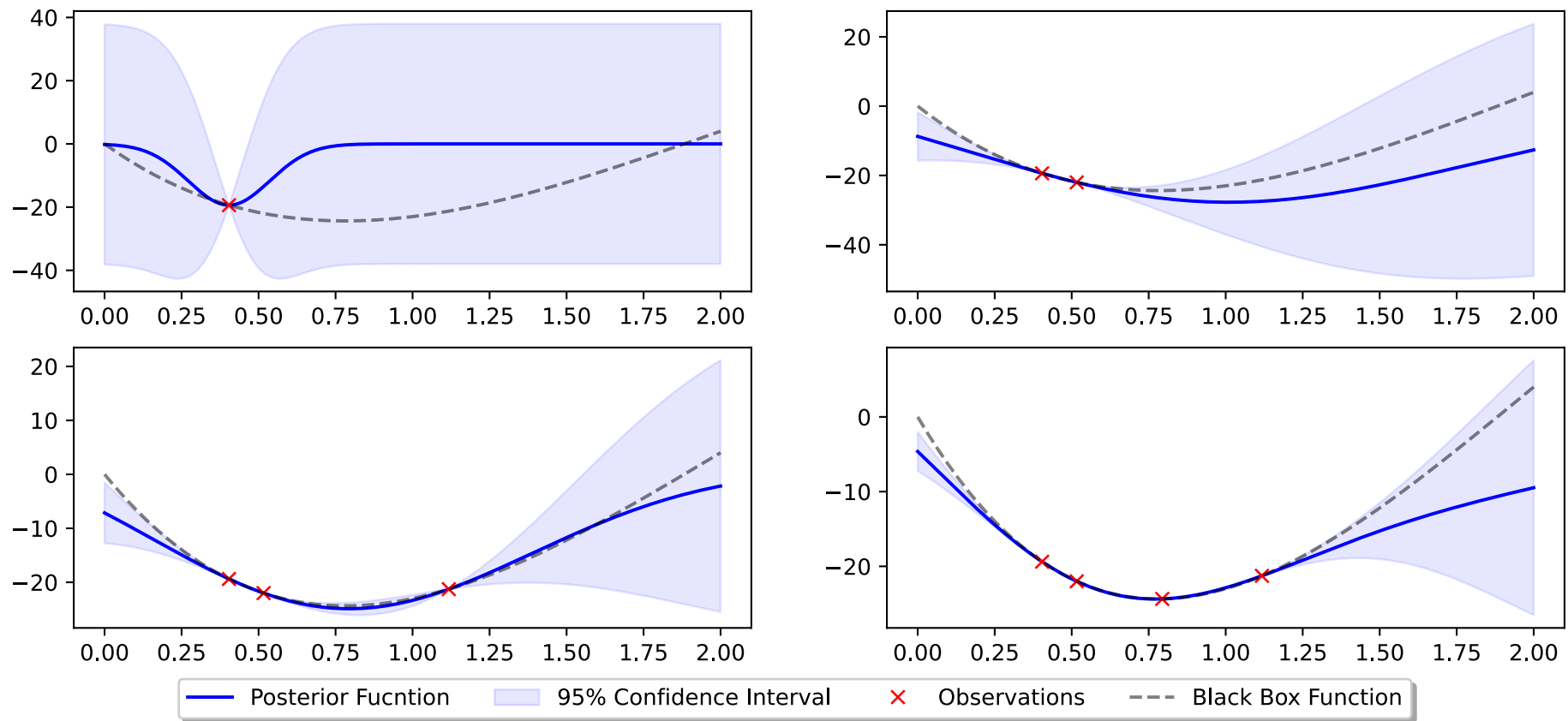


Figure 3: Simple 1D Example. Each plot shows the behavior of posterior function after adding one observation.

Set initialization size to be 5, and iteration size  $n$  to be 5. We only need to evaluate 10 points to get the near-optimal result.

Apparently, this case is easy for Gaussian Process Regression to fit, indicating that Bayesian Optimizer is suitable for few-shots optimization problem.

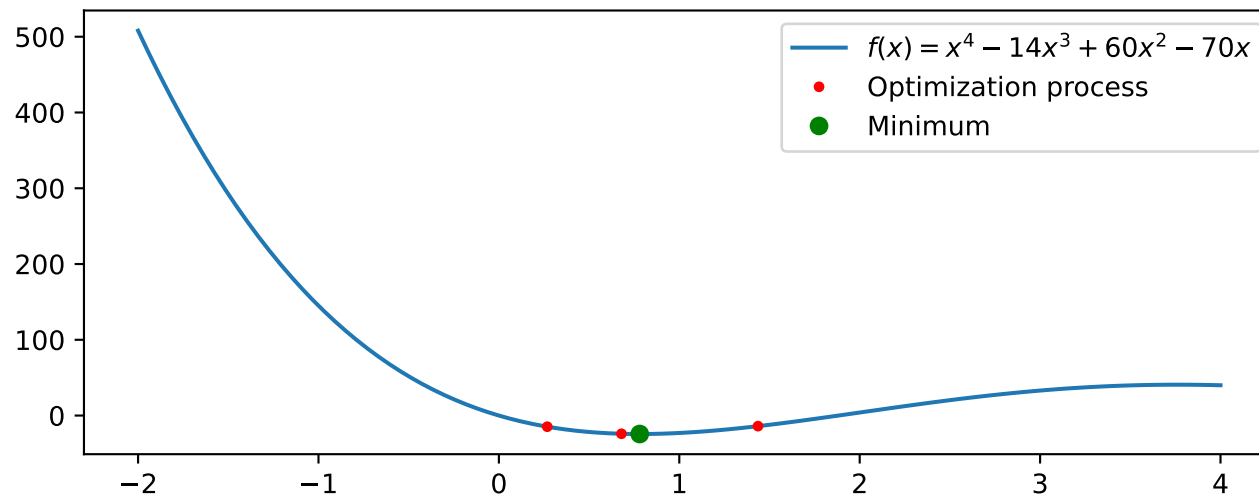


Figure 4: 1-D Case Optimized by Bayesian Optimizer.

For a **maximization problem**, the simplest EI of a newly sampled and unevaluated point  $y$  is defined as

$$\text{EI} = E[(y - M_n)\mathbb{I}(y > M_n)]$$

where  $M_{i=1}^{n_i} y_i$  is the best value observed so far.

Our implementation adds  $\kappa$  as a hyperparameter to balance between exploration and exploitation,

$$\text{EI}(\kappa) = E[(y - M_n + \kappa)\mathbb{I}(y > M_n - \kappa)]$$

When  $\kappa$  is larger, the Bayesian Optimizer tends to explore more in points with poor estimation on surrogate model; vice versa, the Bayesian Optimizer will exploit the current optimum when  $\kappa$  is small.



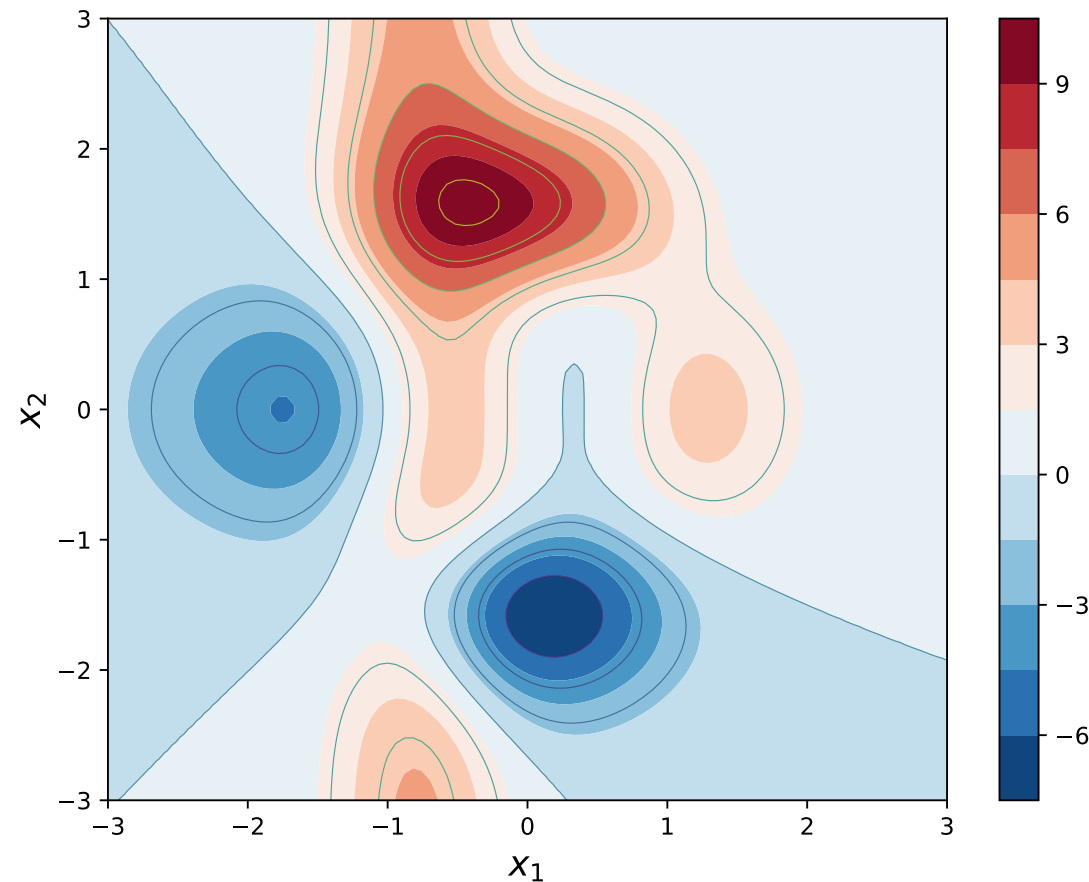


Figure 5: 2 Dimensional Example.

The plot of  $3 \cdot (1 - x)^2 \cdot e^{-x^2 - (x+1)^2} - 10 \cdot \left(\frac{x}{5} - x^3 - y^5\right) \cdot e^{-x^2 - y^2} - 3 \cdot e^{-(x+2)^2 - y^2}$  is shown in Figure 5; this will be used as the 2-D problem in following contexts.

By applying the default setting shown in Figure 1 with 10 iterations, we only budget 110 evaluations and the result is near-perfect.

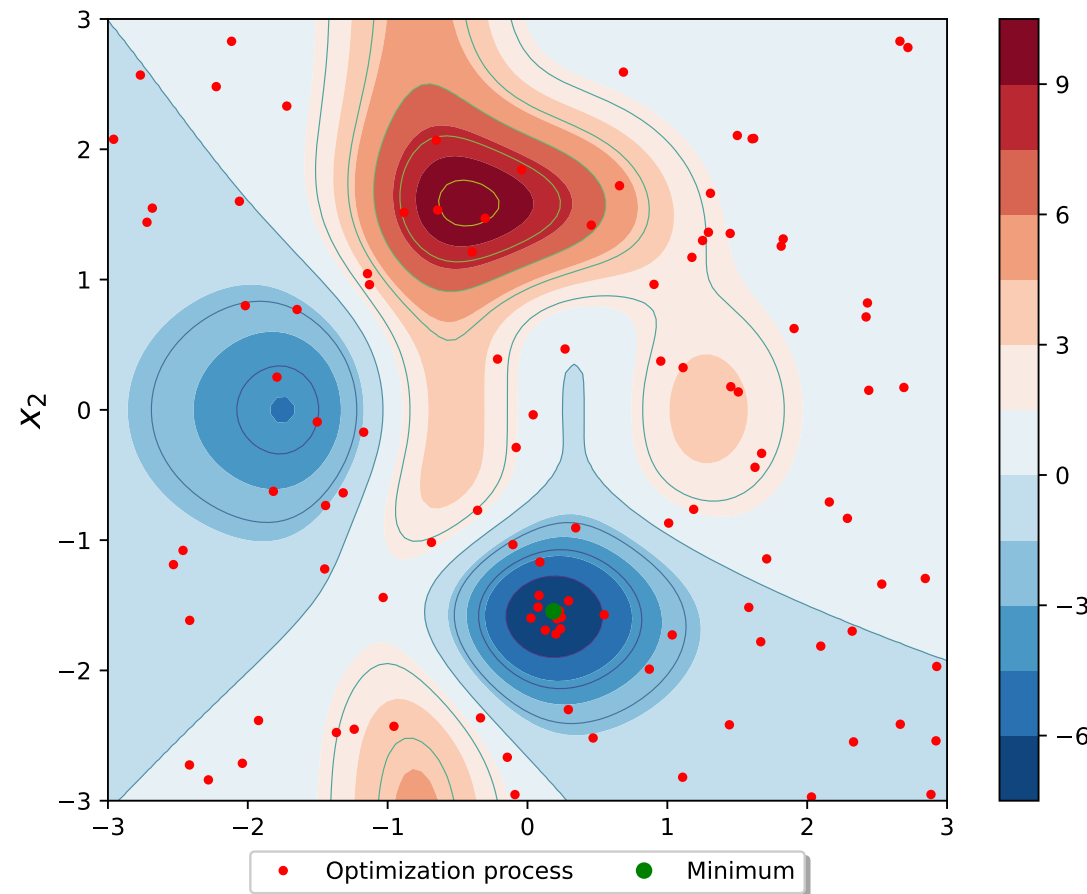


Figure 6: 2-D Case Optimized by Bayesian Optimizer.

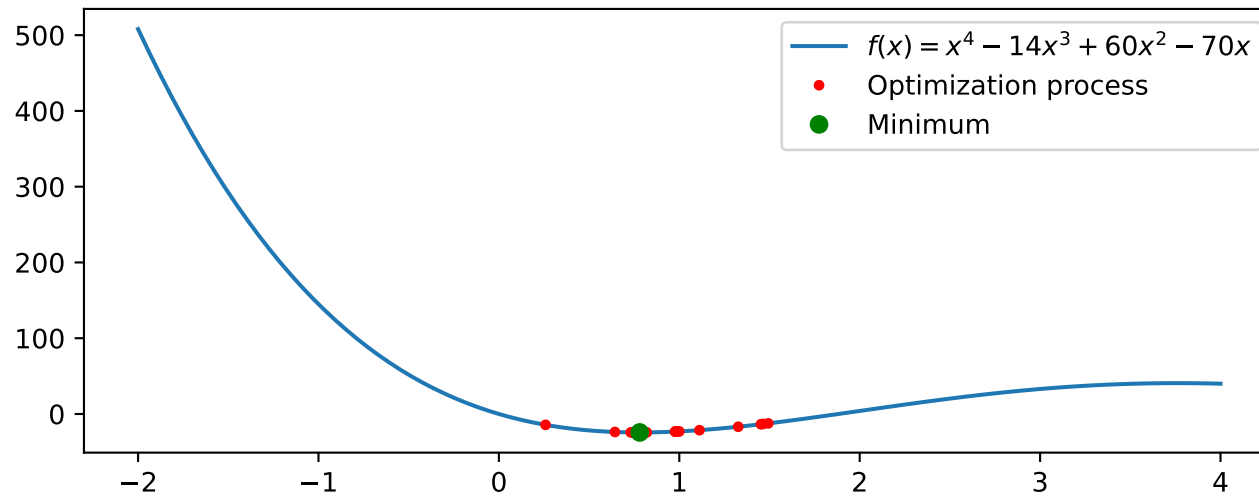


Figure 7: 1-D Case Optimized by Simulated Annealing (SA).

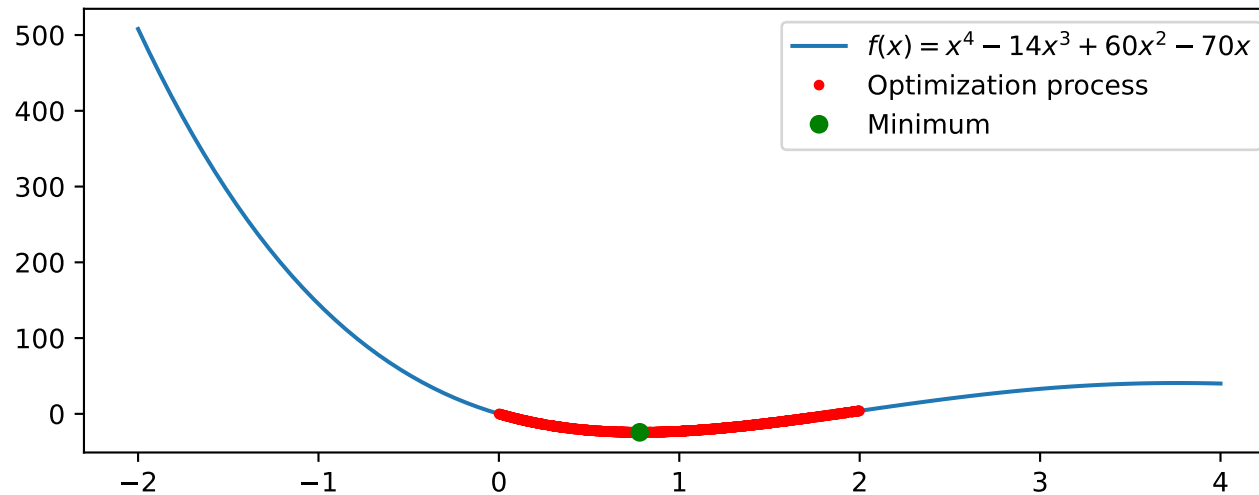


Figure 8: 1-D Case Optimized by Genetic Algorithms (GA).

**Simulated Annealing (SA):** Figure 7 shows the 1-D case optimized by simulated annealing after 1000 evaluations. The temperature function is  $100 \times 0.8^n$ , thus a lot of those points which are evaluated behind the first few ones are clustered around minimum.

**Genetic Algorithms (GA):** Figure 8 shows the 1-D case optimized by genetic algorithm after 1000 evaluations (20 points a generation with 98 gernerations). The evolution strategy includes single-point crossover and mutation by normal distribution sampling.

**Conclusion:** 1-D case is too easy for SA, GA, and Bayesian Optimizer. We need a more complex case to differentiate their performances.

With the same parameter setting to 1-D case, we leverage SA to optimize the aforementioned 2-D case.

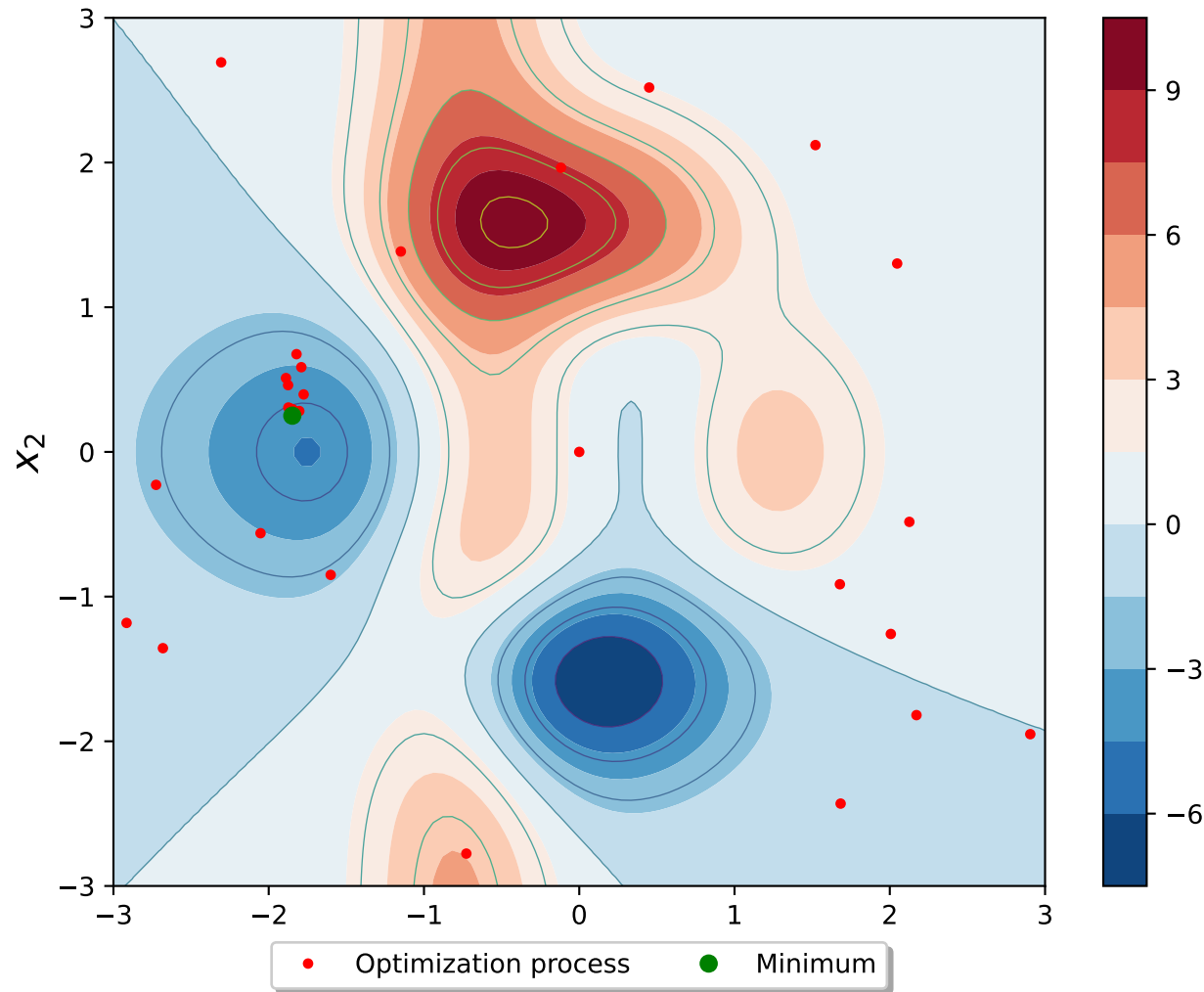


Figure 9: 1-D Case Optimized by Simulated Annealing (SA).

With the same parameter setting to 1-D case, we leverage GA to optimize the aforementioned 2-D case.

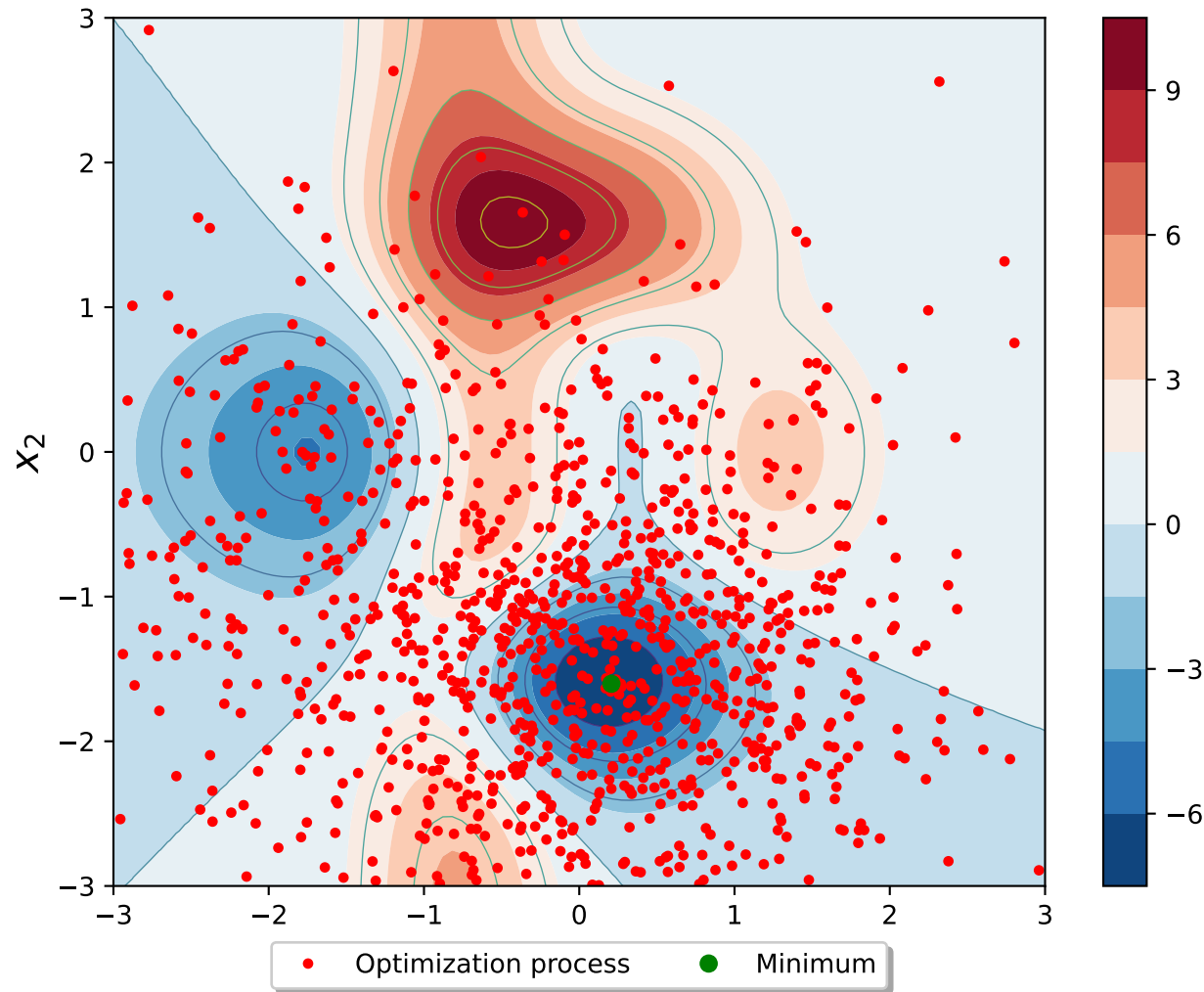


Figure 10: 1-D Case Optimized by Genetic Algorithm (GA).

**Simulated Annealing (SA):** As Figure 9 shows, the optimal point of SA falls into a local minimum and it never approaches the global optimal point. It probably has to do with the temperature and the initial value.

**Genetic Algorithms (GA):** As Figure 10 shows, the evaluated points of GA scatter largely but also have to trend to cluster around the global minimum, by which they finally reach a near perfect result.

**Bayesian Optimizer:** As Figure 6 shows, only budget 10 points, the Bayesian Optimizer can easily optimize to the optimal point.

**Conclusion:** Bayesian optimizers also show excellent optimization with a few shots on 2-dimensional problems, *however, is it really without drawbacks?*

In every iteration, the Bayesian Optimier will make use of all evaluated points  $(\mathbf{x}, \mathbf{f}(\mathbf{x}))$  to conduct inference on newly sampled points  $\mathbf{x}^*$ . We would like to estimate the  $\mathbf{y}^*$  by Gaussian Process Regression. The prior distribution of  $\mathbf{y}^*$  is set as  $N(\boldsymbol{\mu}_y, K_{yy})$ .

Assume that  $(\mathbf{f}(\mathbf{x}), \mathbf{y}^*) \sim \text{Multivariate Gaussian Distribution}$ .

$$\begin{bmatrix} \mathbf{f}(\mathbf{x}) \\ \mathbf{y}^* \end{bmatrix} \sim N \left( \begin{pmatrix} \boldsymbol{\mu}_f \\ \boldsymbol{\mu}_y \end{pmatrix}, \begin{pmatrix} K_{ff} & K_{fy} \\ K_{fy}^T & K_{yy} \end{pmatrix} \right)$$

And  $(\mathbf{y}^* | \mathbf{f}(\mathbf{x})) \sim N(K_{fy}^T K_{ff}^{-1} \mathbf{f}(\mathbf{x}), K_{yy} - K_{fy}^T K_{ff}^{-1} K_{fy})$ . Let  $n$  be the number of all evaluated points and  $m$  be the dimension of input space. We can know from above equation that the computation consumption of  $(\boldsymbol{\mu}, \boldsymbol{\sigma})$  of  $\mathbf{y}^*$  is  $O(n^3 + m^2 \times n^2)$ , which means the Bayesian Optimizer will be hard to fit when it comes to high dimensional problems or when the iteration time is too large.