LEARNING REGION OF INTEREST FOR BAYESIAN OPTIMIZATION WITH ADAPTIVE LEVEL-SET ESTIMATION

Fengxue Zhang¹, Jialin Song², James Bowden⁴, Alexander Ladd³, Yisong Yue⁴, Thomas A. Desautels³, Yuxin Chen¹

¹University of Chicago, ²NVIDIA, Inc, ³Lawrence Livermore National Laboratory, ⁴California Institute of Technology



Motivating Applications Bayesian Optimization on Partition(s) Are you having trouble optimizing your scientific experiments? Is it because your objective is *too complex*? Heuristics including Train GP on the partitions Partitions Partition the Gaussian clustering, Monte-Carlo tree $f_i: (X'_i) \to Y'_i$ Is the global model failing to capture useful locality? search, trust region Search Space Process(es) Where (X'_i, Y'_i) are observed pairs. identification etc. PELTON WHEEL OSCILLATOR & GENERATOR Local Acquisition Function(s) New observation $acq_i(x)$ (x_t, y_t) SEA RETURN. Select candidate to query SEA WATER Global $argmax acq_i(x)$ PISTON Optimization AG FLOW Figure 3: Bayesian Optimization on Partition(s)

Figure 1: Left: Protein-protein Interactions; Right: Water Converter

Previous methods rely on heuristics that normally introduce *additional complexity* to the surrogate models, which could incur challenges in fine-tuning the *hyperparameters* of the heuristics.

Deep Kernel Learning

Algorithmic Details



Figure 2: Deep Kernel Learning (DKL)

- A *latent space mapping* g to project input \mathcal{X} to a latent space \mathcal{Z}
- An objective mapping $h : \mathcal{Z} \to \mathbb{R}$ such that $f(x) \approx h(g(x))$
- $h: \mathcal{Z} \to \mathbb{R}$ often modeled by *Gaussian Processes*
- [ZNC22] initialize the latent space mapping with the encoder of the pre-train auto encoder
- [TDHL20] periodically retrain the auto encoder to improve the latent space.

Global Estimation with DKL

- Learn the global partitioning on **X** with Deep Kernel Learning (DKL) [Wil+16] due to its scalability
- Underlying global function $f_g \coloneqq f$ is assumed to be drawn from a global Gaussian process.
- Maximizing the negative log-likelihood (NLL) $-\log(P(\mathbf{y}_t|X_t, \theta_{f_a,t}))$ [RW05].
- Sample unlabeled dataset from X to pre-trained an Auto-Encoder and use the parameters of its encoder to initialize the neural network q [FCT20].

Regions of Interest Filtering

Algorithm 1 Bayesian Optimization with Adaptive Level-Set Estimation (BALLET) 1: Input:Search space X, initial observation D_0 , horizon T; 2: for t = 1 to T do Learn the global estimation $\mathcal{GP}_{f_g,t}: \theta_{f_g,t} \leftarrow \arg \max_{\theta_{f_g}} - \log(P(\mathbf{y}_t | X_{t-1}, \theta_{f_g}))$ 3: Partition by region of interest filtering: $\hat{\mathbf{X}}_t \leftarrow \{\mathbf{x} \in \mathbf{X} | UCB_{f_g,t}(\mathbf{x}) > LCB_{f_g,t,max}\}$ 4: Partition the historical observation: $\hat{\mathbf{D}} = \{(\mathbf{x}, y) \in \mathbf{D} | \mathbf{x} \in \hat{\mathbf{X}}_t\}.$ 5: Learn the superlevel-set GP: $\mathcal{GP}_{\hat{f},t}: \theta_{\hat{f},t} \leftarrow \arg \max_{\theta_{\hat{f}}} -\log(P(\mathbf{y}_t | \hat{\mathbf{X}}_t, \theta_{\hat{f}}))$ 6: Optimize the superlevel-set acquisition function: $\mathbf{x}_{t+1} \leftarrow \arg \max \alpha_{\hat{f}}(\mathbf{x})$ 7: Update D: $\mathbf{D}_{t+1} \leftarrow \mathbf{D}_t \cup \{(\mathbf{x}_{t+1}, y_{t+1})\}$ 8:

9: **Output**: $\max y_t$

In contrast to existing work, we consider a *non-parametric model* for partitioning the search space, which has shown remarkable performance in real-world tasks while having *few hyperparameters* to maintain.

Bayesian Optimization with Adaptive Level-Set Estimation



With the confidence interval of the global \mathcal{GP}_{f_g} , we could define the upper confidence bound $UCB_{f_q}(\mathbf{x})$ and lower confidence bound $LCB_{f_a}(\mathbf{x})$ We attain the region of interest.

$$\hat{\mathbf{X}} = \{ x \in \mathbf{X} | UCB_{f_g}(\mathbf{x}) > LCB_{f_g,max} \}$$
(1)

- The global \mathcal{GP}_{f_q} enables a filtering on **X** to locate the region of interest $\hat{\mathbf{X}}$.
- It is desired for $\hat{\mathbf{X}}$ that with high probability, $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x})$ is contained by $\hat{\mathbf{X}}$,
- Also $|\hat{\mathbf{X}}| \ll |\mathbf{X}|$, then \hat{f} , the objective function defined on $\hat{\mathbf{X}}$ is therefore of reduced complexity.

The algorithm filters the search space \mathbf{X} using the UCB and LCB estimated by a *global GP*. Then it feeds another GP the *regions of interest* of the search space $\hat{\mathbf{X}}$ and the filtered historical observations $(\hat{\mathbf{X}}_t, \hat{\mathbf{Y}}_t)$. Global optimization at each iteration is conducted on the *adaptive partition* using this GP and its acquisition function.

