

# Monte Carlo Tree Search based Variable Selection for High Dimensional Bayesian Optimization

Lei Song\*, Ke Xue\*, Xiaobin Huang, Chao Qian<sup>†</sup>  
Email: [qianc@lamda.nju.edu.cn](mailto:qianc@lamda.nju.edu.cn)

School of Artificial Intelligence  
Nanjing University, China



# Black-box Optimization

We consider the following problem formulation for black-box optimization (BBO):

$$\max_{x \in \mathcal{X}} f(x)$$

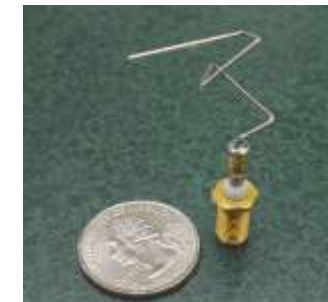
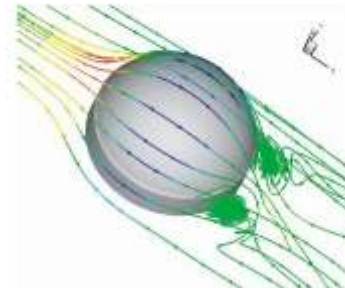
where only the evaluation  $f(x)$  is available and **no additional information** is known

The methods based on first or second order information (e.g., gradient descent) can not be used

Traditional BBO algorithms:

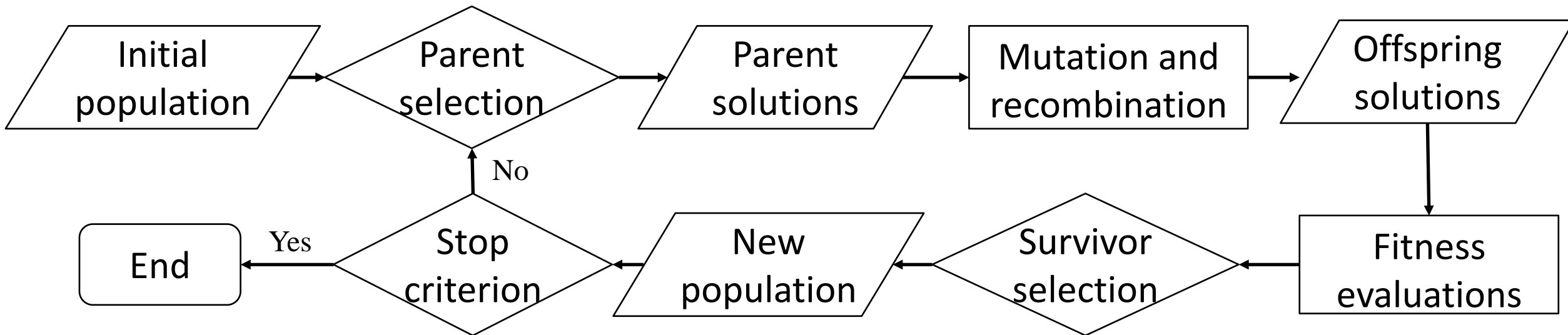
- Evolutionary algorithms
- Evolutionary strategies
- Bayesian optimization

Application:



# Evolutionary Algorithms

**Evolutionary algorithms (EAs)** are a kind of randomized heuristic optimization algorithms, inspired by nature evolution (reproduction with variation + nature selection)



- Population-based and easy to be parallelized
- Discrete inputs

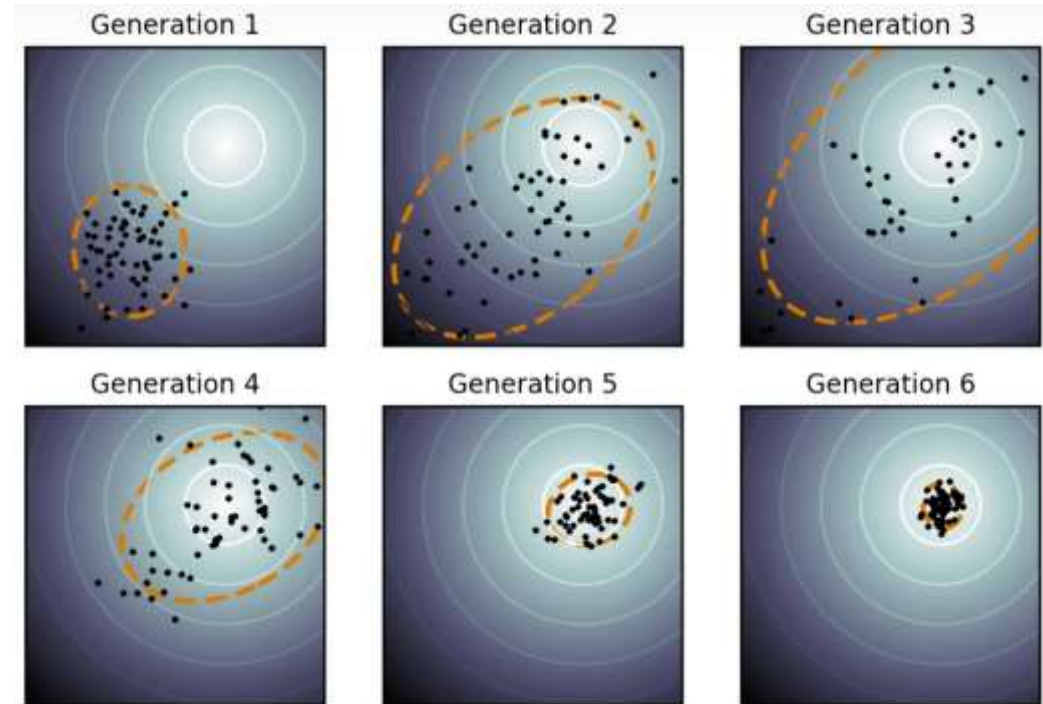
# Evolutionary Strategies

**Evolutionary strategies (ES)** are a popular variant of EAs for continuous problems

CMA-ES is one of the most popular algorithms in ES

CMA-ES samples offspring from  $N(m, \Sigma)$  and uses the adaptive covariance matrix to balance exploration and exploitation

- Hundreds of continuous parameters



# Bayesian Optimization

**Bayesian optimization (BO)** uses surrogate model to approximate  $f$  and obtains the next query point via acquisition function

---

## Algorithm 1 BO Framework

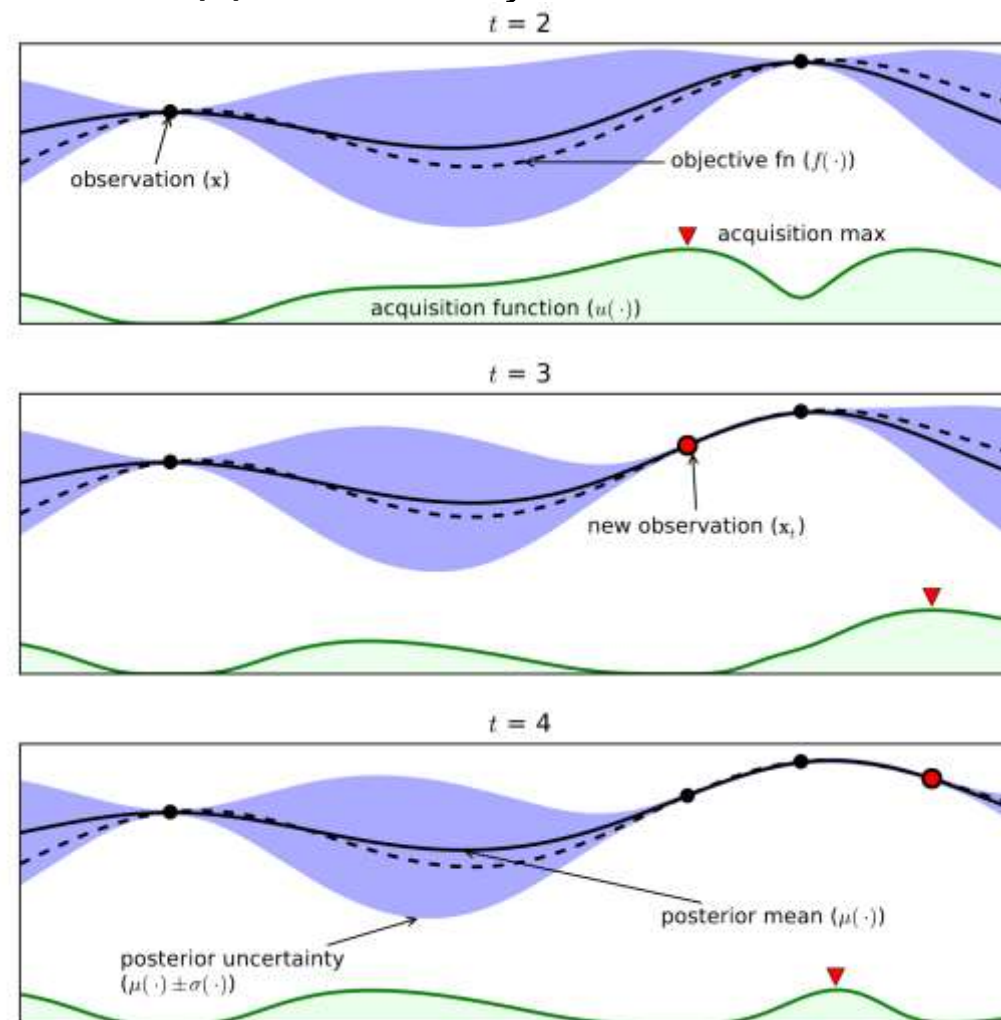
---

**Input:** iteration budget  $T$

**Process:**

- 1: let  $D_0 = \emptyset$ ;
  - 2: **for**  $t = 1 : T$  **do**
  - 3:    $\mathbf{x}_t = \arg \max_{\mathbf{x} \in \mathcal{X}} \text{acq}(\mathbf{x})$ ;   obtain next point
  - 4:   evaluate  $f$  at  $\mathbf{x}_t$  to obtain  $y_t$ ;
  - 5:   augment the data  $D_t = D_{t-1} \cup \{(\mathbf{x}_t, y_t)\}$  and update the GP model   update model
  - 6: **end for**
- 

- Less than 30 parameters
- Expensive evaluation



# Black-box Optimization

---

Algorithms	Runtime	Problem scale	Type of variables
EAs	short	large	discrete
ES	short	large	continuous
BO	long	small	continuous

There are still many challenges for BBO in real-world problems, e.g.,

- Multi-objective problems
- High-dimensional problems
- Complex constraints



# High-dimensional Bayesian Optimization

---

Scaling BO to high-dimensional problems is a challenge:

- Search space increases exponentially
- Computation cost of fitting GP and optimizing the acquisition function is time-consuming

Current approaches usually solve high-dimensional BO in a low-dimensional subspace:

1. Obtain a low-dimensional subspace
2. Optimize in the low-dimensional subspace
3. Project the low-dimensional solution back to the high-dimensional space

# High-dimensional Bayesian Optimization

---

Different approaches are based on different assumptions to obtain the low-dimensional subspace:

- **Decomposition:**  $f$  can be decomposed into the sum of low-dimensional functions
- **Embedding:** only a few dimensions affect  $f$  significantly
- **Variable selection:** only a few *axis-aligned* dimensions affect  $f$  significantly



## Decomposition

---

**Add-GP-UCB** assumes that  $f$  can be decomposed into the sum of **disjoint** low-dimensional functions

$$f(\mathbf{x}) = f^1(\mathbf{x}^{(1)}) + \dots + f^k(\mathbf{x}^{(k)}), \forall i, j, \mathbf{x}^{(i)} \cap \mathbf{x}^{(j)} = \emptyset$$

- Maximize the likelihood to learn a low-dimensional decomposition
- Optimize the low-dimensional functions separately
- Concatenate variables of low-dimensional functions

$$\begin{aligned} f(\mathbf{x}) &= f^{(1)}(x_1, x_3, x_4) \\ &+ f^{(2)}(x_2) \\ &+ f^{(3)}(x_5, x_6) \end{aligned}$$

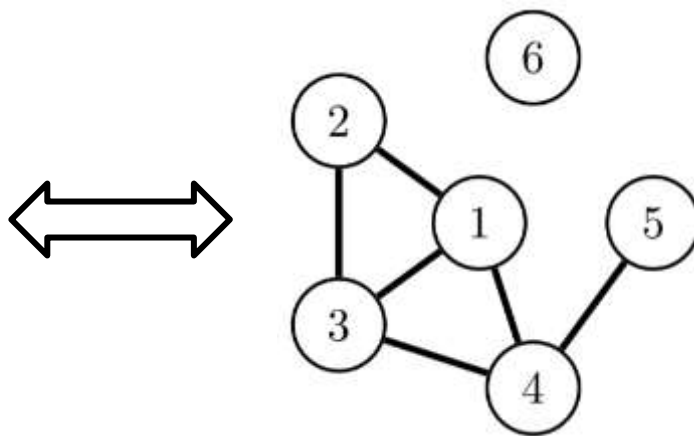
## Decomposition

**Overlapping** generalizes Add-GP-UCB to **overlapping** conditions

$$f(\mathbf{x}) = f^1(\mathbf{x}^{(1)}) + \dots + f^k(\mathbf{x}^{(k)})$$

- Maximize the likelihood to learn a low-dimensional decomposition
- Optimize the low-dimensional functions on the graph similar to message passing
- Concatenate variables of low-dimensional functions

$$\begin{aligned} f(\mathbf{x}) &= f^{(1)}(x_1, x_2, x_3) \\ &+ f^{(2)}(x_1, x_3, x_4) \\ &+ f^{(3)}(x_4, x_5) \\ &+ f^{(4)}(x_6) \end{aligned}$$



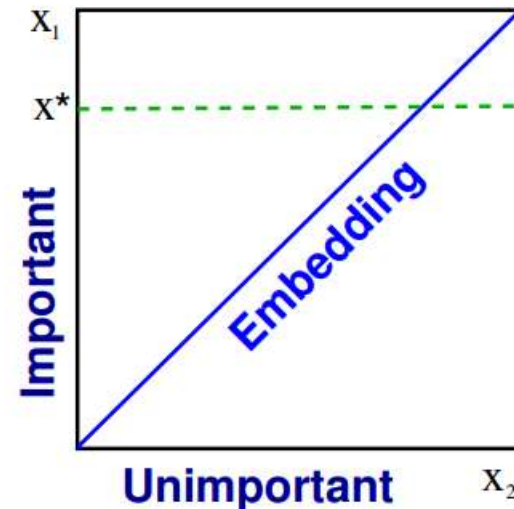
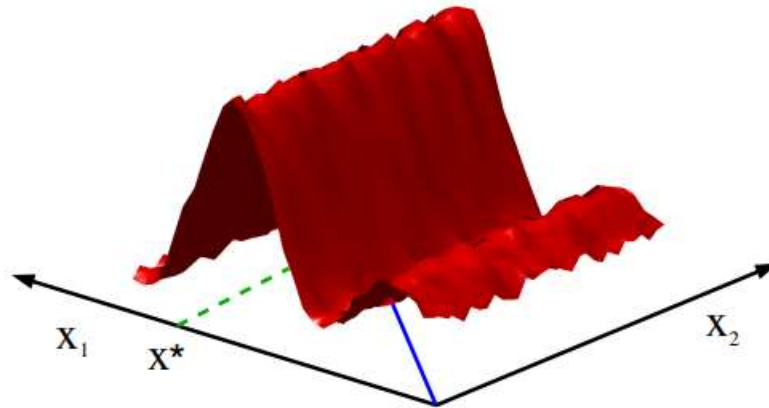
# Embedding

Assuming low effective dimensionality, **REMBO** uses a random embedding matrix to obtain the low-dimensional subspace:

$$\mathbf{M} \in \mathbb{R}^{D \times d}, \mathbf{M}_{ij} \sim N(0, 1)$$

Then, the optimization problem is:

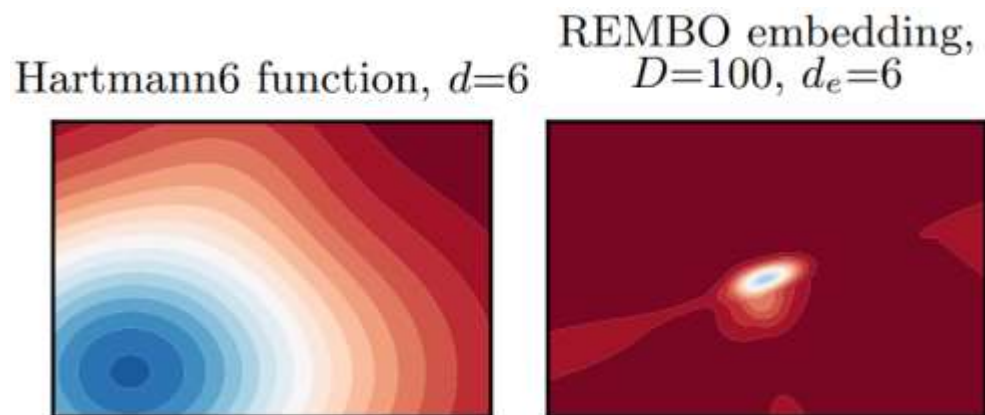
$$\max_{z \in \mathbb{R}^d} f(\mathbf{M}z)$$



# Embedding

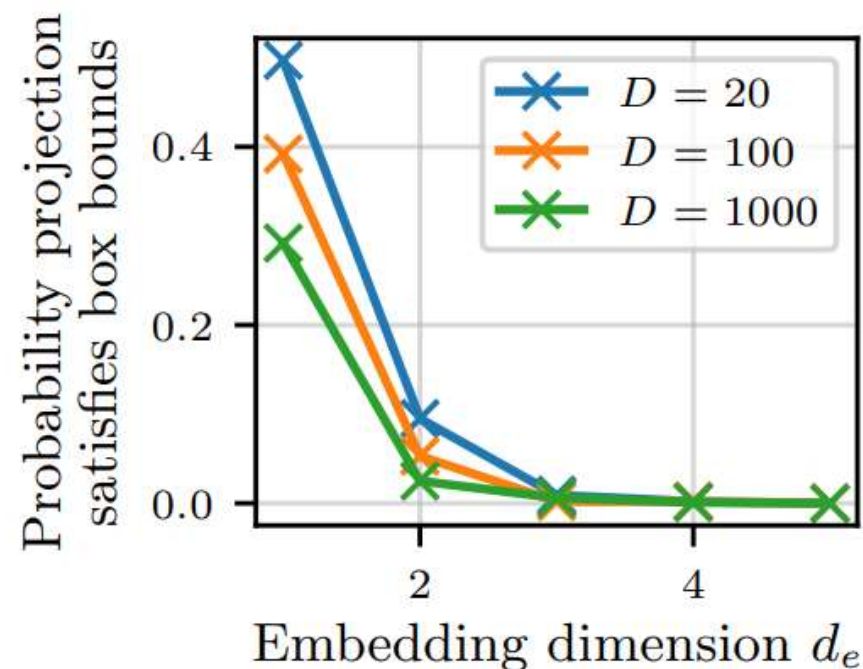
ALEBO improves several misconceptions in previous embedding methods, e.g.,

- The box bounds result in a nonlinear distortion of the search space
- Many points map to the facets



Methods:

- A Mahalanobis kernel
- A constrained acquisition function optimization



## Variable Selection

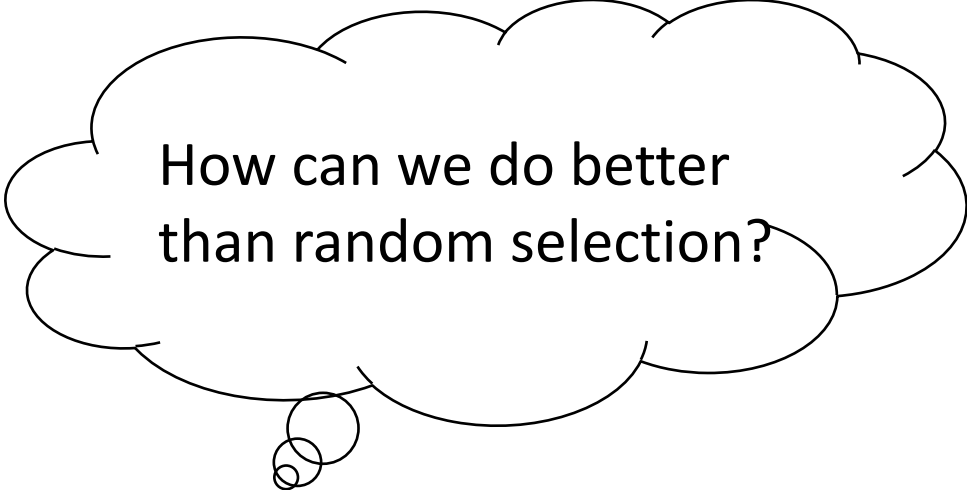
---

Dropout: select  $d$  variables randomly and optimize the selected variables

- Select  $d$  variables randomly
- Optimize the selected variables
- Use “fill-in” strategy to obtain the unselected variables

Advantage:

- Variable selection is much simpler than embedding and can reduce the runtime



How can we do better  
than random selection?

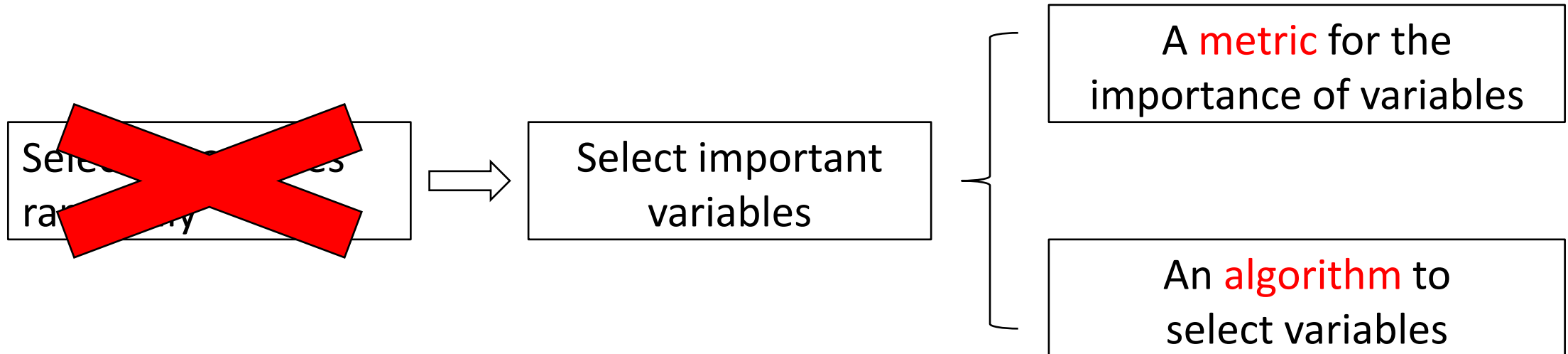
## Motivation

---

The importance of different variables are different

Thus, randomly select variables is inefficient

**We should pay more attention to the **important** variables!**



## MCTS-VS

---

**Monte Carlo Tree Search based Variable Selection (MCTS-VS)** uses MCTS to iteratively select and optimize a subset of important variables, and uses “fill-in” strategy to obtain the unselected variables

- Variable score  $s$  is the metric of the importance of variables
- MCTS is employed to partition the variables into important and unimportant ones, select and optimize the important variables
- “Fill-in” strategy for unselected variables



## Variable Score

---

**Variable score**  $\mathbf{s} \in \mathbb{R}^D$  is a  $D$ -dimensional vector, where the  $i$ -th element represents the importance of the  $i$ -th variable

$$\mathbf{s} = \left( \sum_{(\mathbb{M}, \mathcal{D}) \in \mathbb{D}} \sum_{(x^i, y^i) \in \mathcal{D}} y^i \cdot g(\mathbb{M}) \right) / \left( \sum_{(\mathbb{M}, \mathcal{D}) \in \mathbb{D}} |\mathcal{D}| \cdot g(\mathbb{M}) \right)$$

The sum of query evaluations  
optimizing the variables indexed by  $\mathbb{M}$

The number of queries using  
each variable

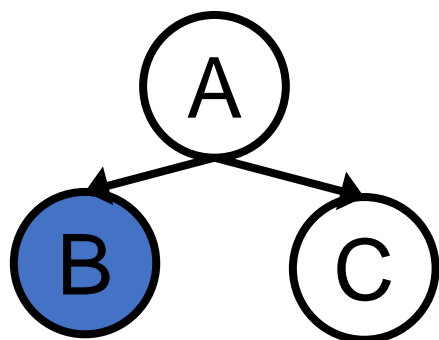
- $(\mathbb{M}, \mathcal{D})$  represents the indices of selected variables and the corresponding samples
  - E.g.,  $\mathbb{M} = \{2, 5, 7\}$  and  $\mathcal{D}$  is obtained by optimizing  $\left\{ ([x_2, x_5, x_7]^i, y^i) \right\}_{i=1}^{t-1}$
- $g: 2^{[D]} \rightarrow \{0, 1\}^D$ , and the  $i$ -th element is 1 if  $i \in \mathbb{M}$ , and 0 otherwise

## A Brief Introduction to MCTS

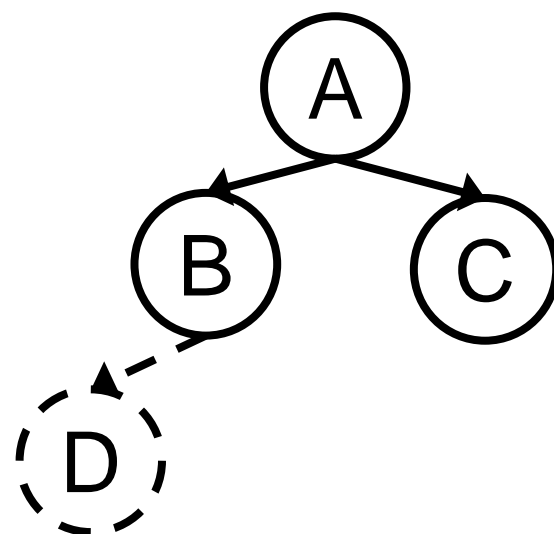
Tree node  $X$  represents the state, and stores  $v_X$  representing its goodness and the number  $n_X$  of visits

UCB is used to select node, balancing the **exploitation** and **exploration**:

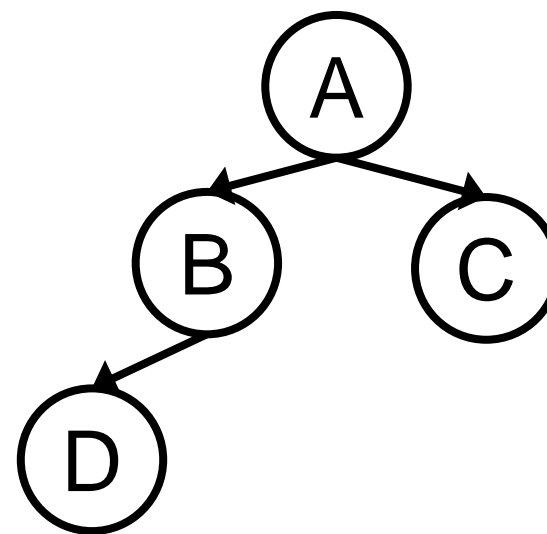
$$v_X + 2C_p \sqrt{2(\log n_p)/n_X}$$



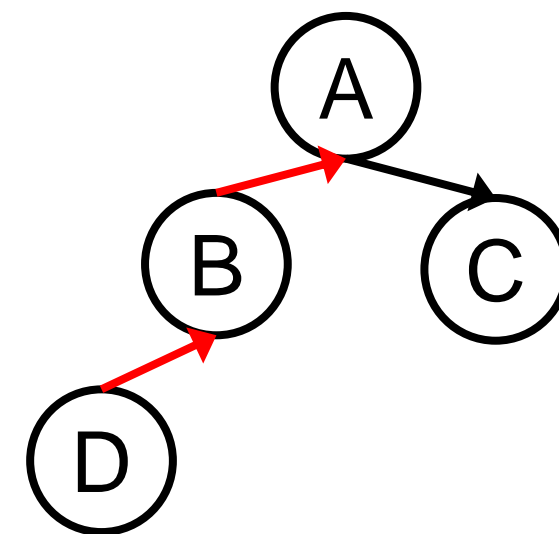
*Selection*



*Expansion*



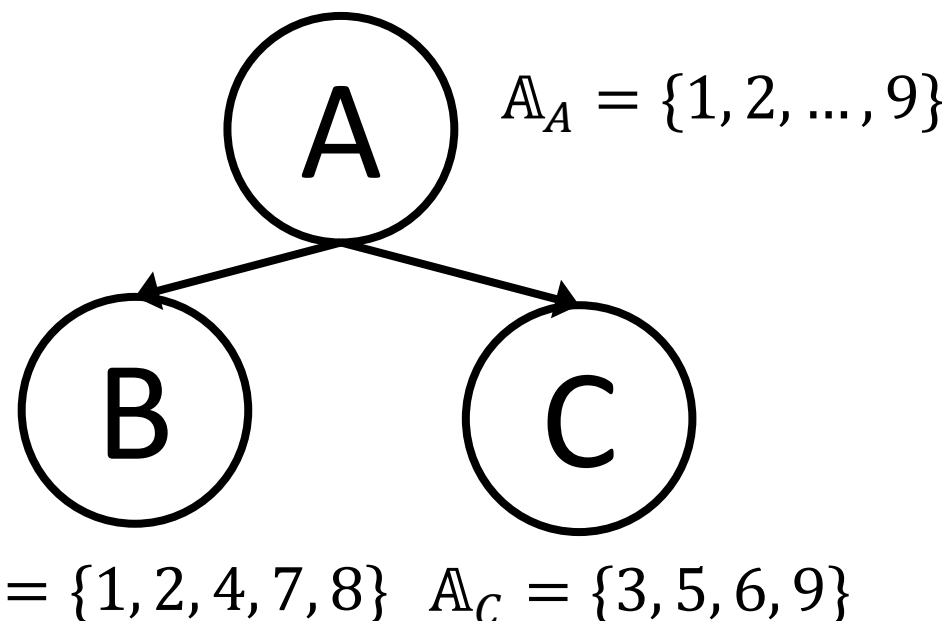
*Simulation*



*Back-propagation*

## Example Illustration of MCTS-VS

Current state of MCTS-VS



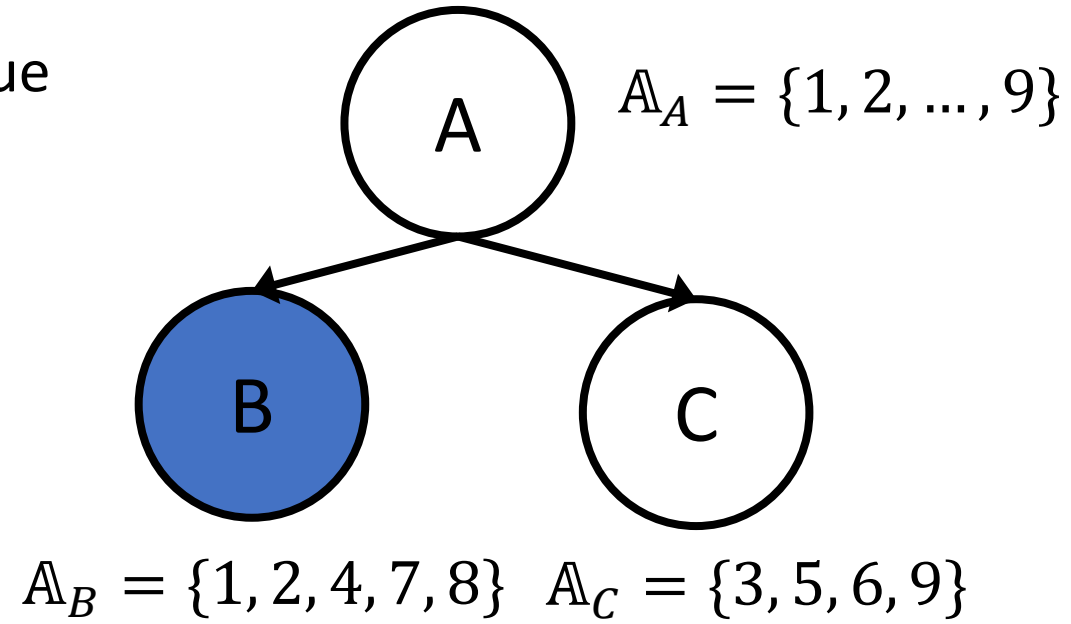
Tree node  $X$  represents a subset of variables, denoted by index set  $A_X \subseteq [D]$

- The root node represents all variables
- $v_X$  is defined as the average score (i.e., importance) of variables contained by  $X$ , which is calculated by  $s \cdot g(A_X)/|A_X|$
- $n_X$  is the number of visits

## Example Illustration of MCTS-VS

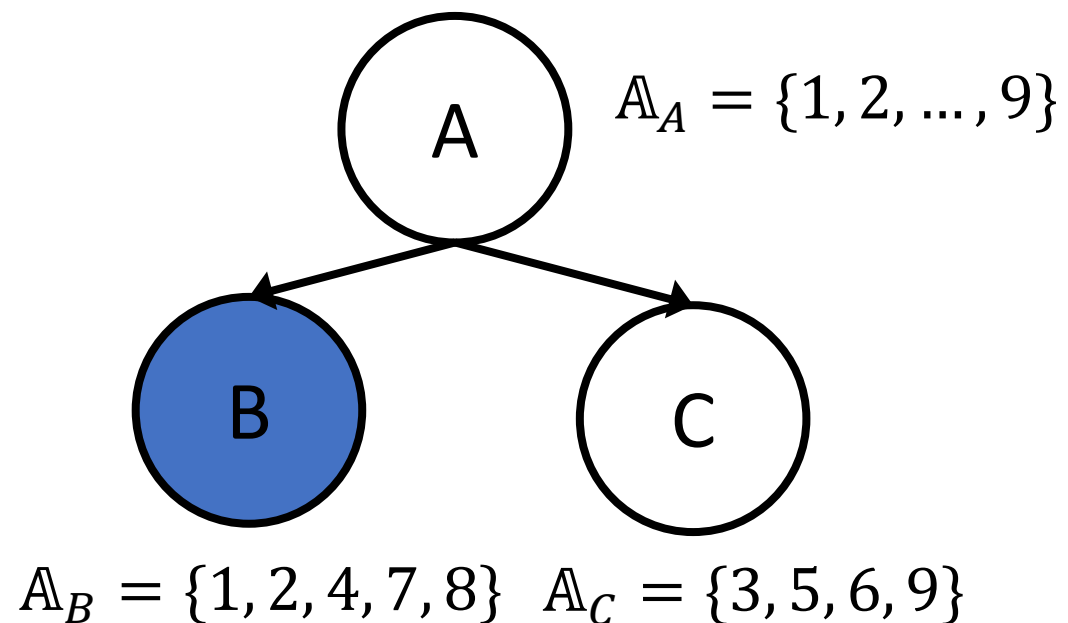
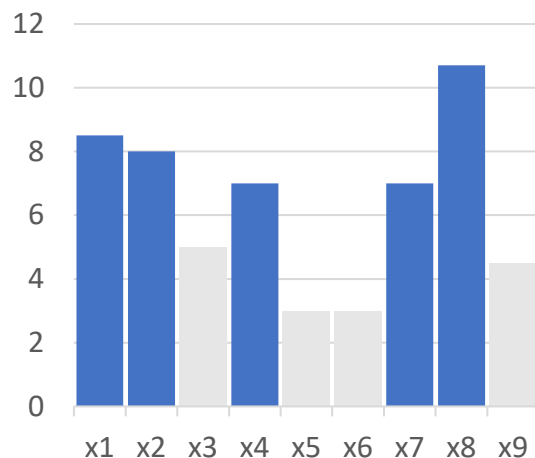
---

Select node  $B$  based on UCB value



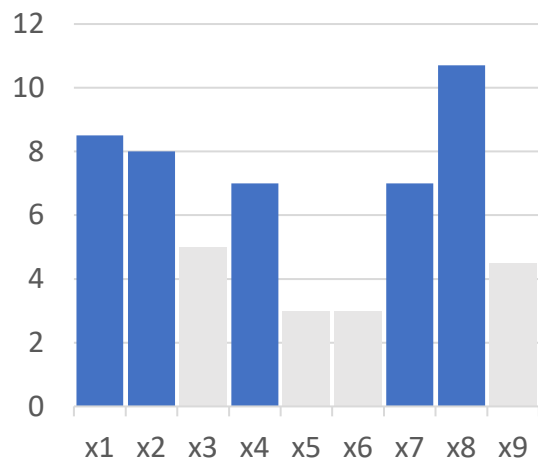
## Example Illustration of MCTS-VS

Variable score  $s$

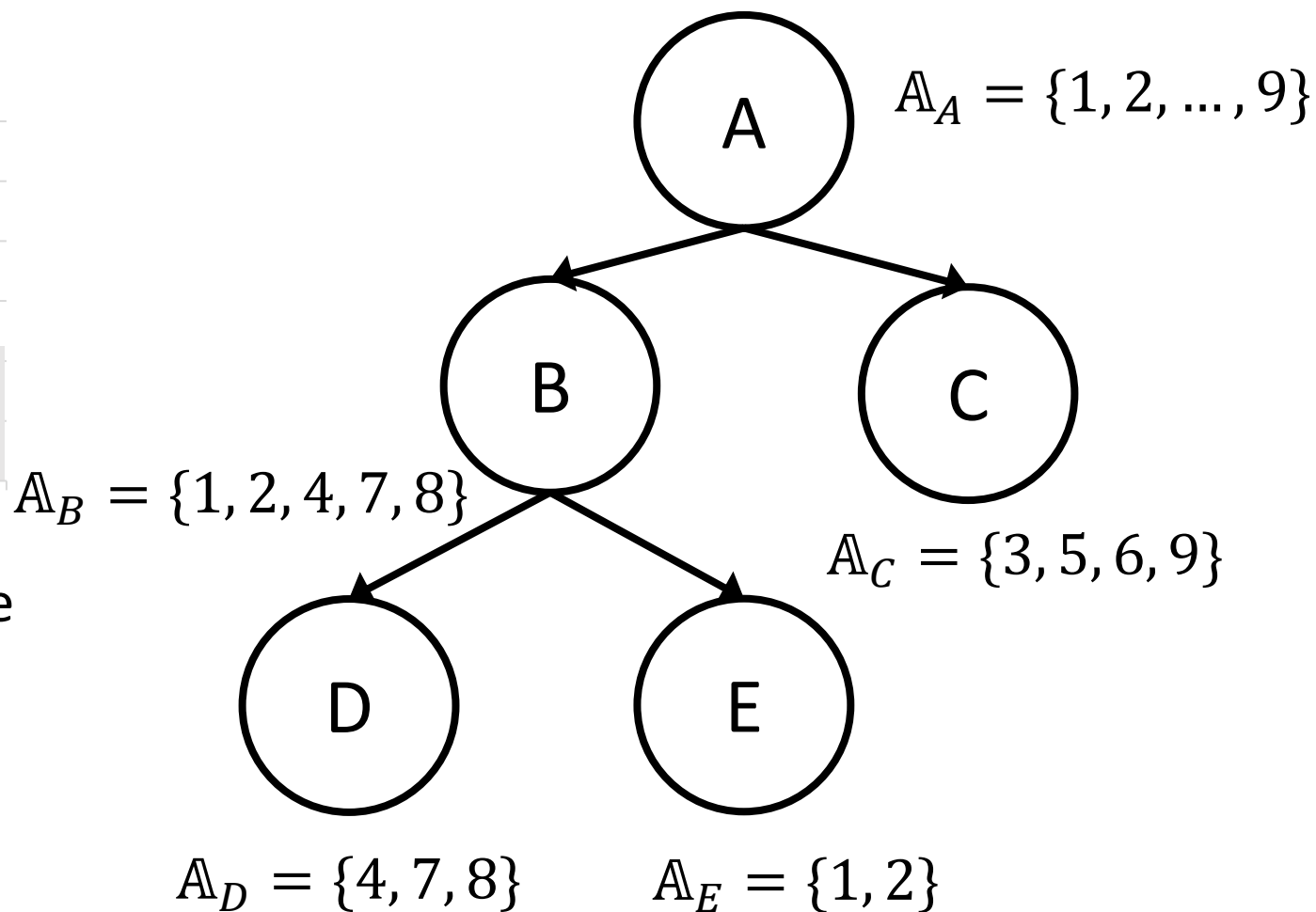
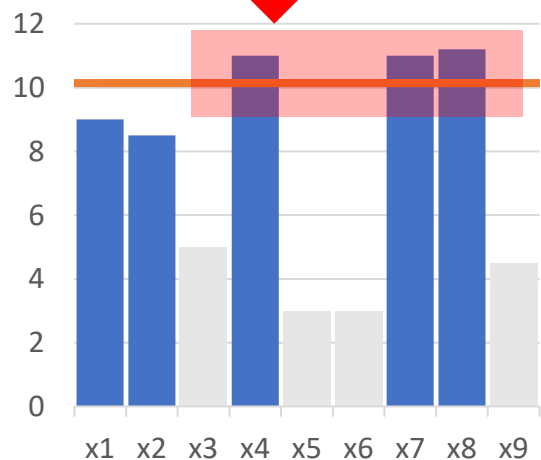


## Example Illustration of MCTS-VS

Variable score  $s$

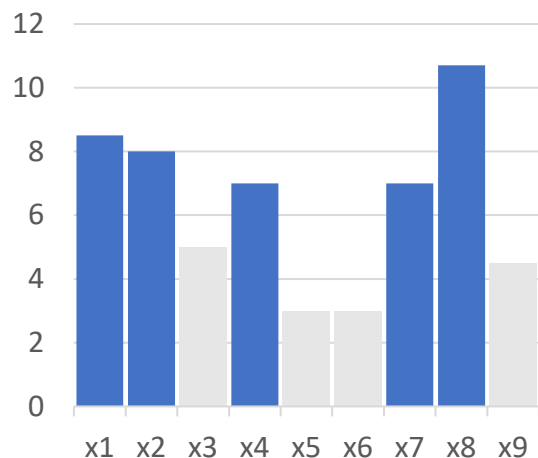


optimize

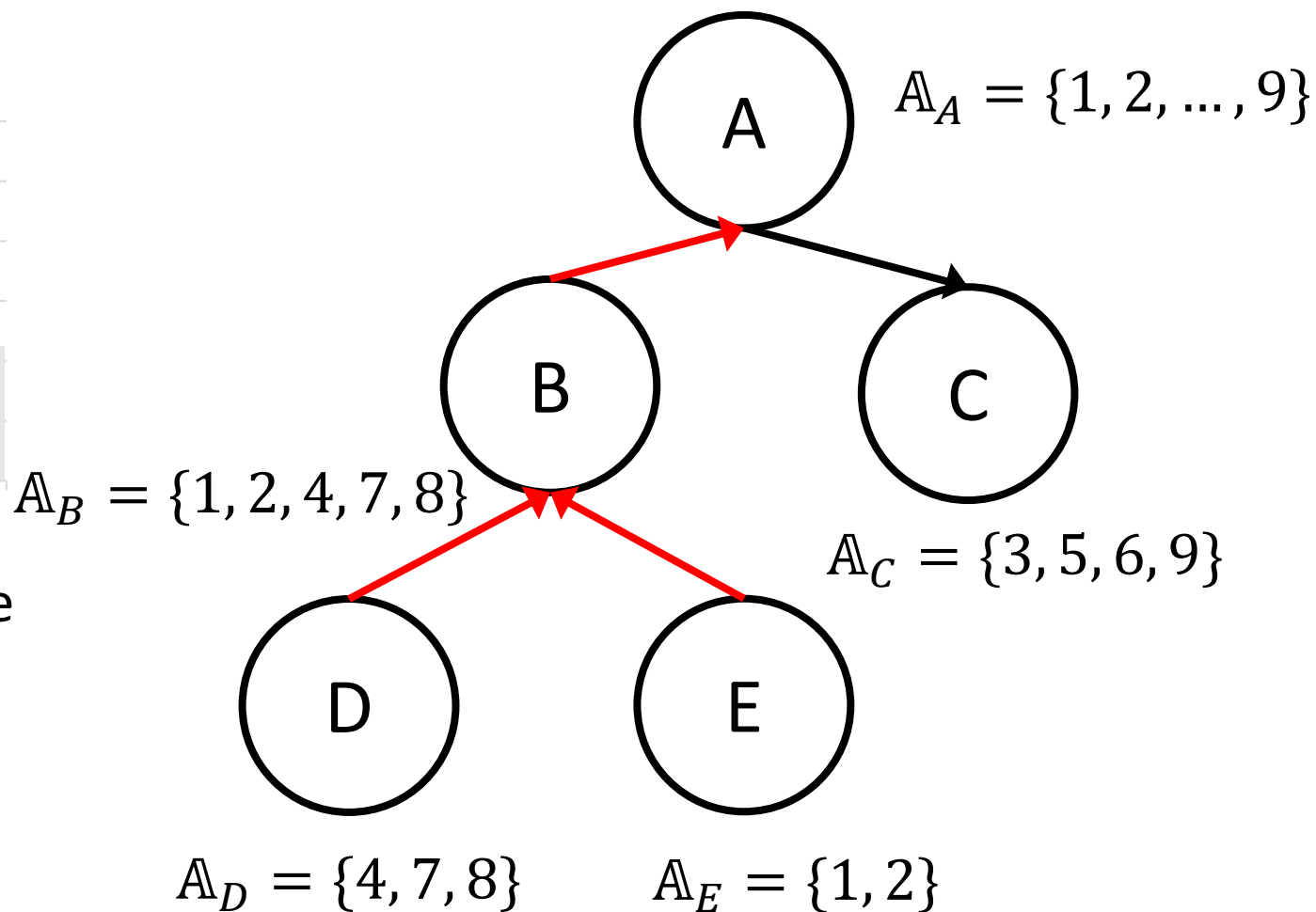
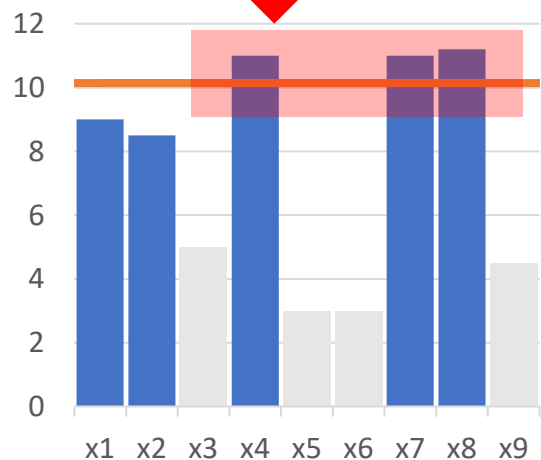


## Example Illustration of MCTS-VS

Variable score  $s$



optimize





## “Fill-in” Strategy

---

The best  $k$  samples:  $\{(\mathbf{x}^{*j}, y^{*j})\}_{j=1}^k$ , for unselected variable  $x_i$ :

- The best- $k$  strategy:  $x_i$  is uniformly selected from  $\{x_i^{*j}\}_{j=1}^k$  at random
- The average best- $k$  strategy:  $x_i$  is the average of  $\{x_i^{*j}\}_{j=1}^k$
- The random strategy:  $x_i$  is sampled from the domain randomly

# Theoretical Analysis

---

- Cumulative regret  $R_T = \sum_{t=1}^T (f(x^*) - f(x^t))$
- Assumption: The function  $f$  is a GP sample path. For some  $a, b > 0$ , given  $L > 0$ , the partial derivatives of  $f$  satisfy that  $\forall i \in [D], \exists \alpha_i \geq 0$ ,

$$P \left( \sup_{x \in \mathcal{X}} \left| \frac{\partial f}{\partial x_i} \right| < \alpha_i L \right) \geq 1 - a e^{-\left(\frac{L}{b}\right)^2}$$

- Theorem:  $\forall \delta \in (0, 1)$ , let  $\beta_t = 2 \log \left( \frac{4\pi_t}{\delta} \right) + 2d_t \log(d_t t^2 b r \sqrt{\log(\frac{4Da}{\delta})})$  and  $L = b \sqrt{\log \frac{4Da}{\delta}}$ , and  $\{\pi_t\}_{t \geq 1}$  satisfies  $\sum_{t \geq 1} \pi_t^{-1} = 1$  and  $\pi_t > 0$ . Let  $\beta_T^* = \max_{1 \leq i \leq T} \beta_t$ . At iteration  $T$ ,

$$R_T \leq \sqrt{C_1 T \beta_T^* \gamma_T} + 2\alpha_{max} + 2 \sum_{t=1}^T \sum_{i \in [D] \setminus \mathbb{M}_t} \alpha_i^* L r$$

# Theoretical Analysis

---

$$R_T \leq \sqrt{C_1 T \beta_T^* \gamma_T} + 2\alpha_{max} + 2 \sum_{t=1}^T \sum_{i \in [D] \setminus \mathbb{M}_t} \alpha_i^* Lr$$

- The regret from optimization
- The regret from unselected variables

Insight:

- The variable selection can reduce the computational complexity while increasing the regret
- A good variable selection algorithm should select as important variables as possible, i.e., variables with as larger  $\alpha_i^*$  as possible

# Experiments

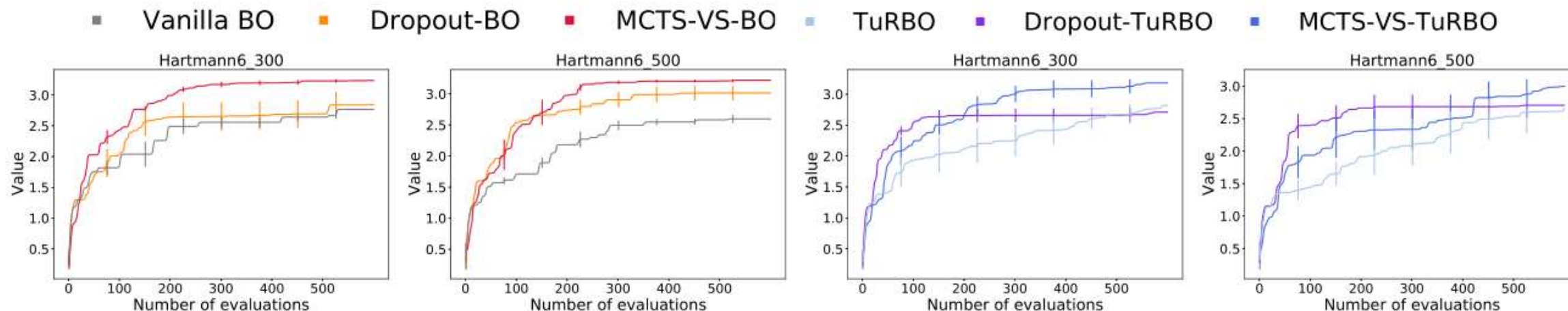
---

We want to know the following research questions (RQs):

- RQ1: Can BO benefit from variable selection?
- RQ2: How does MCTS-VS perform compared with state-of-the-art methods?
- RQ3: How about the runtime of MCTS-VS?
- RQ4: Can MCTS-VS select more important variables than Dropout? (why)
- RQ5: Is MCTS-VS sensitive to the hyper-parameters?

## Experiments – RQ1

Effectiveness of variable selection:

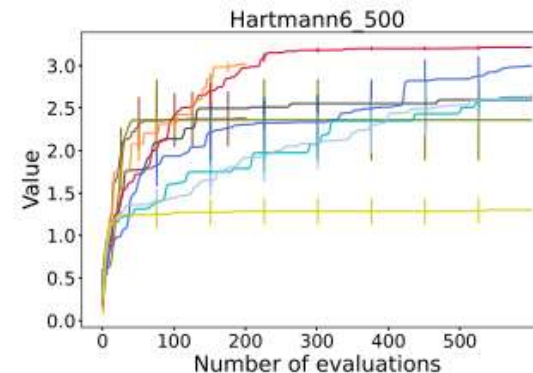
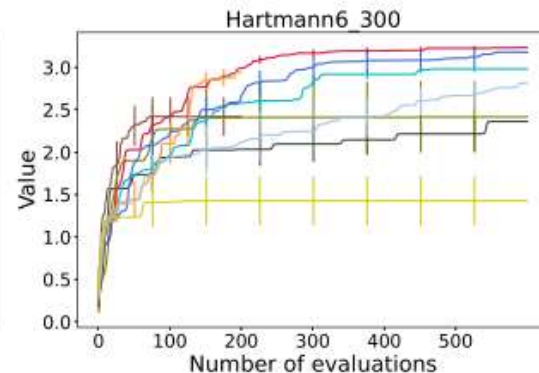
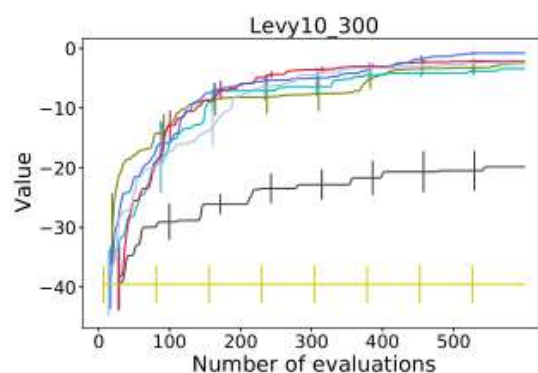
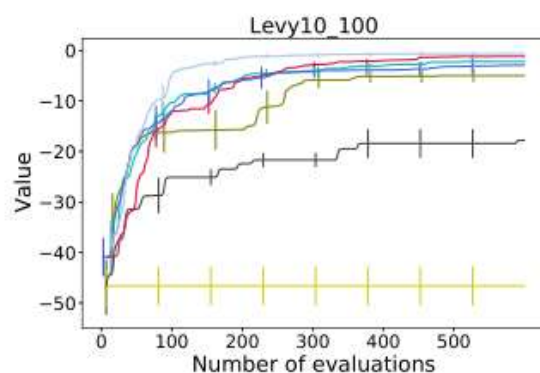


- Dropout is better than BO without variable selection
- MCTS-VS is better than Dropout

# Experiments – RQ2

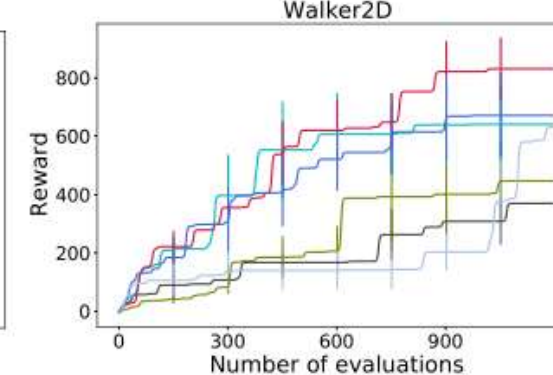
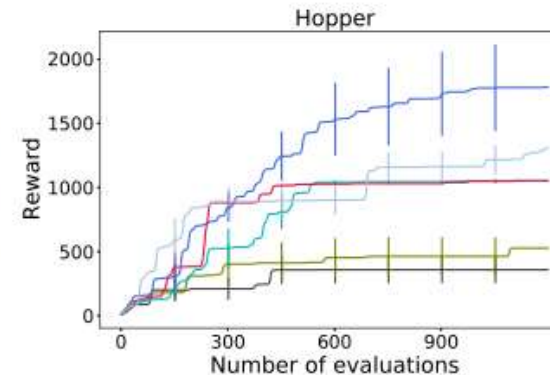
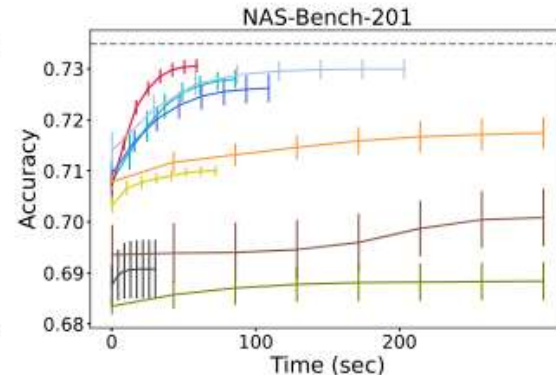
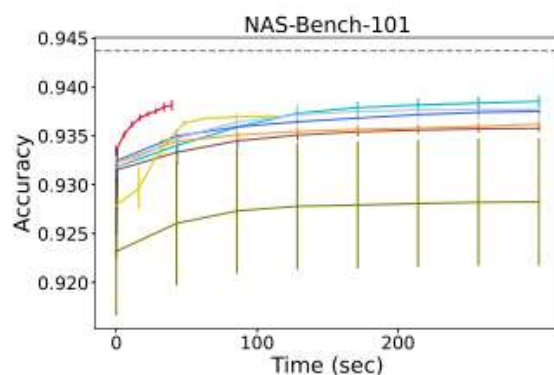
## Synthetic functions:

■ MCTS-VS-BO    ■ MCTS-VS-TuRBO    ■ TuRBO    ■ LA-MCTS-TuRBO  
■ SAASBO    ■ HeSBO    ■ ALEBO    ■ CMA-ES    ■ VAE-BO



## Real-world problems:

MCTS-VS is comparable with sota methods



## Experiments – RQ3

### Runtime comparison:

METHOD	LEVY10_100	LEVY10_300	HARTMANN6_300	HARTMANN6_500
VANILLA BO	3.190	4.140	4.844	5.540
DROPOUT-BO	2.707	3.225	3.237	3.685
MCTS-VS-BO	2.683	3.753	3.711	4.590
TURBO	8.621	9.206	9.201	9.754
LA-MCTS-TURBO	14.431	22.165	25.853	34.381
MCTS-VS-TURBO	4.912	5.616	5.613	5.893
SAASBO	/	/	2185.678	4163.121
HESBO	220.459	185.092	51.678	55.699
ALEBO	/	/	470.714	512.641
CMA-ES	0.030	0.043	0.043	0.045

Variable selection can reduce the runtime



## Experiments – RQ4

---

Recall comparison:

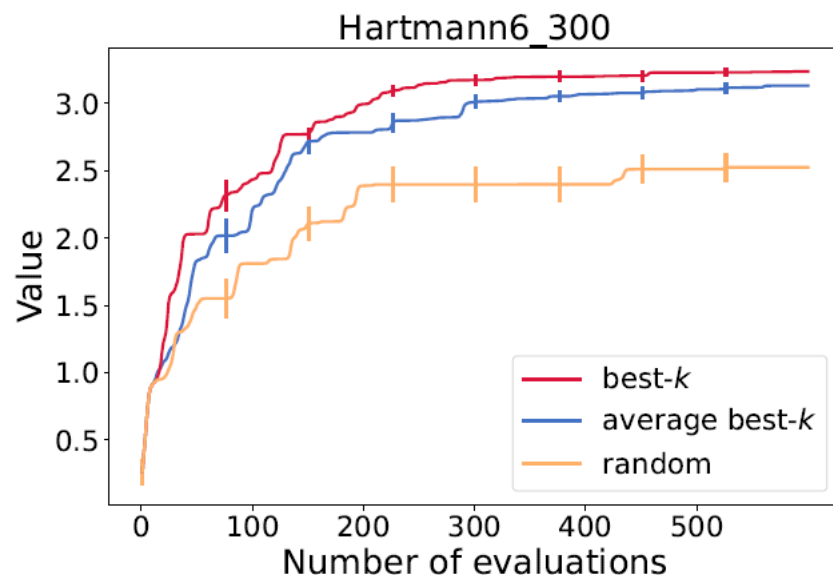
METHOD	LEVY10_100	LEVY10_300	HARTMANN6_300	HARTMANN6_500
DROPOUT	0.100	0.030	0.020	0.012
MCTS-VS	0.429	0.433	0.352	0.350

Recall  $\frac{d_t^*}{d}$  is used to compare the quality of variable selection, where  $d_t^*$  is the number of valid variables selected at iteration  $t$  and  $d$  is the number of valid variables

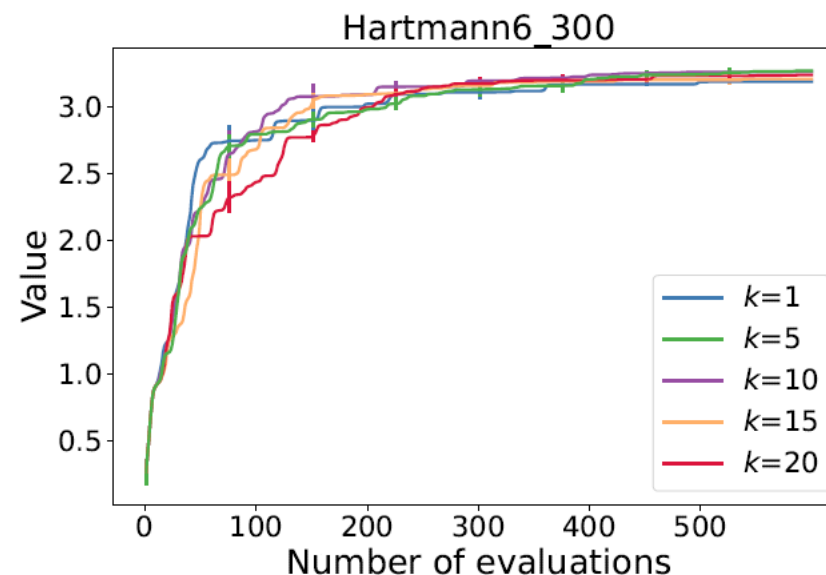
- Dropout:  $\frac{d}{D}$  in expectation
- MCTS-VS: run for 600 evaluations on five different random seeds and calculate the average recall

The recall of MCTS-VS is much larger than Dropout

## Experiments – RQ5



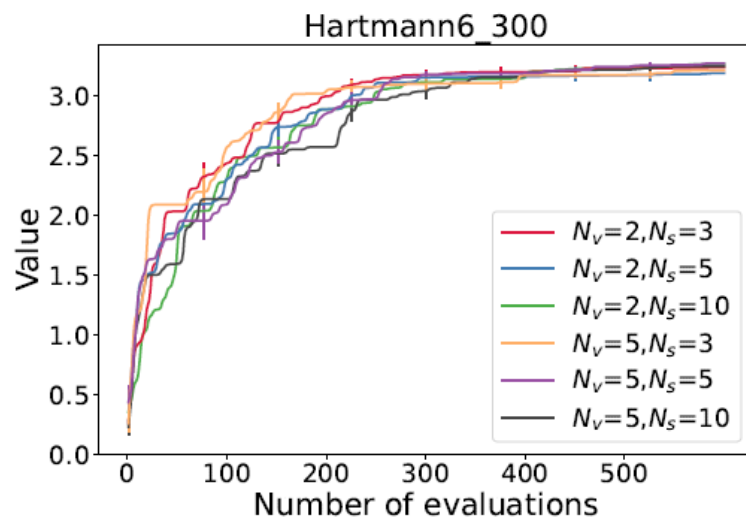
(a) “Fill-in” strategy



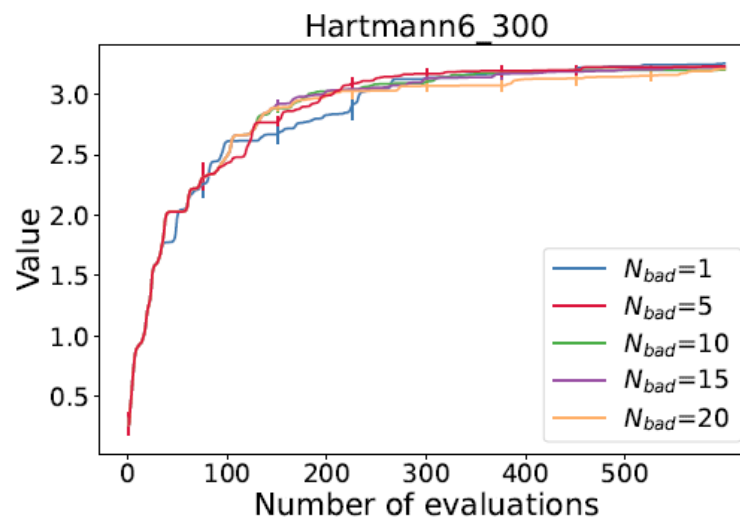
(b) Hyper-parameter  $k$  of the best- $k$  strategy

The best- $k$  strategy is good, and MCTS-VS is **not sensitive** to the selection of  $k$

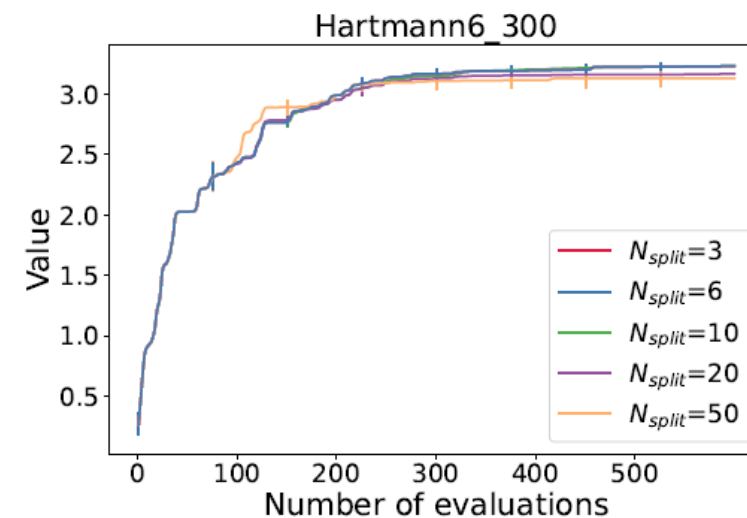
## Experiments – RQ5



(a) Number of samples



(b)  $N_{bad}$



(c)  $N_{split}$

MCTS-VS is **not sensitive** to other hyper-parameters

## Conclusion

---

MCTS-VS uses MCTS to recursively partition the variables into important and unimportant ones, and only optimizes those important variables

Feature work:

- A more well-designed metric for importance
- A specific theoretical analysis for MCTS-VS

*Thank you!*