

# Transfer Learning for Bayesian Optimization with Principal Component Analysis

Masui, Hideyuki; Romeres, Diego; Nikovski, Daniel N.

TR2022-169 December 20, 2022

## Abstract

Bayesian Optimization has been widely used for black-box optimization. Especially in the field of machine learning, BO has obtained remarkable results in hyperparameters optimization. However, the best hyperparameters depend on the specific task and traditionally the BO algorithm needs to be repeated for each task. On the other hand, the relationship between hyperparameters and objectives has similar tendency among tasks. Therefore, transfer learning is an important technology to accelerate the optimization of novel task by leveraging the knowledge acquired in prior tasks. In this work, we propose a new transfer learning strategy for BO. We use information geometry based principal component analysis (PCA) to extract a low-dimension manifold from a set of Gaussian process (GP) posteriors that models the objective functions of the prior tasks. Then, the low dimensional parameters of this manifold can be optimized to adapt to a new task and set a prior distribution for the objective function of the novel task. Experiments on hyperparameters optimization benchmarks show that our proposed algorithm, called BO-PCA, accelerates the learning of an unseen task (less data are required) while having low computational cost.

*International Conference on Machine Learning and Applications (ICMLA) 2022*

© 2022 MERL. This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Research Laboratories, Inc.; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Research Laboratories, Inc. All rights reserved.



# Transfer Learning for Bayesian Optimization with Principal Component Analysis

1<sup>st</sup> Hideyuki Masui

*Mitsubishi Electric Corporation*

Kamakura, Kanagawa, Japan

Masui.Hideyuki@bc.MitsubishiElectric.co.jp

2<sup>nd</sup> Diego Romeres

*Mitsubishi Electric Research Labs*

Cambridge, MA, USA

romeres@merl.com

3<sup>rd</sup> Daniel Nikovski

*Mitsubishi Electric Research Labs*

Cambridge, MA, USA

nikovski@merl.com

**Abstract**—Bayesian Optimization has been widely used for black-box optimization. Especially in the field of machine learning, BO has obtained remarkable results in hyperparameters optimization. However, the best hyperparameters depend on the specific task and traditionally the BO algorithm needs to be repeated for each task. On the other hand, the relationship between hyperparameters and objectives has similar tendency among tasks. Therefore, transfer learning is an important technology to accelerate the optimization of novel task by leveraging the knowledge acquired in prior tasks. In this work, we propose a new transfer learning strategy for BO. We use information geometry based principal component analysis (PCA) to extract a low-dimension manifold from a set of Gaussian process (GP) posteriors that models the objective functions of the prior tasks. Then, the low dimensional parameters of this manifold can be optimized to adapt to a new task and set a prior distribution for the objective function of the novel task. Experiments on hyperparameters optimization benchmarks show that our proposed algorithm, called BO-PCA, accelerates the learning of an unseen task (less data are required) while having low computational cost.

**Index Terms**—Bayesian optimization, Gaussian process, transfer learning, principal component analysis

## I. INTRODUCTION

Bayesian Optimization (BO) is a well-established framework for solving black-box optimization problems of the form:

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad (1)$$

where  $f(\cdot)$  is an expensive-to-evaluate objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$  and  $\mathcal{X} \subset \mathbb{R}^D$  denotes the configuration space. BO has been applied with remarkable results in a wide range of areas, such as hyperparameter optimization [4], [20], robotics [11], [16], materials design [29], medical science [22] and many others. BO searches for the global optimum  $\mathbf{x}^*$  by iteratively selecting input configurations and updating a surrogate model of the objective function  $f$ . In the case of hyperparameter optimization (HPO) task,  $\mathbf{x}$  is a vector of hyperparameters and  $f$  represents a validation error.

In this work, we focus on transfer learning for black-box optimization [2]. We assume that we have already completed  $T$  somewhat similar black-box optimization tasks  $\{f_t\}_{t=1}^T$  defined over the same configuration space  $\mathcal{X}$ . If a new similar black-box optimization task is given, it could be beneficial to leverage the prior tasks as ancillary information when solving

the new one. Hence, our goal is to leverage prior source tasks' data sets to accelerate the optimization of the test task  $f_{T+1}$ .

Our algorithm accomplishes the knowledge transfer by extracting a low-dimensional manifold and parameter vector from the prior tasks' data, based on information geometry, and updating the parameter vector to solve the new task. This approach has several advantages. First, by using information geometry, a low-dimensional parameter vector is learned for improved stability even though a limited number of source tasks and observations. Second, the optimization problem of these few parameters in each iteration of BO can be converted to one of linear least squares. These advantages are important, because in real situations, it is commonly difficult to collect a lot of data from source tasks to train complex models. Furthermore, we often have to use limited computer resources.

**Contributions.** Our main contributions are as follow:

- 1) We introduce a new transfer learning method for BO, called BO-PCA, that uses information geometry-based feature extraction for Gaussian process learning of the objective functions.
- 2) The proposed algorithm is able to update online the transfer learning model in each iteration of the BO algorithm.
- 3) We show the effectiveness of our algorithm over s.o.t.a. approaches, using standard HPO benchmark tasks and data sets.

## II. BACKGROUND

### A. Bayesian Optimization

BO is a class of machine learning-based algorithms that aim to solve the problem defined in (1), when standard optimization algorithms typically fail due to the complexity of the objective function  $f(\mathbf{x})$ . This class of algorithms is typically composed of iterative procedures that consist of two main computational steps with the objective of exploring the configuration space  $\mathcal{X}$  and finding an optimal configuration  $\mathbf{x}^*$ .

The first step consists of constructing a surrogate function of the objective function  $f(\cdot)$  using a set of collected observations  $\mathcal{D} = (\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ . We assume that the vector  $\mathbf{y}$  consists of noisy observations of  $f(\cdot)$ , i.e.,

$$y_i = f(\mathbf{x}_i) + \epsilon, \quad (2)$$

where  $\epsilon$  is an independent zero-mean Gaussian noise with variance  $\beta^{-1}$ . Typically, the surrogate function is estimated using Gaussian process regression (GPR) [25]. In this case, the objective function is modeled a priori as a Gaussian Process (GP),  $f(\cdot) \sim \mathcal{N}(m_0(\cdot), k(\cdot, \cdot))$ , where  $m_0(\cdot)$  is an arbitrary mean function and  $k(\cdot, \cdot)$  is the a priori covariance matrix, defined by means of a kernel. Typical choices of the kernel function are the RBF and the Matérn 5/2 kernel function [25]:

$$k_{5/2}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\sqrt{5}r}{\lambda} + \frac{5r^2}{3\lambda^2}\right) \exp\left(-\frac{\sqrt{5}r}{\lambda}\right), \quad (3)$$

where  $r = |\mathbf{x} - \mathbf{x}'|$  and  $\lambda$  is a length scale parameter called a hyperparameter. Every kernel function is parameterized by some hyperparameters and is usually optimized by maximization of the marginal log-likelihood:

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\lambda}) \propto -\log |\mathbf{K}\boldsymbol{\lambda}| - \mathbf{y}^T \mathbf{K}\boldsymbol{\lambda}^{-1} \mathbf{y}. \quad (4)$$

Remarkably, the posterior distribution of  $f$  for an unseen point  $\mathbf{x}$ , namely  $\hat{f}(\mathbf{x}|\mathcal{D})$ , can be computed analytically, and is normally distributed with mean  $\mu(\mathbf{x})$  and variance  $\sigma^2(\mathbf{x})$  which are defined below:

$$\mu(\mathbf{x}) = m_0(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \beta^{-1}\mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}_0), \quad (5)$$

$$\sigma^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \beta^{-1}\mathbf{I})^{-1} \mathbf{k}(\mathbf{x}), \quad (6)$$

where  $\mathbf{K} = k(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N}$  is defined element-wise through the kernel function  $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ , the a priori covariance matrix, and  $\mathbf{k}(\mathbf{x}) = k(\mathbf{x}, \mathbf{X}) \in \mathbb{R}^{N \times 1}$  is a vector with entries  $k_i = k(\mathbf{x}, \mathbf{x}_i)$  and  $\mathbf{m}_0 := m_0(\mathbf{X}) \in \mathbb{R}^{N \times 1}$  is the mean vector, where each entry is defined as  $m_{0,i} = m_0(\mathbf{x}_i)$  (typically  $m_0(\mathbf{x}) \equiv 0$ ). For later convenience, the posterior distribution  $\hat{f}(\mathbf{x}|\mathcal{D})$  can also be defined based on the GP posterior of the collected observations,  $f(\mathbf{X})$ , which is a multivariate normal distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , with

$$\boldsymbol{\mu} = \mathbf{m}_0 + \mathbf{K} (\mathbf{K} + \beta^{-1}\mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}_0), \quad (7)$$

$$\boldsymbol{\Sigma} = \mathbf{K} - \mathbf{K} (\mathbf{K} + \beta^{-1}\mathbf{I})^{-1} \mathbf{K}, \quad (8)$$

After standard computation, an equivalent definition of (5) and (6) based on  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  is:

$$\mu(\mathbf{x}) = m_0(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} (\boldsymbol{\mu} - \mathbf{m}_0), \quad (9)$$

$$\sigma^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} (\boldsymbol{\Sigma} - \mathbf{K}) \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}). \quad (10)$$

The second step consists in maximizing an acquisition function, which is defined according to an exploration-exploitation criterion. One of the most commonly used function is the expected improvement (EI) [14]:

$$\alpha(\mathbf{x}|\mathcal{D}) = \mathbb{E}_{\hat{f}(\mathbf{x}|\mathcal{D})} \left[ \max(0, \hat{f}(\mathbf{x}^+) - \hat{f}(\mathbf{x})) \right] \quad (11)$$

$$= \sigma(\mathbf{x}) z \Phi(z) + \sigma(\mathbf{x}) \phi(z), \quad (12)$$

where  $\mathbf{x}^+ = \arg \min[\hat{f}(\mathbf{x}_1), \hat{f}(\mathbf{x}_2), \dots, \hat{f}(\mathbf{x}_N)]$  is the current best configuration,  $z = \frac{\hat{f}(\mathbf{x}^+) - \mu(\mathbf{x})}{\sigma(\mathbf{x})}$  is a standard normal distribution, and  $\Phi(\cdot)$  and  $\phi(\cdot)$  denote the CDF and PDF of  $\hat{f}(\cdot)$ , respectively. W.l.o.g. in this paper, we use the EI function

as the acquisition function, but others could be used, e.g. the probability of improvement, the knowledge gradient, etc. The interested reader can refer for example to [6], [19], [21] for a more in depth description of BO.

## B. GP-ePCA and GP-mPCA

Principal component analysis (PCA) is a widely used algorithm to extract a low-dimensional linear subspace from high dimensional data that explains that data well. The case in which each of these data points represents a GP distribution was studied in [1], and used in [10] to propose a meta-learning algorithm for GP based on PCA.

Let us assume that there are  $T$  data sets  $\{\mathcal{D}_t\}_{t=1}^T = \{(\mathbf{X}_t, \mathbf{y}_t)\}_{t=1}^T$  with  $\mathbf{X}_t \in \mathbb{R}^{N_t \times D}$  and  $\mathbf{y}_t \in \mathbb{R}^{N_t \times 1}$ , where  $N_t$  is the number of data points in data set  $t$ . With a slight abuse of notation, some of the quantities defined in the previous section for one data set will be redefined for multiple data sets. Let now  $\mathbf{X} = \bigcup_{t=1}^T \mathbf{X}_t$  be a union set of the input sets. Following Section II-A, the estimated GP posterior of  $\mathbf{X}$  given the  $t$ -th data set is a multivariate normal distribution  $\hat{f}_t(\mathbf{X}|\mathcal{D}_t) \sim (\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ , with

$$\boldsymbol{\mu}_t = \mathbf{m}_0(\mathbf{X}) + \mathbf{K}_t (\mathbf{K}_{tt} + \beta^{-1}\mathbf{I})^{-1} (\mathbf{y}_t - \mathbf{m}_0(\mathbf{X}_t)), \quad (13)$$

$$\boldsymbol{\Sigma}_t = \mathbf{K} - \mathbf{K}_t (\mathbf{K}_{tt} + \beta^{-1}\mathbf{I})^{-1} \mathbf{K}_t^T, \quad (14)$$

where  $\mathbf{K} = k(\mathbf{X}, \mathbf{X})$ ,  $\mathbf{K}_t = k(\mathbf{X}, \mathbf{X}_t)$ ,  $\mathbf{K}_{tt} = k(\mathbf{X}_t, \mathbf{X}_t)$ .

Since the multivariate normal distribution belongs to the exponential family, each GP posterior  $\hat{f}_t(\mathbf{X}|\mathcal{D}_t)$  is a point on a Riemannian flat manifold, here denoted by  $\mathcal{S}$ . A flat manifold is spanned by a coordinate system. In particular, there are two coordinate systems, the e-coordinates and m-coordinates, that define two subspaces linear w.r.t.  $\boldsymbol{\xi}_t = (\boldsymbol{\theta}_t^T, \text{vec}(\boldsymbol{\Theta}_t)^T)^T$  and  $\boldsymbol{\zeta}_t = (\boldsymbol{\eta}_t^T, \text{vec}(\mathbf{H}_t)^T)^T$ , respectively, where:

$$\boldsymbol{\theta}_t = \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\mu}_t, \quad \boldsymbol{\Theta}_t = -\frac{1}{2} \boldsymbol{\Sigma}_t^{-1}, \quad (15)$$

$$\boldsymbol{\eta}_t = \boldsymbol{\mu}_t, \quad \mathbf{H}_t = \boldsymbol{\mu}_t \boldsymbol{\mu}_t^T + \boldsymbol{\Sigma}_t^T. \quad (16)$$

The goal of the GP-ePCA and GP-mPCA methods is to project  $\{\boldsymbol{\xi}_t\}_{t=1}^T$  and  $\{\boldsymbol{\zeta}_t\}_{t=1}^T$ , respectively, into a low dimension linear submanifold  $\mathcal{M} \subset \mathcal{S}$  which is expressed by the set of points:

$$\tilde{\boldsymbol{\xi}}_t = (\tilde{\boldsymbol{\theta}}_t^T, \text{vec}(\tilde{\boldsymbol{\Theta}}_t)^T)^T = \mathbf{U}_\xi \mathbf{w}_t + \mathbf{u}_{\xi,0}, \quad (17)$$

$$\tilde{\boldsymbol{\zeta}}_t = (\tilde{\boldsymbol{\eta}}_t^T, \text{vec}(\tilde{\mathbf{H}}_t)^T)^T = \mathbf{U}_\zeta \mathbf{w}_t + \mathbf{u}_{\zeta,0}. \quad (18)$$

where  $\mathbf{U}_\xi, \mathbf{U}_\zeta \in \mathbb{R}^{(N+N^2) \times L}$  with  $L < T$  and  $\mathbf{w}_t \in \mathbb{R}^{L \times 1}$ . When  $\{\boldsymbol{\xi}_t\}_{t=1}^T$  or  $\{\boldsymbol{\zeta}_t\}_{t=1}^T$  are observed, which is at training time, the task of GP-ePCA and GP-mPCA is to optimize the weights  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T)^T$ , the basis vectors  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_L)$ , and the offset  $\mathbf{u}_0$  by minimizing the sum of the KL divergence:

$$\hat{\mathbf{W}}, \hat{\mathbf{U}}, \hat{\mathbf{u}}_0 = \arg \min_{\mathbf{W}, \mathbf{U}, \mathbf{u}_0} E(\mathbf{W}, \mathbf{U}, \mathbf{u}_0), \quad (19)$$

$$E(\mathbf{W}, \mathbf{U}, \mathbf{u}_0) = \begin{cases} \sum_{t=1}^T D_{KL} \left( \boldsymbol{\xi}_t \parallel \tilde{\boldsymbol{\xi}}_t \right) & \text{(GP-ePCA)} \\ \sum_{t=1}^T D_{KL} \left( \boldsymbol{\zeta}_t \parallel \tilde{\boldsymbol{\zeta}}_t \right) & \text{(GP-mPCA)} \end{cases}.$$

As a result, GP-ePCA, respectively GP-mPCA, obtain a low-dimensional parameter vector  $\mathbf{w}_t$  which characterizes each training data set. Therefore, at prediction time, given the weights  $\mathbf{w}$ , the posterior distribution  $\mathcal{N}(\tilde{\boldsymbol{\mu}}(\mathbf{w}), \tilde{\boldsymbol{\Sigma}}(\mathbf{w}))$  can be restored by inverting (15) (and (16)). The posterior estimator of a test point  $\mathbf{x}$  can be computed with mean and variance:

$$\begin{aligned}\tilde{m}(\mathbf{x}, \mathbf{w}) &= m_0(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}(\tilde{\boldsymbol{\mu}}(\mathbf{w}) - \mathbf{m}_0), \\ \tilde{\sigma}^2(\mathbf{x}, \mathbf{w}) &= k(\mathbf{x}, \mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}(\tilde{\boldsymbol{\Sigma}}(\mathbf{w}) - \mathbf{K})\mathbf{K}^{-1}\mathbf{k}(\mathbf{x}).\end{aligned}\quad (20)$$

where  $\tilde{\boldsymbol{\mu}}(\mathbf{w})$  and  $\tilde{\boldsymbol{\Sigma}}(\mathbf{w})$  are computed as below:

$$\tilde{\boldsymbol{\mu}}(\mathbf{w}) = \begin{cases} -0.5\tilde{\boldsymbol{\Theta}}^{-1}\tilde{\boldsymbol{\theta}} & \text{(GP-ePCA)} \\ \tilde{\boldsymbol{\eta}} & \text{(GP-mPCA)} \end{cases}, \quad (22)$$

$$\tilde{\boldsymbol{\Sigma}}(\mathbf{w}) = \begin{cases} -0.5\tilde{\boldsymbol{\Theta}}^{-1} & \text{(GP-ePCA)} \\ \tilde{\mathbf{H}} - \tilde{\boldsymbol{\eta}}\tilde{\boldsymbol{\eta}}^T & \text{(GP-mPCA)} \end{cases}. \quad (23)$$

### C. Sparse GP-ePCA and GP-mPCA

The major problem of GP-PCA is the calculation cost. If the size of  $\mathbf{X} = \bigcup_{t=1}^T \mathbf{X}_t$  is  $|\mathbf{X}| = N$ , the computation cost of GP-PCA becomes  $\mathcal{O}(N^3)$ . The authors of GP-PCA [10] also proposed sparse GP-PCA to reduce calculation cost which is based on sparse Gaussian process [12]. Let  $\mathbf{Z}$  be a set of inducing points  $\mathbf{Z} = \bigcup_{t=1}^T \mathbf{Z}_t$ ,  $\mathbf{K}_z = k(\mathbf{X}, \mathbf{X}_z)$  and  $\mathbf{K}_{zz} = k(\mathbf{Z}, \mathbf{Z})$ . The covariance matrix can be approximated as below:

$$\mathbf{K} \approx \mathbf{K}_z \mathbf{K}_{zz}^{-1} \mathbf{K}_z^T. \quad (24)$$

In the case of sparse GP-PCA, the posterior distribution of each task corresponds to (13) and (14) is defined as below:

$$\boldsymbol{\mu}_t = \mathbf{m}_0(\mathbf{X}_z) + \mathbf{K}_{zz} \mathbf{A}_{zz}^{-1} \mathbf{K}_{zt} (\mathbf{y}_t - \mathbf{m}_0(\mathbf{X}_t)), \quad (25)$$

$$\boldsymbol{\Sigma}_t = \beta^{-1} \mathbf{K}_{zz} \mathbf{A}_{zz}^{-1} \mathbf{K}_{zz}, \quad (26)$$

where  $\mathbf{A}_{zz} = \beta^{-1} \mathbf{K}_{zz} + \mathbf{K}_z^T \mathbf{K}_z$  and  $\mathbf{K}_{zt} = k(\mathbf{X}_z, \mathbf{X}_t)$ . Let the number of inducing points as  $|\mathbf{Z}_t| = M$  and share this for all tasks,  $\mathbf{Z}_1 = \mathbf{Z}_2 = \dots = \mathbf{Z}_T$ . In this case, the calculation cost of GP-PCA is reduced from  $\mathcal{O}(N^3)$  to  $\mathcal{O}(M^3)$ .

## III. PROBLEM FORMULATION

Consider  $T$  source tasks available at training time, each of this task is defined by an objective function  $\{f_t\}_{t=1}^T$  and for each task  $N_t$  labeled data are collected  $\{\mathcal{D}_t\}_{t=1}^T = \{ \{(\mathbf{x}_{t,n}, y_{t,n})\}_{n=1}^{N_t} \}_{t=1}^T$ . Given the data sets  $\{\mathcal{D}_t\}_{t=1}^T$ , it is possible to estimate  $T$  surrogate functions, using GPR or any other machine learning function approximator. We are interested in a data-efficient transfer learning BO algorithm capable of transferring the knowledge acquired in the  $T$  source tasks to solve a novel test task ( $t = T + 1$ ), i.e.,

$$\mathbf{x}_{T+1}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f_{T+1}(\mathbf{x} | \{\mathcal{D}_t\}_{t=1}^T). \quad (27)$$

The question we want to answer is: ‘‘How can we learn efficiently, in terms of data complexity, a new task, having previously learned  $T$  similar source tasks?’’

## A. Related Work

Standard BO without transfer learning has been already discussed in the previous sections, see [19] for a survey. The use of PCA in BO was proposed in [17] to find a low dimensional space for high dimensional data point where the surrogate of the objective function is learned. This is very different from our approach since PCA was used in the observations and not in the GP posteriors like in GP-PCA. Moreover, in this work we focus on transfer learning of BO from source tasks to a novel unseen task. This problem attracted a lot of research and is strongly connected to the concept of meta-learning [9], [24]. In [2], the authors propose a series of benchmark tasks, data sets as well as metrics to facilitate the comparison among different methods. Moreover, the authors of [5] not only propose a method for transfer learning BO called RGPE, but they also make available the data sets and the implementation of s.o.t.a. methods. We will refer to this implementation and benchmarks to compare our algorithm. The main algorithms for BO transfer learning are summarized in the following. ABLR (adaptive Bayesian linear regression) [15], is a multi-task Bayesian linear regression approach. Transfer learning is achieved by sharing the deep neural network which is trained with source tasks data. Then, separate Bayesian linear regressors are embedded for each task. TST-R [27] (two-stage transfer surrogate using pairwise hyperparameter performance rankings) is an ensemble approach. At the first stage, a surrogate model is trained for each source and test task. Then, at the second stage, the weights are calculated for each model using similarity between the test task and each source task. TAF-R (transfer acquisition function framework using pairwise hyperparameter performance rankings) [28], is a similar ensemble approach to TST-R. The difference between TST-R is that TAF-R forms ensembles of acquisition functions in order to overcome different scales of evaluation scores. Finally, RGPE (ranking-weighted Gaussian process ensemble) [5], is another ensemble surrogate models approach similar to TST-R. The weights are calculated by optimizing a ranking loss using surrogate models and test task data.

## IV. BO-PCA

We propose to solve the transfer learning BO problem (27) by learning the surrogate function with a GP-PCA method (both GP-ePCA and GP-mPCA are possible). We call this approach BO-PCA. GP-PCA is a principled way of encoding the GP posteriors in a low-dimensional manifold, and at training time it defines the space where the posterior of the novel  $T + 1$  task can be (data) efficiently learned by creating transfer learning model that depends on a low dimensional set of parameters easier to update online. The proposed method is described in Algorithm 1.

During training time, we learn the  $\{\hat{f}_t(\mathbf{X} | \mathcal{D}_t)\}_{t=1}^T$  surrogate functions of the source tasks using GP-ePCA (or GP-mPCA), to obtain  $L$  basis  $\hat{\mathbf{U}}$  and the offset  $\hat{\mathbf{u}}_0$ , as described in Section II-B. Then, in order to optimize an unseen test task,  $f_{T+1}$ , we fix the basis functions  $\{\hat{\mathbf{U}}, \hat{\mathbf{u}}_0\}$ , collect a small number of data  $n_0$  from  $f_{T+1}$ , and each time a new  $n$ -th

---

**Algorithm 1** BO with GP-PCA Prior

---

**Input:** source tasks data  $\{\mathcal{D}_t\}_{t=1}^T$ , number of initial points  $n_0$ , budget  $N_{t+1}$ , GP-PCA manifold dimension  $L$ .  
Execute GP-ePCA to obtain  $\hat{\mathbf{W}}_\xi, \hat{\mathbf{U}}_\xi, \hat{\mathbf{u}}_{\xi,0}$  (19)  
(Execute GP-mPCA to obtain  $\hat{\mathbf{W}}_\zeta, \hat{\mathbf{U}}_\zeta, \hat{\mathbf{u}}_{\zeta,0}$  (19))  
Observe  $f_{T+1}$  at  $n_0$  points  $\mathcal{D}_{T+1} = \{\mathbf{x}_{T+1,n}, y_{T+1,n}\}_{n=1}^{n_0}$   
Set  $n = n_0$   
**while**  $n < N_{T+1}$  **do**  
  Fit weights  $\hat{\mathbf{w}}$  of GP-ePCA (GP-mPCA) (28)  
  Update weight  $\mathbf{w}_{T+1} = \hat{\mathbf{w}}$   
  Set prior mean  $m_0(\mathbf{x}) = \tilde{m}(\mathbf{x}, \mathbf{w}_{T+1})$  for  $f_{T+1}$   
  Learn  $\hat{f}_t(\mathbf{X}|\mathcal{D}_{T+1})$  using GPR  
  Optimize  $\mathbf{x}_{T+1,n} = \arg \max_{\mathbf{x} \in \mathcal{X}} \alpha(\mathbf{x}|\mathcal{D}_{T+1})$   
  Observe  $y_{T+1,n} = f_{T+1}(\mathbf{x}_{T+1,n})$   
  Update  $\mathcal{D}_{T+1} = \mathcal{D}_{T+1} \cup \{(\mathbf{x}_{T+1,n}, y_{T+1,n})\}$   
   $n = n + 1$   
**end while**  
**Output:**  $\hat{\mathbf{x}}_{T+1} = \arg \min_{n=1,2,\dots,N_{T+1}} f_{T+1}(\mathbf{x}_{T+1,n})$

---

data point  $(\mathbf{x}_{T+1,n}, y_{T+1,n})$  arrives, the weight vector,  $\mathbf{w}_{T+1}$  is updated by minimizing the least-squares problem:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^L} \|\mathbf{y}_{T+1} - \tilde{\mathbf{m}}(\mathbf{X}_{t+1}, \mathbf{w})\|_2^2, \quad (28)$$

where  $\tilde{\mathbf{m}}(\mathbf{X}_{t+1}, \mathbf{w})$  is a  $n_{T+1} \times 1$  vector where each entry  $\tilde{m}_i(\mathbf{X}_{t+1}, \mathbf{w}) = \tilde{m}(\mathbf{x}_{T+1,i}, \mathbf{w})$ . We use this GP-PCA estimated function  $\tilde{m}(\mathbf{x}, \hat{\mathbf{w}})$  as a prior mean function of  $f_{T+1}$ , that is  $m_0(\mathbf{x}) = \tilde{m}(\mathbf{x}, \mathbf{w}_{T+1})$ . We can therefore use GPR to learn the surrogate objective function of  $\hat{f}_t(\mathbf{X}|\mathcal{D}_{T+1})$ . The subsequent steps of each iteration follow in accordance with standard BO algorithms.

*Remark 1:* Note that after computing the optimal vector of weights  $\hat{\mathbf{w}}$ , we could compute both the mean and covariance GP-PCA functions, see equations (22) and (23), respectively. However, the source tasks commonly have many more data points than the test task, and for this reason, the GP-PCA covariance function tends to significantly affect the GP posterior variance (10), which leads to two practical limitations if used in Algorithm 1. First, given the many data points in the source tasks, the estimated variance tends to be narrow, preventing the BO algorithm to explore sufficiently. Second, if some test task configurations are far away from all source task configurations, the posterior variance tends to be erroneously very large, only because that point was not seen during training. For these reasons, only the GP-PCA mean function is used to transfer the knowledge of the source tasks to the novel task.

#### A. GP-ePCA vs GP-mPCA

Next, we discuss the main differences between GP-ePCA and GP-mPCA that are relevant to implementing our transfer learning BO optimization algorithm. In Algorithm 1, the weights are optimized by solving problem (28) at each iteration (i.e., every time a new data pair becomes available). In the

case of GP-ePCA, the estimation vector  $\tilde{\boldsymbol{\mu}}(\mathbf{w})$  is calculated as shown below, following (22),

$$\tilde{\boldsymbol{\mu}}_{\text{ePCA}}(\mathbf{w}) = -0.5 \left( \hat{\mathbf{U}}_{\Theta} \mathbf{w} + \hat{\mathbf{u}}_{0,\Theta} \right)^{-1} \left( \hat{\mathbf{U}}_{\theta} \mathbf{w} + \hat{\mathbf{u}}_{0,\theta} \right), \quad (29)$$

where,

$$\begin{bmatrix} \tilde{\theta}(\mathbf{w}) \\ \text{vec} \left( \tilde{\Theta}(\mathbf{w}) \right) \end{bmatrix} = \hat{\mathbf{U}}_{\xi} \mathbf{w} + \hat{\mathbf{u}}_{\xi,0} = \begin{bmatrix} \hat{\mathbf{U}}_{\theta} \mathbf{w} + \hat{\mathbf{u}}_{0,\theta} \\ \hat{\mathbf{U}}_{\Theta} \mathbf{w} + \hat{\mathbf{u}}_{0,\Theta} \end{bmatrix}. \quad (30)$$

Analogously, in the case of GP-mPCA:

$$\tilde{\boldsymbol{\mu}}_{\text{mPCA}}(\mathbf{w}) = \hat{\mathbf{U}}_{\eta} \mathbf{w} + \hat{\mathbf{u}}_{0,\eta}, \quad (31)$$

where

$$\begin{bmatrix} \tilde{\eta}(\mathbf{w}) \\ \text{vec} \left( \tilde{\mathbf{H}}(\mathbf{w}) \right) \end{bmatrix} = \hat{\mathbf{U}}_{\zeta} \mathbf{w} + \hat{\mathbf{u}}_0 = \begin{bmatrix} \hat{\mathbf{U}}_{\eta} \mathbf{w} + \hat{\mathbf{u}}_{0,\eta} \\ \hat{\mathbf{U}}_{\mathbf{H}} \mathbf{w} + \hat{\mathbf{u}}_{0,\mathbf{H}} \end{bmatrix}. \quad (32)$$

It follows that  $\tilde{\boldsymbol{\mu}}_{\text{mPCA}}(\mathbf{w})$  is a linear function of the weights  $\mathbf{w}$ , while  $\tilde{\boldsymbol{\mu}}_{\text{ePCA}}(\mathbf{w})$  is a non-linear function of the same weights. For this reason, the optimization problem (28) for GP-ePCA is much harder to solve and in practice standard gradient-based methods were failing to obtain optimal weights in many experiments. To make GP-ePCA working we had to apply derivative-free optimization algorithms such as CMA-ES [8], which are in general slower to compute. Therefore, GP-mPCA is superior to GP-ePCA from the point of view of calculation cost.

#### B. Online GP-mPCA

Following the discussion in the previous section the weight optimization problem (28) for GP-mPCA is linear and is obtained as below by substituting (31) into (20):

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^L} \|\mathbf{y}_{T+1} - k(\mathbf{X}_{t+1}, \mathbf{X}) \mathbf{K}^{-1} \tilde{\boldsymbol{\mu}}(\mathbf{w})\|_2^2, \quad (33)$$

$$= \arg \min_{\mathbf{w} \in \mathbb{R}^L} \|\mathbf{y}'_{T+1} - \phi(\mathbf{X}_{t+1}) \mathbf{w}\|_2^2, \quad (34)$$

where

$$\mathbf{y}'_{T+1} = \mathbf{y}_{T+1} - k(\mathbf{X}_{t+1}, \mathbf{X}) \mathbf{K}^{-1} \hat{\mathbf{u}}_{0,\eta}, \quad (35)$$

$$\phi(\mathbf{X}_{T+1}) = k(\mathbf{X}_{t+1}, \mathbf{X}) \mathbf{K}^{-1} \hat{\mathbf{U}}_{\eta}. \quad (36)$$

Hence, (34) is a least squares problem and can be solved recursively by means of the recursive least squares algorithm [13, Chp. 11]. Therefore, the weights at the  $n$ -th iteration can be updated exactly based on the weights at the  $n-1$ -th iteration, without having to solve the whole optimization problem again, and the update does not depend on the number of data, leading to a fast weight update step.

#### C. Initialization of GP-mPCA

The GP-PCA algorithm uses gradient descent in (19) to optimize  $\mathbf{U}$ ,  $\mathbf{u}_0$  and  $\mathbf{W}$ . That is, it is not guaranteed to converge to the global optimum. Therefore, initialization of  $\mathbf{U}$ ,  $\mathbf{u}_0$  and  $\mathbf{W}$  is very important. The authors of m-PCA (e-PCA) [1] suggest to use m-center (e-center) for  $\mathbf{u}_0$  and to compute

TABLE I  
SEARCH SPACE AND HYPERPARAMETERS DEFINITION FOR EACH BENCHMARK.

Quadratic			AdaBoost			SVM			NN		
Name	Range	Log	Name	Range	Log	Name	Range	Log	Name	Range	Log
$x_0$	$[-5, 5]$	No	# Iterations	$[2, 10^4]$	Yes	Use linear kernel	{No, Yes}	No	Batch size	$[16, 512]$	Yes
$x_1$	$[-5, 5]$	No	# Terms	$[2, 30]$	Yes	Use polynomial kernel	{No, Yes}	No	Learning rate	$[10^{-4}, 10^{-1}]$	Yes
$x_2$	$[-5, 5]$	No				Use RBF kernel	{No, Yes}	No	Momentum	$[0.1, 0.99]$	No
						Trade-off parameter	$[2^{-5}, 2^6]$	No	Weight decay	$[10^{-5}, 10^{-1}]$	Yes
						Degree of the kernel	$[2, 10]$	No	# Layers	$[1, 5]$	No
						Width of the kernel	$[10^{-4}, 10^3]$	Yes	# Units per layer	$[64, 1024]$	Yes
									Dropout	$[0.0, 1.0]$	No

standard PCA to initialize  $\mathbf{U}$  and  $\mathbf{W}$ . The definition of m-center (e-center) is given by the arithmetic mean of samples in m-coordinate (e-coordinate),

$$\hat{\mathbf{u}}_0 = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\eta}_t, \quad \left( \hat{\mathbf{u}}_0 = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\theta}_t \right). \quad (37)$$

In order to leverage the linearity of GP-mPCA, we modify its initialization. In light of the fact that our algorithm only requires the mean  $\tilde{\boldsymbol{\mu}}(\mathbf{w})$  of the restored GP distribution and  $\tilde{\boldsymbol{\mu}}(\mathbf{w})$  is obtained only with  $\hat{\mathbf{U}}_\eta$  and  $\hat{\mathbf{u}}_{0,\eta}$  (and not also  $\hat{\mathbf{U}}_{\mathbf{H}}$  and  $\hat{\mathbf{u}}_{0,\mathbf{H}}$ ) in (31), it is not necessary to extract low-dimension features from  $\{\mathbf{H}_t\}_{t=1}^T$ , defined in (16).

We propose to expand the basis functions of  $\mathbf{U}$  from  $L$  to  $L + T$  which corresponds to also increase the dimensionality of the weights  $\mathbf{W}$  to  $L + T$

$$\mathbf{U}\mathbf{W} + \mathbf{u}_0 = \begin{bmatrix} \mathbf{U}_{\eta,\eta} & \mathbf{U}_{\eta,\mathbf{H}} \\ \mathbf{U}_{\mathbf{H},\eta} & \mathbf{U}_{\mathbf{H},\mathbf{H}} \end{bmatrix} \begin{bmatrix} \mathbf{W}_\eta \\ \mathbf{W}_{\mathbf{H}} \end{bmatrix} + \begin{bmatrix} \mathbf{u}_{0,\eta} \\ \mathbf{u}_{0,\mathbf{H}} \end{bmatrix}, \quad (38)$$

where each matrix and vector is initialized as below.

- $\mathbf{u}_{0,\eta}$ : m-center
- $\mathbf{U}_{\eta,\eta}, \mathbf{W}_\eta$ : PCA
- $\mathbf{U}_{\mathbf{H},\mathbf{H}} = [\text{vec}(\boldsymbol{\Sigma}_1 + \tilde{\boldsymbol{\eta}}_1 \tilde{\boldsymbol{\eta}}_1^T), \dots, \text{vec}(\boldsymbol{\Sigma}_T + \tilde{\boldsymbol{\eta}}_T \tilde{\boldsymbol{\eta}}_T^T)]$
- $\mathbf{U}_{\eta,\mathbf{H}} = \mathbf{O}, \mathbf{U}_{\mathbf{H},\eta} = \mathbf{O}, \mathbf{W}_{\mathbf{H}} = \mathbf{I}, \mathbf{u}_{0,\mathbf{H}} = \mathbf{0}$

This implies that the features  $\mathbf{U}_{\eta,\eta}$  and the weights  $\mathbf{W}_\eta$  are extracted to characterize only the mean functions  $\{\boldsymbol{\mu}_t\}_{t=1}^T$  and not the covariance functions. Moreover, this expansion theoretically means that the covariance function can be reconstructed exactly i.e.,  $\{\boldsymbol{\Sigma}_t\}_{t=1}^T = \{\tilde{\boldsymbol{\Sigma}}(\mathbf{w}_t)\}_{t=1}^T$ .

## V. EXPERIMENTS

In this section, we aim to validate the proposed transfer learning BO algorithm w.r.t. state-of-the-art algorithms in terms of data efficiency, accuracy and computational time. BO-PCA is compared against the transfer learning BO algorithms described in Section III-A, namely ABLR [15], TST-R [27], TAF-R [28], RGPE [5] and Vanilla BO. The implementation of these baselines, and of the benchmark datasets was made available in the AutoML/transfer-hpo-framework repositories<sup>1</sup> [5], and we rely on the same implementation for fairness of comparison. The proposed algorithm is implemented in Python and uses GPy<sup>2</sup> for standard GPR. All computation-time based

experiments were performed with a 8-core i7-4790K CPU with 16 GB of memory on an Ubuntu machine.

### A. HPO Benchmark Problems

We consider 4 well known benchmark functions to evaluate the performance of the algorithms. Table II summarizes the main properties of each benchmark, such as the dimension of the hyperparameters to learn, the number of tasks, the number of sampling data collected for each task, and their respective references.

TABLE II  
SUMMARY OF BENCHMARKS

Name	# Dim	$N_{\text{tasks}}$	$N_t$	Reference
Quadratic Function	3	30	50	[15]
AdaBoost Grid	2	50	50	[18]
SVM Grid	6	50	50	[18]
NN Grid	7	35	50	[30]

The first benchmark consists of parameterized quadratic functions [15]:

$$f_t(\mathbf{x}|a_t, b_t, c_t) = a_t \|\mathbf{x}\|_2^2 + b_t \mathbf{1}^T \mathbf{x} + c_t, \quad (39)$$

where the three coefficients are sampled uniformly from  $a, b, c \sim \text{Uniform}(0.1, 10)$ . Then, we use two HPO benchmarks called Adaboost and SVM defined in [18]. These benchmarks were created using 50 classification datasets chosen randomly from the UCI repository<sup>3</sup>. They are denominated grid benchmarks, because their test accuracy was pre-computed on a discretized grid of hyperparameters, 2 and 6, respectively. Finally, the last HPO benchmark is composed of deep neural networks defined in [30]. This benchmark was created using training the networks in 35 different tasks in the OpenML datasets [7] and accounting for hyperparameters of dimension 7. The test accuracy was pre-computed for 2,000 configurations, which were sampled randomly. The hyperparameters meaning and search space of each benchmark are shown in Table I. The ‘log’ column shows whether the hyperparameters are converted into the logarithmic scale for numerical stability purposes.

<sup>1</sup><https://github.com/automl/transfer-hpo-framework>

<sup>2</sup><https://sheffieldml.github.io/GPy/>

<sup>3</sup><https://archive.ics.uci.edu/ml/index.php>

TABLE III  
NORMALIZED REGRET

	Quadratic [ $\times 10^5$ ]					AdaBoost [ $\times 10^2$ ]				
	10	20	30	40	50	10	20	30	40	50
Vanilla BO	519	10.7	0.82	0.69	0.59	4.07	2.27	<b>1.40</b>	<b>0.74</b>	<b>0.57</b>
ABLR	776	519	287	200	169	4.29	2.43	1.56	1.11	0.68
TST-R	734	704	677	598	536	4.46	2.76	1.76	1.16	<b>0.71</b>
TAF-R	78.8	8.33	1.01	0.84	0.72	<b>3.91</b>	2.36	<b>1.53</b>	<b>0.89</b>	<b>0.59</b>
RGPE	141	83.5	58.2	46.3	31.1	4.75	2.51	<b>1.55</b>	<b>0.92</b>	<b>0.55</b>
<b>BO-ePCA</b>	<b>108</b>	1.39	<b>0.48</b>	0.44	<b>0.40</b>	<b>4.02</b>	<b>2.26</b>	<b>1.59</b>	1.23	0.95
<b>BO-mPCA</b>	<b>76.7</b>	<b>0.79</b>	<b>0.42</b>	<b>0.35</b>	<b>0.34</b>	<b>3.91</b>	<b>2.21</b>	<b>1.52</b>	1.23	1.00
	SVM [ $\times 10^2$ ]					NN [ $\times 10^2$ ]				
	10	20	30	40	50	10	20	30	40	50
Vanilla BO	4.70	2.55	1.56	0.91	0.74	4.64	2.70	1.62	<b>1.08</b>	<b>0.86</b>
ABLR	4.29	2.43	1.56	1.11	0.68	5.12	3.64	3.30	2.73	2.43
TST-R	5.37	3.97	3.06	2.19	1.77	4.73	3.58	2.50	1.81	1.54
TAF-R	4.69	2.52	1.65	0.95	0.66	4.37	2.39	1.39	<b>0.98</b>	<b>0.85</b>
RGPE	5.72	3.94	2.91	1.96	1.43	<b>3.23</b>	<b>1.91</b>	1.50	1.21	1.03
<b>BO-ePCA</b>	<b>3.97</b>	1.77	<b>0.85</b>	<b>0.54</b>	<b>0.33</b>	<b>3.30</b>	<b>1.72</b>	<b>1.21</b>	<b>0.99</b>	<b>0.90</b>
<b>BO-mPCA</b>	<b>3.97</b>	<b>1.54</b>	<b>0.80</b>	<b>0.40</b>	<b>0.27</b>	<b>3.25</b>	<b>1.67</b>	<b>1.23</b>	<b>1.02</b>	<b>0.88</b>

### B. Experimental Design

For each benchmark, we perform the experiments in a leave-one-task-out style. That is, among all possible tasks, we pick one as the test task, and all the others are considered source tasks. This experiment is repeated in turn, so that each and every task becomes a test task.

From each source task,  $N_t = 50$  points were randomly collected and the number of basis functions in BO-PCA was set to  $L = 1$ . We used sparse GP-PCA for all experiments. The number of inducing points selected was  $|\mathbf{Z}| = 30$  for the Quadratic and AdaBoost benchmark, and  $|\mathbf{Z}| = 50$  for the SVM and NN benchmark. All the inducing points were collected by Latin hypercube sampling [23]. At test time, all 5 methods started by collecting the same  $n_0 = 5$  initial points, which were generated by a sequential model-free hyperparameter optimization [26], from the test task function. The number of total data collected for each test task was set to  $N_{T+1} = 50$ , corresponding to 45 iterations of the BO transfer learning algorithms.

After performing these experiments, so that each and every task became a test task, we averaged the performance across test tasks for each benchmark. Two evaluation metrics were computed: the average normalized regret and the average rank, that have been proposed as standard metrics for HPO problems in [2]. First, the average normalized regret represents the average of 0-1 range scaled regret on each task. The definition of average normalized regret after  $N$  iterations is

$$r_N = \frac{1}{N_{\text{tasks}}} \sum_{t=1}^{N_{\text{tasks}}} \min_{n \in [1, \dots, N]} \frac{y_{t,n} - f_t^{\min}}{f_t^{\max} - f_t^{\min}}, \quad (40)$$

where  $N_{\text{tasks}}$  is a number of tasks of the benchmark,  $y_{t,n}$  is the observation of task  $t$  at the  $n$ -th iteration and  $f_t^{\min}$  and

$f_t^{\max}$  are the best and worst global values of the function that are available for evaluation purposes only:

$$f_t^{\min} = \min_{\mathbf{x}} f_t(\mathbf{x}), \quad f_t^{\max} = \max_{\mathbf{x}} f_t(\mathbf{x}). \quad (41)$$

Second, the competing methods are ranked at each iteration accordingly to the normalized regret, going from the best method with score equal 1 to the worst one with increasing score. Every experiment is repeated 15 times and each time a different set of  $N_t = 50$  data point is randomly selected for each source task.

TABLE IV  
COMPUTATION TIME [MS]

Method	Quadratic			AdaBoost		
	10	30	50	10	30	50
ABLR	33K	14K	35K	83K	31K	2.9K
TST-R	13.7	40.1	155	24.1	68.5	234
TAF-R	13.7	40.3	154	24.2	68.4	233
RGPE	63.5	234	762	86.4	334	1072
<b>BO-ePCA</b>	299	310	473	330	335	464
<b>BO-mPCA</b>	<b>1.50</b>	<b>1.59</b>	<b>1.73</b>	<b>1.57</b>	<b>1.57</b>	<b>1.71</b>
Method	SVM			NN		
	10	30	50	10	30	50
ABLR	60K	82K	61K	32K	30K	9.6K
TST-R	29.4	78.3	259	16.9	50.4	188
TAF-R	27.0	76.6	244	16.5	50.0	192
RGPE	100	375	1764	69.9	261	847
<b>BO-ePCA</b>	367	374	616	322	333	471
<b>BO-mPCA</b>	<b>2.68</b>	<b>3.26</b>	<b>2.71</b>	<b>2.28</b>	<b>2.25</b>	<b>2.35</b>

### C. Results

The results of the normalized regret are shown in Table III. For each method, we compared the averaged normalized regret after 10, 20, 30, 40 and 50 iterations of BO. We underline the best value per benchmark at each number of

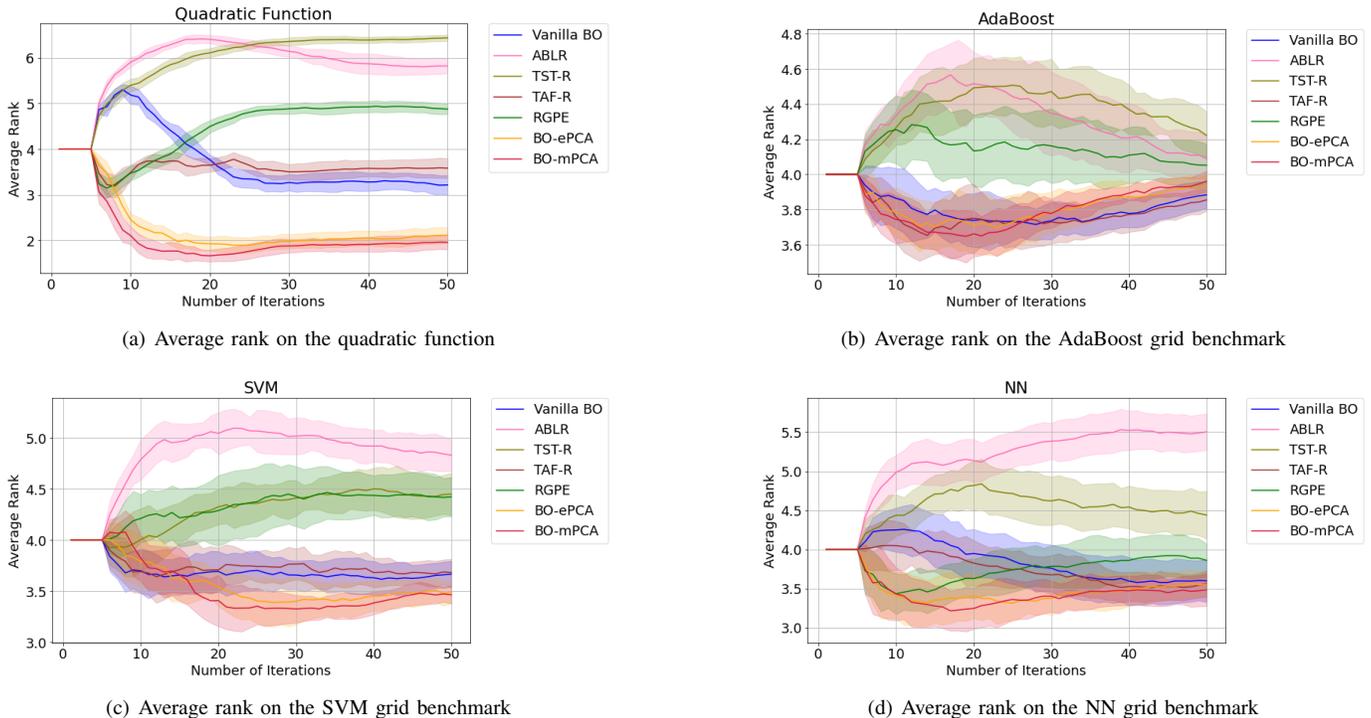


Fig. 1. Average rank

iterations and show in boldface the values which are not significantly different from the best value according to a Wilcoxon signed-rank test [3] with significant level  $\alpha = 0.05$ . BO-PCA outperformed or performed as well as the best of the other transfer learning algorithms in almost all the benchmarks and all the number of iterations. Indeed, when the average normalized regret of BO-PCA does not have the best value, it is not statistically different from the best method. Only in the AdaBoost benchmark at 40 and 50 iterations, the proposed method was outperformed. This first metric tell us that BO-PCA was empirically the most data-efficient algorithm and the most accurate in most of the cases. These results are confirmed also by the rank shown in Figure 1, where each of the figures displays the average rank and its standard deviation for one of the four benchmarks. Also, the rank shows that BO-PCA is either the best or comparable to the best competitor, which sometimes is RGPE and sometimes plain BO, for each iteration of each task.

Finally, we also compare the computational time needed to update the transfer learning models of the 6 methods. This corresponds to the time to update the Bayesian linear regression for ABLR, and to the time to update the weight vector for the other transfer learning methods. Note that this is the discriminatory step, as all the other steps in the BO iteration are equivalent in all the algorithms (except for ABLR that is anyway performing the worst). The results are shown in Table IV.

The weights of GP-mPCA are updated much faster than any other method. This is because it is the only algorithm that can be updated recursively, as described in Section IV-B.

In addition, although the computational time of TST-R, TAF-R and RGPE increases at each iteration, GP-mPCA remains constant, as it does not depend on the number of data. Note that the computational time for ABLR is much slower than the others, because of numerical issues related to matrix inversion that cause long computations for numerical adjustments.

We can conclude that BO-PCA outperforms the s.o.t.a. algorithms in the analyzed benchmarks. Comparing the two variants of this algorithm (using GP-ePCA and GP-mPCA), we obtain similar results in terms of performance, but with GP-mPCA always slightly outperforming GP-ePCA. This supports the efficacy of the proposed initialization method designed specifically for GP-mPCA and described in Section IV-C. Given these results, and the fact that GP-mPCA is much faster in terms of computational time, we can conclude that BO-PCA equipped with GP-mPCA is the most promising algorithm for transfer learning BO.

## VI. DISCUSSION AND CONCLUSIONS

In this work, we proposed an efficient BO algorithm for transfer learning based on GP-PCA, called BO-PCA. It first learns low-dimensional features from a set of Gaussian processes posteriors of the source tasks using GP-PCA. Then, it updates the prior functions by optimizing the low dimensional parameters of the extracted features to fit the target task. The approach could be used with both GP-ePCA and GP-mPCA. In the latter case, the transfer learning procedure becomes a least squares problem which can be solved recursively and independent of the number of data. Therefore, our algorithm requires much lower computational cost compared to previous

methods. We show that our algorithm either works more accurately and is more data efficient or works as well as the previous methods for HPO tasks, even though it incurs a much lower computational cost.

## REFERENCES

- [1] S. Akaho. e-pca and m-pca: Dimension reduction of parameters by information geometry. In *IEEE international Joint Conference on Neural Networks*, volume 1, pages 129–134, 2004.
- [2] S. P. Arango, H. S. Jomaa, M. Wistuba, and J. Grabocka. Hpo-b: A large-scale reproducible benchmark for black-box hpo based on openml. *arXiv preprint arXiv:2106.06257*, 2021.
- [3] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7(1):1–30, 2006.
- [4] S. Falkner, A. Klein, and F. Hutter. Bohb: Robust and efficient hyperparameter optimization at scale. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1437–1446, 2018.
- [5] M. Feurer, B. Letham, and E. Bakshy. Scalable meta-learning for bayesian optimization using ranking-weighted gaussian process ensembles. In *AutoML Workshop at ICML*, volume 7, 2018.
- [6] P. I. Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- [7] P. Gijbbers, E. LeDell, J. Thomas, S. Poirier, B. Bischl, and J. Vanschoren. An open source automl benchmark. *arXiv preprint arXiv:1907.00909*, 2019.
- [8] N. Hansen, D. V. Arnold, and A. Auger. *Evolution Strategies*, pages 871–898. Springer Berlin Heidelberg, 2015.
- [9] F. Hutter, L. Kotthoff, and J. Vanschoren, editors. *Automatic Machine Learning: Methods, Systems, Challenges*. Springer, 2019.
- [10] H. Ishibashi and S. Akaho. Principal component analysis for gaussian process posteriors. *Neural Computation*, 34:1189–1219, 2022.
- [11] N. Jaquier, V. Borovitskiy, A. Smolensky, A. Terenin, T. Asfour, and L. Rozo. Geometry-aware bayesian optimization in robotics using riemannian matern kernels. In *Conference on Robot Learning (CoRL)*, 2021.
- [12] H. Liu, Y.-S. Ong, X. Shen, and J. Cai. When gaussian process meets big data: A review of scalable gps. *IEEE transactions on neural networks and learning systems*, 31(11):4405–4423, 2020.
- [13] L. Ljung. System identification. In *Signal analysis and prediction*, pages 163–173. Springer, 1998.
- [14] J. Mockus, V. Tiesis, and A. Zilinskas. The application of bayesian methods for seeking the extremum. *Towards global optimization*, 2(117-129):2, 1978.
- [15] V. Perrone, R. Jenatton, M. Seeger, and C. Archambeau. Scalable hyperparameter transfer learning. In *Proceedings of the 31st International Conference on Advances in Neural Information Processing Systems*, pages 12751–12761, 2018.
- [16] A. Rai, R. Antonova, F. Meier, and C. Atkeson. Using simulation to improve sample-efficiency of bayesian optimization for bipedal robots. *Journal of Machine Learning Research*, 20:1–24, 2019.
- [17] E. Raponi, H. Wang, M. Bujny, S. Boria, and C. Doerr. High dimensional bayesian optimization assisted by principal component analysis. In *International Conference on Parallel Problem Solving from Nature*, pages 169–183. Springer, 2020.
- [18] N. Schilling, M. Wistuba, and L. Schmidt-Thieme. Scalable hyperparameter optimization with products of gaussian process experts. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 33–48. Springer, 2016.
- [19] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- [20] J. Snoek, H. Larochelle, and R. Adams. Practical bayesian optimization of machine learning algorithms. In *Proceedings of the 25th International Conference on Advances in Neural Information Processing Systems*, pages 2960–2968, 2012.
- [21] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 25, 2012.
- [22] N. Sugihara, H. Sugihara, Y. Otake, R. Krems, H. Nakamura, and S. Fuse. Rapid and mild one-flow synthetic approach to unsymmetrical sulfamides guided by bayesian optimization. *Chemistry-Methods*, 1(11):484–490, 2021.
- [23] B. Tang. Orthogonal array-based latin hypercubes. *Journal of the American statistical association*, 88(424):1392–1397, 1993.
- [24] J. Vanschoren. Meta-learning. In Hutter et al. [9], pages 39–68.
- [25] C. K. Williams and C. E. Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- [26] M. Wistuba, N. Schilling, and L. Schmidt-Thieme. Sequential model-free hyperparameter tuning. In *2015 IEEE international conference on data mining*, pages 1033–1038. IEEE, 2015.
- [27] M. Wistuba, N. Schilling, and L. Schmidt-Thieme. Two-stage transfer surrogate model for automatic hyperparameter optimization. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 199–214. Springer, 2016.
- [28] M. Wistuba, N. Schilling, and L. Schmidt-Thieme. Scalable gaussian process-based transfer surrogates for hyperparameter optimization. *Machine Learning*, 107(1):43–78, 2018.
- [29] Y. Zhang, D. Apley, and W. Chen. Bayesian optimization for materials design with mixed quantitative and qualitative variables. *Scientific Reports*, 10(1):1–13, 2020.
- [30] L. Zimme, M. Lindauer, and F. Hutter. Auto-pytorch tabular: Multifidelity metalearning for efficient and robust autodl. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(9):3079 – 3090, 2021.