

Model Selection based on Regional Error Estimation of Surrogates

Ali Mehmani¹, Souma Chowdhury², Jie Zhang³, Weiyang Tong⁴, Achille Messac⁵

¹ Syracuse University, Syracuse, NY, USA, amehmani@syr.edu

² Syracuse University, Syracuse, NY, USA, sochowdh@syr.edu

³ Syracuse University, Syracuse, NY, USA, jzhang56@syr.edu

⁴ Syracuse University, Syracuse, NY, USA, wtong@syr.edu

⁵ Syracuse University, Syracuse, NY, USA, messac@syr.edu

1. Abstract

The analysis of complex system behavior often demands expensive experiments or computational simulations. Surrogate modeling techniques are often used to provide a tractable and inexpensive approximation of such complex system behavior. Owing to the lack of any general guidelines regarding the suitability of different surrogate models for different applications, model selection approach can be helpful to choose the best surrogate technique. This paper investigates the effectiveness of a recently developed method for surrogate error quantification called, Regional Error Estimation of Surrogate (REES), to select the best surrogate model based on the level of accuracy. The REES method is developed based on the concept that the accuracy of the approximation methods is related to the amount of available resources. In the REES method, intermediate surrogates are iteratively constructed over heuristic subsets of the available sample points (i.e., intermediate training points) and tested over the remaining available sample points (i.e., intermediate test points). The statistical mode of the median and the maximum error distributions are selected to represent the overall and maximum error at each iteration. The estimated modes of the median and maximum error distributions are then represented as functions of the number of intermediate training points using a regression model. The regression models are used to predict the overall and minimum accuracy of the final surrogate. These two error measures are then applied to select the best surrogate. The proposed model selection technique is applied to select the best surrogate among (i) Kriging, (ii) Radial Basis Functions (RBF), (iii) Extended Radial basis Functions (E-RBF), and (iv) Quadratic Response Surface (QRS), for standard test functions and a wind farm power generation function. The REES-based model selection is compared with (i) model-selection based on cross-validation errors and (ii) model-selection based on error estimated on a large set of additional test points; the latter is assumed to provide the correct model selection. The REES-based model selection is found to be significantly more accurate than that based on cross-validation errors.

2. Keywords: Surrogate model; Error quantification; Model selection; Sampling; Decision making;

3. Introduction

Engineering design problems often involve computationally intensive simulation models (high fidelity models) or expensive experiment-based system evaluations. An accurate surrogate model is an effective tool for providing a tractable and an inexpensive approximation of the actual system evaluation. Major surrogate modeling methods include Polynomial Response Surfaces [1], Kriging [2, 3], Moving Least Square [4, 5], Radial Basis Functions (RBF) [6], Neural Networks [7], and hybrid surrogate modeling [8]. These methods have been applied to a wide range of disciplines, such as aerospace design, automotive design, chemistry, and material science [9]. Since there is not unique surrogates suitable for all applications, model selection techniques are used to select the best surrogate based on one or more error measures among available candidate surrogates. In addition, the error measures can be applied as a model parameter selection method (e.g., to find the best value for the shape parameter in RBF). Popular error measures used as model selection criteria include [10]: (i) split sample, (ii) cross-validation, (iii) bootstrapping, and (iv) Akaike's information criterion (AIC) [11]. These error measures either provide limited information regarding the accuracy of surrogates, or require additional system evaluations.

In this paper, we investigate the effectiveness of the REES error measurement method [12] in quantifying the predictive ability of the surrogates. This method is model independent and seeks to be universal in

application. While in this paper we are seeking to apply the REES method as a model selection technique, this method can also be applied as a model parameters selection method. In the following section, the formulation of the REES error measurement method as a model selection approach is introduced in detail.

4. Model Selection based on REES

As a model selection criterion, the REES method is applied to provide information about the overall accuracy of an estimated function without investing additional system evaluations. The REES method predicts the error by modeling the variation of the error with an increasing density of training points. The REES method is detailed in algorithm 1.

Algorithm 1 Model Selection based on REES

Suppose there are J candidates (or surrogates); indexed by j

INPUT:

Set Number of sample points N

Set Number of iterations N^{it} ; indexed by t

Set Size of intermediate training points at each iteration, ${}^t n$, where ${}^t n < {}^{t+1} n$

Set Number of combinations at each iteration, K^t where $K^t \leq \binom{N}{n^t}$; indexed by k

for all candidates, $j = 1, \dots, J$ **do**

X = Experimental Design(N)

$F | X$ = Evaluate System (X)

$\{X\} = \{(X_i, F_i)_{i=1}^N\}$

for $t = 1, \dots, N^{it}$ **do**

for $k = 1, \dots, K^t$ **do**

Choose $\{\beta\} \subset \{X\}$, where $\#\{\beta\} = {}^t n$

Define intermediate training points, $\{X^{TR}\} = \{\beta\}$

Define intermediate test points, $\{X^{TE}\} = \{X\} - \{\beta\}$

Construct intermediate surrogate S_k using $\{X^{TR}\}$

Estimate median and maximum errors;

$E_{med,t}^k = median(e_m)_{m=1, \dots, \#\{X^{TE}\}}$, $E_{max,t}^k = max(e_m)_{m=1, \dots, \#\{X^{TE}\}}$

end for

Fit distributions of the median and the maximum errors over all K^t combinations

Determine the mode of the median and maximum error distributions; Mo_{med}^t and Mo_{max}^t

end for

Construct a final surrogate using $\{X\}$

Train regression functions using Mo_{med}^t and $Mo_{max}^t \quad \forall t$

Predict the overall and maximum errors in the final surrogate (using regression models); $E_j^{REES_{med}}$

and $E_j^{REES_{max}}$

end for

return Best surrogate with the smallest $E_j^{REES_{med}}$ and/or $E_j^{REES_{max}}$

In this method, for each surrogate candidate, the intermediate surrogates are constructed iteratively for all combinations of each iteration using the intermediate training points, and are tested over the intermediate test points. The median and maximum errors are then estimated for each combination. The median is applied since it is a useful measure of central tenancy which is less vulnerable to outliers.

$$E_{med}^k = median(e_m)_{m=1, \dots, \#\{X^{TE}\}}$$

$$E_{max}^k = max(e_m)_{m=1, \dots, \#\{X^{TE}\}}$$

where $\#\{X^{TE}\}$ represents the number of test points at t^{th} iteration; and e represents the Relative

Absolute Error value estimated on intermediate test points.

$$e_m = RAE_m; RAE_m = \left| \frac{y_m - \hat{y}_m}{y_m} \right| \quad (1)$$

where y_m and \hat{y}_m are the actual and predicted values on m^{th} intermediate test point, respectively.

The median and maximum errors estimated over all K^t combinations are used to develop a probabilistic model at each iterations. In this study, the lognormal distribution is selected because there is often orders of magnitude variation in the median and maximum errors. The density function of a lognormal distribution is given by

$$f(x; \mu, \sigma) = \frac{1}{x\sqrt{2\pi\sigma^2}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}} \quad (2)$$

where the σ^2 and μ represents the shape and log-scale parameters of the distribution, respectively. The regression models are then applied to relate the statistical mode of the median and maximum error distributions (Mo_{med}^t and Mo_{max}^t) to the number of the training points ($^t n$) at each iteration. These regression models called the *variation of error with sample density (VESD)*. This models are then used to predict the overall and maximum errors of each surrogate candidate, and to select the best surrogate.

In this study, three types of the regression functions are applied

Type 1 Exponential regression model

$$F(x) = a_0 e^{a_1 x} \quad (3)$$

Type 2 Multiplicative regression model

$$F(x) = a_0 x^{a_1} \quad (4)$$

Type 3 Linear regression model

$$F(x) = a_0 + a_1 x \quad (5)$$

where a_0 and a_1 are unknown coefficients to be determined. The choice of these functions assume a *smooth monotonic decrease of the error with the training point density within that region*. In this paper, the root mean squared error metric is used to select the best-fit regression model.

5. Application of the Model Selection based on REES Method

The effectiveness of the model selection based on the REES method is explored to select the best surrogate between all candidates including (i) Kriging, (ii) RBF, (iii) E-RBF, and (iv) Quadratic Response Surface (QRS) on four benchmark problems and an engineering design problem. The results of the REES method in selecting the best surrogate are compared with the model selection based on actual errors evaluated using additional test points. In algorithm 2, the formulation of evaluating actual errors and using it as a model selection criterion is given in detail. Here, the actual error is defined by the mode of error distribution trained using errors evaluated on additional test points.

To illustrate the potential greater effectiveness of model selection based on the REES error measurement over a prediction sum of square (PRESS), the *normalized* PRESS method based on the *leave-one-out cross-validation* approach is defined in Algorithm 3.

To do a sensible comparison, the *normalized* Root Mean Square Error (RMSE) evaluated on additional test points is used to evaluate the performance of PRESS in selecting the best surrogate among other surrogate candidates.

$$E^{\text{ActualRMSE}} = \sqrt{\frac{1}{N^{\text{Test}}} \sum_{i=1}^{N^{\text{Test}}} (RAE_i)^2} \quad (6)$$

5.1 Benchmark Problems

The performance of the model selection based on REES method is evaluated using the following analytical test problems:

Branin-Hoo function (2 variables)

$$f(x) = \left(x_2 - \frac{5.1x_1^2}{4\pi^2} + \frac{5x_1}{\pi} - 6 \right)^2 + 10 \left(1 - \frac{1}{8\pi} \right) \cos(x_1) + 10$$

where $x_1 \in [-5 \ 10]$, $x_2 \in [0 \ 15]$

Algorithm 2 Model Selection based on the actual error on additional test points

Suppose there are J candidates (or surrogates); indexed by j

INPUT:

Set Number of additional test points N^{test}

for all candidates, $j = 1, \dots, J$ **do**

for $i = 1, \dots, N^{test}$ **do**

 Estimate actual value on i^{th} test point; $y_i = \text{System}(x_i)$

 Estimate predicted value on i^{th} test point; $\hat{y}_i = \text{Surrogate}(x_i)$

 Estimate RAE on i^{th} test point; $RAE_i = \left| \frac{y_i - \hat{y}_i}{y_i} \right|$

end for

end for

Fit a distribution to the errors (RAEs) evaluated on test points

Determine the mode of the error distribution; $E_j^{Actual} = \text{mode}(\text{error distribution})$

return Best surrogate surrogate with the smallest E_j^{Actual}

Algorithm 3 Model Selection based on the normalized PRESS

Suppose there are J candidates (or surrogates); indexed by j

INPUT:

Set Number of training points N

for all candidates, $j = 1, \dots, J$ **do**

for $i = 1, \dots, N$ **do**

 Estimate predicted value on i^{th} point using a final surrogate; $\hat{y}_i = \text{Surrogate}(x_i)$

 Estimate predicted value on i^{th} point using an intermediate surrogate*; $\hat{y}_i^{-i} = \text{Intermediate Surrogate}(x_i)$

 Estimate RAE on i^{th} point; $RAE_i^{CV} = \left| \frac{\hat{y}_i - \hat{y}_i^{-i}}{\hat{y}_i} \right|$

end for

end for

Estimate the root mean square of the errors (RAEs), $E_j^{PRESS} = \sqrt{\frac{1}{N} \sum_{i=1}^N (RAE_i^{CV})^2}$

return Best surrogate with the lowest E_j^{PRESS}

**Intermediate surrogate* is the surrogate constructed using all sample points except the i^{th} point.

Hartmann function (6 design variables)

$$f(x) = -\sum_{i=1}^4 c_i \exp\left\{-\sum_{j=1}^n A_{ij} (x_j - P_{ij})^2\right\} \quad (7)$$

where $x = (x_1 \ x_2 \ \dots \ x_n)$, $x_i \in [0 \ 1]$

In Hartmann-6, $n = 6$; the constants c , A , and P , are a 1×4 vector, a 4×6 matrix, and a 4×6 matrix, respectively [13].

Dixon & Price functions (12 and 18 design variables))

$$f(x) = (x_1 - 1)^2 + \sum_{i=2}^n i (2x_i^2 - x_{i-1})^2 \quad (8)$$

where $x_i \in [-10 \ 10]$, $i = 1, \dots, n$
 $n = 12$ and 18 .

5.2 Wind Farm Power Generation

The effectiveness of model selection based on REES for engineering problems is illustrated using a wind farm power generation model, adopted from the Unrestricted Wind Farm Layout Optimization (UWFLO) framework [14, 15]. Surrogate candidates are developed using Kriging, RBF, E-RBF, and QRS to represent the *power generation of an array-like wind farm*. It is assumed that turbines are arranged in a row-column pattern over the farm site. Hence, the wind farm power generation can be represented as a function of the *streamwise spacing* and *spanwise spacing* between turbines, with respect to the south direction (as shown in Fig. 1). The annual average power generation of a wind farm is a complex and expensive function of the turbine features, the turbine arrangement (or farm layout), and the local wind resource variations. A surrogate model offers a more tractable (and inexpensive) representation of the farm power generation in terms of key design parameters. To train the surrogate model, the actual annual-average wind farm power generation is estimated using an advanced power generation model developed by Chowdhury et al. [14, 15]. The selection of the surrogate model among available surrogate candidates in wind farm design and analysis is critical since further decision highly depends on the surrogate selection.

In this case study, surrogates are constructed to represent the power generation of an array-like *100-turbine wind farm* as a function of the *streamwise spacing* (x_h) and the *spanwise spacing* (x_l) between turbines. The turbines are arranged in a 10×10 patterns in this case. The bivariate normal distributions of wind data obtained for a site in North Dakota[16] is used for this case study. The lower and upper bounds of x_h and x_l , based on the wind turbine rotor diameter (D), are specified as

$$5D < x_h < 30D$$

$$1.1D < x_l < 10D$$

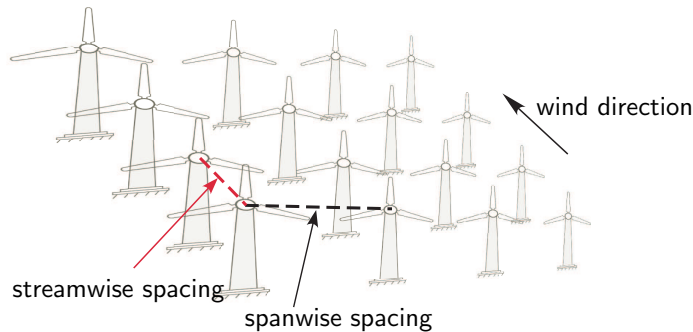


Figure 1: Wind farm array schematic

5.2 Numerical Settings

Numerical settings for the application of REES are provided in Table 1, which lists (i) the number

of variables, (ii) the number of training points, (iii) the number of iterations, and (iv) the size of training points at each iteration (as a function of iteration, t). The Optimal Latin Hypercube based on Translational Propagation algorithm [17] is adopted to determine locations of the full set of sample points $\{X\}$, and additional test points for the benchmark problems. This design of experiments method is accompanied by the modified maximum distance criterion.

Table 1: Numerical setup for test problems

Function	No. of variables	No. of sample points, N	No. of iteration, N^{it}	No. of training points at each iteration, t_n
Branin-Hoo function	2	32	4	$22+2t$
Hartmann	6	100	5	$88+2t$
Dixon & Price	12	240	5	$228+2t$
Dixon & Price	18	400	6	$386+2t$
Wind farm power generation	2	30	4	$19+2t$

To implement the Kriging method, the DACE (design and analysis of computer experiments) package developed by Lophaven et al. [18] is used. The bounds on the correlation parameters in the nonlinear optimization in Kriging, θ_l and θ_u are specified to be 0.1 and 20, respectively. The zero-order polynomial function is used as a regression model. To implement RBF, the multiquadric radial basis function [6] is used where the shape parameter is set to $c = 0.9$. In implementation of E-RBF [19], the shape parameter is set to $c = 0.9$; the λ parameter is set to 4.75; and the order of monomial in non-radial basis functions is fixed at 2.

6 Results and Discussion

In this study, the best surrogate is the most accurate model (based on the overall and/or minimum accuracy) among all candidates over the entire design domain. To do model selection based on the REES, the variation of error with sample density (VESD) regression models for all surrogate candidates are constructed for all test problems. The VESD regression models used to predict the overall and maximum errors in different surrogates for benchmark test problems are illustrated in Figs. 2 and 3, respectively. These regression models for wind farm power generation problem are illustrated in Fig. 4. The VESD models are trained using the statistical mode of the median and maximum error distributions (Mo_{med} and Mo_{max}) at each iteration in different surrogates (as explained in Algorithm 1). In Figs. 2 and 4(a) solid circles represent quantified mode of median errors at each iteration in different surrogates; squares represent predicted overall error in different surrogates. Similarly, solid circles and squares in Figs. 3 and 4(b) represent quantified mode of maximum errors at each iteration and the predicted maximum error in different surrogates.

Table 2 shows the results of the REES method in selecting the best surrogate among all candidates for different test problems. Here, the overall error estimated using REES is applied to select the best surrogate. In this table, the result of model selection based on the RESS is also compared with those models selected based on the actual error evaluated on additional test points (explained in Algorithm 2). It is observed that REES selects the best surrogate correctly. This observation shows that REES can be effectively used as a model selection method to select the best surrogate for different application domains.

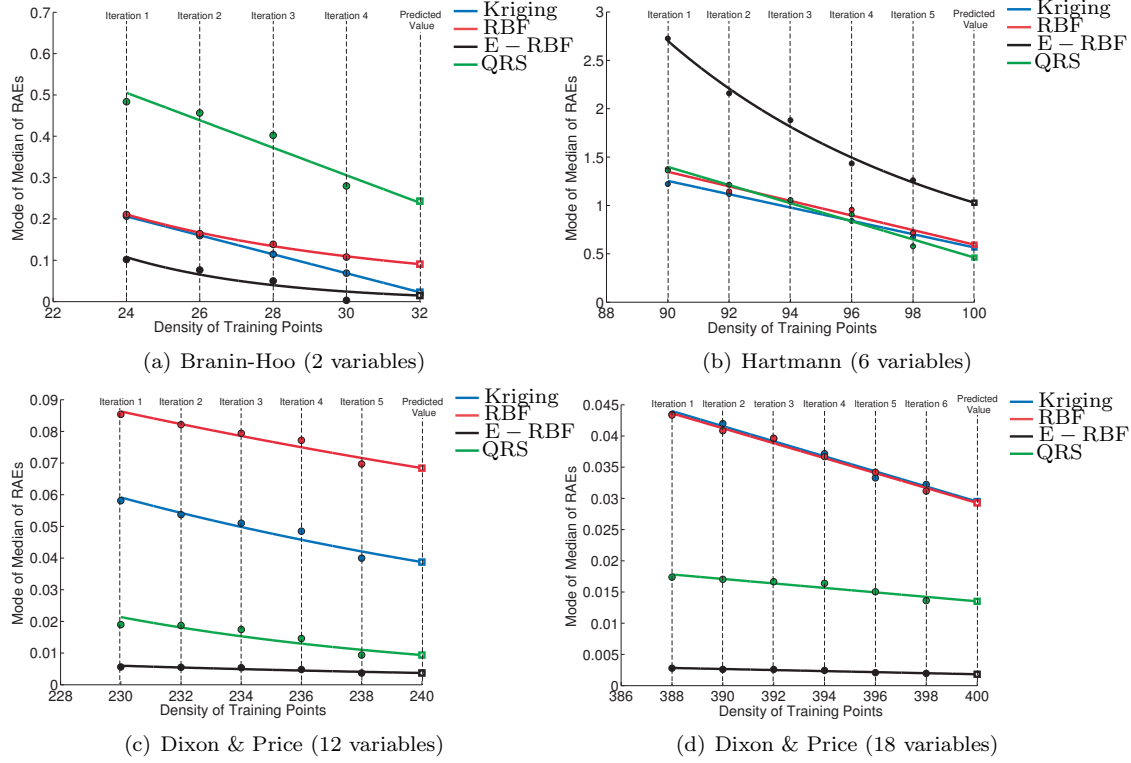


Figure 2: VESD regression models used to predict the overall error

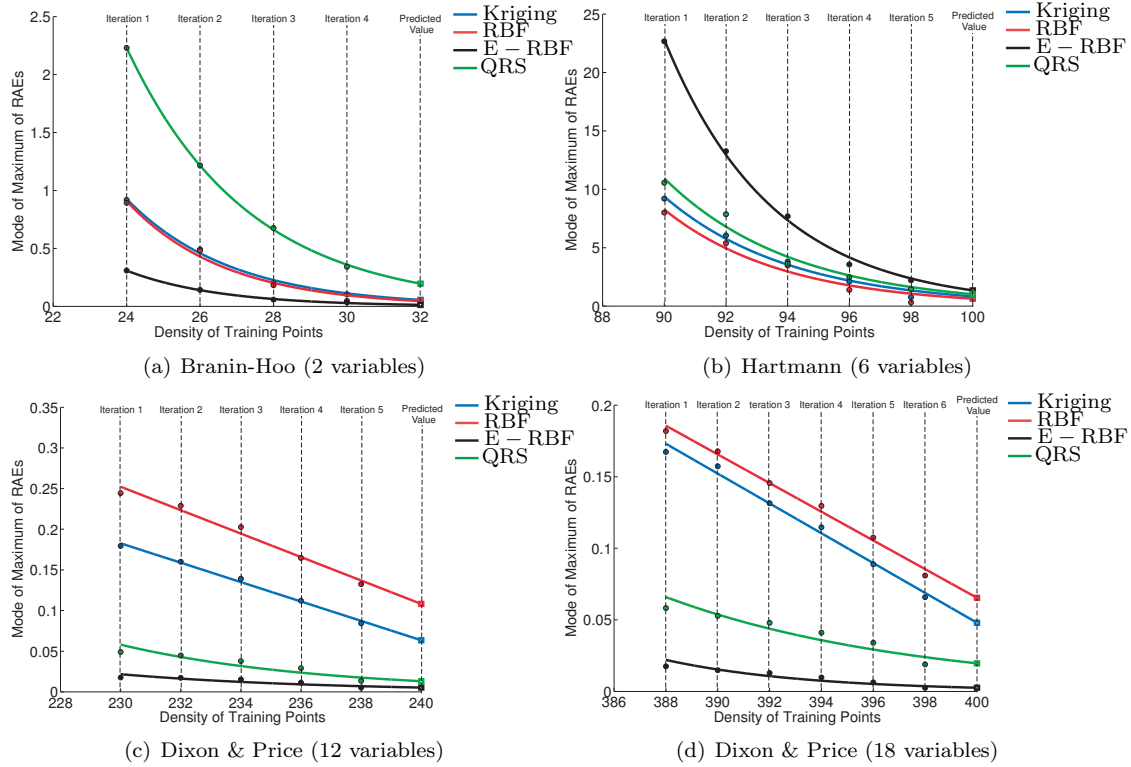


Figure 3: VESD regression models used to predict the maximum error

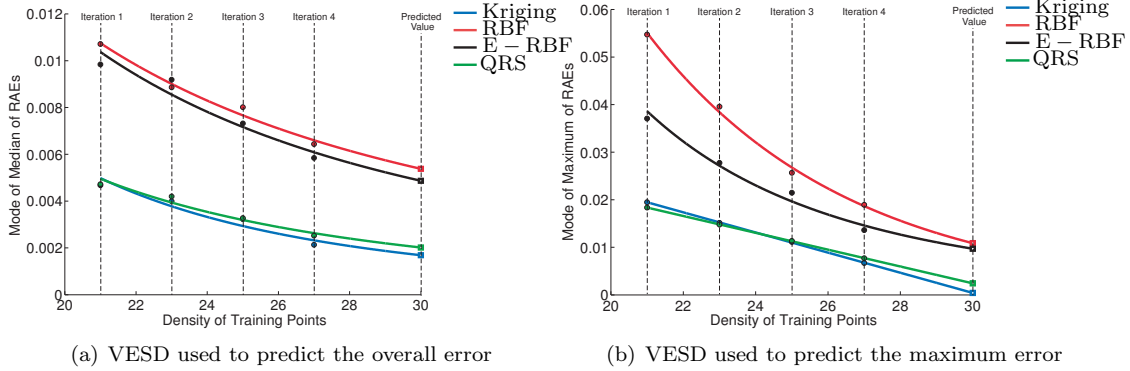


Figure 4: Wind farm power generation

Table 2: Model selection based on the REES and mode of the actual errors evaluated on additional test points (E^{Actual})

Function	REES				E^{Actual}			
	Kriging	RBF	E-RBF	QRS	Kriging	RBF	E-RBF	QRS
Branin-Hoo			X				X	
Hartmann-6				X				X
Dixon & Price-12			X				X	
Dixon & Price-18			X				X	
Wind farm power generation	X				X			

In Table 3, the results of model selection based on the *normalized* PRESS are compared with those based on the mean of the *normalized* mean square error evaluated using a large number of test points. It is observed that in most of the cases, the normalized PRESS could not select the best surrogate correctly. Owing to the relative accuracy of the intermediate surrogates illustrated in Fig. 2 for all test problems; the model selection based on the error evaluated in the last iteration should have the same output as the model selection based on the REES. The issue with model selection based on PRESS is attributed to the use of the *mean* to aggregate the errors evaluated using the *leave-one-out* approach, which makes PRESS vulnerable to the outliers.

Table 3: Model selection based on the *normalize* PRESS and Root Mean Square Errors on additional test points ($E^{Actual_{RMSE}}$)

Function	Normalized PRESS				$E^{Actual_{RMSE}}$			
	Kriging	RBF	E-RBF	QRS	Kriging	RBF	E-RBF	QRS
Branin-Hoo			X			X		
Hartmann-6	X					X		
Dixon & Price-12				X			X	
Dixon & Price-18				X			X	
Wind farm power generation				X				X

In some applications (e.g., evaluating the maximum stress in structural analysis), we need to select the best surrogate based on the maximum error evaluated in the entire domain. To do this end, the maximum error estimated using REES ($E^{REES_{max}}$) is applied. Table 3 represents the results of model selection based on the maximum error evaluated using REES; model selection based on the 75th percentile of the *leave one out cross validation* errors evaluated by considering all combinations; model selection based on 75th percentile of the actual errors evaluated on a large number of test points. It is observed that, in the majority of cases, the REES method selects the best surrogate based on the maximum error correctly.

Table 4: Model selection based on the maximum error evaluated using REES, 75th percentile of the actual errors evaluated on test points, and 75th percentile of the *leave one out cross validation* (CV) errors

Function	Maximum error based on REES				75 th percentile of the CV errors				75 th percentile of the actual errors			
	Kriging	RBF	E-RBF	QRS	Kriging	RBF	E-RBF	QRS	Kriging	RBF	E-RBF	QRS
Branin-Hoo			X					X			X	
Hartmann-6		X						X	X			
Dixon & Price-12			X					X			X	
Dixon & Price-18			X					X			X	
Wind farm power generation	X							X	X			

7 Conclusion

This paper presents a new model selection approach to select the best surrogate among available surrogate models based on the level of accuracy. This approach investigates the effectiveness of a recently developed error measure, Regional Error Estimation of Surrogate (REES). In the REES method, intermediate surrogates are iteratively constructed with heuristic subsets of the available sample points, and remaining points are used to evaluate the error of the estimated function. Regression models are then used to represent the surrogate error as a function of the number of training points. This regression model is used to predict the overall accuracy of the final surrogates, and to select the best surrogate model. The effectiveness of the model selection based on the REES are illustrated using standard test problems and a wind farm power generation problem. The results show that the proposed method selects the best surrogate among all surrogate candidates with a higher level of confidence in comparison to the *normalized* prediction sum of square (PRESS) estimated by the *leave-one-out cross-validation* approach.

References

- [1] R. Jin, W. Chen, and T. W. Simpson. Comparative studies of metamodeling techniques under multiple modeling criteria. *AIAA*, (4801), 2000.
- [2] T. Simpson, J. Korte, T. Mauery, and F. Mistree. Kriging models for global approximation in simulation-based multidisciplinary design optimization. *AIAA Journal*, 39(12):2233–2241, 2001.
- [3] A. Forrester and A. Keane. Recent advances in surrogate-based optimization. *Progress in Aerospace Sciences*, 45(1-3):50–79, 2009.
- [4] K.K. Choi, B.D. Young, and R.J. Yang. Moving least square method for reliability-based design optimization. pages 4–8, Dalian, China, June 2001. 4th World Congress of Structural and Multidisciplinary Optimization.
- [5] V. V. Toropov, U. Schramm, A. Sahai, R. D. Jones, and T. Zeguer. Design optimization and stochastic analysis based on the moving least squares method. Rio de Janeiro, 30 May - 03 June 2005. 6th World Congresses of Structural and Multidisciplinary Optimization.
- [6] R. L. Hardy. Multiquadric equations of topography and other irregular surfaces. *Journal of Geophysical Research*, 76:1905–1915, 1971.
- [7] B. Yegnanarayana. *Artificial Neural Networks*. PHI Learning Pvt. Ltd., 2004.
- [8] J. Zhang, S. Chowdhury, and A. Messac. An adaptive hybrid surrogate model. *Structural and Multidisciplinary Optimization*, 46(2):223–238, 2012.
- [9] G. Wang and S. Shan. Review of metamodeling techniques in support of engineering design optimization. *Journal of Mechanical Design*, 129(4):370–381, 2007.
- [10] N. Queipo, R. Haftka, W. Shyy, T. Goel, R. Vaidyanathan, and P. Tucker. Surrogate-based analysis and optimization. *Progress in Aerospace Sciences*, 41(1):1–28, 2005.
- [11] H. Bozdogan. Akaike’s information criterion and recent developments in information complexity. *Journal of Mathematical Psychology*, 44:62–91, 2000.

- [12] A. Mehmani, S. Chowdhury, J. Zhang, W. Tong, and A. Messac. Quantifying regional error in surrogates by modeling its relationship with sample density. In *54th AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics and Materials Conference*, Boston, MA, USA, April 2013.
- [13] J. Zhang, S. Chowdhury, and A. Messac. An adaptive hybrid surrogate model. *Structural and Multidisciplinary Optimization (accepted)*, 46(2):223–238, August 2012.
- [14] S. Chowdhury, J. Zhang, A. Messac, and L. Castillo. Unrestricted wind farm layout optimization (uwflo): Investigating key factors influencing the maximum power generation. *Renewable Energy*, 38(1):16–30, 2012.
- [15] S. Chowdhury, J. Zhang, A. Messac, and L. Castillo. Optimizing the arrangement and the selection of turbines for a wind farm subject to varying wind conditions. *Renewable Energy*, 52:273–282, 2013.
- [16] NDAWN. The north dakota agricultural weather network, 2010. <http://ndawn.ndsu.nodak.edu/>.
- [17] F. A. C. Viana, G. Venter, and V. Balabanov. An algorithm for fast optimal latin hypercube design of experiments. *International Journal for Numerical Methods in Engineering*, 82(2):135156, 2010.
- [18] S. N. Lophaven, H. B. Nielsen, and J. Sondergaard. Dace - a matlab kriging toolbox, version 2.0. Technical Report IMM-REP-2002-12, Informatics and Mathematical Modelling Report, Technical University of Denmark, 2002.
- [19] A. Mullur and A. Messac. Extended radial basis functions: More flexible and effective metamodeling. *AIAA Journal*, 43(6):1306–1315, 2005.