

Combinatorial Bayesian Optimization using the Graph Cartesian Product

1. Introduction

- Black-box optimization:

$$x_{opt} \argmin_{x \in X} f(x)$$

- Non-differentiable f
- Expensive to query f
- Noisy f

- Typically, the search space X is continuous and the function f is assumed to be smooth.

- A more challenging problem is when the search space is **combinatorial** (e.g., variables are categorical or ordinal).

- There are two main challenges for combinatorial problems:

- How to define a smooth function on combinatorial objects?
→ **kernel (prior on smoothness)**
- How to efficiently select next points in a combinatorial space?
→ **acquisition function optimization**

2. Smoothness

- A space of combinatorial variables -

$$C_1, C_2, \dots, C_{d-1}, C_d$$

↓

- A graph corresponding to the space -

$$G_i = G(C_i) \text{ for } i = 1, \dots, d$$

$$G = G_1 \boxtimes G_2 \boxtimes \dots \boxtimes G_{d-1} \boxtimes G_d$$

↓

- The concept of smoothness on functions on a graph -

Graph Fourier Analysis using graph Laplacian $L(G)$

Eigendecomposition of $\{(\lambda_i, u_i)\}$

Smaller $\lambda \Rightarrow$ Smoother u

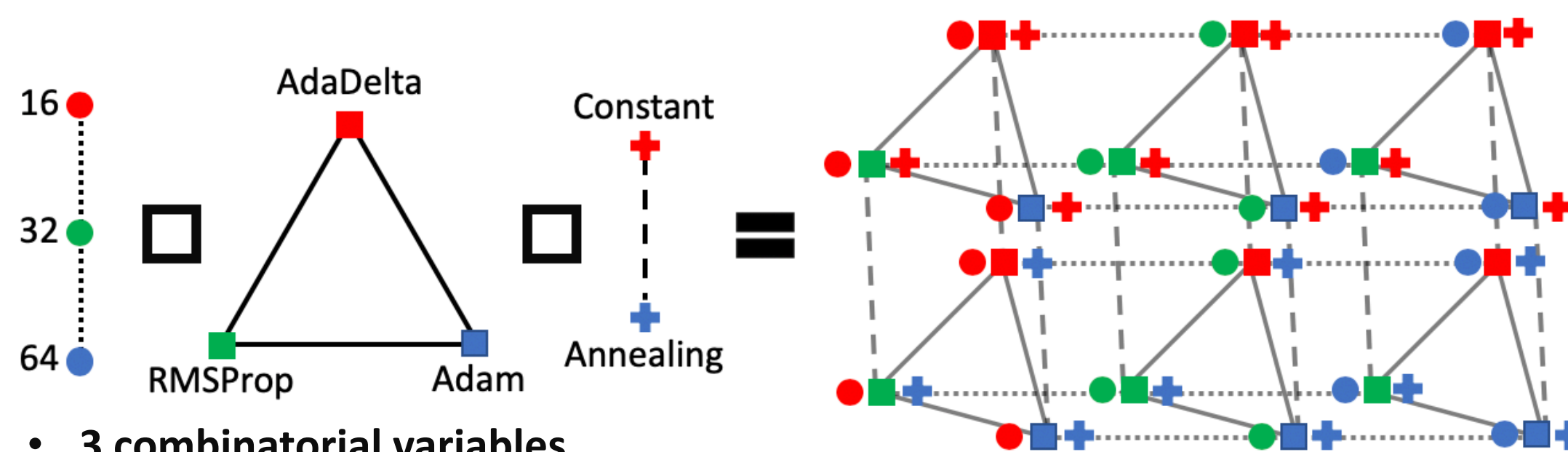
↓

- Smoothness of functions on combinatorial variables -

A kernel on a space of combinatorial variables

$$K((c_1, \dots, c_d), (c_1', \dots, c_d'))$$

3. Combinatorial graph (Search space)



- 3 combinatorial variables

- Batch size $C_1 = \{16, 32, 64\}$
- Optimizer $C_2 = \{AdaDelta, RMSProp, Adam\}$
- Learning rate annealing $C_3 = \{Constant, Annealing\}$

- 3 subgraphs

$$G_1 = G(C_1), G_2 = G(C_2), G_3 = G(C_3)$$

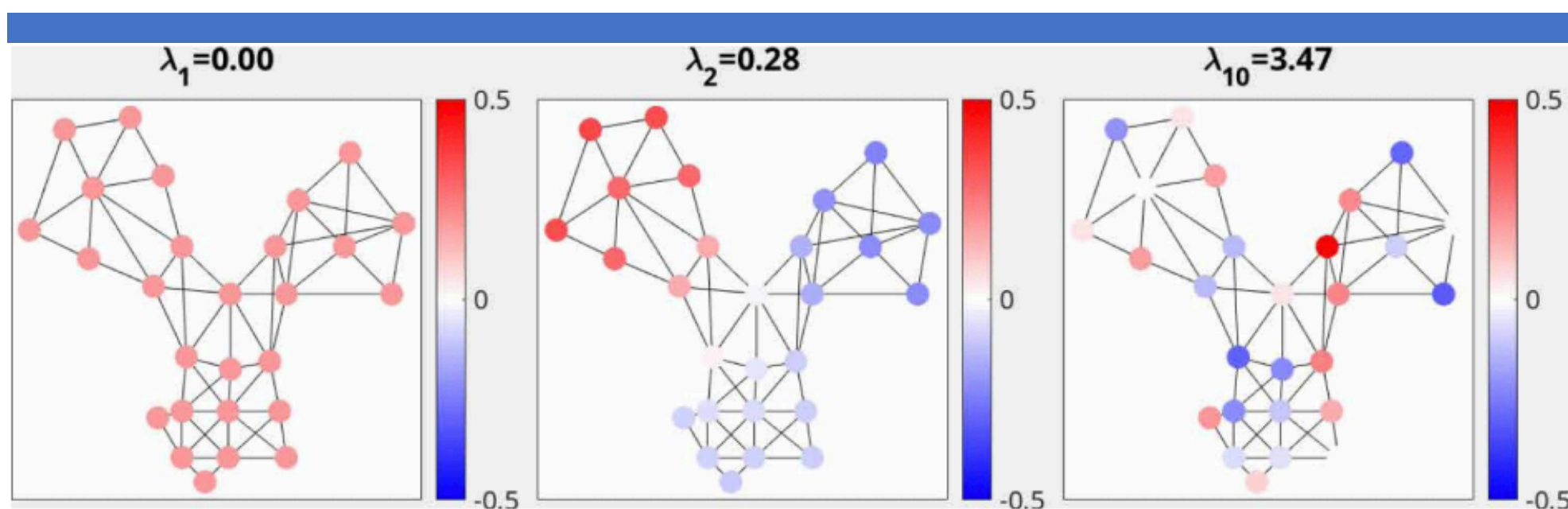
- Complete graphs for categorical variables
- Path graphs for ordinal variables

- Combinatorial graph : Graph Cartesian product of subgraphs

$$G = G_1 \boxtimes G_2 \boxtimes G_3$$

- Vertices : sets of specific choices of categorical/ordinal variables.
- Edges : similarities between 2 sets of choices

4. Graph Signal Processing



Graph Signal Processing: Overview, Challenges, and Applications, Ortega et. al., IEEE.

- Graph Fourier Transform

- Graph Laplacian $L(G) = D - A$, D : degree mat., A : adjacency mat.
- Eigendecomposition of graph Laplacian $L(G) : \{(\lambda_i, u_i)\}_{i=1, \dots, |V|}$
- λ_i represents smoothness(energy) of u_i

- Approximate a function with eigenfunctions with small eigenvalues.

→ $f \approx \sum c_i u_i$ while trying to make c_i small if λ_i is large

- Diffusion kernel → GP **nonparametric** approach

$$K_G(v, \tilde{v} | \beta) = [e^{-\beta L(G)}]_{v, \tilde{v}} = [U e^{-\beta \Lambda} U^T]_{v, \tilde{v}}$$

5. Graph Cartesian Product

- Graph Cartesian product and Kronecker sum

$$\begin{aligned} L(G_1 \boxtimes G_2) &= L(G_1) \oplus L(G_2) \\ &= L(G_1) \otimes I_2 + I_1 \otimes L(G_2) \end{aligned}$$

- Kronecker sum and matrix exponential

$$\begin{aligned} K_{G_1 \boxtimes G_2} &= e^{-\beta L(G_1 \boxtimes G_2)} = e^{-\beta (L(G_1) \oplus L(G_2))} \\ &= e^{-\beta (L(G_1) \otimes I_2 + I_1 \otimes L(G_2))} \\ &= e^{-\beta L(G_1)} \otimes e^{-\beta L(G_2)} \\ &= K_{G_1} \otimes K_{G_2} \end{aligned}$$

- Efficient computation of diffusion kernels

$$O(\prod_{i=1}^d |V_i|^3) \rightarrow O(\sum_{i=1}^d |V_i|^3)$$

- Able to handle large graphs

$$|V| \in \{2^{24}, 5^{21}, 2^{28}, 2^{43}, 2^{60}, 2^{32}\}$$

6. ARD Diffusion Kernel

- Diffusion kernel has a single kernel parameter β

$$K_G(v, \tilde{v} | \beta) = \bigotimes_i K_{G_i}(v_i, \tilde{v}_i | \beta_i)$$

- A multiplicand in the Kronecker product corresponds to a sub-graph

- Each sub-graph corresponds to a variable

- Variable-wise kernel parameter → **ARD diffusion kernel**

$$K_G(v, \tilde{v} | \beta) = \bigotimes_i K_{G_i}(v_i, \tilde{v}_i | \beta_i)$$

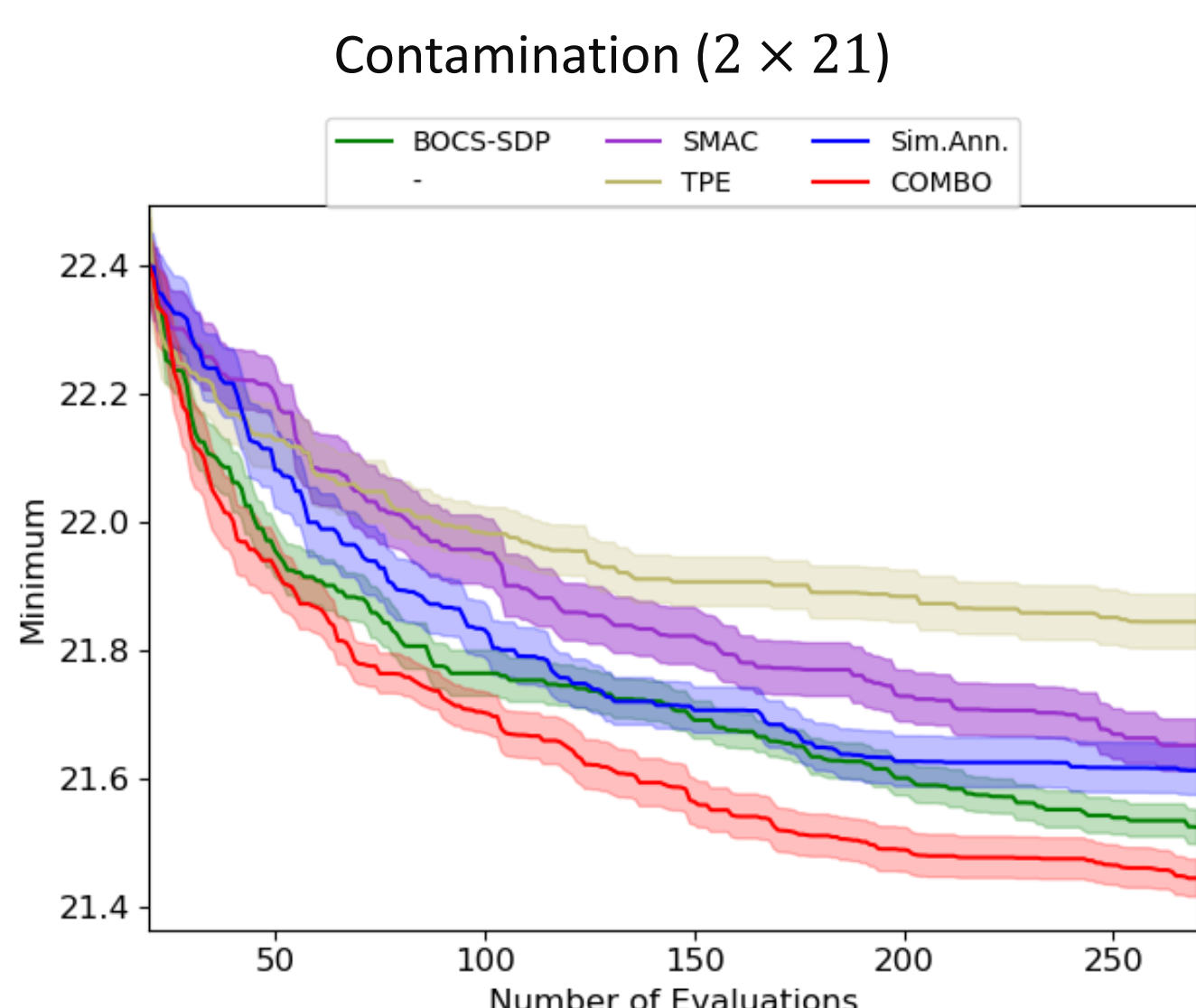
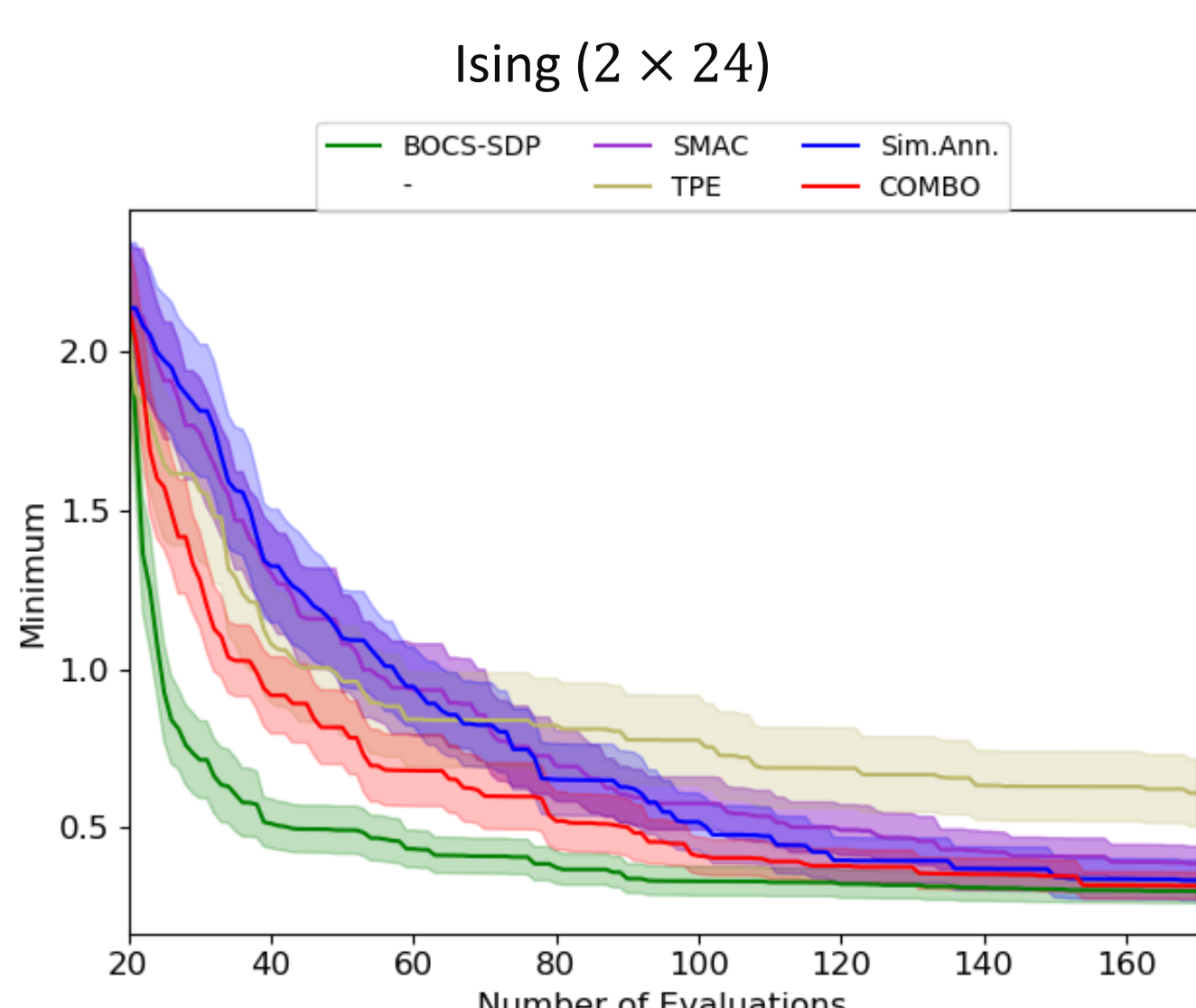
- Relevant variables can be selected automatically

- Horseshoe priors on $\{\beta_i\}_{i=1, \dots, d}$ promotes more effective feature selection

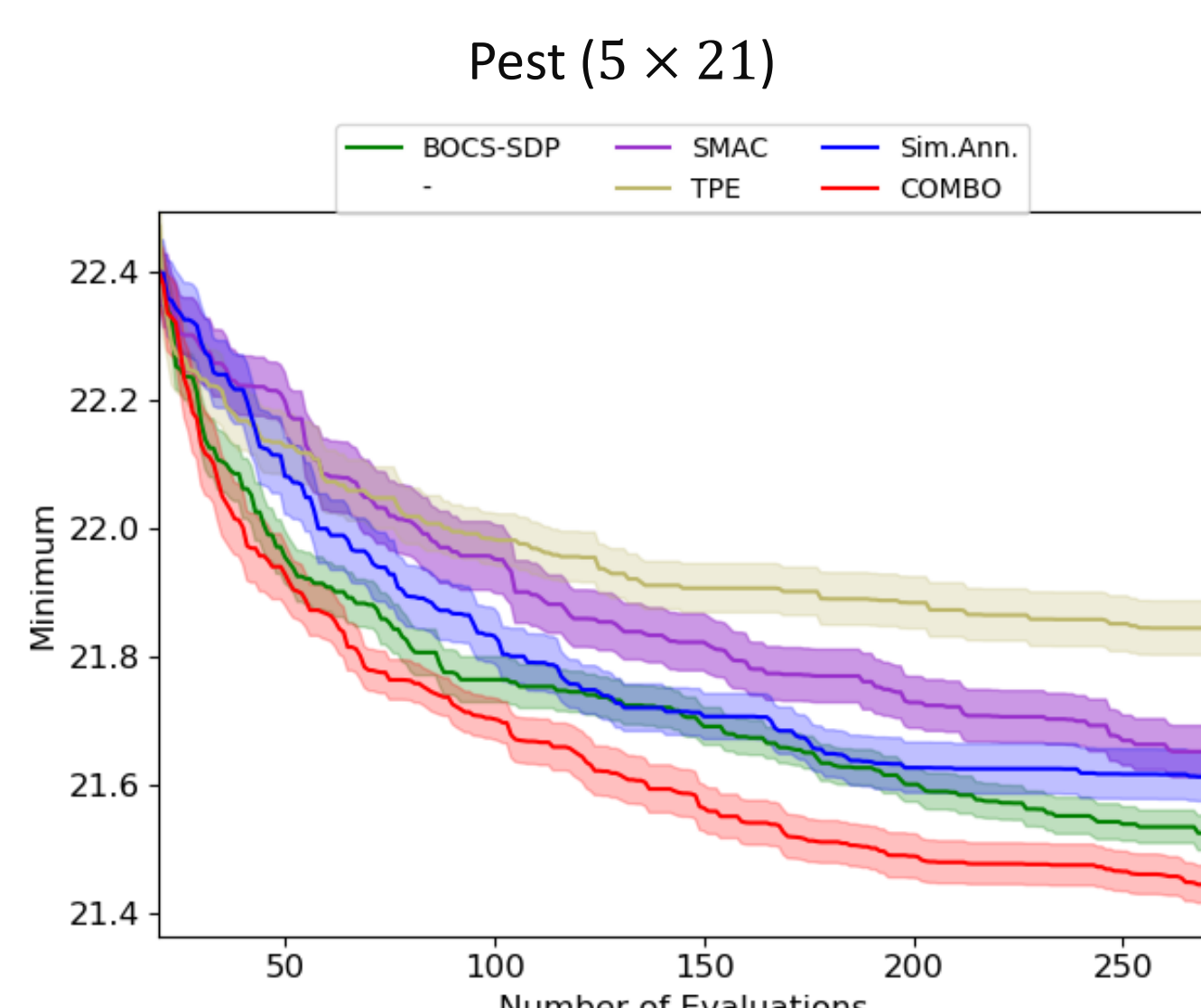
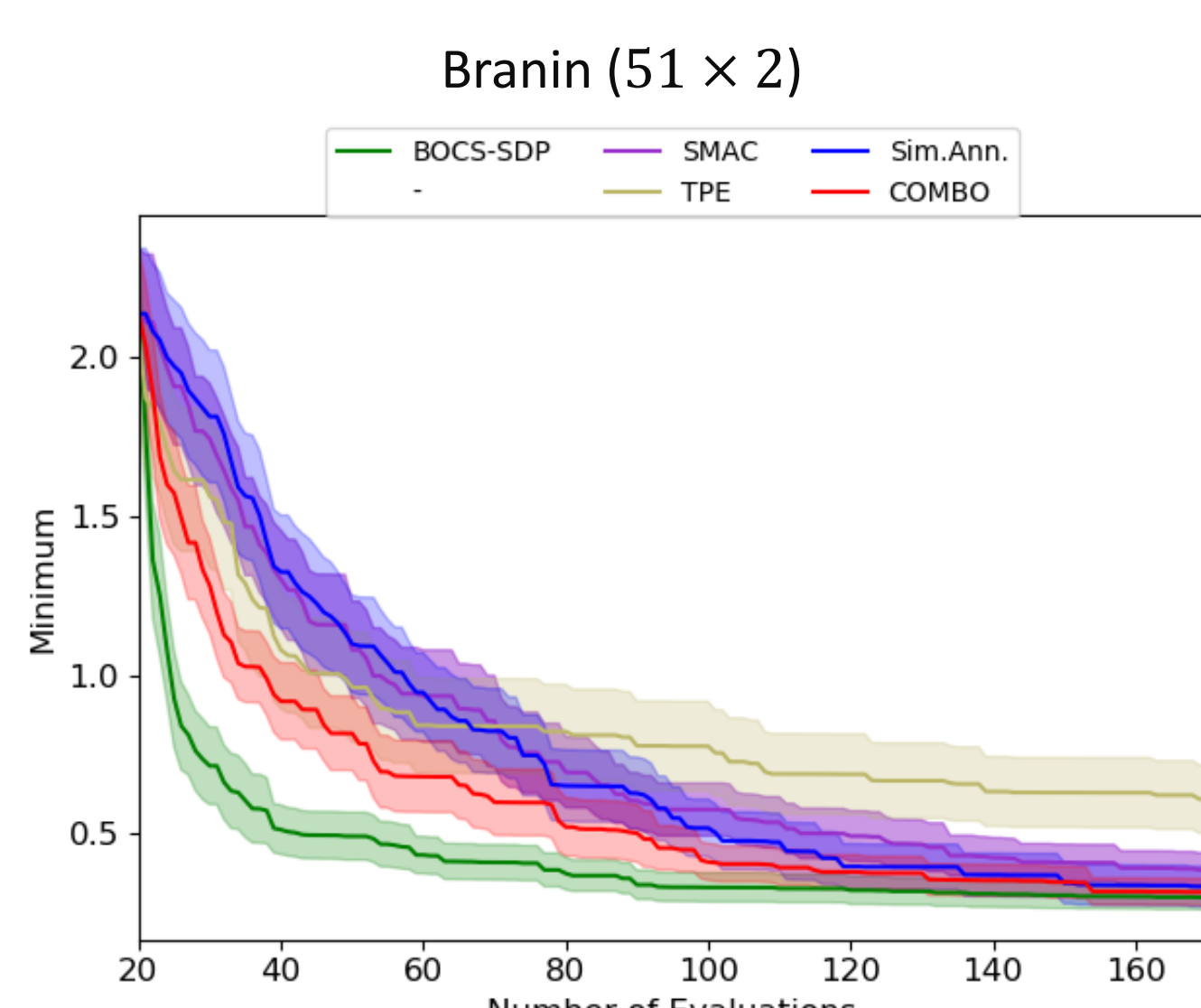
- We use slide sampling to sample $\{\beta_i\}_{i=1, \dots, d}$.

7. Experiments

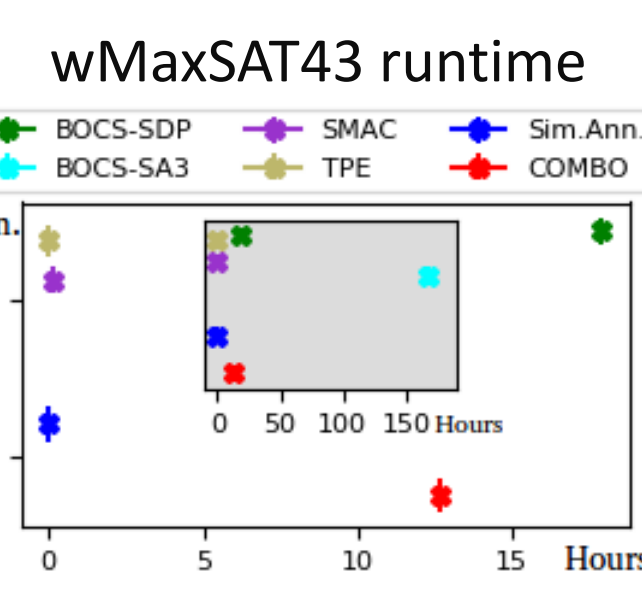
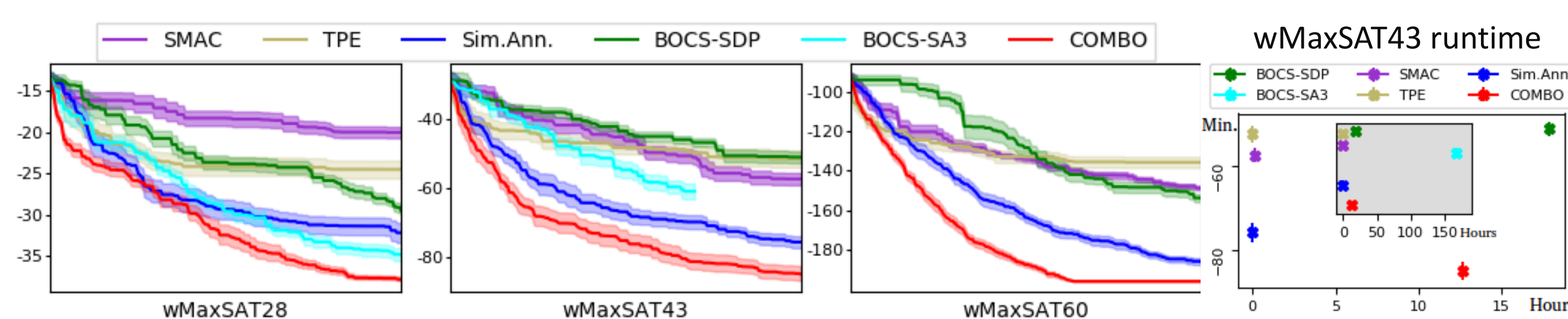
Binary



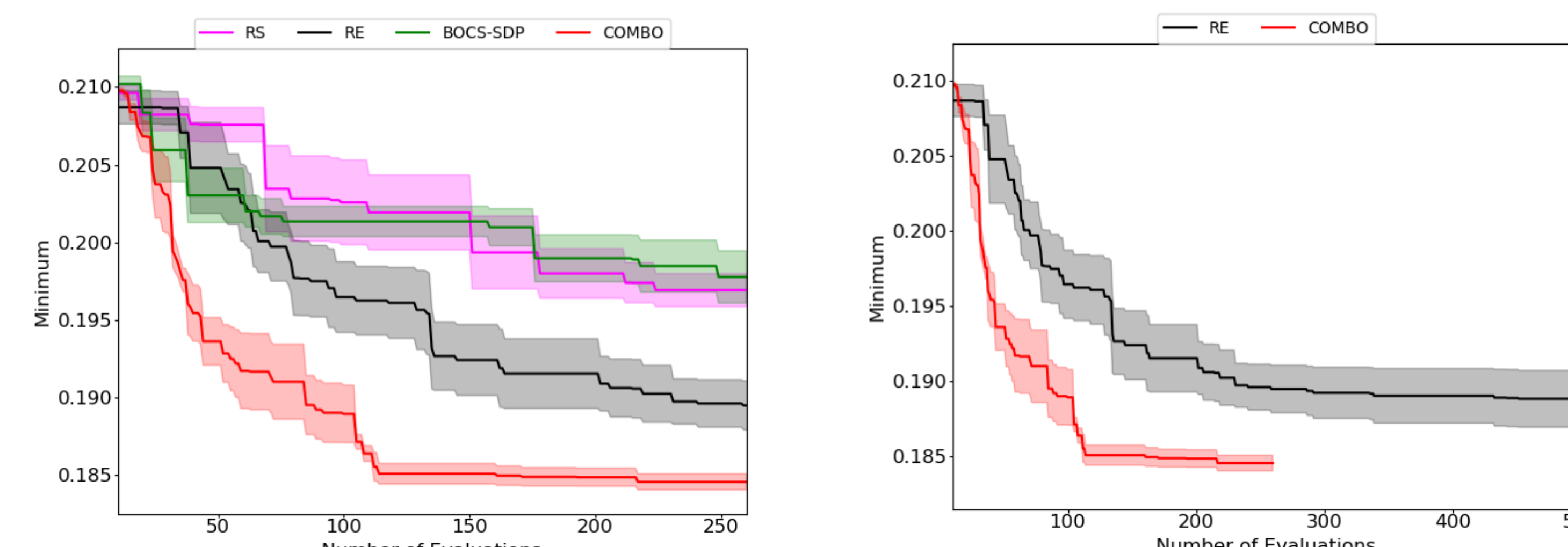
Ordinal & Multi-cate.



Weighted MaxSAT



Architecture Search



8. Conclusions

- We propose COMBO, a Bayesian Optimization for combinatorial search spaces using Gaussian Processes.

- The application of the graph Cartesian product allows to reduce exponential complexity to a linear complexity.

- The ARD diffusion kernel allows to better model complex functions on combinatorial objects.

- We show supremacy of COMBO on various combinatorial optimization problems.