

# A VARIABLE SELECTION METHOD IN A MULTIVARIATE NONPARAMETRIC REGRESSION MODEL: APPLICATION TO GEOSCIENCE

Mary E. Savino <sup>1</sup> & Céline Lévy-Leduc <sup>2</sup>

<sup>1</sup> *Andra, 92290 Châtenay-Malabry, France and Université Paris-Saclay, AgroParisTech, INRAE, UMR MIA Paris-Saclay, 91120 Palaiseau, France*

*mary.savino@agroparistech.fr*

<sup>2</sup> *Université Paris-Saclay, AgroParisTech, INRAE, UMR MIA Paris-Saclay, 91120 Palaiseau, France*

*celine.levy-leduc@agroparistech.fr*

**Résumé.** Nous présentons ici une nouvelle méthode de sélection de variables dans un modèle de régression non-paramétrique multivarié et reposant sur des données afin d’identifier les variables dont dépend réellement la fonction de régression. Cette méthode consiste à approcher la fonction sous-jacente par une combinaison linéaire de B-splines d’ordre  $M$  ainsi que par la combinaison de leurs interactions deux-à-deux. Les coefficients de cette combinaison linéaire sont estimés en minimisant le critère des moindres carrés pénalisé par la somme des normes  $\ell_2$  des dérivées partielles par rapport à chaque variable dont dépend la fonction. Nous montrons que la méthode proposée peut être reformulée sous la forme d’un critère de type Group Lasso. Nous validons notre approche à travers différentes expériences numériques en faisant notamment varier le nombre d’observations, le niveau de bruit et le nombre total de variables. Nous la comparons également à deux autres méthodes de l’état de l’art et une application à un système géochimique réel est présentée. A travers ces différentes applications, notre approche démontre de meilleures performances statistiques que les autres méthodes auxquelles nous l’avons comparée. Notre méthode est implémentée dans le package R `absorber` qui sera bientôt disponible sur le “Comprehensive R Archive Network” (CRAN).

**Mots-clés.** Sélection de variables, régression non paramétrique, B-splines, Group Lasso

**Summary.** In this presentation, we introduce a novel data-driven variable selection approach in a multivariate nonparametric regression model designed to capture only the variables on which the regression function depends. The underlying idea of our method consists in approximating the function by a linear combination of B-splines of order  $M$  and their pairwise interactions. The coefficients of this linear combination are estimated by minimizing the penalized least-squares criterion. The penalization consists of the sum of the  $\ell_2$ -norms of the partial derivatives with respect to the different variables on which the function depends. We show that our proposed method can be reformulated as a Group Lasso problem. We investigate the statistical performance of our approach through numerical experiments varying the number of observations, the noise level and the total number of variables. We also compare it to two other state-of-the-art methods. An application to a geochemical system is also proposed. In these different frameworks, our approach exhibits better performance than the other methods. Our completely data-driven method is implemented in the `absorber` R package which will be soon available on the Comprehensive R Archive Network (CRAN).

**Keywords.** Variable selection, nonparametric regression, B-splines, Group Lasso

# 1 Introduction

The simulation of geochemical models that incorporate precipitation and dissolution reactions of minerals coupled to other physical processes represents a challenging task. Reactive transport modeling (RTM) serves as an illustration, striving to simultaneously consider geochemical reactions, fluid flow, heat transfer and solute transport. This challenge has led to the development of Machine Learning (ML)-based approaches aimed at estimating real solutions for full simulation models through the use of surrogate models. The main idea here consists in solving the transport equations explicitly and approximating solutions for geochemical reactions at equilibrium using surrogate models at each time step. A wealth of reviews and surveys on surrogate models for RTM is available in the works of Razavi et al. (2012); Asher et al. (2015); Jatnieks et al. (2016). Another approach to improve the surrogate model accuracy while reducing the CPU times is to reduce the number of input variables to consider in the model. This can be reformulated as a variable selection problem in the following framework. Let us consider that we have  $n$  observations satisfying the following nonparametric regression model:

$$Y_i = f(x_i) + \varepsilon_i, \quad x_i = (x_i^{(1)}, \dots, x_i^{(p)}) \in \mathbb{R}^p, \quad 1 \leq i \leq n \quad (1)$$

where  $f$  is an unknown real-valued function and where the  $\varepsilon_i$ 's are i.i.d centered random variables of variance  $\sigma^2$ . We will also assume that  $f$  actually depends only on  $\tilde{d}$  variables instead of  $p$ , with  $\tilde{d} < p$ , which means that there exists a real-valued function  $\tilde{f}$  such that  $f(x) = \tilde{f}(\tilde{x})$ , where  $x \in \mathbb{R}^p$  and  $\tilde{x} \in \mathbb{R}^{\tilde{d}}$ . Variable selection consists in identifying the components of  $\tilde{x}$ .

We propose a novel method for variable selection motivated by Radchenko and James (2010) using a multivariate nonparametric regression model to retrieve the  $\tilde{d}$  relevant variables on which  $f$  in (1) depends. Our approach, presented in Section 2, consists in approximating  $f$  using a linear combination of B-splines and their pairwise interactions. Additionally, drawing inspiration from the methodology of Rosasco et al. (2010), the coefficients of the linear combination are estimated by minimizing the usual least-squares criterion penalized by the sum of the  $\ell_2$ -norms of the partial derivatives with respect to the different variables on which  $f$  depends. We show that our proposed method can be reformulated as a Group Lasso problem defined by Yuan and Lin (2006). Two different approaches to choose the penalization parameter are presented. The statistical performance of our approach are investigated in Section 3 and a geochemical application is given in Section 4.

## 2 Method

### 2.1 Approximation of $f$ using B-splines

Let  $\mathbf{t}_\ell = (t_{\ell,1}, \dots, t_{\ell,K})$  be a set of  $K$  points called knots and let  $\mathcal{S}_\ell$  be a compact subset of  $\mathbb{R}$ . Following De Boor (1978, p. 89-90) and Hastie et al. (2009, p. 160), the augmented knot sequence  $\boldsymbol{\tau}_\ell$  is defined as follows:

$$\begin{aligned} \tau_{\ell,1} &= \dots = \tau_{\ell,M} = x_{min}^{(\ell)}, \\ \tau_{\ell,j+M} &= t_{\ell,j}, \quad j = 1, \dots, K, \\ \tau_{\ell,K+M+1} &= \dots = \tau_{\ell,K+2M} = x_{max}^{(\ell)}, \\ \boldsymbol{\tau}_\ell &= (\tau_{\ell,1}, \dots, \tau_{\ell,K+2M}) = \underbrace{(x_{min}^{(\ell)}, \dots, x_{min}^{(\ell)})}_{M \text{ times}}, \underbrace{(t_{\ell,1}, \dots, t_{\ell,K})}_{\mathbf{t}_\ell}, \underbrace{(x_{max}^{(\ell)}, \dots, x_{max}^{(\ell)})}_{M \text{ times}}, \end{aligned}$$

where  $x_{min}^{(\ell)}$  and  $x_{max}^{(\ell)}$  are the lower and upper bounds of  $\mathcal{S}_\ell$ , respectively.

Denoting by  $B_{k,m}^{(\ell)}$  the  $k$ th B-spline basis function of order  $m$  with  $m \leq M$  for the knot sequence  $\boldsymbol{\tau}_\ell$  and for the dimension  $\ell$ , B-splines are defined by the following recursion:

$$B_{k,1}^{(\ell)}(x^{(\ell)}) = \begin{cases} 1 & \text{if } \tau_{\ell,k} \leq x^{(\ell)} < \tau_{\ell,k+1} \\ 0 & \text{otherwise} \end{cases} \quad \text{for } k = 1, \dots, K + 2M - 1, \quad (2)$$

and for  $2 \leq m \leq M$ ,

$$B_{k,m}^{(\ell)}(x^{(\ell)}) = \frac{x^{(\ell)} - \tau_{\ell,k}}{\tau_{\ell,k+m-1} - \tau_{\ell,k}} B_{k,m-1}^{(\ell)}(x^{(\ell)}) + \frac{\tau_{\ell,k+m} - x^{(\ell)}}{\tau_{\ell,k+m} - \tau_{\ell,k+1}} B_{k+1,m-1}^{(\ell)}(x^{(\ell)}), \quad (3)$$

for  $k = 1, \dots, (K + 2M - m)$ .

Inspired by Radchenko and James (2010), we propose approximating the function  $f(x^{(1)}, \dots, x^{(p)})$  appearing in (1) by a linear combination of B-splines of each variable  $x^{(1)}, \dots, x^{(p)}$  and of pairwise interaction of them as follows:

$$F(x^{(1)}, \dots, x^{(p)}) = \sum_{\ell=1}^p \sum_{k=1}^{K+M} \beta_k^{(\ell)} B_k^{(\ell)}(x^{(\ell)}) + \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \left( \sum_{k=1}^{K+M} \sum_{q=1}^{K+M} \beta_{k,q}^{(\ell,j)} B_k^{(\ell)}(x^{(\ell)}) B_q^{(j)}(x^{(j)}) \right), \quad (4)$$

where  $B_k^{(\ell)} = B_{k,M}^{(\ell)}$  is defined in (2) and (3) and where  $\beta_k^{(\ell)}$  and  $\beta_{k,q}^{(\ell,j)}$  are unknown coefficients. Observe that the column vector  $(F(x_i^{(1)}, \dots, x_i^{(p)}))_{1 \leq i \leq n}$  (4) can be rewritten as follows:

$$\sum_{\ell=1}^p \Psi_\ell \boldsymbol{\beta}_\ell + \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \Phi_{\ell j} \boldsymbol{\beta}_{\ell,j}. \quad (5)$$

where  $\Psi_\ell$  is a  $n \times (K+M)$  matrix such that its  $i$ th row is equal to  $(B_1^{(\ell)}(x_i^{(\ell)}), \dots, B_{K+M}^{(\ell)}(x_i^{(\ell)}))$  and  $\beta_\ell = (\beta_1^{(\ell)} \dots \beta_{K+M}^{(\ell)})^T$  for  $1 \leq \ell \leq p$ ,  $A^T$  denoting the transpose of the matrix  $A$ . Moreover,  $\Phi_{\ell j}$  is an  $n \times (K+M)^2$  matrix such that its  $i$ th row satisfies  $(\Phi_{\ell j})_{i,\bullet} = ((\Psi_\ell)_{i,\bullet} \otimes (\Psi_j)_{i,\bullet})$ ,  $\otimes$  denoting the Kronecker product,  $(\Psi_\ell)_{i,\bullet}$  denoting the  $i$ th row of  $\Psi_\ell$  and  $\beta_{\ell,j} = (\beta_{1,1}^{(\ell,j)} \beta_{1,2}^{(\ell,j)} \dots \beta_{K+M,K+M}^{(\ell,j)})^T$  for  $1 \leq \ell < j \leq p$ .

## 2.2 Description of our variable selection method

Inspired by the methodology of Rosasco et al. (2010), we propose selecting the variables on which  $f$  depends by estimating the coefficients  $\beta_\ell$  and  $\beta_{\ell,j}$  appearing in (5) through the minimization of the following regularized criterion:

$$\begin{aligned} & (\widehat{\beta}_1(\lambda), \dots, \widehat{\beta}_p(\lambda), \widehat{\beta}_{1,2}(\lambda), \dots, \widehat{\beta}_{(p-1),p}(\lambda)) \\ &= \underset{\substack{(\beta_1, \dots, \beta_p) \\ (\beta_{1,2}, \dots, \beta_{(p-1),p})}}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \sum_{\ell=1}^p \Psi_\ell \beta_\ell - \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \Phi_{\ell j} \beta_{\ell,j} \right\|_2^2 + \lambda \sum_{\ell=1}^p \sqrt{\sum_{i=1}^n \partial_\ell F(x_i)^2} \right), \end{aligned}$$

where  $\mathbf{Y} = (Y_1, \dots, Y_n)$ , the  $Y_i$ 's being defined in (1),  $\partial_\ell F(x_i)$  denotes the  $\ell$ th partial derivative of  $F$  defined in (4) at some observation point  $x_i = (x_i^{(1)}, \dots, x_i^{(p)})$  and  $\|y\|_2^2 = \sum_{i=1}^n y_i^2$ . Note that the idea underlying this criterion is that when a function does not depend on a variable its partial derivative with respect to this variable is equal to zero.

Using the definition of  $F$  given in (5), the criterion can be rewritten as follows:

$$\begin{aligned} & (\widehat{\beta}_1(\lambda), \dots, \widehat{\beta}_p(\lambda), \widehat{\beta}_{1,2}(\lambda), \dots, \widehat{\beta}_{(p-1),p}(\lambda)) \\ &= \underset{\substack{(\beta_1, \dots, \beta_p) \\ (\beta_{12}, \dots, \beta_{(p-1)p})}}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \sum_{\ell=1}^p \Psi_\ell \beta_\ell - \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \Phi_{\ell j} \beta_{\ell,j} \right\|_2^2 \right. \\ & \quad \left. + \lambda \sum_{\ell=1}^p \left\| \Psi'_\ell \beta_\ell + \sum_{j=\ell+1}^p (\partial_\ell \Phi_{\ell j}) \beta_{\ell,j} + \sum_{1 \leq j < \ell} (\partial_\ell \Phi_{j\ell}) \beta_{j,\ell} \right\|_2 \right), \end{aligned} \tag{6}$$

where  $\Psi'_\ell$  is the  $n \times (K+M)$  matrix such that  $(\Psi'_\ell)_{i,k} = B_k^{(\ell)'}(x_i^{(\ell)})$ ,  $B_k^{(\ell)'}$  denoting the first derivative of  $B_k^{(\ell)}$ . The  $i$ th row of  $(\partial_\ell \Phi_{\ell j})$  (resp.  $(\partial_\ell \Phi_{j\ell})$ ) is defined by  $(\partial_\ell \Phi_{\ell j})_{i,\bullet} = ((\Psi'_\ell)_{i,\bullet} \otimes (\Psi_j)_{i,\bullet})$  (resp.  $(\partial_\ell \Phi_{j\ell})_{i,\bullet} = ((\Psi_j)_{i,\bullet} \otimes (\Psi'_\ell)_{i,\bullet})$ ). By denoting  $(\partial_\ell \Phi_{\ell\bullet}) = ((\partial_\ell \Phi_{\ell(\ell+1)}) \dots (\partial_\ell \Phi_{\ell p}))$ ,  $(\partial_\ell \Phi_{\bullet\ell}) = ((\partial_\ell \Phi_{1\ell}) \dots (\partial_\ell \Phi_{(\ell-1)\ell}))$ ,  $\beta_{\ell\bullet} = (\beta_{\ell,(\ell+1)}^T \dots \beta_{\ell,p}^T)^T$  and  $\beta_{\bullet\ell} = (\beta_{1,\ell}^T \dots \beta_{(\ell-1),\ell}^T)^T$ , the penalty term can be written as:

$$\lambda \sum_{\ell=1}^p \left\| \Psi'_\ell \beta_\ell + (\partial_\ell \Phi_{\ell\bullet}) \beta_{\ell\bullet} + (\partial_\ell \Phi_{\bullet\ell}) \beta_{\bullet\ell} \right\|_2 =: \lambda \sum_{\ell=1}^p \left\| (\partial_\ell \Theta_\ell) \gamma_\ell \right\|_2, \tag{7}$$

where  $\boldsymbol{\gamma}_\ell = (\boldsymbol{\beta}_\ell^T \boldsymbol{\beta}_{\ell\bullet}^T \boldsymbol{\beta}_{\bullet\ell}^T)^T$ . The least-squares term can be rewritten as follows:

$$\begin{aligned}
& \left\| \mathbf{Y} - \sum_{\ell=1}^p \Psi_\ell \boldsymbol{\beta}_\ell - \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \Phi_{\ell j} \boldsymbol{\beta}_{\ell,j} \right\|_2^2 \\
&= \left\| \mathbf{Y} - \sum_{\ell=1}^p \Psi_\ell \boldsymbol{\beta}_\ell - \frac{1}{2} \left( \sum_{\ell=1}^{p-1} \sum_{j=\ell+1}^p \Phi_{\ell j} \boldsymbol{\beta}_{\ell,j} + \sum_{\ell=2}^p \sum_{j=1}^{\ell-1} \Phi_{j\ell} \boldsymbol{\beta}_{j,\ell} \right) \right\|_2^2 \\
&=: \left\| \mathbf{Y} - \sum_{\ell=1}^p \Theta_\ell \boldsymbol{\gamma}_\ell \right\|_2^2.
\end{aligned} \tag{8}$$

Equation (8) comes by setting  $\Theta_1 = \left( \Psi_1 \quad \frac{1}{2} \Phi_{1\bullet} \right)$  and  $\Theta_p = \left( \Psi_p \quad \frac{1}{2} \Phi_{\bullet p} \right)$ , where  $\Phi_{\ell\bullet} = (\Phi_{\ell(\ell+1)} \dots \Phi_{\ell p})$  and  $\Phi_{\bullet\ell} = (\Phi_{1\ell} \dots \Phi_{(\ell-1)\ell})$ . Combining (7) and (8), (6) can be rewritten as:

$$(\hat{\boldsymbol{\gamma}}_1(\lambda), \dots, \hat{\boldsymbol{\gamma}}_p(\lambda)) = \underset{(\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p)}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \sum_{\ell=1}^p \Theta_\ell \boldsymbol{\gamma}_\ell \right\|_2^2 + \lambda \sum_{\ell=1}^p \left\| (\partial_\ell \Theta_\ell) \boldsymbol{\gamma}_\ell \right\|_2 \right). \tag{9}$$

By defining  $\boldsymbol{\alpha}_\ell = (\partial_\ell \Theta_\ell) \boldsymbol{\gamma}_\ell$  and  $\tilde{\mathbf{X}}_\ell = \Theta_\ell (\partial_\ell \Theta_\ell)^+$ ,  $A^+$  being the Moore-Penrose inverse of matrix  $A$ , (9) can be rewritten as:

$$(\hat{\boldsymbol{\alpha}}_1(\lambda), \dots, \hat{\boldsymbol{\alpha}}_p(\lambda)) = \underset{(\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_p)}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \sum_{\ell=1}^p \tilde{\mathbf{X}}_\ell \boldsymbol{\alpha}_\ell \right\|_2^2 + \lambda \sum_{\ell=1}^p \left\| \boldsymbol{\alpha}_\ell \right\|_2 \right). \tag{10}$$

The last formulation of our variable selection criterion (10) can be seen as a group lasso problem introduced by Yuan and Lin (2006), where the size  $p_\ell$  of each group  $\ell$  belonging to  $\{1, \dots, p\}$  is equal to  $n$ . The coefficients  $\hat{\boldsymbol{\gamma}}_\ell(\lambda)$  are thus obtained as follows:

$$\hat{\boldsymbol{\gamma}}_\ell(\lambda) = (\partial_\ell \Theta_\ell)^+ \hat{\boldsymbol{\alpha}}_\ell(\lambda). \tag{11}$$

Thus, we define the active variables for each  $\lambda$  belonging to a given set  $\Lambda$  as follows:

$$\mathcal{V}_\lambda = \left\{ \ell, \sum_{k \geq 1} |\hat{\boldsymbol{\gamma}}_{\ell,k}(\lambda)| \neq 0 \right\}, \tag{12}$$

where  $\hat{\boldsymbol{\gamma}}_{\ell,k}(\lambda)$  is the  $k$ th coefficient of  $\hat{\boldsymbol{\gamma}}_\ell(\lambda)$ . We also introduce the set  $\mathcal{V}_f$  of the indices of the  $d$  relevant variables on which  $f$  in (1) actually depends that have to be selected among the  $p$  variables. We also denote the set  $\overline{\mathcal{V}}_f$  of the indices of the irrelevant variables on which  $f$  does not depend.

### 3 Numerical experiments

In this section, we will investigate the statistical performance of our method called ABSORBER and implemented in the `absorber` R package when the variance of the noise  $\sigma^2$

increases as well as the number of observations  $n$ . We will also study how this novel method behaves when the number of variables  $p$  grows. To demonstrate the efficiency of our method, we will compare it to two state-of-the-art methods for feature selection: LassoNet introduced in Lemhadri et al. (2021) and the widely used Random Forests (RF) introduced by Breiman (2001), using their default parameters.

### 3.1 Metrics and selection criteria

We first introduce metrics to assess the efficiency of our method:

- **True Positive Rate (TPR)** and the **False Positive Rate (FPR)**, for each  $\lambda$ :

$$\text{TPR}(\lambda) = \frac{\text{TP}(\lambda)}{d} = \frac{|\mathcal{V}_\lambda \cap \mathcal{V}_f|}{d} \quad \text{and} \quad \text{FPR}(\lambda) = \frac{\text{FP}(\lambda)}{p-d} = \frac{|\mathcal{V}_\lambda \cap \overline{\mathcal{V}_f}|}{p-d},$$

where  $d < p$ ,  $|\mathcal{A}|$  is the cardinality of the set  $\mathcal{A}$ ,  $\text{TP}(\lambda)$  and  $\text{FP}(\lambda)$  are the number of true selected variables and the number of false selected variables for  $\lambda$ , respectively.  $\mathcal{V}_\lambda$ ,  $\mathcal{V}_f$  and  $\overline{\mathcal{V}_f}$  and are introduced in (12) and in the text following this equation.

We also propose two criteria to choose  $\lambda$ . One allows the user to choose a threshold of percentage of selection and the other leverages the Akaike Information Criterion (AIC), introduced in Akaike (1973) to automatically choose  $\lambda$ . Both are defined as follows:

- **Percentage of variable selection**, for each variable  $\ell$  belonging to  $\{1, \dots, p\}$ :

$$P_\ell = \frac{100}{|\Lambda|} \sum_{\lambda \in \Lambda} \mathbb{1}\{\ell \in \mathcal{V}_\lambda\}, \quad (13)$$

where  $|\Lambda|$  is the total number of parameters in the set  $\Lambda$ ,  $\mathbb{1}\{A\} = 1$  if the event  $A$  holds and 0 if not and  $\mathcal{V}_\lambda$  is defined in (12).

- **AIC**:

$$\text{AIC}(\lambda) = n \ln \left( \frac{\text{RSS}(\lambda)}{n} \right) + 2T_\lambda, \quad (14)$$

where  $T_\lambda$  is the number of terms appearing in (5) by keeping only the variables selected with  $\lambda$  and  $\text{RSS}(\lambda)$  is the residual sum of squares defined as follows:

$$\text{RSS}(\lambda) = \left\| \mathbf{Y} - \widehat{\mathbf{Y}}(\lambda) \right\|_2^2 \quad \text{with} \quad \widehat{\mathbf{Y}}(\lambda) = \sum_{\ell=1}^p \Theta_\ell \widehat{\gamma}_\ell(\lambda), \quad (15)$$

where  $\widehat{\gamma}_\ell(\lambda)$  is defined in (11). Then, the chosen  $\lambda = \lambda_{\text{AIC}}$  is such that:

$$\lambda_{\text{AIC}} = \underset{\lambda \in \Lambda}{\text{argmin}} (\text{AIC}(\lambda)). \quad (16)$$

### 3.2 Results

We apply ABSORBER, LassoNet and RF on noisy observations satisfying (1) with  $f = f_1$  defined as:

$$f_1(x^{(1)}, \dots, x^{(10)}) = 1.8 \cos(x^{(1)}) \sin(x^{(7)} + 1) - 5 \ln(x^{(3)} + 1) - \frac{0.9}{x^{(10)^2 + 1}, \quad (17)$$

$$(x^{(1)}, \dots, x^{(10)}) \in [0, 1]^{10}.$$

and we calculate the percentage of selection defined in the previous section. Here,  $\mathcal{V}_{f_1} = \{1, 3, 7, 10\}$  are the relevant variables to be selected. Note that we use  $K = 1$  evenly spaced knots in the B-splines bases involved in ABSORBER. We refer the reader to Savino and Lévy-Leduc (2024) for a numerical investigation justifying this choice.

Similarly to ABSORBER, the percentage of selection can be computed for each penalization parameter of LassoNet. Concerning the RF method, we convert the percentage of increased mean square error for each variable as the model excludes them one-by-one into a percentage of selection for each of them.

The results obtained for  $n = 700$  and  $n = 2000$  observations with ten random samplings of the set are displayed in Figure 1. Despite being sensitive to the noise level  $\sigma$  in the obser-

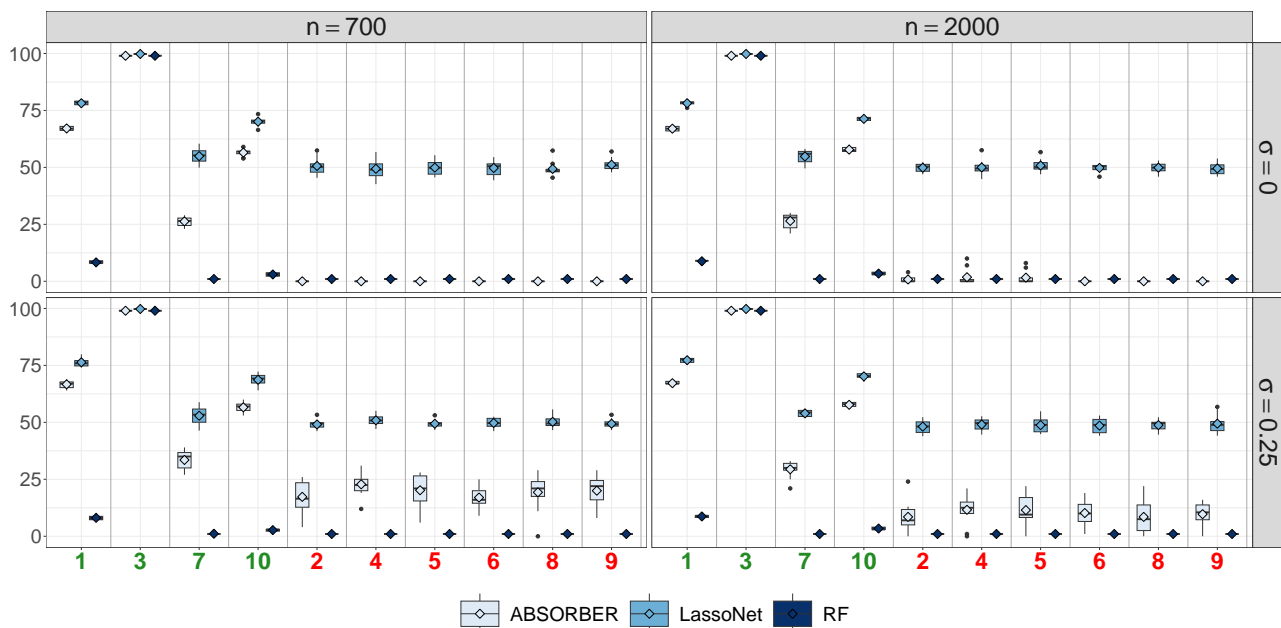


Figure 1: Percentage of selection of each variable of  $f_1$  with three different methods: ABSORBER, LassoNet and RandomForests (RF) with an increasing number of observations  $n$  (left to right) and of the value of  $\sigma$  (top to bottom). 10 random samplings of  $\mathbf{Y}$  were used to obtain these results. The empty diamonds inside the boxplots correspond to the mean value and the plain bullets outside the boxplots are the extreme values.

variations, ABSORBER succeeds in selecting the relevant variables with a satisfying percentage (between 25% and 40%) and allows us to choose a threshold at the visible gap (25 %) to

select the correct variables. In contrast, LassoNet and Random Forests either select irrelevant variables at a high percentage (50%) or fail to select the relevant variables, respectively. Increasing the number of observations allows us to reduce the percentage of selection for irrelevant variables with ABSORBER to 12% while maintaining the minimum percentage of relevant variables up to 30%. The two other methods appear to be unaffected by changes in  $n$ . These conclusions emphasize that our method outperforms those two methods for variable selection while requiring only a few parameters to choose.

We then apply our variable selection method using  $\lambda_{\text{AIC}}$  as described in (16) and we calculate the value of  $\text{TPR}(\lambda_{\text{AIC}})$  and  $\text{FPR}(\lambda_{\text{AIC}})$ . We can observe that increasing  $\sigma$  slightly alters the performance of our method. However, for smaller noise levels ( $\sigma < 0.25$ ) and with  $n \geq 700$  our method enables  $\text{TPR}(\lambda_{\text{AIC}}) = 1$  while maintaining  $\text{FPR}(\lambda_{\text{AIC}})$  to a value smaller than 0.05. This means that almost no irrelevant variables are chosen demonstrating the efficiency of our variable selection procedure.

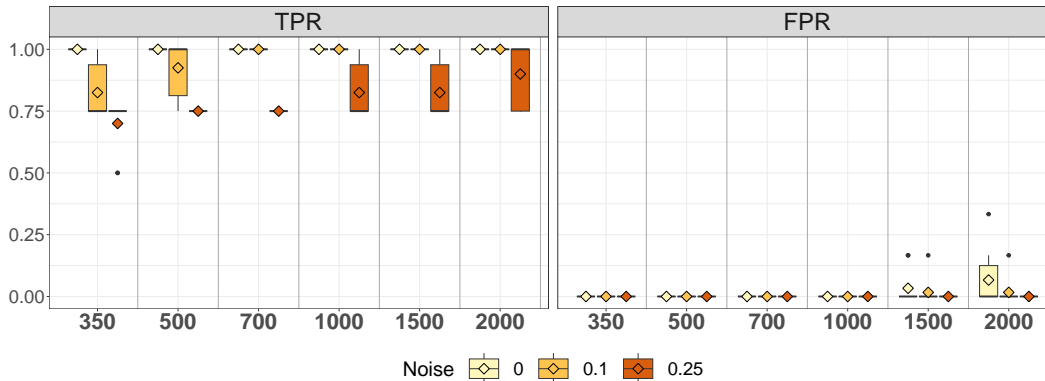


Figure 2:  $\text{TPR}(\lambda)$  and  $\text{FPR}(\lambda)$  values by choosing  $\lambda = \lambda_{\text{AIC}}$  for  $f_1$  with three noise levels in the observation sets. 10 random samplings of  $\mathbf{Y}$  were used to obtain these results. The empty diamonds inside the boxplots correspond to the mean value and the plain bullets outside the boxplots are the extreme values.

## 4 A geochemical application

We propose applying our method ABSORBER to a real geochemical system derived from the calcite dissolution and precipitation study in Kolditz et al. (2012). Our observation sets are here generated using the geochemical solver PHREEQC as in Parkhurst and Appelo (2013). For the purposes of this study, we specifically focus on the calcite precipitation/dissolution such that we can define a function  $f_2$  depending on normalized concentrations and quantities of elements:

$$\begin{aligned} \text{Calcite} &= f_2(\text{C}^*, \text{Ca}^*, \text{K}^*, \text{Cl}^*, \text{Calcite}^*, x^{(6)}, x^{(7)}, x^{(8)}, x^{(9)}, x^{(10)}) \\ &= \tilde{f}_2(\text{C}^*, \text{Ca}^*, \text{Calcite}^*), \end{aligned}$$



where Calcite denotes the amount of produced calcite and  $C^*$ ,  $Ca^*$ ,  $K^*$ ,  $Cl^*$ ,  $Calcite^*$  are the normalized concentrations and quantities initially present of C, Ca, K, Cl and Calcite, respectively. The normalization of each variable is done by taking into account the minimal and the maximal bound of the values so that each variable belongs to  $[0, 1]$ .  $x^{(6)}, x^{(7)}, x^{(8)}, x^{(9)}, x^{(10)}$  are synthetic irrelevant variables obtained through uniform sampling between 0 and 1. Hereafter,  $\mathcal{V}_{f_2} = \{C^*, Ca^*, Calcite^*\}$  and  $\overline{\mathcal{V}}_{f_2} = \{K^*, Cl^*, x^{(6)}, x^{(7)}, x^{(8)}, x^{(9)}, x^{(10)}\}$ . We aim at retrieving  $\mathcal{V}_{f_2}$  by applying our method.

The results for the application of ABSORBER, LassoNet and RF to  $f_2$  are displayed in the left part of Figure 3. Our method consistently succeeds in selecting the relevant variables belonging to  $\mathcal{V}_{f_2}$  (62.5% of selection) while discarding the irrelevant ones (almost 0% of selection) even with a small dataset size  $n = 350$ . In contrast, LassoNet selects all the variables and fails to discriminate the relevant from the irrelevant ones. Random Forests, on the other hand, tends to detect only one variable of  $\mathcal{V}_{f_2}$ . This shows once again that our method outperforms the other two in this geochemical case.

Furthermore, we used the AIC to select the parameter  $\lambda$  and to automatically choose the relevant variables. The corresponding results are shown in right part of Figure 3. This statistical criterion proves to be highly efficient as evidenced by  $TPR(\lambda_{AIC}) = 1$  and  $FPR(\lambda_{AIC}) = 0$ , regardless of  $n$ .

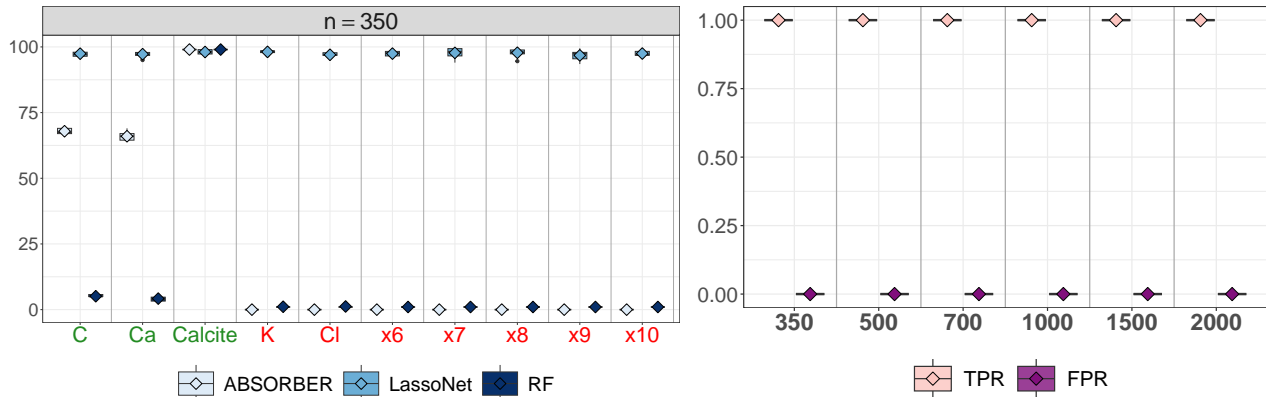


Figure 3: (Left) Percentage of selection of each variable of  $f_2$  with three different methods: ABSORBER, LassoNet and Random Forests (RF) with  $n = 350$  observations. (Right)  $TPR(\lambda)$  and  $FPR(\lambda)$  values by choosing  $\lambda = \lambda_{AIC}$  for  $f_2$  with an increasing number of observations  $n$ . 10 random samplings of  $\mathbf{Y}$  were used to obtain these results. The empty diamonds inside the boxplots correspond to the mean value and the plain bullets outside the boxplots are the extreme values.

## References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Proceedings of the 2nd International Symposium on Information Theory*, pp. 267–281. In B.N.Petrov, F.Csaki (Eds.).

- Asher, M. J., B. F. Croke, A. J. Jakeman, and L. J. Peeters (2015). A review of surrogate models and their application to groundwater modeling. *Water Resources Research* 51(8), 5957–5973.
- Breiman, L. (2001). Random forests. *Machine learning* 45, 5–32.
- De Boor, C. (1978). *A practical guide to splines*, Volume 27. Springer-Verlag New York.
- Hastie, T., R. Tibshirani, and J. Friedman (2009). *The Elements of Statistical Learning: Data mining, inference, and prediction*. New York, NY, USA: Springer New York Inc.
- Jatnieks, J., M. De Lucia, D. Dransch, and M. Sips (2016). Data-driven surrogate model approach for improving the performance of reactive transport simulations. *Energy Procedia* 97, 447–453.
- Kolditz, O., U.-J. Görke, H. Shao, and W. Wang (2012). *Thermo-hydro-mechanical-chemical processes in porous media: benchmarks and examples*, Volume 86. Springer Science & Business Media.
- Lemhadri, I., F. Ruan, L. Abraham, and R. Tibshirani (2021). LassoNet: A neural network with feature sparsity. *The Journal of Machine Learning Research* 22(1), 5633–5661.
- Parkhurst, D. L. and C. Appelo (2013). *Description of input and examples for PHREEQC version 3: a computer program for speciation, batch-reaction, one-dimensional transport, and inverse geochemical calculations*. U.S.G.S. Techniques and Methods.
- Radchenko, P. and G. M. James (2010). Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association* 105(492), 1541–1553.
- Razavi, S., B. A. Tolson, and D. H. Burn (2012). Review of surrogate modeling in water resources. *Water Resources Research* 48(7).
- Rosasco, L., M. Santoro, S. Mosci, A. Verri, and S. Villa (2010). A regularization approach to nonlinear variable selection. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 653–660. JMLR Workshop and Conference Proceedings.
- Savino, M. E. and C. Lévy-Leduc (2024). A novel variable selection method in a nonlinear multivariate model using B-splines with an application to geoscience. hal-04434820.
- Yuan, M. and Y. Lin (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 68(1), 49–67.