Information-Theoretic Bayesian Optimization for Bilevel Optimization Problems

Anonymous authors Paper under double-blind review

ABSTRACT

A bilevel optimization problem consists of two optimization problems nested as an upper- and a lower-level problem, in which the optimality of the lower-level problem defines a constraint for the upper-level problem. This paper considers Bayesian optimization (BO) for the case that both the upper- and lower-levels involve expensive black-box functions. Because of its nested structure, bilevel optimization has a complex problem definition and, compared with other standard extensions of BO such as multi-objective or constraint settings, it has not been widely studied. We propose an information-theoretic approach that considers the information gain of both the upper- and lower-optimal solutions and values. This enables us to define a unified criterion that measures the benefit for both level problems, simultaneously. Further, we also show a practical lower bound based approach to evaluating the information gain. We empirically demonstrate the effectiveness of our proposed method through several benchmark datasets.

1 Introduction

000

001

002 003 004

006

007 008 009

010 011

012

013

014

016

017

018

019

021

023

024

026027028

029

031

033

034

037

040

041

042

043

044

046

047

048

049

050 051

052

The bilevel optimization is a standard formulation for a decision making problem that has a hierarchical structure. It consists of two optimization problems nested as an upper- and a lower-level problem, in which the optimality of the lower-level problem defines a constraint for the upper-level problem. For example, in the computational materials design, a target property should be optimized under the constraint of the energy minimization. Bilevel optimization techniques is applicable to hierarchical decision makings in a variety of contexts such as inverse optimal control (Suryan et al., 2016), chemical reaction optimization (Abbassi et al., 2021), and shape optimization (Herskovits et al., 2000).

We particularly focus on the case both level problems are defined by expensive black-box functions. Most of BO studies for bilevel optimization consider applying BO only to the upper-level problem (e.g., Kieffer et al., 2017; Dogan & Prestwich, 2023) as pointed out by (Chew et al., 2025). Typically, under a selected query for the upper-level problem, repeated queries to the lower-level problem is required, and further, the gradient of the lower-level problem is often assumed (e.g., Fu et al., 2024). Islam et al. (2018) and Wang et al. (2021) consider BO in both levels, but repeated queries on lower-level is still required. These approaches are not fully suitable when both levels are expensive black-boxes in which the gradient is not available. On the other hand, recently a few methods without those limitations have also been studied. Ekmekcioglu et al. (2024) combine the Thompson sampling on the upper-level query and a knowledge gradient-based extension of multi-task BO on the lowerlevel, but the theoretical justification for the combination of these two different criteria has not been revealed. Further, Chew et al. (2025) propose the well-known GP upper confidence bound (UCB) based approach to bilevel BO, called BILBO. Although BILBO has a theoretical regret guarantee, in general, the performance of GP-UCB based methods depend on the selection of the balancing parameter of the exploitation and exploration, because the theoretically recommended value often does not provide the best performance (Srinivas et al., 2010).

We propose an information-theoretic approach that considers the simultaneous information gain for both the upper- and lower-optimal solutions and values, which we call bilevel information gain. This enables us to define a unified criterion that measures the benefit for both level problems simultaneously, which is not necessarily common in the case of bilevel methods as mentioned in

the previous paragraph. Although the effectiveness of information-theoretic BO has been shown in several different contexts (e.g., Hennig & Schuler, 2012; Hernández-Lobato et al., 2014; Hoffman & Ghahramani, 2015; Wang & Jegelka, 2017; Hernández-Lobato et al., 2015; Hernandez-Lobato et al., 2016; Suzuki et al., 2020; Takeno et al., 2022a;b; Hvarfner et al., 2022; Tu et al., 2022), it has not been combined with bilevel optimization, to our knowledge. We first define bilevel information gain by extending the idea of the joint entropy search (Hvarfner et al., 2022; Tu et al., 2022). Unfortunately, the original definition of bilevel information gain is computationally intractable, and we show that a natural extension of the truncation based approximation, which has been widely employed in information-theoretic BO (e.g., Wang & Jegelka, 2017), can be derived. By combining the truncation based approximation and a variational lower bound (Takeno et al., 2022b), we obtain our criterion called Bilevel optimization via Lower-bound based Joint Entropy Search (BLJES). Further, while we mainly consider 'coupled setting' in which upper- and lower-level observations are obtained simultaneously, 'decoupled setting', in which a separate observation for each level is available, is also discussed. For example, in the case that the objective function values are outputs of some simulators (e.g., physical simulation), if a common simulator provides both level observations simultaneously, coupled setting is suitable, while if the upper- and the lower-level observations are from different simulators, decoupled setting can be more appropriate. We further propose an extension for the case that each level problem has inequality constraints (i.e., each level problem is a constraint problem).

Our contributions are summarized as follows.

- We show an information-theoretic formulation of bilevel BO, which has never been explored, to our knowledge. Bilevel information gain is defined to measure the benefit for both level problems.
- We derive a lower bound based approximation of bilevel information gain. We extend
 the standard truncation based approach in the single-level information-theoretic BO to the
 bilevel problem.
- We further propose extensions for decoupled setting and constraint problems. We show that our framework can handle these settings by a natural extension of bilevel information gain.

We demonstrate effectiveness of BLJES through functions generated from Gaussian processes and several benchmark problems.

2 Preliminaries

Bilevel optimization. Let $f: X \times \Theta \to \mathbb{R}$ and $g: X \times \Theta \to \mathbb{R}$ denote the upper- and the lower-level objective functions, respectively, both of which are assumed to be costly black-box functions. The upper- and the lower-level variables are denoted by $x \in X$ and $\theta \in \Theta$, respectively, where $X \subset \mathbb{R}^{d_X}$ and $\Theta \subset \mathbb{R}^{d_\Theta}$. The bilevel optimization problem is formulated as:

$$\max_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x}))$$
s.t. $\boldsymbol{\theta}^*(\boldsymbol{x}) = \underset{\boldsymbol{\theta} \in \Theta}{\arg \max} g(\boldsymbol{x}, \boldsymbol{\theta}),$ (1)

where $\theta^*(x)$ represents the optimal solution of the lower-level problem for a given upper-level variable x. For simplicity, we assume the lower-level optimum $\theta^*(x)$ is uniquely determined for each x (called optimistic setting (Sinha et al., 2020)). The bilevel optimal solution is denoted by (x^*, θ^*) , while the lower-level optimum corresponding to a given x is written as $(x, \theta^*(x))$, noting that $\theta^* = \theta^*(x^*)$. The upper- and the lower-level optimal values are denoted by $f^* := f(x^*, \theta^*)$ and $g^* := g(x^*, \theta^*)$, respectively. Figure 1 shows an illustration. Observations of the objective functions contain additive Gaussian noise $y^f_{(x,\theta)} := f(x,\theta) + \epsilon^f$, $\epsilon^f \sim \mathcal{N}(0, \{\sigma^f_{\text{noise}}\}^2)$, and $y^g_{(x,\theta)} := g(x,\theta) + \epsilon^g$, $\epsilon^g \sim \mathcal{N}(0, \{\sigma^g_{\text{noise}}\}^2)$, where σ^f_{noise} are the noise standard deviations, respectively. Let $\mathcal{D}_t := \{(x_i,\theta_i,y_i^f,y_i^g)\}_{i=1}^n$ be the dataset that we have at the t-th iteration of BO, where $y^f_i := y^f_{(x_i,\theta_i)}$ and $y^g_i := y^g_{(x_i,\theta_i)}$. The number of observed points n is $t + n_0$, where n_0 is the number of the initial observations.

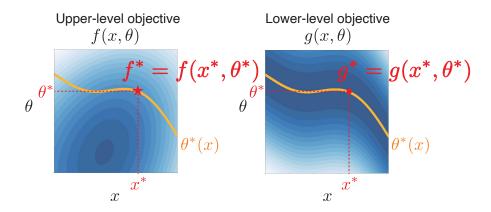


Figure 1: Example of Bilevel optimization ($d_X = 1, d_{\Theta} = 1$) and its optimal solution. For each upper-level variable x, the feasible solution is defined by $\theta^*(x) = \arg \max_{\theta} g(x, \theta)$, shown as the orange line. The optimal solution (x^*, θ^*) is the maximizer of f on the orange line.

Gaussian process. The upper- and the lower-level objective functions are each modeled by independent Gaussian processes (GPs) with kernel functions $k^f((x,\theta),(x',\theta'))$ and $k^g((x,\theta),(x',\theta'))$, respectively. Given the dataset \mathcal{D}_t , the predictive distribution of an objective function $h \in \{f,g\}$ at a point (x,θ) is expressed as:

$$h(\boldsymbol{x}, \boldsymbol{\theta}) \mid \mathcal{D}_{t} \sim \mathcal{N}(\mu_{t}^{h}(\boldsymbol{x}, \boldsymbol{\theta}), \{\sigma_{t}^{h}(\boldsymbol{x}, \boldsymbol{\theta})\}^{2}),$$

$$\mu_{t}^{h}(\boldsymbol{x}, \boldsymbol{\theta}) = \boldsymbol{k}^{h \top} (\boldsymbol{K}^{h} + \{\sigma_{\text{noise}}^{h}\}^{2} \boldsymbol{I})^{-1} \boldsymbol{y}^{h},$$

$$\{\sigma_{t}^{h}(\boldsymbol{x}, \boldsymbol{\theta})\}^{2} = k^{h}((\boldsymbol{x}, \boldsymbol{\theta}), (\boldsymbol{x}, \boldsymbol{\theta})) - \boldsymbol{k}^{h \top} (\boldsymbol{K}^{h} + \{\sigma_{\text{noise}}^{h}\}^{2} \boldsymbol{I})^{-1} \boldsymbol{k}^{h},$$

$$(2)$$

where $\mathbf{y}^h = (y_1^h, \dots, y_n^h)^\top$, $\mathbf{k}^h = (k^h((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}_1, \boldsymbol{\theta}_1)), \dots, k^h((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}_t, \boldsymbol{\theta}_n)))^\top$, and $\mathbf{K}^h \in \mathbb{R}^{n \times n}$ is the kernel matrix with an entry $k^h((\mathbf{x}_i, \boldsymbol{\theta}_i), (\mathbf{x}_j, \boldsymbol{\theta}_j))$ at a position (i, j). Here, $\mathbf{I} \in \mathbb{R}^{n \times n}$ denotes the identity matrix.

Bayesian optimization. We consider BO for the bilevel optimization problem (1). Bayesian optimization is a method for efficiently optimizing black-box functions with a limited number of samples. At step t, GPs are fitted to the dataset \mathcal{D}_t , and the next query point is determined as $\arg\max_{x,\theta} \alpha_t(x,\theta)$, where $\alpha_t(x,\theta)$ denotes the acquisition function. After sampling the query point, the newly obtained data are added to the dataset, and the GPs are refitted.

3 BILEVEL OPTIMIZATION VIA LOWER-BOUND BASED JOINT ENTROPY SEARCH

We consider bilevel BO based on the information gain for the optimal solutions and values $(\boldsymbol{x}^*, \boldsymbol{\theta}^*, f^*, \text{and } g^*)$ achieved by next observations $y^f_{(\boldsymbol{x},\boldsymbol{\theta})}$ and $y^g_{(\boldsymbol{x},\boldsymbol{\theta})}$, which we call bilevel information gain. Note that we regard the optimal $(\boldsymbol{x}^*, \boldsymbol{\theta}^*, f^*, g^*)$ as random variables defined by the predictive distributions of the objective functions $f(\boldsymbol{x},\boldsymbol{\theta})$ and $g(\boldsymbol{x},\boldsymbol{\theta})$. Our approach combines the concept of entropy search (Hennig & Schuler, 2012), in particular, joint entropy search (Hvarfner et al., 2022; Tu et al., 2022), and a variational lower bound based approximation of mutual information (MI). We refer to our proposed method as *Bilevel optimization via Lower-bound based Joint Entropy Search* (BLJES).

3.1 Lower Bound of Mutual Information

Bilevel information gain is represented as the MI between the candidate observations $(y_{(\mathbf{x},\theta)}^f, y_{(\mathbf{x},\theta)}^g)$ and the set of the optimal solutions and their upper- and lower-objective values $\{f^*, g^*, \mathbf{x}^*, \theta^*\}$:

$$\mathrm{MI}(y_{(\boldsymbol{x},\boldsymbol{\theta})}^f,y_{(\boldsymbol{x},\boldsymbol{\theta})}^g\;;\;f^*,g^*,\boldsymbol{x}^*,\boldsymbol{\theta}^*\mid\mathcal{D}_t).$$

This criterion naturally allows simultaneous consideration of both the upper- and the lower-objectives. Since the direct evaluation of this MI is difficult, we employ a lower bound based approximation.

Let $\Omega := \{y_{(x,\theta)}^f, y_{(x,\theta)}^g, f^*, g^*, x^*, \theta^*\}$. Our lower bound of the MI is derived by a technique that is often used in the context of the variational approximation (e.g., Poole et al., 2019):

$$MI(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g}; f^{*}, g^{*}, x^{*}, \theta^{*} \mid \mathcal{D}_{t}) = \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] \\
= \mathbb{E}_{f^{*}, g^{*}, x^{*}, \theta^{*}} \left[\mathbb{E}_{y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}} \left[\log \frac{q(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] \\
+ KL \left(p(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}) \parallel q(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}) \right) \right] \\
\geq \mathbb{E}_{\Omega} \left[\log \frac{q(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f}, y_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] =: LB(x, \theta), \tag{3}$$

where KL is Kullback-Leibler (KL) divergence and $q(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t)$ is a variational distribution (q can be any density function as far as the KL divergence can be defined). The inequality of the last line can be taken because the KL divergence is non-negative (the equality holds when $p(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t) = q(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t)$). Similar lower bounds of the MI have been used in information-theoretic multi-objective and constraint BO (Ishikura & Karasuyama, 2025; Takeno et al., 2022b).

The variational distribution q is an approximation of $p(y_{(x,\theta)}^f, y_{(x,\theta)}^g, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t)$ for which an exact analytical representation is difficult to know. The difficulty is in the conditioning by the optimal solutions and values, for which the most widely accepted approach in information-theoretic BO is to use truncated distributions (e.g., Wang & Jegelka, 2017; Suzuki et al., 2020; Tu et al., 2022). For example, in the case of well-known max-value entropy search (MES), proposed by (Wang & Jegelka, 2017) for the standard single-level problem $\max_x f(x)$, the predictive distribution conditioning on the max-value $f_{\text{unc}}^* := \arg\max_x f(x)$ is approximated by the truncated normal distribution, i.e., $p(f(x) \mid f_{\text{unc}}^*) \approx p(f(x) \mid f(x) \leq f_{\text{unc}}^*)$. When f_{unc}^* is given, $f(x') \leq f_{\text{unc}}^*$ should hold for any x' (and there should exist at least one x' such that $f(x') = f_{\text{unc}}^*$), while MES simplifies this condition so that $f(x) \leq f_{\text{unc}}^*$ holds only for the current x. Similar simplifications have been employed by most of information-theoretic BO algorithms and shown superior performance.

We extend the truncation based approach to our bilevel problem as follows.

$$q(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t) := p(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f(x, \theta^*(x)) \le f^*, g(x^*, \theta) \le g^*, \mathcal{D}_t^+),$$
(4)

where $\mathcal{D}_t^+ = \mathcal{D}_t \cup \{(\boldsymbol{x}^*, \boldsymbol{\theta}^*, f^*, g^*)\}$ is the dataset augmented by the optimal point $(\boldsymbol{x}^*, \boldsymbol{\theta}^*, f^*, g^*)$. The right hand side has the three conditions, each of which can be interpreted as follows.

- When f^* is given, $f(x', \theta^*(x')) \le f^*$ should hold for $\forall x'$. However, this condition is computationally intractable as mentioned for the case of MES. Based on a similar idea of MES, the condition $f(x, \theta^*(x)) \le f^*$ is only imposed on the current x.
- When g^* is given, $g(x^*, \theta') \le g^*$ should hold for $\forall \theta'$. We replace it with $g(x^*, \theta) \le g^*$ in which the inequality is only imposed on the current θ .
- In the right hand side of (4), \mathcal{D}_t is replaced \mathcal{D}_t^+ . This condition can impose that the GPs satisfy $f(\mathbf{x}^*, \mathbf{\theta}^*) = f^*$ and $g(\mathbf{x}^*, \mathbf{\theta}^*) = g^*$ by adding $(\mathbf{x}^*, \mathbf{\theta}^*, f^*, g^*)$ into the training data.

By substituting (4) into (3) and using the conditional independence of $y_{(x,\theta)}^f$ and $y_{(x,\theta)}^g$ in the right hand side of (4), we see

$$LB(\boldsymbol{x}, \boldsymbol{\theta}) = \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f}, y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid f(\boldsymbol{x}, \boldsymbol{\theta}^{*}(\boldsymbol{x})) \leq f^{*}, g(\boldsymbol{x}^{*}, \boldsymbol{\theta}) \leq g^{*}, \mathcal{D}_{t}^{+})}{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f}, y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})} \right]$$

$$= \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f} \mid f(\boldsymbol{x}, \boldsymbol{\theta}^{*}(\boldsymbol{x})) \leq f^{*}, \mathcal{D}_{t}^{+})}{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f} \mid \mathcal{D}_{t})} + \log \frac{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid g(\boldsymbol{x}^{*}, \boldsymbol{\theta}) \leq g^{*}, \mathcal{D}_{t}^{+})}{p(y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})} \right]. \quad (5)$$

The inside of the expectation (5) can be analytically derived. For both the first and second terms, the denominators are the predictive distribution of the GPs, whose density can be obtained from (2). Next, we consider the numerator of the first term of (5). Although the truncation is imposed on $f(x, \theta^*(x))$, the distribution is for $y_{(x,\theta)}^f$ that has different input point (x,θ) from the truncated point, unlike the case of MES. The following theorem shows that the analytical representation can still be derived even with this difference:

Theorem 3.1. Let (m_1^f, s_1^f) , (m_2^f, s_2^f) , and (m_3^f, s_3^f) be the mean and standard deviation of $p(f(\mathbf{x}, \boldsymbol{\theta}^*(\mathbf{x})) \mid y_{(\mathbf{x}, \boldsymbol{\theta})}^f, \mathcal{D}_t^+)$, $p(f(\mathbf{x}, \boldsymbol{\theta}^*(\mathbf{x})) \mid \mathcal{D}_t^+)$, and $p(y_{(\mathbf{x}, \boldsymbol{\theta})}^f \mid \mathcal{D}_t^+)$, respectively. Then,

$$p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid f(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x})) \leq f^{*}, \mathcal{D}_{t}^{+}) = \begin{cases} \Phi\left(\frac{f^{*}-m_{1}^{f}}{s_{1}^{f}}\right) \phi\left(\frac{y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f}-m_{3}^{f}}{s_{3}^{f}}\right) / \left\{\Phi\left(\frac{f^{*}-m_{2}^{f}}{s_{2}^{f}}\right) s_{3}^{f}\right\} & if \, \boldsymbol{x} \neq \boldsymbol{x}^{*}, \\ \phi\left(\frac{y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f}-m_{3}^{f}}{s_{3}^{f}}\right) / s_{3}^{f} & otherwise, \end{cases}$$

$$(6)$$

where ϕ and Φ are the probability density function (PDF) and cumulative density function (CDF) of the standard normal distribution, respectively.

The proof is in Appendix A.1. Note that all of $\{(m_i^f, s_i^f)\}_{i=1}^3$ can be analytically calculated from the GP posterior of f for which details are also show in Appendix A.1. For the numerator in the second term of (5) can be reduced to the similar analytical form, which is shown in Appendix A.2. As a result, we obtain an analytical form of the inside of the expectation (5).

3.2 Computations

The expectation of (5) is approximated by the Monte-Carlo method by sampling $\Omega := \{y_{(x,\theta)}^f, y_{(x,\theta)}^g, f^*, g^*, x^*, \theta^*\}$. The sample of all the elements of Ω can be obtained through the sample of the objective functions f and g. We use random Fourier feature (RFF) (Rahimi & Recht, 2008), by which the GP posterior can be approximated by the Bayesian linear model. For a D-dimensional RFF vector $\phi(x,\theta) \in \mathbb{R}^D$, the linear model $\mathbf{w}^{h\top}\phi(x,\theta)$ can be constructed, where $\mathbf{w}^h \in \mathbb{R}^D$ is a parameter vector and $h \in \{f,g\}$ represents one of objective functions. By sampling \mathbf{w}^h from the posterior, approximate sample paths of f and g are obtained, which denoted as f and g, respectively. Then, the sample of $(f^*, g^*, x^*, \theta^*)$ is obtained through $\max_x \tilde{f}(x, \tilde{\theta}^*(x))$ s.t. $\tilde{\theta}^*(x) = \arg\max_\theta \tilde{g}(x, \theta)$, which can be seen as a white-box bilevel optimization problem. Since both f and g are represented by the Bayesian liner model, they are differentiable. Then, the gradient $\partial \tilde{f}(x, \theta^*(x))/\partial x$ can be obtained through the implicit function theorem (see Appendix B for detail), by which standard gradient based optimization methods can be applied. The samples of $y_{(x,\theta)}^f$ and $y_{(x,\theta)}^g$ can be obtained by adding the random noise from $\mathcal{N}(0, \{\sigma_{\text{noise}}^f\}^2)$ and $\mathcal{N}(0, \{\sigma_{\text{noise}}^g\}^2)$ to the sampled $f(x,\theta)$ and $g(x,\theta)$, respectively.

Let K be the number of samplings of Ω . In a variety of contexts of information-theoretic BO (e.g., Wang & Jegelka, 2017), the superior performance has been repeatedly shown with small K settings (e.g., 10). After obtaining K samples of Ω , the Monte-Carlo approximation of (5) can be calculated for any (x, θ) (Note that these K samples are reused during the acquisition function optimization without regenerating for each candidate (x, θ)). Since (5) contains $\theta^*(x)$, the maximization of the approximate (5) is also bilevel optimization, for which gradient-based methods can also be applied through the implicit function theorem (details are also in Appendix B).

4 Extensions

We here describe two extensions of BLJES, which are for decoupled setting and constraint problems.

4.1 DECOUPLED SETTING

We mainly consider the setting in which $y_{(x,\theta)}^f$ and $y_{(x,\theta)}^g$ are observed simultaneously, which we call 'coupled' setting. On the other hand, only one of $y_{(x,\theta)}^f$ or $y_{(x,\theta)}^g$ can be separately observed in some scenarios. In this paper, this setting is called 'decoupled' setting, inspired by the similar setting in multi-objective BO (Hernandez-Lobato et al., 2016).

A natural criterion for decoupled setting is information gain obtained by only one of $y_{(x,\theta)}^f$ or $y_{(x,\theta)}^g$, for which the lower bounds can be derived by the almost same way as (3):

$$\operatorname{MI}(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f}; f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*} \mid \mathcal{D}_{t}) \geq \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid f(\boldsymbol{x}, \boldsymbol{\theta}^{*}(\boldsymbol{x})) \leq f^{*}, \mathcal{D}_{t}^{+})}{p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid \mathcal{D}_{t})} \right], \tag{7}$$

$$\operatorname{MI}(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{g}; f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*} \mid \mathcal{D}_{t}) \geq \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid g(\boldsymbol{x}^{*}, \boldsymbol{\theta}) \leq g^{*}, \mathcal{D}_{t}^{+})}{p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})} \right]. \tag{8}$$

The derivation is in Appendix C. For both the inside of the expectation of (7) and (8), the analytical calculations shown in section 3.1 can be used. The expectation is approximated by Monte-Carlo sampling of Ω , which is also same as coupled setting. As a result, the decision making not only for selecting (x, θ) , but also selecting the upper- or the lower-observation (i.e., $y_{(x,\theta)}^f$ or $y_{(x,\theta)}^g$) can be performed.

4.2 Incorporating Constraint Problems

In a more general formulation of bilevel optimization, constraints are imposed on both of the upperand the lower-level problems. When we have N and M inequality constraints for the upper- and the lower-level problems, respectively, the bilevel optimization problem is written as

$$\begin{aligned} \max_{\boldsymbol{x} \in X} f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) \\ \text{s.t. } c_n^U(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) \geq 0, n = 1, \dots, N \\ \boldsymbol{\theta}^*(\boldsymbol{x}) = \underset{\boldsymbol{\theta} \in \Theta}{\arg\max} \{ g(\boldsymbol{x}, \boldsymbol{\theta}) \mid c_m^L(\boldsymbol{x}, \boldsymbol{\theta}) \geq 0, m = 1, \dots, M \}, \end{aligned}$$

where $c^U: \mathbb{R}^{d_X+d_\Theta} \to \mathbb{R}$ and $c^L: \mathbb{R}^{d_X+d_\Theta} \to \mathbb{R}$ are constraint functions. We assume that c^U and c^L are also expensive black-box functions and modeled by the independent GPs.

For constraint BO, Takeno et al. (2022b) show an information-theoretic approach based on a lower bound with a truncated variational distribution. By combining the truncation shown by (Takeno et al., 2022b) and our bilevel information gain, we can extend BLJES to the bilevel constraint problem. In particular, the conditioning on the predictive distributions by the optimal points are required to extend. For example, if f^* is given, the inequality $f(x, \theta^*(x)) \leq f^*$ is imposed only when the constraints $c_n^U(x, \theta^*(x)) \geq 0, n = 1, \ldots, N$ hold (if the constraints are not satisfied, $f(x, \theta^*(x))$ is not truncated). Details are in Appendix D.

5 Experiments

We evaluated the performance of BLJES by using sample path functions from the GP prior and several benchmark functions. For baselines, we employed Random selection and BILBO. The initial number of observations was set $n_0 = 5$ random points. The both level observations contain an additive noise whose mean is 0 and standard deviation is 10^{-3} . Each experiment was performed 10 times with different initial points. We used the Gaussian kernel for both the GPs of f and g, in which the prior mean, the kernel length-scale, the output scale, and the noise variance are optimized

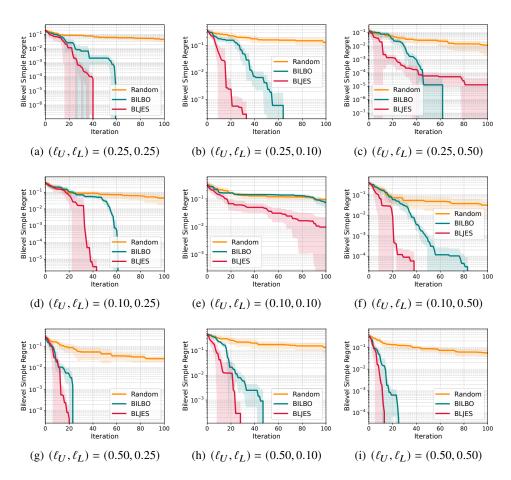


Figure 2: Regret comparison on functions from the GP prior.

by the marginal likelihood at every iteration. In BLJES, the number of Mont-Carlo samples was set as K = 30. We here employed the pool setting (query candidates are finite grid points defined later) because BILBO is proposed for the finite domain setting.

For the metric at the t-th iteration, we used the following criterion, denoted as bilevel simple regret:

$$\min_{i \in [n0+t]} \max_{h \in \{f,g\}} r_h(\mathbf{x}_i, \boldsymbol{\theta}_i), \tag{9}$$

where

$$\begin{split} r_f(\boldsymbol{x}_i, \boldsymbol{\theta}_i) &= \max(0, f^* - f(\boldsymbol{x}_i, \boldsymbol{\theta}_i)) / (f^* - \min_{\boldsymbol{x}, \boldsymbol{\theta}} f(\boldsymbol{x}, \boldsymbol{\theta})), \\ r_g(\boldsymbol{x}_i, \boldsymbol{\theta}_i) &= (g(\boldsymbol{x}_i, \boldsymbol{\theta}^*(\boldsymbol{x}_i)) - g(\boldsymbol{x}_i, \boldsymbol{\theta}_i)) / (g(\boldsymbol{x}_i, \boldsymbol{\theta}^*(\boldsymbol{x}_i)) - \min_{\boldsymbol{\theta}} g(\boldsymbol{x}_i, \boldsymbol{\theta})). \end{split}$$

Our metric (9) takes the larger value between $r_f(x_i, \theta_i)$, which represents the regret of the upperlevel problem, and $r_g(x_i, \theta_i)$, which represents those of the lower-level problem. Since $f(x_i, \theta_i)$ can be larger than f^* , the 'max' operation is taken to guarantee $r_f(x_i, \theta_i) \geq 0$, while the numerator of $r_g(x_i, \theta_i)$ is non-negative without 'max' from the definition of $\theta^*(x_i)$. The denominators of $r_f(x_i, \theta_i)$ and $r_g(x_i, \theta_i)$ are for absorbing the scale difference of two objectives. In (9), we employed the best value obtained during the entire search procedure by taking the minimum with respect to observed points.

We first provide the results on coupled setting for GP sample path functions (section 5.1) and benchmark functions (section 5.2). Further, the results on decoupled setting (section 5.3) and different *K* settings (section 5.4) are also reported. Appendix presents other details of the settings (Appendix E.1) and results such as a larger noise setting (Appendix E.2), the continuous domain (Appendix E.5), and constraint problems (Appendix E.6).

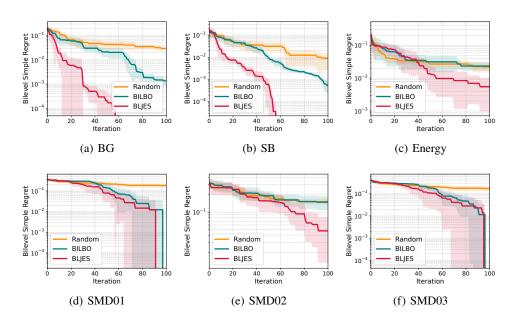


Figure 3: Regret comparison on benchmark problems.

5.1 Sample Path from GP Prior

We first used the sample path from the GP prior as the true objective functions, i.e., $f \sim \mathcal{GP}(0,k)$ and $g \sim \mathcal{GP}(0,k)$, in which k is the Gaussian kernel $k((x,\theta),(x',\theta')) = \exp\{-(\|x-x'\|^2 + \|\theta-\theta'\|^2)/(2\ell^2)\}$. For the length scale ℓ , we use different values $\ell_U \in \{0.25,0.10,0.50\}$ for f and $\ell_L \in \{0.25,0.10,0.50\}$ for g, respectively. The input space is $d_X = d_{\Theta} = 1$ and [0,1] for each. The candidate points are a combination of 100 grid points in each dimension (100² points).

The results are shown in Fig. 2. Overall, BLJES shows superior performance for a variety of the length scale functions. Only for $(\ell_U, \ell_L) = (0.25, 0.50)$, BILBO decreased the regret to 0 faster at about 60 iterations, but BLJES also reached the small value (10^{-4}) around that iterations.

5.2 Benchmark Functions

We here used six benchmark problems. Two functions are created by combining benchmark functions of single-level optimization. In the first problem, denoted as BG, the upper objective is BraninHoo $(d_X = 1)$ and the lower objective Goldstein-price $(d_{\Theta} = 1)$, which was used in (Chew et al., 2025). In the second problem, denoted as SB, the upper objective is SixHumpCamel $(d_X = 1)$ and the lower objective BraninHoo $(d_{\Theta} = 1)$, which was used in (Ekmekcioglu et al., 2024). The third problem, denoted as Energy, is a simulator based energy market problem $(d_X = 2 \text{ and } d_{\Theta} = 2)$ introduced by (Chew et al., 2025), in which this data is regarded as a real-world dataset. From the fourth to the sixth problems, denoted as SMD01, 02, and 03 $(d_X = 2 \text{ and } d_{\Theta} = 2)$, are test problems specifically designed for bilevel optimization benchmark (Sinha et al., 2014). The number of grid points in each dimension is 100 for GB and SB (100^2 points) , and 10 for Energy and SMD (10^4 points) .

The results are shown in Fig. 3. BLJES has obviously superior performance in BG, SB, Energy, and SMD02. For SMD01 and SMD03, similar performance is shown in BLJES and BILBO, both of which rapidly decrease the regret compared with Random.

5.3 Decoupled Setting

We here evaluate performance on decoupled setting, for which regret comparison is shown in Fig. 4. The objective functions are the same functions used before. We see that MLJES shows smaller regret values for most of problems except only for SMD02 in which BILBO shows better performance. The results indicate that our MI based criterion is effective also for decoupled setting.

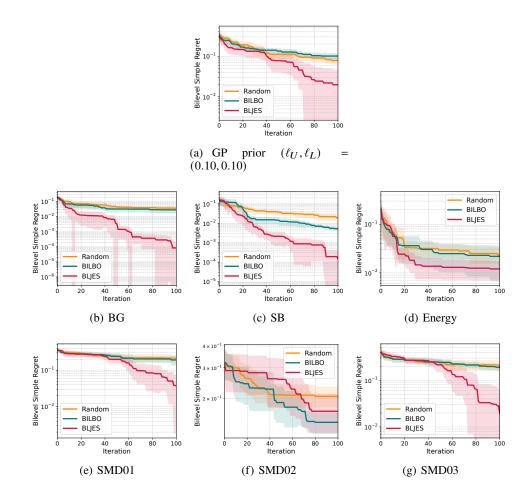


Figure 4: Regret comparison on benchmark problems under decoupled setting.

5.4 Effect of the Number of Samplings

We evaluate the effect of the number of samplings K on the performance. Figure 5 shows the regret of BLJES with K=10,20,30, and 50 on the BG benchmark problem. Note that the result of K=30 is same as Fig. 3 (a). Although K=50 was slightly better than other settings in the end of the optimization, we do not see large differences. Similar tendency has been reported in information-theoretic BO studies (Wang & Jegelka, 2017; Takeno et al., 2022a). See in Appendix E.4 for the results on other problems.

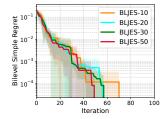


Figure 5: BLJES with different *K* on the BG benchmark.

6 Conclusion

We propose an information-theoretic approach to bilevel Bayesian optimization, called Bilevel optimization via Lower-bound based Joint Entropy Search (BLJES). BLJES considers information gain of optimal points and values of both the upper- and lower- level problems simultaneously, by which we can define a unified criterion that measures the benefit for both the problems. We derive a lower bound based approximation of bilevel information gain, which can be seen as a natural extension of the single level information-theoretic Bayesian optimization. Further, we also propose extensions for decoupled setting and constraint problems. The effectiveness of BLJES is demonstrated through sample path functions from Gaussian processes and benchmark functions.

REFERENCES

- Malek Abbassi, Abir Chaabani, Lamjed Ben Said, and Nabil Absi. An approximation-based chemical reaction algorithm for combinatorial multi-objective bi-level optimization problems. In 2021 IEEE Congress on Evolutionary Computation (CEC), pp. 1627–1634, 2021.
- Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A framework for efficient monte-carlo bayesian optimization. In *Advances in Neural Information Processing Systems 33*, 2020. URL http://arxiv.org/abs/1910.06403.
- Ruth Wan Theng Chew, Quoc Phong Nguyen, and Bryan Kian Hsiang Low. BILBO: Bilevel Bayesian optimization. In *Proceedings of the 42th International Conference on Machine Learning*, 2025.
- Vedat Dogan and Steven Prestwich. Bilevel optimization by conditional bayesian optimization. In *Machine Learning, Optimization, and Data Science: 9th International Conference, LOD 2023, Grasmere, UK, September 22–26, 2023, Revised Selected Papers, Part I,* pp. 243–258, Berlin, Heidelberg, 2023. Springer-Verlag. ISBN 978-3-031-53968-8.
- Omer Ekmekcioglu, Nursen Aydin, and Juergen Branke. Bayesian optimization of bilevel problems, 2024.
- Shi Fu, Fengxiang He, Xinmei Tian, and Dacheng Tao. Convergence of bayesian bilevel optimization. In *The Twelfth International Conference on Learning Representations*, 2024.
- Philipp Hennig and Christian J. Schuler. Entropy search for information-efficient global optimization. *Journal of Machine Learning Research*, 13(57):1809–1837, 2012.
- Daniel Hernandez-Lobato, Jose Hernandez-Lobato, Amar Shah, and Ryan Adams. Predictive entropy search for multi-objective Bayesian optimization. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pp. 1492–1501. PMLR, 2016.
- José Miguel Hernández-Lobato, Matthew W. Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. In *Advances in Neural Information Processing Systems* 27, pp. 918–926. Curran Associates, Inc., 2014.
- José Miguel Hernández-Lobato, Michael A. Gelbart, Matthew W. Hoffman, Ryan P. Adams, and Zoubin Ghahramani. Predictive entropy search for Bayesian optimization with unknown constraints. In *Proceedings of the 32th International Conference on Machine Learning*, volume 37, pp. 1699–1707. PMLR, 2015.
- J. Herskovits, A. Leontiev, G. Dias, and G. Santos. Contact shape optimization: a bilevel programming approach. *Structural and Multidisciplinary Optimization*, 20(3):214–221, Nov 2000. ISSN 1615-1488.
- Matthew W. Hoffman and Zoubin Ghahramani. Output-space predictive entropy search for flexible global optimization. In *NIPS Workshop on Bayesian Optimization*, 2015.
- Carl Hvarfner, Frank Hutter, and Luigi Nardi. Joint entropy search for maximally-informed bayesian optimization. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2022. Curran Associates Inc.
- M. Ishikura and M. Karasuyama. Pareto-frontier entropy search with variational lower bound maximization. In *Proceedings of the 42th International Conference on Machine Learning*, 2025.
- Md Monjurul Islam, Hemant Kumar Singh, and Tapabrata Ray. Efficient global optimization for solving computationally expensive bilevel optimization problems. In 2018 IEEE Congress on Evolutionary Computation (CEC), pp. 1–8, 2018.
- Emmanuel Kieffer, Grégoire Danoy, Pascal Bouvry, and Anass Nagih. Bayesian optimization approach of general bi-level problems. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, GECCO '17, pp. 1614–1621, New York, NY, USA, 2017. Association for Computing Machinery.

- Victor Picheny, Tobias Wagner, and David Ginsbourger. A benchmark of kriging-based infill criteria for noisy optimization. *Structural and multidisciplinary optimization*, 48(3):607–626, 2013.
- Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational bounds of mutual information. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 5171–5180. PMLR, 2019.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20*, pp. 1177–1184. Curran Associates, Inc., 2008.
- Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. Test problem construction for single-objective bilevel optimization. *Evolutionary Computation*, 22(3):439–477, 2014.
- Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. A review on bilevel optimization: From classical to evolutionary approaches and applications, 2020.
- N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pp. 1015–1022. Omnipress, 2010.
- Varun Suryan, Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. Handling inverse optimal control problems using evolutionary bilevel optimization. In 2016 IEEE Congress on Evolutionary Computation (CEC), pp. 1893–1900, 2016.
- Shinya Suzuki, Shion Takeno, Tomoyuki Tamura, Kazuki Shitara, and Masayuki Karasuyama. Multi-objective Bayesian optimization using Pareto-frontier entropy. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 9279–9288. PMLR, 2020.
- Shion Takeno, Hitoshi Fukuoka, Yuhki Tsukada, Toshiyuki Koyama, Motoki Shiga, Ichiro Takeuchi, and Masayuki Karasuyama. A generalized framework of multi-fidelity max-value entropy search through joint entropy. *Neural Computation*, 2022a. To appear.
- Shion Takeno, Tomoyuki Tamura, Kazuki Shitara, and Masayuki Karasuyama. Sequential and parallel constrained max-value entropy search via information lower bound. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 20960–20986. PMLR, 2022b.
- Ben Tu, Axel Gandy, Nikolas Kantas, and Behrang Shafei. Joint entropy search for multi-objective Bayesian optimization. *Advances in Neural Information Processing Systems*, 35:9922–9938, 2022.
- Bing Wang, Hemant Kumar Singh, and Tapabrata Ray. Comparing expected improvement and kriging believer for expensive bilevel optimization. In 2021 IEEE Congress on Evolutionary Computation (CEC), pp. 1635–1642, 2021.
- Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient Bayesian optimization. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 3627–3635. PMLR, 2017.

A Derivation of Lower Bound

A.1 Proof of Theorem 3.1

From Bayes theorem,

$$p(y_{(x,\theta)}^{f} \mid f_{(x,\theta^{*}(x))} \leq f^{*}, \mathcal{D}_{t}^{+}) = \frac{p(f(x,\theta^{*}(x)) \leq f^{*} \mid y_{(x,\theta)}^{f}, \mathcal{D}_{t}^{+})p(y_{(x,\theta)}^{f} \mid \mathcal{D}_{t}^{+})}{p(f(x,\theta^{*}(x)) \leq f^{*} \mid \mathcal{D}_{t}^{+})}.$$
 (10)

All the three densities in the right hand side, the analytical representations can be derived as follows.

 • The probability $p(f(x, \theta^*(x)) \le f^* \mid y_{(x,\theta)}^f, \mathcal{D}_t^+)$ is calculated by the density

$$f_{(x,\theta^*(x))} \mid y_{(x,\theta)}^f, \mathcal{D}_t^+ \sim \mathcal{N}(m_1^f, \{s_1^f\}^2),$$

for which the mean m_1^f and variance $\{s_1^f\}^2$ can be derived by considering the conditional density of the joint posterior of $f(x, \theta^*(x)), y_{(x, \theta)}^f$, and $f(x^*, \theta^*)$ as

$$\begin{split} m_1^f &= \mu_t^f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) \\ &+ \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta})) \\ \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^f(\boldsymbol{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\operatorname{noise}}^f\}^2 & \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} y_{(\boldsymbol{x}, \boldsymbol{\theta})}^f - \mu_t^f(\boldsymbol{x}, \boldsymbol{\theta}) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}, \boldsymbol{\theta})) & \{\sigma_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)\}^2 \end{bmatrix}^{-1} \begin{bmatrix} y_{(\boldsymbol{x}, \boldsymbol{\theta})}^f - \mu_t^f(\boldsymbol{x}, \boldsymbol{\theta}) \\ f_{(\boldsymbol{x}^*, \boldsymbol{\theta}^*)} - \mu_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*) \end{bmatrix} \\ \{s_1^f\}^2 &= \{\sigma_t^f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta}))\}^2 \\ &- \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta})) \\ \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^f(\boldsymbol{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\operatorname{noise}}^f\}^2 & \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta})) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*) \\ \operatorname{Cov}_t^f((\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*), (\boldsymbol{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}^{-1} \end{bmatrix}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*, \boldsymbol{x}), (\boldsymbol{x}^*, \boldsymbol{\theta}$$

Note that $Cov_t^f((x, \theta), (x', \theta'))$ is the posterior covariance between $f(x, \theta)$ and $f(x', \theta')$, given \mathcal{D}_t . By using m_1^f and s_1^f , we have

$$p(f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) \le f^* \mid y_{(\boldsymbol{x}, \boldsymbol{\theta})}^f, \mathcal{D}_t^+) = \begin{cases} \Phi\left(\frac{f^* - m_1^f}{s_1^f}\right) & \text{if } \boldsymbol{x} \ne \boldsymbol{x}^*, \\ 1 & \text{otherwise,} \end{cases}$$
(11)

where Φ is the cumulative density function (CDF) of the standard normal distribution.

• Next, to calculate the denominator $p(f_{(x,\theta^*(x))} \le f^* \mid \mathcal{D}_t^+)$, we consider the density

$$f_{(\boldsymbol{x},\boldsymbol{\theta}^*(\boldsymbol{x}))} \mid \mathcal{D}_t^+ \sim \mathcal{N}(m_2^f, \{s_2^f\}^2),$$

for which the mean m_2^f and variance $\{s_2^f\}^2$ can be derived by considering the conditional density of the joint posterior of $f_{(x,\theta^*(x))}$ and $f_{(x^*,\theta^*)}$ as

$$\begin{split} m_2^f &= \mu_t^f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) + \frac{\operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*))}{\left\{\sigma_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)\right\}^2} (f(\boldsymbol{x}^*, \boldsymbol{\theta}^*) - \mu_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)) \\ \{s_2^f\}^2 &= \sigma_t^f(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) - \frac{\left\{\operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\left\{\sigma_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)\right\}^2} \end{split}$$

Then, we obtain

$$p(f_{(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x}))} \leq f^{*} \mid \mathcal{D}_{t}^{+}) = \begin{cases} \Phi\left(\frac{f^{*}-m_{2}^{f}}{s_{2}^{f}}\right) & \text{if } \boldsymbol{x} \neq \boldsymbol{x}^{*} \\ 1 & \text{otherwise} \end{cases}$$
(12)

• The density $p(y_{(\mathbf{x}, \boldsymbol{\theta})}^f \mid \mathcal{D}_t^+)$ can also be derived by a similar approach as

$$y_{(\mathbf{r},\boldsymbol{\theta})}^f \mid \mathcal{D}_t^+ \sim \mathcal{N}(m_3^f, s_3^f),$$

where

$$m_3^f = \mu_t^f(\boldsymbol{x}, \boldsymbol{\theta}) + \frac{\operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*))}{\left\{\sigma_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)\right\}^2} (f(\boldsymbol{x}^*, \boldsymbol{\theta}^*) - \mu_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*))$$

$$\{s_3^f\}^2 = \{\sigma_t^f(\boldsymbol{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\text{noise}}^f\}^2 - \frac{\left\{\operatorname{Cov}_t^f((\boldsymbol{x}, \boldsymbol{\theta}), (\boldsymbol{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\left\{\sigma_t^f(\boldsymbol{x}^*, \boldsymbol{\theta}^*)\right\}^2}.$$

Therefore,

$$p(y_{(\mathbf{x},\theta)}^f \mid \mathcal{D}_t^+) = \phi \left(\frac{y_{(\mathbf{x},\theta)}^f - m_3^f}{s_3^f} \right) / s_3^f.$$
 (13)

where ϕ is the density function of the standard normal distribution.

By substituting (11), (12), and (13) into (10), we obtain (6).

A.2 Analytical Representation of $p(y_{(\mathbf{x}, \boldsymbol{\theta})}^g \mid g(\mathbf{x}^*, \boldsymbol{\theta}) \leq g^*, \mathcal{D}_t^+)$

In the case of $p(y_{(\mathbf{x},\theta)}^g \mid g(\mathbf{x}^*,\theta) \leq g^*, \mathcal{D}_t^+)$, almost the same derivation can be applied as (6). Therefore, we here only show the final result

$$p(y_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid g(\boldsymbol{x}^*,\boldsymbol{\theta}) \leq g^*, \mathcal{D}_t^+) = \begin{cases} \Phi\left(\frac{g^* - m_1^g}{s_1^g}\right) \phi\left(\frac{y_{(\boldsymbol{x},\boldsymbol{\theta})}^g - m_3^g}{s_3^g}\right) / \Phi\left(\frac{g^* - m_2^g}{s_2^g}\right) s_3^g & \text{if } \boldsymbol{\theta} \neq \boldsymbol{\theta}^*, \\ \phi\left(\frac{y_{(\boldsymbol{x},\boldsymbol{\theta})}^g - m_3^g}{s_3^g}\right) / s_3^g & \text{otherwise,} \end{cases}$$

where

$$\begin{split} m_1^g &= \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}) + \begin{bmatrix} \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta})) \\ \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\text{noise}}^g\}^2 & \text{Cov}_t^g((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} y_{(\mathbf{x}, \boldsymbol{\theta})}^g - \mu_t^g(\mathbf{x}, \boldsymbol{\theta}) \\ \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}^*), (\mathbf{x}, \boldsymbol{\theta})) & \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2 \end{bmatrix}^{-1} \begin{bmatrix} y_{(\mathbf{x}, \boldsymbol{\theta})}^g - \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*) \\ g_{(\mathbf{x}^*, \boldsymbol{\theta}^*)} - \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}, \\ \{s_1^g\}^2 &= \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta})) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\text{noise}}^g\}^2 & \text{Cov}_t^f((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta})) \\ \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta})\}^2 + \{\sigma_{\text{noise}}^g\}^2 & \text{Cov}_t^f((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^{-1} \begin{bmatrix} \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta})) \\ \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^\top \begin{bmatrix} \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta}^*), (\mathbf{x}, \boldsymbol{\theta}^*), (\mathbf{x}, \boldsymbol{\theta}^*), (\mathbf{x}^*, \boldsymbol{\theta}^*) \end{bmatrix}^{-1} \end{bmatrix} \begin{bmatrix} \text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}, \boldsymbol{\theta}^*)) \end{bmatrix}, \\ m_2^g &= \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}) + \frac{\text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*))}{\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2} (g(\mathbf{x}^*, \boldsymbol{\theta}^*) - \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)) \end{bmatrix}^2, \\ \{s_2^g\}^2 &= \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta})\}^2 - \frac{\left\{\text{Cov}_t^g((\mathbf{x}^*, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2} (g(\mathbf{x}^*, \boldsymbol{\theta}^*) - \mu_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)), \\ \{s_2^g\}^2 &= \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta})\}^2 + \left\{\sigma_{\text{noise}}^g\}^2 - \frac{\left\{\text{Cov}_t^g((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2}, \\ \{s_3^g\}^2 &= \{\sigma_t^g(\mathbf{x}, \boldsymbol{\theta})\}^2 + \left\{\sigma_{\text{noise}}^g\}^2 - \frac{\left\{\text{Cov}_t^g((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2}, \\ \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2 - \frac{\left\{\text{Cov}_t^g((\mathbf{x}, \boldsymbol{\theta}), (\mathbf{x}^*, \boldsymbol{\theta}^*))\right\}^2}{\left\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}, \\ \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2 - \frac{\left\{\text{Cov}_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}{\left\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}, \\ \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2 - \frac{\left\{\text{Cov}_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}{\left\{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}, \\ \{\sigma_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\}^2 - \frac{\left\{\text{Cov}_t^g(\mathbf{x}^*, \boldsymbol{\theta}^*)\right\}^2}{\left\{\sigma_t$$

and $Cov_t^g((x,\theta),(x',\theta'))$ is the posterior covariance between $g(x,\theta)$ and $g(x',\theta')$ given \mathcal{D}_t .

B Detail of Gradient Computations

First, we consider the gradient for $\partial \tilde{f}(x, \theta^*(x))/\partial x$, which is required to obtain the sample of x^*, θ^*, f^* , and g^* . For $\tilde{\theta}^*(x) = \arg \max_{\theta} \tilde{g}(x, \theta)$, the implicit function theorem derives

$$\frac{\partial \tilde{\boldsymbol{\theta}}^*(\boldsymbol{x})}{\partial \boldsymbol{x}^\top} = -\left\{ \frac{\partial^2 \tilde{\boldsymbol{g}}(\boldsymbol{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}^*(\boldsymbol{x})} \right\}^{-1} \left. \frac{\partial^2 \tilde{\boldsymbol{g}}(\boldsymbol{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{x}^\top} \bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}^*(\boldsymbol{x})},$$

from which we can calculate

$$\frac{\partial \tilde{f}(x,\tilde{\boldsymbol{\theta}}^*(x))}{\partial x} = \left. \frac{\partial \tilde{f}(x,\boldsymbol{\theta})}{\partial x} \right|_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}^*(x)} + \left\{ \frac{\partial \tilde{\boldsymbol{\theta}}^*(x)}{\partial x^\top} \right\}^\top \left. \frac{\partial \tilde{f}(x,\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}^*(x)}.$$

Next, we consider the acquisition function maximization. Let

$$\tilde{a}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\theta}') \coloneqq \log \frac{p(\tilde{y}_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f} \mid \tilde{f}(\boldsymbol{x}, \boldsymbol{\theta}') \leq \tilde{f}^{*}, \tilde{\mathcal{D}}_{t}^{+})}{p(\tilde{y}_{(\boldsymbol{x}, \boldsymbol{\theta})}^{f} \mid \mathcal{D}_{t})} + \log \frac{p(\tilde{y}_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid \tilde{g}(\tilde{\boldsymbol{x}}^{*}, \boldsymbol{\theta}) \leq \tilde{g}^{*}, \tilde{\mathcal{D}}_{t}^{+})}{p(\tilde{y}_{(\boldsymbol{x}, \boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})}$$

be the inside of the expectation of (5) in which $\theta^*(x)$ is replaced with θ' , and variables in Ω is replaced by a sample, denoted with ". Note that $\tilde{D}_t^+ = \mathcal{D}_t \cup (\tilde{x}^*, \tilde{\theta}^*, \tilde{f}^*, \tilde{g}^*)$. Then, the gradient with respect to x can be written as

$$\frac{\partial \tilde{a}(x,\theta,\tilde{\theta}^*(x))}{\partial x} = \frac{\partial \tilde{a}(x,\theta,\theta')}{\partial x}\bigg|_{\theta'=\tilde{\theta}^*(x)} + \left\{\frac{\partial \tilde{\theta}^*(x)}{\partial x^\top}\right\}^\top \frac{\partial \tilde{a}(x,\theta,\theta')}{\partial \theta'}\bigg|_{\theta'=\tilde{\theta}^*(x)}$$

The gradient with respect θ can be obtained through the usual derivative.

C Lower Bound of Decoupled Setting

The lower bound of the information gain for the upper-level observation is

$$\begin{split} & \text{MI}(y_{(x,\theta)}^{f}; \ f^{*}, g^{*}, x^{*}, \theta^{*} \mid \mathcal{D}_{t}) = \mathbb{E}_{\Omega} \left[\log \frac{p(y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f} \mid \mathcal{D}_{t})} \right] \\ & = \mathbb{E}_{f^{*}, g^{*}, x^{*}, \theta^{*}} \left[\mathbb{E}_{y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}} \left[\log \frac{q(y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f} \mid \mathcal{D}_{t})} \right] \right. \\ & + \text{KL} \left(p(y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}) \parallel q(y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t}) \right) \right] \\ & \geq \mathbb{E}_{\Omega} \left[\log \frac{q(y_{(x,\theta)}^{f} \mid f^{*}, g^{*}, x^{*}, \theta^{*}, \mathcal{D}_{t})}{p(y_{(x,\theta)}^{f} \mid \mathcal{D}_{t})} \right] =: \text{LB}^{f}(x, \theta). \end{split}$$

By setting the variational distribution as

$$q(y_{(\mathbf{x},\boldsymbol{\theta})}^f \mid f^*, g^*, \mathbf{x}^*, \boldsymbol{\theta}^*, \mathcal{D}_t) \coloneqq p(y_{(\mathbf{x},\boldsymbol{\theta})}^f \mid f(\mathbf{x}, \boldsymbol{\theta}^*(\mathbf{x})) \le f^*, g(\mathbf{x}^*, \boldsymbol{\theta}) \le g^*, \mathcal{D}_t^+)$$
$$= p(y_{(\mathbf{x},\boldsymbol{\theta})}^f \mid f(\mathbf{x}, \boldsymbol{\theta}^*(\mathbf{x})) \le f^*, \mathcal{D}_t^+),$$

we obtain the lower bound (7).

D EXTENSION FOR CONSTRAINT PROBLEMS

Let

$$\begin{aligned} & \boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta})}^f \coloneqq (f(\boldsymbol{x},\boldsymbol{\theta}), c_1^U(\boldsymbol{x},\boldsymbol{\theta}), \dots, c_N^U(\boldsymbol{x},\boldsymbol{\theta}))^\top, \\ & \boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta})}^g \coloneqq (g(\boldsymbol{x},\boldsymbol{\theta}), c_1^L(\boldsymbol{x},\boldsymbol{\theta}), \dots, c_M^L(\boldsymbol{x},\boldsymbol{\theta}))^\top, \end{aligned}$$

be the vectors in which the objective function and the constraint functions are concatenated for the upper- and the lower-level problems, respectively, and

$$\mathbf{y}_{(\mathbf{x},\theta)}^{f} \coloneqq (y_{(\mathbf{x},\theta)}^{f}, y_{(\mathbf{x},\theta)}^{c_{1}^{U}}, \dots, y_{(\mathbf{x},\theta)}^{c_{N}^{U}})^{\mathsf{T}},$$
$$\mathbf{y}_{(\mathbf{x},\theta)}^{g} \coloneqq (y_{(\mathbf{x},\theta)}^{g}, y_{(\mathbf{x},\theta)}^{c_{1}^{L}}, \dots, y_{(\mathbf{x},\theta)}^{c_{M}})^{\mathsf{T}},$$

are the counterparts of noisy observations, where $y_{(x,\theta)}^{c_n^L} \coloneqq c_n^U(x,\theta) + \epsilon^{c_n^U}$, $\epsilon^{c_n^U} \sim \mathcal{N}(0, \{\sigma_{\text{noise}}^{c_n^L}\}^2)$ and $y_{(x,\theta)}^{c_m^U} \coloneqq c_m^L(x,\theta) + \epsilon^{c_m^L}$, $\epsilon^{c_m^U} \sim \mathcal{N}(0, \{\sigma_{\text{noise}}^{c_m^L}\}^2)$. We observe $(y_{x,\theta}^f, y_{x,\theta}^g)$ at every BO iteration for selected (x,θ) , i.e., $\mathcal{D}_t = \{(x_i,\theta_i,y_{x_i,\theta_i}^f,y_{x_i,\theta_i}^g)\}_{i=1}^n$ in the constraint setting. In addition to f and g, the independent GPs are also fitted to c_n^U and c_m^L , for which the posteriors given \mathcal{D}_t are written as $\mathcal{N}(\mu_t^{c_n^U}(x,\theta), \{\sigma_t^{c_n^U}(x,\theta)\}^2)$ and $\mathcal{N}(\mu_t^{c_m^L}(x,\theta), \{\sigma_t^{c_m^U}(x,\theta)\}^2)$, respectively.

D.1 LOWER BOUND

 The MI and its lower bound can be derived by the same approach as (3):

$$\begin{aligned} & \text{MI}(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \; ; \; f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*} \mid \mathcal{D}_{t}) = \mathbb{E}_{\Omega} \left[\log \frac{p(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t})}{p(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] \\ & = \mathbb{E}_{f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}} \left[\mathbb{E}_{\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t}} \left[\log \frac{q(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t})}{p(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] \\ & + \text{KL} \left(p(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t}) \parallel q(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t}) \right) \right] \\ & \geq \mathbb{E}_{\Omega} \left[\log \frac{q(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid f^{*}, g^{*}, \boldsymbol{x}^{*}, \boldsymbol{\theta}^{*}, \mathcal{D}_{t})}{p(\boldsymbol{y}_{(x,\theta)}^{f}, \boldsymbol{y}_{(x,\theta)}^{g} \mid \mathcal{D}_{t})} \right] =: \text{LB}_{c}(\boldsymbol{x}, \boldsymbol{\theta}), \end{aligned}$$

where here $\mathcal{D}_t^+ := \mathcal{D}_t \cup \{(\boldsymbol{x}^*, \boldsymbol{\theta}^*, f^*, g^*)\}.$

To define the variational distribution q, we follow the same approach as information-theoretic constraint BO proposed by (Takeno et al., 2022b). Let $\mathcal{A}^f = \{(c_0, \mathbf{c}) \mid c_0 \geq f^*, \mathbf{c} \geq 0, c_0 \in \mathbb{R}, \mathbf{c} \in \mathbb{R}^N\}$ and $\mathcal{A}^g = \{(c_0, \mathbf{c}) \mid c_0 \geq g^*, \mathbf{c} \geq 0, c_0 \in \mathbb{R}, \mathbf{c} \in \mathbb{R}^M\}$. When f^* is given, \mathcal{A}^f is the region that $h^f_{x,\theta^*(x)}$ cannot exist for $\forall x$. When g^* is given, \mathcal{A}^g is the region that $h^g_{x^*,\theta}$ cannot exist for $\forall \theta$. Based on the same simplification of the conditioning discussed in section 3.1, we define the variational distribution as

$$q(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid f^*, g^*, x^*, \theta^*, \mathcal{D}_t) \coloneqq p(y_{(x,\theta)}^f, y_{(x,\theta)}^g \mid \boldsymbol{h}_{(x,\theta^*(x))}^f \in \bar{\mathcal{A}^f}, \boldsymbol{h}_{(x^*,\theta)}^g \in \bar{\mathcal{A}^g}, \mathcal{D}_t^+),$$

where $\bar{\mathcal{A}}^f$ and $\bar{\mathcal{A}}^g$ are the complement sets of \mathcal{A}^f and \mathcal{A}^g , respectively. As a result, we see

$$\begin{split} \operatorname{LB}_{c}(\boldsymbol{x},\boldsymbol{\theta}) &= \mathbb{E}_{\Omega} \left[\log \frac{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f},\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x}))}^{f} \in \bar{\mathcal{A}}^{f}, \boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g}, \mathcal{D}_{t}^{+})}{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f},\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})} \right] \\ &= \mathbb{E}_{\Omega} \left[\log \frac{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid \boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x}))}^{f} \in \bar{\mathcal{A}}^{f}, \mathcal{D}_{t}^{+})}{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid \mathcal{D}_{t})} + \log \frac{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g}, \mathcal{D}_{t}^{+})}{p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t})} \right]. \end{split}$$

D.2 ANALYTICAL REPRESENTATION OF VARIATIONAL DISTRIBUTION

From Bayes theorem,

$$p(\mathbf{y}_{(\mathbf{x},\theta)}^{f} \mid \mathbf{h}_{(\mathbf{x},\theta^{*}(\mathbf{x}))}^{f} \in \bar{\mathcal{A}}^{f}, \mathcal{D}_{t}^{+}) = \frac{p(\mathbf{h}_{(\mathbf{x},\theta^{*}(\mathbf{x}))}^{f} \in \bar{\mathcal{A}}^{f} \mid \mathbf{y}_{(\mathbf{x},\theta)}^{f}, \mathcal{D}_{t}^{+})p(\mathbf{y}_{(\mathbf{x},\theta)}^{f} \mid \mathcal{D}_{t}^{+})}{p(\mathbf{h}_{(\mathbf{x},\theta^{*}(\mathbf{x}))}^{f} \in \bar{\mathcal{A}}^{f} \mid \mathcal{D}_{t}^{+})}.$$
 (14)

The density $p(\boldsymbol{h}_{(x,\theta^*(x))}^f \mid \boldsymbol{y}_{(x,\theta)}^f, \mathcal{D}_t^+)$ is an (N+1)-dimensional independent Gaussian distribution, for which the first dimension is $\mathcal{N}(m_1^f, \{s_1^f\}^2)$ shown in Appendix A.1 and from the second to the (N+1)-th dimension is $\mathcal{N}(m^{c_N^U}, \{s_1^{c_N^U}\}^2)$ where

$$\begin{split} m^{c_n^U} &= \mu_t^{c_n^U}(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) + \frac{\text{Cov}_t^{c_n^U}((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta}))}{\left\{\sigma_t^{c_n^U}(\boldsymbol{x}, \boldsymbol{\theta})\right\}^2} (y_{(\boldsymbol{x}, \boldsymbol{\theta})}^{c_n^U} - \mu_t^{c_n^U}(\boldsymbol{x}, \boldsymbol{\theta})), \\ &\{s^{c_n^U}\}^2 = \sigma_t^{c_n^U}(\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})) - \frac{\left\{\text{Cov}_t^{c_n^U}((\boldsymbol{x}, \boldsymbol{\theta}^*(\boldsymbol{x})), (\boldsymbol{x}, \boldsymbol{\theta}))\right\}^2}{\left\{\sigma_t^{c_n^U}(\boldsymbol{x}, \boldsymbol{\theta})\right\}^2}. \end{split}$$

As a result, we can derive

$$\begin{split} p(\boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x}))}^{f} \in \bar{\mathcal{A}}^{f} \mid \boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f}, \mathcal{D}_{t}^{+}) &= 1 - (1 - \Phi(\frac{f^{*} - m_{1}^{f}}{s_{1}^{f}})) \prod_{n=1}^{N} (1 - \Phi(\frac{0 - m_{n}^{c_{n}^{U}}}{s_{n}^{c_{n}^{U}}})), \\ p(\boldsymbol{h}_{(\boldsymbol{x},\boldsymbol{\theta}^{*}(\boldsymbol{x}))}^{f} \in \bar{\mathcal{A}}^{f} \mid \mathcal{D}_{t}^{+}) &= 1 - (1 - \Phi(\frac{f^{*} - m_{2}^{f}}{s_{2}^{f}})) \prod_{n=1}^{N} (1 - \Phi(\frac{0 - \mu_{t}^{c_{n}^{U}}(\boldsymbol{x}, \boldsymbol{\theta}^{*}(\boldsymbol{x}))}{\sigma_{t}^{c_{n}^{U}}(\boldsymbol{x}, \boldsymbol{\theta}^{*}(\boldsymbol{x}))})), \\ p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} \mid \mathcal{D}_{t}^{+}) &= \phi(\frac{\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{f} - m_{3}^{f}}{s_{3}^{f}}) \prod_{n=1}^{N} \phi(\frac{\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{c_{n}^{U}} - \mu_{t}^{c_{n}^{U}}(\boldsymbol{x}, \boldsymbol{\theta})}{\sigma_{t}^{c_{n}^{U}}(\boldsymbol{x}, \boldsymbol{\theta})})/(s_{3}^{f} \sigma_{t}^{c_{n}^{U}}(\boldsymbol{x}, \boldsymbol{\theta})). \end{split}$$

Similarly, for the lower-level density, Bayes theorem transforms

$$p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g}, \mathcal{D}_{t}^{+}) = \frac{p(\boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g} \mid \boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g}, \mathcal{D}_{t}^{+})p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t}^{+})}{p(\boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g} \mid \mathcal{D}_{t}^{+})}$$

Here again, the density of the first dimension of $p(\boldsymbol{h}_{(\boldsymbol{x}^*,\boldsymbol{\theta})}^g \mid \boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^g, \mathcal{D}_t^+)$ is $\mathcal{N}(m_1^g, \{s_1^g\}^2)$ shown in Appendix A.1 and from the second to the (M+1)-th dimension is $\mathcal{N}(m^{c_m^L}, \{s^{c_m^L}\}^2)$ where

$$\begin{split} m^{c_m^L} &= \mu_t^{c_m^L}(x^*, \theta) + \frac{\text{Cov}_t^{c_m^L}((x^*, \theta), (x, \theta))}{\left\{\sigma_t^{c_m^L}(x, \theta)\right\}^2} (y^{c_m^L}(x, \theta) - \mu_t^{c_m^L}(x, \theta)), \\ &\{s^{c_m^L}\}^2 = \sigma_t^{c_m^L}(x^*, \theta) - \frac{\left\{\text{Cov}_t^{c_m^L}((x^*, \theta), (x, \theta))\right\}^2}{\left\{\sigma_t^{c_m^L}(x, \theta)\right\}^2}. \end{split}$$

As a result

$$\begin{split} p(\boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g} \mid \boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g}, \mathcal{D}_{t}^{+}) &= 1 - \left(1 - \Phi\left(\frac{g^{*} - m_{1}^{g}}{s_{1}^{g}}\right)\right) \prod_{m=1}^{M} \left(1 - \Phi\left(\frac{0 - m_{m}^{c_{m}^{L}}}{s_{m}^{c_{m}^{L}}}\right)\right) \\ p(\boldsymbol{h}_{(\boldsymbol{x}^{*},\boldsymbol{\theta})}^{g} \in \bar{\mathcal{A}}^{g} \mid \mathcal{D}_{t}^{+}) &= 1 - \left(1 - \Phi\left(\frac{g^{*} - m_{2}^{g}}{s_{2}^{g}}\right)\right) \prod_{m=1}^{M} \left(1 - \Phi\left(\frac{0 - \mu_{t}^{c_{m}^{L}}(\boldsymbol{x}^{*}, \boldsymbol{\theta})}{\sigma_{t}^{c_{m}^{L}}(\boldsymbol{x}^{*}, \boldsymbol{\theta})}\right)\right) \\ p(\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} \mid \mathcal{D}_{t}^{+}) &= \phi\left(\frac{\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{g} - m_{3}^{g}}{s_{3}^{g}}\right) \prod_{m=1}^{M} \phi\left(\frac{\boldsymbol{y}_{(\boldsymbol{x},\boldsymbol{\theta})}^{c_{m}} - \mu_{t}^{c_{m}^{L}}(\boldsymbol{x}, \boldsymbol{\theta})}{\sigma_{t}^{c_{m}^{L}}(\boldsymbol{x}, \boldsymbol{\theta})}\right) / (s_{3}^{g} \sigma_{t}^{c_{m}^{L}}(\boldsymbol{x}, \boldsymbol{\theta})) \end{split}$$

E SUPPLEMENTARY FOR EXPERIMENTS

E.1 Other Details of Experimental Settings

We used the SingleTaskGP model of BoTorch (Balandat et al., 2020) to define the GPs. The output of each benchmark function are transformed by signed log1p function $\operatorname{sign}(y) \log(1 + |y|)$ except for BraninHoo and Goldstein-price for which the transformation shown by (Picheny et al., 2013) was used. For benchmark functions, the input space is scaled to $[0,1]^{d_X+d_\Theta}$ from the original input domain.

E.2 Larger Noise Setting

Figure 6 shows results with a stronger noise setting (the noise standard deviation is set as 10^{-1}). We do not see large difference for the relative performance among compared methods compared with the small noise setting shown in Fig. 3.

E.3 Additional Results on Decoupled Setting

Figure 7 shows regret in decoupled setting for all different length scale settings of the GP prior, which generates true objectives. We obviously see that BLJES was superior or comparable to BILBO.

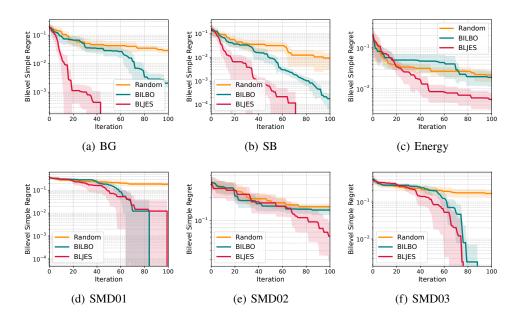


Figure 6: Regret comparison with 10^{-1} noise standard deviation.

E.4 Additional Results on Effect of the Number of Samplings

Figure 8 shows results of BLJES for different K. We do not see particularly large differences among different K settings in these benchmarks.

E.5 Continuous Domain

Figure 9 shows the regret in the case of \mathcal{X} and Θ are the continuous space. We employed gradient based optimizers for both of the bilevel problem defined by sample paths and the acquisition function maximization (gradient of a bilevel problem is discussed in Appendix B). Here, BILBO is not performed because Chew et al. (2025) only discuss the finite domain. We see that BLJES efficiently decreases the regret even in the continuous space. Only in SMD02, BLJES was not efficient compared with the random selection.

E.6 Constraint Problems

For empirical evaluation, we employed problems from the bilevel optimization benchmark (Sinha et al., 2014), denoted as SDM09 ($d_X = 2$, $d_{\Theta} = 2$, N = 1, M = 1), 10 ($d_X = 2$, $d_{\Theta} = 2$, N = 2, M = 1), 11 ($d_X = 2$, $d_{\Theta} = 2$, N = 1, M = 1), and 12 ($d_X = 2$, $d_{\Theta} = 2$, N = 3, M = 2). The number of grid points in each dimension is 10 (10^4 points). The evaluation metric is

$$\min_{i \in [n_0 + t]} \max_{h \in \{f, g, c_1^U, \dots, c_N^U, c_1^L, \dots, c_M^L\}} r_h(\mathbf{x}_i, \boldsymbol{\theta}_i)$$

where r_h become r_f and r_g shown in section 5 if h = f or g, and

$$r_c(\mathbf{x}_i, \mathbf{\theta}_i) = \max(0, -c(\mathbf{x}_i, \mathbf{\theta}_i)) / \max_{\mathbf{x}, \mathbf{\theta}} (\max(0, -c(\mathbf{x}, \mathbf{\theta}))), c \in \{c_1^U, \dots, c_N^U, c_1^L, \dots, c_M^L\}$$

if $h \in \{c_1^U, \dots, c_N^U, c_1^L, \dots, c_M^L\}$. The other settings are same as described in the beginning of section 5.

The results are in Fig. 10. For SMD09 and SMD11, BLJES shows faster decrease of the regret. For SMD12, BLJES and BILBO are comparable and both of them are much better than Random. For SMD10, BLJES rapidly decreased the regret, while BLJES also quickly decreased the regret (the difference is in small scale values).

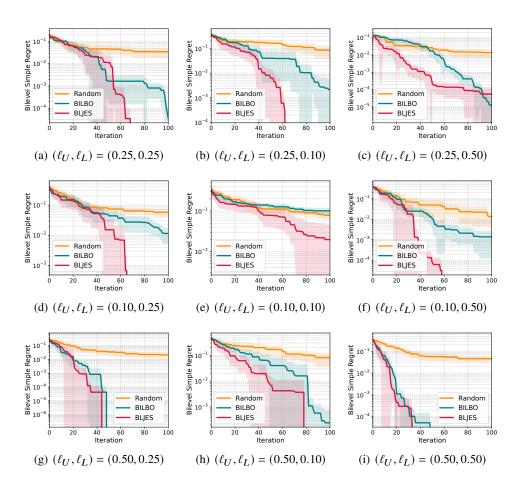


Figure 7: Regret comparison on functions from the GP prior under decoupled setting.

F LLM Usage

In this manuscript, LLM was only used to polish writing.

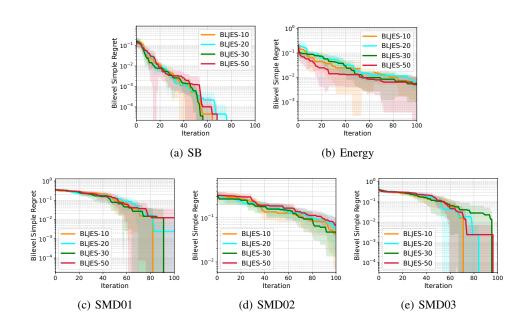


Figure 8: BLJES with different *K*.

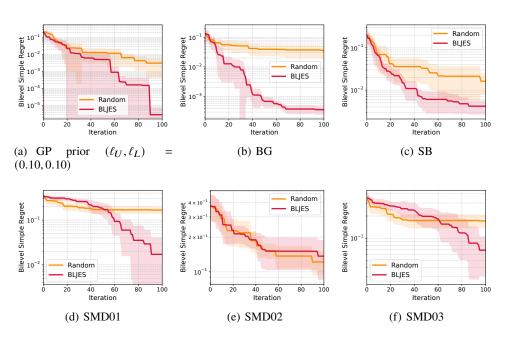


Figure 9: Regret comparison on continuous input domain.

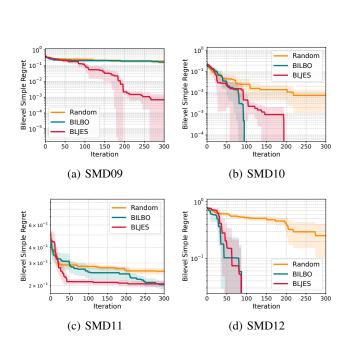


Figure 10: Regret comparison in constraint problems.