

# Which Surrogate Works for Empirical Performance Modelling? A Case Study with Differential Evolution

Ke Li<sup>†\*</sup>, Zilin Xiang<sup>\*</sup>, Kay Chen Tan<sup>‡</sup>

<sup>†</sup>*Department of Computer Science, University of Exeter  
EX4 4QF, Exeter, UK*

<sup>\*</sup>*College of Computer Science and Engineering, University of Electronic Science and Technology of China  
611731, Chengdu, China*

<sup>‡</sup>*Department of Computer Science, City University of Hong Kong  
Hong Kong SAR, China  
k.li@exeter.ac.uk, zilin.xiang@hotmail.com, kaytan@cityu.edu.hk*

**Abstract**—It is not uncommon that meta-heuristic algorithms contain some intrinsic parameters, the optimal configuration of which is crucial for achieving their peak performance. However, evaluating the effectiveness of a configuration is expensive, as it involves many costly runs of the target algorithm. Perhaps surprisingly, it is possible to build a cheap-to-evaluate surrogate that models the algorithm’s empirical performance as a function of its parameters. Such surrogates constitute an important building block for understanding algorithm performance, algorithm portfolio/selection, and the automatic algorithm configuration. In principle, many off-the-shelf machine learning techniques can be used to build surrogates. In this paper, we take the differential evolution (DE) as the baseline algorithm for proof-of-concept study. Regression models are trained to model the DE’s empirical performance given a parameter configuration. In particular, we evaluate and compare four popular regression algorithms both in terms of how well they predict the empirical performance with respect to a particular parameter configuration, and also how well they approximate the parameter *versus* the empirical performance landscapes.

**Index Terms**—Empirical performance modelling, parameter configuration, landscape analysis, differential evolution.

## I. INTRODUCTION

Meta-heuristic algorithms are normally accompanied by some parameters which can influence their search behaviour on various optimisation problems. Parameter optimisation (PO) aims to find a best possible parameter configuration  $\theta^*$  from the parameter space  $\Theta$ , which consists of all possible configurations, of the target algorithm and helps it achieve its peak performance on a black-box optimisation problem. Formally, given an algorithm, PO can be defined as the following black-box meta-optimisation problem:

$$\begin{aligned} & \text{minimize } \mathcal{L}(f(\mathbf{x}), \theta) \\ & \text{subject to } \theta \in \Theta \end{aligned} \quad (1)$$

where  $f(\mathbf{x})$  is the optimisation problem under consideration, and  $\mathbf{x} \in \mathbb{R}^d$  is a decision variable.  $\mathcal{L}(f(\mathbf{x}), \theta)$  is the performance measure associated with a configuration  $\theta$  of the target algorithm. In particular, it can either be the runtime cost (e.g. the CPU wall time and/or the number of function evaluations) or the error of the solution found by the target algorithm.

PO is a challenging black-box meta-optimisation problem. First, its landscape is complex and change with the target algorithm when solving different problems. Second, the parameters associated with the target algorithm can have various types (e.g. numerical, integer and categorical) and the number of parameters can be potentially large depending on the algorithm specification. In addition, PO is intrinsically expensive as it requires to explore  $\Theta$  by running the target algorithm with different configurations, where evaluating the effectiveness of a configuration will in turn cost a large amount of function evaluations and/or CPU wall time. In the evolutionary computation (EC) community, constructing a cheap-to-evaluate surrogate in lieu of calling the physically expensive objective function has been widely accepted as an effective way for expensive optimisation [1]. The design and analysis of computer experiments in statistics also uses surrogate models to either fit a global model of the overall landscape or sequentially identify the global optimum of the underlying function [2]. In the automatic parameter configuration field, sequential model-based Bayesian optimisation methods [3]–[5] have shown strong performance in PO, compared to some traditional methods like grid search and random search [6] and can compete or even surpass the results tuned by experienced human experts. Moreover, regression models have been extensively used in meta-learning to predict the algorithm performance across various datasets [7]. Note that all these lines of research need to construct surrogate models of a computationally expensive and complex function in order to inform an active learning criterion that identifies new inputs to evaluate.

The problem of PO has a long history dating back to the 90s [8]. Recently, it becomes increasingly popular in both meta-heuristics (e.g. [3], [4], [9], [10]) and machine learning (e.g. [5], [11], [12]) communities, especially with the

This work was supported by UKRI Future Leaders Fellowship under grant MR/S017062/1

development of emerging automated machine learning [13]. In this paper, instead of developing new algorithms for PO, we focus on studying surrogate models, which sit in the core of the model-based PO framework. We take the differential evolution (DE) [14], one of the most popular black-box optimiser in the EC community, as the baseline algorithm. To obtain the empirical performance data on a given optimisation problem, we evaluate the performance of DE with respect to 5,940 parameter configurations in an expensive offline phase. The collected performance data are used to train a regression model and to validate its generalisation ability for predicting empirical performance of unseen parameter configurations. Here we consider four off-the-shelf regression algorithms for empirical performance modelling. In particular, we evaluate and compare their abilities in terms of how well they predict the empirical performance with respect to a particular parameter configuration, and also how well they approximate the parameter configuration *versus* the empirical performance landscapes. We envisage that this aspect will shed light on the study of the characteristics of surrogate models in future.

The rest of this paper is organised as follows. Section II describes the methodologies that we used to setup the experiments. Section III presents and analyses the experimental results. Finally, Section IV concludes this paper and provides some future directions.

## II. METHODOLOGY

This section mainly describes the benchmark problems chosen in our empirical studies, the baseline algorithm DE and its corresponding parameters, the performance measure used to evaluate the quality of a particular parameter configuration, the method used to collect the algorithm performance data, and the regression algorithms used to build surrogates for modelling the empirical performance.

### A. Benchmark Problems

In this paper, we consider choosing six widely used elementary test problems (i.e. sphere, ellipsoid, rosenbrock, ackley, griewank and rastrigin) and the first fourteen test problems (i.e. excluding those hybrid composite functions) from the CEC 2005 competition [15] to constitute the benchmark problems. To facilitate the notation in Section III, the six elementary functions are denoted as F1 to F6 and those from the CEC 2005 competition are denoted as F7 to F20. Note that these test problems have various characteristics. In particular, F1, F2 and F7 to F11 are unimodal functions while the others are multi-modal functions. All test problems have analytically defined continuous objective functions with a known global optimum. The number of variables of each test problem varies from 2 to 30 (in particular  $d \in \{2, 10, 30\}$ ) and the range of variables is set according to their original paper.

### B. DE and its Parameters

DE [14] is one of the most popular black-box optimisation algorithm in the EC community. One of the major reasons that contributes to its success is its simple structure. For a

vanilla DE, an offspring solution  $\mathbf{x}^c$  is generated by a two-step procedure. First, a trial vector  $\bar{\mathbf{x}}$  is generated as:

$$\bar{\mathbf{x}} = \mathbf{x}^1 + F \times (\mathbf{x}^2 - \mathbf{x}^3) \quad (2)$$

where  $F \in (0, 3]$ , known as the evolution step size, is a parameter of DE.  $\mathbf{x}^1$ ,  $\mathbf{x}^2$  and  $\mathbf{x}^3$  are randomly chosen from the parent population. Afterwards,  $\mathbf{x}^c$  is generated as:

$$x_i^c = \begin{cases} \bar{x}_i & \text{if } (\text{rand} < CR) \vee (i = j) \\ x_i & \text{otherwise} \end{cases} \quad (3)$$

where  $i \in \{1, \dots, d\}$ ,  $j$  is an integer randomly chosen from 1 to  $d$ .  $\mathbf{x}$  is the parent solution under consideration.  $\text{rand}$  is a random number chosen from 0 to 1, and  $CR \in [0, 1]$ , known as the crossover rate, is another parameter of DE. In addition, the population size  $NP \in \mathbb{N}$  is also a parameter.

Many studies have demonstrated that the performance of DE is highly sensitive to its parameter settings [16]. During the past decade, many efforts have been devoted to the development of advanced DE variants that are able to adaptively set the parameters on the fly [17]–[19] and/or find a good configuration in an offline manner [20]. Since the major purpose of this paper is to investigate the ability of building the surrogate for modelling the empirical performance of an algorithm with respect to its corresponding parameter configurations, we focus on the vanilla DE [14] which is simple yet without losing the generality of the observations. Obviously,  $NP$  is an integer parameter, while  $F$  and  $CR$  are numerical parameters.

### C. Performance Measure

As the global optimum of each test problem is known a priori, this paper uses the approximation error to evaluate the empirical performance of a particular parameter configuration. Specifically, it is computed as:

$$\Psi(f(\mathbf{x}), \theta) = f(\mathbf{x}) - f(\mathbf{x}^*) \quad (4)$$

where  $\theta$  is a parameter configuration of DE,  $\mathbf{x}$  is the best-so-far solution found by the DE with the parameter configuration  $\theta$ , and  $\mathbf{x}^*$  is the global optimum. Since DE is a stochastic algorithm, each parameter configuration needs to be repeated more than one time in practice. Thus, the performance of a parameter configuration  $\theta$  is measured as an averaged approximation error:

$$\mathcal{L}(f(\mathbf{x}), \theta) = \frac{1}{n} \sum_{i=1}^n \Psi_i(f(\mathbf{x}), \theta) \quad (5)$$

where  $\Psi_i(f(\mathbf{x}), \theta)$  is the approximation error of a configuration  $\theta$  at the  $i$ -th run and  $n$  is the number of repetitions of experiments with  $\theta$  where we set  $n = 31$  in our experiments.

### D. Data Collection

In principle, algorithm performance data used to construct the surrogate model of an algorithm's empirical performance can be obtained by any means. Since this paper aims to investigate the overall surrogate modelling ability of an algorithm's performance with respect to its parameter space, we

are interested in every corner of the space. To this end, the parameter space is sampled in a grid manner, where we chose 9 different  $NP$  settings, i.e.  $NP = i \times d$ ,  $i \in \{2, \dots, 10\}$ , 60 different values for  $F \in (0, 3]$  with a step size 0.05, and 11 different values for  $CR \in [0, 1]$  with a step size 0.1. Therefore, there are 5,940 different parameter configurations in total.

#### E. Regression Algorithms for Surrogate Modelling

In this paper, four regression algorithms, i.e. Gaussian process (GP), random forest (RF), support vector machine for regression (SVR), radial basis function networks (RBFN), are considered as the candidates for surrogate modelling of DE's empirical performance. Note that these regression algorithms have been widely used in the model-based PO in the algorithm configuration literature [21].

To construct a surrogate model on a particular problem instance, each of these four models is trained on the performance data (only 70% of them are used for training while the remaining 30% are used for testing) collected by running the DE algorithm with various parameter configurations on each problem instance as introduced in Section II-D. Note that learning a surrogate model is no free lunch, as each regression algorithm also requires some hyper-parameters to be tuned. To identify the best possible configurations for each regression algorithm, we apply the random search [6] to explore the hyper-parameter space. Specifically, as for GP, we need to choose an appropriate kernel among RBF, rational quadratic and Matérn; as for RF, the number of trees in a forest is chosen from 2 to 100, the minimum number of samples required to split an internal node is chosen from 2 to 11, the number of features to consider when looking for the best split is set in the range  $[0.001, 1]$ , the criterion used to measure the quality of a split is either mean squared error or mean absolute error and the minimum number of samples required to be at a leaf node is chosen from 1 to 11; as for SVR, the kernel is chosen between RBF and Sigmoid, the maximal margin  $\epsilon$  is chosen from  $[0.01, 1]$ , the regularisation parameter  $C$  is set in between 1 and 10, and  $\gamma$  is chosen from  $[0.01, 1]$  if RBF is used as the kernel. A 5-fold cross-validation (using 80% of the training data for training and the remaining 20% data for testing) is used to evaluate the training performance of a particular hyper-parameter configuration of a regression algorithm. To have a fair comparison, all surrogate modelling procedures are implemented by `scikit-learn`, a machine learning toolbox in Python<sup>1</sup>.

### III. EXPERIMENTS AND RESULTS

In this section, we will present and compare experimental evaluations of the quality of surrogates constructed by different regression algorithms introduced in Section II-E. The experimental results are analysed according to the following three research questions (RQs).

**RQ1:** *Which surrogate model works best for empirical performance modelling on various kinds of benchmark problems?*

<sup>1</sup><https://scikit-learn.org/stable/>

**RQ2:** *Does the empirical performance predicted by a surrogate model follow the order as the ground truth?*

**RQ3:** *How does the empirical performance landscape fit by a surrogate model compare with the ground truth?*

#### A. Comparisons of Different Surrogate Models

Bearing the RQ1 in mind, this section empirically compares the generalisation performance of four regression algorithms on unseen parameter configurations. In particular, the root mean square error (RMSE) is used to measure the generalisation performance and it is calculated as:

$$RMSE = \sqrt{\frac{\sum_{i=1}^{\hat{n}} (\hat{L}(f(\mathbf{x}), \theta_i) - L(f(\mathbf{x}), \theta_i))^2}{\hat{n}}} \quad (6)$$

where  $\hat{L}(f(\mathbf{x}), \theta_i)$  is the approximation error of a parameter configuration  $\theta_i$  estimated by a surrogate model; while  $L(\mathbf{x}, \theta_i)$  is the observed approximation error of  $\theta_i$ ,  $i \in \{1, \dots, \hat{n}\}$  and  $\hat{n}$  is the number of data in the testing set.

From the results shown in Tables I to III, we clearly see that GP and RF are the best regression algorithms to build the surrogate for modelling the empirical performance. RBFN is slightly worse than GP and RF, while SVR is the worst choice except on F14 when  $d = 2$ . Note that our observations of promising performance of GP and RF are also in line with some results reported in the contemporary algorithm configuration literature [21]. Furthermore, we find that the performance of different regression algorithms are consistent across different dimensions. This makes sense as a surrogate model is built upon the parameter configurations themselves, which are independent from the problem instances. In addition, we find that the RMSE dramatically increases with the dimensionality of the underlying problem. This can be explained as the significant degeneration of the performance of DE with the dimensionality which in term largely increases the approximation errors.

To have a better understanding of the generalisation performance of different surrogate models (especially the relationship between the predicted performance and its ground truth given a particular parameter configuration), we calculate the Pearson correlation coefficient (PCC) of the results:

$$PCC = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad (7)$$

where  $X$  represents the set of observed approximation errors of all parameter configurations in the testing set while  $Y$  is the set of approximation errors estimated by a surrogate model.  $\text{cov}(X, Y)$  is the covariance of  $X$  and  $Y$ ,  $\sigma_X$  and  $\sigma_Y$  are the standard deviations of  $X$  and  $Y$ . In particular, a higher PCC indicates a better correlation between the predicted performance and the ground truth.

From the results shown in Tables I to III, we can see that the observations are in line with the RMSE. The performance of GP and RF are the most competitive regression algorithms in almost all cases, where the correlation between the predicted



performance and its ground truth is relatively high. The performance of RBFN is very close to those of GP and RF, while the PCC obtained by SVR is the worst. To have a visual understanding of this point, we also provide the scatter plots of *ground truth vs predicted performance* in Figures 1 to 3<sup>2</sup>. According to the observations from these figures and Tables I to III, we summarise our findings as follows.

- As shown in Tables I to III, the RMSEs of all four regression algorithms are huge (over  $10^7$ ) on F9 and F12. This is because the performance of DE are miserable on these two test problems with almost all sampled 5,940 parameter configurations. Accordingly, the deviations of the predicted empirical performance are in a relatively large scale. This also explains the increase of RMSEs with the problem dimensionality. However, according to PCCs, we find that the correlation between the predicted empirical performance and the ground truth of GP, RBFN and RF are acceptable.
- The RMSEs of the first six elementary test problems (i.e. F1 to F6), which are relatively simple, are better than those from CEC 2005 competition. Accordingly, the deviations between the predicted performance and the ground truth are small. This indicates that most parameter configurations are able to lead to an acceptable performance of DE. In other words, DE is not sensitive to its configurations on these problems.
- As shown in Fig. 2, we find that SVR largely underestimates the approximation error on F8. Similar observations can be found on F7, F9, F10, F12 and F18 as shown in the supplementary document.
- As shown in Fig. 3, we find that scatter plots are crowded in the middle region of the diagonal line. This implies that all parameter configurations fail to lead to a decent result. Similar observations can be found on F13 and F20 when the number of variables becomes large in the supplementary document.

Based on the above discussions, we come up with the following response to RQ1:

**Response to RQ1: GP and RF are the best regression algorithms for building the surrogate model of empirical performance. In addition, the quality of the surrogate model depend on the quality of the performance data.**

#### B. Comparisons of Performance Ranks Obtained by Different Surrogate Models

When using a surrogate in a sequential model-based PO, the prediction accuracy of this model is not utterly important. Instead, reliably differentiating the promising ones with respect to their unpromising counterparts can also provide useful information to guide the optimisation process. In other words, for a set of parameter configurations, we expect that the ranks (or the order) of the empirical performance predicted by a surrogate model can follow those of the ground truth. To this

<sup>2</sup>More comprehensive figures are moved to the supplementary document, which can be downloaded from <http://coda-group.github.io/cec19-sup.pdf>.

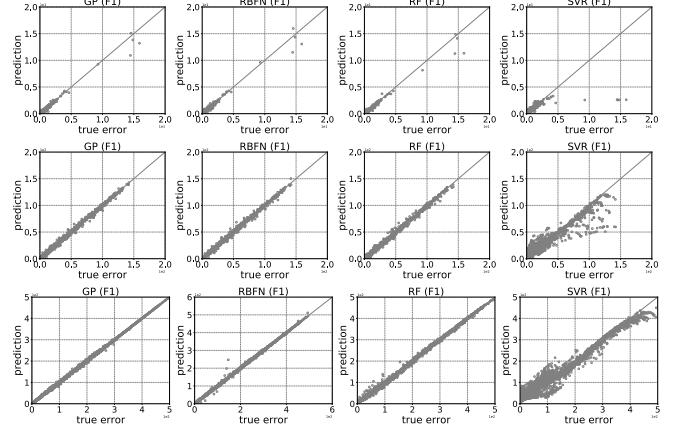


Fig. 1: Scatter plots of the empirical performance predicted by a surrogate model *vs* the observed empirical performance on the testing set (i.e. unseen parameter configurations). In particular, three rows respectively represent results on F1 where  $d = 2, 10, 30$ .

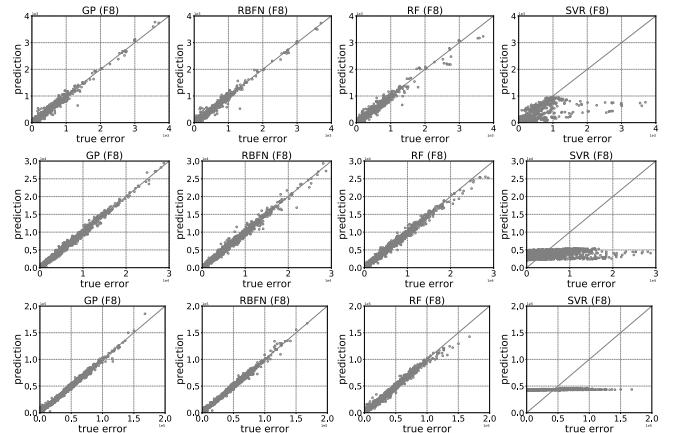


Fig. 2: Scatter plots of the empirical performance predicted by a surrogate model *vs* the observed empirical performance on the testing set (i.e. unseen parameter configurations). In particular, three rows respectively represent results on F8 where  $d = 2, 10, 30$ .

end, we consider using the Spearman's rank correlation coefficient (SRCC) to measure the statistical dependence between the ranks of the predicted performance and the ground truth. Note that the calculation of SRCC is almost the same as that of PCC, except that the raw data is replaced by the corresponding ranks.

$$SRCC = \frac{\text{cov}(r_X, r_Y)}{\sigma_{r_X} \sigma_{r_Y}} \quad (8)$$

where  $r_X$  indicates the ranks of the observed approximation errors of all parameters configurations in the testing set while  $r_Y$  is the ranks of those estimated approximation errors. A higher SRCC indicates a better dependency between the predicted performance and the ground truth.

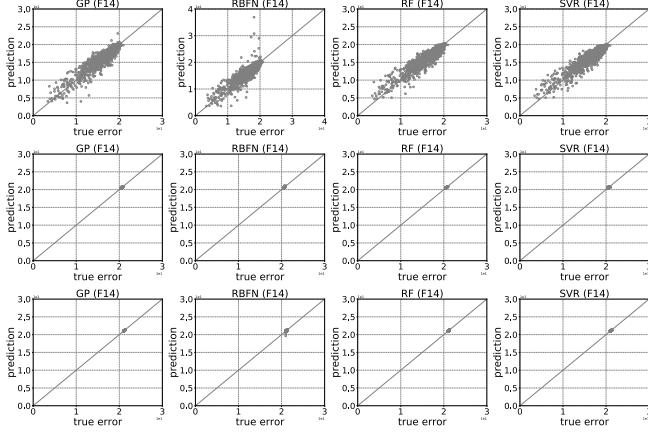


Fig. 3: Scatter plots of the empirical performance predicted by a surrogate model *vs* the observed empirical performance on the testing set (i.e. unseen parameter configurations). In particular, three rows respectively represent results on F14 where  $d = 2, 10, 30$ .

From the results shown in Tables I to III, we can still come up with the conclusion that GP and RF are the most reliable regression algorithms for building the surrogate model of the empirical performance. They almost dominate the top two positions in terms of SRCC. It is interesting to note that the SRCCs obtained by SVR are not as poor as its performance on RMSE and PCC. It is even comparable with GP and RF in some cases, e.g. on F20. This suggests that the prediction made by SVR has a decent chance to differentiate the order between two parameter configurations. In this case, SVR might be useful in a model-based PO process where it can be used as a comparison-based surrogate [22]. Furthermore, we also notice that RBFN does not show a good performance on SRCC. It is even sometimes worse than SVR. This indicates that although the prediction made by RBFN is numerically close to the ground truth, it may still mislead a model-based PO as it messes up the order of similar parameter configurations.

Based on the above discussion, we come up with the following response to RQ2:

**Response to RQ2: GP and RF are able to preserve the order of the empirical performance of different parameter configurations. In particular, SVR, which performs poorly on predicting the empirical performance, shows comparable performance for order preservation.**

#### C. Comparisons of Landscape Approximation

In previous subsections, we mainly focus on investigating the quality of surrogate models from the approximation accuracy perspective. For the last RQ, we plan to study of the quality of surrogate models from a landscape analysis perspective. Considering the testing data set, we compare the landscapes of the empirical performance predicted by different regression algorithms to the landscape of the ground truth. To this end, we use the kernel density estimation (KDE)

method<sup>3</sup> to estimate a probability density function (PDF) of the empirical performance. To have a visual comparison, Figs 4 to 6 shows the plots of the estimated PDFs of four different regression algorithms and the ground truth. From these figures, we can see that the prediction made by GP, RF and RBFN almost fit the distribution of the ground truth. In contrast, the estimated PDF of SVR deviates from the ground truth in many cases. This becomes more evident when the dimensionality of the underlying problem becomes large.

Since the surrogate model considered in this paper is a mapping between a parameter configuration and its corresponding empirical performance, it is interesting to consider a more complex landscape that is a joint probability distribution of parameter configuration and empirical performance. As it is non-trivial to visualise a multi-dimensional distribution, we try to understand the proximity of the landscape approximated by the surrogate model and that of the ground truth from a statistical distance perspective. To this end, we apply the earth mover's distance (EMD) [23], also known as Wasserstein metric, to evaluate the dissimilarity between two multi-dimensional distributions. Generally speaking, given two distributions, the EMD measures the minimum cost of turning one distribution into the other. In our context, similar landscapes are expected to have a relatively small EMD whereas large EMD values will imply that the landscapes are significantly different from each other. Due to the page limit, we do not intend to elaborate the calculation procedure of EMD, interested readers can refer to [23] for more details. From the comparison results of EMD values shown in Table IV, we find that GP, RF and RBFN have the same level of approximation to the ground truth whereas the divergence values obtained by SVR are relatively large in almost all cases. All these observations are also in line with the RMSEs discussed in Section III-A.

Based on the above discussion, we come up with the following response to RQ3:

**Response to RQ3: The landscapes of the empirical performance predicted GP, RF and RBFN well approximate the ground truth; while the landscapes obtained by SVR deviate from the ground truth to a certain extent.**

#### IV. CONCLUSIONS AND FUTURE DIRECTIONS

It is not uncommon a meta-heuristic algorithm is accompanied by some parameters, the settings of which largely influence its performance on various problems. Tweaking the parameter configuration of a meta-heuristic algorithm to achieve its peak performance on a certain problem can be treated as an optimisation process, as known as PO. Due to the stochastic property of most meta-heuristic algorithms, evaluating the quality of a particular parameter configuration usually requires to run the target algorithms several times. Therefore, it is inarguably that PO is computationally expensive. Building a cheap-to-evaluate surrogate model in lieu of a computationally expensive experiment has been widely accepted as a major approach for expensive optimisation.

<sup>3</sup><https://uk.mathworks.com/help/stats/ksdensity.html>

TABLE IV: Comparisons of EMD between the surrogate model built by four regression algorithms and the ground truth

Problem	$d$	GP	RBFN	RF	SVR	Problem	$d$	GP	RBFN	RF	SVR
F1	2	3.9123E-2	4.1131E-2	<b>3.8218E-2</b>	2.1449E-1	F11	2	<b>1.6881E+0</b>	2.0925E+0	2.1808E+0	1.0413E+1
	10	7.4359E-1	<b>7.0765E-1</b>	8.3284E-1	4.3693E+0		10	1.8064E+1	<b>1.6778E+1</b>	1.8437E+1	2.7649E+2
	30	<b>1.4450E+0</b>	1.8732E+0	2.7342E+0	1.0298E+1		30	1.2526E+2	<b>9.1694E+1</b>	1.6699E+2	5.3056E+3
F2	2	7.7335E-1	8.1710E-1	<b>7.6167E-1</b>	1.9532E+0	F12	2	<b>1.2890E+6</b>	1.3150E+6	2.1687E+6	2.3287E+7
	10	<b>1.7967E+1</b>	1.9940E+2	2.4552E+2	3.6687E+3		10	<b>1.2284E+7</b>	1.2908E+7	1.7883E+7	4.9220E+8
	30	1.8648E+3	<b>1.5890E+3</b>	2.9031E+3	2.4368E+5		30	1.7468E+8	<b>1.7299E+8</b>	2.3724E+8	7.7375E+9
F3	2	2.5151E+1	2.6593E+1	<b>1.7417E+1</b>	3.0550E+1	F13	2	<b>4.0126E-1</b>	6.3972E-1	8.5819E-1	3.8741E+0
	10	<b>8.5900E+2</b>	1.1733E+3	9.8913E+2	1.5110E+4		10	<b>6.4477E+0</b>	1.8491E+1	7.6609E+0	1.6785E+2
	30	8.8047E+3	<b>7.7415E+3</b>	1.2159E+4	7.9610E+5		30	4.5265E+1	<b>6.6990E+1</b>	2.6507E+1	7.9910E+2
F4	2	2.2263E-1	2.5603E-1	<b>1.2124E-1</b>	4.1061E-1	F14	2	<b>3.4411E-1</b>	3.5736E-1	3.9440E-1	4.1817E-1
	10	<b>2.4946E-1</b>	2.6379E-1	3.5331E-1	1.1068E+0		10	2.3350E-2	7.1733E-2	<b>2.0398E-2</b>	4.7365E-2
	30	<b>7.6201E-2</b>	1.5856E-1	1.4125E-1	6.4510E-1		30	1.6935E-2	7.1612E-2	<b>1.4588E-2</b>	6.7035E-2
F5	2	1.0446E-2	1.1772E-2	<b>1.0122E-2</b>	1.1300E-2	F15	2	<b>2.3576E-1</b>	2.4196E-1	2.6502E-1	3.1265E-1
	10	2.7905E-2	2.7986E-2	<b>2.7347E-2</b>	9.4473E-2		10	<b>1.0227E+0</b>	1.0634E+0	1.2373E+0	3.0674E+0
	30	<b>3.3163E-2</b>	3.6446E-2	4.8344E-2	2.5692E-1		30	3.2431E+0	4.1342E+0	5.0520E+0	3.4736E+1
F6	2	2.1599E-1	2.3193E-1	<b>2.0715E-1</b>	3.0823E-1	F16	2	<b>2.8404E-1</b>	3.1241E-1	3.2114E-1	3.9340E-1
	10	<b>1.3799E+0</b>	1.5950E+0	4.9727E+0	1.4138E+0		10	1.3076E+0	<b>1.2970E+0</b>	1.5427E+0	4.6604E+0
	30	<b>2.6476E+0</b>	2.8931E+0	3.7200E+0	1.5060E+1		30	<b>4.3627E+0</b>	5.3324E+0	6.6650E+0	6.0207E+1
F7	2	5.6405E-2	<b>5.5008E-2</b>	5.7562E+0	6.5047E+1	F17	2	6.5516E-2	7.2802E-2	<b>5.8908E-2</b>	6.8641E-2
	10	<b>6.8407E+1</b>	6.8467E+1	8.7470E+1	1.8612E+3		10	<b>1.8298E-1</b>	1.9186E-1	1.9132E-1	2.3036E-1
	30	5.1245E+2	<b>4.4561E+2</b>	5.4400E+2	1.7949E+4		30	<b>4.1992E-1</b>	5.7923E-1	4.9740E-1	1.8528E+0
F8	2	<b>6.4421E-2</b>	8.2598E+0	1.0246E+1	7.5015E+1	F18	2	7.9171E+0	8.8528E+0	9.5477E+0	5.0562E+1
	10	7.3016E+1	<b>6.8773E+1</b>	9.2672E+1	2.8420E+3		10	<b>3.8751E+2</b>	3.9933E+2	6.9504E+2	1.9780E+4
	30	5.8338E+2	<b>4.7055E+2</b>	7.4258E+2	2.1170E+4		30	6.8663E+3	<b>5.9223E+3</b>	9.8100E+3	3.7050E+5
F9	2	<b>1.2270E+6</b>	1.4303E+6	1.5548E+6	2.0123E+7	F19	2	2.9912E-1	2.8090E-1	<b>2.6204E-1</b>	8.0894E-1
	10	<b>4.1771E+5</b>	4.7615E+5	5.7221E+5	8.1209E+6		10	<b>7.5440E+0</b>	7.6484E-1	8.2718E-1	2.5802E+0
	30	3.9609E+6	<b>3.8603E+6</b>	5.1305E+6	1.6692E+8		30	1.37388E+1	<b>1.2999E+1</b>	1.5938E+1	2.7295E+2
F10	2	9.5085E+0	<b>6.9783E+0</b>	1.2881E+1	6.7627E+1	F20	2	3.5455E-2	3.9468E-2	<b>3.3128E-2</b>	3.7422E-2
	10	7.8455E+1	<b>7.6291E+1</b>	8.3358E+1	2.9098E+3		10	4.0848E-2	4.3134E-2	<b>3.3434E-2</b>	5.6918E-2
	30	5.7160E+2	<b>4.3714E+2</b>	7.0516E+2	2.0198E+4		30	<b>5.6556E-2</b>	7.5748E-2	5.7657E-2	1.8357E-1

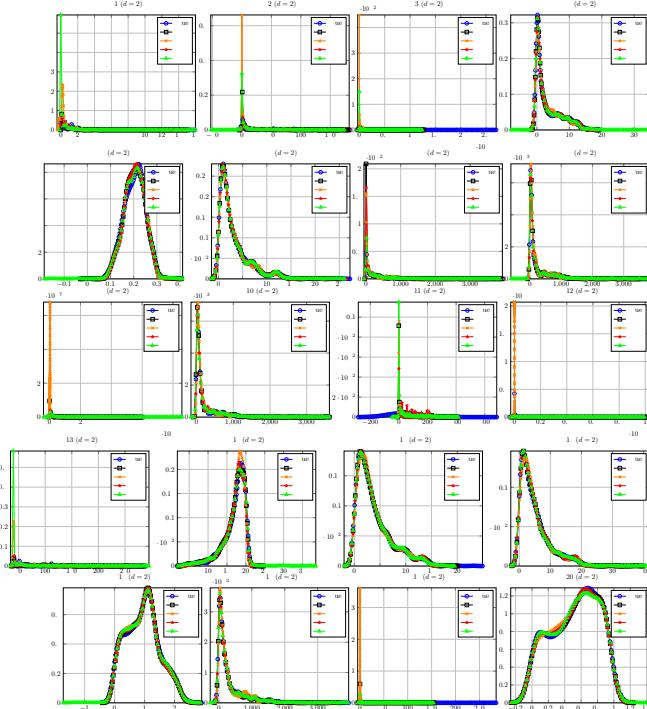


Fig. 4: Estimated probability density distribution of the empirical performance predicted by four different regression algorithms and the ground truth ( $d = 2$ ).

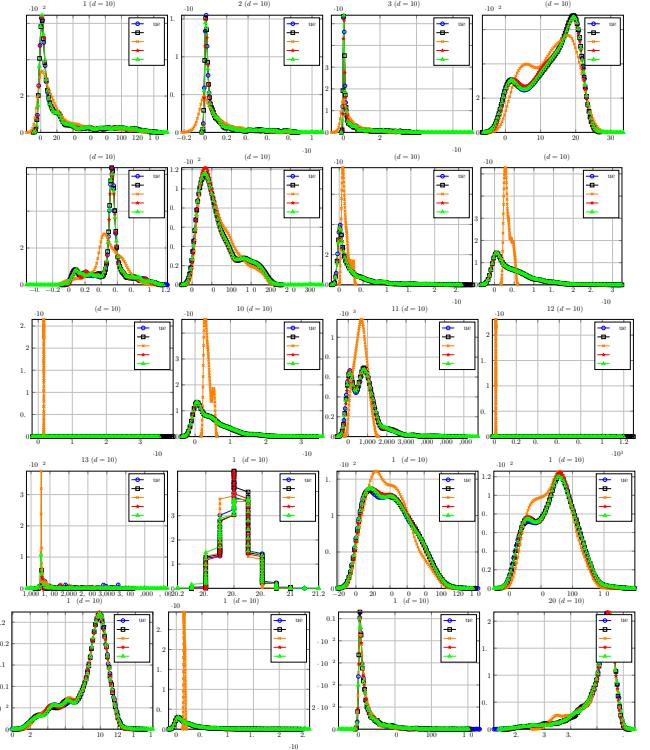


Fig. 5: Estimated probability density distribution of the empirical performance predicted by four different regression algorithms and the ground truth ( $d = 10$ ).

Instead of developing a new algorithm for PO, this paper aims to study a fundamental issue — investigating the ability of four prevalent regression algorithms for building a surrogate model of empirical performance. From our extensive experiments, we find that surrogate models built by GP and RF have shown promising generalisation ability for predicting the empirical performance of unseen parameter configurations. In particular, the prediction accuracy depends on the quality of the original

performance data. This implies that it needs to be careful to use a surrogate model in the early stage of a PO process. Furthermore, we find that although SVR does not show a promising performance for predicting the approximation error of a parameter configuration, it is able to differentiate the order of two parameter configurations.

Generally speaking, we hope this work will be useful to

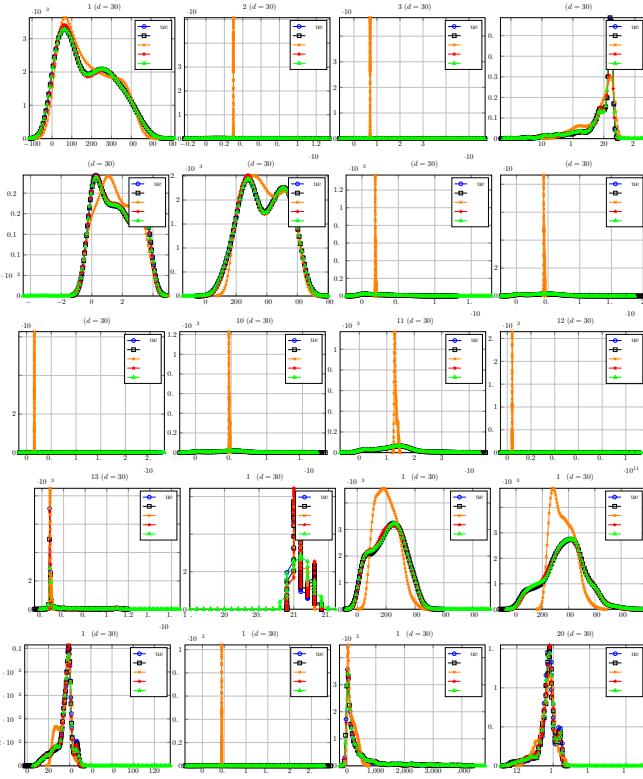


Fig. 6: Estimated probability density distribution of the empirical performance predicted by four different regression algorithms and the ground truth ( $d = 30$ ).

a wide variety of researchers who seek to model algorithm performance for algorithm analysis, scheduling, algorithm portfolio construction, automated algorithm configuration, and other applications. As for the coming next step, we plan to explore the following three aspects.

- We would like to apply the regression algorithms investigated in this paper in the context of model-based PO. Although using design and analysis of computer experiments in the context of PO has already been studied in some previous work (e.g. sequential PO [3]), it is still worthwhile to see whether the observations in the offline training are directly applicable to online PO.
- Since collecting a performance data in PO is computationally expensive, it might be interesting to use the offline trained surrogate models to generate pseudo data. In this rigour, semi-supervised learning can be useful to address a small data challenge.
- Here we set the PO as a per-instance scenario. In the prevalent algorithm configuration literature [4], it is more interesting to combine the problem feature into the surrogate modelling process so that we can generalise the PO to a range of similar problems.

## REFERENCES

- [1] Y. Jin, “Surrogate-assisted evolutionary computation: Recent advances and future challenges,” *Swarm and Evol. Comput.*, vol. 1, no. 2, pp. 61–70, 2011.
- [2] T. J. Santner, B. J. Williams, and W. I. Notz, *The Design and Analysis of Computer Experiments*. Springer, 2003.
- [3] T. Bartz-Beielstein, C. Lasarczyk, and M. Preuss, “Sequential parameter optimization,” in *CEC'05: Proc. of the 2005 IEEE Congress on Evol. Comput.*, 2005, pp. 773–780.
- [4] F. Hutter, H. H. Hoos, and K. Leyton-Brown, “Sequential model-based optimization for general algorithm configuration,” in *LION'11: Proc. of 5th International Conference on Learning and Intelligent Optimization*, 2011, pp. 507–523.
- [5] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown, “Auto-WEKA: combined selection and hyperparameter optimization of classification algorithms,” in *KDD'13: Proc. of 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2013, pp. 847–855.
- [6] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *J. Machine Learning Research*, vol. 13, pp. 281–305, 2012.
- [7] M. Reif, F. Shafait, M. Goldstein, T. Breuel, and A. Dengel, “Automatic classifier selection for non-experts,” *Pattern Anal. Appl.*, vol. 17, no. 1, pp. 83–96, 2014.
- [8] R. Kohavi and G. H. John, “Automatic parameter selection by minimizing estimated error,” in *ICML'95: Proc. of 12th International Conference on Machine Learning*, 1995, pp. 304–312.
- [9] A. Blot, H. H. Hoos, L. Jourdan, M. Kessaci-Marmion, and H. Trautmann, “Mo-paramil: A multi-objective automatic algorithm configuration framework,” in *LION'16: Proc. of 10th International Conference on Learning and Intelligent Optimization*, 2016, pp. 32–47.
- [10] M. López-Ibáñez, J. Dubois-Lacoste, L. Pérez Cáceres, T. Stützle, and M. Birattari, “The irace package: Iterated racing for automatic algorithm configuration,” *Oper. Res. Perspectives*, vol. 3, pp. 43–58, 2016.
- [11] J. Snoek, H. Larochelle, and R. P. Adams, “Practical bayesian optimization of machine learning algorithms,” in *NIPS'12: Proc. of 26th Annual Conference on Neural Information Processing Systems*, 2012, pp. 2960–2968.
- [12] S. Sanders and C. G. Giraud-Carrier, “Informing the use of hyperparameter optimization through metalearning,” in *ICDM'17: Proc. of 2017 IEEE International Conference on Data Mining*, 2017, pp. 1051–1056.
- [13] “Neurips 2018 challenge: The 3rd automl challenge: Automl for lifelong machine learning,” <https://www.4paradigm.com/competition/nips2018>.
- [14] R. Storn and K. V. Price, “Differential evolution - A simple and efficient heuristic for global optimization over continuous spaces,” *J. Global Optimization*, vol. 11, no. 4, pp. 341–359, 1997.
- [15] P. N. Suganthan, N. Hansen, K. Deb, J. J. Liang, Y.-P. Chen, A. Anger, and S. Tiwari, “Problem definitions and evaluation criteria for the CEC 2005 special session on real-parameter optimization,” NTU and IIT Kanpur, Technical Report 2005005, 2005.
- [16] S. Das and P. N. Suganthan, “Differential evolution: A survey of the state-of-the-art,” *IEEE Trans. Evol. Comput.*, vol. 15, no. 1, pp. 4–31, 2011.
- [17] J. Brest, S. Greiner, B. Boskovic, M. Mernik, and V. Zumer, “Self-adapting control parameters in differential evolution: A comparative study on numerical benchmark problems,” *IEEE Trans. Evol. Comput.*, vol. 10, no. 6, pp. 646–657, 2006.
- [18] A. K. Qin, V. L. Huang, and P. N. Suganthan, “Differential evolution algorithm with strategy adaptation for global numerical optimization,” *IEEE Trans. Evol. Comput.*, vol. 13, no. 2, pp. 398–417, 2009.
- [19] K. Li, Á. Fialho, and S. Kwong, “Multi-objective differential evolution with adaptive control of parameters and operators,” in *LION'11: Proc. of 5th International Conference on Learning and Intelligent Optimization*, 2011, pp. 473–487.
- [20] N. Belkhir, J. Dréo, P. Savéant, and M. Schoenauer, “Feature based algorithm configuration: A case study with differential evolution,” in *PPSN'16: Proc. of 14th International Conference on Parallel Problem Solving from Nature - PPSN XIV*, 2016, pp. 156–166.
- [21] F. Hutter, L. Xu, H. H. Hoos, and K. Leyton-Brown, “Algorithm runtime prediction: Methods & evaluation,” *Artif. Intell.*, vol. 206, pp. 79–111, 2014.
- [22] I. Loshchilov, M. Schoenauer, and M. Sebag, “Comparison-based optimizers need comparison-based surrogates,” in *PPSN'10: Proc. of 11th International Conference on Parallel Problem Solving from Nature*, 2010, pp. 364–373.
- [23] Y. Rubner, C. Tomasi, and L. J. Guibas, “The earth mover's distance as a metric for image retrieval,” *International Journal of Computer Vision*, vol. 40, no. 2, pp. 99–121, 2000.