
Prior-guided Bayesian Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

While Bayesian Optimization (BO) is a very popular method for optimizing expensive black-box functions, it fails to leverage the knowledge of domain experts. This causes BO to waste function evaluations on bad design choices (e.g., machine learning hyperparameters) that the expert already knows to work poorly. To address this issue, we introduce Prior-guided Bayesian Optimization (PrBO). PrBO allows users to transfer their knowledge into the optimization process in the form of priors about which parts of the input space will yield the best performance, rather than BO’s standard priors over functions (which are much less intuitive for users). PrBO then combines these priors with BO’s standard probabilistic model to form a pseudo-posterior used to select which points to evaluate next. We show that PrBO is around 12x faster than state-of-the-art methods without user priors and $10,000\times$ faster than random search on a common suite of benchmarks. PrBO also converges faster even if the user priors are not entirely accurate and robustly recovers from misleading priors.

1 Introduction

Bayesian Optimization (BO) is a data-efficient method for the joint optimization of design choices that gained great popularity in recent years. It is impacting a wide range of areas, including hyperparameter optimization (Snoek et al., 2012; Falkner et al., 2018), AutoML (Feurer et al., 2015a; Hutter et al., 2018), Computer Go (Chen et al., 2018), hardware design (Koeplinger et al., 2018; Nardi et al., 2019), and many others. It promises greater automation so as to increase both product quality and human productivity. As a result, BO is also established in many large tech companies, e.g., with Google Vizier (Golovin et al., 2017) and Facebook BoTorch (Balandat et al., 2019).

Nevertheless domain experts often have substantial prior knowledge that standard BO cannot incorporate. Users can incorporate prior knowledge by narrowing the search space; however, this type of hard prior can lead to poor performance by missing important regions. BO also supports a prior over functions $p(f)$, e.g., via a kernel function. However, this is not the prior experts have: users often know which ranges of hyperparameters tend to work best, and are able to specify a probability distribution $p_{\text{best}}(x)$ to quantify these priors. E.g., many users of the Adam optimizer (Kingma & Ba, 2015) know that its best learning rate is often in the vicinity of $1e-3$. Similarly, Navruzyan et al. (2019) derived neural network hyperparameter priors for image datasets based on their experience with five datasets. In these cases, users know potentially good values for a new application, but cannot be certain about them.

As a result, many competent users instead revert to manual search, which can fully incorporate their prior knowledge. A recent survey showed that most NeurIPS 2019 and ICLR 2020 papers that reported having tuned hyperparameters used manual search, with only a very small fraction using BO (Bouthillier & Varoquaux, 2020). In order for BO to be adopted widely, and help facilitate faster progress in the ML community by tuning hyperparameters faster and better, it is therefore crucial to devise a method that allows experts to fully transfer their knowledge into the optimization. In this

paper, we introduce Prior-guided Bayesian Optimization (PrBO), a novel BO variant that allows users to transfer their knowledge into BO in the form of priors. PrBO then combines this prior knowledge with its learning model in order to learn better and faster where to find promising hyperparameter configurations. Our technical contributions with PrBO are:

1. For the first time, user prior knowledge can be combined with standard BO probabilistic models, such as Gaussian Processes (GPs), Random Forests (RFs), and Bayesian Neural Networks (BNNs), leading to improved performance.
2. Unlike previous approaches, PrBO is flexible w.r.t. how the prior is defined, allowing previously impossible-to-inject (e.g. decay and exponential) priors.
3. PrBO gives more importance to the model as iterations progress, gradually forgetting the prior knowledge and ensuring robustness against misleading priors.

We demonstrate the effectiveness of PrBO on a common suite of benchmarks, showing that accurate prior knowledge helps PrBO to achieve similar performance to current state-of-the-art on average $12\times$ faster. PrBO also achieves better final performance in all but one of the benchmarks tested.

2 Background: Tree-structured Parzen Estimator

Whereas the standard probabilistic model in BO directly models $p(y|\mathbf{x})$, the Tree-structured Parzen Estimator (TPE) approach of Bergstra et al. (2011) models $p(\mathbf{x}|y)$ and $p(y)$ instead. This is done by constructing two parametric densities, $g(\mathbf{x})$ and $l(\mathbf{x})$, which are computed using the observations with function value above and below a given threshold, respectively. The separating threshold y^* is defined as a quantile of the observed function values. TPE then defines $p(\mathbf{x}|y)$ as:

$$p(\mathbf{x}|y) = l(\mathbf{x})I(y < y^*) + g(\mathbf{x})(1 - I(y < y^*)), \quad (1)$$

where $I(y < y^*)$ is 1 when $y < y^*$ and 0 otherwise. The parametrization of the generative model $p(\mathbf{x}, y) = p(\mathbf{x}|y)p(y)$ facilitates the computation of EI as it leads to $EI_{y^*}(\mathbf{x}) \propto l(\mathbf{x})/g(\mathbf{x})$ and, thus, $\arg \max_{\mathbf{x} \in \mathcal{X}} EI_{y^*}(\mathbf{x}) = \arg \max_{\mathbf{x} \in \mathcal{X}} l(\mathbf{x})/g(\mathbf{x})$.

3 Bayesian Optimization with Priors

We propose a BO approach dubbed PrBO that allows field experts to transfer prior knowledge into the optimization in the form of priors. PrBO combines this user-defined prior with a probabilistic model that captures the likelihood of the observed data $(\mathbf{x}_i, y_i)_{i=1}^n$. PrBO is independent from the probabilistic model being used, i.e., it can be freely combined with, e.g., GPs, RFs or BNNs.

3.1 Priors

PrBO allows users to transfer prior knowledge into BO. This is done via a prior distribution that informs where in the input space \mathcal{X} we expect to find good $f(\mathbf{x})$ values. A point is considered “good” if it leads to low function values. We denote the prior distribution $P_g(\mathbf{x})$, where g denotes that this is a prior on good points and $\mathbf{x} \in \mathcal{X}$ is a given point. Similarly, we define a prior on where in the input space we expect to have “bad” points. Although we could have a user-defined probability distribution $P_b(\mathbf{x})$, we aimed to keep the load on users low and thus, for simplicity, compute $P_b(\mathbf{x}) = 1 - P_g(\mathbf{x})$ ($P_g(\mathbf{x})$ is normalized to $[0, 1]$ by min-max scaling before computing $P_b(\mathbf{x})$).¹

In practice, \mathbf{x} contains several dimensions but it is difficult for experts to provide a multivariate distribution $P_g(\mathbf{x})$. Users can easily specify, e.g., draw, a univariate or bivariate probability distribution for continuous dimensions or provide a list of probabilities for discrete dimensions. In PrBO, users are free to define a complex multi-variate distribution, but we expect the standard use case to be that users only want to specify univariate distributions, implicitly assuming a prior that factors as $P_g(\mathbf{x}) = \prod_{i=1}^D P_g(x_i)$, where D is the number of dimensions in \mathcal{X} , x_i is the i -th input dimension of \mathcal{X} . We show examples of continuous and discrete priors in Appendices A and E, respectively. We use

¹We note that for continuous spaces, this $P_b(\mathbf{x})$ is not a probability distribution, and therefore only a pseudo-prior, as it does not integrate to 1. For discrete spaces, we normalize $P_b(\mathbf{x})$ so that it sums to 1 and therefore is a proper probability distribution and prior.

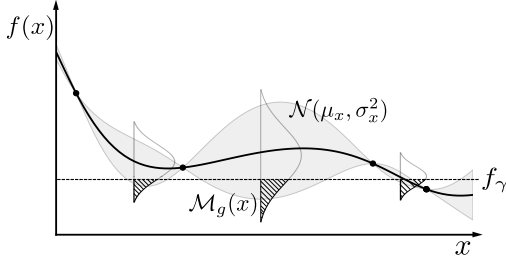


Figure 1: Our model is composed by a probabilistic model and the probability of improving over the threshold f_γ , i.e., right tail of the Gaussian. The black curve is the probabilistic model’s mean and the shaded area is the model’s variance.

Algorithm 1 PrBO Algorithm. \mathcal{D} keeps track of all function evaluations so far: $(\mathbf{x}_i, y_i)_{i=1}^t$.

```

1: Input: Input space  $\mathcal{X}$ , user-defined prior
   distributions  $P_g(\mathbf{x})$  and  $P_b(\mathbf{x})$ , quantile  $\gamma$ 
   and BO budget  $B$ .
2: Output: Optimized point  $\mathbf{x}_{inc}$ .
3:  $\mathcal{D} \leftarrow \text{Initialize}(\mathcal{X})$ 
4: for  $t = 1$  to  $B$  do
5:    $\mathcal{M}_g(\mathbf{x}) \leftarrow \text{fit\_model\_good}(\mathcal{D})$ 
6:    $\mathcal{M}_b(\mathbf{x}) \leftarrow \text{fit\_model\_bad}(\mathcal{D})$ 
7:    $g(\mathbf{x}) \leftarrow P_g(\mathbf{x}) \cdot \mathcal{M}_g(\mathbf{x})^{\frac{t}{\beta}}$ 
8:    $b(\mathbf{x}) \leftarrow P_b(\mathbf{x}) \cdot \mathcal{M}_b(\mathbf{x})^{\frac{t}{\beta}}$ 
9:    $\mathbf{x}_t \in \arg \max_{\mathbf{x} \in \mathcal{X}} EI_{f_\gamma}(\mathbf{x})$ 
10:   $y_t \leftarrow f(\mathbf{x}_t)$ 
11:   $\mathcal{D} = \mathcal{D} \cup (\mathbf{x}_t, y_t)$ 
12: end for
13:  $\mathbf{x}_{inc} \leftarrow \text{ComputeBest}(\mathcal{D})$ 
14: return  $\mathbf{x}_{inc}$ 

```

factorized priors in our experiments to mimic what we expect most users will provide. In Appendix F, we show that these factorized priors lead to similar performance compared to multivariate priors.

3.2 Model

Whereas the standard probabilistic model in BO, e.g., a GP, quantifies $p(y|\mathbf{x})$ directly, that model is hard to combine with the user-defined prior $P_g(\mathbf{x})$. We therefore introduce a method to translate the standard probabilistic model $p(y|\mathbf{x})$ into a model that is easier to combine with this prior. Similar to the TPE work (see Sec. 2), our model combines $p(\mathbf{x}|y)$ and $p(y)$ instead of directly modeling $p(y|\mathbf{x})$.

The computation we perform for this translation is to quantify the probability that a given input \mathbf{x} is “good” under our standard probabilistic model $p(y|\mathbf{x})$. As in TPE, we define settings as “good” if their observed y -value is below a certain quantile γ of the observed function values (so that $p(y < f_\gamma) = \gamma$). We in addition exploit the fact that our standard probabilistic model $p(y|\mathbf{x})$ has a Gaussian form, and under this Gaussian prediction we can compute the probability $\mathcal{M}_g(\mathbf{x})$ of the function value lying below a certain quantile using the standard closed-form formula for Probability of Improvement (PI, Kushner (1964)):

$$\mathcal{M}_g(\mathbf{x}) = p(f(\mathbf{x}) < f_\gamma | \mathbf{x}) = \Phi\left(\frac{f_\gamma - \mu_{\mathbf{x}}}{\sigma_{\mathbf{x}}}\right), \quad (2)$$

where $\mu_{\mathbf{x}}$ and $\sigma_{\mathbf{x}}$ are the mean and standard deviation of the probabilistic model at \mathbf{x} , and Φ is the standard normal CDF, see Figure 1. Likewise, we compute a probability $\mathcal{M}_b(\mathbf{x})$ of \mathbf{x} being bad.

3.3 Pseudo-posterior

PrBO combines the prior in Section (3.1) and the model in Eq. (2) into a pseudo-posterior. It represents the updated beliefs on where we can find good points, based on the prior and data that has been observed. The pseudo-posterior is computed as the product of the prior and the model:

$$g(\mathbf{x}) \propto P_g(\mathbf{x}) \mathcal{M}_g(\mathbf{x})^{\frac{t}{\beta}}, \quad (3)$$

where t is the current iteration, β is an optimization hyperparameter, $\mathcal{M}_g(\mathbf{x})$ is defined in Eq. (2), and $P_g(\mathbf{x})$ is the prior defined in Sec 3.1, rescaled to $[0, 1]$. We note that this pseudo-posterior is not normalized, but this suffices for PrBO to determine the next \mathbf{x}_t as the normalization constant cancels out (c.f. Section 3.4). Since $g(\mathbf{x})$ is not normalized and we add a t/β exponent to Eq. 3, we refer to $g(\mathbf{x})$ as a pseudo-posterior, to emphasize that it is not a standard posterior probability distribution.

107 The t/β fraction in Eq. (3) controls how much weight is given to the model. As the optimization
 108 progresses, more weight is given to the model over the prior. Intuitively, we put more emphasis on
 109 the model as it observes more data and becomes more accurate. We do this under the assumption
 110 that the model will eventually be better than the user at predicting where to find good points. This
 111 also allows to recover from misleading priors as we show in Appendix A; similar to, and inspired by
 112 Bayesian models, the data ultimately washes out the prior. The β hyperparameter defines the balance
 113 between prior and model, with higher β values giving more importance to the prior and requiring
 114 more data to overrule it.

115 We note that computing Equation (3) directly can lead to numerical issues. Namely, the pseudo-
 116 posterior can reach extremely low values if the prior and model probabilities are low, especially as
 117 t/β grows. To prevent this, in practice, PrBO uses the logarithm of the pseudo-posterior instead:
 118 $\log(g(\mathbf{x})) \propto \log(P_g(\mathbf{x})) + \frac{t}{\beta} \cdot \log(\mathcal{M}_g(\mathbf{x}))$. Once again, we also define an analogous pseudo-
 119 posterior distribution on bad \mathbf{x} : $b(\mathbf{x})$, and use these quantities to define a density model $p(\mathbf{x}|y)$:

$$p(\mathbf{x}|y) \propto \begin{cases} g(\mathbf{x}) & \text{if } y < f_\gamma \\ b(\mathbf{x}) & \text{if } y \geq f_\gamma. \end{cases} \quad (4)$$

120 3.4 Acquisition Function

121 We adopt the EI formulation used in (Bergstra et al., 2011) by replacing their Adaptive Parzen
 122 Estimators with our computation of the pseudo-posterior in Eq. (3). Namely, we compute EI as:

$$EI_{f_\gamma}(\mathbf{x}) := \int_{-\inf}^{\inf} \max(f_\gamma - y, 0) p(y|\mathbf{x}) dy = \int_{-\inf}^{f_\gamma} (f_\gamma - y) \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} dy \propto \left(\gamma + \frac{b(\mathbf{x})}{g(\mathbf{x})}(1 - \gamma) \right)^{-1}. \quad (5)$$

123 The full derivation of Eq. (5) is shown in Appendix B. Eq. (5) shows that to maximize improvement
 124 we would like points \mathbf{x} with high probability under $g(\mathbf{x})$ and low probability under $b(\mathbf{x})$, i.e.,
 125 minimizing the ratio $b(\mathbf{x})/g(\mathbf{x})$. We note that the point that minimizes the ratio for our unnormalized
 126 pseudo-posteriors will be the same that minimizes the ratio for the normalized pseudo-posterior and,
 127 thus, the computation of the normalized pseudo-posteriors is unnecessary.

128 The dynamics of PrBO can be understood in terms of the following proposition:

129 **Proposition 1** *Given f_γ , $P_g(\mathbf{x})$, $P_b(\mathbf{x})$, $\mathcal{M}_g(\mathbf{x})$, $\mathcal{M}_b(\mathbf{x})$, $g(\mathbf{x})$, $b(\mathbf{x})$, $p(\mathbf{x}|y)$, and β as above, then*

$$\lim_{t \rightarrow \infty} \arg \max_{\mathbf{x} \in \mathcal{X}} EI_{f_\gamma}(\mathbf{x}) = \lim_{t \rightarrow \infty} \arg \max_{\mathbf{x} \in \mathcal{X}} \mathcal{M}_g(\mathbf{x}),$$

130 *where $\mathcal{M}_g(\mathbf{x})$ and EI_{f_γ} are as defined in Eqs. 2 and 5, respectively.*

131 In early BO iterations the prior will have a predominant role, but in later BO iterations the model will
 132 grow more important, and as Proposition 1 shows, if PrBO is run long enough the prior washes out
 133 and PrBO *only* trusts the probabilistic model informed by the data.

134 3.5 Putting It All Together

135 Algorithm 1 shows the PrBO algorithm, based on the components defined in the previous sections. In
 136 Line 3, PrBO starts with a design of experiments (DoE) phase, where it randomly samples a number
 137 of points from the user-defined prior $P_g(\mathbf{x})$. After initialization, the BO loop starts at Line 4. In each
 138 loop iteration, PrBO fits the probabilistic model on the previously evaluated points (lines 5 and 6)
 139 and computes the pseudo-posteriors $g(\mathbf{x})$ and $b(\mathbf{x})$ (lines 7 and 8 respectively). The EI acquisition
 140 function is computed next, using the pseudo-posteriors, and the point that maximizes EI is selected as
 141 the next point to evaluate at line 9. The black-box function evaluation is performed at Line 10. This
 142 BO loop is repeated for a pre-defined number of iterations, according to the user-defined budget B .

143 4 Experiments

144 We implement both GPs and RFs as predictive models and use GPs in our experiments unless stated
 145 otherwise. We set the model weight $\beta = 10$ and the model quantile to $\gamma = 0.05$, see our sensitivity

146 studies in Appendices I and J. Before starting the main BO loop in PrBO, we randomly sample $D + 1$
147 points from the prior. We optimize EI using a multi-start local search, similar to SMAC (Hutter et al.,
148 2011). We start with four synthetic benchmarks: Branin, SVM, FC-NET and XGBoost, which are 2,
149 2, 6 and 8 dimensional, respectively. The last three are part of the Profet benchmarks (Klein et al.,
150 2019), generated by a generative model built using performance data on OpenML or UCI datasets.
151 See Appendix C for more details on the experimental setup. Due to space constraints, we defer
152 the (qualitatively similar) results for the SVM benchmark to Appendix D. For the same reason, we
153 defer to Appendix E the study of a real-world application, a programming language and compiler
154 named *Spatial* for the design of application accelerators, i.e., FPGAs. We apply PrBO to three
155 *Spatial* benchmarks, namely, 7D shallow and deep CNNs, and a 10D molecular dynamics grid
156 application. We optimize design runtime constrained to the design fitting a target FPGA. Compared
157 to the previous state-of-the-art, PrBO converges on average $1.49\times$ faster on two benchmarks and
158 achieves $1.28\times$ better final performance on the third. For context, this is a significant improvement
159 in the FPGA field, where a 10% improvement could qualify for acceptance in a top-tier conference.

160 4.1 Prior Selection

161 In this section we study the effect of choosing a prior. A suitable property of the prior is that, by
162 selecting a tighter prior around an optimum, we would expect sampling from the prior to have an
163 increased performance. To the limit, if the prior is composed by only one point which is one of
164 the global optima, then the first sample (and all of them) from the prior will hit the optimum. To
165 have a sanity check of this property, we build an artificial prior in a controlled way. We rely on an
166 automated computation of the prior by computing a univariate Kernel Density Estimation (KDE)
167 using a Gaussian kernel on the synthetic benchmarks introduced above. We note that the goal of
168 these synthetic priors is to have an unbiased prior for our experiments, whereas manual priors would
169 be biased by our own expertise of these benchmarks. In practice, users will manually define these
170 priors without needing additional experiments.

171 We experiment with an array of varying quality priors. We select a constant $10D$ points in each prior
172 and vary the size of the random sample dataset so that we can make the priors more sharply peaked
173 around the optima in a controlled environment. We use the best performing $10D$ samples to create
174 the prior from a uniform random sample dataset size of $10D \frac{100}{x}$; we refer to this prior as $x\%$ in
175 Figure 2. As an example the XGBoost benchmark has $d = 8$, so, 100% means we sample 80 points
176 and use all 80 to create the prior, 10% means we sample 800 points and use the best performing 80
177 to create the prior, 1% means we sample 8,000 and use the best 80 to create the prior, and so on.

178 Figure 2 shows the performance of purely sampling from the prior and running PrBO, respectively,
179 after $10D$ function evaluations with different priors. A bigger random sample dataset and a smaller
180 percentage leads to a tighter prior around the optimum, making the argument for a stronger prior.
181 This is confirmed by Figure 2, where a sharply peaked prior (right side of the figure) leads to a better
182 performance in both scenarios. In addition we observe that in contrast to sampling from the prior,
183 PrBO achieves a smaller regret by being able to evolve from the initial prior and making independent
184 steps towards better values of the objective function. More extensive experiments with a similar trend,
185 including the rest of the benchmarks, are in Appendix H.

186 4.2 Comparison Against Strong Baselines

187 Figure 3 compares PrBO to other optimizers using the log simple regret on five runs (mean and std
188 error reported) on the synthetic benchmarks. We compare the results of PrBO with and without priors
189 (both weak and strong) to $10,000\times$ random search (RS, i.e., for each BO sample we draw 10,000
190 uniform random samples), sampling from the strong prior only, and Spearmint (Snoek et al., 2012)
191 which is a well-adopted BO approach using GPs and the EI acquisition function.

192 PrBO prior beats $10,000\times$ RS and PrBO weak prior on all benchmarks. It also either outperforms or
193 matches the performance of sampling from the prior; this is expected because prior sampling cannot
194 recover from a non-ideal prior. The two methods are identical up to the initialization phase because
195 they both sample from the same prior in that phase.

196 PrBO Prior is more sample efficient and finds better or equal results than Spearmint on three out
197 of the four benchmarks. On XGBoost, PrBO leads the performance until 139 BO iterations, where
198 Spearmint catches up and achieves slightly better results in the end (function values of 8.986 for

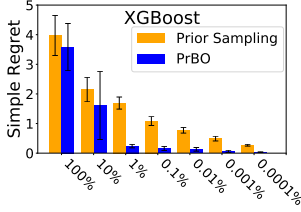


Figure 2: Regret of prior sampling and PrBO with different priors ($\mu \pm \sigma$ on 5 reps.).

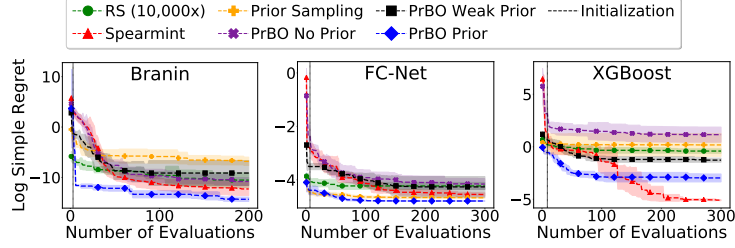


Figure 3: Log regret comparison of PrBO with and without priors, RS, and Spearmint (mean \pm one std on 5 repetitions). We run the benchmarks for 100D iterations, capped at 300.

199 Spearmint vs. 9.026 for PrBO after 300 iterations). Importantly, in all our experiments, PrBO with a
200 good prior consistently shows tremendous speedups in the early phases of the optimization process,
201 typically only requiring on average 8.25 iterations to reach the performance that Spearmint reaches
202 after 100 iterations ($12.12\times$ faster). Thus in comparison to traditional BO approaches, PrBO makes
203 use of the best of both worlds, leveraging prior knowledge and efficient optimization based on BO.

204 5 Related Work

205 TPE by Bergstra et al. (2011) supports limited hand-designed priors in the form of normal or log-
206 normal distributions. We make three technical contributions that make PrBO more flexible than TPE.
207 First, we generalize over the TPE approach by allowing more flexible priors; second, our approach is
208 model-agnostic (i.e., PrBO is not limited to the TPE model; we use GPs and RFs in our experiments);
209 and third, PrBO is inspired by Bayesian models that give more importance to the data as iterations
210 progress. We also show that PrBO outperforms HyperOpt’s TPE in Appendix G.

211 In parallel work, Li et al. (2020) also allow users to specify priors via a probability distribution. Their
212 two-level approach first samples a number of configurations by maximizing Thompson samples from
213 a GP posterior and then chooses the configuration that has the highest prior as the next to evaluate. In
214 contrast, our method leverages the information from the prior more directly and ensures that the prior
215 gets washed out as we collect more data, enabling PrBO to overcome misspecified priors.

216 Similarly, black-box optimization tools, such as SMAC (Hutter et al., 2011) or iRace (López-Ibáñez
217 et al., 2016) also support simple hand-designed priors, e.g. log-transformations. However, these are
218 not properly reflected in the predictive models and both cannot explicitly recover from bad priors.

219 Our work also relates to other meta-learning for BO approaches (Vanschoren, 2019), where BO is
220 applied to many similar optimization problems in a sequence such that knowledge about the general
221 problem structure can be exploited in future optimization problems. In contrast to these approaches,
222 PrBO is the first method that allows human experts to explicitly specify their priors. Furthermore,
223 PrBO does not depend on any meta-features (Feurer et al., 2015b) and only incorporates a single
224 prior instead of many priors from different experiments (Lindauer & Hutter, 2018).

225 6 Conclusions and Future Work

226 We have proposed a novel BO variant, PrBO, that allows users to transfer their expert knowledge
227 into the optimization in the form of priors about which parts of the input space will yield the best
228 performance. These are different than common priors over functions which are much less intuitive
229 for users. BO failed so far to leverage the knowledge of domain experts, not only causing inefficiency
230 but also getting users away from applying BO approaches because they could not exploit their prior
231 knowledge. PrBO addresses this issue and will therefore facilitate the adoption of BO. We showed
232 that PrBO is $12.12\times$ more sample efficient than state-of-the-art methods, and $10,000\times$ faster than
233 random search, on a common suite of benchmarks, and also achieves better final performance in all
234 but one of the benchmarks tested. PrBO also robustly recovers from misleading priors. In future work,
235 we will study how PrBO can be used to leverage other types of prior knowledge from meta-learning,
236 to boost BO’s performance even further.

References

- Maximilian Balandat, Brian Karrer, Daniel R Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. Botorch: Programmable bayesian optimization in pytorch. *arXiv preprint arXiv:1910.06403*, 2019.
- James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pp. 2546–2554, 2011.
- Xavier Bouthillier and Gaël Varoquaux. Survey of machine-learning experimental methods at NeurIPS2019 and ICLR2020. Research report, Inria Saclay Ile de France, January 2020. URL <https://hal.archives-ouvertes.fr/hal-02447823>.
- Yutian Chen, Aja Huang, Ziyu Wang, Ioannis Antonoglou, Julian Schrittwieser, David Silver, and Nando de Freitas. Bayesian optimization in alphago. *CoRR*, abs/1812.06855, 2018.
- Laurence Charles Ward Dixon. The global optimization problem: an introduction. *Toward global optimization*, 2:1–15, 1978.
- Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: robust and efficient hyperparameter optimization at scale. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 1436–1445, 2018.
- Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems* 28, pp. 2962–2970. Curran Associates, Inc., 2015a.
- Matthias Feurer, Jost Tobias Springenberg, and Frank Hutter. Initializing bayesian hyperparameter optimization via meta-learning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pp. 1128–1135, 2015b.
- Jacob R Gardner, Matt J Kusner, Zhixiang Eddie Xu, Kilian Q Weinberger, and John P Cunningham. Bayesian optimization with inequality constraints. In *Proceedings of the 31st International Conference on Machine Learning, ICML*, 2014.
- Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and D. Sculley. Google vizier: A service for black-box optimization. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017.
- GPy. GPy: A gaussian process framework in python. <http://github.com/SheffieldML/GPy>, since 2012.
- F. Hutter, L. Xu, H. Hoos, and K. Leyton-Brown. Algorithm runtime prediction: Methods & evaluation. *Artificial Intelligence*, 206:79–111, 2014.
- Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pp. 507–523. Springer, 2011.
- Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren (eds.). *Automated Machine Learning: Methods, Systems, Challenges*. Springer, 2018. In press, available at <http://automl.org/book>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR*, 2015.
- Aaron Klein, Zhenwen Dai, Frank Hutter, Neil D. Lawrence, and Javier Gonzalez. Meta-surrogate benchmarking for hyperparameter optimization. In *Advances in Neural Information Processing Systems NeurIPS*, pp. 6267–6277, 2019.
- David Koeplinger, Matthew Feldman, Raghu Prabhakar, Yaqi Zhang, Stefan Hadjis, Ruben Fiszal, Tian Zhao, Luigi Nardi, Ardavan Pedram, Christos Kozyrakis, and Kunle Olukotun. Spatial: A Language and Compiler for Application Accelerators. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, June 2018.

284 Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in
285 the presence of noise. *Journal of Basic Engineering*, 86(1):97–106, 1964.

286 Cheng Li, Sunil Gupta, Santu Rana, Vu Nguyen, Antonio Robles-Kelly, and Svetha Venkatesh.
287 Incorporating expert prior knowledge into experimental design via posterior sampling. *arXiv*
288 *preprint arXiv:2002.11256*, 2020.

289 Marius Lindauer and Frank Hutter. Warmstarting of model-based algorithm configuration. In
290 *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pp. 1355–1362,
291 2018.

292 Manuel López-Ibáñez, Jérémie Dubois-Lacoste, Leslie Pérez Cáceres, Thomas Stützle, and Mauro
293 Birattari. The irace package: Iterated racing for automatic algorithm configuration. *Operations*
294 *Research Perspectives*, 3:43–58, 2016.

295 Luigi Nardi, David Koeplinger, and Kunle Olukotun. Practical design space exploration. In *27th*
296 *IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecom-*
297 *munication Systems, MASCOTS*, 2019.

298 Arshak Navruzian, Frank Sharp, Jeremy Howard, and Antoine Saliou. Optimizing hy-
299 perparams for image datasets in fastai. [https://platform.ai/blog/page/1/](https://platform.ai/blog/page/1/optimizing-hyperparams-for-image-datasets-in-fastai/)
300 [optimizing-hyperparams-for-image-datasets-in-fastai/](https://platform.ai/blog/page/1/optimizing-hyperparams-for-image-datasets-in-fastai/), 2019.

301 Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business
302 Media, 2012.

303 Andrei Paleyes, Mark Pullin, Maren Mahsereci, Neil Lawrence, and Javier González. Emulation
304 of physical processes with emukit. In *Second Workshop on Machine Learning and the Physical*
305 *Sciences, NeurIPS*, 2019.

306 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
307 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
308 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
309 12:2825–2830, 2011.

310 Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine
311 learning algorithms. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp.
312 2960–2968, 2012.

313 Joaquin Vanschoren. Meta-learning. In *Automated Machine Learning - Methods, Systems, Challenges*,
314 pp. 35–61. 2019.