

Figure 4: PrBO on the 1D Branin function. The leftmost column shows the exponential prior. The other columns show the model and the log pseudo-posterior after 0 (RS only), 10, and 20 BO iterations. PrBO forgets the wrong prior on the local optimum and converges to the global optimum.

A Prior Forgetting

In this section, we show that PrBO can recover from a misleading prior, thanks to our predictive model and the t/β parameter in the pseudo-posterior computation Eq. 3. As BO progresses, the predictive model becomes more accurate and receives more weight, guiding optimization away from the wrong prior and towards better values of the objective function. Figure 4 shows PrBO on the 1D Branin function with an exponential prior. Columns (b), (c), and (d) show PrBO after $D + 1 = 2$ initial samples and 0, 10, 20 BO iterations, respectively. After initialization, as shown in Column (b), the pseudo-posterior is nearly identical to the exponential prior and guides PrBO towards the region of the space on the right, which is towards the local optimum. This happens until the predictive model becomes certain there will be no more improvement from sampling that region (Columns (c) and (d)). After that, the predictive model guides the pseudo-posterior towards exploring regions with high uncertainty. Once the global minimum region is found, the pseudo-posterior starts balancing exploiting the global minimum and exploring regions with high uncertainty, as shown in 4d (bottom). Notably, the pseudo-posterior after $x > 4$ falls to 0 in 4d (top), as the predictive model is certain there will be no improvement from sampling the region of the local optimum.

B EI Derivation

Here, we provide a full derivation of Eq. (5):

$$EI_{f_\gamma}(\mathbf{x}) := \int_{-\inf}^{\inf} \max(f_\gamma - y, 0) p(y|\mathbf{x}) dy = \int_{-\inf}^{f_\gamma} (f_\gamma - y) \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} dy. \quad (6)$$

As defined in Section 3.2, $p(y < f_\gamma) = \gamma$ and γ is a quantile of the observed objective values $\{y^{(i)}\}$. Then $p(\mathbf{x}) = \int_{\mathbb{R}} p(\mathbf{x}|y)p(y) dy = \gamma g(\mathbf{x}) + (1 - \gamma)b(\mathbf{x})$, where $g(\mathbf{x})$ and $b(\mathbf{x})$ are the posteriors introduced in Section 3.3. Therefore

$$\int_{-\inf}^{f_\gamma} (f_\gamma - y) p(\mathbf{x}|y)p(y) dy = g(\mathbf{x}) \int_{-\inf}^{f_\gamma} (f_\gamma - y) p(y) dy = \gamma f_\gamma g(\mathbf{x}) - g(\mathbf{x}) \int_{-\inf}^{f_\gamma} y p(y) dy, \quad (7)$$

so that finally

$$EI_{f_\gamma}(\mathbf{x}) = \frac{\gamma f_\gamma g(\mathbf{x}) - g(\mathbf{x}) \int_{-\inf}^{f_\gamma} y p(y) dy}{\gamma g(\mathbf{x}) + (1 - \gamma)b(\mathbf{x})} \propto \left(\gamma + \frac{b(\mathbf{x})}{g(\mathbf{x})}(1 - \gamma) \right)^{-1}. \quad (8)$$

C Experimental Setup

We use a combination of publicly available implementations for our predictive models. For our Gaussian Process (GP) model, we use GPy’s (GPy, since 2012) GP implementation with the Matérn5/2 kernel. We use different length-scales for each input dimensions, learned via Automatic Relevance Determination (ARD) (Neal, 2012). For our Random Forests (RF), we use scikit-learn’s RF implementation (Pedregosa et al., 2011). We set the fraction of features per split to 0.5, the minimum number of samples for a split to 5 and disable bagging. We also adapt our RF implementation to use the same split selection approach as Hutter et al. (2014).

For our constrained Bayesian Optimization (cBO) approach, we use scikit-learn’s RF classifier, trained on previously explored configurations, to predict the probability of a configuration being feasible. We then weight our EI acquisition function by this probability of feasibility, as proposed by Gardner et al. (2014). We normalize our EI acquisition function before considering the probability of feasibility, to ensure both values are in the same range. This cBO implementation is used in the Spatial use-case as in Nardi et al. (2019).

For all experiments, we set the model weight to $\beta = 10$ and the model quantile to $\gamma = 0.05$, see Appendices J and I. Before the main BO loop, PrBO is initialized by random sampling $D + 1$ points from the prior, where D is the number of input dimensions. We use the public implementation of Spearmin³, which uses 2 random samples for initialization. We set the bandwidth of our KDE priors to $100n^{-\frac{1}{D}}$, where D is the number of input dimensions, see Appendix H. We normalize our KDE priors before computing the pseudo-posterior, to ensure they are in the same range as our model. We also implement interleaving which randomly samples a point to explore during BO with 10% chance.

We optimize EI using a multi-start local search, similar to the one used in SMAC (Hutter et al., 2011). Namely, we start local searches on the 10 best points evaluated in previous BO iterations, on the 10 best performing points from a set of 10,000 random samples and on the 10 best performing points from 10,000 random samples drawn from the prior. To compute the neighbors of each of these 30 total points, we normalize the range of each objective to $[0, 1]$ and randomly sample four neighbors from a truncated Gaussian centered at the original value and with standard deviation $\sigma = 0.2$.

We use four synthetic benchmarks in our experiments.

Branin. The Branin function is a well-known synthetic benchmark for optimization problems (Dixon, 1978). The Branin function has two input dimensions and three global minima.

SVM. Hyperparameter-optimization benchmark in 2D based on Profet (Klein et al., 2019). This benchmark is generated by a generative meta-model built using a set of SVM classification models trained on 16 OpenML tasks. The benchmark has two input parameters, corresponding to SVM hyperparameters.

FC-Net. Hyperparameter and architecture optimization benchmark in 6D based on Profet. The FC-Net benchmark is generated by a generative meta-model built using a set of feed-forward neural networks trained on the same 16 OpenML tasks as the SVM benchmark. The benchmark has six input parameters corresponding to network hyperparameters.

XGBoost. Hyperparameter-optimization benchmark in 8D based on Profet. The XGBoost benchmark is generated by a generative meta-model built using a set of XGBoost regression models in 11 UCI datasets. The benchmark has eight input parameters, corresponding to XGBoost hyperparameters.

The search spaces for each benchmark are summarized in Table 1. For the Profet benchmarks, we report the original ranges and whether or not a log scale was used. However, in practice, Profet’s generative model transforms the range of all hyperparameters to a linear $[0, 1]$ range. We use Emukit’s public implementation for these benchmarks (Paley et al., 2019).

D SVM Regret Comparison

In addition to the experiments in Section 4.2, we show the performance of PrBO on the SVM benchmark. Figure 5 shows a log regret comparison of PrBO, Spearmin, Prior Sampling and $10,000 \times$ RS. We note that the results are similar to the other benchmarks in Figure 3. Namely, PrBO

³<https://github.com/HIPS/Spearmin>

Table 1: Search spaces for our synthetic benchmarks. For the Profet benchmarks, we report the original ranges and whether or not a log scale was used.

Benchmark	Parameter name	Parameter values	Log scale
Branin	x_1	$[-5, 10]$	-
	x_2	$[0, 15]$	-
SVM	C	$[e^{-10}, e^{10}]$	✓
	γ	$[e^{-10}, e^{10}]$	✓
FCNet	learning rate	$[10^{-6}, 10^{-1}]$	✓
	batch size	$[8, 128]$	✓
	units layer 1	$[16, 512]$	✓
	units layer 2	$[16, 512]$	✓
	dropout rate l1	$[0.0, 0.99]$	-
	dropout rate l2	$[0.0, 0.99]$	-
XGBoost	learning rate	$[10^{-6}, 10^{-1}]$	✓
	gamma	$[0, 2]$	-
	L1 regularization	$[10^{-5}, 10^3]$	✓
	L2 regularization	$[10^{-5}, 10^3]$	✓
	number of estimators	$[10, 500]$	-
	subsampling	$[0.1, 1]$	-
	maximum depth	$[1, 15]$	-
	minimum child weight	$[0, 20]$	-

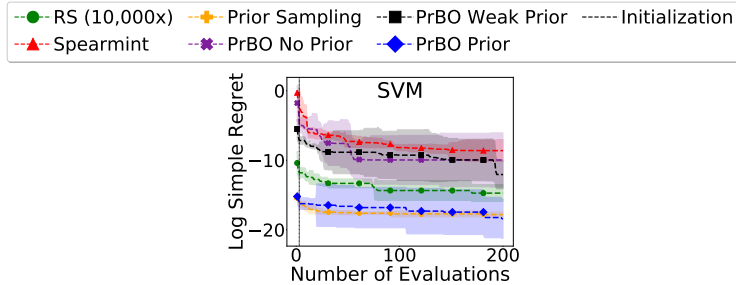


Figure 5: Log regret comparison of PrBO with and without priors, $10,000\times$ RS, and Spearmint (mean \pm one std on 5 repetitions). We run the benchmark for 200 iterations.

with a strong prior outperforms RS and spearmint. PrBO also outperforms Spearmint with a weak prior and even with a uniform prior.

E Spatial Real-world Application

Spatial (Koeplinger et al., 2018) is a programming language and corresponding compiler for the design of application accelerators on reconfigurable architectures, e.g. field-programmable gate arrays (FPGAs). These reconfigurable architectures are a type of logic chip that can be reconfigured via software to implement different applications. *Spatial* provides users with a high-level of abstraction for hardware design, so that they can easily design their own applications on FPGAs. It allows users to specify parameters that do not change the behavior of the application, but impact the runtime and resource-usage (e.g. logic units) of the final design. During compilation, the *Spatial* compiler estimates the ranges of these parameters and estimates the resource-usage and runtime of the application for different parameter values. These parameters can then be optimized during compilation in order to achieve the design with the fastest runtime. We fully integrate PrBO as a pass in *Spatial*'s compiler, so that *Spatial* can automatically use PrBO for the optimization during compilation. This enables *Spatial* to seamlessly call PrBO during the compilation of any new application to guide the search towards the best design on an application-specific basis.

In our experiments, we introduce for the first time the automatic optimization of three *Spatial* real-world applications, namely, 7D shallow and deep CNNs, and a 10D molecular dynamics grid application. Previous work by Nardi et al. (2019) had applied automatic optimization of *Spatial*

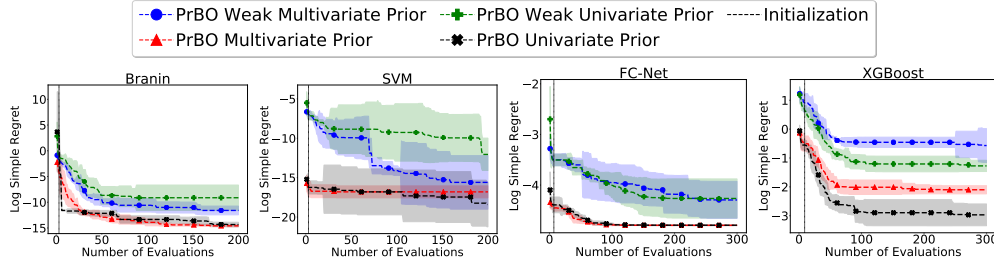


Figure 7: Log regret comparison of PrBO with multivariate and univariate KDE priors (mean \pm one std on 5 repetitions). We run the benchmarks for $100D$ iterations, capped at 300.

We compare the performance of PrBO to RS, manual optimization, and HyperMapper (Nardi et al., 2019), the current state-of-the-art BO solution for *Spatial*. For a fair comparison, both PrBO and HyperMapper use RFs as predictive probabilistic model. The priors for these *Spatial* applications take the form of a list of probabilities, containing the probability of each ordinal or categorical value being good. Each benchmark also has a default configuration, which ensures all methods start with at least one feasible configuration. The priors, expert configuration, and default configurations for these benchmarks were provided once by an unbiased *Spatial* developer, who is not an author of this paper, and kept unchanged during the entire project. The search space, priors, and expert configuration for each application are presented in Tables 2, 3, and 4.

Figure 6 shows the log regret on the *Spatial* benchmarks. PrBO vastly outperforms RS in all benchmarks; notably, RS does not improve over the default configuration in MD Grid. PrBO is also able to leverage the expert’s prior and outperform the expert’s configuration in all benchmarks ($2.68\times$, $1.06\times$, and $10.4\times$ speedup for shallow CNN, deep CNN, and MD Grid, respectively). Compared to HyperMapper, PrBO achieves better performance in the MD Grid benchmark ($1.28\times$ speedup). For context, this is a significant improvement in the FPGA field, where a 10% improvement could qualify for acceptance in a top-tier conference. In the CNN benchmarks, PrBO converges to the minima regions faster than HyperMapper ($1.58\times$ and $1.4\times$ faster for shallow and deep respectively). Thus, PrBO leverages the best of both worlds (the expert’s prior knowledge and BO) to provide a new state of the art for *Spatial*.

F Multivariate Prior Comparison

Figure 7 shows a log regret comparison of PrBO with univariate and multivariate KDE priors. We show results for univariate and multivariate versions of our weak and strong KDE priors. We use the best $10D$ points out of $1,000D$ and $10,000,000D$ randomly sampled points to create our weak and strong priors, respectively. We use the same points to create the univariate and multivariate priors. We recall that the goal of these synthetic priors is to have an unbiased prior for our experiments, whereas manual priors would be biased by our own expertise of these benchmarks. In practice, users will manually define these priors without needing additional experiments.

In all cases PrBO achieves similar performance with univariate and multivariate priors. For the Branin and SVM benchmarks, the weak multivariate prior leads to slightly better results than the weak univariate prior. However, the difference is small, in the order of 10^{-4} and 10^{-6} , respectively.

Surprisingly, for the XGBoost benchmark, the univariate version for both the weak and strong priors lead to better results than their respective multivariate counterparts, though, once again, the difference in performance is small, around 0.2 and 0.03 for the weak and strong prior, respectively, whereas the XGBoost benchmark can reach values as high as $f(\mathbf{x}) = 600$. Our hypothesis is that this difference comes from the bandwidth estimator ($100n^{-\frac{1}{D}}$), which leads to larger bandwidths, consequently, smoother priors, when a multivariate prior is constructed.

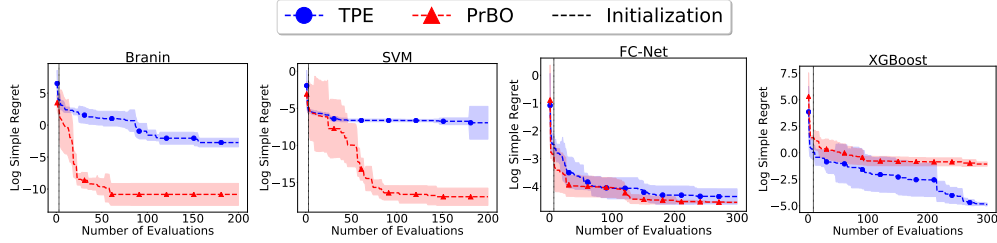


Figure 8: Log regret comparison of PrBO and TPE (mean +/- one std on 5 repetitions). We run the benchmarks for 100D iterations, capped at 300.

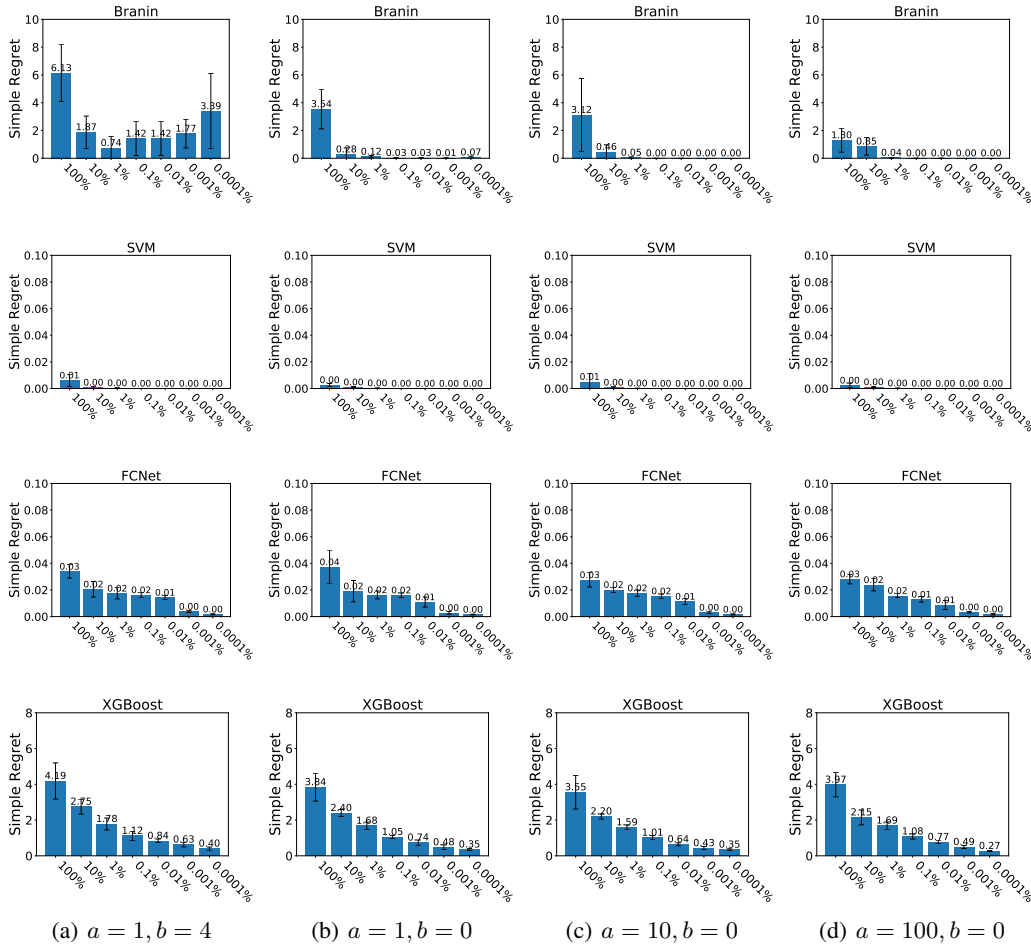


Figure 9: Simple regret of sampling from the prior with different priors for our synthetic benchmarks (mean +/- one std on 5 repetitions). A more informative prior gives better results in all benchmarks.

G Comparison to TPE

We compare PrBO to the TPE approach of Bergstra et al. (2011) on our four synthetic benchmarks. We use Hyperopt’s implementation⁴ of TPE, which defines priors as one of a list of supported distributions, including Uniform, Normal, and Lognormal distributions. Since it is not possible to input the KDE priors introduced in Section 4.1 into the TPE algorithm, we instead use manually defined priors in

⁴<https://github.com/hyperopt/hyperopt>

the format supported by the Hyperopt implementation. We note that this is straightforward in PrBO, as PrBO supports any form of probability distribution as a prior. We are then able to perform a fair comparison between the two approaches that use the same exact prior.

We define the prior for each input parameter as a Gaussian distribution with mean at the optimum and with standard deviation equal to half of the parameters range. For the Branin prior, we arbitrarily choose one of the optima, i.e., the $(\pi, 2.275)$ optimum. For the Profet benchmarks, we use the minimum out of 10,000,000 random samples as an approximation of the optimum. We note that using Hyperopt’s Gaussian priors leads to an unbounded search space, which sometimes leads TPE to suggest parameter configurations outside the allowed parameter range. To prevent these values from being evaluated, we convert values outside the parameter range to be equal to the upper or lower range limit, depending on which limit was exceeded.

Figure 8 shows a log regret comparison between PrBO and TPE on our four synthetic benchmarks. PrBO outperforms TPE in three out of four benchmarks, namely, Branin, SVM, and FCNet. We note, however, that the good performance of TPE on XGBoost may be an artifact of the approach of clipping values to its maximal or minimal values as mentioned above. In fact, the clipping nudges TPE towards promising configurations in this case, since XGBoost has low function value near the edges of the search space. Overall, the better performance of PrBO is expected, since PrBO is able to combine prior knowledge with more sample-efficient surrogates, which leads to better performance.

H Prior Bandwidth Selection

We show the effect of different bandwidth sizes on the univariate KDE prior. For that, we compare the performance of sampling from the prior and PrBO with different bandwidth sizes. We use scipy’s Gaussian KDE implementation and modify its bandwidth size with four variations of Scott’s Rule $an^{-\frac{1}{D+b}}$. We experiment with $a = 1, b = 4$ (scipy’s default); $a = 1, b = 0$; $a = 10, b = 0$; and $a = 100, b = 0$. Note that larger values for a and smaller values for b lead to smaller bandwidths. For each bandwidth size, we show results for an array of varying quality priors. We select a constant $10D$ points in each prior and vary the size of the uniform random sample dataset. We follow the following rule: we use the best performing $10D$ samples to create the prior from a random sample dataset size of $10D \frac{100}{x}$; we refer to this prior as $x\%$. We experiment with dataset sizes varying from $10D$ to $10^7 D$.

Figure 9 shows the performance of purely sampling from the prior. We note that, in most cases, using a larger dataset leads to better results. This is expected, sampling more points means we find more points near the optima and, therefore, the prior will be built with points closer to the optima. Likewise, we note that smaller bandwidths often lead to better results, especially as more points are sampled. This is also expected, since a smaller bandwidth means the prior distribution will be more peaked around the optima. However, there are a couple of exceptions to these trends. First, for the Branin, sampling more points does not lead to a better prior when we use $a = 1, b = 4$, this is likely because the multiple minima of the Branin and the bigger bandwidth lead the prior to be oversmoothed, missing the peaks near the optima. Second, smaller bandwidths do not always lead to better performance for smaller random sample datasets. This happens because we find points farther from the optima in these datasets and end up computing priors peaked at points that are farther from the optima, i.e., our priors become misleading. The effects of these misleading priors can be especially noticed for the 100% random samples dataset. Based on these results, we set our KDE priors to $100n^{-\frac{1}{D}}$, where D is the number of input dimensions.

Figure 10 shows the performance of PrBO for different priors. The same observations from Figure 9 hold here. Namely, sampling more points and using smaller bandwidths lead to better performance. Also, the 100% dataset once again leads to inconsistent results, since it is a misleading prior for PrBO. Based on these results, we use the smallest bandwidth and largest dataset in our experiments, i.e. $a = 100, b = 0$, and 0.0001%. Intuitively, this is a reasonable choice, since these priors will be our closest approximation to an ideal prior that is centered exactly at the optima, where sampling from the prior always leads to the optimum. Our results in Figures 9 and 10 shows that this combination leads to the best results in all benchmarks, as expected.

Figure 11 shows a performance comparison between PrBO and sampling from the prior. For these results, we use $a = 100, b = 0$ and compare the regret of PrBO and sampling from the prior for

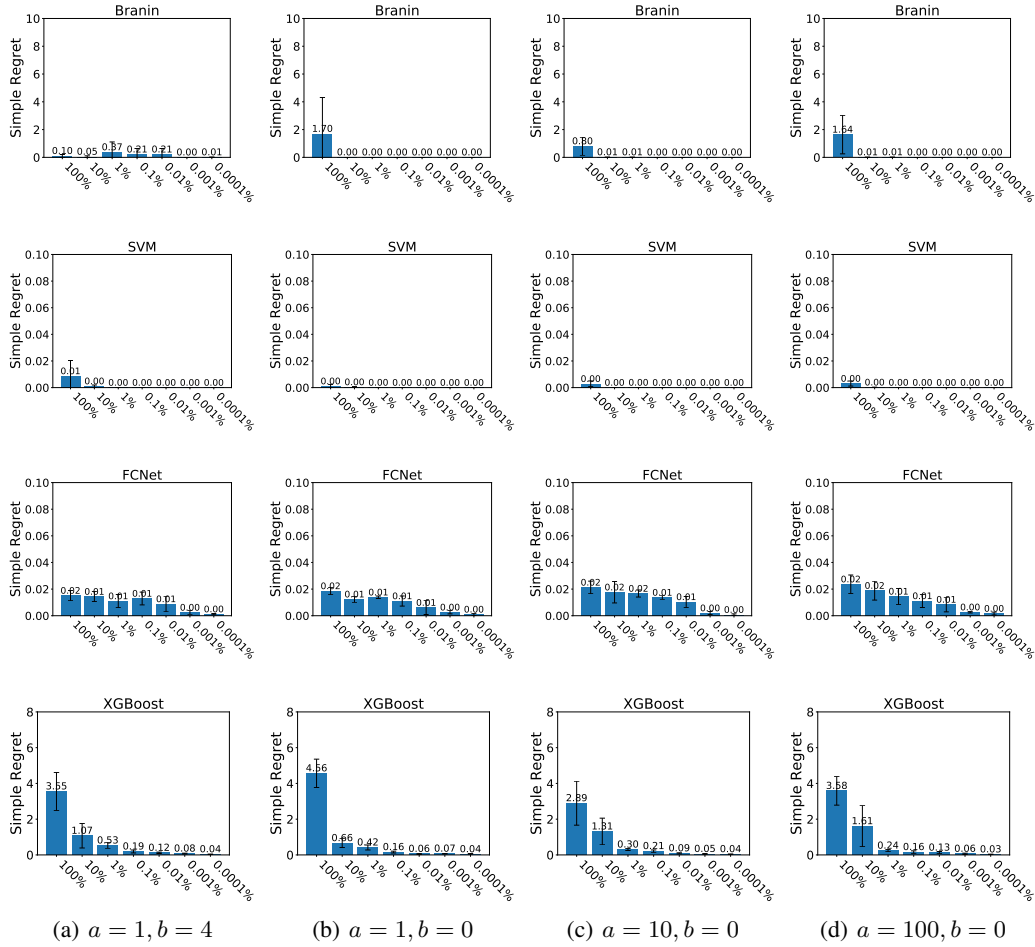


Figure 10: Simple regret of PrBO with different priors for our synthetic benchmarks (mean +/- one std on 5 repetitions). A more informative prior gives better results in all benchmarks.

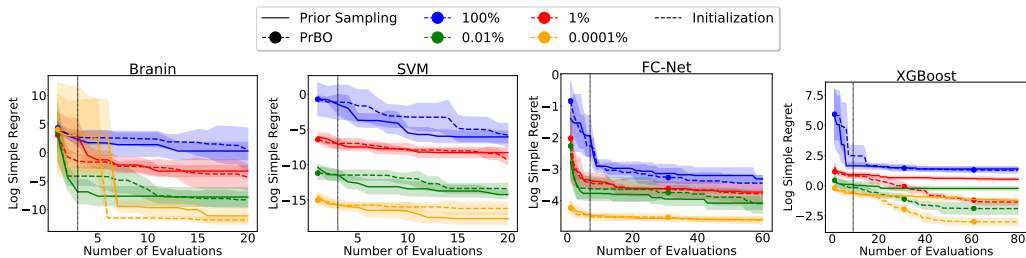


Figure 11: Log simple regret comparison between PrBO and sampling from the prior. The shaded lines are mean +/- one std error.

different dataset sizes. PrBO performs better for nearly all dataset sizes and benchmarks. This is expected as PrBO complements the prior with its probabilistic model, learning which regions within the prior are better to explore and also recovering from misleading priors. There are two exceptions on the SVM benchmark, where sampling from the prior performs slightly better for 0.01% and 0.0001% datasets. We note, however, that the difference in performance is extremely small, in the order of 10^{-5} and 10^{-7} , respectively.

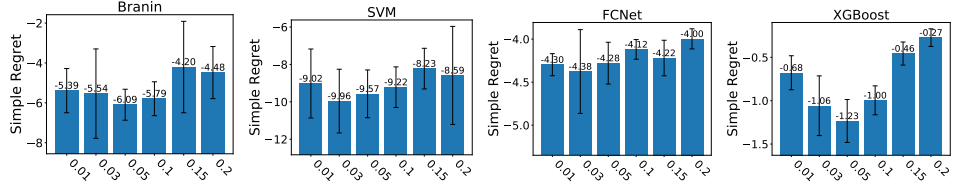


Figure 12: Comparison of PrBO with a weak KDE prior and different values for the γ hyperparameter on our four synthetic benchmarks (mean \pm one std on 5 repetitions). We run PrBO with a budget of $10D$ function evaluations.

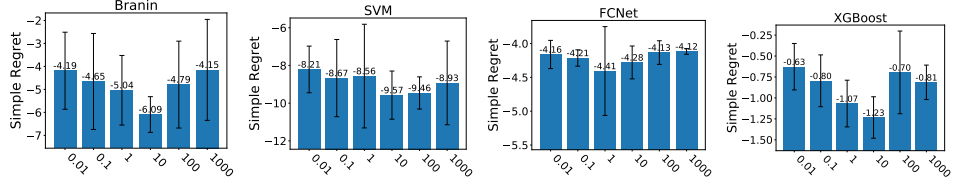


Figure 13: Comparison of PrBO with a weak KDE prior and different values for the β hyperparameter on our four synthetic benchmarks (mean \pm one std on 5 repetitions). We run PrBO with a budget of $10D$ function evaluations.

I γ -Sensitivity Study

We show the effect of the γ hyperparameter introduced in Section 3.2 for the quantile identifying the points considered to be good. To show this, we compare the performance of PrBO with a weak KDE prior and different γ values. We use our weak prior as it leads to greater variation in performance, which helps to visualize better the impact of the γ hyperparameter. For all experiments, we initialize PrBO with $D + 1$ random samples and then run PrBO until it reaches $10D$ function evaluations. For each γ value, we run PrBO five times and report mean and standard deviation.

Figure 12 shows the results of our comparison. We first note that values near the lower and higher extremes lead to degraded performance, this is expected, since these values will lead to an excess of either exploitation or exploration. Further, we note that PrBO achieves similar performance for all values of γ , however, $\gamma = 0.03$ and $\gamma = 0.05$ consistently lead to better performance, with $\gamma = 0.05$ usually leading to lower deviation.

J β -Sensitivity Study

We show the effect of the β hyperparameter introduced in Section 3.3 for controlling the influence of the prior over time. To show the effects of β , we compare the performance of PrBO with a weak KDE prior and different β values on our four synthetic benchmarks. We use our weak prior as it leads to greater variation in performance, which helps to visualize better the impact of the β hyperparameter. For all experiments, we initialize PrBO with $D + 1$ random samples and then run PrBO until it reaches $10D$ function evaluations. For each β value, we run PrBO five times and report mean and standard deviation.

Figure 13 shows the results of our comparison. We note that values of β that are too low (near 0.01) or too high (near 1,000) lead to lower performance. This shows that putting too much emphasis on the model or the prior will lead to degraded performance, as expected. Further, we note that $\beta = 10$ lead to the best performance in three out of our four benchmarks. This result is reasonable, as $\beta = 10$ means PrBO will put more emphasis on the prior in early iterations, when the predictive model is still not accurate, and slowly shift towards putting more emphasis on the model as the model sees more data and becomes more accurate.