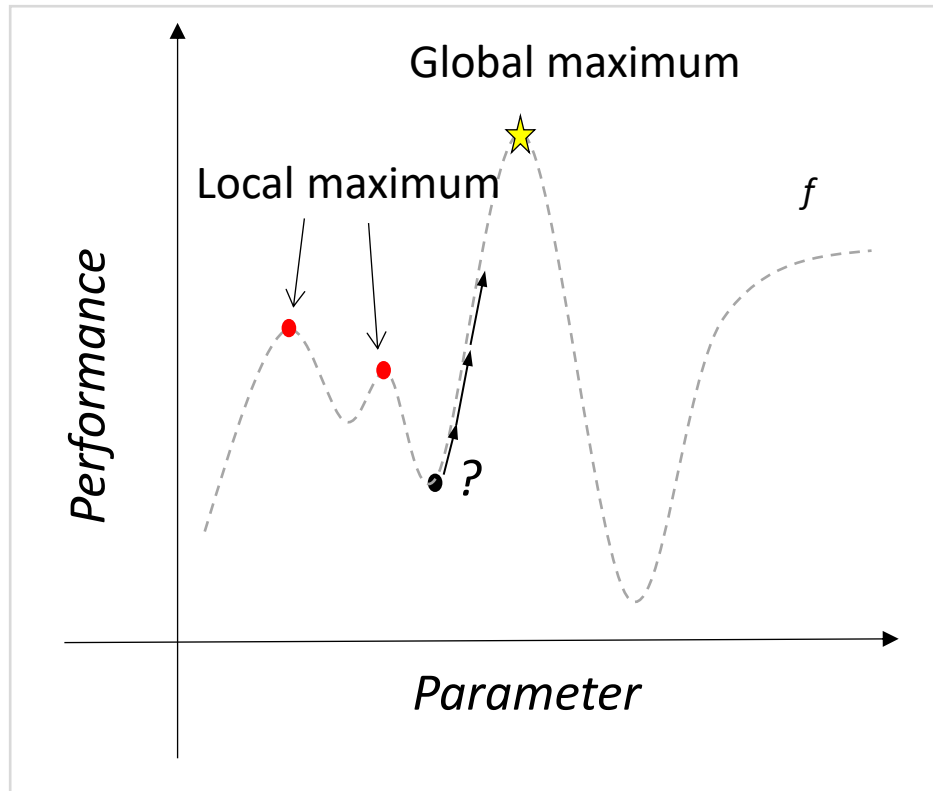


Bayesian Optimization

Felix Luong



➔ **Bayesian Optimization (BO)**

- Gradient-based methods
(Newton's method, BFGS with multi-star, etc.)

➔ **Unable to optimize without gradient** ✗

- Derivative-free methods
(DIRECT, Genetic Algorithm, Particle Swarm Optimization, etc.)

➔ **Sample inefficient due to costly evaluations** ✗

- ✓ Black-box optimization
- ✓ Sample efficiency
- ✓ Global optimization

Introduction

Automotive & Car design



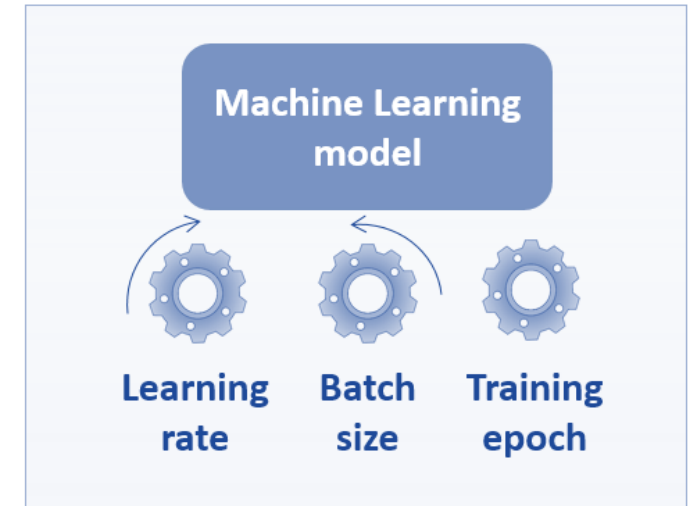
Maximize driving safety
Minimize fuel usage

Heat treatment



Maximize alloy's strength
Minimize cost

Hyper-parameter tuning



Maximize accuracy
Minimize training time

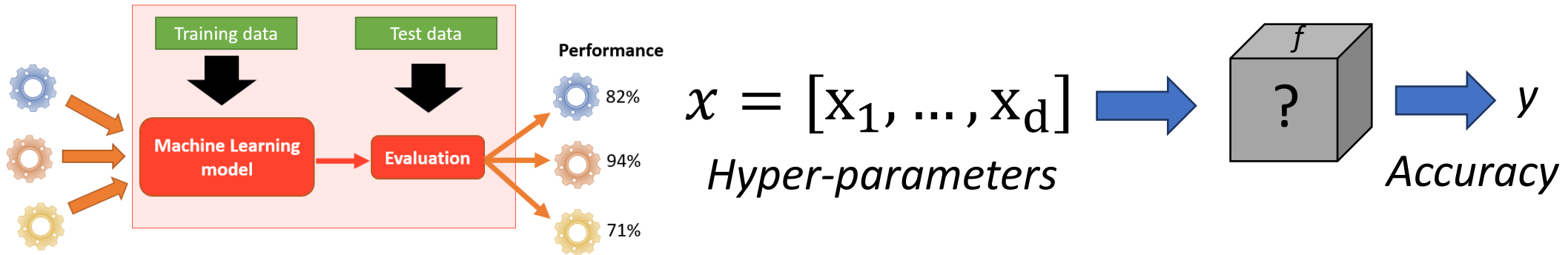
Images:

<https://www.akka-technologies.com/sectors/automotive>

<https://www.reliance-foundry.com/blog/tempering-steel#gref>

Introduction

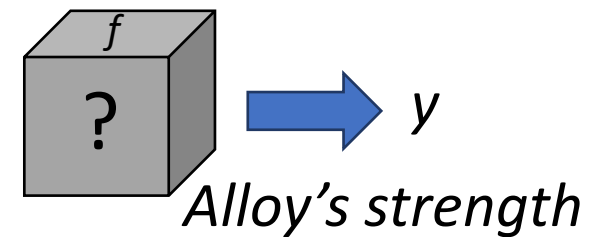
Real-world optimizations are **costly** and **black-box**.



- Require training and testing the model.



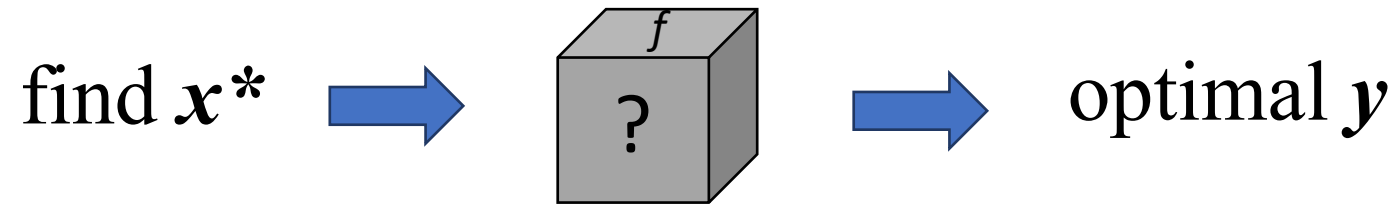
$x = [x_1, \dots, x_d]$
Temperature, time, etc.



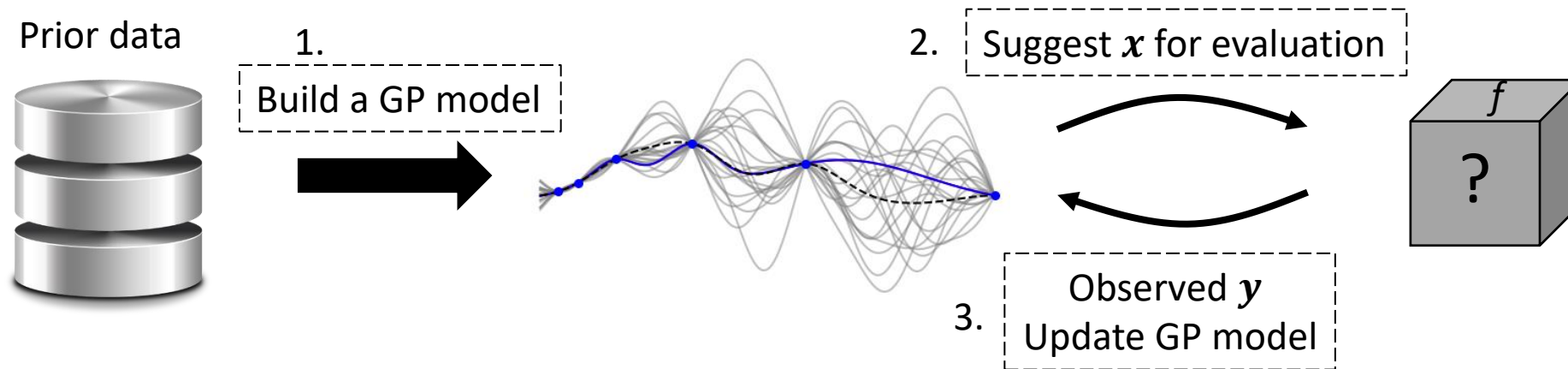
- Need real materials and time for experiments.

Introduction

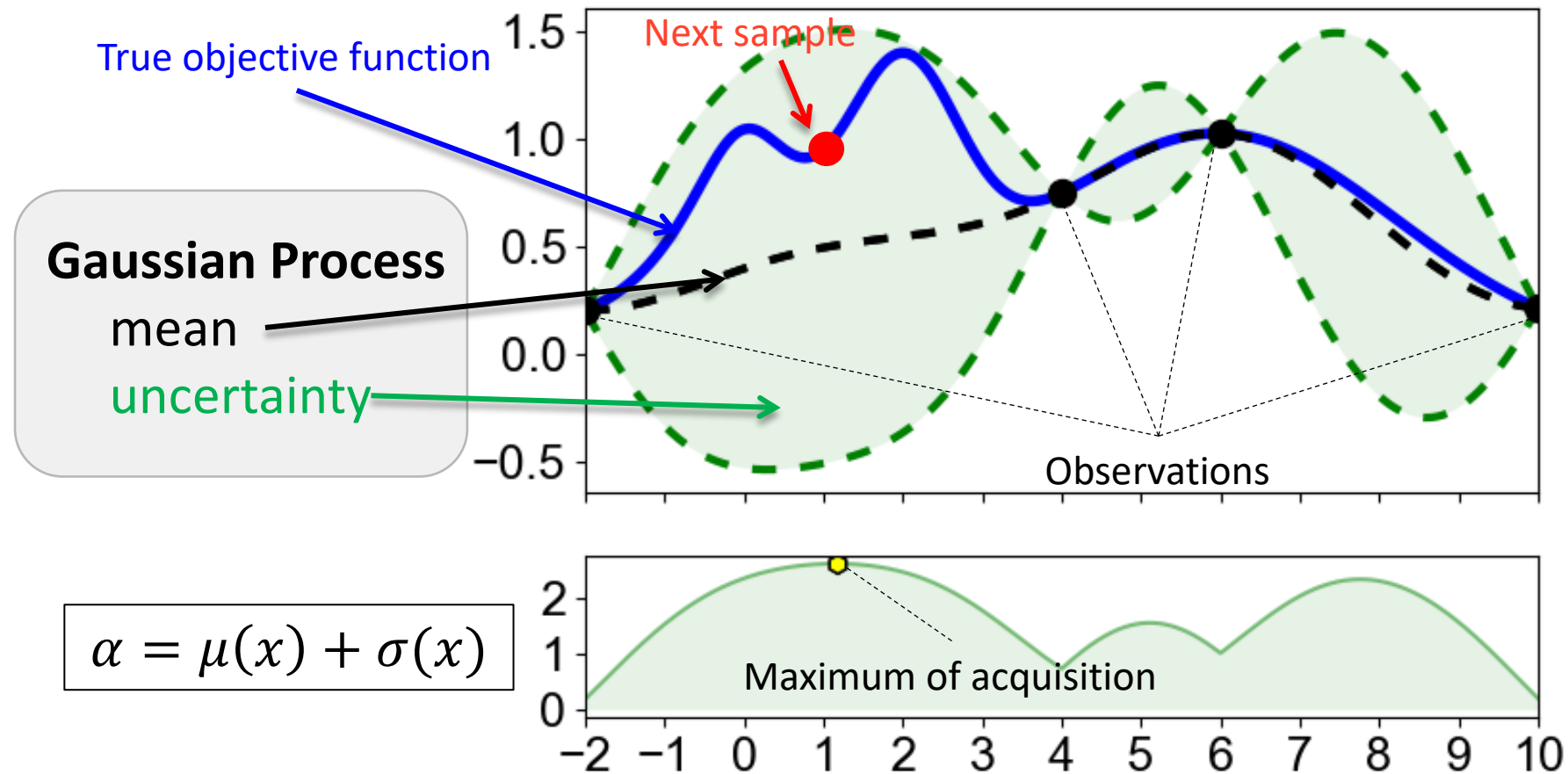
- BO is an efficient tools for automatically finding the optimal solution of expensive black-box functions.



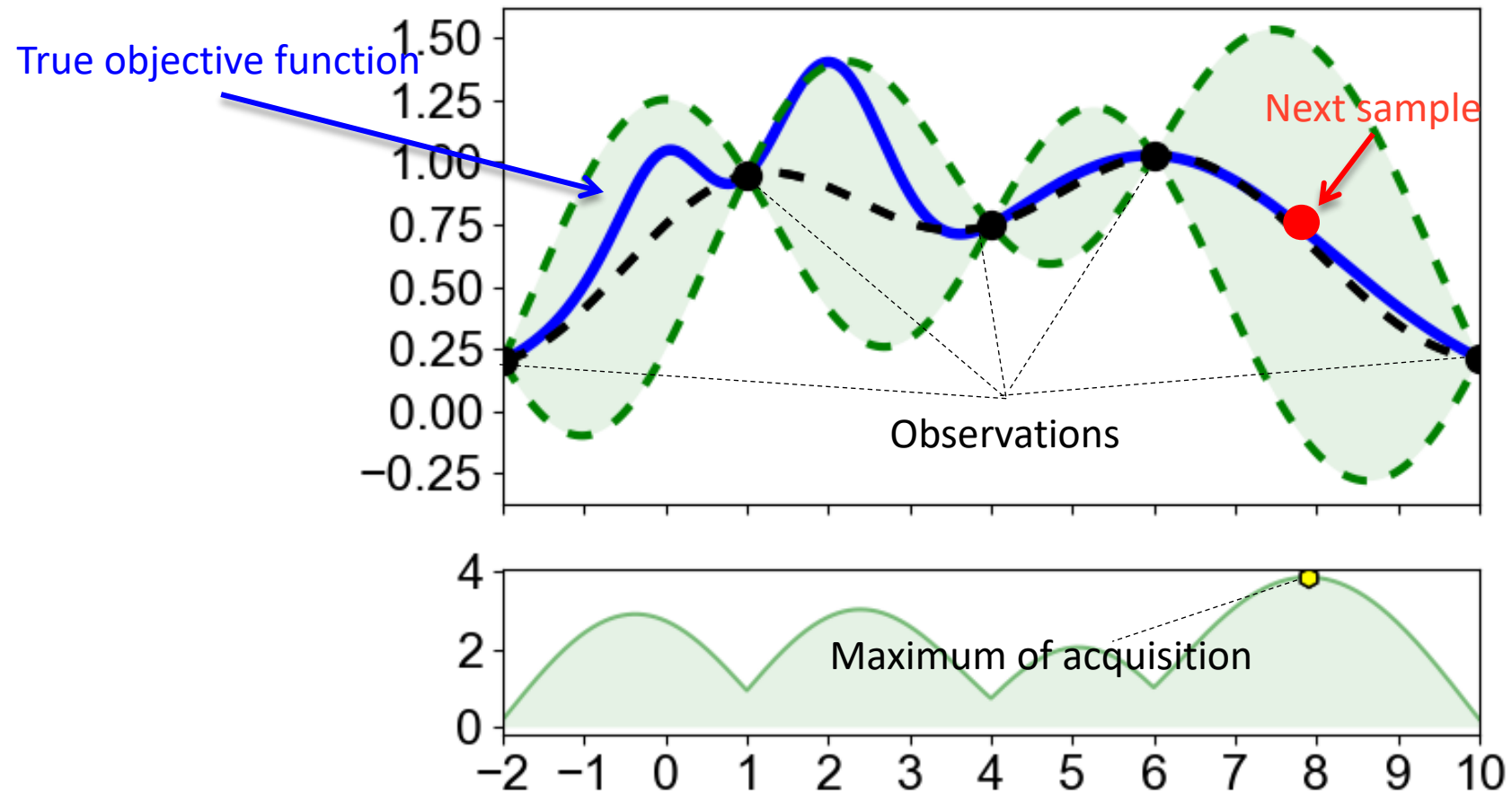
- BO uses a Gaussian process to model the black-box objective function.



Background



Background



Common acquisition functions

- Probability of Improvement (PI)

$$\alpha_{\text{PI}} = P(f(x) > f(x^+)) = \Phi\left(\frac{\mu(x) - f(x^+) - \xi}{\sigma(x)}\right)$$

- Expected improvement (EI)

$$\alpha_{\text{EI}} = \begin{cases} z \Phi\left(\frac{z}{\sigma(x)}\right) + \sigma(x) \phi\left(\frac{z}{\sigma(x)}\right) & \text{if } \sigma(x) > 0 \\ 0 & \text{if } \sigma(x) = 0 \end{cases}$$

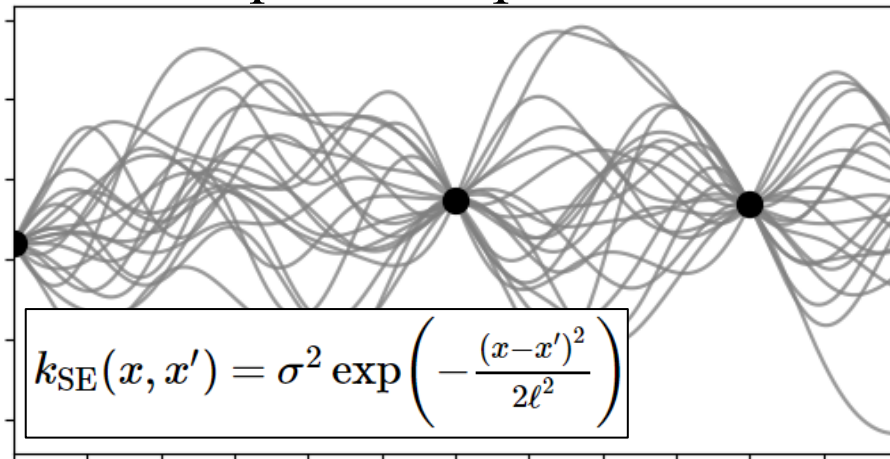
- Upper Confidence Bound (UCB)

$$\alpha_{\text{UCB}} = \mu(x) + \beta \sigma(x)$$

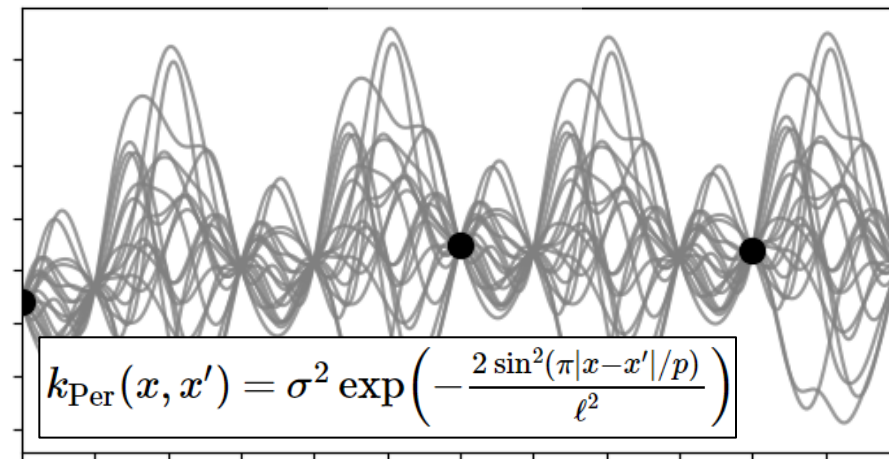
Kernel functions

- Important part of the Gaussian process used **to model the objective function**.
- Capture the **smoothness** and **correlation** between points.

Squared Exponential



Periodic



$$x, x' \in \mathbb{R}^d$$

Problem statement

- Consider tuning hyper-parameters of a machine learning model as an expensive black-box function.
- Use BO to find the best configuration of hyper-parameters.
- Some important hyper-parameters:
 - **Continuous**: Learning rate, momentum, dropout rate, ..
 - **Discrete**: Number of layer/hidden units, batch size, epochs, ..
- BO incorporates a Gaussian Process to sample a next configuration by maximizing an acquisition function in **a continuous domain**.

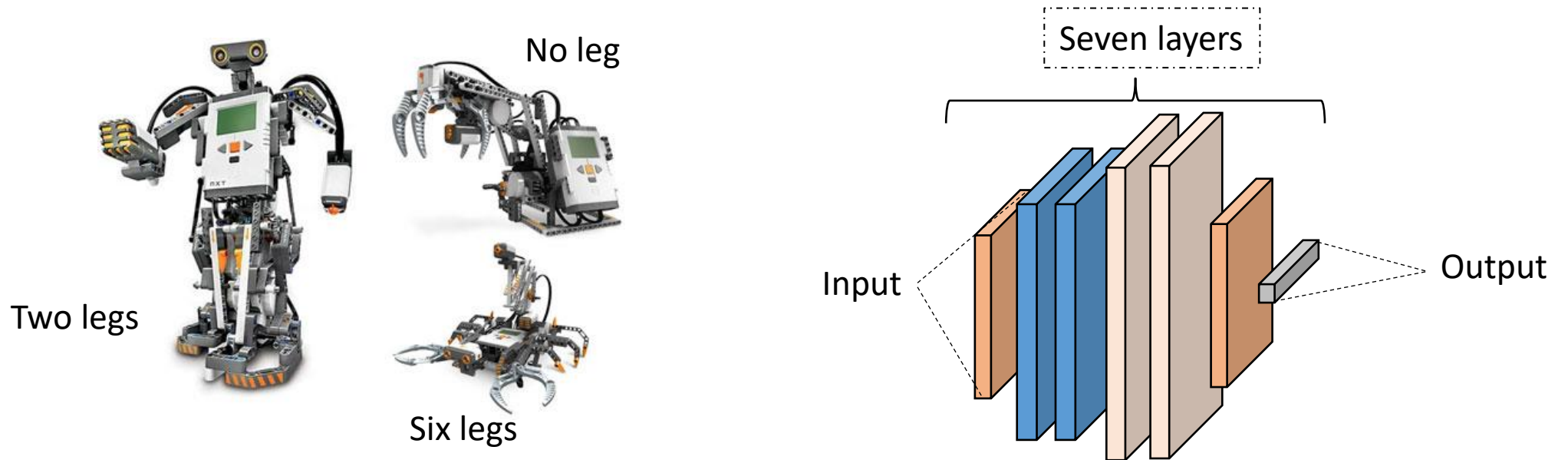


How to deal with discrete variables?

Discrete variables

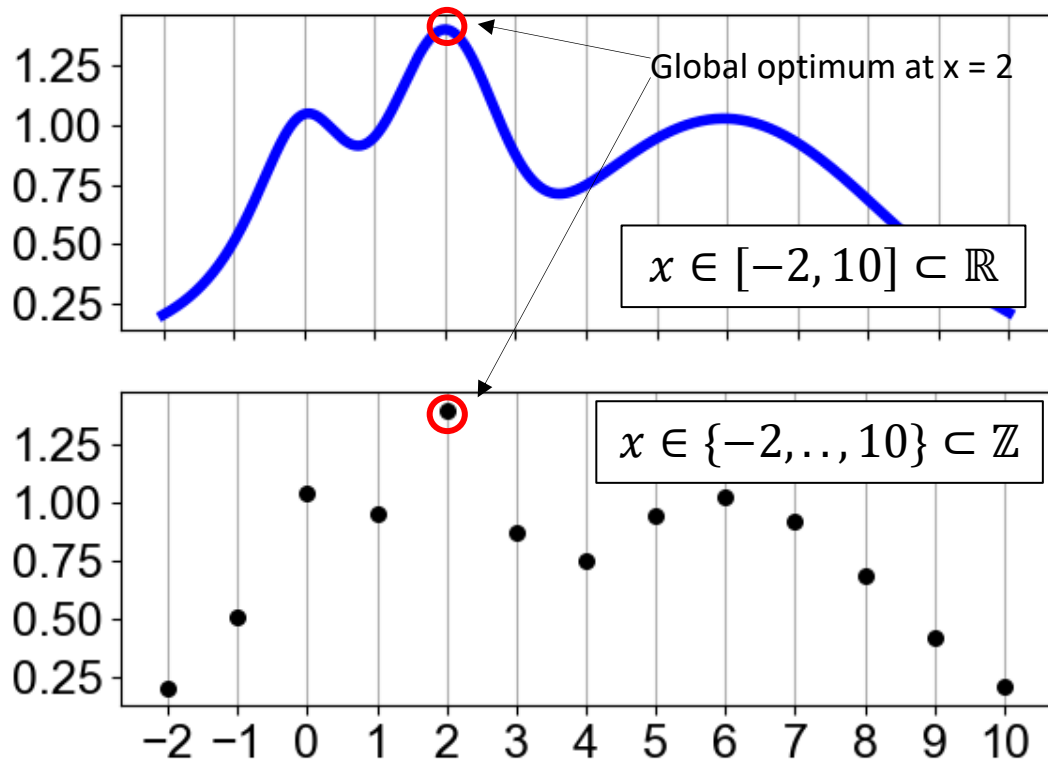
- **Discrete variable:** numeric variables. Have a countable number of values between any two values.

Example: number of wheels/legs/arms of a robot;
number of hidden units/layers of a neural network model.



Optimization with discrete variables

Assume a discrete variable x has a relation to the function evaluation y .



➔ We want to find the optimal point among **numerous** discrete points.

For example:

- Tune a neural network's hyper-parameters:
 - Hidden layers: 1 – 10
 - Hidden units: 1 – 1000

- ➔
- The number of discrete points grows **exponentially**.
 - **Expensive** to try all settings.

Example function: $y = f(x) = e^{-(x-2)^2} + e^{\left(\frac{-(x-6)^2}{10}\right)} + \frac{1}{e^{x^2} + 1}$

Optimization with discrete variables

Existing methods for dealing with discrete inputs in optimization:

- Bayesian Optimization with Naïve Rounding
- Transformation of inputs [1]
- Sequential model-based optimization (SMAC) [2]
- Tree-Parzen Estimators (TPE) [3]

[1] Garrido-Merchán, E. C., & Hernández-Lobato, D. (2018). **Dealing with categorical and integer-valued variables in Bayesian optimization with Gaussian processes.** *arXiv preprint arXiv:1805.03463*.

[2] Hutter, F., Hoos, H.H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: LION11. pp. 507{523. Springer (2011)

[3] Bergstra, J. S., Bardenet, R., Bengio, Y., & Kégl, B. (2011). **Algorithms for hyper-parameter optimization.** In *Advances in neural information processing systems* (pp. 2546-2554).


Existing method 1: BO with Naïve Rounding

- Treat **discretes variables** as **continuous**.
- Need to **round** the continuous value of a discrete variable to the nearest discrete value.

BO with **Naïve Rounding** Algorithm

1. **for** $t = 1, 2, \dots, T$ **do**
2. Build a model GP with D_t
3. Select new x_{t+1} by optimizing acquisition α
 $x_{t+1} = \operatorname{argmax}_{x \in X \subseteq \mathbb{R}^d} \alpha_{UCB}(\beta_t)$
4. **$x_{t+1} = \text{round}(x_{t+1})$**
5. Query the objective function to obtain y_{t+1}
6. Augment data $D_{t+1} = \{D_t, (x_{t+1}, y_{t+1})\}$
7. **end for**

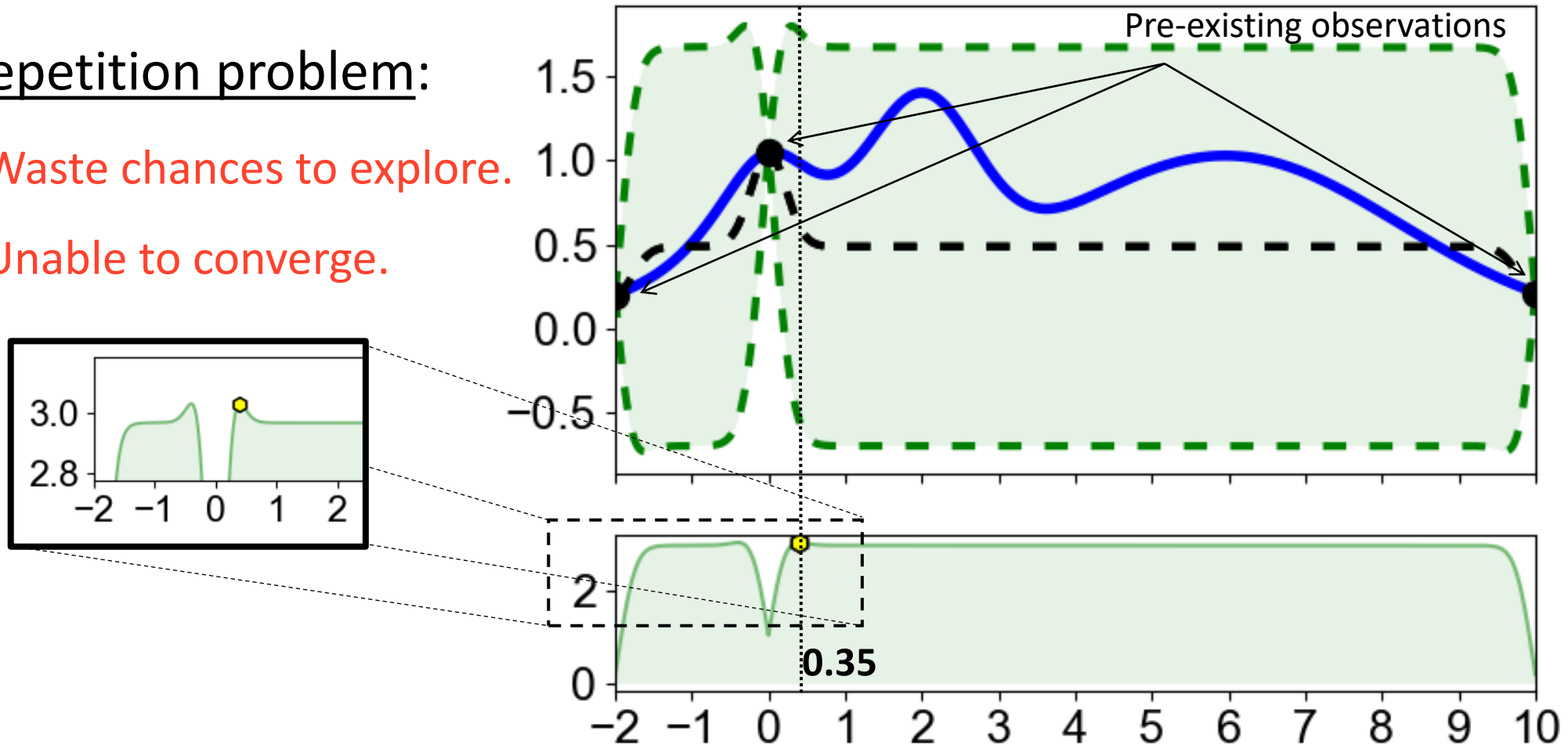
Round the continuous value
of discrete variables



Existing method 1: BO with Naïve Rounding

Repetition problem:

- Waste chances to explore.
- Unable to converge.



Existing method 2: Transformation of inputs [1]

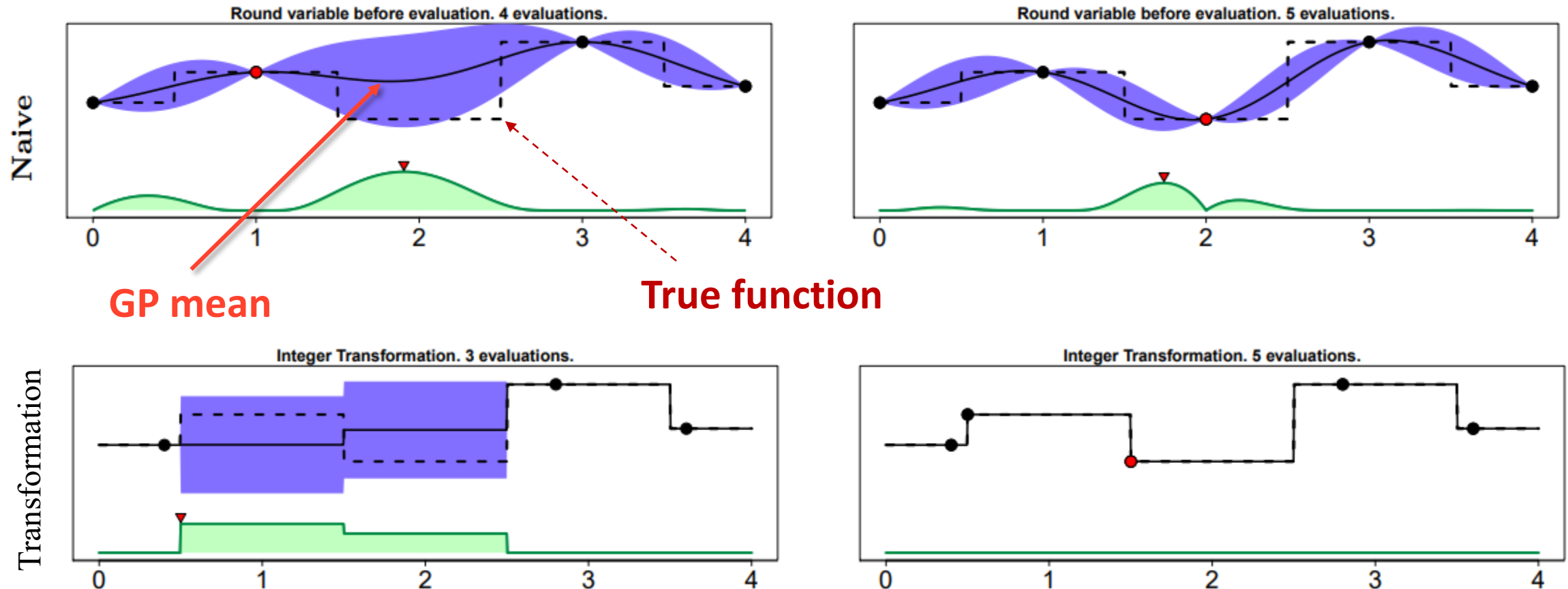
Transform the input points to \mathbf{k} to obtain an alternative covariance.

$$k'(x, x') = k(T(x), T(x')),$$

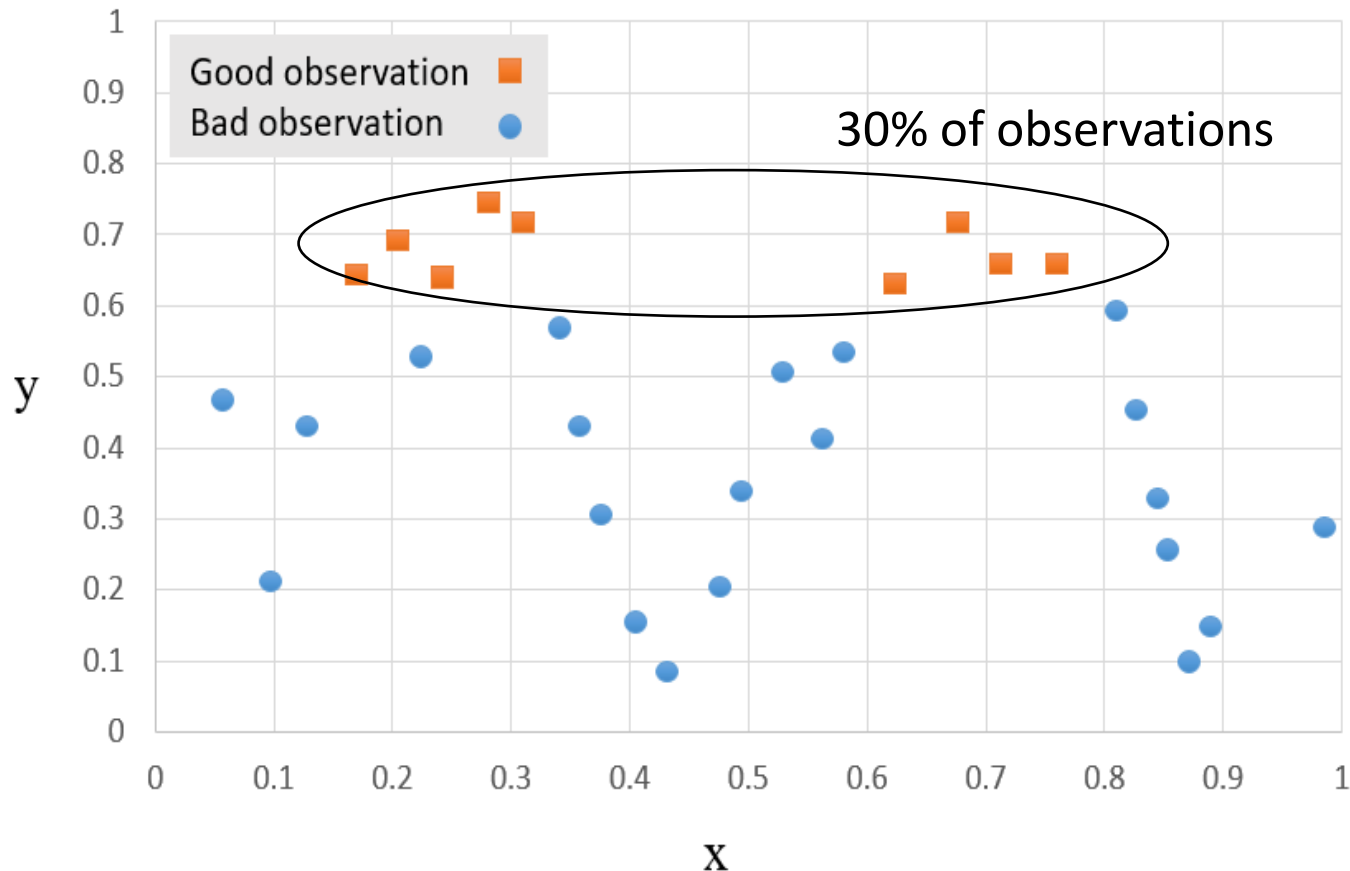
where $\mathbf{T}(\mathbf{x})$ is a transformation of all discrete inputs of $f(\cdot)$.

The values of discrete variables **are rounded to the closest discrete value**.

Comparison between Naïve and Transformation methods

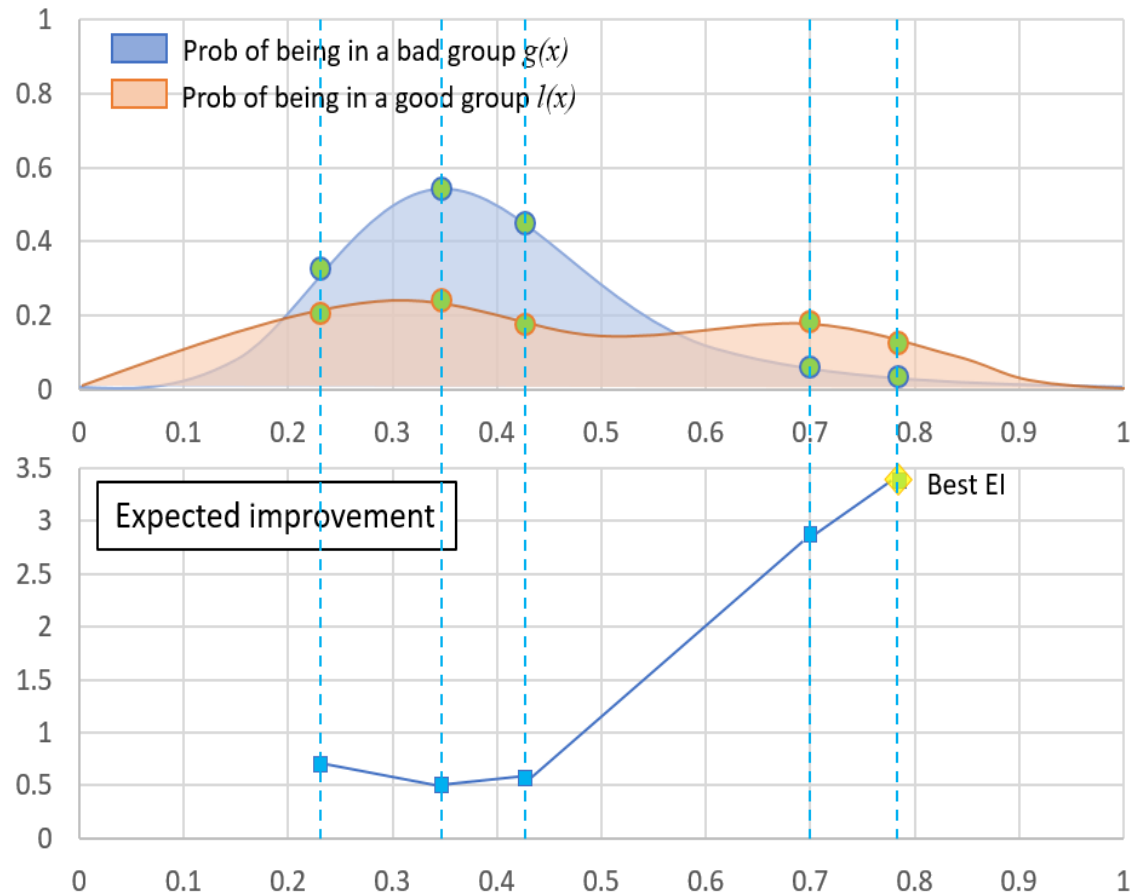


Existing method 3: Tree-Parzen Estimators [2]



1. Sort observations by evaluation.
2. Take a percentage of **best observations** for **good** group.
3. Remaining observations belong to **bad** group.

Existing method 3: Tree-Parzen Estimators [2]



$$EI_{y^*} = \left(\gamma + \frac{g(x)}{l(x)} (1 - \gamma) \right)^{-1}$$

Advantages:

- Sample from the discrete distribution.

Disadvantages:

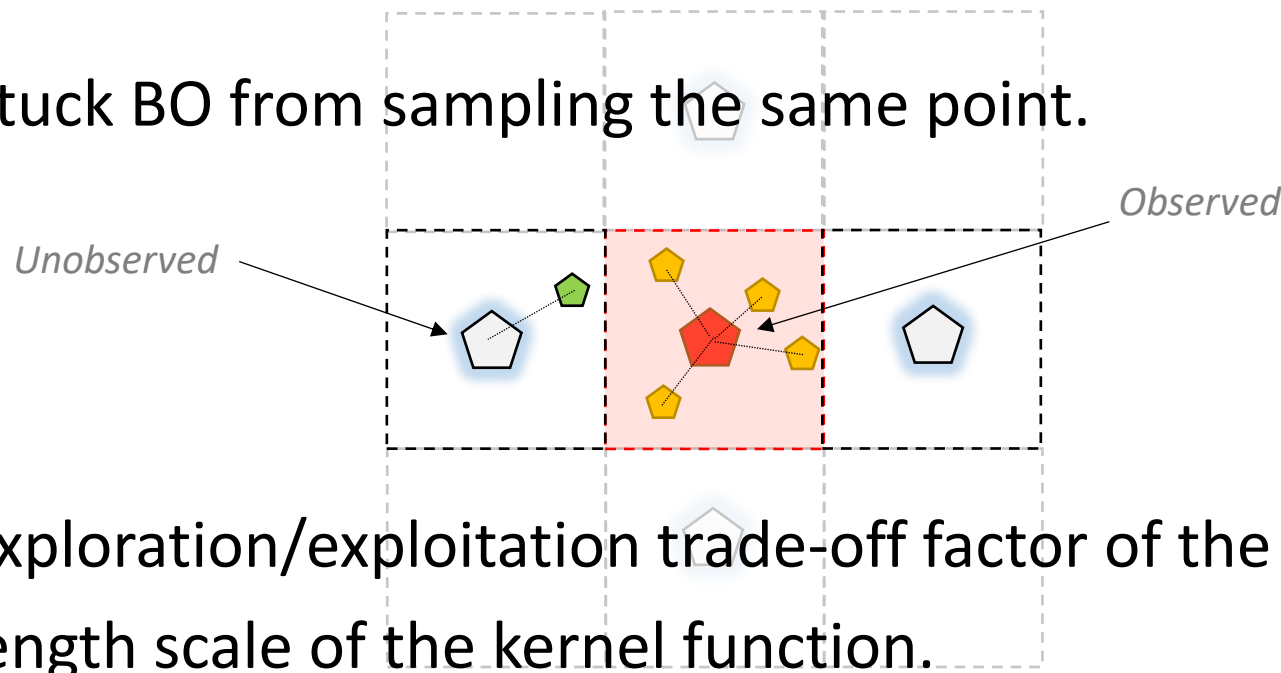
- Need sufficient amount of observations.
- Randomly draw from distributions.

Proposed method

❑ Question:

*“How to **avoid repetition** in optimizing functions with **discrete inputs**?”*

❑ Idea: Unstuck BO from sampling the same point.



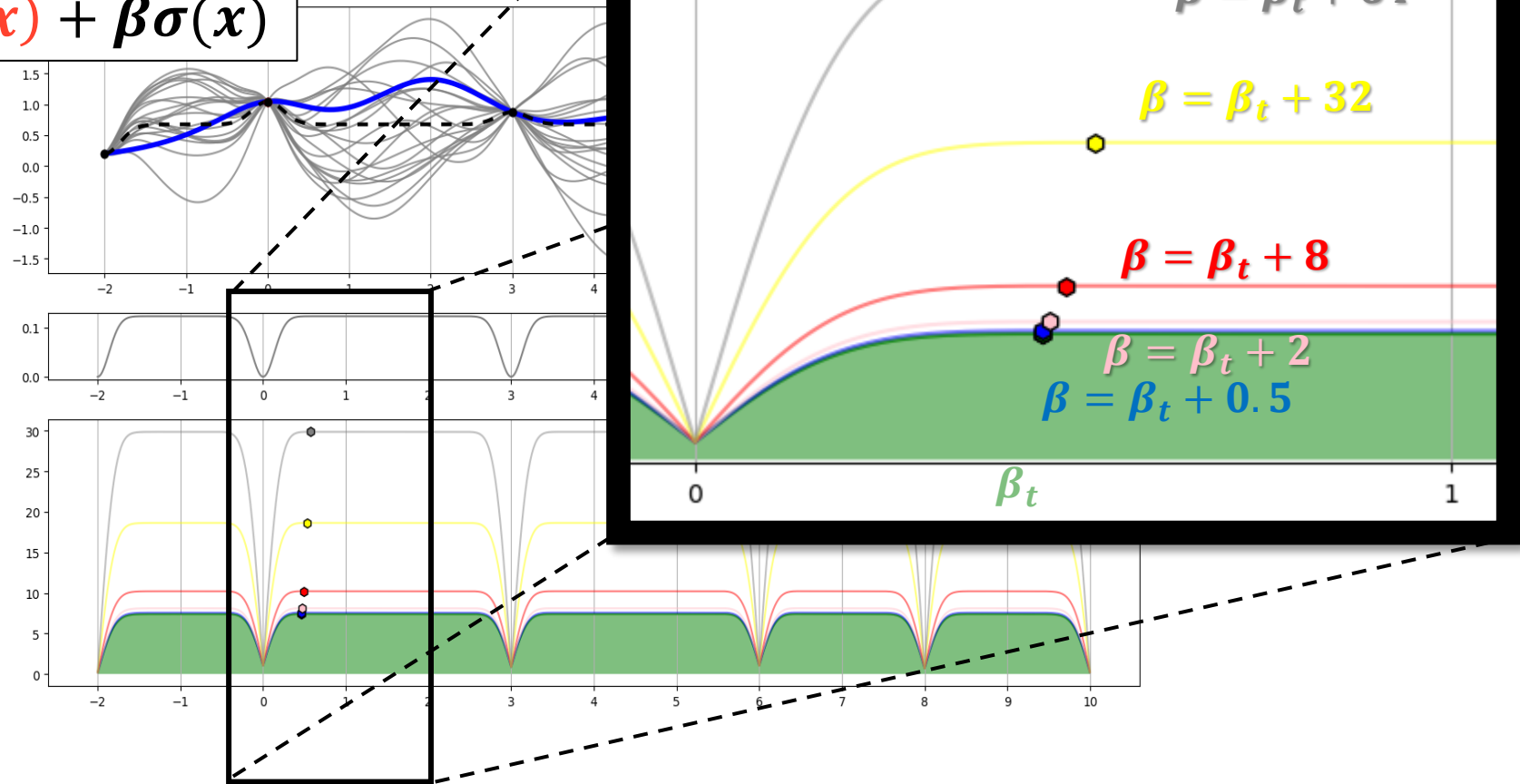
❑ Solution:

- Adjust exploration/exploitation trade-off factor of the acquisition function.
- Adjust length scale of the kernel function.

Proposed method - Adjust β

Dominated by mean:

$$\alpha_{UCB} = \mu(x) + \beta\sigma(x)$$

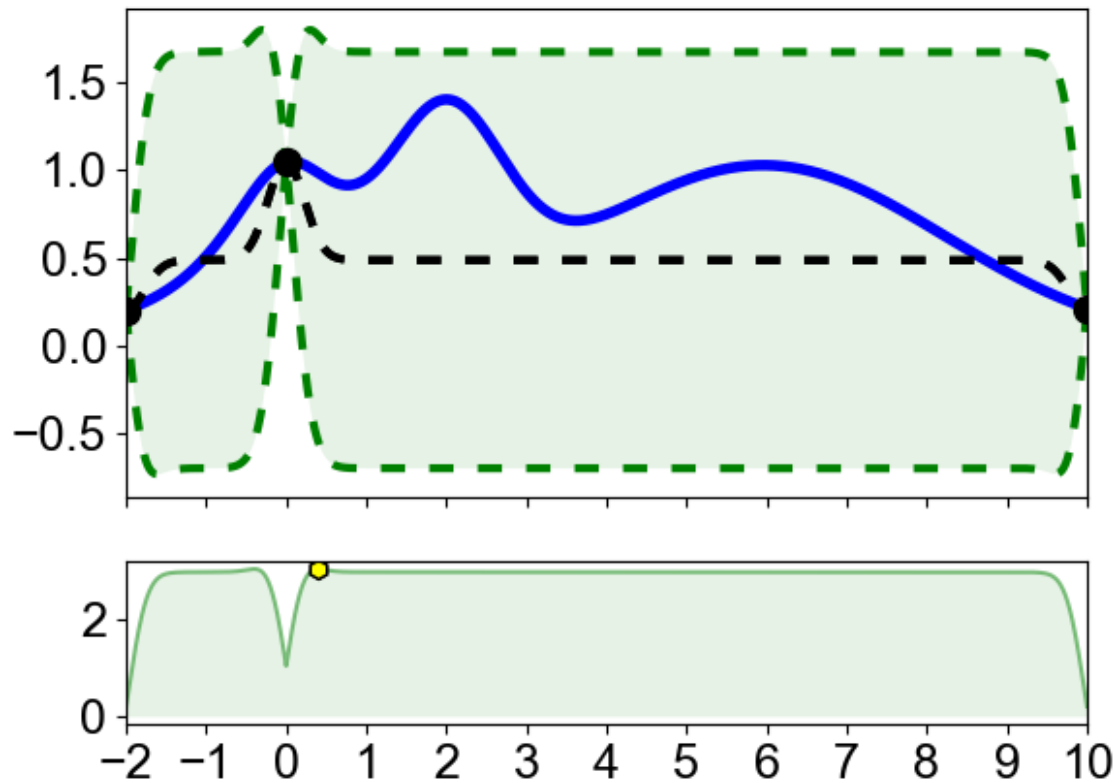


Proposed method - Adjust l

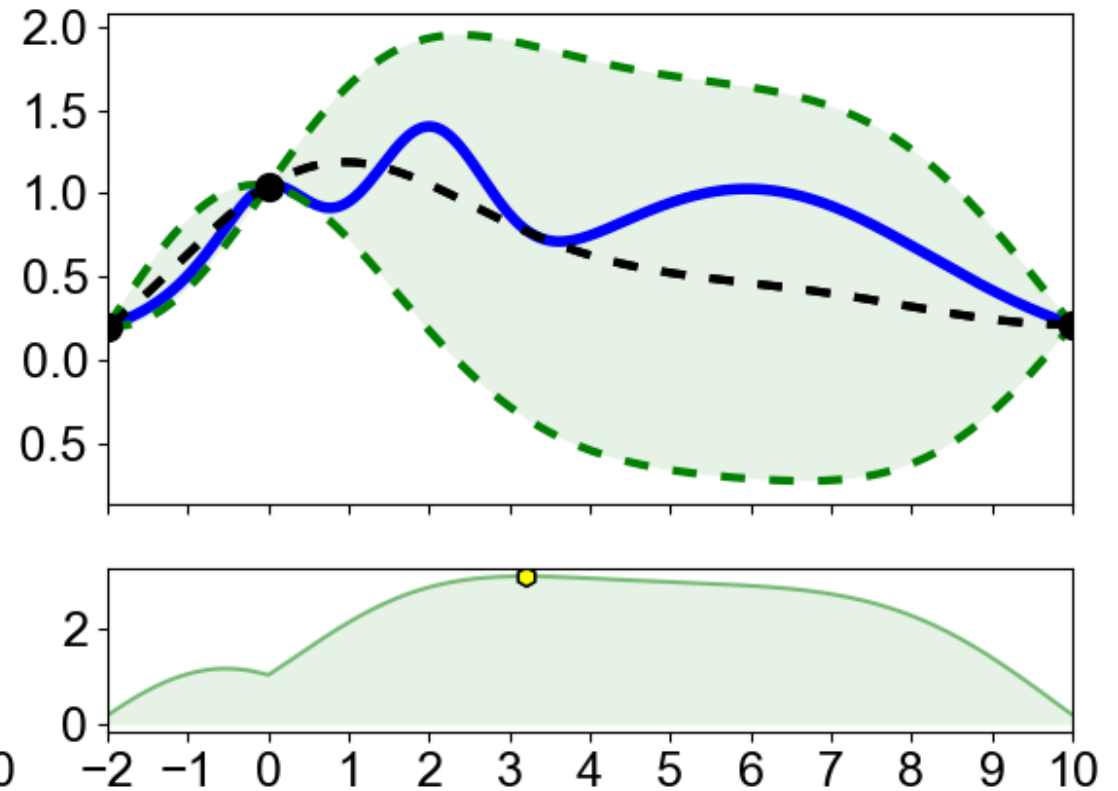
Dominated by variance:

$$\alpha_{UCB} = \mu(x) + \beta \sigma(x)$$

Small length scale



Large length scale



Proposed method

Discrete BO Algorithm

1. **for** $t = 1, 2, \dots, T$ **do**
2. Build a model GP with D_t
3. Sample the next point $x_{t+1} = \operatorname{argmax}_{x \in X \subseteq \mathbb{R}^d} \alpha_{UCB}(\beta_t)$

4. $x_{t+1} = \operatorname{round}(x_{t+1})$ #

5. **if** $x_{t+1} \in D_t$:

Find β, l just

6. **Search for opti**

7. $x_{t+1} = \operatorname{argmax}_{x \in X \subseteq \mathbb{R}^d} \alpha_{UCB}(\beta_t)$ with β, l

8. $x_{t+1} = \operatorname{round}(x_{t+1})$

9. Query the objective function to obtain y_{t+1}

10. Augment data $D_{t+1} = \{D_t, (x_{t+1}, y_{t+1})\}$

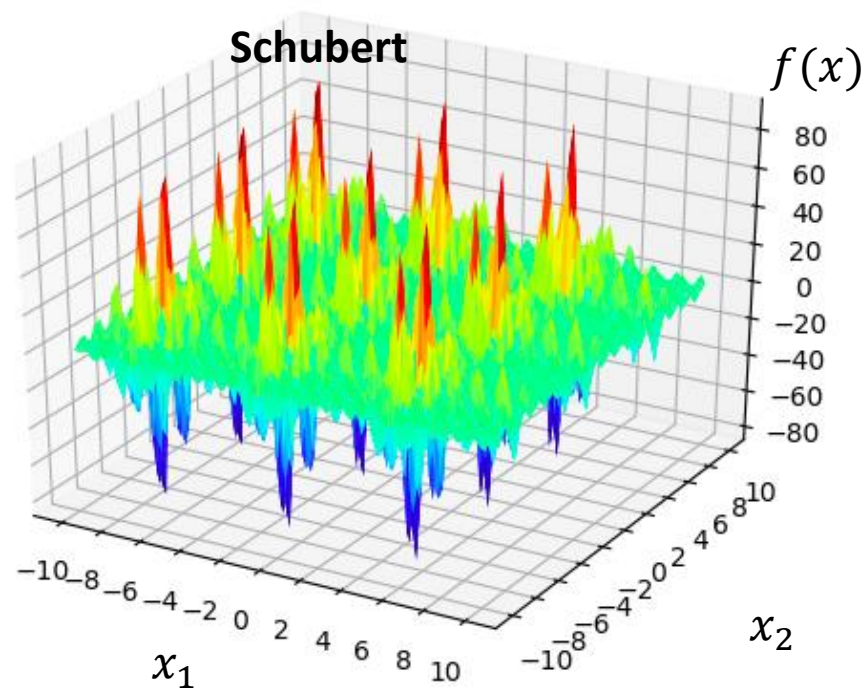
11. **end for**

$$\beta^*, l^* = \operatorname{argmin} g(\beta_t + \Delta\beta, l)$$

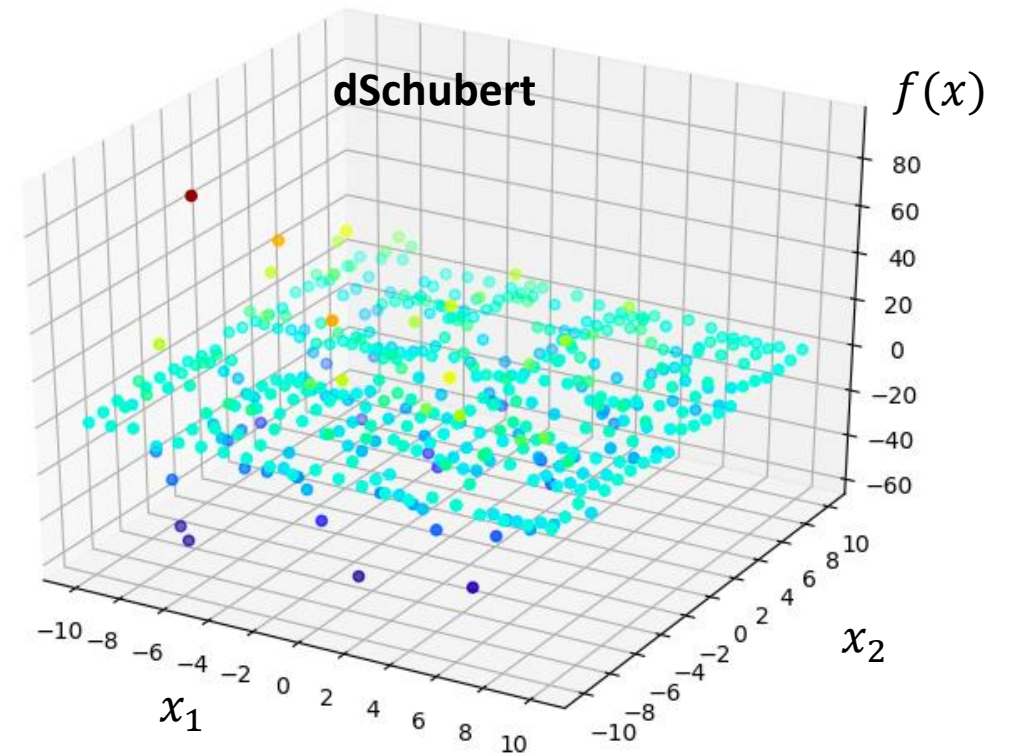
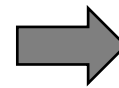
$$g(\beta_t + \Delta\beta, l) = \Delta\beta + \|x_{t+1} - x'_{t+1}\|_2 + P(x'_{t+1})$$

Result – Synthetic functions

Synthetic function	Dimension	Range
dSchubert 2D	2	$x_1, x_2 \in \{-10, \dots, 10\} \subset \mathbb{Z}$

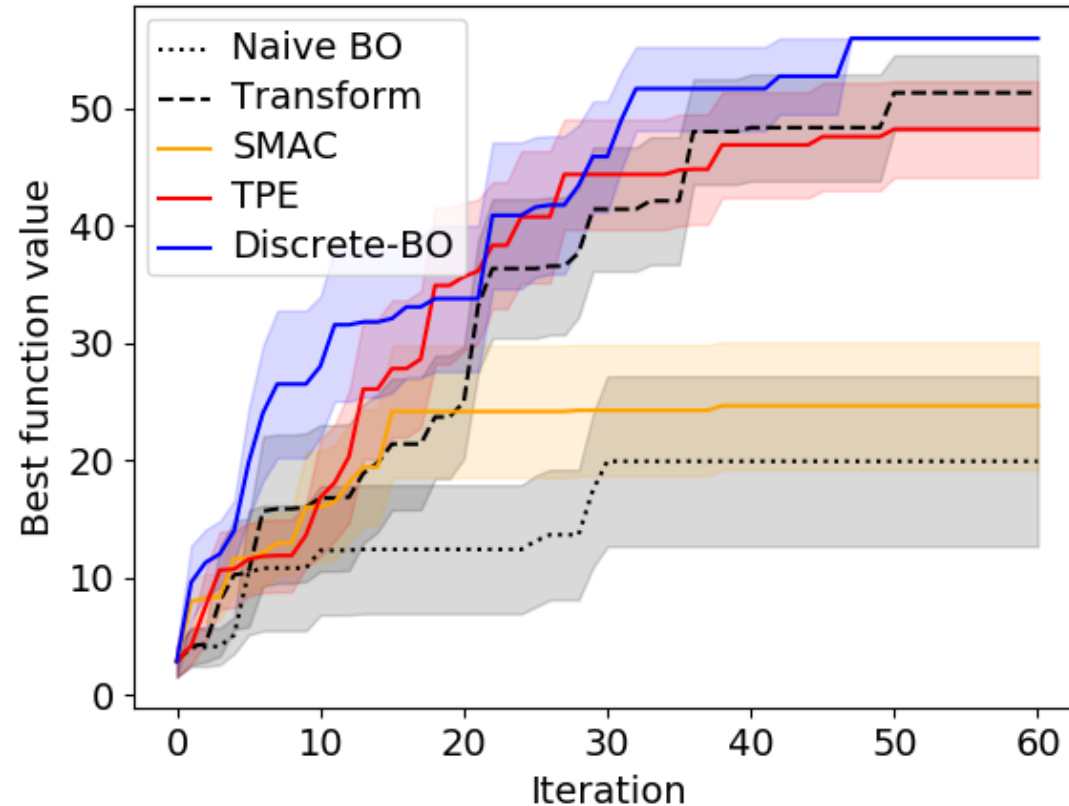


Discretize

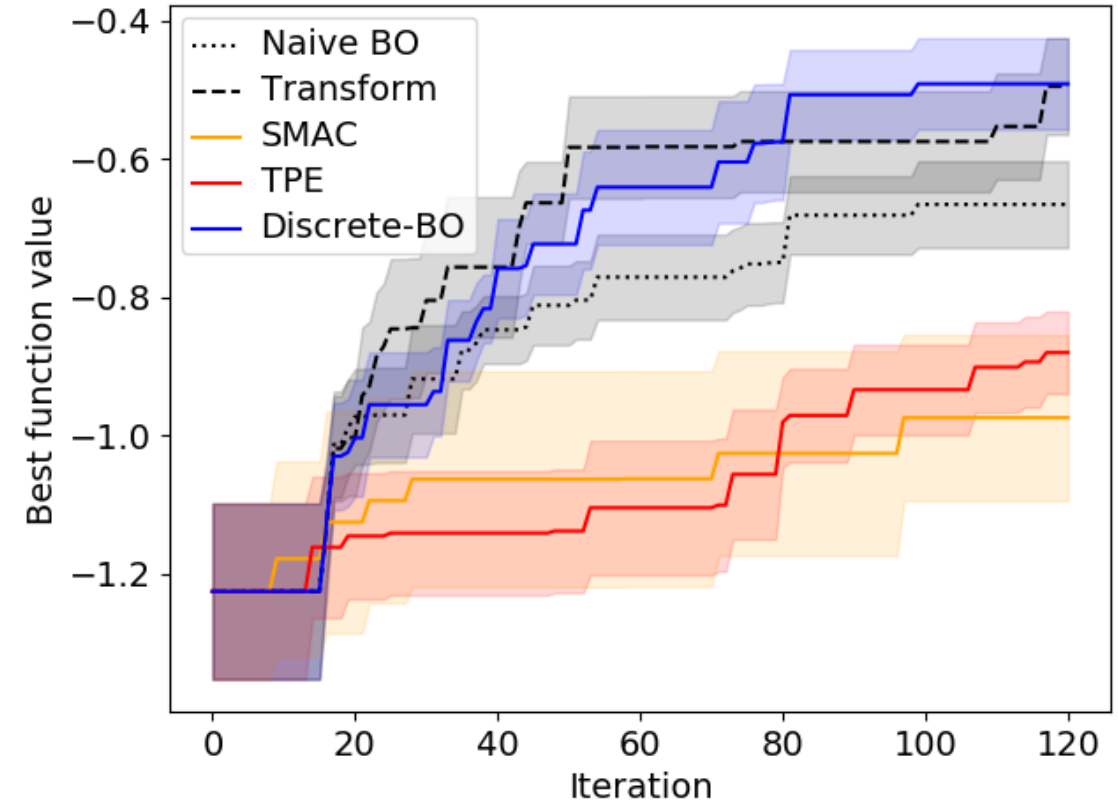


Result – Synthetic functions

Schubert2d

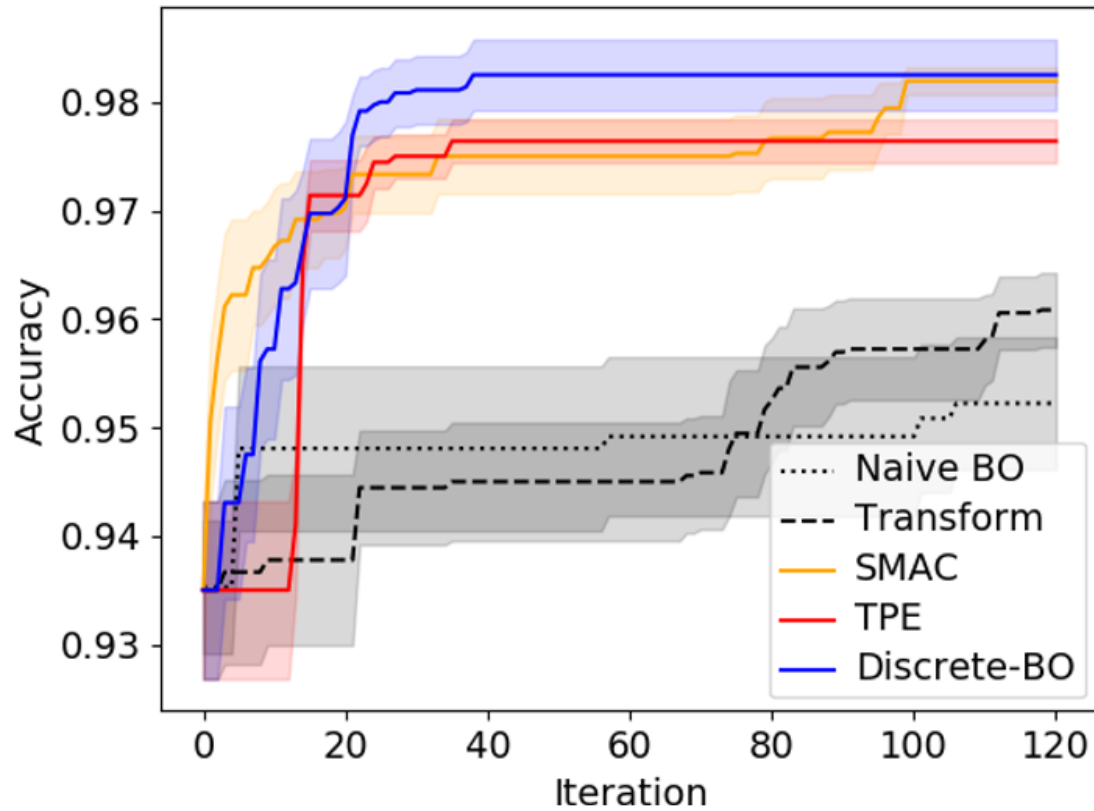


Griewank3d

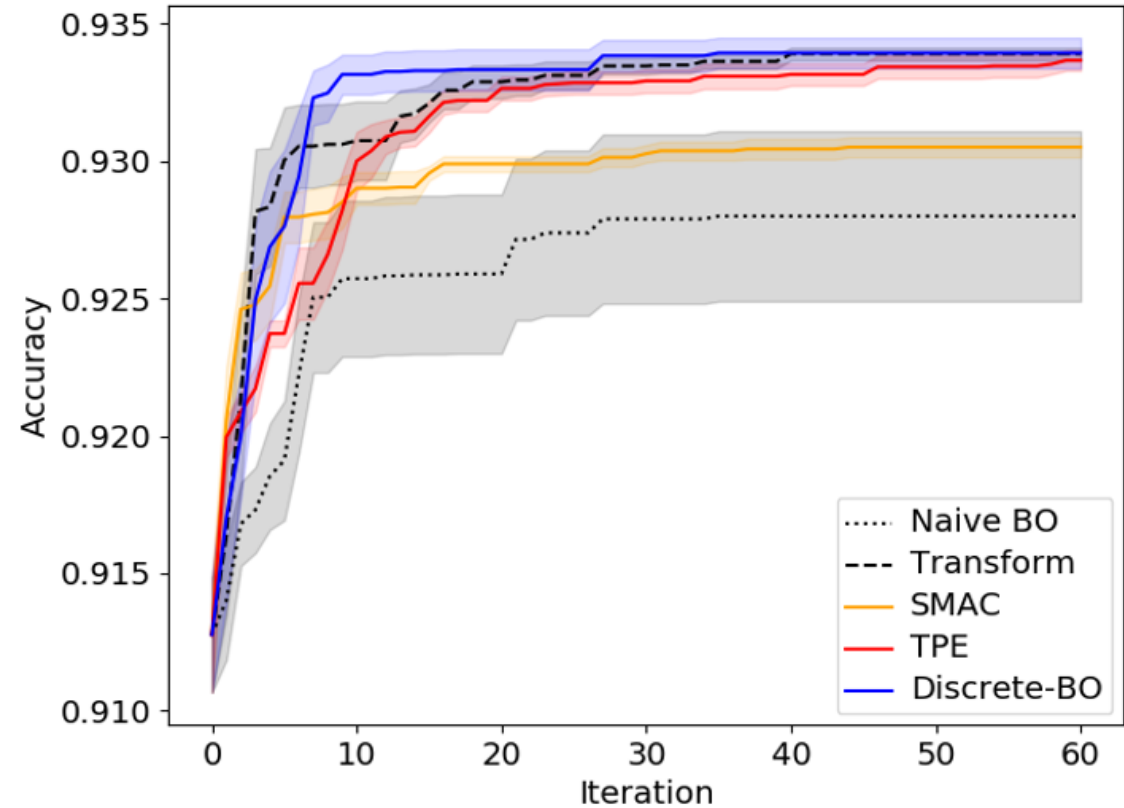


Result – Classifier models

ANN - Digits



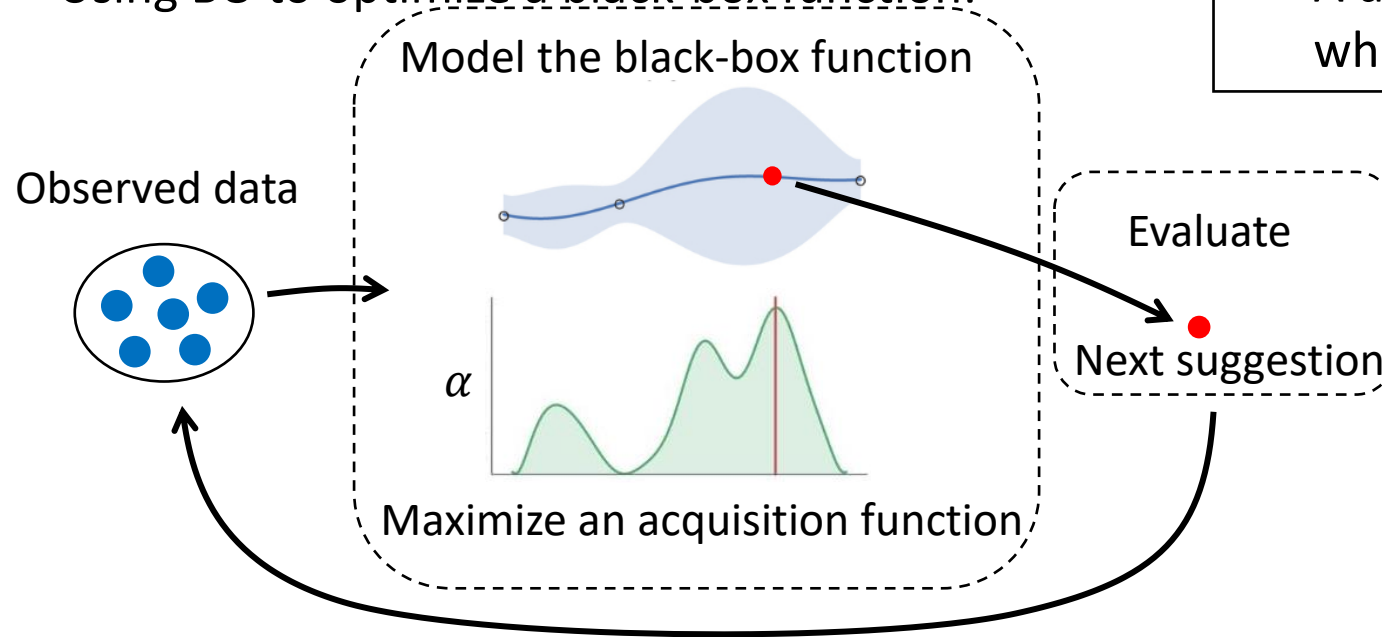
Random forest



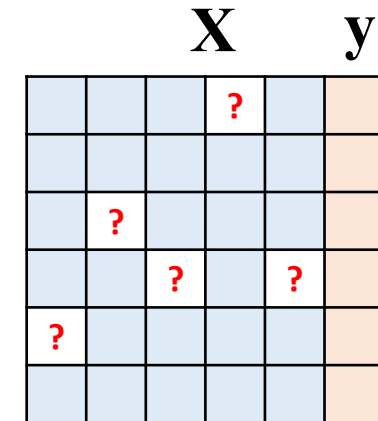
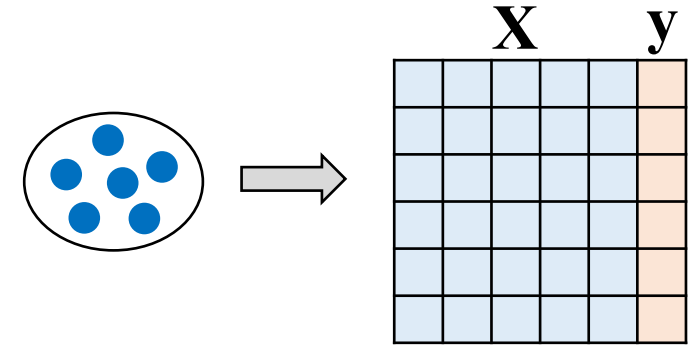
Missing Inputs

Using BO to optimize a black-box function:

- A data point $(\mathbf{x}, y) = \{x_1, x_2, \dots, x_d, y\}$ where $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d, y \in \mathbb{R}$



$$\alpha^{UCB} = \mu(\mathbf{x}) + \beta\sigma(\mathbf{x})$$



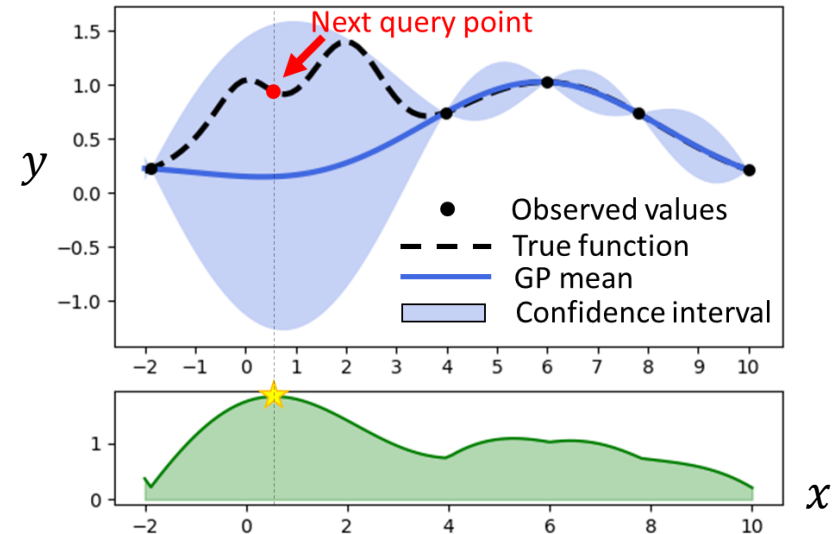
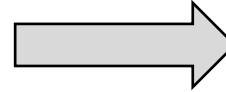
➡ What if there are missing values in observed data?

Missing Inputs

❖ Case M1: Missing values already in historical data

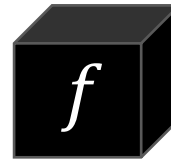
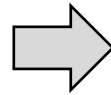
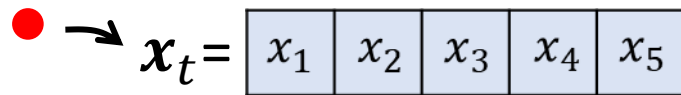
X					y
		?		?	
			?		
	?				
				?	
?					

How to
build the
GP?



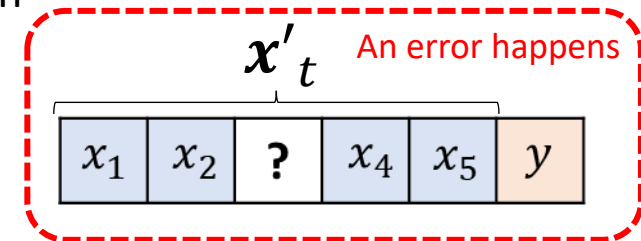
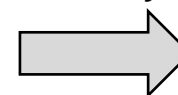
❖ Case M2: Missing values in new observed point x' at iteration t

Next query point



Obtain function

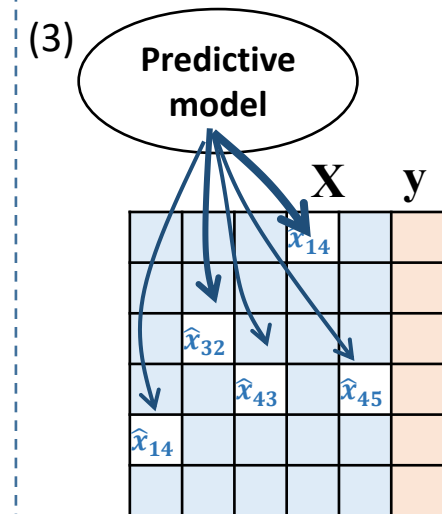
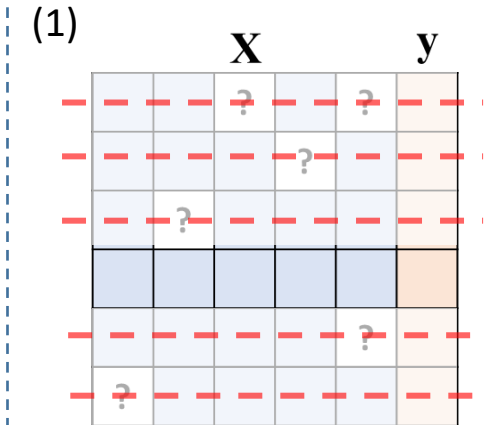
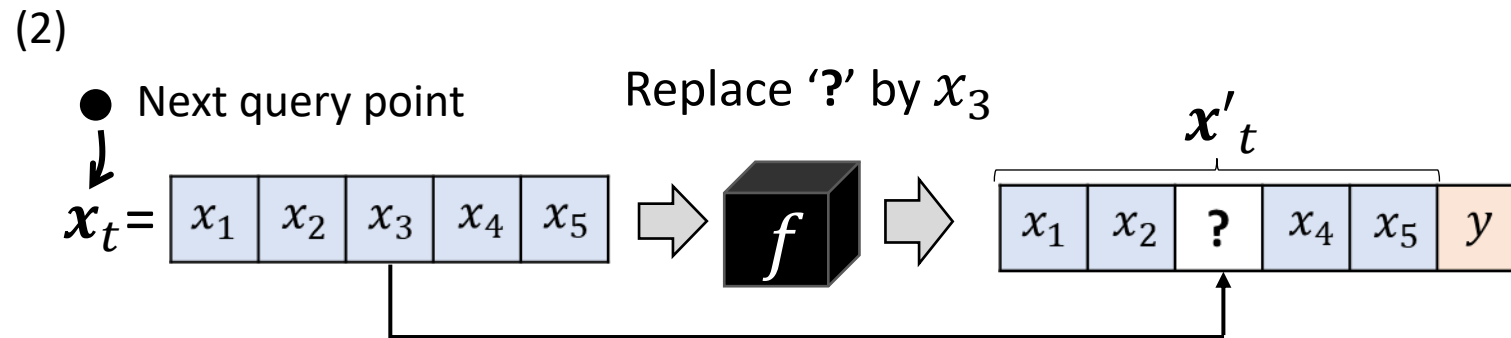
value y



Optimization with Missing Inputs

What if we have missing values?

1. **Remove** observations which have missing values.
2. In case *M2*, we might **use the suggested point** obtained from maximizing the acquisition α (e.g., EI, UCB).
3. **Predict** missing values.



Optimization with Missing Inputs

(1) Remove observations (*DropBO*)

(1)

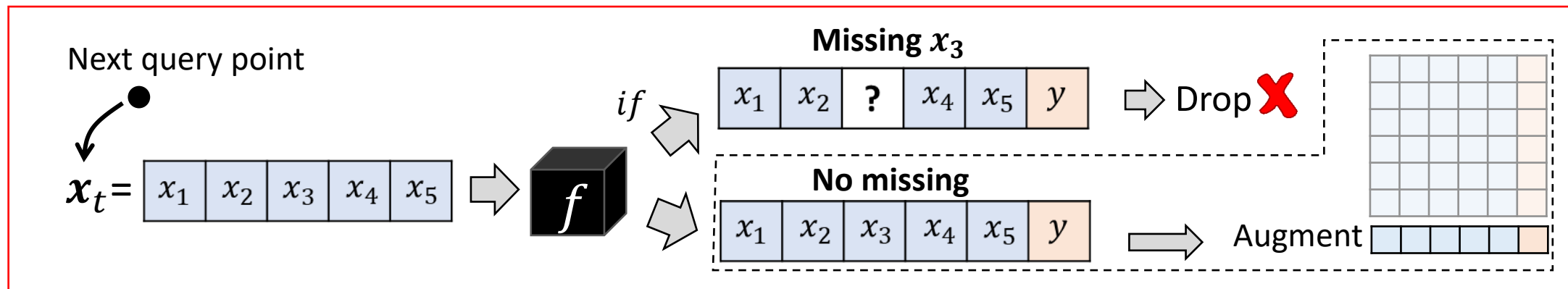
X					y
?					

Pros:

- Very accurate data in observations set

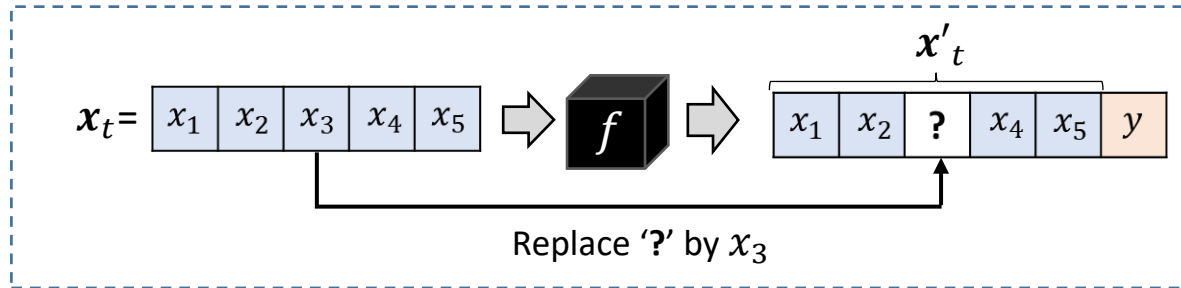
Cons:

- Have less number of observations
- Inefficient if missing rate is high



Optimization with Missing Inputs

(2) Use suggested point (*SuggestBO*)



Pros:

- Perform well when missing noise is *small*
- No need to train a predictive model

Cons:

- Inefficient if missing noise is *high*
- Unable to solve case *M1* (missing values in historical data)

Given $\mathbf{x}_t \in \mathcal{X}^d$, when evaluate \mathbf{x}_t , we obtain $y_t = f(\mathbf{x}'_t)$ instead of $y_t = f(\mathbf{x}_t)$. *SuggestBO* performs well only when $\|\mathbf{x}_t - \mathbf{x}'_t\|_2$ is small.

Optimization with Missing Inputs

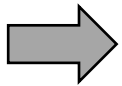
(3) Predict missing values

How to correctly predict missing values?

Imputations

- Mean/mode: does not consider the **correlation** between values.
- K-nearest neighbors: depends on **available data** and **distance metric**.
- Regression (e.g., random forest, neural network): requires **many training data**.

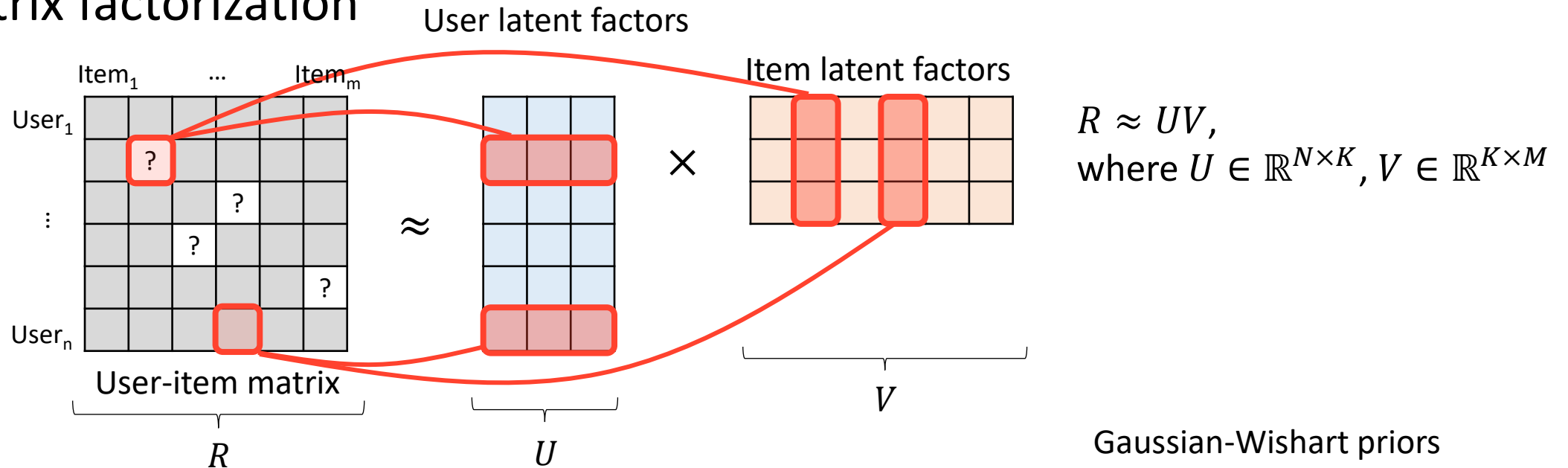
BO under uncertain inputs (BO-uGP^[1]): unable to work if **actual value is unknown**.



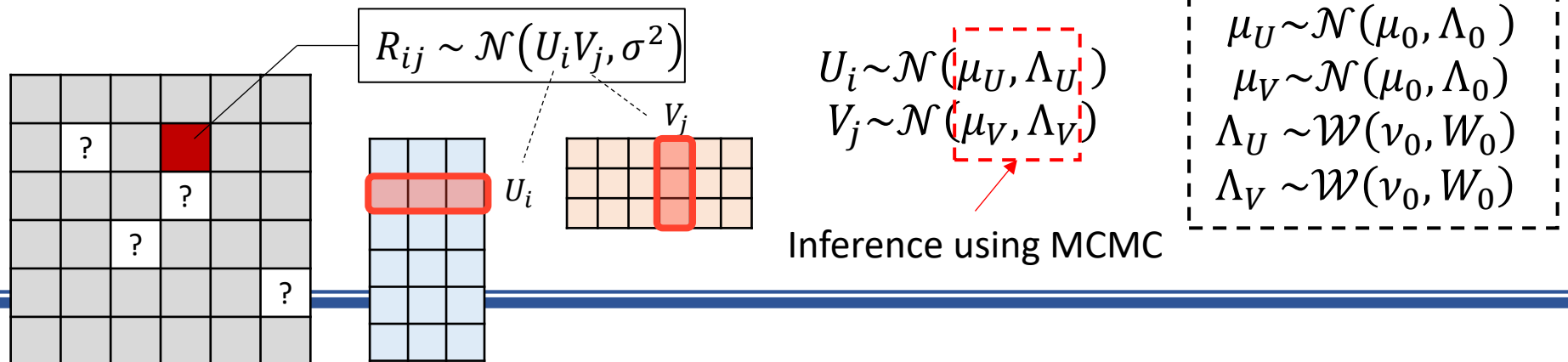
Find a way which takes into account the **correlation** between values
to build the **distribution of missing data**

Bayesian Probabilistic Matrix Factorization (BPMF)

Matrix factorization

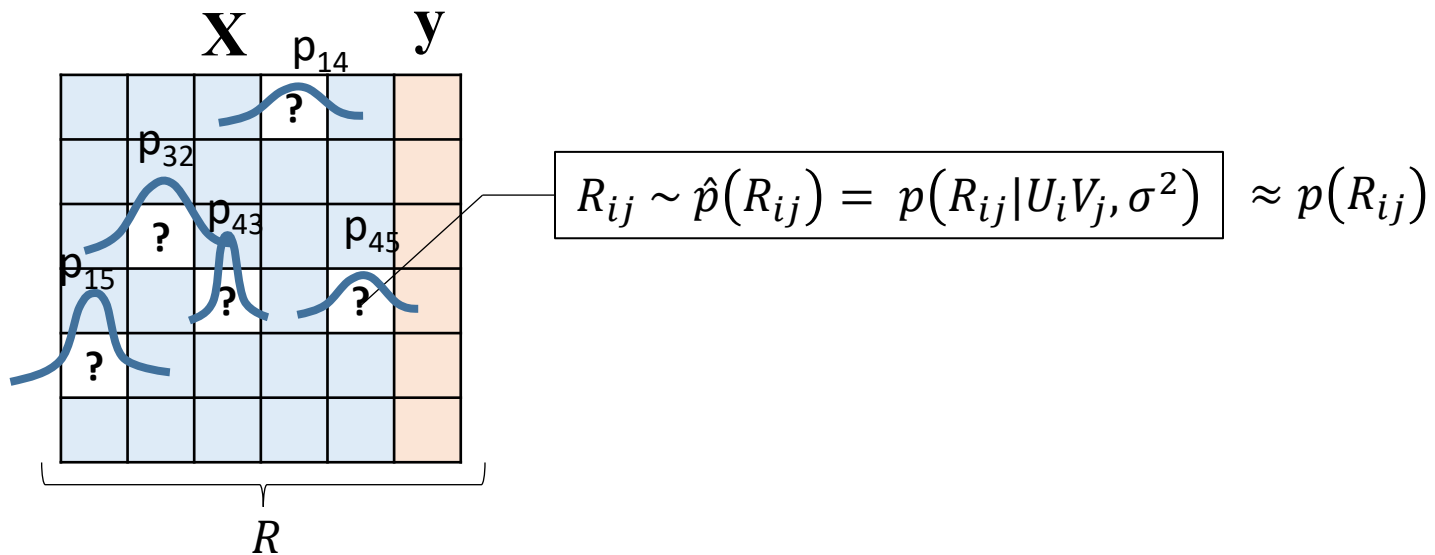


BPMF



Our method

Consider observations set $[\mathbf{X}, \mathbf{y}]$ as the matrix R , we can compute the probability distribution of each value R_{ij} in the matrix in the same way as BPMF:



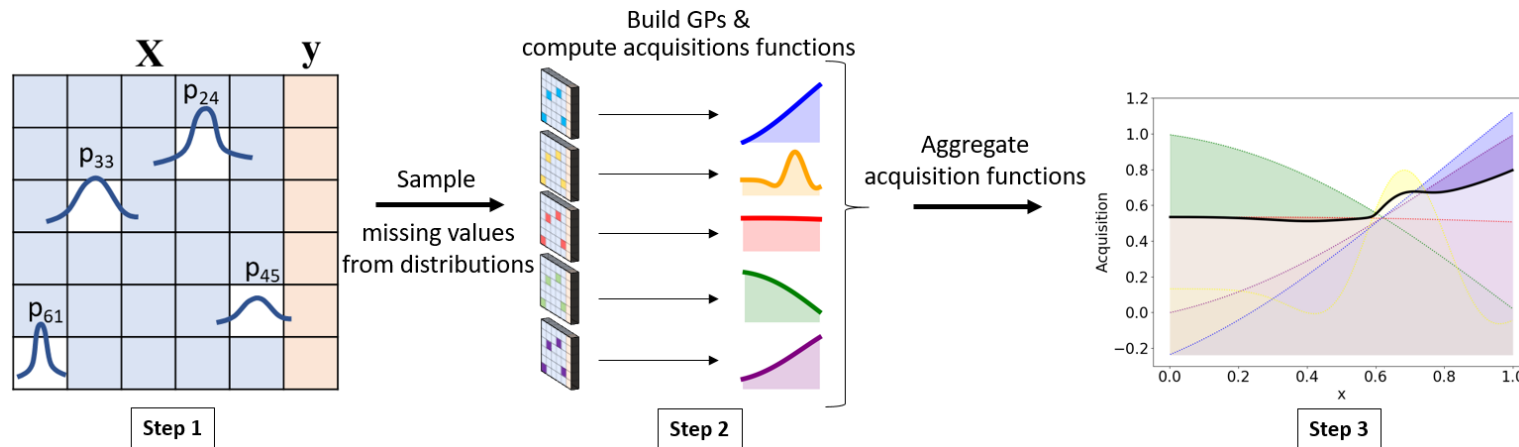
New values R_{ij} can be **directly sampled** from the distributions. However, it is **uncertain** to randomly sample one value as a substitution to the missing value.

Our method

Instead of sampling **one** value for each missing location, we can sample **n** values.

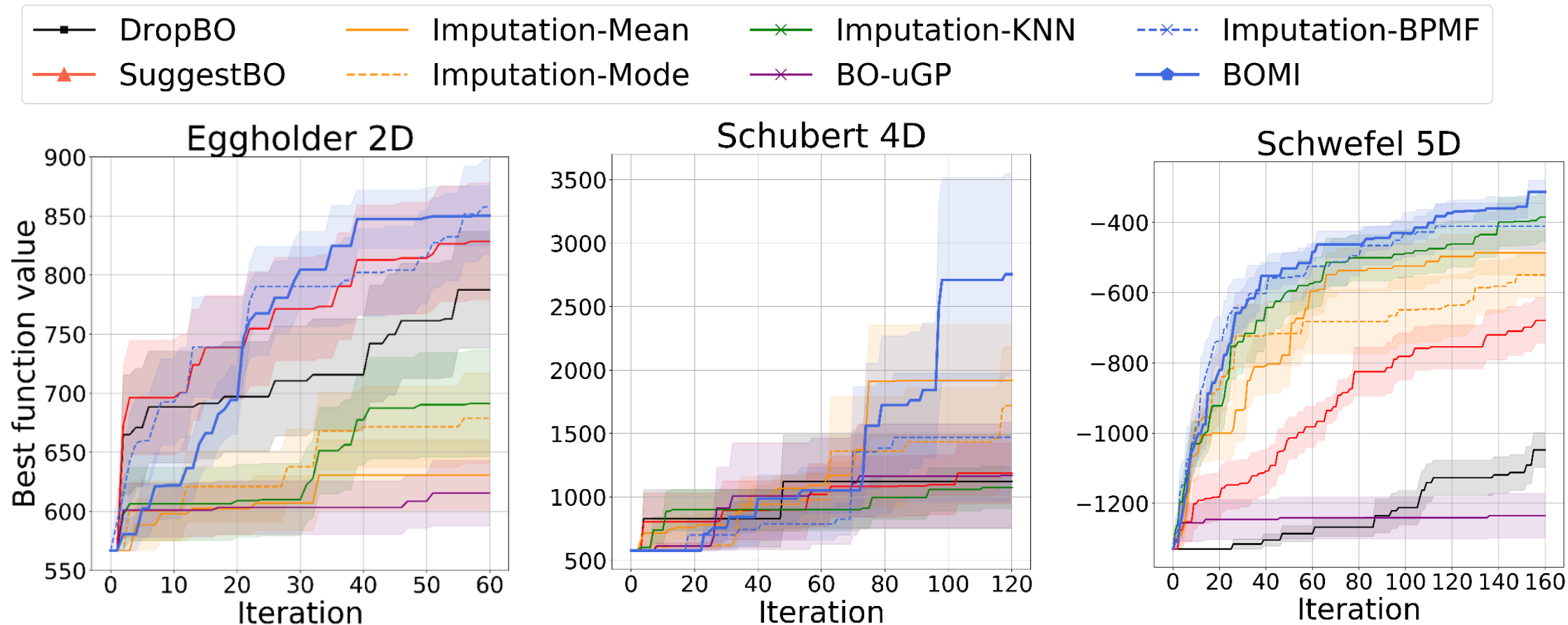
We use the following acquisition function: Agreement tradeoff term

$$\begin{aligned}\alpha^{UCB-MI} &= \mu_{\alpha}(\alpha^{UCB}(x)) + \beta_{\alpha} \sigma_{\alpha}(\alpha^{UCB}(x)) \\ &= \frac{1}{Q} \sum_{q=1}^Q (\alpha_q^{UCB}(x)) + \beta_{\alpha} \sqrt{\frac{\sum_{q=1}^Q (\alpha_q^{UCB}(x) - \frac{1}{Q} \sum_{q=1}^Q \alpha_q^{UCB}(x))^2}{Q-1}}\end{aligned}$$



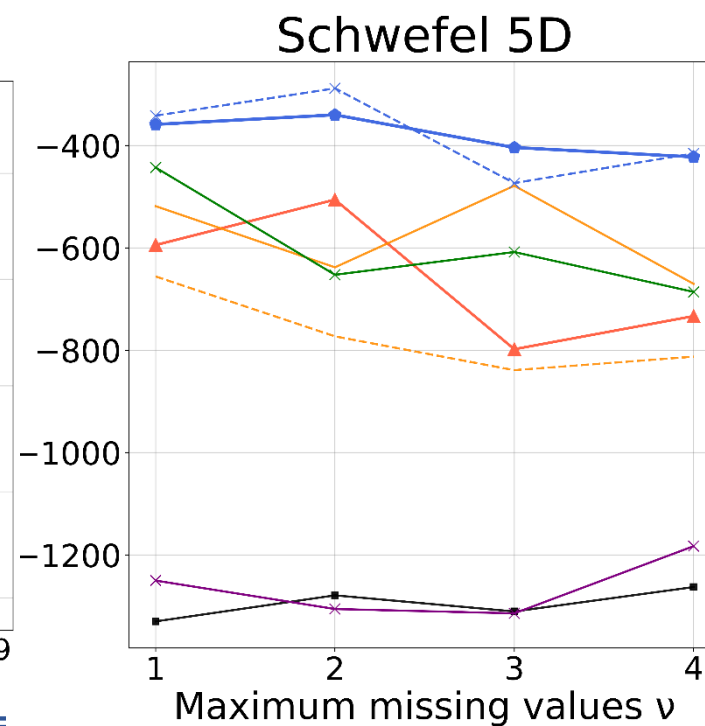
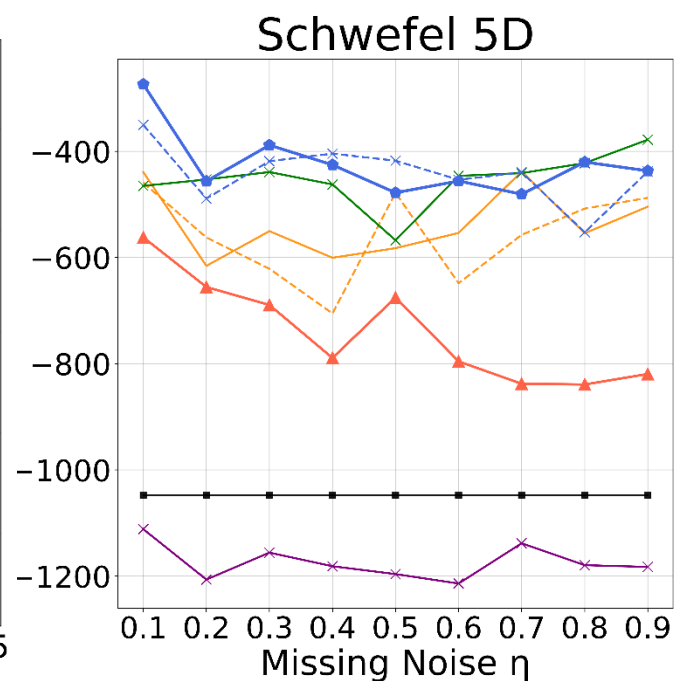
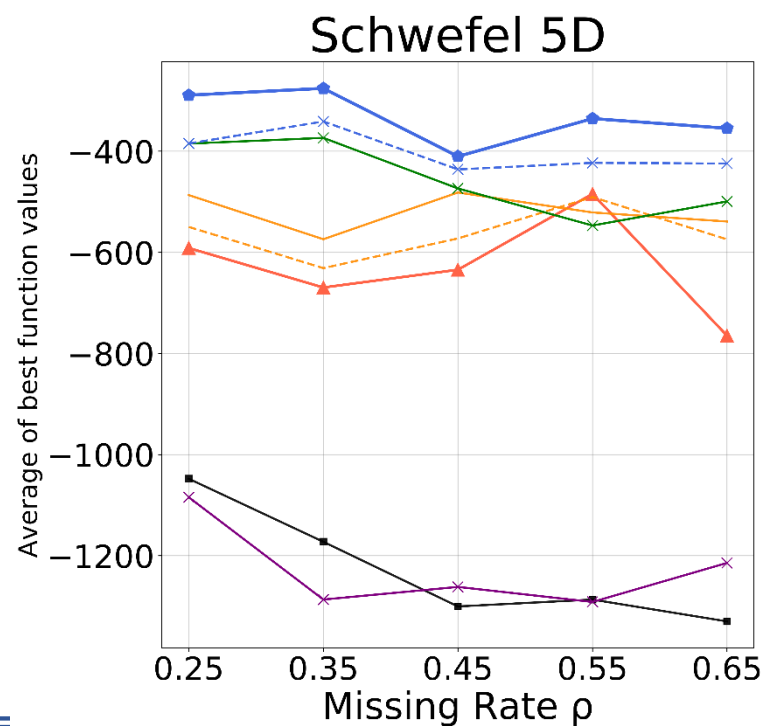
Experiments

Comparison on synthetic functions

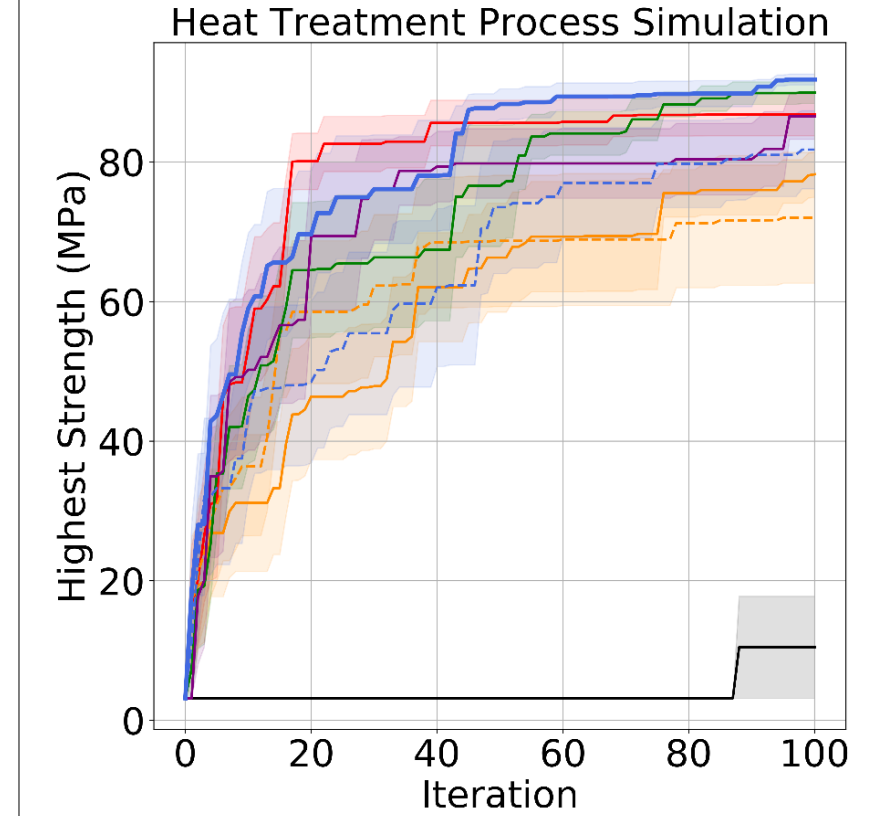
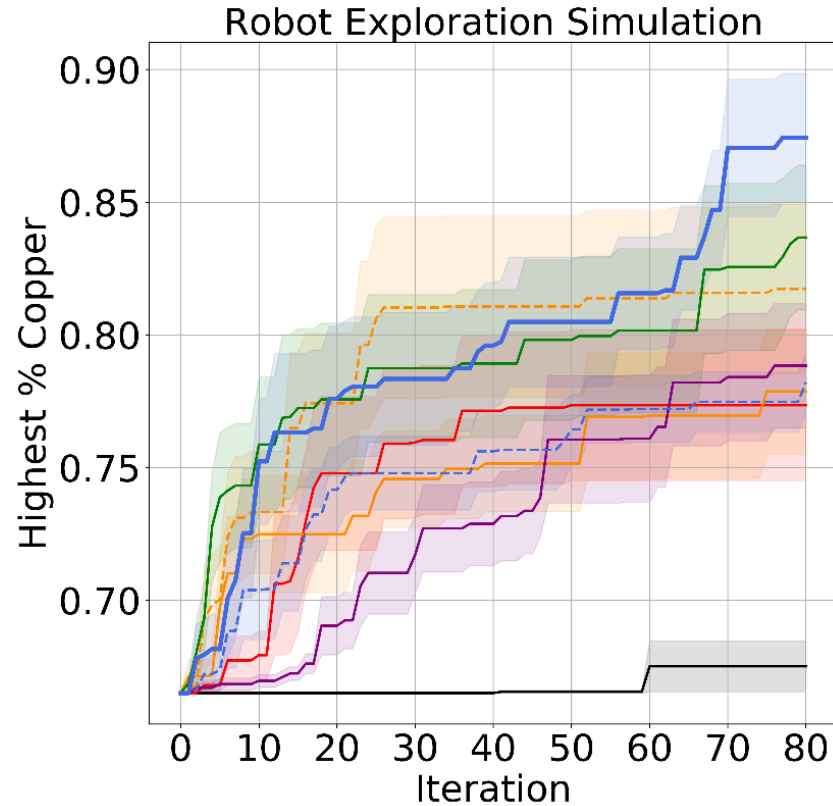
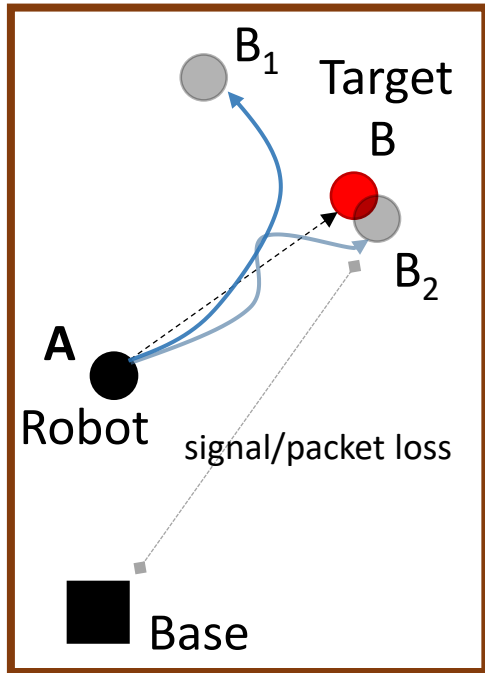


Experiments

Stability



Experiments – Simulation of real problems



THANK YOU!!
