

New Approach by Kriging Models to Problems in QSAR

Kai-Tai Fang,^{*,†,‡} Hong Yin,^{†,‡} and Yi-Zeng Liang[§]

Department of Mathematics, Hong Kong Baptist University, Hong Kong, China, College of Mathematics and Statistics, Wuhan University, Wuhan 430072, P. R. China, and College of Chemistry and Chemical Engineering, Central South University, Changsha 410083, P. R. China

Received June 23, 2004

Most models in quantitative structure and activity relationship (QSAR) research, proposed by various techniques such as ordinary least squares regression, principal components regression, partial least squares regression, and multivariate adaptive regression splines, involve a linear parametric part and a random error part. The random errors in those models are assumed to be independently identical distributed. However, the independence assumption is not reasonable in many cases. Some dependence among errors should be considered just like Kriging. It has been successfully used in computer experiments for modeling. The aim of this paper is to apply Kriging models to QSAR. Our experiments show that the Kriging models can significantly improve the performances of the models obtained by many existing methods.

1. INTRODUCTION

The molecular descriptors include various topological indices, quantum chemical descriptors, physicochemical parameters, and so on. They all give structure descriptions of chemical compounds. Chemometrics, especially, quantitative structure activity relationship (QSAR) and quantitative structure–property relationship (QSPR), attempt to correlate physical, chemical, and biological activities or properties with structural descriptors of compounds and find a suitable model, called metamodel, to establish relationships between molecule descriptors and activities or properties. The results are useful in theoretical and computational chemistry, biochemistry, pharmacology, and environment research.

The common approach for building a structure–activity or property relationship consists of the following steps: (1) to develop (or to select) the descriptors for the molecular structure; (2) to apply proper mathematical methods to construct a metamodel (a metamodel is an approximate model to the true model); and (3) to evaluate the model built. This study is concerned with the last two steps.

Techniques in multivariate analysis and data mining, such as ordinary least squares regression, principal components regression,¹ partial least squares regression,² multivariate adaptive regression splines,³ and multivariate additive regression tree,^{4,5} are useful tools for modeling. Metamodels generated by these methods, basically are linear models with independently identical distributed (iid) random errors, i.e.,

$$y_i = \sum_{j=1}^q \beta_j f_j(\mathbf{x}_i) + \epsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

where $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{ip}]$ denotes the i th sample which is described by p molecule descriptors, $\mathbf{y} = [y_1, y_2, \dots, y_n]'$ (y_i

$\in \mathcal{R}$) is a corresponding vector of n molecular activity or property values and ϵ_i 's are realizations of Gaussian independent random variables with mean zero and variance σ_ϵ^2 . The functions f_j 's are chosen from a given set of basis functions, such as polynomials basis, splines,^{6,7} wavelets,⁸ principal components, etc. The first part $\sum_{j=1}^q \beta_j f_j(\mathbf{x}_i)$ in model (1) is called *parametric item* and the second part is about *error*.

How to design a model (1)? Most of authors paid their attentions to the parametric item and put their work on building a good parametric item from given basis functions. However, in their study, they are faced with many difficulties: (1) *collinearity*: there are high correlations among variables; (2) *sparsity*: the number of observations is relatively small compared to the number of variables; and (3) *the curse of dimensionality*: many methods cannot be implemented due to computational complexity.

There are a lot of discussions on this direction in the literature. Some authors suggested to employ sliced inverse regression,^{9,10} projection pursuit regression,¹¹ lasso,¹² and so on to overcome the existing questions above. In this paper, we do not discuss how to optimize the *parametric item*. We may steer clear of the obstruction to look for another direction.

It seems possible for us to improve the performance of model (1) from the other aspect. The assumption of independent and identical distributed errors in model (1) is not always true. It might not be reasonable in many cases in chemometrics. For instance, many examples show that there can still be unacceptably large residuals compared to measurement errors^{13,14} in many models of QSAR/QSPR research. The reason for this may be diverse. The simplest and the most natural reflection on our mind is that the unaccepted residuals could be dependent, since a limited number of numeric molecular descriptors could not capture most of the information of molecular structures. So it becomes more reasonable to consider some kind of dependence of errors between any pairs of observations. These dependent errors will present more information than an independent

* Corresponding author phone: 852-3411-7025; fax: 852-3411-5811; e-mail: ktfang@hkbu.edu.hk.

[†] Hong Kong Baptist University.

[‡] Wuhan University.

[§] Center South University.

Table 1. OLSR, Kriging Model, and Their Prediction Results for Test Data

variable names	OLSR $\hat{\beta}$	Kriging $\hat{\beta}$	Kriging $\hat{\theta}$
constant	135.057	0.009	
<i>Chi</i> ₁	0.824	0.041	1.085
<i>Chi</i> ₂	47.804	1.109	0.500
<i>Chi</i> ₃	4.949	0.083	0.442
<i>Chi</i> ₄	-0.344	0.007	1.042
<i>Chi</i> ₅	-3.898	-0.073	0.346
<i>Chi</i> ₆	-8.763	-0.202	2.354
<i>Chi</i> ₇	-1.982	-0.051	9.809
<i>hari</i> ₁	-1.575	-0.069	3.007
<i>hidi</i> ₁	-3.066	-0.073	1.277
<i>hwi</i> ₁	-1.222	0.022	0.543
<i>idi</i> ₁	6.617	0.208	26.089
<i>infi</i> ₁	-8.675	0.024	0.500
<i>infi</i> ₂	-2.944	-0.088	0.240
<i>infi</i> ₃	-2.856	-0.320	0.850
<i>infi</i> ₄	7.223	0.212	0.055
test RMSE	6.363	5.811	8.68% (improved)
test AAE	4.439	3.686	16.96% (improved)

Table 2. Correlation Function References

name	$R_{\theta}^{(j)}(x_{ij} - x_{kj})$
EXP	$\exp(-\theta_j d_j)$
GUASS	$\exp(-\theta_j d_j^2)$
LIN	$\max\{0, 1 - \theta_j d_j \}$
SPLINE	$1 - 3\xi_j^2 + 2\xi_j^3; \xi_j = \min\{1, \theta_j d_j \}$

situation. For instance, we might use a Gaussian process $\{z(\mathbf{x}_i), i = 1, 2, \dots, n\}$ instead of independent random variables ϵ_i 's. In this case the underlying model would be

$$y(\mathbf{x}) = \sum_{j=1}^q \beta_j f_j(\mathbf{x}) + z(\mathbf{x}) \quad (2)$$

where $y(\mathbf{x})$ is the response at \mathbf{x} and $z(\mathbf{x})$ follows a Gaussian process with mean zero, $E(z(\mathbf{x})) = 0$, and covariance function

$$\text{Cov}[z(\mathbf{x}), z(\mathbf{x}^*)] = \sigma_z^2 R_{\theta}(\mathbf{x}; \mathbf{x}^*) \quad (3)$$

where σ_z^2 is called process variance and $R_{\theta}(\cdot, \cdot)$ is the correlation function of two samples with some unknown parameter vector $\theta = [\theta_1, \theta_2, \dots, \theta_p]$. When the process $\{z(\mathbf{x})\}$ is stationary, $R_{\theta}(\mathbf{x}, \mathbf{x}^*) = R_{\theta}(d(\mathbf{x}, \mathbf{x}^*))$ depends only on the distance, $d(\mathbf{x}, \mathbf{x}^*)$, between \mathbf{x} and \mathbf{x}^* . The larger the distance, the smaller the correlation between two samples. Many referenced correlation functions can be found in Table 2. Applying model (2) into the data set we have

$$y_i = \sum_{j=1}^q \beta_j f_j(\mathbf{x}_i) + z(\mathbf{x}_i), i = 1, 2, \dots, n \quad (4)$$

In fact, model (2) is just the Kriging model that has been widely used for the design and analysis of computer experiments.¹⁵ So the main objective of this paper is to introduce what we believe to be a new and promising approach for prediction in QSAR/QSPR research.

For easily understanding, now let us see an example. We consider the boiling points of 530 saturated hydrocarbons with 2–10 carbon atoms consulting from Rücker and Rücker.¹⁶ There are 15 topological indices on our hands: 7 *Chi*'s, 1 *hari*, 1 *hidi*, 1 *hwi*, 1 *idi*, and 4 *infi*'s. Generally,

ordinary least squares regression (OLSR) is the simplest and widely applied model in QSAR. We can first compare the prediction results of OLSR with those of Kriging. If input variables are not changed, the 15 topological indices can be looked on as a set of basic functions. The main difference between OLSR and Kriging is that the former postulates the independence of errors, while the latter prefer the dependence to independence of errors.

One group results can be displayed when the size of training data is $\mathcal{L} = 424$ sampled randomly from 530 saturated hydrocarbons and the number of test data is $\mathcal{T} = 106$. OLSR and its corresponding Kriging model sharing the same basic functions with OLSR are obtained based on 424 training data. The least squares estimators of $\beta_i, i = 0, \dots, 15$ under model (1) and model (4) are listed in Table 1. Figure 1 gives plots of fittings and residuals for 424 training data and 106 test data computed by OLSR and the corresponding Kriging model. We found that the fittings obtained by the Kriging model are almost equal to the actual observations which practically verifies that Kriging is an interpolator. For comparing prediction ability of these two models, root-mean-squared error (RMSE) and average absolute error (AAE) of test set are employed as criteria to evaluate model's accuracy. Their computational formulas are

$$RMSE = \sqrt{\frac{1}{t} \sum_{i=1}^t (y_i - \hat{y}_i)^2} \quad (5)$$

and

$$AAE = \frac{1}{t} \sum_{i=1}^t |y_i - \hat{y}_i| \quad (6)$$

respectively, where t is the number of test set and y_i is observation of the response, and \hat{y}_i is the corresponding prediction value. Our calculation shows that $RMSE(OLSR) = 6.363$, $RMSE(Kriging) = 5.811$; $AAE(OLSR) = 4.439$, and $AAE(Kriging) = 3.686$. For the sake of comparisons, we define the degree of improvement using 'ratio' = $(old - new)/old$ ('new' denotes the prediction results of Kriging and 'old' is those of other metamodels). So the improvement degrees of Kriging against OLSR for RMSE and AAE are 8.68% and 16.96%, respectively. It shows that the Kriging model can significantly improve its original model with independent errors.

The above experiment shows the usefulness of the Kriging model in QSAR. However, how to apply the Kriging model to data in QSAR, several issues need to be raised:

(1) Choice of Basic Functions. Many existing techniques in multivariate analysis and data mining can help us to choose a suitable model based on model (1), whose basic functions can be used as reference for model (2). So in this paper we consider only the ordinary least squares regression (OLSR), the principal components regression (PCR), and the partial least squares regression (PLSR) for generating metamodels. In fact, other techniques can also be applied to find more choices of basic functions for model (2). For our experience the improvement made by Kriging will be more significant if the parametric item in referenced metamodels is not far from the true model.

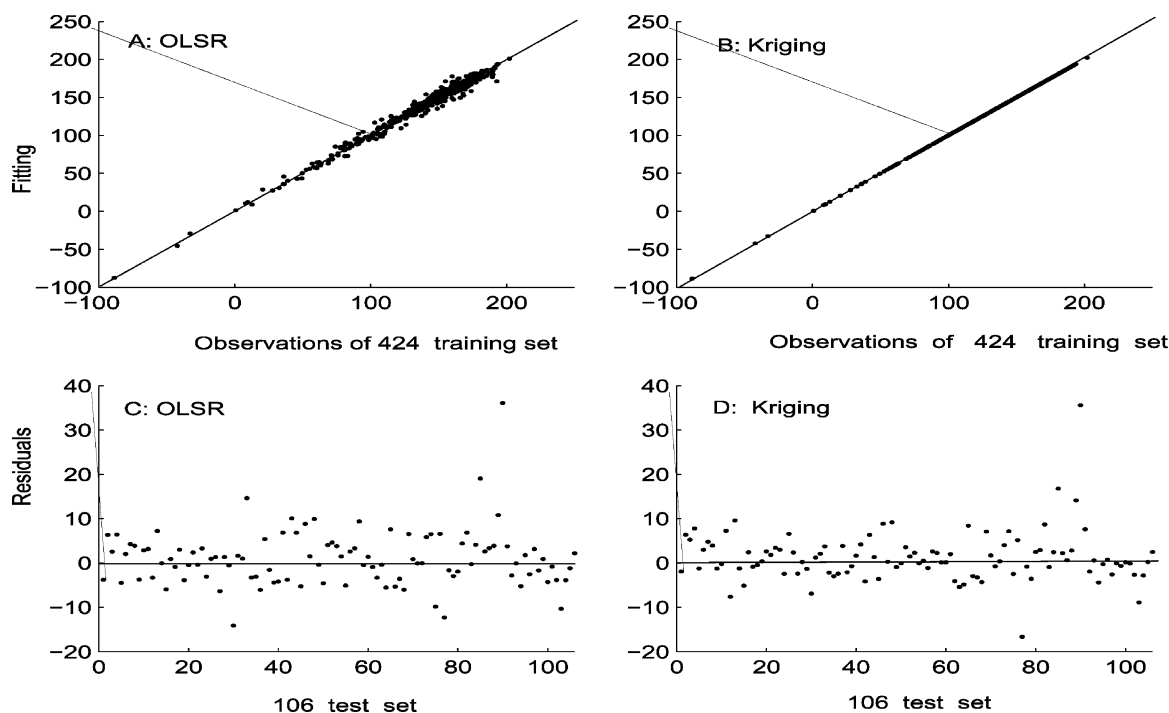


Figure 1. The fittings and residuals for training data and test data (A) and (C) the OLSR model and (B) and (D) the Kriging model.

(2) Training Data and Test Data. How to split the original data set into a training set and a test set is an important issue in the use of the Kriging model. In this paper we consider four cases and want to know which splitting way is the best.

(3) Comparisons. The criteria RMSE and AAE will be employed for comparing Kriging models with other meta-models in prediction ability.

The following sections will deal with these issues in detail.

The paper is organized as follows. Kriging is described in Section 2. Three methods, OLSR, PCR, and PLSR, are briefly introduced in Section 3, and they are used to be comparative models with Kriging. The Kriging model is also called semiparametric since the model contains a parametric item (or linear regression part) and a nonparametric part, the latter is considered as the realization of a random process. We shall show that if the process is stationary the Kriging model can improve the performance of model (1) under a suitable choice of basic functions. Some applications of this new approach to a real data of chemometrics are analyzed in Section 3. The detailed experimental process and the analysis of results are also presented. Our experiments show that the new approach can significantly improve the results obtained by other models where the errors are independent. Some important conclusions and the future work will be addressed in Section 4.

2. KRIGING MODELS

The word “Kriging” is synonymous with optimal spatial prediction. It has been named after a South African mining engineer with the name Krige, who first popularized stochastic methods for spatial predictions. A brief overview of Kriging is given in this section. For a comprehensive review readers can refer to the literature.^{17,18}

Consider the Kriging model (4) covariance

$$\text{Cov}[z(\mathbf{x}_i), z(\mathbf{x}_k)] = \sigma_z^2 R_\theta(\mathbf{x}_i, \mathbf{x}_k) = \sigma_z^2 \prod_{j=1}^p R_{\theta_j}^{(j)}(x_{ij} - x_{kj}) \quad (7)$$

where the correlation function $R_{\theta_j}^{(j)}$ can be chosen from Table 2 ($d_j = x_{ij} - x_{kj}$).

So the unknown parameters in the Kriging model (4) are the coefficient vector $\beta = [\beta_1, \beta_2, \dots, \beta_q]'$ and covariance's parameters $\sigma_z^2, \theta = [\theta_1, \theta_2, \dots, \theta_p]$. For estimations of these unknown parameters, the unbiased linear predictor $\hat{y}(\mathbf{x}) = \mathbf{c}'\mathbf{y}$ is favorable, where \mathbf{y} is the response vector of the training data and \hat{y} is the fitting result of the training data or the prediction result of the test data. It is easily shown that minimizing the mean square error of this predictor under the unbiasedness condition $\mathbf{f}(\mathbf{x}) = \mathbf{c}'\mathbf{F}$ one gets the best linear unbiased estimation of $y(\mathbf{x})$

$$\hat{y}(\mathbf{x}) = \hat{\mathbf{c}}'\mathbf{y} = \mathbf{f}(\mathbf{x})\hat{\beta} + \mathbf{r}(\mathbf{x})'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{F}\hat{\beta}) \quad (8)$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_q(\mathbf{x})], \mathbf{r}(\mathbf{x}) = [R_\theta(\mathbf{x}_1, \mathbf{x}), R_\theta(\mathbf{x}_2, \mathbf{x}), \dots, R_\theta(\mathbf{x}_n, \mathbf{x})]' \quad (9)$$

$$\mathbf{R} = [\mathbf{r}(\mathbf{x}_1), \mathbf{r}(\mathbf{x}_2), \dots, \mathbf{r}(\mathbf{x}_n)], \hat{\beta} = (\mathbf{F}'\mathbf{R}^{-1}\mathbf{F})^{-1}\mathbf{F}'\mathbf{R}^{-1}\mathbf{y} \quad (10)$$

and \mathbf{F} is an expanded design matrix with elements $f_{ij} = f_j(\mathbf{x}_i)$ ($1 \leq i \leq n, 1 \leq j \leq q$). When \mathbf{x} is the i th training sample \mathbf{x}_i , from (8) we have

$$\hat{y}(\mathbf{x}_i) = \mathbf{f}(\mathbf{x}_i)\hat{\beta} + \mathbf{r}(\mathbf{x}_i)'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{F}\hat{\beta}) = \mathbf{f}(\mathbf{x}_i)\hat{\beta} + \mathbf{e}_i'(\mathbf{y} - \mathbf{F}\hat{\beta}) = y_i \quad (11)$$

which shows theoretically that the Kriging predictor interpolates the training data. Therefore, the fitting process involves two stages: obtaining the weighted least squares estimate $\hat{\beta}$ and then interpolating the residuals $\mathbf{y} - \mathbf{F}\hat{\beta}$.

Another approach for estimations of β , σ_z^2 , and θ is by the maximum likelihood method. From the normality assumption of the Gaussian Kriging model, the density of \mathbf{y} is given by

$$(2\pi)^{-n/2} (\sigma_z^2)^{-n} |\mathbf{R}|^{-1/2} \exp \left\{ -\frac{1}{2\sigma_z^2} (\mathbf{y} - \mathbf{F}\beta)' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{F}\beta) \right\} \quad (12)$$

After dropping a constant, the log-likelihood function of the training data becomes

$$l(\beta, \sigma_z^2, \theta) = -\frac{1}{2} \left[n \ln(\sigma_z^2) + \ln(|\mathbf{R}|) + \frac{1}{\sigma_z^2} (\mathbf{y} - \mathbf{F}\beta)' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{F}\beta) \right] \quad (13)$$

For a given θ , by differentiation $l(\beta, \sigma_z^2, \theta)$ with respect to β and σ_z^2 , respectively, the maximum likelihood estimator of the β is just the weighted least squares estimator, and the maximum likelihood estimator of σ_z^2 is

$$\hat{\sigma}_z^2 = \frac{1}{n} (\mathbf{y} - \mathbf{F}\hat{\beta})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{F}\hat{\beta}) \quad (14)$$

The parameters $[\theta_1, \dots, \theta_p]$ in \mathbf{R} are obtained by minimizing $|\mathbf{R}|^{1/n} \hat{\sigma}_z^2$, and the latter is a function including only correlation parameters $[\theta_1, \dots, \theta_p]$. There is an algorithm¹⁹ by which we can compute estimator of θ efficiently.

The formulations above should not be used directly for practical computation because design matrix \mathbf{F} and correlation matrix \mathbf{R} may be nearly singular and σ_z^2 could be very small. So we have to modify the algorithm and get another one which is described in the Appendix of this paper.

3. NEW APPROACH AND EXPERIMENTS

We still consider the boiling points of 530 saturated hydrocarbons with 15 topological indices mentioned in the Introduction.

3.1. Parametric Item. In the Kriging model (4), the parametric item $\sum_{j=1}^q \beta_j f_j(\mathbf{x}_i)$ is a linear combination of a set of basis functions. In our experiments the basis functions in a parametric item is chosen by three methods: ordinary least squares regression (OLSR), principal components regression (PCR), and partial least squares regression (PLSR).

Ordinary Least Squares Regression. Given input columns $\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(p)}$, the linear model

$$\hat{\mathbf{y}} = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j \mathbf{x}_{(j)} \quad (15)$$

has been often employed to find the relationship between response variable and input variables because the linear model is simple and can obtain a good result in many cases. When the dimension of inputs is very large, one may use selection of variables to reduce the number of variables. In

our experiment, we simply use all variables ($q = p = 15$) as a set of basis functions.

Principal Components Regression. When a regression model has a large number of input variables, the collinearity between variables becomes a serious problem. One approach to overcome this difficulty is to transfer the original variables to new orthogonal variables, and we hope that only a few of the new variables will accommodate most of the variations of the data set. This approach has been long employed in the statistical literature. For example, the principal components regression accords with this idea. It “models out” the main part of the variation from the original data set \mathbf{X} by

$$\mathbf{z}_{(j)} = \mathbf{X} \mathbf{v}_j, j = 1, 2, \dots, q \leq p \quad (16)$$

where vector \mathbf{v}_j is the eigenvector associated to the j th largest eigenvalue of the covariance matrix $\text{Cov}(\mathbf{X})$. Here, $\mathbf{z}_{(j)}$ can be regarded as observations of the j th new variable. Since the principal components $\mathbf{z}_{(j)}$'s are orthogonal, this regression is just a sum of univariate regressions

$$\hat{\mathbf{y}} = \bar{y} + \sum_{j=1}^q \theta_j \mathbf{z}_{(j)} \quad (17)$$

where \bar{y} is the mean of observed responses of the training data. The choice of q will be discussed in Section 3.2. In our Kriging model, the basis functions in parametric item are just the principal components which are selected by PCR.

Partial Least Squares Regression. Partial least squares regression (PLSR) is another approach aiming for dimension reduction which was introduced by Wold in 1975. The idea of PLSR is similar to that of PCR but with the modification that both \mathbf{y} and \mathbf{X} are considered in the process. We assume that \mathbf{y} is centered, and each $\mathbf{x}_{(j)}$ is standardized to have mean 0 and variance 1. PLS begins by computing the univariate regression coefficient $\hat{\gamma}_{1j}$ of \mathbf{y} on each $\mathbf{x}_{(j)}$, that is, $\hat{\gamma}_{1j} = \mathbf{x}_{(j)}' \mathbf{y}$. So PLS constructs the derived input $\mathbf{z}_{(1)} = \sum_{j=1}^p \hat{\gamma}_{1j} \mathbf{x}_{(j)}$ which is the first partial least squares direction. Hence in the construction of each $\mathbf{z}_{(j)}$, the inputs are weighted by the strength of their univariate effect on \mathbf{y} .

Fitting the first new feature $\mathbf{z}_{(1)}$ to the response \mathbf{y} , the least-squares estimate $\hat{\theta}_1$ on $\mathbf{z}_{(1)}$ can be easily found. Then we orthogonalize $\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(p)}$ with respect to $\mathbf{z}_{(1)}$ to get residual matrix $\mathbf{X}^{(1)} = [\mathbf{x}_{(1)}^{(1)}, \mathbf{x}_{(2)}^{(1)}, \dots, \mathbf{x}_{(p)}^{(1)}]$. Just like the process producing the first PLS direction $\mathbf{z}_{(1)}$, we generate the second PLS direction $\mathbf{z}_{(2)}$ based on $\mathbf{X}^{(1)}$ and \mathbf{y} . Repeat the process until q ($q \leq p$) directions have been computed. The choice of q depends on the specific problem. We shall discuss this in Section 3.2. So the parametric item in the Kriging model can employ the q PLS directions selected by the above technique as a set of basis function.

3.2. Cross-Validation and Assessing Model Accuracy.

Cross-Validation. Cross-validation has been widely used in data modeling. Data are split into a training set and a test set. The training set is for model development, while the test set is for making a judgment of the performance of the selected model or the metamodel. When the number of observations is small, the k -fold cross-validation has been popularly used. In our practice we choose 5-fold cross-validation. In fact, how many fold cross-validation is not an important point in this paper. You can choose other fold cross-validation and get the same conclusions with us. In

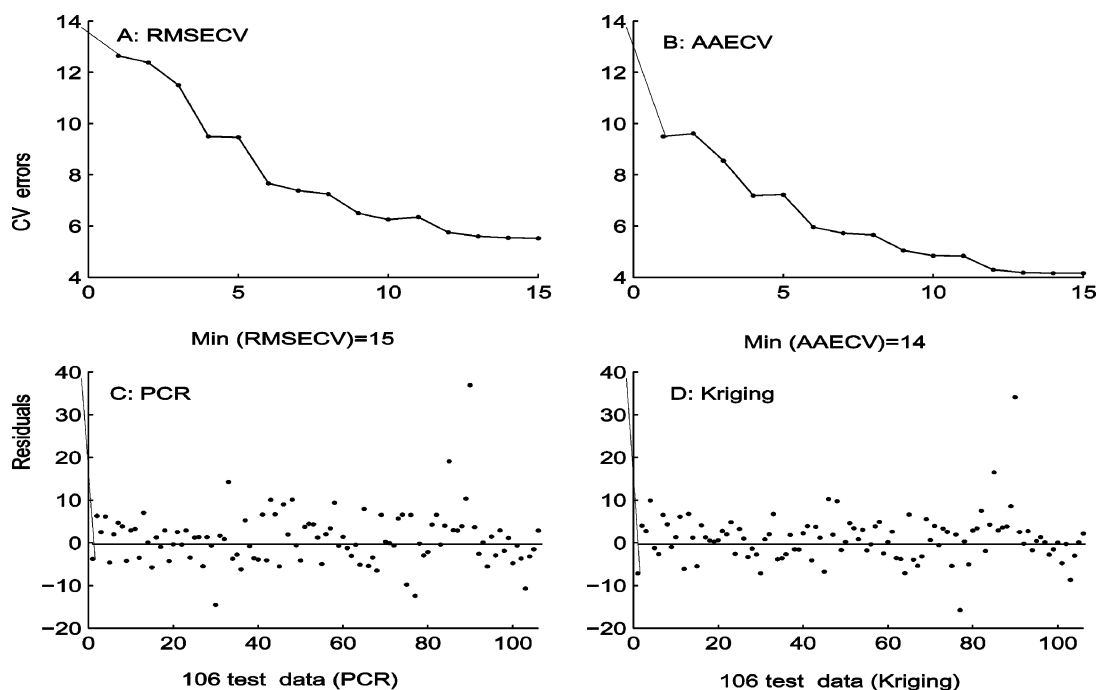


Figure 2. The results of cross-validation for training data and predictions of PCR and Kriging for test data.

our case, both PCR and PLSR have an unknown parameter q which can be chosen to minimize the estimate of prediction error based on 5-fold cross-validation. Briefly, 5-fold cross-validation works by dividing the training data randomly into 5 equal parts. The learning method is fit to four-fifths of the training data, and the prediction error is computed on the remaining one-fifth. This is done in turn for each one-fifth of the training data, and 5 prediction error estimates are averaged; finally we choose q to minimize the averaged estimate of prediction errors.

Assessing Model Accuracy. For comparing metamodelling in the cross-validation process, root-mean-squared error (RMSE) and average absolute error (AAE) can be used as criteria.

3.3. Experiments and Comparisons. The detailed procedure of our experiments on the boiling data is given as follows:

(1) The data were randomly divided into a training set \mathcal{T} and a test set \mathcal{S} . For studying the affect of the sample size of the training set, we consider four cases, i.e., each takes 20%, 40%, 60%, and 80% of the original data as a training data. As a consequence, the four training sets have $l_1 = 530 \times 0.2 = 106$, $l_2 = 530 \times 0.4 = 212$, $l_3 = 530 \times 0.6 = 318$, and $l_4 = 530 \times 0.8 = 424$ observations, respectively, so the size of corresponding test sets are $t_1 = 424$, $t_2 = 318$, $t_3 = 212$, and $t_4 = 106$.

(2) A metamodel was constructed by the three methods (OLSR, PCR, and PLSR) based on each of the training sets via the use of 5-fold cross-validation. For OLSR, this step is just to compute the estimation of coefficients β . For PCR and PLSR two estimators of q can be obtained. Let \hat{q}_1 be estimated under the RMSE and \hat{q}_2 be estimated under the AAE by the use of the method PCR/PLSR. The final estimator of q is suggested by $\hat{q} = \min(\hat{q}_1, \hat{q}_2)$. In this step we can have $12 = 3 \times 4$ metamodelling, and the basis functions of each metamodel are employed for the parametric item in the Kriging model that will be developed in the next step.

(3) For each metamodel obtained in the previous step the corresponding Kriging model can be established by the method introduced in Section 2.

(4) Calculate RMSE and AAE for each of the 12 metamodelling and their corresponding Kriging models. For fair comparisons we repeat the above process 100 times and use the average of respective 100 RMSEs and 100 AAEs to stand for predictive ability of the corresponding model.

We can also browse the fitting and prediction results of once sampling for PCR and PLSR and their corresponding Kriging models advance.

PCR and the Corresponding Kriging Models' Results Analysis. The same 424 training data and 106 test data are used in PCR and the corresponding Kriging model. First, we should carry out PCA with the training data and determine the number of components, that is \hat{q} , in the final PCR and corresponding Kriging by 5-fold cross-validation. Once q is estimated, the prediction results of test data for two metamodelling can be compared. Figure 2 shows the error results of 5-fold cross-validation of PCR for 424 training data and prediction results of PCR and the relevant Kriging model for 106 test data:

(A) and (B) in Figure 2 show the RMSE and AAE results of 5-fold cross-validation (RMSECV and AAECV) for the PCR model. The number of principal components (basis functions) in the final PCR and the Kriging model is 14 because of $\min(\hat{q}_1, \hat{q}_2) = \min(15, 14) = 14$. RMSE and AAE of test data for PCR and the corresponding Kriging model are as follows: RMSE (PCR) = 6.400, RMSE (Kriging) = 5.635, AAE (PCR) = 4.425, and AAE (Kriging) = 3.800. The improvement degrees of Kriging against PCR are 11.95% and 14.12% respectively.

PLSR and the Corresponding Kriging Models' Results Analysis. The building process of these two metamodelling are the same with that of PCR and its corresponding Kriging model. The final number of basis functions through 5-fold cross-validation is 12 ($\hat{q} = \min(\hat{q}_1, \hat{q}_2) = 12$). RMSE and

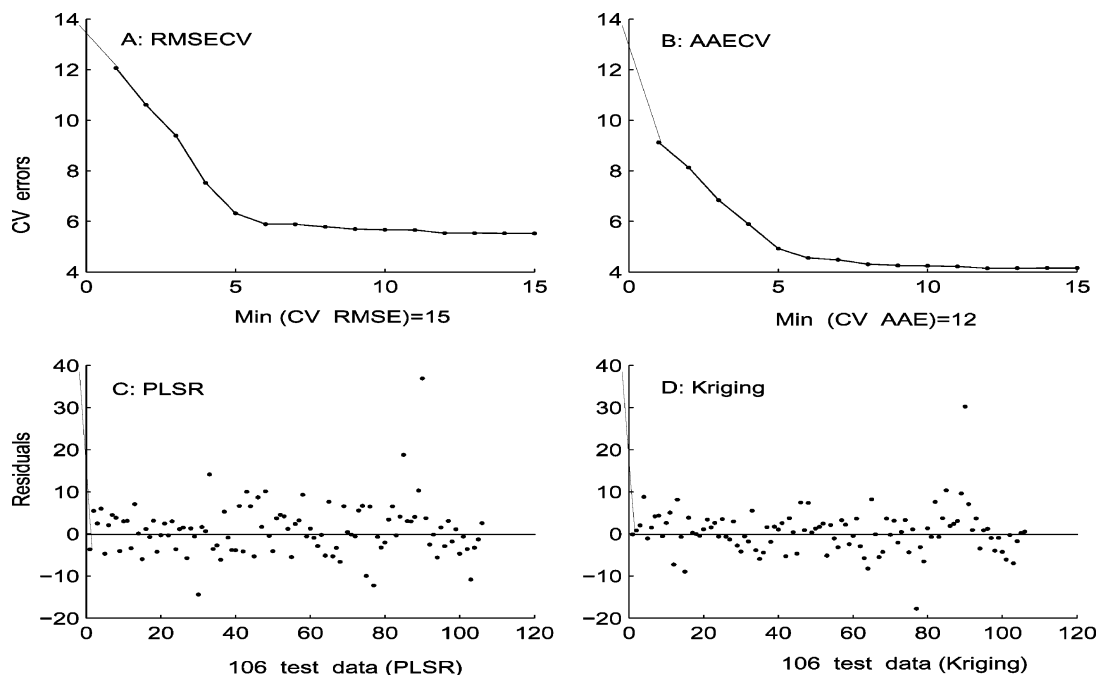


Figure 3. The results of cross-validation for training data and predictions of PLSR and Kriging for test data.

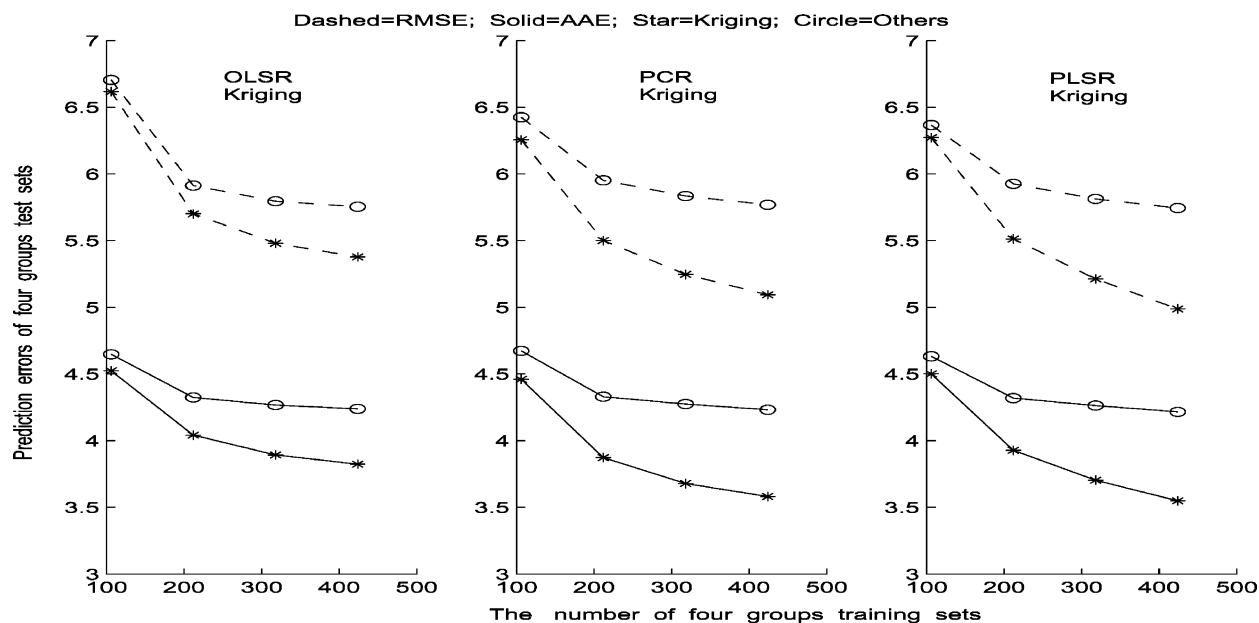


Figure 4. Prediction errors and degrees of improvement of four groups test sets by Kriging models against other metamodels.

AAE of the test data for PLSR and its corresponding Kriging are as follows: RMSE (PLSR) = 6.365, RMSE (Kriging) = 5.216, AAE (PLSR) = 4.395, and AAE (Kriging) = 3.499. The improvement degrees of Kriging against PLSR are 18.05% and 20.39%. Figure 3 is the error results of 5-fold cross-validation for 424 training data and prediction results of PLSR and relevant Kriging for 106 test data.

Above are the results of once random sampling for a sort of training data size. Now we will show the results of four groups training data with different sizes on 100 average. So we have 24 models in all: 12 metamodels are suggested by OLSR, PCR, and PLSR, and 12 their corresponding Kriging models. Figure 4 gives comparisons of each metamodel and its corresponding Kriging model based on two criteria, RMSE and AAE, where there are dashed lines for RMSE and solid lines for AAE. Abscissa in Figure 4 denotes the

number of four training set groups, and ordinate shows the averaged prediction errors of the corresponding four test set groups. Circles are the results of OLSR, PCR, and PLSR, and stars denote the performances of Kriging models. It is clear that Kriging is superior to all corresponding metamodels. Table 3 lists averaged AAEs and RMSEs for the 24 models and the degrees of improvement. From Table 3 we will address the following remarks from several aspects:

Sample Size. We note that all models accuracy improves when training samples size increases, and the larger training set size is, the more improvement by the Kriging model is.

RMSE and AAE. The trends of RMSE and AAE are similar, but AAE lines are always below the RMSE lines.

The Number of Variables in Final Models. The first column of Table 3 lists the number of variables in each final models. The models of the first subtable use all original

Table 3. Prediction Results and Improvements of Four Groups Test Sets by Kriging against OLSR, PCR, and PLSR

no. of var. q	test set T	AAE		RMSE		increase		
		OLSR	Kriging	OLSR	Kriging	Kriging (%)	vs	OLSR (%)
$\hat{q} = 15$	$t_1 = 424$	4.6476	4.5238	6.7040	6.6164	2.67		1.33
$\hat{q} = 15$	$t_2 = 318$	4.3243	4.0416	5.9123	5.7006	6.54		3.57
$\hat{q} = 15$	$t_3 = 212$	4.2675	3.8936	5.7948	5.4802	8.76		5.46
$\hat{q} = 15$	$t_4 = 106$	4.2389	3.8243	5.7537	5.3776	9.65		6.51

no. of var. q	test set T	AAE		RMSE		increase		
		PCR	Kriging	PCR	Kriging	Kriging (%)	vs	PCR (%)
$\hat{q} = 11.89$	$t_1 = 424$	4.6729	4.4599	6.4244	6.2556	4.48		2.57
$\hat{q} = 13.25$	$t_2 = 318$	4.3295	3.8732	5.9518	5.5009	10.52		7.60
$\hat{q} = 13.79$	$t_3 = 212$	4.2748	3.6795	5.8336	5.2472	13.87		10.09
$\hat{q} = 13.81$	$t_4 = 106$	4.2318	3.5816	5.7684	5.0939	15.22		11.70

no. of var. \hat{q}	test set T	AAE		RMSE		increase		
		PLSR	Kriging	PLSR	Kriging	Kriging (%)	vs	PLSR (%)
$\hat{q} = 9.6$	$t_1 = 424$	4.6324	4.5015	6.3660	6.2714	2.85		1.49
$\hat{q} = 11.56$	$t_2 = 318$	4.3193	3.9284	5.9256	5.5119	9.04		7.02
$\hat{q} = 12.36$	$t_3 = 212$	4.2637	3.7037	5.8124	5.2128	13.12		10.36
$\hat{q} = 12.6$	$t_4 = 106$	4.2164	3.5494	5.7443	4.9875	15.72		13.21

variables. \hat{q} 's of the second and the third subtable are the averaged estimations based on 5-fold cross-validation after 100 duplications. We find the number of variables listed in the third subtable are less than those of the other two subtables, which testifies a recognized conclusion that PLSR closely reach the PCR prediction results using fewer components than those of PCR.

Degrees of Improvement. The degrees of improvement by Kriging against OLSR are very minor relative to other improvement results, but it becomes obvious when the training set size increases. The degrees of improvement by Kriging against PCR and PLSR are very apparent, and they also increase when the number of training data are being added.

Impact of Basis Functions. In fact, we have constructed three Kriging models based on three different sets of basis functions. If you want to obtain perfect prediction results and better improvements, Kriging based on principal components can be chosen. On the other hand, if you want to pick out a model which does not have better prediction ability and improvements, then the most parsimonious model, Kriging based on PLS directions, may be a better choice.

4. CONCLUSION

This boiling points data have been utilized to demonstrate that the Kriging method is a promising model approximation technique in QSAR research. There are several research issues to address for the application of Kriging and our future work:

1. Since the Kriging model interpolates the training data, if we select the training set by experimental designs, the prediction results may be better.

2. In our example, we use the original variables or a set of basis functions computed by PCR and PLSR as inputs. Particular basis functions such as piecewise polynomials, splines, and wavelets bases also can be employed.

3. The *parametric item* $\sum_{j=1}^q \beta_j f_j(\mathbf{x}_i)$ is a global reflection. If it is not designed correctly, the results will be affected.

Future work on the boiling points includes adding more design variables, changing the basis functions, and employing

other experiments design techniques and comparing with more approximation approaches.

ACKNOWLEDGMENT

The work is partially supported by the Hong Kong UGC grant RGC/HKBU 2044/02P and Statistics Research and Consultancy Centre, Hong Kong Baptist University.

APPENDIX

1. Compute \mathbf{R} 's Cholesky factor \mathbf{C}

$$\mathbf{R} = \mathbf{C}\mathbf{C}' \quad (18)$$

where \mathbf{C} is an upper triangular matrix, and it can be obtained in Matlab language $\mathbf{C} = \text{chol}(\mathbf{R})$.

2. Carry out "decorrelation" transformation:

$$\tilde{\mathbf{F}} = \mathbf{C}^{-1}\mathbf{F}, \tilde{\mathbf{y}} = \mathbf{C}^{-1}\mathbf{y} \quad (19)$$

3. Compute QR factorization of $\tilde{\mathbf{F}}$

$$\tilde{\mathbf{F}} = \mathbf{Q}\mathbf{G}' \quad (20)$$

where \mathbf{Q} has orthonormal columns and \mathbf{G}' is upper triangular. They can be computed in Matlab language $[\mathbf{Q}, \mathbf{G}', \mathbf{E}] = \text{qr}(\tilde{\mathbf{F}})$.

4. Compute the estimations of β and σ_z^2 :

$$\hat{\beta} = \mathbf{G}^{-1}\mathbf{Q}'\tilde{\mathbf{y}}, \hat{\sigma}_z^2 = \frac{1}{n}(\tilde{\mathbf{y}} - \tilde{\mathbf{F}}\hat{\beta})'(\tilde{\mathbf{y}} - \tilde{\mathbf{F}}\hat{\beta}) \quad (21)$$

5. Compute the objective function

$$\hat{\sigma}_z^2 |\mathbf{R}|^{1/n} = \hat{\sigma}_z^2 \Pi(c_{jj})^{2/n} \quad (22)$$

where c_{jj} 's are the diagonal elements of the Cholesky factor \mathbf{C} .

6. Simple pattern search to optimize θ .

REFERENCES AND NOTES

- (1) Massy, W. F. Principal components regression in exploratory statistical research. *J. Am. Stat. Assoc.* **1965**, *60*, 234–246.

- (2) Wold, H. Soft modelling by latent variables: the nonlinear iterative partial least squares (NIPALS) approach. *Perspect. Probability Stat.* In Honor of M. S. Bartlett, **1975**, 117–144.
- (3) Friedman, J. Multivariate adaptive regression splines (with discussion). *Ann. Statistics* **1991**, 19(1), 1–141.
- (4) Friedman, J. *Stochastic gradient boosting*; Techniquial report; Stanford University, 1999.
- (5) Friedman, J. Greedy function approximation: the gradient boosting machine. *Ann. Statistics* **2001**, 29(5).
- (6) Green, P.; Silverman, B. *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty approach*; Chapman and Hall: London, 1994.
- (7) Wahba, G. *Spline Models for Observational Data*; SIAM: Philadelphia, 1990.
- (8) Daubechies, I. *Ten Lectures in Wavelets*; Society for Industrial and Applied Mathematics: Philadelphia, PA, 1992.
- (9) Li, K. C. Sliced inverse regression for dimension reduction. *J. Am. Stat. Assoc.* **1991**, 86, 316–327.
- (10) He, P.; Fang, K. T.; Xu, C. J. The classification tree combined with SIR and its applications to classification of mass spectra. *J. Data Sci.* **2003**, 1, 425–445.
- (11) Friedman, J.; Tukey, J. A projection pursuit algorithm for exploratory data analysis. *IEEE Trans. Comput. Ser. C* **1974**, 23, 881–889.
- (12) Tibshirani, R. Regression shrinkage and selection via the lasso. *J. Royal Statist. Soc. B* **1996**, 58, 267–288.
- (13) Rohrbaugh, R. H.; Jurs, P. C. Prediction of Gas Chromatographic Retention Indexes of Selected Olefins. *Anal. Chem.* **1985**, 57, 2770–3.
- (14) Du, Y. P.; Liang, Y. Z.; Yun, D. Data Mining for Seeking an Accurate Quantitative Relationship between Molecular Structure and GC Retention Indices of Alkenes by Projection Pursuit. *J. Chem. Inf. Comput. Sci.* **2002**, 42, 1283–1292.
- (15) Sacks, J.; Welch, W. J.; Mitchell, T. J.; Wynn, H. P. Design and analysis of computer experiments. *Stat. Sci.* **1989**, 4, 409–435.
- (16) Rücker, G.; Rücker, C. On topological indices, boiling points, and cycloalkanes. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 788–802.
- (17) Cressie, N. *Statistics for Spatial Data*, revised edition; Wiley: New York, 1993.
- (18) Rivoirard, J. *Introduction to Disjunctive Kriging and Nonlinear Geostatistics*; Oxford University Press: Oxford, 1994.
- (19) Lophaven, S. N.; Nielsen, H. B.; Søndergaard, J. Aspects of the Matlab Toolbox DACE. Report IMM-REP-2002-13, *Informatics and Mathematical Modelling*. DTU, **2002**, 44 pages.

CI049798M