

An Information-Theoretic Framework for Unifying Active Learning Problems

Quoc Phong Nguyen,¹ Bryan Kian Hsiang Low,¹ Patrick Jaillet²

¹Dept. of Computer Science, National University of Singapore, Republic of Singapore

²Dept. of Electrical Engineering and Computer Science, MIT, USA
{qphong, lowkh}@comp.nus.edu.sg, jaillet@mit.edu

Abstract

This paper presents an information-theoretic framework for unifying active learning problems: *level set estimation* (LSE), *Bayesian optimization* (BO), and their generalized variant. We first introduce a novel active learning criterion that subsumes an existing LSE algorithm and achieves state-of-the-art performance in LSE problems with a continuous input domain. Then, by exploiting the relationship between LSE and BO, we design a competitive information-theoretic acquisition function for BO that has interesting connections to upper confidence bound and *max-value entropy search* (MES). The latter connection reveals a drawback of MES which has important implications on not only MES but also on other MES-based acquisition functions. Finally, our unifying information-theoretic framework can be applied to solve a generalized problem of LSE and BO involving multiple level sets in a data-efficient manner. We empirically evaluate the performance of our proposed algorithms using synthetic benchmark functions, a real-world dataset, and in hyperparameter tuning of machine learning models.

1 Introduction

Level set estimation (LSE) is about determining a level set of an unknown function or, alternatively, a *superlevel set* of the function (i.e., a region of inputs where the function values are of at least a known threshold) given a finite budget of expensive (possibly noisy) function evaluations (Gotovos et al. 2013). It has important applications in environmental sensing/monitoring where the unknown function is a (spatial) field of some quantity of interest (e.g., pH, temperature, and solar radiation) (Galland, Réfrégier, and Germain 2004). On the other hand, *Bayesian optimization* (BO) has gained significant recognition in science and engineering fields (Brochu, Cora, and de Freitas 2010; Calandra et al. 2014; Krause and Ong 2011; Shahriari et al. 2015; Snoek, Larochelle, and Adams 2012) for its effectiveness in optimizing a black-box objective function (i.e., without a closed-form expression/derivative) using a finite budget of expensive (possibly noisy) function evaluations. At first glance, one may straightforwardly regard BO as LSE by setting the threshold as the maximum value of the objective function, i.e., the superlevel set is reduced to a set of maximizers.

However, the *unknown* maximum value in BO does not satisfy the requirement of a *known* threshold in LSE. This poses the challenge of whether it is possible to develop a framework to unify LSE and BO. Though the work of Bogunovic et al. (2016) has developed such a unified approach called truncated variance reduction, it is demonstrated mainly on problems with a discrete input domain and requires enumerating over all inputs in a set of “unclassified points”, which can be prohibitively large in practice (or infinite when the input domain is continuous and not discretized). In contrast, our work here proposes the *first information-theoretic* framework for unifying both LSE and BO that can empirically outperform the state-of-the-art LSE criteria and scale to real-world problems with a continuous input domain.

To shed light on the connection between LSE and BO, we propose to view BO as an active learning problem that involves actively estimating the superlevel set of the objective function with respect to an estimate of its maximum value; such a problem reduces to LSE when the maximum value is known instead. Improving the estimation of the superlevel set in turn refines the estimate of the maximum value. As the estimate approaches the true maximum value, the superlevel set becomes a set of the maximizers of the objective function. Unfortunately, existing LSE criteria cannot be directly applied to BO since they either impose a noiseless assumption (Low et al. 2012) or cannot handle an unknown threshold (Bryan et al. 2006).

A key contribution of our work here therefore lies in introducing a novel information-theoretic active learning criterion for LSE (Sec. 3) that can be exploited for designing a new acquisition function for BO with interesting connections to *upper confidence bound* (UCB) and *max-value entropy search* (MES) (Sec. 4). The latter connection reveals a drawback of MES (Remark 3), which has important implications on not only MES, but also on other MES-based acquisition functions such as those handling multiple objectives (Belakaria, Deshwal, and Doppa 2019; Suzuki et al. 2020) or fidelities (Takeno et al. 2020).

The other main contribution of our work is to show how our proposed unifying information-theoretic framework can be applied to solve a generalized problem of LSE and BO involving multiple level sets/thresholds in a data-efficient manner. This problem, namely *implicit LSE* (Sec. 5), is about identifying a region of inputs whose function values

differ from the (unknown) maximum value by at most a specified *tolerance*. It is motivated from the estimation of *hotspots* in environmental fields, which correspond to regions with large field measurements (Gotovos et al. 2013). In summary, the specific contributions of our work include:

- A novel information-theoretic active learning criterion for LSE problems with a continuous input domain (Sec. 3), which subsumes an existing LSE criterion (Low et al. 2012) and empirically outperforms state-of-the-art LSE criteria (Bryan et al. 2006; Low et al. 2012) on synthetic benchmark functions and a real-world dataset (Sec. 6.1);
- A new information-theoretic acquisition function for BO problems with interesting connections to UCB and MES; the latter connection reveals a drawback of using MES (Sec. 4). We empirically evaluate the performance of our proposed BO algorithm using several synthetic benchmark functions, a real-world dataset, and in hyperparameter tuning of a logistic regression model and a convolutional neural network for image classification with MNIST and CIFAR-10 datasets (Sec. 6.2);
- Applying our unifying information-theoretic framework to solve the implicit LSE problem in a data-efficient manner (Sec. 5).

2 Gaussian Process (GP)

Let the unknown objective function be denoted as $f : \mathcal{X} \rightarrow \mathbb{R}$ over a bounded input domain $\mathcal{X} \subset \mathbb{R}^d$. An LSE/BO algorithm repeatedly selects an input query $\mathbf{x} \in \mathcal{X}$ for evaluating f to obtain a noisy observation $y_{\mathbf{x}} \triangleq f(\mathbf{x}) + \epsilon_{\mathbf{x}}$ of its function value $f(\mathbf{x})$ corrupted by an additive Gaussian noise $\epsilon_{\mathbf{x}} \sim \mathcal{N}(0, \sigma_n^2)$ with noise variance σ_n^2 . Since it is expensive to evaluate f , the goal of LSE (BO) is to strategically select input queries for finding the level/superlevel set (global maximizer(s)) as rapidly as possible. To achieve this, we model f using a GP: Let $\{f(\mathbf{x}')\}_{\mathbf{x}' \in \mathcal{X}}$ denote a GP, i.e., every finite subset of $\{f(\mathbf{x}')\}_{\mathbf{x}' \in \mathcal{X}}$ follows a multivariate Gaussian distribution (Rasmussen and Williams 2006). Then, the GP is fully specified by its *prior* mean $\mathbb{E}[f(\mathbf{x}')]$ and covariance $k_{\mathbf{x}'\mathbf{x}''} \triangleq \text{cov}[f(\mathbf{x}'), f(\mathbf{x}'')]$ for all $\mathbf{x}', \mathbf{x}'' \in \mathcal{X}$; the latter can be defined by, for example, the widely-used squared exponential kernel $k_{\mathbf{x}'\mathbf{x}''} \triangleq \sigma_s^2 \exp(-0.5(\mathbf{x}' - \mathbf{x}'')^\top \Lambda^{-2}(\mathbf{x}' - \mathbf{x}''))$ where $\Lambda \triangleq \text{diag}(\ell_1, \dots, \ell_d)$ and σ_s^2 are its length-scale and signal variance hyperparameters, respectively. For notational simplicity (and w.l.o.g.), the prior mean is assumed to be zero. Given a column vector $\mathbf{y}_{\mathcal{D}} \triangleq (y_{\mathbf{x}'})_{\mathbf{x}' \in \mathcal{D}}$ of noisy observations from evaluating f at a set \mathcal{D} of input queries selected in previous iterations, the GP posterior belief of the function value at any input query \mathbf{x} is a Gaussian $p(f(\mathbf{x})|\mathbf{y}_{\mathcal{D}}) = \mathcal{N}(f(\mathbf{x})|\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2)$ with the following *posterior* mean $\mu_{\mathbf{x}}$ and variance $\sigma_{\mathbf{x}}^2$:

$$\begin{aligned} \mu_{\mathbf{x}} &\triangleq \mathbf{K}_{\mathbf{x}\mathcal{D}}(\mathbf{K}_{\mathcal{D}\mathcal{D}} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}_{\mathcal{D}} \\ \sigma_{\mathbf{x}}^2 &\triangleq k_{\mathbf{x}\mathbf{x}} - \mathbf{K}_{\mathbf{x}\mathcal{D}}(\mathbf{K}_{\mathcal{D}\mathcal{D}} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{K}_{\mathcal{D}\mathbf{x}} \end{aligned} \quad (1)$$

where $\mathbf{K}_{\mathbf{x}\mathcal{D}} \triangleq (k_{\mathbf{x}\mathbf{x}'})_{\mathbf{x}' \in \mathcal{D}}$, $\mathbf{K}_{\mathcal{D}\mathcal{D}} \triangleq (k_{\mathbf{x}'\mathbf{x}''})_{\mathbf{x}', \mathbf{x}'' \in \mathcal{D}}$, $\mathbf{K}_{\mathcal{D}\mathbf{x}} \triangleq \mathbf{K}_{\mathbf{x}\mathcal{D}}^\top$, and \mathbf{I} is an identity matrix. Then, $p(y_{\mathbf{x}}|\mathbf{y}_{\mathcal{D}}) = \mathcal{N}(y_{\mathbf{x}}|\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2 + \sigma_n^2)$.

3 Binary Entropy Search (BES) for Level Set Estimation (LSE)

An LSE algorithm repeatedly selects the next input query $\mathbf{x} \in \mathcal{X}$ for evaluating f to maximize some active learning criterion based on the GP posterior belief of f given the observations $\mathbf{y}_{\mathcal{D}}$ obtained in previous iterations such that the superlevel set $\mathcal{X}_{f_o}^+ \triangleq \{\mathbf{x}' \in \mathcal{X} | f(\mathbf{x}') \geq f_o\}$ of f w.r.t. a given threshold f_o can be found as rapidly as possible.

In this section, we propose an information-theoretic active learning criterion for LSE which measures the information gain on the superlevel set $\mathcal{X}_{f_o}^+$ from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$. Let $\gamma_{\mathbf{x}}^\circ$ denote an indicator variable of label -1 if $\mathbf{x} \in \mathcal{X}_{f_o}^+$ (i.e., superlevel set), and label 1 otherwise (i.e., \mathbf{x} is in sublevel set $\mathcal{X}_{f_o}^- \triangleq \{\mathbf{x}' \in \mathcal{X} | f(\mathbf{x}') < f_o\}$). We can view $\gamma_{\mathbf{x}}^\circ$ as a class label of \mathbf{x} and LSE as a binary classification problem that classifies whether each $\mathbf{x} \in \mathcal{X}$ is in the superlevel set $\mathcal{X}_{f_o}^+$ or the sublevel set $\mathcal{X}_{f_o}^-$. Let $\gamma_{\mathcal{X}}^\circ \triangleq (\gamma_{\mathbf{x}'}^\circ)_{\mathbf{x}' \in \mathcal{X}}$. The active learning criterion can therefore be measured as the mutual information $I(y_{\mathbf{x}}; \gamma_{\mathcal{X}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o)$ which cannot be evaluated tractably with a continuous \mathcal{X} . So, we simplify it to the information gain on class label $\gamma_{\mathbf{x}}^\circ$ from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$:

$$\begin{aligned} \alpha_{\text{BES}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}}) &\triangleq I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o) \\ &= H(p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o)) - \mathbb{E}_{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_o)} [H(p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}} \cup \{\mathbf{x}\}, f_o))] \end{aligned} \quad (2)$$

where the *prior entropy* of $\gamma_{\mathbf{x}}^\circ$ is defined as

$$H(p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o)) \triangleq - \sum_{\gamma_{\mathbf{x}}^\circ} p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o) \log p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o)$$

and the *posterior entropy* $H(p(\gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}} \cup \{\mathbf{x}\}, f_o))$ of $\gamma_{\mathbf{x}}^\circ$ given $y_{\mathbf{x}}$ is defined in a similar manner. Since $\gamma_{\mathbf{x}}^\circ$ is binary, our active learning criterion $\alpha_{\text{BES}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}})$ is named *binary entropy search* (BES). Note that BES (2) can be interpreted as the expected reduction in the uncertainty (entropy) of $\gamma_{\mathbf{x}}^\circ$ from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$. Though replacing $\gamma_{\mathcal{X}}^\circ$ with $\gamma_{\mathbf{x}}^\circ$ appears to be a simplification, BES demonstrates state-of-the-art performance in our experiments (Sec. 6.1). Such a simplification is also commonly adopted by existing acquisition functions for BO (e.g., (Suzuki et al. 2020; Wang and Jegelka 2017)). BES (2) can be evaluated as follows:

$$\begin{aligned} \alpha_{\text{BES}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}}) &\triangleq I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^\circ | \mathbf{y}_{\mathcal{D}}, f_o) \\ &= \mathbb{E}_{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})} \left[\sum_{\gamma_{\mathbf{x}}^\circ} \Psi(\gamma_{\mathbf{x}}^\circ g_{\mathbf{x}}(y_{\mathbf{x}}, f_o)) \log \frac{\Psi(\gamma_{\mathbf{x}}^\circ g_{\mathbf{x}}(y_{\mathbf{x}}, f_o))}{\Psi(\gamma_{\mathbf{x}}^\circ h_{\mathbf{x}}(f_o))} \right] \end{aligned} \quad (3)$$

where $g_{\mathbf{x}}(y_{\mathbf{x}}, f_o) \triangleq (\sigma_+^2 f_o - \sigma_n^2 \mu_{\mathbf{x}} - \sigma_{\mathbf{x}}^2 y_{\mathbf{x}}) / (\sigma_{\mathbf{x}} \sigma_n \sigma_+)$, $h_{\mathbf{x}}(f_o) \triangleq (f_o - \mu_{\mathbf{x}}) / \sigma_{\mathbf{x}}$, and Ψ denotes the c.d.f. of the standard Gaussian distribution. Its derivation is shown in Appendix A. BES (3) can thus be optimized w.r.t. input query \mathbf{x} via stochastic gradient ascent.

Fig. 1a shows LSE with the threshold $f_o = 0$ being viewed as a binary classification problem that classifies whether each $\mathbf{x} \in [0, 10]$ is in $\mathcal{X}_{f_o}^+$ or $\mathcal{X}_{f_o}^-$ and the level set w.r.t. $f_o = 0$ is likely to be found on the decision boundary

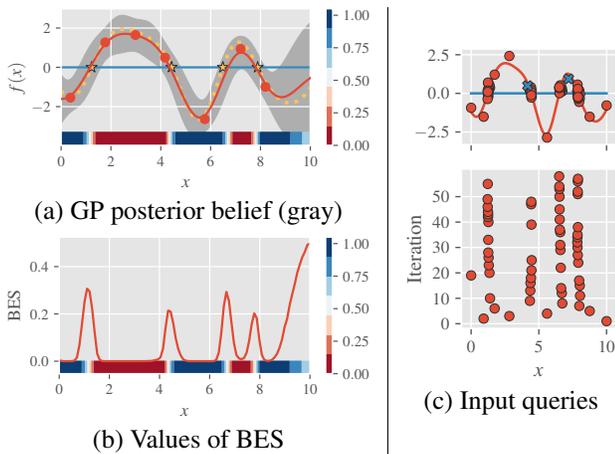


Figure 1: LSE with the threshold $f_o = 0$ as a binary classification problem that classifies if each $\mathbf{x} \in [0, 10]$ is in $\mathcal{X}_{f_o}^+$ or $\mathcal{X}_{f_o}^-$; $p(\mathbf{x} \in \mathcal{X}_{f_o}^+ | \mathbf{y}_{\mathcal{D}})$ is shown on x -axes of left plots. (a) The objective function f , level set, 7 observations, GP posterior mean (1), and $f_o = 0$ are plotted as a yellow dotted line, yellow stars, red circles, a red curve, and a blue line, respectively. (b) Plot of values of BES based on GP posterior belief in Fig. 1a. The bottom plot in (c): input queries (red circles) selected by BES w.r.t. iteration no.; the top plot in (c): $f_o = 0$, 2 prior observations, and observations plotted as a blue line, blue crosses, and red circles, respectively.

(i.e., white regions on the x -axis). Fig. 1b shows large values of BES on the decision boundary that is likely to contain the level set, which is desirable. Fig. 1c shows BES using about 10 observations to explore and find roughly the level set w.r.t. $f_o = 0$. Then, BES exploits by distributing its observations on the decision boundary (i.e., level set).

Remark 1 (Special case of BES) When the observation is noiseless (i.e., $y_{\mathbf{x}} = f(\mathbf{x})$ or $\sigma_n^2 = 0$), $\gamma_{\mathbf{x}}^o$ is fully determined by the values of f_o and $y_{\mathbf{x}} = f(\mathbf{x})$. Then, $p(\gamma_{\mathbf{x}}^o | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, f_o)$ is either 0 or 1 and the posterior entropy term in (2) thus becomes 0. So, BES reduces to the prior entropy term in (2) and the resulting active learning algorithm: $\max_{\mathbf{x} \in \mathcal{X}} H(p(\gamma_{\mathbf{x}}^o | \mathbf{y}_{\mathcal{D}}, f_o))$ is called *entropy maximization* (EM), as proposed by Low et al. (2012). EM can therefore be viewed as a special case of BES due to its noiseless observations. In other words, BES subsumes EM.

4 BES for Maximum Value Prediction (BES-MP) in Bayesian Optimization (BO)

A BO algorithm repeatedly selects the next input query $\mathbf{x} \in \mathcal{X}$ for evaluating f to maximize some acquisition function based on the GP posterior belief of f given the observations $\mathbf{y}_{\mathcal{D}}$ obtained in previous iterations such that the maximizer(s) of f can be found as rapidly as possible.

Given an estimate f_* of the maximum value of f , the superlevel set $\mathcal{X}_{f_*}^+ \triangleq \{\mathbf{x} \in \mathcal{X} | f(\mathbf{x}) \geq f_*\}$ w.r.t. the threshold f_* can be regarded as a set of potential maximizers. So, BO can be viewed as an active learning problem that

involves actively estimating the superlevel set $\mathcal{X}_{f_*}^+$, which corresponds to an LSE problem. Therefore, we exploit our proposed BES criterion for LSE (Sec. 3) to design a new acquisition function for BO, specifically, the information gain on class label $\gamma_{\mathbf{x}}^*$ (i.e., indicator variable of label -1 if $\mathbf{x} \in \mathcal{X}_{f_*}^+$, and label 1 otherwise) from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$: $I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)$ which can be optimized via stochastic gradient ascent by replacing $\gamma_{\mathbf{x}}^o$ and f_o with $\gamma_{\mathbf{x}}^*$ and f_* in (3), respectively. However, since the maximum value of f is unknown, we estimate it with a set \mathcal{F}_* of samples of the maximum value of functions drawn from the GP posterior belief (1). These functions are drawn by applying the random Fourier feature approximation to GP (Rahimi and Recht 2008), which is widely used in existing information-theoretic acquisition functions (Hernández-Lobato, Hoffman, and Ghahramani 2014; Hoffman and Ghahramani 2015; Wang and Jegelka 2017). Then, we propose the acquisition function called *BES for maximum value prediction* (BES-MP) by averaging our BES criterion (for LSE) over the set \mathcal{F}_* of maximum value samples:

$$\alpha_{\text{BES-MP}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}}) \triangleq |\mathcal{F}_*|^{-1} \sum_{f_* \in \mathcal{F}_*} I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*). \quad (4)$$

At first glance, it may not seem straightforward to justify averaging BES over \mathcal{F}_* in (4). To do so, we have proven in Appendix B that the average of BES over \mathcal{F}_* (4) is in fact the mutual information between $y_{\mathbf{x}}$ and the jointly distributed random variables $(\gamma_{\mathbf{x}}^*, f_*)$ ¹

$$\alpha_{\text{BES-MP}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}}) = I(y_{\mathbf{x}}; (\gamma_{\mathbf{x}}^*, f_*) | \mathbf{y}_{\mathcal{D}}) \quad (5)$$

where we overload the notation f_* to denote a discrete uniform random variable on the support \mathcal{F}_* whose distribution approximates that of the unknown maximum value of f .

In the rest of this section, we will investigate the connections between BES-MP and existing acquisition functions: UCB (Srinivas et al. 2010) and MES (Wang and Jegelka 2017). Our result below reveals that UCB can, interestingly, be derived from BES-MP by choosing a deterministic estimate of the maximum value of f , as proven in Appendix C:

Theorem 1 (Connection to UCB) *Define acquisition function of UCB as $\alpha_{\text{UCB}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}}) \triangleq \mu_{\mathbf{x}} + \beta \sigma_{\mathbf{x}}$ ($\beta > 0$) and $\mathbf{x}_{\text{UCB}} \triangleq \arg \max_{\mathbf{x} \in \mathcal{X}} \alpha_{\text{UCB}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}})$. If observation $y_{\mathbf{x}}$ is noiseless and the estimate of maximum value of f is chosen deterministically: $f_* = \alpha_{\text{UCB}}(\mathbf{x}_{\text{UCB}}, \mathbf{y}_{\mathcal{D}})$, then BES-MP selects the same input queries as that selected by UCB.*

For noisy observation $y_{\mathbf{x}}$, though both BES-MP and MES employ a set \mathcal{F}_* of samples of the maximum value of f , BES-MP differs significantly from MES in both its interpretation and model of noisy observation, as explained in the two remarks below:

Remark 2 (Interpretation as information gain) BES-MP (5) can be interpreted as information gain on both the class

¹An alternative acquisition function would be the mutual information $I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}})$ where f_* is marginalized out. But, its empirical performance does not differ much from that of (4). So, we focus on (4) which can be seamlessly unified with BES for LSE.

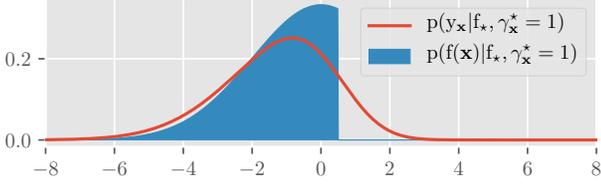


Figure 2: Plots of $p(y_{\mathbf{x}}|f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1)$ vs. a truncated Gaussian distribution of $p(f(\mathbf{x})|f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1)$.

label $\gamma_{\mathbf{x}}^{\star}$ and the threshold $f_{\star} \in \mathcal{F}_{\star}$ inducing the superlevel set $\mathcal{X}_{f_{\star}}^+$ (of potential maximizers) from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$. In contrast, MES measures the information gain on maximum value from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$. BES-MP (4) is also closely related to BES (2), thus allowing our unifying information-theoretic framework for BO and LSE to be established.

Remark 3 (Model of noisy observation $y_{\mathbf{x}}$) Another key distinction between BES-MP and MES lies in how they model $p(y_{\mathbf{x}}|y_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1)$. MES assumes that given $f_{\star} \geq f(\mathbf{x})$, the observation $y_{\mathbf{x}}$ at an input query \mathbf{x} must be at most f_{\star} (Wang and Jegelka 2017), which leads to an (upper-tail) truncated Gaussian distribution of $p(y_{\mathbf{x}}|y_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1)$ and its closed-form expression. However, due to noise $\epsilon_{\mathbf{x}}$, $y_{\mathbf{x}} = f(\mathbf{x}) + \epsilon_{\mathbf{x}}$ can be larger than f_{\star} even though $f_{\star} \geq f(\mathbf{x})$, as shown in Fig. 2. This issue can be interpreted as MES assuming to observe the noiseless $f(\mathbf{x})$ when in fact, only the noisy $y_{\mathbf{x}}$ is observed, which implies that MES overestimates the information gain on the maximum value from observing a noisy $y_{\mathbf{x}}$. This overestimation is significant when the noise variance σ_n^2 is large relative to the posterior variance $\sigma_{\mathbf{x}}^2$ (1) of $f(\mathbf{x})$. Such an issue also plagues the other MES-based acquisition functions such as those handling multiple objectives (Belakaria, Deshwal, and Doppa 2019; Suzuki et al. 2020) or fidelities (Takeno et al. 2020). On the other hand, BES-MP models $p(y_{\mathbf{x}}|y_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star})$ accurately, which may suggest an improvement to these other MES-based acquisition functions to be considered for future work.

5 Implicit Level Set Estimation (LSE)

Implicit LSE is about finding the superlevel set w.r.t. an unknown threshold that differs from the maximum value of f by a specified *tolerance*. It is motivated from the estimation of *hotspots* (i.e., superlevel sets) in environmental fields, which are regions of locations (i.e., inputs) whose field measurements (i.e., function values) are of at least a threshold. Since such measurements may vary throughout the year, it is desirable to define the threshold based on the (unknown) maximum value of the environmental field, which explains the term of *implicit level set*. For example, farmers are interested to identify the regions of their farms with high (or low) phosphorus level. Recall that LSE aims to find the superlevel set w.r.t. a known threshold while BO aims to find the maximizer(s) of the objective function, i.e., the superlevel set w.r.t. the unknown maximum value. Therefore, our LSE and BO algorithms cannot be directly applied to solve the implicit LSE problem.

A variant of an implicit LSE problem with a discrete input domain has been introduced in (Gotovos et al. 2013) where the threshold is expressed as a percentage of $\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. However, in this paper, we prefer our above definition as it accounts meaningfully for negative function values. Note that existing works only consider problems with a discrete input domain (Gotovos et al. 2013) while our work here addresses problems with a continuous input domain such as those in our experiments.

Let $\alpha \geq 0$ be the specified tolerance. The threshold in implicit LSE is then $\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) - \alpha$ which is not known due to the unknown maximum value: $\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. So, the implicit LSE problem is about finding the superlevel set w.r.t. $\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) - \alpha$. It is a generalized variant of BO and LSE as it reduces to BO when $\alpha = 0$ and to LSE when the maximum value of f is known.

Following the design of BES-MP in Sec. 4, one may be tempted to solve the implicit LSE problem by averaging BES over the set $\mathcal{F}_{\alpha} \triangleq \{f_{\star} - \alpha | f_{\star} \in \mathcal{F}_{\star}\}$ where \mathcal{F}_{\star} is a set of samples of the maximum value of f defined in Sec. 4 previously; $f_{\alpha} \in \mathcal{F}_{\alpha}$ is then an estimate of the unknown threshold in implicit LSE. Define the superlevel set $\mathcal{X}_{f_{\alpha}}^+ \triangleq \{\mathbf{x} \in \mathcal{X} | f(\mathbf{x}) \geq f_{\alpha}\}$ w.r.t. f_{α} . Let $\gamma_{\mathbf{x}}^{\alpha}$ denote an indicator variable of label -1 if $\mathbf{x} \in \mathcal{X}_{f_{\alpha}}^+$, and label 1 otherwise. Similar to (4), the active learning criterion of BES-MP for implicit LSE can be written as

$$(1/|\mathcal{F}_{\alpha}|) \sum_{f_{\alpha} \in \mathcal{F}_{\alpha}} I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\alpha} | y_{\mathcal{D}}, f_{\alpha}). \quad (6)$$

Like (5), (6) can also be expressed as $I(y_{\mathbf{x}}; (\gamma_{\mathbf{x}}^{\alpha}, f_{\alpha}) | y_{\mathcal{D}})$ which can be interpreted as the information gain on both the class label $\gamma_{\mathbf{x}}^{\alpha}$ and the threshold $f_{\alpha} \in \mathcal{F}_{\alpha}$ inducing the superlevel set $\mathcal{X}_{f_{\alpha}}^+$ from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$.² We can optimize (6) in the same manner as (4).

Unfortunately, the above BES-MP only actively estimates the decision boundaries between $\mathcal{X}_{f_{\alpha}}^+$ and $\mathcal{X}_{f_{\alpha}}^-$ for $f_{\alpha} \in \mathcal{F}_{\alpha}$. Since these decision boundaries can be far from the maximizer(s) (e.g., when α is large), it is unlikely that BES-MP queries at the maximizer(s), hence yielding poor estimates f_{\star} of the maximum value. For example, Fig. 3a shows that BES-MP has only 1 input query near to the maximizer of f . The poor estimates f_{\star} entail poor estimates $f_{\alpha} = f_{\star} - \alpha$ (i.e., dashed blue lines in Fig. 3a) and hence the poor performance of BES-MP in implicit LSE.

To improve the performance of BES-MP in implicit LSE, we consider a generalization of LSE to the *k-level set estimation* (*k-LSE*) problem (i.e., with multiple thresholds). It is an active learning problem that involves actively estimating the k level sets where the threshold of the i -th level set is represented by b_i . Let $\mathbf{b} \triangleq (b_i)_{i=1}^k$ denote a vector of thresholds in ascending order, i.e., $b_i < b_j$ if $i < j$. The k -LSE is equivalent to a $(k+1)$ -class classification problem that classifies each $\mathbf{x} \in \mathcal{X}$ into $k+1$ classes. Let $\gamma_{\mathbf{x}}^k \in \{0, 1, \dots, k\}$ denote the class label of an input \mathbf{x} such that it is of label 0 if $f(\mathbf{x}) \in (-\infty, b_1)$, and label i if $f(\mathbf{x}) \in [b_i, b_{i+1})$

²We also overload the notation f_{α} to denote a discrete uniform random variable on the support \mathcal{F}_{α} whose distribution approximates that of the unknown threshold in implicit LSE.

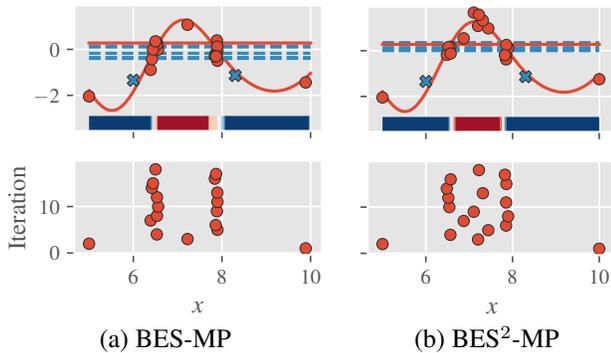


Figure 3: Input queries of (a) BES-MP and (b) BES²-MP in an implicit LSE problem. The notations are the same as those in Fig. 1 except that the ground truth f_α is plotted as a solid red line and the 5 estimates of f_α given \mathcal{D} (i.e., after 20 queries) are plotted as dashed blue lines.

and $1 \leq i \leq k$ where $b_{k+1} \triangleq \infty$. Similar to the design of BES, we propose an active learning criterion for k -LSE called BES^k that measures the information gain on class label $\gamma_{\mathbf{x}}^k$ from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$:

$$\alpha_{\text{BES}^k}(\mathbf{x}, \mathcal{Y}_{\mathcal{D}}) \triangleq I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^k | \mathcal{Y}_{\mathcal{D}}, \mathbf{b})$$

which can be expressed in a form that can be optimized via stochastic gradient ascent (Appendix D).

Implicit LSE can be viewed as a k -LSE problem such that the vector of thresholds is unknown (due to the unknown maximum value of f). So, we can exploit our BES^k criterion for k -LSE to design an active learning criterion for implicit LSE called BES²-MP (i.e., $k = 2$) by averaging BES^k over a set \mathcal{B} of estimates $\mathbf{b} = (f_* - \alpha, f_*)^\top$ for $f_* \in \mathcal{F}_*$:

$$\alpha_{\text{BES}^2\text{-MP}}(\mathbf{x}, \mathcal{Y}_{\mathcal{D}}) \triangleq (1/|\mathcal{B}|) \sum_{\mathbf{b} \in \mathcal{B}} I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^k | \mathcal{Y}_{\mathcal{D}}, \mathbf{b}).$$

Similar to BES-MP (6), BES²-MP can also be expressed as $I(y_{\mathbf{x}}; (\gamma_{\mathbf{x}}^k, \mathbf{b}) | \mathcal{Y}_{\mathcal{D}})$ which can be interpreted as the information gain on both the class label $\gamma_{\mathbf{x}}^k$ and the threshold vector \mathbf{b} inducing the 2 level sets from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$.³ Fig. 3b shows that BES²-MP uses several input queries to determine the maximum value of f but BES-MP (Fig. 3a) does not. As a result, BES²-MP can estimate $f_* - \alpha$ (i.e., f_α) more accurately than BES-MP, which can be observed from Fig. 3 by comparing the dashed blue lines representing f_α samples with the solid red line representing the ground truth threshold.

Remark 4 (A unifying framework) We introduce a unifying framework for LSE, BO, and implicit LSE problems by interpreting our proposed active learning criteria or acquisition function as information gain on the class label and the threshold vector \mathbf{b} of length k from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$. By setting $k = 1$, our unifying framework encompasses BES for LSE when the threshold is

³We also overload the notation \mathbf{b} to denote a discrete uniform random variable on the support \mathcal{B} whose distribution approximates that of the vector of unknown thresholds.

known (Sec. 3) and BES-MP for BO when the threshold is unknown (Sec. 4). By setting $k = 2$, our unifying framework encompasses BES²-MP for implicit LSE when the threshold vector is unknown.

6 Experiments and Discussion

This section empirically evaluates the performance of our proposed LSE (Sec. 6.1), BO (Sec. 6.2), and implicit LSE (Sec. 6.3) algorithms against that of state-of-the-art methods using synthetic benchmark functions, a real-world dataset, and in hyperparameter tuning of machine learning models. The code is available at <https://github.com/qphong/bes-mp>.

6.1 Level Set Estimation (LSE)

In this subsection, we empirically compare the performance of BES against that of the state-of-the-art EM (Low et al. 2012) (Remark 1) and *straddle* (STRDL) heuristic (Bryan et al. 2006) in the LSE problem. The methods of Bogunovic et al. (2016) and Gotovos et al. (2013) are demonstrated mainly on problems with a discrete input domain and hence not directly applicable to our experiments with a continuous input domain. Furthermore, STRDL is empirically shown to achieve comparable performance to these methods. So, STRDL is chosen as a direct competitor with BES while other methods (Bogunovic et al. 2016; Gotovos et al. 2013) are not empirically compared here. Since LSE is a binary classification problem (see Sec. 3) in a continuous domain \mathcal{X} , we use the log loss as the performance metric:

$$-(1/|\mathcal{X}'|) \sum_{\mathbf{x} \in \mathcal{X}'} \log p(c_{\mathbf{x}}^\circ (f(\mathbf{x}) - f_\circ) < 0 | \mathcal{Y}_{\mathcal{D}}) \quad (7)$$

where \mathcal{X}' is a set of 7000 uniformly sampled inputs from \mathcal{X} and $c_{\mathbf{x}}^\circ$ is an indicator variable of label -1 if $\mathbf{x} \in \mathcal{X}_{f_\circ}^+$, and label 1 otherwise. Each experiment is repeated 30 times to account for the randomness in the observation and the optimization. Results of the log 10 of the average of the log loss are presented.

As EM assumes noiseless observations (Remark 1), our experiments are performed with observations of both small ($\sigma_n^2 = 0.0001$) and large ($\sigma_n^2 = 0.09$) noise variances. The GP hyperparameters are learned using *maximum likelihood estimation* (MLE) (Rasmussen and Williams 2006). Regarding the synthetic functions, the function values are normalized and shifted to ensure a zero prior mean.

Results for the synthetic benchmark objective functions⁴ are shown in Figs. 4a to 4h. We can observe that (a) BES outperforms the other active learning criteria for both noise variance values, (b) EM outperforms STRDL when the noise variance is small ($\sigma_n^2 = 0.0001$), as shown in Figs. 4a, 4e, and 4g, and (c) the performance of EM deteriorates when the noise variance is large ($\sigma_n^2 = 0.09$) as it is outperformed by STRDL, as shown in Figs. 4b, 4d, and 4f. The last observation can be explained by the assumption of EM about noiseless observations (Remark 1).

Fig. 4i shows the results for an LSE problem on an estimated real-world phosphorus field (Webster and Oliver

⁴Details of the synthetic functions are available at <https://www.sfu.ca/~ssurjano/optimization.html>.

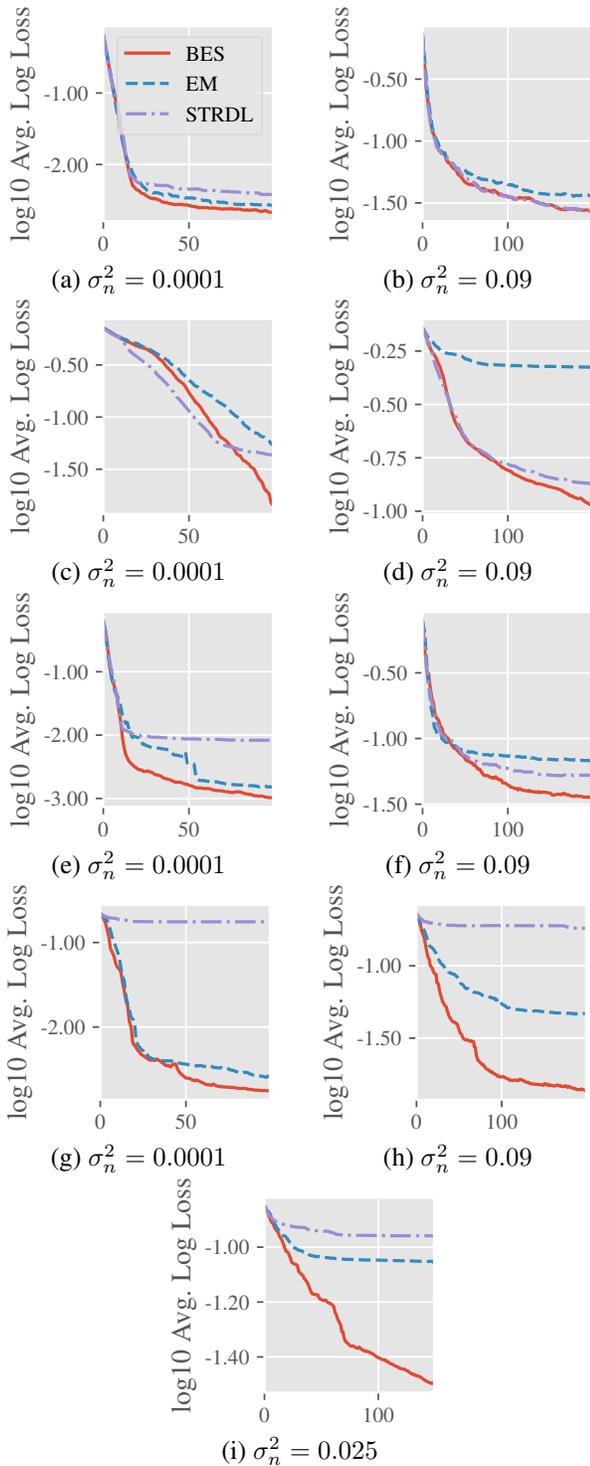


Figure 4: The log 10 of the average of the log loss for LSE experiments with synthetic functions: functions drawn from GP with (a-b) $l = 1/3$, (c-d) $l = 0.125$, (e-f) Branin, (g-h) Michaelwicz; and (i) an estimated phosphorus field.

2007). The noise variance is $\sigma_n^2 = 0.025$ which is learned from the dataset using MLE. It can be observed that BES

outperforms EM and STRDL significantly, while EM outperforms STRDL. The standard deviation (SD) of the log loss is shown in Table 1 in Appendix E.

6.2 Bayesian Optimization (BO)

This subsection evaluates the empirical performance of BES-MP against that of the existing acquisition functions: PES, MES, UCB, and EI in optimizing synthetic benchmark functions like Michaelwicz, Hartmann-3d, and Goldstein (the negative values of functions are used), and an estimated environmental field from the phosphorus dataset (see Sec 6.1). The noise variance in the experiments with the synthetic benchmark functions is 0.01. The GP hyperparameters are learned using MLE and $|\mathcal{F}_*|$ is set to 5.

We also use BO to tune the hyperparameters of 2 machine learning models. Firstly, we train a logistic regression model on the MNIST dataset which consists of 28×28 grayscale images of 10 handwritten digits. The hyperparameters include the L2 regularization weight (in $[10^{-6}, 1]$), the batch size (in $[20, 500]$), and the learning rate (in $[10^{-3}, 1]$). So, the input dimension of BO is 3. The objective function is the validation accuracy on a validation set of 14K images. Secondly, we train a CNN on the CIFAR-10 dataset which consists of 50K 32×32 color images in 10 classes. The CNN includes a convolutional layer followed by a dense layer. The hyperparameters include the batch size (in $[32, 512]$), the learning rate (in $[10^{-6}, 10^{-2}]$) and the learning rate decay (in $[10^{-7}, 10^{-3}]$) of the RMSprop optimization method, the convolutional filter size (in $[128, 256]$), and the number of hidden neurons in the dense layer (in $[64, 256]$). So, the input dimension of BO is 5. The objective function is the validation accuracy on a validation set of 10K images. We normalize the inputs in these experiments.

Following the work of Bogunovic et al. (2016), the performance metric is the regret of the best input query so far, i.e., $(\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})) - (\max_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}))$. The regret is averaged over 10 random runs to account for the randomness in the stochastic optimization and the noisy observation.

Fig. 5 shows that BES-MP outperforms the other acquisition functions in most of the experiments. In the other plots, BES-MP demonstrates a comparable performance to that of EI or PES. On the other hand, the performance of MES is not stable: for example, it does not perform well in Figs. 5a, 5b, and 5d. This can be explained by Remark 3. The SD of the regret is shown in Table 2 in Appendix E.

6.3 Implicit Level Set Estimation (LSE)

This subsection empirically illustrates the advantage of BES²-MP over BES-MP in implicit LSE problems which include several synthetic benchmark functions and an estimated phosphorus field (see Sec. 6.1). The tolerance α is specified as 0.2. The noise variance σ_n^2 in the observations of the synthetic functions is 0.0001. The GP hyperparameters are optimized using MLE. The number $|\mathcal{F}_*|$ of maximum value samples is 5. Similar to Sec. 6.1, the performance metric is the log loss. Unlike (7), since the threshold is unknown, it is marginalized out in the log loss expression:

$$-|\mathcal{X}'|^{-1} \sum_{\mathbf{x} \in \mathcal{X}'} \log(p(c_{\mathbf{x}}^\alpha (f(\mathbf{x}) - f_* + \alpha) < 0 | \mathcal{Y}_{\mathcal{D}}))$$

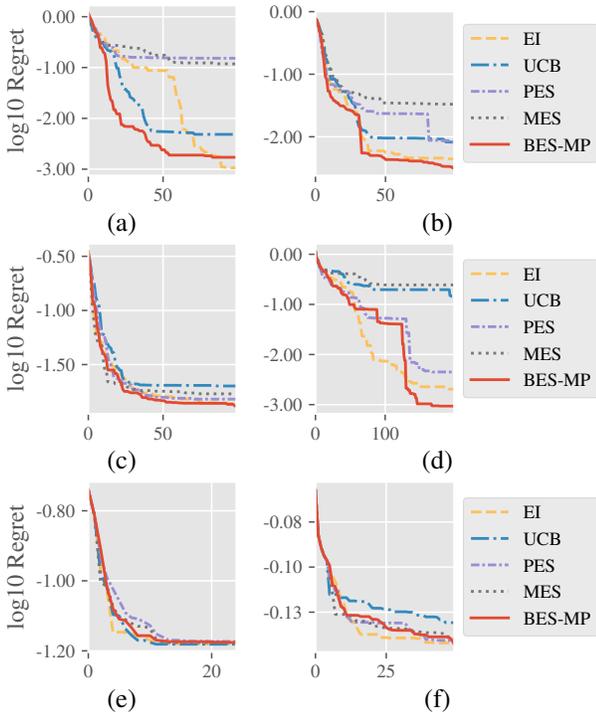


Figure 5: BO experiments with synthetic functions: (a) Michaelwicz, (b) Hartmann-3d, (c) Goldstein; real-world optimization problems: (d) an estimated phosphorus field; in hyperparameter tuning for training (e) a logistic regression model on MNIST, and (f) a CNN on CIFAR-10.

where f_* is marginalized: $p(c_x^\alpha (f(\mathbf{x}) - f_* + \alpha) < 0 | \mathbf{y}_D) = |\mathcal{F}_*|^{-1} \sum_{f_* \in \mathcal{F}_*} p(c_x^\alpha (f(\mathbf{x}) - f_* + \alpha) < 0 | \mathbf{y}_D, f_*)$; c_x^α is an indicator variable of label -1 if $\mathbf{x} \in \mathcal{X}_{f_*}^+$, and label 1 otherwise. Each experiment is repeated 30 times. Results of the log 10 of the average of the log loss are presented. We also reduce these implicit LSE problems to LSE problems by providing the threshold (i.e., $\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) - 0.2$) to the active learning criteria for LSE: BES, EM, and STRDL, and plotting their log losses. As the threshold is known, these methods serve as baselines that should outperform implicit LSE algorithms, i.e., BES²-MP and BES-MP.

Fig. 6 shows that BES²-MP outperforms BES-MP in all experiments, as expected from our discussion in Sec. 5. Besides, BES-MP does not converge in Figs. 6a and 6d as BES-MP does not gather observations to learn about the maximum value of f (Sec. 5). Regarding the baselines with known thresholds (i.e., active learning criteria for LSE: BES, EM, and STRDL), BES achieves the best performance. However, BES²-MP outperforms EM in Figs. 6a and 6e likely due to noisy observations. Surprisingly, even with known thresholds, STRDL is still outperformed by our BES²-MP and BES-MP in several experiments. It is different from the work of Gotovos et al. (2013) where baselines with known thresholds are empirically shown to outperform all methods with unknown thresholds. The SD of the log loss is shown in Table 3 in Appendix E.

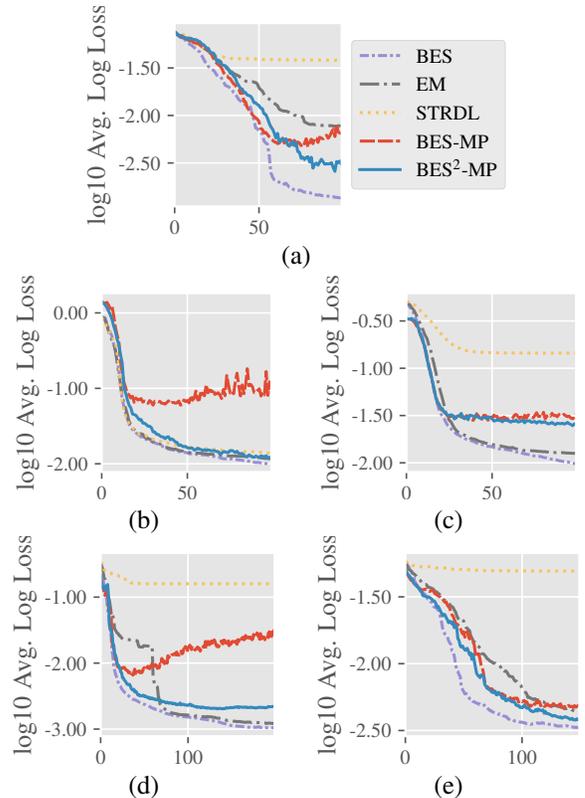


Figure 6: Implicit LSE experiments with synthetic functions: (a) a function sampled from GP with $l = 0.125$, (b) Branin, (c) Goldstein, (d) Hartmann-3d; and (e) an estimated phosphorus field.

7 Conclusion

This paper describes an information-theoretic framework for unifying the LSE, BO, and implicit LSE problems. We propose the first active learning criteria based on mutual information for LSE and implicit LSE problems, which yield the state-of-the-art empirical performance in estimating the level set of synthetic benchmark functions and an environmental field with a continuous input domain. By exploiting the relationship between LSE and BO, we design an information-theoretic acquisition function and study its connections to UCB and MES. It highlights a critical issue in modeling the noisy observation among the MES-based acquisition functions, which implies their overestimation of the information gain on the maximum value from the noisy observation. Our proposed acquisition function achieves a competitive performance in comparison with existing acquisition functions for BO in optimizing synthetic benchmark functions, an environmental field, and in hyperparameter tuning of logistic regression model and CNN. We will consider generalizing our framework to nonmyopic BO (Kharkovskii, Ling, and Low 2020; Ling, Low, and Jaillet 2016), batch BO (Daxberger and Low 2017), high-dimensional BO (Hoang, Hoang, and Low 2018), and multi-fidelity BO (Zhang, Dai, and Low 2019) settings.

Broader Impact

From our perspective, the societal benefits of the proposed framework outweigh its negative impact.

Our LSE and implicit LSE algorithms can be used for developing methods to monitor/locate hotspots (i.e., regions where environmental field measurements exceed a threshold) in an environmental field (e.g., over lakes and farms), which has potential applications in agriculture, aquaculture, and pollution control. While some people believe that this development can have a negative impact by reducing the salary of the related jobs, the long-term benefits are more significant. For example, high-yield and low-cost agriculture can help to sustain the growing population and reduce the food price, which benefits the whole society.

BO is well-known for a wide range of applications such as automated machine learning. With the comparison between our proposed BES-MP and other information-theoretic acquisition functions, other researchers can have a better understanding of BES-MP to employ/enhance it in their own research. Furthermore, our comparison can help engineers to understand and improve existing systems implemented with MES through the clarification of its drawback, for example, by correcting the approximation in Remark 3 if the observation noise is noticeable.

Acknowledgments. This research/project is supported by the National Research Foundation, Prime Minister’s Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) program, Singapore-MIT Alliance for Research and Technology (SMART) Future Urban Mobility (FM) IRG. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore.

References

- Belakaria, S.; Deshwal, A.; and Doppa, J. R. 2019. Max-value entropy search for multi-objective Bayesian optimization. In *Proc. NeurIPS*, 7825–7835.
- Bogunovic, I.; Scarlett, J.; Krause, A.; and Cevher, V. 2016. Truncated variance reduction: A unified approach to Bayesian optimization and level-set estimation. In *Proc. NeurIPS*, 1507–1515.
- Brochu, E.; Cora, V. M.; and de Freitas, N. 2010. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. arXiv:1012.2599.
- Bryan, B.; Nichol, R. C.; Genovese, C. R.; Schneider, J.; Miller, C. J.; and Wasserman, L. 2006. Active learning for identifying function threshold boundaries. In *Proc. NeurIPS*, 163–170.
- Calandra, R.; Seyfarth, A.; Peters, J.; and Deisenroth, M. P. 2014. An experimental comparison of Bayesian optimization for bipedal locomotion. In *Proc. ICRA*, 1951–1958.
- Daxberger, E. A.; and Low, K. H. 2017. Distributed Batch Gaussian process optimization. In *Proc. ICML*, 951–960.
- Galland, F.; Réfrégier, P.; and Germain, O. 2004. Synthetic aperture radar oil spill segmentation by stochastic complexity minimization. *IEEE Geoscience and Remote Sensing Letters* 1(4): 295–299.
- Gotovos, A.; Casati, N.; Hitz, G.; and Krause, A. 2013. Active learning for level set estimation. In *Proc. IJCAI*, 1344–1350.
- Hernández-Lobato, J. M.; Hoffman, M. W.; and Ghahramani, Z. 2014. Predictive entropy search for efficient global optimization of black-box functions. In *Proc. NeurIPS*, 918–926.
- Hoang, T. N.; Hoang, Q. M.; and Low, K. H. 2018. Decentralized high-dimensional Bayesian optimization with factor graphs. In *Proc. AAAI*, 3231–3238.
- Hoffman, M. W.; and Ghahramani, Z. 2015. Output-space predictive entropy search for flexible global optimization. In *Proc. NeurIPS Workshop on Bayesian Optimization*.
- Kharkovskii, D.; Ling, C. K.; and Low, K. H. 2020. Nonmyopic Gaussian process optimization with macro-actions. In *Proc. AISTATS*, 4593–4604.
- Kingma, D. P.; and Ba, J. 2015. Adam: A method for stochastic optimization. In *Proc. ICLR*.
- Kingma, D. P.; and Welling, M. 2013. Auto-encoding variational Bayes. arXiv:1312.6114.
- Krause, A.; and Ong, C. S. 2011. Contextual Gaussian process bandit optimization. In *Proc. NeurIPS*, 2447–2455.
- Ling, C. K.; Low, K. H.; and Jaillet, P. 2016. Gaussian process planning with Lipschitz continuous reward functions: Towards unifying Bayesian optimization, active learning, and beyond. In *Proc. AAAI*, 1860–1866.
- Low, K. H.; Chen, J.; Dolan, J. M.; Chien, S.; and Thompson, D. R. 2012. Decentralized active robotic exploration and mapping for probabilistic field classification in environmental sensing. In *Proc. AAMAS*, 105–112.
- Rahimi, A.; and Recht, B. 2008. Random features for large-scale kernel machines. In *Proc. NeurIPS*, 1177–1184.
- Rasmussen, C. E.; and Williams, C. K. I. 2006. *Gaussian Processes for Machine Learning*. MIT Press.
- Shahriari, B.; Swersky, K.; Wang, Z.; Adams, R.; and de Freitas, N. 2015. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE* 104(1): 148–175.
- Snoek, J.; Larochelle, H.; and Adams, R. 2012. Practical Bayesian optimization of machine learning algorithms. In *Proc. NeurIPS*, 2951–2959.
- Srinivas, N.; Krause, A.; Kakade, S.; and Seeger, M. 2010. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proc. ICML*, 1015–1022.
- Suzuki, S.; Takeno, S.; Tamura, T.; Shitara, K.; and Karasuyama, M. 2020. Multi-objective Bayesian optimization using Pareto-frontier entropy. In *Proc. ICML*.
- Takeno, S.; Fukuoka, H.; Tsukada, Y.; Koyama, T.; Shiga, M.; Takeuchi, I.; and Karasuyama, M. 2020. Multi-fidelity

Bayesian optimization with max-value entropy search and its parallelization. In *Proc. ICML*.

Wang, Z.; and Jegelka, S. 2017. Max-value entropy search for efficient Bayesian optimization. In *Proc. ICML*, 3627–3635.

Webster, R.; and Oliver, M. 2007. *Geostatistics for Environmental Scientists*. John Wiley & Sons, Inc., 2nd edition.

Zhang, Y.; Dai, Z.; and Low, K. H. 2019. Bayesian optimization with binary auxiliary information. In *Proc. UAI*.

A Derivation of (3)

It is known that $I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})$ is the *Kullback-Leibler (KL) divergence* between $p(y_{\mathbf{x}}, \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})$ and $p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})$. So,

$$\begin{aligned} & I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) \\ &= \sum_{\gamma_{\mathbf{x}}^{\circ}} \int p(y_{\mathbf{x}}, \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) \log \frac{p(y_{\mathbf{x}}, \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})}{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})} dy_{\mathbf{x}} \\ &= \mathbb{E}_{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})} \left[\sum_{\gamma_{\mathbf{x}}^{\circ}} p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, f_{\circ}) \log \frac{p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, f_{\circ})}{p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})} \right]. \end{aligned} \quad (8)$$

Note that

$$\begin{aligned} p(\gamma_{\mathbf{x}}^{\circ} = -1 | \mathbf{y}_{\mathcal{D}}, f_{\circ}) &= 1 - p(\gamma_{\mathbf{x}}^{\circ} = 1 | \mathbf{y}_{\mathcal{D}}, f_{\circ}) \\ p(\gamma_{\mathbf{x}}^{\circ} = 1 | \mathbf{y}_{\mathcal{D}}, f_{\circ}) &= p(f(\mathbf{x}) < f_{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) = \Psi \left(\frac{f_{\circ} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \right) \end{aligned}$$

where $\Psi((f_{\circ} - \mu_{\mathbf{x}})/\sigma_{\mathbf{x}})$ is the *cumulative density function* (c.d.f.) of the standard Gaussian distribution at $(f_{\circ} - \mu_{\mathbf{x}})/\sigma_{\mathbf{x}}$. Then,

$$p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) = \Psi \left(\gamma_{\mathbf{x}}^{\circ} \frac{f_{\circ} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \right) = \Psi(\gamma_{\mathbf{x}}^{\circ} h_{\mathbf{x}}(f_{\circ})) \quad (9)$$

where $h_{\mathbf{x}}(f_{\circ}) \triangleq (f_{\circ} - \mu_{\mathbf{x}})/\sigma_{\mathbf{x}}$.

We can evaluate $p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, f_{\circ})$ in the same manner as $p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ})$ by computing the GP posterior belief $p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}})$ (1) with all the observations $\mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}$, which incurs $\mathcal{O}((|\mathcal{D}| + 1)^3)$ time. On the other hand, we can compute $p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}})$ via an incremental update of $p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}) = \mathcal{N}(f(\mathbf{x}) | \mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2)$ with the new observation $y_{\mathbf{x}}$ as follows:

$$\begin{aligned} p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}) &= \frac{p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}) p(y_{\mathbf{x}} | f(\mathbf{x}))}{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})} \\ &= \mathcal{N} \left(f(\mathbf{x}) \left| \frac{\sigma_{\mathbf{x}}^2 y_{\mathbf{x}} + \sigma_n^2 \mu_{\mathbf{x}}}{\sigma_+^2}, \frac{\sigma_{\mathbf{x}}^2 \sigma_n^2}{\sigma_+^2} \right. \right) \end{aligned}$$

where $\sigma_+^2 = \sigma_{\mathbf{x}}^2 + \sigma_n^2$ is previously defined in Sec. 2. As a result,

$$\begin{aligned} & p(\gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, f_{\circ}) \\ &= \Psi \left(\gamma_{\mathbf{x}}^{\circ} \left(f_{\circ} - \frac{\sigma_{\mathbf{x}}^2 y_{\mathbf{x}} + \sigma_n^2 \mu_{\mathbf{x}}}{\sigma_+^2} \right) \left/ \sqrt{\frac{\sigma_{\mathbf{x}}^2 \sigma_n^2}{\sigma_+^2}} \right. \right) \quad (10) \\ &= \Psi \left(\gamma_{\mathbf{x}}^{\circ} \frac{\sigma_+^2 f_{\circ} - \sigma_n^2 \mu_{\mathbf{x}} - \sigma_{\mathbf{x}}^2 y_{\mathbf{x}}}{\sigma_{\mathbf{x}} \sigma_n \sigma_+} \right) \\ &= \Psi(\gamma_{\mathbf{x}}^{\circ} g_{\mathbf{x}}(y_{\mathbf{x}}, f_{\circ})) \end{aligned}$$

where $g_{\mathbf{x}}(y_{\mathbf{x}}, f_{\circ}) \triangleq (\sigma_+^2 f_{\circ} - \sigma_n^2 \mu_{\mathbf{x}} - \sigma_{\mathbf{x}}^2 y_{\mathbf{x}}) / (\sigma_{\mathbf{x}} \sigma_n \sigma_+)$. By plugging (9) and (10) into (8),

$$\begin{aligned} & I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\circ} | \mathbf{y}_{\mathcal{D}}, f_{\circ}) \\ &= \mathbb{E}_{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})} \left[\sum_{\gamma_{\mathbf{x}}^{\circ}} \Psi(\gamma_{\mathbf{x}}^{\circ} g_{\mathbf{x}}(y_{\mathbf{x}}, f_{\circ})) \log \frac{\Psi(\gamma_{\mathbf{x}}^{\circ} g_{\mathbf{x}}(y_{\mathbf{x}}, f_{\circ}))}{\Psi(\gamma_{\mathbf{x}}^{\circ} h_{\mathbf{x}}(f_{\circ}))} \right]. \end{aligned}$$

B Proof of (5)

In this subsection, we overload the notation f_{\star} to denote a discrete uniform random variable on the support \mathcal{F}_{\star} , i.e., $p(f_{\star}) = 1/|\mathcal{F}_{\star}|$ for all $f_{\star} \in \mathcal{F}_{\star}$. We will prove that

$$I(y_{\mathbf{x}}; (\gamma_{\mathbf{x}}^{\star}, f_{\star}) | \mathbf{y}_{\mathcal{D}}) = \frac{1}{|\mathcal{F}_{\star}|} \sum_{f_{\star} \in \mathcal{F}_{\star}} I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\star} | \mathbf{y}_{\mathcal{D}}, f_{\star})$$

where the RHS is the definition of $\alpha_{\text{BES-MP}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}})$ in (4) and the LHS is the mutual information in (5) which allows BES-MP to be interpreted as the information gain on both the class label $\gamma_{\mathbf{x}}^{\star}$ and the threshold $f_{\star} \in \mathcal{F}_{\star}$ inducing the superlevel set $\mathcal{X}_{f_{\star}}^+$ (of potential maximizers) from evaluating f at input query \mathbf{x} to observe $y_{\mathbf{x}}$ (Remark 2).

Firstly, we show that $f(\mathbf{x})$ and f_{\star} are conditionally independent if $\gamma_{\mathbf{x}}^{\star}$ is unobserved. We know that

$$\begin{aligned} & p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}) \\ &= p(f(\mathbf{x}), \gamma_{\mathbf{x}}^{\star} = 1 | \mathbf{y}_{\mathcal{D}}, f_{\star}) + p(f(\mathbf{x}), \gamma_{\mathbf{x}}^{\star} = -1 | \mathbf{y}_{\mathcal{D}}, f_{\star}) \\ &= p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1) p(\gamma_{\mathbf{x}}^{\star} = 1 | \mathbf{y}_{\mathcal{D}}, f_{\star}) \\ &\quad + p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = -1) p(\gamma_{\mathbf{x}}^{\star} = -1 | \mathbf{y}_{\mathcal{D}}, f_{\star}). \end{aligned}$$

We observe that $p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1)$ is a truncated Gaussian probability density function on the support $(-\infty, f_{\star})$:

$$p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = 1) = \frac{\mathbb{I}_{f(\mathbf{x}) < f_{\star}} p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}})}{p(\gamma_{\mathbf{x}}^{\star} = 1 | \mathbf{y}_{\mathcal{D}}, f_{\star})}.$$

Similarly,

$$p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}, \gamma_{\mathbf{x}}^{\star} = -1) = \frac{\mathbb{I}_{f(\mathbf{x}) \geq f_{\star}} p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}})}{p(\gamma_{\mathbf{x}}^{\star} = -1 | \mathbf{y}_{\mathcal{D}}, f_{\star})}.$$

Therefore,

$$p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}) = p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}})$$

which implies that $f(\mathbf{x})$ and f_{\star} are conditionally independent if $\gamma_{\mathbf{x}}^{\star}$ is unobserved. Consequently, $y_{\mathbf{x}}$ and f_{\star} are conditionally independent if $\gamma_{\mathbf{x}}^{\star}$ is unobserved:

$$\begin{aligned} p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}) &= \int p(y_{\mathbf{x}} | f(\mathbf{x})) p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}, f_{\star}) dy_{\mathbf{x}} \\ &= \int p(y_{\mathbf{x}} | f(\mathbf{x})) p(f(\mathbf{x}) | \mathbf{y}_{\mathcal{D}}) dy_{\mathbf{x}} = p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}). \end{aligned}$$

It follows that we can express the prior entropy as follows:

$$\begin{aligned} & H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})) \\ &= - \int p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}) \log p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}) dy_{\mathbf{x}} \\ &= - \sum_{f_{\star} \in \mathcal{F}_{\star}} p(f_{\star}) \int p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}) \log p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}) dy_{\mathbf{x}} \\ &= \mathbb{E}_{p(f_{\star})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}))]. \end{aligned}$$

Therefore,

$$\begin{aligned} & I(y_{\mathbf{x}}; (\gamma_{\mathbf{x}}^{\star}, f_{\star}) | \mathbf{y}_{\mathcal{D}}) \\ &= H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})) - \mathbb{E}_{p(\gamma_{\mathbf{x}}^{\star}, f_{\star} | \mathbf{y}_{\mathcal{D}})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, \gamma_{\mathbf{x}}^{\star}, f_{\star}))] \\ &= \mathbb{E}_{p(f_{\star})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}))] - \mathbb{E}_{p(\gamma_{\mathbf{x}}^{\star}, f_{\star} | \mathbf{y}_{\mathcal{D}})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, \gamma_{\mathbf{x}}^{\star}, f_{\star}))] \\ &= \mathbb{E}_{p(f_{\star})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, f_{\star}))] - \mathbb{E}_{p(\gamma_{\mathbf{x}}^{\star} | \mathbf{y}_{\mathcal{D}}, f_{\star})} [H(p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, \gamma_{\mathbf{x}}^{\star}, f_{\star}))] \\ &= \mathbb{E}_{p(f_{\star})} [I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^{\star} | \mathbf{y}_{\mathcal{D}}, f_{\star})]. \end{aligned}$$

Since f_* follows a discrete uniform distribution on the support \mathcal{F}_* , $p(f_*) = 1/|\mathcal{F}_*|$. So,

$$\mathbb{E}_{p(f_*)}[I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)] = \frac{1}{|\mathcal{F}_*|} \sum_{f_* \in \mathcal{F}_*} I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*).$$

C Proof of Theorem 1

If the observation is noiseless (i.e., $\sigma_n^2 = 0$) and $f_* = \alpha_{\text{UCB}}(\mathbf{x}_{\text{UCB}}, \mathbf{y}_{\mathcal{D}})$, then BES-MP reduces to only the prior entropy $H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*))$, as explained in Remark 1. We will prove that BES-MP selects the same input queries as that selected by UCB:

$$\operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)) = \mathbf{x}_{\text{UCB}}.$$

We adapt a proof from that of Low et al. (2012) to show that maximizing $H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*))$ is equivalent to minimizing $|f_* - \mu_{\mathbf{x}}|/\sigma_{\mathbf{x}}$ w.r.t. $\mathbf{x} \in \mathcal{X}$:

$$\begin{aligned} & \max_{\mathbf{x} \in \mathcal{X}} H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)) \\ &= \max_{\mathbf{x} \in \mathcal{X}} - \sum_{\gamma_{\mathbf{x}}^*} p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*) \log p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*) \\ &= \min_{\mathbf{x} \in \mathcal{X}} \left(p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*) \log p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*) \right. \\ & \quad \left. + p(\gamma_{\mathbf{x}}^* = 1 | \mathbf{y}_{\mathcal{D}}, f_*) \log p(\gamma_{\mathbf{x}}^* = 1 | \mathbf{y}_{\mathcal{D}}, f_*) \right) \\ &= \min_{\mathbf{x} \in \mathcal{X}} \left(p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*) \log p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*) \right. \\ & \quad \left. + (1 - p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*)) \log(1 - p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*)) \right) \\ &= \min_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{2} - p(\gamma_{\mathbf{x}}^* = -1 | \mathbf{y}_{\mathcal{D}}, f_*) \right| \\ &= \min_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{2} - p(f(\mathbf{x}) \geq f_* | \mathbf{y}_{\mathcal{D}}, f_*) \right| \\ &= \min_{\mathbf{x} \in \mathcal{X}} \left| \operatorname{erf} \left(\frac{f_* - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}} \sqrt{2}} \right) \right| \\ &= \min_{\mathbf{x} \in \mathcal{X}} \frac{|f_* - \mu_{\mathbf{x}}|}{\sigma_{\mathbf{x}}}. \end{aligned}$$

That is,

$$\operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{|f_* - \mu_{\mathbf{x}}|}{\sigma_{\mathbf{x}}}. \quad (11)$$

Since $\mathbf{x}_{\text{UCB}} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \alpha_{\text{UCB}}(\mathbf{x}, \mathbf{y}_{\mathcal{D}})$ and $\beta > 0$,

$$\begin{aligned} \mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}} &\geq \mu_{\mathbf{x}} + \beta \sigma_{\mathbf{x}} \\ \mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}} - \mu_{\mathbf{x}} &\geq \beta \sigma_{\mathbf{x}} \geq 0 \end{aligned}$$

for all $\mathbf{x} \in \mathcal{X}$. It follows that since $f_* = \alpha_{\text{UCB}}(\mathbf{x}_{\text{UCB}}, \mathbf{y}_{\mathcal{D}}) = \mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}}$, we can bound $|f_* - \mu_{\mathbf{x}}|/\sigma_{\mathbf{x}}$ from below:

$$\begin{aligned} \frac{|f_* - \mu_{\mathbf{x}}|}{\sigma_{\mathbf{x}}} &= \frac{|\mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}} - \mu_{\mathbf{x}}|}{\sigma_{\mathbf{x}}} \\ &= \frac{\mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \geq \frac{\beta \sigma_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \geq \beta. \end{aligned}$$

Furthermore, since $f_* = \alpha_{\text{UCB}}(\mathbf{x}_{\text{UCB}}, \mathbf{y}_{\mathcal{D}}) = \mu_{\mathbf{x}_{\text{UCB}}} + \beta \sigma_{\mathbf{x}_{\text{UCB}}}$, when $\mathbf{x} = \mathbf{x}_{\text{UCB}}$,

$$\frac{|f_* - \mu_{\mathbf{x}_{\text{UCB}}}|}{\sigma_{\mathbf{x}_{\text{UCB}}}} = \beta.$$

Therefore,

$$\operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{|f_* - \mu_{\mathbf{x}}|}{\sigma_{\mathbf{x}}} = \mathbf{x}_{\text{UCB}}. \quad (12)$$

From (11) and (12), we have shown that when $f_* = \alpha_{\text{UCB}}(\mathbf{x}_{\text{UCB}}, \mathbf{y}_{\mathcal{D}})$, $\operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} H(p(\gamma_{\mathbf{x}}^* | \mathbf{y}_{\mathcal{D}}, f_*)) = \mathbf{x}_{\text{UCB}}$.

D Alternative Form of BES^k

It is known that $I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b})$ is the KL divergence between $p(y_{\mathbf{x}}, \gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b})$ and $p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}}, \mathbf{b}) p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b})$. So, we can obtain a similar expression to (8) (Appendix A):

$$\begin{aligned} & I(y_{\mathbf{x}}; \gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b}) \\ &= \mathbb{E}_{p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})} \left[\sum_{\gamma_{\mathbf{x}}^k} p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, \mathbf{b}) \log \frac{p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, \mathbf{b})}{p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b})} \right]. \end{aligned} \quad (13)$$

Let $b_0 \triangleq -\infty$ and $b_{k+1} \triangleq \infty$. Then, $p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b})$ and $p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, \mathbf{b})$ can be expressed as follows:

$$\begin{aligned} & p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D}}, \mathbf{b}) \\ &= p(b_{\gamma_{\mathbf{x}}^k} \leq f(\mathbf{x}) < b_{\gamma_{\mathbf{x}}^k+1} | \mathbf{y}_{\mathcal{D}}, \mathbf{b}) \\ &= \Psi \left(\frac{b_{\gamma_{\mathbf{x}}^k+1} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \right) - \Psi \left(\frac{b_{\gamma_{\mathbf{x}}^k} - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}} \right) \\ &= \Psi(h_{\mathbf{x}}(b_{\gamma_{\mathbf{x}}^k+1})) - \Psi(h_{\mathbf{x}}(b_{\gamma_{\mathbf{x}}^k})) \end{aligned}$$

and

$$p(\gamma_{\mathbf{x}}^k | \mathbf{y}_{\mathcal{D} \cup \{\mathbf{x}\}}, \mathbf{b}) = \Psi(g_{\mathbf{x}}(y_{\mathbf{x}}, b_{\gamma_{\mathbf{x}}^k+1})) - \Psi(g_{\mathbf{x}}(y_{\mathbf{x}}, b_{\gamma_{\mathbf{x}}^k}))$$

where $h_{\mathbf{x}}$ and $g_{\mathbf{x}}$ are previously defined in the line after (9) and (10), respectively.

We can optimize (13) via stochastic gradient ascent by reparameterizing the GP posterior belief $p(y_{\mathbf{x}} | \mathbf{y}_{\mathcal{D}})$ to a standard Gaussian distribution (Kingma and Welling 2013).

E Further Experimental Results

In this subsection, we present both the mean/average and the standard deviation of the log loss for LSE (Sec. 6.1) and implicit LSE (Sec. 6.3) experiments and the regret for BO experiments (Sec. 6.2) in the last iteration. The results are shown in Tables 1, 2, and 3 below:

Table 1: Mean/average and standard deviation of the log loss for the LSE experiments (Sec. 6.1).

Experiment	σ_n^2	BES	EM	STRDL
GP sample ($l = 1/3$)	0.0001	0.0022 \pm 0.0011	0.0027 \pm 0.0015	0.0038 \pm 0.0012
	0.09	0.0270 \pm 0.0114	0.0360 \pm 0.0100	0.0265 \pm 0.0065
GP sample ($l = 0.125$)	0.0001	0.0136 \pm 0.0046	0.0535 \pm 0.0276	0.0436 \pm 0.0185
	0.09	0.1067 \pm 0.0192	0.4722 \pm 0.1403	0.1339 \pm 0.0264
Branin	0.0001	0.0010 \pm 0.0004	0.0015 \pm 0.0006	0.0083 \pm 0.0140
	0.09	0.0354 \pm 0.0208	0.0673 \pm 0.0232	0.0522 \pm 0.0236
Michaelwicz	0.0001	0.0017 \pm 0.0004	0.0026 \pm 0.0008	0.1758 \pm 0.1035
	0.09	0.0136 \pm 0.0041	0.0467 \pm 0.0813	0.1815 \pm 0.0871
Phosphorus	0.0251	0.0318 \pm 0.0019	0.0870 \pm 0.0430	0.1100 \pm 0.0438

Table 2: Mean/average and standard deviation of the regret for the BO experiments (Sec. 6.2).

Experiment	BES-MP	PES	EI	UCB	MES
Michaelwicz	0.0017 \pm 0.0017	0.1524 \pm 0.2943	0.0011 \pm 0.0009	0.0048 \pm 0.0052	0.1178 \pm 0.2275
Hartmann-3d	0.0031 \pm 0.0017	0.0083 \pm 0.0044	0.0044 \pm 0.0033	0.0082 \pm 0.0060	0.0332 \pm 0.0749
Goldstein	0.0131 \pm 0.0033	0.0152 \pm 0.0037	0.0139 \pm 0.0029	0.0200 \pm 0.0118	0.0170 \pm 0.0042
Phosphorus	0.0012 \pm 0.0013	0.0034 \pm 0.0029	0.0011 \pm 0.0009	0.1886 \pm 0.3068	0.2447 \pm 0.3182
MNIST	0.0667 \pm 0.0030	0.0669 \pm 0.0030	0.0659 \pm 0.0001	0.0659 \pm 0.0000	0.0660 \pm 0.0002
CIFAR-10	0.376 \pm 0.022	0.378 \pm 0.018	0.377 \pm 0.010	0.395 \pm 0.023	0.381 \pm 0.022

Table 3: Mean/average and standard deviation of the log loss for the implicit LSE experiments (Sec. 6.3).

Experiment	Unknown f_* (i.e., implicit LSE)		Known f_* (i.e., reducing implicit LSE to LSE)		
	BES ² -MP	BES-MP	BES	EM	STRDL
GP sample	0.0016 \pm 0.0010	0.0264 \pm 0.0251	0.0010 \pm 0.0004	0.0022 \pm 0.0018	0.0380 \pm 0.0224
Branin	0.0125 \pm 0.0047	0.1355 \pm 0.2009	0.0100 \pm 0.0024	0.0115 \pm 0.0025	0.0140 \pm 0.0028
Goldstein	0.0251 \pm 0.0094	0.0344 \pm 0.0308	0.0097 \pm 0.0026	0.0125 \pm 0.0028	0.1442 \pm 0.2235
Hartmann-3d	0.0023 \pm 0.0006	0.0331 \pm 0.0365	0.0010 \pm 0.0003	0.0012 \pm 0.0003	0.1598 \pm 0.2098
Phosphorus	0.0032 \pm 0.0007	0.0045 \pm 0.0018	0.0029 \pm 0.0008	0.0037 \pm 0.0018	0.0491 \pm 0.186