
Batch Bayesian optimisation via density-ratio estimation with guarantees

Anonymous Author(s)

Affiliation

Address

email

1 Appendix

2 This appendix complements the main paper with proofs, experiment details and additional experiments
3 and discussions. [Appendix A](#) presents full proofs for the main theoretical results in the paper.
4 In [Appendix B](#), we discuss an approach to derive alternative regret bounds for BORE under a
5 Thompson sampling perspective. We discuss the theoretical analysis of BORE with its non-constant
6 approximation for the observations quantile τ in [Appendix C](#). In [Appendix D](#), we present further
7 details on the experiments setup. Finally, [Appendix E](#) presents additional experiments on higher-
8 dimensional settings.

9 A Proofs

10 This section presents proofs for the main theoretical results in the paper. We start with a few auxiliary
11 results from the GP-UCB literature [[1](#), [2](#)], following up with the proofs for the main theorems.

12 A.1 Auxiliary results

13 **Lemma A.1** (Srinivas et al. [[1](#), Lemma 5.3]). *The information gain for a sequence of $N \geq 1$*
14 *observations $\{\mathbf{x}_i, z_i\}_{i=1}^N$, where $z_i = f(\mathbf{x}_i) + \nu_i$, $\nu_i \sim \mathcal{N}(0, \lambda)$, can be expressed in terms of*
15 *the predictive variances. Namely, if $f \sim \mathcal{GP}(m, k)$, then the information gain provided by the*
16 *observations is such that:*

$$I(\mathbf{z}_N, \mathbf{f}_N | \mathcal{X}_N) = \frac{1}{2} \sum_{i=1}^N \log(1 + \lambda^{-1} \sigma_{i-1}^2(\mathbf{x}_i)), \quad (\text{A.1})$$

17 where $\mathbf{f}_N := [f(\mathbf{x}_i)]_{i=1}^N$ and $\mathcal{X}_N := \{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{X}$.

18 **Lemma A.2** (Chowdhury and Gopalan [[2](#), Lemma 4]). *Following the setting of [Lemma A.1](#), the sum*
19 *of predictive standard deviations at a sequence of N points is bounded in terms of the maximum*
20 *information gain:*

$$\sum_{i=1}^N \sigma_{i-1}(\mathbf{x}_i) \leq \sqrt{4(N+2)\xi_N}. \quad (\text{A.2})$$

21 **Lemma A.3.** *Let $\mathcal{A} \subset \mathcal{X}$ be a finite set of points where a function $f \sim \mathcal{GP}(m, k)$ was evaluated, so*
22 *that the GP posterior covariance function and the corresponding variance are given by:*

$$k_{\mathcal{A}}(\mathbf{x}, \mathbf{x}') := k(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathcal{A})^{\top} (\mathbf{K}(\mathcal{A}) + \eta \mathbf{I})^{-1} k(\mathcal{A}, \mathbf{x}') \quad (\text{A.3})$$

$$\sigma_{\mathcal{A}}^2(\mathbf{x}) := k_{\mathcal{A}}(\mathbf{x}, \mathbf{x}), \quad \mathbf{x}, \mathbf{x}' \in \mathcal{X}, \quad (\text{A.4})$$

23 where $k(\mathbf{x}, \mathcal{A}) := [k(\mathbf{x}, \mathbf{a})]_{\mathbf{a} \in \mathcal{A}}$ and $\mathbf{K}(\mathcal{A}) := [k(\mathbf{a}, \mathbf{a}')]_{\mathbf{a}, \mathbf{a}' \in \mathcal{A}}$. Then, for any given set $\mathcal{B} \supset \mathcal{A}$ of
24 evaluations of f , we have:

$$\sigma_{\mathcal{B}}^2(\mathbf{x}) \leq \sigma_{\mathcal{A}}^2(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}. \quad (\text{A.5})$$

25 *Proof.* The result follows by observing that the GP posterior given observations at \mathcal{A} is a prior for
 26 the GP with the new observations at the complement $\mathcal{C} := \mathcal{B} \setminus \mathcal{A}$. Then we obtain, for all $\mathbf{x} \in \mathcal{X}$:

$$\begin{aligned}\sigma_{\mathcal{B}}^2(\mathbf{x}) &:= k(\mathbf{x}, \mathbf{x}) - k(\mathbf{x}, \mathcal{B})(\mathbf{K}(\mathcal{B}) + \eta\mathbf{I})^{-1}k(\mathcal{B}, \mathbf{x}) \\ &= \sigma_{\mathcal{A}}^2(\mathbf{x}) - k_{\mathcal{A}}(\mathbf{x}, \mathcal{C})(\mathbf{K}_{\mathcal{A}}(\mathcal{C}) + \eta\mathbf{I})^{-1}k_{\mathcal{A}}(\mathcal{C}, \mathbf{x}) \\ &\leq \sigma_{\mathcal{A}}^2(\mathbf{x}),\end{aligned}\tag{A.6}$$

27 since $k_{\mathcal{A}}(\mathbf{x}, \mathcal{C})(\mathbf{K}_{\mathcal{A}}(\mathcal{C}) + \eta\mathbf{I})^{-1}k_{\mathcal{A}}(\mathcal{C}, \mathbf{x})$ is non-negative. \square

28 A.2 Main proofs

29 A.2.1 Proof of Theorem 1

30 *Proof of Theorem 1.* The proof follows by a simple application of Durand et al. [3, Thm. 1] on
 31 GP-UCB to our settings, as $\pi \in \mathcal{H}$ and the stochastic process defining the query locations \mathbf{x}_t and
 32 observation noise $\nu_t := z_t - \pi(\mathbf{x}_t)$ satisfies their assumptions. \square

33 A.2.2 Proof of Theorem 2

34 To prove Theorem 2, we will follow the procedure of GP-UCB proofs [1, 2] by bounding the
 35 approximation error $|\pi(\mathbf{x}) - \hat{\pi}_t(\mathbf{x})|$ via a confidence bound (Theorem 1) and then applying it to the
 36 instant regret. From the instant regret to the cumulative regret, the bounds are extended by means of
 37 the maximum information gain ξ_T introduced in the main text. One of the differences with our proof,
 38 however, is that BORE with a PLS classifier is not following the optimal UCB policy, but instead
 39 a pure-exploitation approach by following the maximum of the mean estimator $\hat{\pi}_t$, which does not
 40 account for uncertainty.

41 *Proof of Theorem 2.* Recalling the classifier-based bound in Section 4 and that for any $\tau \in \mathbb{R}$ the
 42 result in Lemma 1 holds, we have:

$$\begin{aligned}r_t &= f(\mathbf{x}_t) - f(\mathbf{x}^*) \\ &\leq L_\epsilon(\pi(\mathbf{x}^*) - \pi(\mathbf{x}_t))\end{aligned}\tag{A.7}$$

43 According to Theorem 1, working with the confidence bounds on $\pi(\mathbf{x})$, we then have that the instant
 44 regret is bounded with probability at least $1 - \delta$ by:

$$\begin{aligned}\forall t \geq 1, \quad r_t &\leq L_\epsilon(\hat{\pi}_{t-1}(\mathbf{x}^*) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}^*) - \pi(\mathbf{x}_t)) \\ &\leq L_\epsilon(\hat{\pi}_{t-1}(\mathbf{x}^*) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}^*) - \hat{\pi}_{t-1}(\mathbf{x}_t) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}_t)) \\ &\leq L_\epsilon\beta_{t-1}(\delta)(\sigma_{t-1}(\mathbf{x}^*) + \sigma_{t-1}(\mathbf{x}_t)),\end{aligned}\tag{A.8}$$

45 since $\hat{\pi}_{t-1}(\mathbf{x}^*) \leq \max_{\mathbf{x} \in \mathcal{X}} \hat{\pi}_{t-1}(\mathbf{x}) = \hat{\pi}_{t-1}(\mathbf{x}_t)$. Now we can apply Lemma A.2, yielding with
 46 probability at least $1 - \delta$:

$$\begin{aligned}R_T &:= \sum_{t=1}^T r_t \leq L_\epsilon\beta_T(\delta) \sum_{t=1}^T (\sigma_{t-1}(\mathbf{x}_t) + \sigma_{t-1}(\mathbf{x}^*)) \\ &\leq L_\epsilon\beta_T(\delta) \left(\sqrt{4(T+2)\xi_T} + \sum_{t=1}^T \sigma_{t-1}(\mathbf{x}^*) \right)\end{aligned}\tag{A.9}$$

47 since $\beta_t(\delta) \leq \beta_{t+1}(\delta)$ for all $t \geq 1$. This concludes the proof. \square

48 A.2.3 Proof of Theorem 3

49 Again, we will be following standard GP-UCB proofs for this result using the bound in Theorem 1.

50 *Proof of Theorem 3.* Extending the bound in Equation A.7 with Theorem 1, we have with probability
 51 at least $1 - \delta$:

$$\begin{aligned}\forall t \geq 1, \quad r_t &\leq L_\epsilon(\hat{\pi}_{t-1}(\mathbf{x}^*) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}^*) - \pi_{t-1}^*(\mathbf{x}_t)) \\ &\leq L_\epsilon(\hat{\pi}_{t-1}(\mathbf{x}^*) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}^*) - \hat{\pi}_{t-1}(\mathbf{x}_t) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}_t)) \\ &\leq 2L_\epsilon\beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}_t),\end{aligned}\tag{A.10}$$

52 since $\hat{\pi}_{t-1}(\mathbf{x}^*) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}^*) \leq \max_{\mathbf{x} \in \mathcal{X}} \hat{\pi}_{t-1}(\mathbf{x}) + \beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}) = \hat{\pi}_{t-1}(\mathbf{x}_t) +$
 53 $\beta_{t-1}(\delta)\sigma_{t-1}(\mathbf{x}_t)$). Turning our attention to the cumulative regret, with the same probability, we have:
 54

$$R_T := \sum_{t=1}^T r_t \leq 2L_\epsilon \beta_T(\delta) \sum_{t=1}^T \sigma_{t-1}(\mathbf{x}_t) \leq 4L_\epsilon \beta_T(\delta) \sqrt{(T+2)\xi_T}, \quad (\text{A.11})$$

55 which concludes the proof. \square

56 A.2.4 Proof of Theorem 4

57 *Proof.* Starting with the regret definition, we can define a bound in terms of the discrepancy between
 58 the two sampling distributions:

$$\begin{aligned} r_t &:= \mathbb{E}_{\mathbf{x} \sim \hat{p}_t}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \ell}[f(\mathbf{x})] \\ &\leq L_\epsilon (\mathbb{E}_{\mathbf{x} \sim \ell}[\pi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \hat{p}_t}[\pi(\mathbf{x})]) \\ &\leq L_\epsilon \|\pi\|_\infty \int_{\mathcal{X}} |\ell(\mathbf{x}) - q_{t-1}(\mathbf{x})| d\mathbf{x} \\ &\leq L_\epsilon \|\pi\|_\infty \sqrt{\frac{1}{2} D_{\text{KL}}(q_{t-1} \|\ell)}, \quad \forall t \geq 1, \end{aligned} \quad (\text{A.12})$$

59 where the last line is due to Pinsker's inequality [4] applied to the total variation distance between \hat{p}_t
 60 and ℓ (third line).

61 To obtain a bound on $D_{\text{KL}}(\hat{p}_t \|\ell)$, starting from the definition of the terms, with probability at least
 62 $1 - \delta$, we have that:

$$\begin{aligned} \forall t \geq 0, \quad D_{\text{KL}}(\hat{p}_t \|\ell) &= \mathbb{E}_{\mathbf{x} \sim \hat{p}_t}[\log \hat{p}_t(\mathbf{x}) - \log \ell(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim \hat{p}_t}[\log(\hat{\pi}_t(\mathbf{x}) + \beta_t(\delta)\sigma_t(\mathbf{x})) - \log \pi(\mathbf{x}) + \log \eta_\pi - \log \gamma] \\ &\leq \mathbb{E}_{\mathbf{x} \sim \hat{p}_t}[\log(\hat{\pi}_t(\mathbf{x}) + \beta_t(\delta)\sigma_t(\mathbf{x})) - \log \pi(\mathbf{x})], \end{aligned} \quad (\text{A.13})$$

63 which follows from $\eta_t := \int_{\mathcal{X}} (\hat{\pi}_t(\mathbf{x}) + \beta_t(\delta)\sigma_t(\mathbf{x}))p(\mathbf{x}) d\mathbf{x} \geq \int_{\mathcal{X}} \pi(\mathbf{x})p(\mathbf{x}) d\mathbf{x} =: \gamma$. Now, by the
 64 mean value theorem [5], for all $t \geq 0$, we have that the following holds with the same probability:

$$\begin{aligned} |\log(\hat{\pi}_t(\mathbf{x}) + \beta_t(\delta)\sigma_t(\mathbf{x})) - \log \pi(\mathbf{x})| &\leq L_\pi |\hat{\pi}_t(\mathbf{x}) + \beta_t(\delta)\sigma_t(\mathbf{x}) - \pi(\mathbf{x})| \\ &\leq 2L_\pi \beta_t(\delta)\sigma_t(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \end{aligned} \quad (\text{A.14})$$

65 since $\frac{d \log(s)}{ds} < L_\pi < \infty$ for all $s > \min_{\mathbf{x} \in \mathcal{X}} \pi(\mathbf{x}) > 0$, and $|\hat{\pi}_t(\mathbf{x}) - \pi(\mathbf{x})| \leq \beta_t(\delta)\sigma_t(\mathbf{x})$ by
 66 Theorem 1. The first result in the theorem then follows.

67 For the second part of the result, we first note that:

$$\forall T \geq 1, \quad \min_{t \leq T} D_{\text{KL}}(\hat{p}_t \|\ell) \leq \frac{1}{T} \sum_{t=1}^T D_{\text{KL}}(\hat{p}_{t-1} \|\ell) \quad (\text{A.15})$$

68 Following the previous derivations, it holds with probability at least $1 - \delta$ that:

$$\begin{aligned} \sum_{t=1}^T D_{\text{KL}}(\hat{p}_t \|\ell) &\leq 2L_\pi \sum_{t=1}^T \beta_{t-1}(\delta) \mathbb{E}_{\tilde{\mathbf{x}}_t \sim \hat{p}_t}[\sigma_{t-1}(\tilde{\mathbf{x}}_t)] \\ &\leq 2L_\pi \beta_T(\delta) \sum_{t=1}^T \mathbb{E}_{\tilde{\mathbf{x}}_t \sim q_t}[\sigma_{t-1}(\tilde{\mathbf{x}}_t)] \\ &\leq 2L_\pi \beta_T(\delta) \mathbb{E}_{\tilde{\mathbf{x}}_1 \sim q_1, \dots, \tilde{\mathbf{x}}_T \sim q_T} \left[\sum_{t=1}^T \sigma_{t-1}(\tilde{\mathbf{x}}_t) \right], \end{aligned} \quad (\text{A.16})$$

69 since $\beta_T \geq \beta_t$, for all $t \leq T$, and expectations are linear operations. Considering the predictive
 70 variances above, recall that, at each iteration $t \geq 1$, the algorithm selects a batch of i.i.d. points
 71 $\mathcal{B}_t := \{\mathbf{x}_{t,i}\}_{i=1}^M$, sampled from \hat{p}_t , where to evaluate the objective function f . The predictive

72 variance σ_{t-1}^2 is conditioned on all previous observations, which are grouped by batches. We can
 73 then decompose, for any $t \geq 1$:

$$\sigma_t^2(\mathbf{x}) = \sigma_{t-1}^2(\mathbf{x}) - k_{t-1}(\mathbf{x}, \mathcal{B}_t)(\mathbf{K}_{t-1}(\mathcal{B}_t) + \eta\mathbf{I})^{-1}k_{t-1}(\mathcal{B}_t, \mathbf{x}), \quad (\text{A.17})$$

74 where we use the notation introduced in [Lemma A.3](#), and:

$$k_t(\mathbf{x}, \mathbf{x}') = k_{t-1}(\mathbf{x}, \mathbf{x}') - k_{t-1}(\mathbf{x}, \mathcal{B}_t)(\mathbf{K}_{t-1}(\mathcal{B}_t) + \eta\mathbf{I})^{-1}k_{t-1}(\mathcal{B}_t, \mathbf{x}'), \quad t \geq 1, \quad (\text{A.18})$$

$$k_0(\mathbf{x}, \mathbf{x}') := k(\mathbf{x}, \mathbf{x}'). \quad (\text{A.19})$$

75 Therefore, the predictive variance of the batched algorithm is not the same as the predictive variance
 76 of a sequential algorithm, and we cannot directly apply [Lemma A.2](#) to bound the last term in
 77 [Equation A.16](#).

78 [Lemma A.3](#) tells us that the predictive variance given a set of observations is less than the predictive
 79 variance given a subset of observations. Selecting only the first point from within each batch and
 80 applying [Lemma A.3](#), we get, for $t \geq 1$:

$$\sigma_t^2(\mathbf{x}) \leq s_t^2(\mathbf{x}) := k(\mathbf{x}, \mathbf{x}) - k(\mathbf{x}, \mathcal{X}_t)(\mathbf{K}(\mathcal{X}_t) + \eta\mathbf{I})^{-1}k(\mathcal{X}_t, \mathbf{x}), \quad (\text{A.20})$$

81 where $\mathcal{X}_t := \{\mathbf{x}_{i,1}\}_{i=1}^t$, with $\mathbf{x}_{i,1} \in \mathcal{B}_i$, $i \in \{1, \dots, t\}$. Note that the right-hand side of the equation
 82 above is simply the non-batched GP predictive variance. Furthermore, sample points within a batch
 83 are i.i.d., so that $\mathbf{x}_{t,1} \sim q_t$ and $\tilde{\mathbf{x}}_t \sim q_t$ are identically distributed. We can now apply [Lemma A.2](#),
 84 yielding:

$$\mathbb{E}_{\tilde{\mathbf{x}}_1 \sim q_1, \dots, \tilde{\mathbf{x}}_T \sim q_T} \left[\sum_{t=1}^T \sigma_{t-1}(\tilde{\mathbf{x}}_t) \right] \leq \mathbb{E}_{\tilde{\mathbf{x}}_1 \sim q_1, \dots, \tilde{\mathbf{x}}_T \sim q_T} \left[\sum_{t=1}^T s_{t-1}(\tilde{\mathbf{x}}_t) \right] \leq 2\sqrt{(T+2)\xi_T}. \quad (\text{A.21})$$

85 Combining this result with [Equation A.16](#), we obtain:

$$\sum_{t=1}^T D_{\text{KL}}(\hat{p}_t || \ell) \leq 4L_\pi \beta_T(\delta) \sqrt{(T+2)\xi_T} \in \mathcal{O}(\beta_T(\delta) \sqrt{T\xi_T}). \quad (\text{A.22})$$

86 Lastly, from the definition of $\beta_t(\delta)$, we have:

$$\beta_T(\delta) := b + \sigma_\nu \sqrt{2\lambda^{-1} \log(|\mathbf{I} + \lambda^{-1}\mathbf{K}_{\mathcal{D}_T}|^{1/2}/\delta)}, \quad (\text{A.23})$$

87 where:

$$\log(|\mathbf{I} + \lambda^{-1}\mathbf{K}_{\mathcal{D}_T}|^{1/2}) = I(\mathbf{z}_{N_T}, \mathbf{h}_{N_T}) \leq \xi_{N_T} = \xi_{MT}, \quad (\text{A.24})$$

88 for $h \sim \mathcal{GP}(m, k)$. Therefore, the cumulative sum of divergences is such that:

$$\sum_{t=1}^T D_{\text{KL}}(\hat{p}_t || \ell) \in \mathcal{O}(\sqrt{T}(b\sqrt{\xi_T} + \sqrt{\xi_T \xi_{MT}})). \quad (\text{A.25})$$

89 which concludes the proof. \square

90 B Bayesian regret bounds for BORE as Thompson sampling

91 Although in our main results we considered BORE using an optimal classifier according to a least-
 92 squares loss, we may instead consider that, in practice, the trained classifier might be sub-optimal due
 93 to training via gradient descent. In particular, in the case of stochastic gradient descent, Mandt et al.
 94 [\[6\]](#) showed that parameters learnt this way can be seen as approximate samples of a Bayesian posterior
 95 distribution. This is, therefore, the case of Thompson (or posterior) sampling [\[7\]](#). If we consider
 96 that the posterior over the model's function space is Gaussian, e.g., as in the case of infinitely-wide
 97 deep neural networks [\[8, 9\]](#), we may instead analyse the original BORE as a GP-based Thompson
 98 sampling algorithm. We can then apply theoretical results from Russo and Van Roy [\[7\]](#) to use general
 99 GP-UCB approximation guarantees [\[1, 10\]](#) to bound BORE'S Bayesian regret. Note, however, that
 100 this is a different type of analysis compared to the one presented in this paper, which considered a
 101 frequentist setting where the objective function is fixed, but unknown.

102 C Analysis with a non-constant quantile approximation

103 Our main theoretical results so far relied upon the quantile τ being fixed throughout all iterations
 104 $t \in \{1, \dots, T\}$, though in practice we have to approximate the quantile based on the empirical
 105 observations distribution up to time $t \geq 1$. In this section, we discuss the plausibility of the theoretical
 106 results under this practical scenario.

107 The main impact of a time-varying quantile τ_t , and the corresponding classifier $\pi_t(\mathbf{x}) := p(y \leq \tau_t | \mathbf{x})$,
 108 in theoretical results is in the UCB approximation error (Theorem 1). This result depends on the
 109 observation noise $\nu_{t,i} := z_{t,i} - \pi_t(\mathbf{x}_i)$ as perceived by a GP model with observations $z_{t,i} := \mathbb{I}[y_i \leq$
 110 $\tau_t]$, $i \in \{1, \dots, t\}$, to be sub-Gaussian when conditioned on the history. Hence, a few challenges
 111 originate from there. Firstly, the past observations in the vector $\mathbf{z}_t := [z_{t,i}]_{i=1}^t$ are changing across
 112 iterations, due to the update in τ_t . Secondly, the latent function π_t is stochastic, as the quantile τ_t
 113 depends on the current set of observations \mathbf{y}_t . Lastly, it is not very clear how to define a filtration
 114 for the resulting stochastic process such that the GP noise $\nu_{t,i}$ is sub-Gaussian. Nevertheless, as the
 115 number of observations increases, τ converges to a fixed value, making our asymptotic results valid.

116 D Experiment details

117 This section presents details of our experiments setup. We used PyTorch [11] for our implementation
 118 of batch BORE and BORE++, which we plan to make publicly available in the future.

119 D.1 Theory assessment

120 For this experiment, we generated a random classifier as an element of the RKHS \mathcal{H} defined by the
 121 kernel k as:

$$\pi^* := \sum_{i=1}^F \alpha_i k(\cdot, \mathbf{x}_i^*) \in \mathcal{H}, \quad (\text{D.1})$$

122 where $\{\mathbf{x}_i^*\}_{i=1}^F$ and the weights $\{\alpha_i\}_{i=1}^F$ were i.i.d. sampled from a unit uniform distribution $\mathcal{U}(0, 1)$,
 123 with $F := 5$. The norm of π^* is given by:

$$\|\pi^*\|_k = \sqrt{\boldsymbol{\alpha}_F^\top \mathbf{K}_F \boldsymbol{\alpha}_F}, \quad (\text{D.2})$$

124 where $\mathbf{K} := [\mathbf{x}_i^*, \mathbf{x}_j^*]_{i,j=1}^F \in \mathbb{R}^{F \times F}$ and $\boldsymbol{\alpha}_F := [\alpha_1, \dots, \alpha_F]^\top \in \mathbb{R}^F$. To ensure $\pi^*(\mathbf{x}) \leq 1$, we
 125 normalised the weights according to the classifier norm, i.e., $\boldsymbol{\alpha} := \frac{1}{\|\pi^*\|} \boldsymbol{\alpha}$, so that $\|\pi^*\| = 1$, and
 126 consequently $\pi^*(\mathbf{x}) \leq k(\mathbf{x}, \mathbf{x}) \|\pi^*\| = \|\pi^*\| = 1$, for all $\mathbf{x} \in \mathcal{X}$. The kernel was set as the squared
 127 exponential (RBF) with a length-scale of 0.1.

128 Given a threshold $\tau \in \mathbb{R}$, the objective function corresponding to π^* satisfies:

$$f(\mathbf{x}) := \tau - \Phi_\epsilon^{-1}(\pi^*(\mathbf{x})), \quad \forall \mathbf{x} \in \mathcal{X}. \quad (\text{D.3})$$

129 For this experiment, we set $\tau := 0$. Each trial had a different objective function generated for
 130 it. An example of classifier and objective function pair is presented in Figure 1b (main paper).
 131 Observations were composed as function evaluations corrupted by zero-mean Gaussian noise with
 132 variance $\sigma_\epsilon^2 := 0.01$.

133 The search space was configured as a finite set $\mathcal{X} := \{\mathbf{x}_i\}_{i=1}^{N_{\mathcal{X}}} \subset [0, 1]$ by sampling $N_{\mathcal{X}}$ points from
 134 a unit uniform distribution. The number of points in the search space was set as $N_{\mathcal{X}} := 100$. As the
 135 search space is finite, we also know $\gamma := p(y \leq \tau) = \int_{\mathcal{X}} \pi(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \frac{1}{N_{\mathcal{X}}} \sum_{\mathbf{x} \in \mathcal{X}} \pi^*(\mathbf{x})$.

136 Regarding algorithm hyper-parameters, although any upper bound $b \geq \|\pi^*\|$ would work for setting
 137 up β_t , BORE++ was configured with the RKHS norm π^* as the first term in the setting for β_t (see
 138 Theorem 1). To configure GP-UCB according to its theoretical settings [3, Thm. 1], we computed the
 139 RKHS norm of the resulting f in the RKHS. We can compute the norm of f as an element of \mathcal{H} by
 140 solving the following constrained optimisation problem:

$$\|f\|_k = \min_{h \in \mathcal{H}} \|h\|_k, \quad \text{s.t.} \quad h(\mathbf{x}) = f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}. \quad (\text{D.4})$$

Parameter	Setting
Optimiser	Adam
Batch size	64
Steps	100*

Table D.1: Stochastic gradient descent training settings for batch BORE. (*) For the Six-hump Camel problem, we applied 200 steps.

Parameter	Setting
Step size	0.001
Decay rate	0.9
Number of steps	1000*

Table D.2: SVGD settings for batch BORE. (*) For the Hartmann 3D problem, we used 500 steps.

141 As the search space is finite, the solution to this problem is available in closed form as:

$$\|f\|_k = \sqrt{\mathbf{f}_{\mathcal{X}}^T \mathbf{K}_{\mathcal{X}}^{-1} \mathbf{f}_{\mathcal{X}}}, \quad (\text{D.5})$$

142 where $\mathbf{f}_{\mathcal{X}} := [f(\mathbf{x})]_{\mathbf{x} \in \mathcal{X}} \in \mathbb{R}^{N_{\mathcal{X}}}$, and $\mathbf{K}_{\mathcal{X}} := [k(\mathbf{x}, \mathbf{x}')]_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}}$. We set $\delta := 0.1$. For both BORE++
 143 and GP-UCB, the information gain was recomputed at each iteration. Lastly, the regularisation factor
 144 λ was set as $\lambda := \sigma_{\epsilon}^2$ for GP-UCB and as $\lambda := 0.025$, which was found by grid search.

145 D.2 Global optimisation benchmarks

146 For each problem, all methods used 10 initial points uniformly sampled from the search space. As
 147 performance indicator, we measured the simple regret:

$$r_t^* := \min_{i \leq t} r_i = \min_{i \leq t} f(\mathbf{x}_i) - f(\mathbf{x}^*), \quad t \geq 1, \quad (\text{D.6})$$

148 as the global minimum of each of the considered benchmark functions is known. All objective
 149 function evaluations were provided free of noise to the algorithms.

150 Batch BORE was run with a percentile $\gamma := 0.25$, which was applied to estimate the empirical
 151 quantile τ at every iteration $t \in \{1, \dots, T\}$. The method’s classifier model was composed of a
 152 multilayer perceptron neural network model with 2 hidden layers of 32 units each, which was trained
 153 to minimise the binary cross-entropy loss. The activation function was set as the rectified linear
 154 unit (ReLU) with exception for the Hartmann 3D and the Six-hump Camel problem, which were
 155 run with an exponential linear unit (ELU), instead. Training for the neural networks was performed
 156 via stochastic gradient descent, whose settings are presented in Table D.1. SVGD was run applying
 157 Adadelta to configure its steps according to the settings in Table D.2. The SVGD kernel was set as
 158 the squared exponential (RBF) using the median trick to adjust its lengthscale [12].

159 LP-EI [13] was run using L-BFGS [14] to optimise its acquisition function. The optimisation settings
 160 were kept as the default for GPpyOpt [15].

161 The q -EI method [16] was run using the BoTorch implementation [17]. The acquisition function was
 162 optimised via multi-start optimisation with L-BFGS [14] using 10 random restarts. Monte Carlo
 163 integration for q -EI used 256 samples.

164 E Additional experiments

165 We compared batch BORE against the batch BORE++ algorithm on a synthetic optimisation problem
 166 with the Rosenbrock function. The dimensionality of the search space was varied. The cumulative
 167 regret curves for each algorithm are presented in Figure E.1.

168 Both algorithms were configured with a Bayesian logistic regression classifier applying random
 169 Fourier features [18] as feature maps based on the squared-exponential kernel. The number of
 170 features was set as 300, and the classifier was trained via expectation maximisation. Observations
 171 were corrupted by additive Gaussian noise with zero mean and a small noise variance $\sigma_{\epsilon}^2 = 10^{-4}$, and

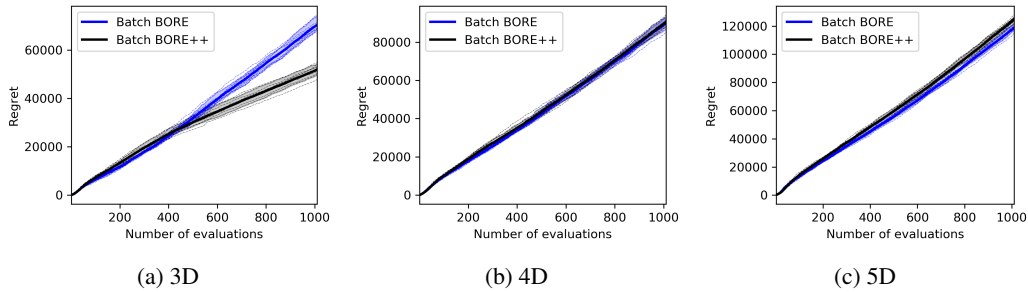


Figure E.1: BORE vs. BORE++ in the batch setting tested on the Rosenbrock function at varying search space dimensionalities. The plots compare the cumulative regret of each algorithm averaged over 10 runs. Shaded areas correspond to the 95% confidence interval.

172 each model was set accordingly. To demonstrate the practicality of the method, the UCB parameter
 173 for BORE++ was fixed at $\beta_t := 3$ across all iterations $t \geq 1$, instead of applying the theoretical setup.
 174 SVGD was configured as its second-order version [19] applying L-BFGS to adjust its steps [14].

175 As the results show in Figure E.1, batch BORE++ has a clear advantage over batch BORE in low
 176 dimensions. However, the performance gains become less obvious at higher dimensionalities and
 177 eventually deteriorate. One of the factors explaining this behaviour is that, as the dimensionality
 178 increases, uncertainty estimates become less useful. Distances between data points largely increase
 179 and affect the posterior variance estimates provided by translation-invariant kernels, such as the
 180 squared-exponential kernel our feature maps were based on. Other classification models may lead to
 181 different behaviours, and their investigation is left for future work.

182 References

- 183 [1] Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias Seeger. Gaussian Process
 184 Optimization in the Bandit Setting: No Regret and Experimental Design. In *Proceedings of the*
 185 *27th International Conference on Machine Learning (ICML 2010)*, pages 1015–1022, 2010.
- 186 [2] Sayak Ray Chowdhury and Aditya Gopalan. On Kernelized Multi-armed Bandits. In *Proceed-*
 187 *ings of the 34th International Conference on Machine Learning (ICML)*, Sydney, Australia,
 188 2017.
- 189 [3] Audrey Durand, Odalric-Ambrym Maillard, and Joelle Pineau. Streaming kernel regression with
 190 provably adaptive mean, variance, and regularization. *Journal of Machine Learning Research*,
 191 19(1):650–683, 2018.
- 192 [4] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A*
 193 *Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- 194 [5] James Raymond Munkres. *Topology: a first course*. Prentice Hall, Edgewood Cliffs, NJ, 1975.
- 195 [6] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. Stochastic gradient descent as
 196 approximate Bayesian inference. *Journal of Machine Learning Research*, 18, 2017.
- 197 [7] Daniel Russo and Benjamin Van Roy. An Information-Theoretic Analysis of Thompson
 198 Sampling. *Journal of Machine Learning Research (JMLR)*, 17:1–30, 2016.
- 199 [8] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
 200 generalization in neural networks. In *Advances in Neural Information Processing Systems*,
 201 Montreal, Canada, 2018.
- 202 [9] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin
 203 Ghahramani. Gaussian Process Behaviour in Wide Deep Neural Networks. In *International*
 204 *Conference on Learning Representations*, Vancouver, Canada, 2018. OpenReview.net.

- 205 [10] Steffen Grünewälder, Jean Yves Audibert, Manfred Opper, and John Shawe-Taylor. Regret
206 bounds for Gaussian process bandit problems. In *Journal of Machine Learning Research*,
207 volume 9, pages 273–280, 2010.
- 208 [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
209 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
210 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
211 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-
212 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché
213 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,
214 pages 8024–8035. Curran Associates, Inc., 2019.
- 215 [12] Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose Bayesian
216 inference algorithm. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- 217 [13] Javier Gonzalez, Zhenwen Dai, Philipp Hennig, and Neil D. Lawrence. Batch Bayesian
218 optimization via local penalization. In *International Conference on Artificial Intelligence and
219 Statistics (AISTATS)*, pages 648–657, Cadiz, Spain, 2016.
- 220 [14] Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyou Zhu. A limited memory algorithm
221 for bound constrained optimization. *SIAM Journal on scientific computing*, 16(5):1190–1208,
222 1995.
- 223 [15] The GPyOpt authors. GPyOpt: A Bayesian optimization framework in python. [http://
224 github.com/SheffieldML/GPyOpt](http://github.com/SheffieldML/GPyOpt), 2016.
- 225 [16] Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine
226 learning algorithms. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors,
227 *Advances in Neural Information Processing Systems 25*, pages 2951–2959. Curran Associates,
228 Inc., 2012.
- 229 [17] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham,
230 Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo
231 Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
- 232 [18] Ali Rahimi and Ben Recht. Random features for large-scale kernel machines. In *Advances in
233 Neural Information Processing (NIPS)*, 2007.
- 234 [19] Gianluca Detommaso, Tiangang Cui, Alessio Spantini, Youssef Marzouk, and Robert Scheichl.
235 A Stein variational Newton method. In *32nd Conference on Neural Information Processing
236 Systems (NeurIPS 2018)*, Montréal, Canada, 2018.