

A comparison of eight metamodeling techniques for the simulation of N₂O fluxes and N leaching from corn crops

Nathalie Villa-Vialaneix^{a,b,*}, Marco Follador^c, Marco Ratto^d, Adrian Leip^c

^a*IUT de Perpignan (Dpt STID, Carcassonne), Univ. Perpignan Via Domitia, France*

^b*Institut de Mathématiques de Toulouse, Université de Toulouse, France*

^c*European Commission, Institute for Environment and Sustainability, CCU, Ispra, Italy*

^d*European Commission, Econometrics and Applied Statistics Unit, Ispra, Italy*

Abstract

The environmental costs of intensive farming activities are often underestimated or not traded by the market, even though they play an important role in addressing future society's needs. The estimation of nitrogen (N) dynamics is thus an important issue which demands detailed simulation based methods and their integrated use to correctly represent complex and non-linear interactions into cropping systems. To calculate the N₂O flux and N leaching from European arable lands, a modeling framework has been developed by linking the CAPRI agro-economic dataset with the DNDC-EUROPE bio-geo-chemical model. But, despite the great power of modern calculators, their use at continental scale is often too computationally costly. By comparing several statistical methods this paper aims to design a metamodel able to approximate the expensive code of the detailed modeling approach, devising the best compromise between estimation performance and simulation speed. We describe the use of two parametric (linear) models and six nonparametric approaches: two methods based on splines (ACOSSO and SDR), one method based on kriging (DACE), a neural networks method (multilayer perceptron, MLP), SVM and a bagging method (random forest, RF). This analysis shows that, as long as few data are available to train the model, splines approaches lead to best results, while when the size of training dataset increases, SVM and RF provide faster and more accurate solutions.

*Corresponding author.

Email address: nathalie.villa@math.univ-toulouse.fr (Nathalie Villa-Vialaneix)

Keywords: Metamodeling, splines, SVM, neural network, random forest, N₂O flux, N leaching, agriculture

1. Introduction

The impact of modern agriculture on the environment is well documented (Power, 2010; Tilman et al., 2002; Scherr and Sthapit, 2009; FAO, 2007, 2005; Singh, 2000; Matson et al., 1997). Intensive farming has a high consumption of nitrogen, which is often in-efficiently used, particularly in livestock production systems (Leip et al., 2011b; Webb et al., 2005; Oenema et al., 2007; Chadwick, 2005). This leads to a large surplus of nitrogen which is lost to the environment. Up to 95% of ammonia emission in Europe have their origin in agricultural activities (Kirchmann et al., 1998; Leip et al., 2011a) contributing to eutrophication, loss of biodiversity and health problems. Beside NH₃, nitrate leaching below the soil root zone and entering the groundwater poses a particular problem for the quality of drinking water (van Grinsven et al., 2006). Additionally, agricultural sector is the major source of anthropogenic emissions of N₂O from the soils, mainly as a consequence of the application of mineral fertilizer or manure nitrogen (Del Grosso et al., 2006; Leip et al., 2011c; European Environment Agency, 2010; Leip et al., 2005). N₂O is a potent greenhouse gas (GHG) contributing with each kilogram emitted about 300 times more to global warming than the same mass emitted as CO₂, on the basis of a 100-years time horizon (Intergovernmental Panel on Climate Change, 2007).

Various European legislations attempt to reduce the environmental impact of the agriculture sector, particularly the Nitrates Directive (European Council, 1991) and the Water Framework Directive (European Council, 2000). Initially, however, compliance to these directives was poor (Oenema et al., 2009; European Commission, 2002). Therefore, with the last reform of the Common Agricultural Policy (CAP) in the year 2003 (European Council, 2003), the European Union introduced a compulsory Cross-Compliance (CC) mechanism to improve compliance with 18 environmental, food safety, animal welfare, and animal and plant health standards (Statutory Management Requirements, SMRs) as well as with requirements to maintain farmlands in good agricultural and environmental condition (Good Agricultural and Environment Condition requirements, GAECs), as prerequisite for receiving direct payments (European Union Commission, 2004; European Council,

34 2009; European Union Commission, 2009; Dimopoulos et al., 2007; Jongeneel
35 et al., 2007). The SMRs are based on pre-existing EU Directives and Reg-
36 ulations such as Nitrate Directives. The GAECs focus on soil erosion, soil
37 organic matter, soil structure and a minimum level of maintenance; for each
38 of these issues a number of standards are listed (Alliance Environnement,
39 2007).

40 It remains nevertheless a challenge to monitor compliance and to assess
41 the impact of the cross-compliance legislations not only on the environment,
42 but also on animal welfare, farmer's income, production levels etc. In or-
43 der to help with this task, the EU-project Cross-Compliance Assessment
44 Tool (CCAT) developed a simulation platform to provide scientifically sound
45 and regionally differentiated responses to various farming scenarios (Elbersen
46 et al., 2010; Jongeneel et al., 2007).

47 CCAT integrates complementary models to assess changes in organic car-
48 bon and nitrogen fluxes from soils (De Vries et al., 2008). Carbon and ni-
49 trogen turnover are very complex processes, characterized by a high spatial
50 variability and a strong dependence on environmental factors such as mete-
51 orological conditions and soils (Shaffer and Ma, 2001; Zhang et al., 2002).
52 Quantification of fluxes, and specifically a meaningful quantification of the
53 response to mitigation measures at the regional level requires the simulation
54 of farm management and the soil/plant/atmosphere continuum at the high-
55 est possible resolution (Anderson et al., 2003; Leip et al., 2011c). For the
56 simulation of N₂O fluxes and N-leaching, the process-based biogeochemistry
57 model DNDC-EUROPE (Leip et al., 2008; Li et al., 1992; Li, 2000) was used.
58 As DNDC-EUROPE is a complex model imposing high computational costs,
59 the time needed to obtain simulation results in large scale applications (such
60 as the European scale) can be restrictive. In particular, the direct use of the
61 deterministic model is prohibited to extract efficiently estimations of the evo-
62 lution of N₂O fluxes and N-leaching under changing conditions. Hence, there
63 is a need for a second level of abstraction, modeling the DNDC-EUROPE
64 model itself, which is called a *meta-model* (see Section 2 for a more specific
65 definition of the concept of metamodeling). Metamodels are defined from a
66 limited number of deterministic simulations for specific applications and/or
67 scenario and allow to obtain fast estimations.

68 This issue is a topic of high interest that has previously been tackled in
69 several papers: among others, (Bouzaher et al., 1993) develop a parametric
70 model, including spatial dependency, to model water pollution. (Krysanova
71 and Haberlandt, 2002; Haberlandt et al., 2002) describe a two-steps approach

72 to address the issue of N leaching and water pollution: they use a process-
73 based model followed by a location of the results with a fuzzy rule. More
74 recently, (Pineros Garcet et al., 2006) compare RBF neural networks with
75 kriging modeling to build a metamodel for a deterministic N leaching model
76 called WAVE (Vanclooster et al., 1996). The present article compares in
77 detail different modeling tools in order to select the most reliable one to
78 meta-model the DNDC-EUROPE tasks in the CCAT project Follador and
79 Leip (2009). This study differs from the work of Vanclooster et al. (1996)
80 because of the adopted European scale and of the analysis of 8 meta-modeling
81 approaches (also including a kriging and a neural network method). The
82 comparison has been based on the evaluation of meta-model performances,
83 in terms of accuracy and computational costs, with different sizes of the
84 training dataset.

85 The rest of the paper is organized as follows: Section 2 introduces the
86 general principles and advantages of using a meta-model; Section 3 reviews
87 in details the different types of metamodels compared in this study; Sec-
88 tion 4 explains the Design Of the Experiments (DOE) and show the results
89 of the comparison, highlighting how the availability of the training data can
90 play an important role in the selection of the best type and form of the
91 approximation. The supplementary material of this paper can be found at:
92 <http://afoludata.jrc.ec.europa.eu/index.php/dataset/detail/232>.

93 **2. From model to metamodel**

94 A model is a simplified representation (abstraction) of reality developed
95 for a specific goal; it may be deterministic or probabilistic. An integrated
96 use of simulation based models is necessary to approximate our perception
97 of complex and nonlinear interactions existing in human-natural systems by
98 means of mathematical input-output (I/O) relationships. Despite the con-
99 tinuous increase of computer performance, the development of large simula-
100 tion platforms remains often prohibited because of computational needs and
101 parametrization constraints. More precisely, every model in a simulation
102 platform such as DNDC-EUROPE, is characterized by several parameters,
103 whose near-optimum set is defined during the calibration. A constraint ap-
104 plies restrictions to the kind of data that the model can use or to specific
105 boundary conditions. The flux of I/O in the simulation platform can thus
106 be impeded by the type of data/boundaries that constraints allow - or not
107 allow - for the models at hand.

108 The use of this kind of simulation platform is therefore not recommended
 109 for all the applications which require many runs, such as sensitivity analysis
 110 or what-if studies. To overcome this limit, the process of abstraction can
 111 be applied to the model itself, obtaining a model of the model (2nd level of
 112 abstraction from reality) called meta-model (Blanning, 1975; Kleijnen, 1975;
 113 Sacks et al., 1989; van Gighc, 1991; Santner et al., 2003). A metamodel is
 114 an approximation of detailed model I/O transformations, built through a
 115 moderate number of computer experiments.

116 Replacing a detailed model with a metamodel generally brings some pay-
 117 offs (Britz and Leip, 2009; Simpson et al., 2001):

- 118 • easier integration into other processes and simulation platforms;
- 119 • faster execution and reduced storage needs to estimate one specific
 120 output;
- 121 • easier applicability across different spatial and/or temporal scales and
 122 site-specific calibrations, as long as data corresponding to the new sys-
 123 tem parametrization are available.

124 As a consequence, a higher number of simulation runs become possible: using
 125 its interpolatory action makes a thorough sensitivity analysis more convenient
 126 and leads to a better understanding of I/O relationships. Also it offers usually
 127 a higher flexibility and can quickly be adapted to achieve a wide range of
 128 goals (prediction, optimization, exploration, validation). However, despites
 129 these advantages, they suffer from a few drawbacks: internal variables or
 130 outputs not originally considered can not be inspected and the prediction
 131 for input regimes outside the training/test set is impossible. Hence, a good
 132 metamodeling methodology should be able to provide fast predictions. But,
 133 considering that limitations, it also must have a low computational cost to be
 134 able to build a new metamodel from a new data set including new variables
 135 and/or a different range for these input variables.

Let (\mathbf{X}, \mathbf{y}) be the dataset consisting of N row vectors of input/output
 pairs (\mathbf{x}_i, y_i) , where $\mathbf{x}_i = (x_i^1, \dots, x_i^d)^T \in \mathbb{R}^d$ ($i = 1, \dots, N$) are the model
 input and $y_i \in \mathbb{R}$ ($i = 1, \dots, N$) are the model responses for N experimental
 runs of the simulation platform. The mathematical representation of I/O
 relationships described by the detailed model can be written as

$$y_i = f(\mathbf{x}_i) \quad i = 1, \dots, N \quad (1)$$

which corresponds to a first abstraction from the real system. From the values of \mathbf{X} and \mathbf{y} , also called *training set*, f is approximated by a function $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$, called metamodel, whose responses can be written as

$$\hat{y}_i = \hat{f}(\mathbf{x}_i).$$

136 and that correspond to a second abstraction from the reality. In this second
137 abstraction, some of the input variables of Eq. (1) might not be useful and
138 one of the issue of metamodeling can be to find the smallest subset of input
139 variables relevant to achieve a good approximation of model (1).

140 Finally, the differences between the real system and the metamodel re-
141 sponse, will be the sum of two approximations (Simpson et al., 2001): the
142 first one introduced by the detailed model (1st abstraction) and the second
143 one due to metamodeling (2nd abstraction). Of course, the validity and ac-
144 curacy of a metamodel are conditioned by the validity of the original model:
145 in the following, it is then supposed that the 1st level of abstraction induces
146 a small error compared to reality. Then, in this paper, we only focus on
147 the second error, $|\hat{y}_i - y_i|$, to assess the performance of different metamodels
148 vs. the detailed DNDC-EUROPE model in order to select the best statisti-
149 cal approach to approximate the complex bio-geo-chemical model at a lower
150 computational cost. Defining a correct metamodeling strategy is very impor-
151 tant to provide an adequate fitting to the model, as suggested by (Kleijnen
152 and Sargent, 2000; Meckesheimer et al., 2002).

153 Recent work, such as (Forrester and Keane, 2009; Wang and Shan, 2007),
154 review the most widely used metamodeling methods: splines based methods
155 (e.g., MARS, kriging...) (Wahba, 1990; Friedman, 1991; Cressie, 1990), neu-
156 ral networks (Bishop, 1995), kernel methods (SVM, SVR...) (Vapnik, 1998;
157 Christmann and Steinwart, 2007), Gaussian Process such as GEM (Kennedy
158 and O'Hagan, 2001), among others. Some of these metamodeling strategies
159 were selected and others added to be compared in this paper. The compar-
160 ison is made on a specific case study related to N leaching and N₂O fluxes
161 prediction which is described in Section 4. The next section briefly describes
162 each of the metamodels compared in this paper.

163 **3. Review of the selected metamodels**

164 Several methods were developed and compared to assess their perfor-
165 mance according to increasing dataset sizes. We provide a brief description

166 of the approaches studied in this paper: two linear models (Section 3.1) and
 167 six nonparametric methods (two based on splines, in Sections 3.2.1 and 3.2.2,
 168 one based on a kriging approach, in Section 3.2.3, which is known to be ef-
 169 ficient when analyzing computer experiments, a neural network method, in
 170 Section 3.2.4, SVM, in Section 3.2.5 and random forest, in Section 3.2.6).

171 3.1. Linear methods

The easiest way to handle the estimation of the model given in Eq. (1) is to suppose that f has a simple parametric form. For example, the *linear model* supposes that $f(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{x} + \beta_0$ where $\boldsymbol{\beta} \in \mathbb{R}^d$ is a vector and β_0 is a real number, both of them have to be estimated from the observations $((\mathbf{x}_i, y_i))_i$. An estimate is given by minimizing the sum of the square errors

$$\sum_{i=1}^N (y_i - (\boldsymbol{\beta}^T \mathbf{x}_i + \beta_0))^2$$

172 which leads to $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ and $\hat{\beta}_0 = \bar{y} - \hat{\boldsymbol{\beta}}^T \bar{\mathbf{X}}$ with $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$
 173 and $\bar{\mathbf{X}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$.

174 In this paper two linear models were used:

- 175 • in the first one, the explanatory variables were the 11 inputs described
 176 in Section 4.2. This model is referred as “LM1”;
- 177 • the second one has been developed starting from the work of (Britz
 178 and Leip, 2009), that includes the 11 inputs of Section 4.2 but also
 179 their non linear transformations (square, square root, logarithm) and
 180 interaction components. A total of 120 coefficients were involved in
 181 this approach which is denoted by “LM2”. Including transformations
 182 and combinations of the 11 inputs has been designed in an attempt to
 183 better model a possible nonlinear phenomenon of the original model.

184 In the second case, due to the large number of explanatory variables, the
 185 model can be over-specified, especially if the training set is small. Actually,
 186 if the dimensionality of the matrix of explanatory variables, \mathbf{X} , has a large
 187 dimension, $\mathbf{X}^T \mathbf{X}$ can be not invertible or ill-conditioned (leading to numerical
 188 instability). Hence, a stepwise selection based on the AIC criterion (Akaike,
 189 1974) has been used to select an optimal subset of explanatory variables
 190 during the training step in order to obtain an accurate solution having a
 191 small number of parameters. This has been performed by using the `stepAIC`
 192 function of the **R** package `MASS`.

193 *3.2. Nonparametric methods*

194 In many modeling problems, linear methods are not enough to catch the
 195 complexity of the phenomenon which is, *per se*, nonlinear. In these situations,
 196 nonparametric are often more suited to obtain accurate approximations of
 197 the phenomenon under study. In this section, six nonparametric approaches
 198 are described: they are compared in Section 4 to model N₂O fluxes and N
 199 leaching.

200 *3.2.1. ACOSSO*

Among nonparametric estimation approach, the smoothing splines (Wahba, 1990; Gu, 2002) is one of the most famous and widely used. Recently, (Storlie et al., 2011) presented the ACOSSO, an adaptive approach based on the COSSO method (Lin and Zhang, 2006) which is in the same line as smoothing splines: it is described as “a new regularization method for simultaneous model fitting and variable selection in nonparametric regression models in the framework of smoothing spline ANOVA”. This method penalizes the sum of component norms, instead of the squared norm employed in the traditional smoothing spline method. More precisely, in splines meta-modeling, it is useful to consider the ANOVA decomposition of f into terms of increasing dimensionality:

$$f(\mathbf{x}) = f(x^1, x^2, \dots, x^d) = f_0 + \sum_j f^{(j)} + \sum_{k>j} f^{(jk)} + \dots + f^{(12\dots d)} \quad (2)$$

201 where x^j is the j -th explanatory variable and where each term is a function
 202 only of the factors in its index, i.e. $f^{(j)} = f(x^j)$, $f^{(jk)} = f(x^j, x^k)$ and
 203 so on. The terms $f^{(j)}$ represent the additive part of the model f , while
 204 all higher order terms $f^{(jk)} \dots f^{(12\dots d)}$ are denoted as “interactions”. The
 205 simplest example of smoothing spline ANOVA model is the additive model
 206 where only $(f^{(j)})_{j=0,\dots,d}$ are used.

To estimate f , we make the usual assumption that $f \in \mathcal{H}$, where \mathcal{H} is a RKHS (Reproducing Kernel Hilbert Space) (Berlinet and Thomas-Agnan, 2004). The space \mathcal{H} can be written as an orthogonal decomposition $\mathcal{H} = \{1\} \oplus \{\bigoplus_{j=1}^q \mathcal{H}_j\}$, where each \mathcal{H}_j is itself a RKHS, \oplus is the direct sum of Hilbert spaces and $j = 1, \dots, q$ spans ANOVA terms of various orders. Typically q includes the main effects plus relevant interaction terms. f is then estimated by \hat{f} that minimizes a criterion being a trade-off between accuracy to the data (empirical mean squared error) and a penalty which

aims at minimizing each ANOVA term:

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(\mathbf{x}_i))^2 + \lambda_0 \sum_{j=1}^q \frac{1}{\theta_j} \|P^j \hat{f}\|_{\mathcal{H}}^2 \quad (3)$$

207 where $P^j \hat{f}$ is the orthogonal projection of \hat{f} onto \mathcal{H}_j and the q -dimensional
 208 vector θ_j of smoothing parameters needs to be tuned somehow, in such a way
 209 that each ANOVA component has the most appropriate degree of smooth-
 210 ness.

211 This statistical estimation problem requires the tuning of the d hyper-
 212 parameters θ_j (λ_0/θ_j are also denoted as smoothing parameters). Various
 213 ways of doing that are available in the literature, by applying generalized
 214 cross-validation (GCV), generalized maximum likelihood procedures (GML)
 215 and so on (Wahba, 1990; Gu, 2002). But, in Eq. (3), q is often large and
 216 the tuning of all θ_j is a formidable problem, implying that in practice the
 217 problem is simplified by setting θ_j to 1 for any j and only λ_0 is tuned. This
 218 simplification, however, strongly limits the flexibility of the smoothing spline
 219 model, possibly leading to poor estimates of the ANOVA components.

Problem (3) also poses the issue of selection of \mathcal{H}_j terms: this is tackled
 rather effectively within the COSSO/ACOSSO framework. The COSSO (Lin
 and Zhang, 2006) penalizes the sum of norms, using a LASSO type penalty
 (Tibshirani, 1996) for the ANOVA model: LASSO penalties are L_1 penalties
 that lead to sparse parameters (i.e., parameters whose coordinates are all
 equal to zero except for a few ones). Hence, using this kind of penalties
 allows us to automatically select the most informative predictor terms \mathcal{H}_j
 with an estimate of \hat{f} that minimizes

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(\mathbf{x}_i))^2 + \lambda \sum_{j=1}^Q \|P^j \hat{f}\|_{\mathcal{H}} \quad (4)$$

220 using a single smoothing parameter λ , and where Q includes *all* ANOVA
 221 terms to be potentially included in \hat{f} , e.g. with a truncation at 2^{nd} or 3^{rd}
 222 order interactions.

It can be shown that the COSSO estimate is also the minimizer of

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(\mathbf{x}_i))^2 + \sum_{j=1}^Q \frac{1}{\theta_j} \|P^j \hat{f}\|_{\mathcal{H}}^2 \quad (5)$$

223 subject to $\sum_{j=1}^Q 1/\theta_j < M$ (where there is a 1-1 mapping between M and
 224 λ). So we can think of the COSSO penalty as the traditional smoothing
 225 spline penalty plus a penalty on the Q smoothing parameters used for each
 226 component. This can also be framed into a linear-quadratic problem, i.e. a
 227 quadratic objective (5) plus a linear constraint on $1/\theta_j$. The LASSO type
 228 penalty has the effect of setting some of the functional components (\mathcal{H}_j 's)
 229 equal to zero (e.g. some variables x^j and some interactions (x^j, x^k) are not
 230 included in the expression of \hat{f}). Thus it “automatically” selects the appropriate
 231 subset q of terms out of the Q “candidates”. The key property of
 232 COSSO is that with one single smoothing parameter (λ or M) it provides
 233 estimates of all θ_j parameters in one shot: therefore it improves considerably
 234 the simplified problem (3) by setting $\theta_j = 1$ (still with one single smoothing
 235 parameter λ_0) and is much more computationally efficient than the full problem
 236 (3) with optimized θ_j 's. An additional improvement from the COSSO
 237 is that the single smoothing parameter λ can be tuned to minimize the BIC
 238 (Bayesian Information Criterion) (Schwarz, 1978), thus allowing to target
 239 the most appropriate degree of parsimony of the metamodel. This is done
 240 by a simple grid-search algorithm as follows (see (Lin and Zhang, 2006) for
 241 details):

- 242 1. for each trial λ value, the COSSO estimate provides the corresponding
- 243 values for θ_j and subsequently its BIC;
- 244 2. the grid-search algorithm will provide the $\hat{\lambda}$ with the smallest BIC.

The adaptive COSSO (ACOSSO) of (Storlie et al., 2011) is an improvement of the COSSO method: in ACOSSO, $\hat{f} \in \mathcal{H}$ minimizes

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(\mathbf{x}_i))^2 + \lambda \sum_{j=1}^q w_j \|P^j \hat{f}\|_{\mathcal{H}} \quad (6)$$

245 where $0 < w_j \leq \infty$ are weights that depend on an initial estimate, $\hat{f}^{(0)}$,
 246 of f , either using (3) with $\theta_j = 1$ or the COSSO estimate (4). The
 247 adaptive weights are obtained as $w_j = \|P^j \hat{f}^{(0)}\|_{L_2}^{-\gamma}$, typically with $\gamma = 2$
 248 and the L_2 norm $\|P^j \hat{f}^{(0)}\|_{L_2} = (\int (P^j \hat{f}^{(0)}(\mathbf{x}))^2 d\mathbf{x})^{1/2}$. The use of adaptive
 249 weights improves the predictive capability of ANOVA models with respect
 250 to the COSSO case: in fact it allows for more flexibility in estimating
 251 important functional components while giving a heavier penalty to unimportant
 252 functional components. The **R** scripts for ACOSSO can be found

253 at <http://www.stat.lanl.gov/staff/CurtStorlie/index.html>. In the
254 present paper we used a MATLAB translation of such R script. The algo-
255 rithm for tuning the hyper-parameters is then modified as follows:

- 256 1. an initial estimate of the ANOVA model $\hat{f}^{(0)}$ is obtained either using
257 (3) with $\theta_j = 1$ or the COSSO estimate (4);
- 258 2. given this trial ANOVA model $\hat{f}^{(0)}$, the weights are computed as $w_j =$
259 $\|P^j \hat{f}^{(0)}\|_{L_2}^{-\gamma}$;
- 260 3. given w_j and for each trial λ value, the ACOSSO estimate (6) provides
261 the corresponding values for θ_j and subsequently its BIC;
- 262 4. the grid-search algorithm will provide the $\hat{\lambda}$ with the smallest BIC.

263 3.2.2. SDR-ACOSSO

264 In a “parallel” stream of research with respect to COSSO-ACOSSO, us-
265 ing the *state-dependent parameter regression* (SDR) approach of (Young,
266 2001), (Ratto et al., 2007) have developed a non-parametric approach, very
267 similar to smoothing splines and kernel regression methods, based on recur-
268 sive filtering and smoothing estimation (the Kalman filter combined with
269 “fixed interval smoothing”). Such a recursive least-squares implementa-
270 tion has some key characteristics: (a) it is combined with optimal maxi-
271 mum likelihood estimation, thus allowing for an estimation of the smooth-
272 ing hyper-parameters based on the estimation of a quality criterion rather
273 than on cross-validation and (b) it provides greater flexibility in adapt-
274 ing to local discontinuities, heavy non-linearity and heteroscedastic error
275 terms. Recently, (Ratto and Pagano, 2010) proposed a unified approach
276 to smoothing spline ANOVA models that combines the best of SDR and
277 ACOSSO: the use of the recursive algorithms in particular can be very ef-
278 fective in *identifying* the important functional components and in providing
279 good estimates of the weights w_j to be used in (6), adding valuable infor-
280 mation in the ACOSSO framework and allowing in many cases to improving
281 ACOSSO performance. The Matlab script for this method can be found at
282 http://eemc.jrc.ec.europa.eu/Software-SS_ANOVA_R.htm.

283 We summarize here the key features of Young’s recursive algorithms of
284 SDR, by considering the case of $d = 1$ and $f(x^1) = f^{(1)}(x^1) + e$, with $e \sim$
285 $N(0, \sigma^2)$. To do so, we rewrite the smoothing problem as $y_i = s_i^1 + e_i$,
286 where $i = 1, \dots, N$ and s_i^1 is the estimate of $f^{(1)}(x_i^1)$. To make the recursive
287 approach meaningful, the MC sample needs to be sorted in ascending order

288 with respect to x^1 : i.e. i and $i - 1$ subscripts are adjacent elements under
 289 such ordering, implying $x_1^1 < x_2^1 < \dots < x_i^1 < \dots < x_N^1$.

To recursively estimate the s_i^1 in SDR it is necessary to characterize it in some stochastic manner, borrowing from non-stationary time series processes (Young and Ng, 1989; Ng and Young, 1990). In the present context, the integrated random walk (IRW) process provides the same smoothing properties of a cubic spline, in the overall State-Space formulation:

$$\begin{aligned} \text{Observation Equation: } y_i &= s_i^1 + e_i \\ \text{State Equations: } s_i^1 &= s_{i-1}^1 + d_{i-1}^1 \\ d_i^1 &= d_{i-1}^1 + \eta_i^1 \end{aligned} \quad (7)$$

290 where d_i^1 is the ‘‘slope’’ of s_i^1 , $\eta_i^1 \sim N(0, \sigma_{\eta_1}^2)$ and η_i^1 (‘‘system disturbance’’
 291 in systems terminology) is assumed to be independent of the ‘‘observation
 292 noise’’ $e_i \sim N(0, \sigma^2)$.

293 Given the ascending ordering of the MC sample, s_i^1 can be estimated by
 294 using the recursive Kalman Filter (KF) and the associated recursive Fixed
 295 Interval Smoothing (FIS) algorithm (see e.g. (Kalman, 1960; Young, 1999)
 296 for details). First, it is necessary to optimize the hyper-parameter associated
 297 with the state space model (7), namely the Noise Variance Ratio (NVR),
 298 where $\text{NVR}_1 = \sigma_{\eta_1}^2 / \sigma^2$. This is accomplished by maximum likelihood opti-
 299 mization (ML) using prediction error decomposition (Schweppe, 1965). The
 300 NVR plays the inverse role of a smoothing parameter: the smaller the NVR,
 301 the smoother the estimate of s_i^1 . Given the NVR, the FIS algorithm then
 302 yields an estimate $\hat{s}_{i|N}^1$ of s_i^1 at each data sample and it can be seen that the
 303 $\hat{s}_{i|N}^1$ from the IRW process is the equivalent of $\hat{f}^{(1)}(x_i^1)$ in the cubic smooth-
 304 ing spline model. At the same time, the recursive procedures provide, in a
 305 natural way, standard errors of the estimated $\hat{s}_{i|N}^1$, that allow for the test-
 306 ing of their relative significance. Finally, it can be easily verified (Ratto and
 307 Pagano, 2010) that by setting $\lambda/\theta_1 = 1/(\text{NVR}_1 \cdot N^4)$, and with evenly spaced
 308 x_i^1 values, the $\hat{f}^{(1)}(x_i^1)$ estimate in the cubic smoothing spline model equals
 309 the $\hat{s}_{i|N}^1$ estimate from the IRW process.

310 The most interesting aspect of the SDR approach is that it is not limited
 311 to the univariate case, but can be effectively extended to the most relevant
 312 multivariate one. In the general additive case, for example, the recursive
 313 procedure needs to be applied, in turn, for each term $f^{(j)}(x_i^j) = \hat{s}_{i|N}^j$, requiring
 314 a different sorting strategy for each $\hat{s}_{i|N}^j$. Hence the ‘‘back-fitting’’ procedure
 315 is applied, as described in (Young, 2000) and (Young, 2001). This procedure

316 provides both ML estimates of all NVR_j 's and the smoothed estimates of the
 317 additive terms $\hat{s}_{i|N}^j$. So, the estimated NVR_j 's can be converted into λ_0/θ_j
 318 values using $\lambda_0/\theta_j = 1/(\text{NVR}_j \cdot N^4)$, allowing us to put the additive model
 319 into the standard cubic spline form.

320 In the SDR context, (Ratto and Pagano, 2010) formalized an interaction
 321 function as the product of two states $s_1 \cdot s_2$, each of them characterized by
 322 an IRW stochastic process. Hence the estimation of a single interaction term
 323 $f(\mathbf{x}_i) = f^{(12)}(x_i^1, x_i^2) + e_i$ is expressed as:

$$\begin{aligned} \text{Observation Equation:} \quad y_i^* &= s_{1,i}^I \cdot s_{2,i}^I + e_i \\ \text{State Equations: } (j = 1, 2) \quad s_{j,i}^I &= s_{j,i-1}^I + d_{j,i-1}^I \\ &d_{j,i}^I = d_{j,i-1}^I + \eta_{j,i}^I \end{aligned} \quad (8)$$

324 where y^* is the model output after having taken out the main effects, $I =$
 325 $1, 2$ is the multi-index denoting the interaction term under estimation and
 326 $\eta_{j,i}^I \sim N(0, \sigma_{\eta_j^I}^2)$. The two terms $s_{j,i}^I$ are estimated iteratively by running the
 327 recursive procedure in turn.

328 The SDR recursive algorithms are usually very efficient in identifying
 329 in the most appropriate way each ANOVA component individually, hence
 330 (Ratto and Pagano, 2010) proposed to exploit this in the ACOSSO framework
 331 as follows.

332 We define $\mathcal{K}_{\langle j \rangle}$ to be the reproducing kernel (r.k.) of an additive term \mathcal{F}_j
 333 of the ANOVA decomposition of the space \mathcal{F} . In the cubic spline case, this
 334 is constructed as the sum of two terms $\mathcal{K}_{\langle j \rangle} = \mathcal{K}_{01\langle j \rangle} \oplus \mathcal{K}_{1\langle j \rangle}$ where $\mathcal{K}_{01\langle j \rangle}$ is
 335 the r.k. of the parametric (linear) part and $\mathcal{K}_{1\langle j \rangle}$ is the r.k. of the purely
 336 non-parametric part. The second order interaction terms are constructed as
 337 the tensor product of the first order terms, for a total of four elements, i.e.

$$\begin{aligned} \mathcal{K}_{\langle i,j \rangle} &= (\mathcal{K}_{01\langle i \rangle} \oplus \mathcal{K}_{1\langle i \rangle}) \otimes (\mathcal{K}_{01\langle j \rangle} \oplus \mathcal{K}_{1\langle j \rangle}) \\ &= (\mathcal{K}_{01\langle i \rangle} \otimes \mathcal{K}_{01\langle j \rangle}) \oplus (\mathcal{K}_{01\langle i \rangle} \otimes \mathcal{K}_{1\langle j \rangle}) \oplus (\mathcal{K}_{1\langle i \rangle} \otimes \mathcal{K}_{01\langle j \rangle}) \oplus (\mathcal{K}_{1\langle i \rangle} \otimes \mathcal{K}_{1\langle j \rangle}) \end{aligned} \quad (9)$$

This suggested that a natural use of the SDR identification and estimation
 in the ACOSSO framework is to apply specific weights to each element of
 the r.k. $\mathcal{K}_{\langle \cdot, \cdot \rangle}$ in (9). In particular the weights are the L_2 norms of each of
 the four elements estimated in (8):

$$\hat{s}_i^I \cdot \hat{s}_j^I = \hat{s}_{01\langle i \rangle}^I \hat{s}_{01\langle j \rangle}^I + \hat{s}_{01\langle i \rangle}^I \hat{s}_{1\langle j \rangle}^I + \hat{s}_{1\langle i \rangle}^I \hat{s}_{01\langle j \rangle}^I + \hat{s}_{1\langle i \rangle}^I \hat{s}_{1\langle j \rangle}^I, \quad (10)$$

338 As shown in (Ratto and Pagano, 2010), this choice can lead to a significant
 339 improvement in the accuracy of ANOVA models with respect to the original

340 ACOSSO approach. Overall, the algorithm for tuning the hyper-parameters
 341 in the combined SDR-ACOSSO reads:

- 342 1. the recursive SDR algorithm is applied to get an initial estimate of each
 343 ANOVA term in turn (back-fitting algorithm);
- 344 2. the weights are computed as the L_2 norms of the parametric and non-
 345 parametric parts of the cubic splines estimates;
- 346 3. given w_j and for each trial λ value, the ACOSSO estimate (6) provides
 347 the corresponding values for θ_j and subsequently its BIC;
- 348 4. the grid-search algorithm will provide the $\hat{\lambda}$ with the smallest BIC.

349 3.2.3. Kriging metamodel: DACE

350 DACE (Lophaven et al., 2002) is a Matlab toolbox used to construct
 351 kriging approximation models on the basis of data coming from computer
 352 experiments. Once we have this approximate model, we can use it as a meta-
 353 model (emulator, surrogate model). We briefly highlight the main features
 354 of DACE. The kriging model can be expressed as a regression

$$\hat{f}(\mathbf{x}) = \beta_1\phi^1(\mathbf{x}) + \dots + \beta_q\phi^q(\mathbf{x}) + \zeta(\mathbf{x}) \quad (11)$$

where $\phi^j, j = 1, \dots, q$ are deterministic regression terms (constant, linear, quadratic, etc.), β_j are the related regression coefficients and ζ is a zero mean random process whose variance depends on the process variance ω^2 and on the correlation $\mathcal{R}(v, w)$ between $\zeta(v)$ and $\zeta(w)$. In kriging, correlation functions are typically used, defined as:

$$\mathcal{R}(\theta, v - w) = \prod_{j=1:d} \mathcal{R}_j(\theta_j, w_j - v_j).$$

In particular, for the generalized exponential correlation function, used in the present paper, one has

$$\mathcal{R}_j(\theta_j, w_j - v_j) = \exp(-\theta_j|w_j - v_j|^{\theta_{d+1}})$$

Then, we can define \mathbf{R} as the correlation matrix at the training points (i.e., the matrix with coordinates $r_{i,j} = \mathcal{R}(\theta, \mathbf{x}_i, \mathbf{x}_j)$) and the vector $\mathbf{r}_\mathbf{x} = [\mathcal{R}(\theta, \mathbf{x}_1, \mathbf{x}), \dots, \mathcal{R}(\theta, \mathbf{x}_N, \mathbf{x})]$, \mathbf{x} being an untried point. Similarly, we define the vector $\phi_\mathbf{x} = [\phi^1(\mathbf{x}) \dots \phi^q(\mathbf{x})]^T$ and the matrix $\Phi = [\phi_{\mathbf{x}_1} \dots \phi_{\mathbf{x}_N}]^T$ (i.e.,

Φ stacks in matrix form all values of $\phi_{\mathbf{x}}$ at the training points). Then, considering the linear regression problem $\Phi\boldsymbol{\beta} \approx \mathbf{y}$ coming from Eq. (11), with parameter $\boldsymbol{\beta} = [\beta_1, \dots, \beta_q]^T \in \mathbb{R}^q$, the GLS solution is given by:

$$\boldsymbol{\beta}^* = (\Phi^T \mathbf{R}^{-1} \Phi)^{-1} \Phi^T \mathbf{R}^{-1} \mathbf{y}$$

which gives the predictor at untried \mathbf{x}

$$\hat{f}(\mathbf{x}) = \phi_{\mathbf{x}}^T \boldsymbol{\beta}^* + \mathbf{r}_{\mathbf{x}}^T \boldsymbol{\gamma}^*,$$

355 where $\boldsymbol{\gamma}^*$ is the N -dimensional vector computed as $\boldsymbol{\gamma}^* = \mathbf{R}^{-1}(\mathbf{y} - \Phi\boldsymbol{\beta}^*)$.

356 The proper estimation of the kriging metamodel requires, of course, to
 357 optimize the hyper-parameters θ in the correlation function: this is typi-
 358 cally performed by maximum likelihood. It is easy to check that the kriging
 359 predictor *interpolates* \mathbf{x}_j , if the latter is a training point.

360 It seems useful to underline that one major difference between DACE and
 361 ANOVA smoothing is the absence of any “observation error” in (11). This is
 362 a natural choice when analyzing computer experiments and it aims to exploit
 363 the “zero-uncertainty” feature of this kind of data. This, in principle, makes
 364 the estimation of kriging metamodels very efficient, as confirmed by the many
 365 successful applications described in literature and justifies the great success
 366 of this kind of metamodels among practitioners. It also seems interesting to
 367 mention the so-called “nugget” effect, which is also used in the kriging liter-
 368 ature (Montès, 1994; Kleijnen, 2009). This is nothing other than a “small”
 369 error term in (11) and it often reduces some numerical problems encountered
 370 in the estimation of the kriging metamodels to the form of (11). The addi-
 371 tion of a nugget term leads to kriging metamodels that smooth, rather than
 372 interpolate, making them more similar to other metamodels presented here.

373 3.2.4. Multilayer perceptron

374 “Neural network” is a general name for statistical methods dedicated to
 375 data mining. They comprise of a combination of simple computational el-
 376 ements (neurons or nodes) densely interconnected through synapses. The
 377 number and organization of the neurons and synapses define the network
 378 topology. One of the most popular neural network class is the “multilayer
 379 perceptrons” (MLP) commonly used to solve a wide range of classification
 380 and regression problems. In particular, MLP are known to be able to approx-
 381 imate any (smooth enough) complex function (Hornik, 1991). Perceptrons

382 were introduced at the end of the 50s by Rosenblatt but they started be-
 383 coming very appealing more recently thanks to the soaring computational
 384 capacities of computers. The works of (Ripley, 1994) and (Bishop, 1995)
 385 provide a general description of these methods and their properties.

386 For the experiments presented in Section 4.3, one-hidden-layer percep-
 387 trons were used. They can be expressed as a function of the form

$$f_{\mathbf{w}} : \mathbf{x} \in \mathbb{R}^p \rightarrow g_1 \left(\sum_{i=1}^Q w_i^{(2)} g_2 \left(\mathbf{x}^T \mathbf{w}_i^{(1)} + w_i^{(0)} \right) + w_0^{(2)} \right)$$

388 where:

- 389 • $\mathbf{w} := \left[(w_i^{(0)})_i, ((\mathbf{w}_i^{(1)})^T)_i, w_0^{(2)}, (w_i^{(2)})_i \right]^T$ are parameters of the model,
 390 called *weights*. They have to be learned in $(\mathbb{R})^Q \times (\mathbb{R}^p)^Q \times \mathbb{R} \times (\mathbb{R})^Q$
 391 during the training;
- 392 • Q is a hyper-parameter indicating the number of neurons on the hidden
 393 layer;
- 394 • g_1 and g_2 are the activation functions of the neural networks. Generally,
 395 in regression cases (when the outputs to be predicted are real values
 396 rather than classes), g_1 is the identity function (hence the outputs are
 397 a linear combination of the neurons on the hidden layer) and g_2 is the
 398 logistic activation function $z \rightarrow \frac{e^z}{1+e^z}$.

The weights are learned in order to minimize the mean square error on the training set:

$$\hat{\mathbf{w}} := \arg \min \sum_{i=1}^n \|y_i - f_{\mathbf{w}}(\mathbf{x}_i)\|^2. \quad (12)$$

399 Unfortunately this error is not a quadratic function of w and thus no exact
 400 algorithm is available to find the global minimum of this optimization prob-
 401 lem (and the existence of such a global minimum is not even guaranteed).
 402 Gradient descent based approximation algorithms are usually computed to
 403 find an approximate solution, where the gradient of $\mathbf{w} \rightarrow f_{\mathbf{w}}(\mathbf{x}_i)$ is calculated
 404 by the back-propagation principle (Werbos, 1974).

Moreover, to avoid overfitting, a penalization strategy, called *weight de-
 cay* (Krogh and Hertz, 1992), was introduced. It consists of replacing the

minimization problem (12) by its penalized version:

$$\hat{\mathbf{w}} := \arg \min \sum_{i=1}^n \|y_i - f_{\mathbf{w}}(\mathbf{x}_i)\|^2 + C \|\mathbf{w}\|^2$$

405 where C is the penalization parameter. The solution of this penalized mean
 406 square error is designed to be smoother than that given by Eq. (12). The `nnet`
 407 `R` function, provided in the `R` package `nnet` (Venables and Ripley, 2002), was
 408 used to train and test the one-hidden-layer MLP. As described in Section 4.3,
 409 a single validation approach was used to tune the hyper-parameters Q and
 410 C which were selected on a grid search ($Q \in \{10, 15, 20, 25, 30\}$ and $C \in$
 411 $\{0, 0.1, 1, 5, 10\}$).

412 3.2.5. SVM (Support Vector Machines)

SVM were introduced by (Boser et al., 1992) originally to address clas-
 sification problems. Subsequently (Vapnik, 1995) presented an application
 to regression problems to predict dependent real valued variables from given
 inputs. In SVM, the estimate \hat{f} is chosen among the family of functions

$$f : \mathbf{x} \in \mathbb{R}^d \rightarrow \langle \mathbf{w}, \phi(\mathbf{x}) \rangle_{\mathcal{H}} + b$$

where ϕ is a function from \mathbb{R}^d into a given Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$, here a
 RKHS, $\mathbf{w} \in \mathcal{H}$ and $b \in \mathbb{R}$ are parameters to be learned from the training
 dataset. Despite several strategies were developed to learn the parameters \mathbf{w}
 and b (Steinwart and Christmann, 2008), we opted for the original approach
 which consists of using the ϵ -insensitive loss function as a quality criterion
 for the regression:

$$L_{\epsilon}(\mathbf{X}, \mathbf{y}, \hat{f}) = \sum_{i=1}^N \max(|\hat{f}(\mathbf{x}_i) - y_i| - \epsilon, 0).$$

This loss function has the property to avoid considering the error when it
 is small enough (smaller than ϵ). His main interest, compared to the usual
 squared error, is its robustness (see (Steinwart and Christman, 2008) for
 a discussion). The SVM regression is based on the minimization of this
 loss function on the learning sample while penalizing the complexity of the
 obtained \hat{f} . More precisely, the idea of SVM regression is to find \mathbf{w} and b
 solutions of:

$$\arg \min_{\mathbf{w}, b} L_{\epsilon}(\mathbf{X}, \mathbf{y}, \hat{f}) + \frac{1}{C} \|\mathbf{w}\|_{\mathcal{H}}^2 \quad (13)$$

413 where the term $\|w\|_{\mathcal{H}}^2$ is the regularization term that controls the complexity
 414 of \hat{f} and C is the regularization parameter: when C is small, \hat{f} is allowed to
 415 make bigger errors in favor of a smaller complexity; if the value of C is high,
 416 \hat{f} makes (almost) no error on the training data but it could have a large
 417 complexity and thus not be able to give good estimations for new observa-
 418 tions (e.g., those of the test set). A good choice must devise a compromise
 419 between the accuracy required by the project and an acceptable metamodel
 420 complexity.

(Vapnik, 1995) demonstrates that, using the Lagrangian and Karush-Kuhn-Tucker conditions, w takes the form

$$\mathbf{w} = \sum_{i=1}^N (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$

421 where α_i and α_i^* solve the so-called *dual optimization problem*:

$$\begin{aligned} \arg \max_{\alpha_i, \alpha_i^*} & \left(-\frac{1}{2} \sum_{i,j=1}^N (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}} \right. & (14) \\ & \left. - \epsilon \sum_{i=1}^N (\alpha_i + \alpha_i^*) + \sum_{i=1}^N y_i (\alpha_i - \alpha_i^*) \right) \\ \text{subject to: } & \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0 \text{ and } \alpha_i, \alpha_i^* \in [0, C]. \end{aligned}$$

422 This is a classical quadratic optimization problem that can be explicitly
 423 solved. (Keerthi et al., 2001) provide a detailed discussion on the way to
 424 compute b once w is found; for the sake of clarity, in this paper we skip the
 425 full description of this step.

In Eq. (14), ϕ is only used through the dot products $(\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}})_{i,j}$. Hence, ϕ is never explicitly given but only accessed through the dot product by defining a kernel, \mathcal{K} :

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}}. \quad (15)$$

426 This is the so-called *kernel trick*. As long as $\mathcal{K} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is symmet-
 427 ric and positive, it is ensured that an underlying Hilbert space \mathcal{H} and an
 428 underlying $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ exist satisfying the relation of Eq. (15). The very

429 common *Gaussian kernel*, $\mathcal{K}_\gamma(u, v) = e^{-\gamma\|u-v\|^2}$ for a $\gamma > 0$, was used in the
 430 simulations.

431 Finally, three hyper-parameters have to be tuned to use SVM regression:

- 432 • ϵ of the loss function;
- 433 • C , the regularization parameter of the SVM;
- 434 • γ , the parameter of the Gaussian kernel.

435 As described in Section 4.3, a single validation approach was used to
 436 tune the hyper-parameters C and γ which were selected on a grid search
 437 ($C \in \{10, 100, 1000, 2000\}$ and $\gamma \in \{0.001, 0.01, 0.1\}$). To reduce the compu-
 438 tational costs and also to limit the number of hyperparameters to the same
 439 value as in MLP case (and thus to prevent the global method from being too
 440 flexible), we avoided tuning ϵ by setting it equal to 1, which corresponds ap-
 441 proximately to the second decile of the target variable for each scenario and
 442 output. This choice fitted the standard proposed by (Mattera and Haykin,
 443 1998) which suggests having a number of Support Vectors smaller than 50%
 444 of the training set. Simulations were done by using the function `svm` from
 445 the **R** package `e1071` based on the `libsvm` library (Chang and Lin, 2001).

446 3.2.6. *Random Forest*

Random forests (RF) were first introduced by (Breiman, 2001) on the
 basis of his studies on bagging and of the works of (Amit and Geman, 1997;
 Ho, 1998) on features selection. Basically, bagging consists of computing
 a large number of elementary regression functions and of averaging them.
 In random forest, elementary regression functions involved in the bagging
 procedure are regression trees (Breiman et al., 1984). Building a regression
 tree aims at finding a series of *splits* deduced from one of the d variables, x^k
 (for a $k \in \{1, \dots, d\}$), and a threshold, τ , that divides the training set into
 two subsamples, called *nodes*: $\{i : x_i^k < \tau\}$ and $\{i : x_i^k \geq \tau\}$. The split of
 a given node, \mathcal{N} , is chosen, among all the possible splits, by minimizing the
 sum of the *homogeneity* of the two corresponding child nodes, \mathcal{N}_c^1 and \mathcal{N}_c^2 ,
 as follows:

$$\sum_{i \in \mathcal{N}_c^i} (y_i - \bar{y}^{\mathcal{N}_c^i})^2$$

447 where $\bar{y}^{\mathcal{N}_c^i} = \frac{1}{|\mathcal{N}_c^i|} \sum_{i \in \mathcal{N}_c^i} y_i$ is the mean value of the output variable for the
 448 observations belonging to \mathcal{N}_c^i (i.e., the intra-node variance).

449 The growth of the tree stops when the child nodes are homogeneous
450 enough (for a previously fixed value of homogeneity) or when the number
451 of observations in the child nodes is smaller than a fixed number (generally
452 chosen between 1 and 5). The prediction obtained for new inputs, \mathbf{x} , is then
453 simply the mean of the outputs, y_i , of the training set that belong to the
454 same terminal node (a leaf). The pros of this method are its easy readability
455 and interpretability; the main drawback is its limited flexibility, especially
456 for regression problems. To overcome this limit, random forests combine a
457 large number (several hundreds or several thousands) of regression trees, T .
458 In the forest, each tree is built sticking to the following algorithm that is
459 made of random perturbations of the original procedure to make the tree
460 under-efficient (i.e., so that none of the tree in the forest is the optimal one
461 for the training dataset):

- 462 1. A given number of observations, m , are randomly chosen from the
463 training set: this subset is called *in-bag* sample whereas the other ob-
464 servations are called *out-of-bag* and are used to check the error of the
465 tree;
- 466 2. For each node of the tree, a given number of variables, q , are randomly
467 selected among all the possible explanatory variables. The best split
468 is then calculated on the basis of these q variables for the m chosen
469 observations.

470 All trees in the forest are fully learned: the final leafs all have homogeneity
471 equal to 0. Once having defined the T regression trees, $\mathcal{T}_1, \dots, \mathcal{T}_T$, the re-
472 gression forest prediction for new input variables, \mathbf{x} , is equal to the mean of
473 the individual predictions obtained by each tree of the forest for \mathbf{x} .

474 Several hyper-parameters can be tuned for random forests such as the
475 number of trees in the final forest or the number of variables randomly se-
476 lected to build a given split. But, as this method is less sensitive to parameter
477 tuning than the other ones (i.e., SVM and MLP), we opted for leaving the
478 default values implemented in the **R** package `randomForest` based on useful
479 heuristics: 500 trees were trained, each defined from a bootstrap sample built
480 with replacement and having the size of the original dataset. Each node was
481 defined from three randomly chosen variables and the trees were grown until
482 the number of observations in each node was smaller than five. Moreover,
483 the full learning process always led to a stabilized out-of-bag error.

484 4. Simulations and results

485 4.1. Application to the Cross Compliance Assessment Tool

486 As described above in the Section 1, the impact assessment of Cross Com-
487 pliance (CC) measures on the EU27 farmlands, required the development of a
488 simulation platform called Cross Compliance Assessment Tool (CCAT). The
489 CCAT framework integrates different models, such as Miterra (Velthof et al.,
490 2009), DNDC-EUROPE (Follador et al., 2011), EPIC (van der Velde et al.,
491 2009) and CAPRI (Britz and Witzke, 2008; Britz, 2008), in order to guarantee
492 an exhaustive evaluation of the effects of agro-environmental standards for
493 different input, scenario assumptions, compliance rates and space-time reso-
494 lutions (Elbersen et al., 2010; De Vries et al., 2008). The simulated outputs
495 are indicators for nitrogen (N) and carbon (C) fluxes, biodiversity and land-
496 scape, market response and animal welfare. The selection of the CC scenarios
497 as well as the definition of the environmental indicators to be considered in
498 this project, are described by (Jongeneel et al., 2008). The CCAT tool eval-
499 uates the effect of agricultural measures on N₂O fluxes and N leaching by
500 means of the meta-model of the mechanistic model DNDC-EUROPE (Fol-
501 lador et al., 2011). N₂O is an important greenhouse gas (Intergovernmental
502 Panel on Climate Change, 2007). Agriculture and in particular agricultural
503 soils are contributing significantly to anthropogenic N₂O emissions (Euro-
504 pean Environment Agency, 2010). N₂O fluxes from soils are characterized
505 by a high spatial variability and the accuracy of estimates can be increased if
506 spatially explicit information is taken into consideration (Leip et al., 2011a).
507 Similarly, leaching of nitrogen from agricultural soils is an important source
508 of surface and groundwater pollution (European Environment Agency, 1995).

509 The main limits of using DNDC-EUROPE directly in the CCAT platform
510 are the high computational costs and memory requirements, due to the large
511 size of input datasets and the complexity and high number of equations to
512 solve. To mitigate this problem, making the integration easier, we decided
513 to develop a metamodel of DNDC-EUROPE (Follador and Leip, 2009). The
514 choice of the best meta-modeling approach has been based on the analysis
515 of performance of different algorithms, as described in details in Section 4.4.
516 The best metamodel is expected to have low computational costs and an
517 acceptable accuracy for all the dataset sizes.

518 *4.2. Input and Output data description*

519 The set of training observations (around 19 000 observations) used to de-
520 fine a metamodel \hat{f} was created by linking the agro-economic CAPRI dataset
521 with the bio-geochemical DNDC-EUROPE model at Homogeneous Spatial
522 Mapping Unit (HSMU) resolution, as described in (Leip et al., 2008). We
523 opted for corn cultivation as case study, since it covers almost 4.6% of UAA
524 (utilized agricultural area) in EU27, playing an important role in human and
525 animal food supply (European Union Commission, 2010)¹ and representing
526 one of the main cropping system in Europe. To obtain a representative
527 sample of situations for the cultivation of corn in EU27, we selected about
528 19,000 HSMUs on which at least 10% of the agricultural land was used for
529 corn (Follador et al., 2011).

530 The input observations used to train the metamodels were drawn from
531 the whole DNDC-EUROPE input database (Leip et al., 2008; Li et al., 1992;
532 Li, 2000), in order to meet the need of simplifying the I/O flux of information
533 between models in the CCAT platform. This screening was based on a pre-
534 liminary sensitivity analysis of input data through the *importance function*
535 of the **R** package `randomForest`, and subsequently it was refined by expert
536 evaluations (Follador et al., 2011; Follador and Leip, 2009). At last, 11 input
537 variables were used:

- 538 • Variable related to N input [$\text{kgN ha}^{-1}\text{yr}^{-1}$], such as mineral fertil-
539 izer (`N_FR`) and manure (`N_MR`) amendments, N from biological fixation
540 (`N_fix`) and N in crop residue (`N_res`);
- 541 • variables related to soil: soil bulk density, `BD`, [g cm^{-3}], topsoil organic
542 carbon, `SOC`, [mass fraction], clay content, `clay`, [fraction] and topsoil
543 pH, `pH`;
- 544 • variables related to climate: annual precipitation `Rain`, [mm yr^{-1}],
545 annual temperature `Tmean` [$^{\circ}\text{C}$] and N in rain, `Nr`, [ppm].

546 They refer to the main driving forces taking part in the simulation of N_2O and
547 N leaching with DNDC-EUROPE, such as farming practices, soil attributes
548 and climate information. In this contribution we only show the results for
549 the corn baseline scenario - that is the conventional corn cultivation in EU27,

¹<http://epp.eurostat.ec.europa.eu>

550 as described by (Follador et al., 2011). Note that a single metamodel was
551 developed for each CC scenario and for each simulated output in CCAT, as
552 described in (Follador and Leip, 2009). Figure 1 summarizes the relations
553 between the DNDC-EUROPE model and the metamodel.

554 As the number of input variables was not large, they were all used in all
555 the metamodeling methods described in Section 3, without additional vari-
556 able selection. The only exception is the second linear model (Section 3.1)
557 which uses a more complete list of input variables obtained by various com-
558 binations of the original 11 variables and thus includes a variable selection
559 process to avoid collinearity issues.

560 Two output variables were studied: the emissions of N_2O ($[\text{kg N yr}^{-1}$
561 $\text{ha}^{-1}]$ for each HSMU), a GHG whose reduction is a leading matter in cli-
562 mate change mitigation strategies, and the nitrogen leaching ($[\text{kg N yr}^{-1}$
563 $\text{ha}^{-1}]$ for each HSMU), which has to be monitored to meet the drinking
564 water quality standards (Askegaard et al., 2005). A metamodel was devel-
565 oped for each single output variable. The flux of information through the
566 DNDC-EUROPE model and its relationship with the metamodel’s one are
567 summarized in Figure 1. The data were extracted using a high performance
568 computer cluster and the extraction process took more than one day for all
569 the 19 000 observations.

570 [Figure 1 about here.]

571 4.3. Training, validation and test approach

572 The training observations were randomly partitioned (without replace-
573 ment) into two groups: 80% of the observations (i.e., $N_L \simeq 15\,000$ HSMU)
574 were used for training (i.e., for defining a convenient \hat{f}) and the 20% re-
575 maining observations (i.e., $N_T \simeq 4\,000$ HSMU) were used for validating the
576 metamodels (i.e., for calculating an error score). Additionally, in order to
577 understand the impact of the training dataset on the goodness of the esti-
578 mations (\hat{y}_i) and to compare the different metamodel performance according
579 to the data availability, we randomly selected from the entire training set a
580 series of subsets, having respectively $N_L = 8\,000, 4\,000, 2\,000, 1\,000, 500,$
581 200 and 100 observations, each consecutive training subset being a subset of
582 the previous one.

583 The methodology used to assess the behavior of different metamodels
584 under various experimental conditions (size of the dataset and nature of the
585 output) are summarized in Description 1.

Description 1 Methodology used to compare the metamodels under various experimental conditions

- 1: **for** Each metamodel, each output and each size N_L **do**
 - 2: {**Train** the metamodel with the N_L training observations \rightarrow definition of \hat{f} ;
 - 3: **Estimate the outputs** for the $N_T \simeq 4\,000$ inputs of the test set from $\hat{f} \rightarrow$ calculation of \hat{y}_i ;
 - 4: **Calculate the test error** by comparing the estimated outputs, \hat{y}_i , vs. the outputs of the DNDC-EUROPE model for the same test observations, y_i .}
 - 5: **end for**
-

586 More precisely, for some metamodels, Step 2 requires the tuning of some
587 hyper-parameters (e.g., SVM have three hyper-parameters, see Section 3).
588 These hyper-parameters were tuned by:

- 589 • for ACOSSO and SDR: a grid-search to minimize BIC plus an algorithm
590 to get the weights w_j : in these cases, an efficient formula, that does
591 not require to compute each leave-one-out estimate of f , can be used to
592 compute the BIC; moreover the COSSO penalty provides all θ_j given
593 λ and w_j in a single shot. In the SDR identification steps, a maximum
594 likelihood strategy is applied to optimize NVR's;
- 595 • for DACE, a maximum likelihood strategy;
- 596 • for MLP, SVM and RF, a *simple validation strategy* preferred to a
597 cross validation strategy to reduce the computational time especially
598 with the largest training datasets): half of the data were used to define
599 several metamodels depending on the values of hyper-parameters on a
600 grid search and the remaining data were used to select the best set of
601 hyper-parameters by minimizing a mean square error criterion.

602 Hence, depending on which features are the most interesting (easy tuning of
603 the hyperparameters, size of the training dataset, size of the dataset need-
604 ing new prediction...), the use of one method is more or less recommended.
605 Table 1 summarizes the main characteristics of the training and validation
606 steps of each method as well as the characteristics to do new predictions. For
607 instance, linear models are more

Table 1: Summary of the main features for training, validation (hyperparameters tuning) and test steps of each method.

Method	Training characteristics	Validation characteristics	New predictions characteristics
LM1	Very fast to train.	There is no hyperparameter to tune.	Very fast.
LM2	Fast to train but much slower than LM1 because of the number of parameters to learn.	There is no hyperparameter to tune.	Very fast.
ACOSSO	Fast to train only if the number of observations is very low: the dimension of the kernel matrix is $N_L \times N_L$ and it is obtained as the sum of the kernels of each $[N_L \times N_L]$ ANOVA term, which can be long to calculate.	One hyperparameter (λ) is tuned twice by minimizing BIC: the first time to get the weights w_j the second to get the final estimate (given λ and w_j the COSSO penalty provides automatically in a single shot all θ_j).	The time needed to obtain new predictions can be high depending on the sizes of both the training dataset and the test dataset. It requires to compute a kernel matrix having dimension $N_L \times N_T$.
<i>Continued on next page</i>			

Table 1 – Continued from previous page

Method	Training characteristics	Validation characteristics	Test characteristics
SDR	<p>Fast to train only if the number of observations is very low: the dimension of the kernel matrix is $N_L \times N_L$ and it is obtained as the sum of the kernels of each $[N_L \times N_L]$ ANOVA term, which can be long to calculate.</p>	<p>As for ACOSSO, the single hyperparameter (λ) is tuned by minimizing BIC: the SDR identification step to provide w_j also optimizes hyperparameters for each ANOVA component but this can be done efficiently by the SDR recursive algorithms (given λ and w_j the COSSO penalty provides automatically in a single shot all θ_j).</p>	<p>The time needed to obtain new predictions can be high depending on the sizes of both the training dataset and the test dataset. It requires to compute a kernel matrix having dimension $N_L \times N_T$.</p>
DACE	<p>Fast to train only if the number of observations is very low: the dimension of the kernel matrix is $N_L \times N_L$, and the inversion of a matrix $N_L \times N_L$ is required in the GLS procedure.</p>	<p>$d + 1$ hyperparameters are tuned by ML, which becomes intractable already for moderate d: each step of the optimization a matrix $N_L \times N_L$ has to be inverted.</p>	<p>The time needed to obtain new predictions can be high depending on the sizes of both the training dataset and the test dataset. It requires to compute a kernel matrix having dimension $N_L \times N_T$.</p>

Continued on next page

Table 1 – Continued from previous page

Method	Training characteristics	Validation characteristics	Test characteristics
MLP	Hard to train: because the error to minimize is not quadratic, the training step faces local minima problems and has thus to be performed several times with various initialization values. It is also very sensitive to the dimensionality of the data (that strongly increases the number of weights to train) and, to a lesser extent, to the number of observations.	2 hyperparameters have to be tuned but one is discrete (number of neurons on the hidden layer) which is easier. Nevertheless, cross validation is not suited: tuning is performed by simple validation and can thus be less accurate. It can be time consuming .	The time needed to obtain new predictions is very low : it depends on the number of predictions.
SVM	Fast to train if the number of observations is low: SVM are almost insensitive to the dimensionality of the data but the dimension of the kernel matrix is $N_L \times N_l$ and can be long to calculate.	Three hyperparameters have to be tuned and in the case where the size of the training dataset is large, cross validation is not suited. Tuning is performed by simple validation and can thus be less accurate. It is also time consuming .	The time needed to obtain new predictions can be high depending on the sizes of both the training dataset and the test dataset. It requires to compute a kernel matrix having dimension $N_L \times N_T$.

Continued on next page

Table 1 – Continued from previous page

Method	Training characteristics	Validation characteristics	Test characteristics
RF	Fast to train: almost insensitive to the size or the dimensionality of the training dataset thanks to the random selections of observations and variables. Most of the time needed to train is due to the number of trees required to stabilize the algorithm, that can sometimes be large.	Almost insensitive to hyperparameters so no extensive tuning is required.	The time needed to obtain new predictions is low : it depends on the number of predictions to do and also on the number of trees in the forest.

608

609 In Step 4, the test quality criterion was evaluated by calculating several
610 quantities:

- the *Mean Squared Error* (MSE):

$$\text{MSE} = \frac{1}{N_T} \sum_{i=1}^{N_T} (\hat{y}_i - y_i)^2$$

611 where y_i and \hat{y}_i are, respectively, the model outputs in the test dataset
612 and the corresponding approximated outputs given by the metamodel.

- the R^2 coefficient:

$$R^2 = 1 - \frac{\sum_{i=1}^{N_T} (\hat{y}_i - y_i)^2}{\sum_{i=1}^{N_T} (\hat{y}_i - \bar{y})^2} = 1 - \frac{\text{MSE}}{\text{Var}(y)}$$

613 where \bar{y} and $\text{Var}(y)$ are the mean and the variance of all y_i in the test
614 dataset. R^2 is equal to 1 if the predictions are perfect and thus gives
615 a way to quantify the accuracy of the predictions to the variability of
616 the variable to predict.

- 617
618
619
 • the *standard deviation of the SE* and the *maximum value of the SE* were also computed to give an insight on the variability of the performance and not only on its mean.

620 *4.4. Results and discussion*

621
622
623
624
625
626
627
 This section includes several ways to compare the methods on the problem described in 4.2. First, Section 4.4.1 compares the accuracy of the predictions for various methods and various training dataset sizes. Then, Section 4.4.2 gives a comparison of the computational times needed either to train the model (with the maximum dataset size) and to make new predictions. Finally, Section 4.4.3 describes the model itself and gives an insight about its physical interpretation.

628 *4.4.1. Accuracy*

629
630
631
632
633
634
635
636
 The performance on the test set is summarized in Tables 2 to 5: they include characteristics about the mean values of the squared errors (MSE and R^2) in Tables 2 and 3, respectively for N₂O and N leaching predictions, as well as characteristics related to the variability of the performance (standard deviations of the squared errors and maximum values of the squared errors) in Tables 4 and 5, respectively for N₂O and N leaching predictions. Note that, in almost all cases, the minimum values of the squared errors were equal or close to 0.

637 [Table 1 about here.]

638 [Table 2 about here.]

639 [Table 3 about here.]

640 [Table 4 about here.]

641
642
643
 The evolution of R^2 on the test set in function of the size of the training set is displayed in Figures 2 (N₂O prediction) and 3 (N leaching) for each method.

644 [Figure 2 about here.]

645 [Figure 3 about here.]

646 From these results, several facts clearly appeared:

- 647 • Even for small datasets, the metamodeling approach behaves correctly
648 with R^2 always greater than 80% for the best approaches. Note that the
649 poorest results (those that are the closest to 80%) are obtained for small
650 training dataset sizes (100 or 200). This means that, in the case where
651 several metamodels are needed to model various assumptions of the
652 input variables ranges, crude but acceptable estimates can be obtained
653 at a very low computational cost. For more efficient predictions, larger
654 datasets are more suited and achieve R^2 values greater than 90%.
- 655 • Predicting N leaching seems an easier task than predicting N_2O fluxes
656 with greater performance for almost any training dataset size. This is
657 not surprising because N_2O is generated as an intermediate product in
658 the denitrification chain, being produced by the reduction of nitrate,
659 but being consumed by N_2O denitrifiers. As a consequence, N_2O fluxes
660 are the result of a fragile equilibrium between those processes which are
661 both highly sensitive on environmental conditions such as pH, oxygen
662 availability, substrate availability (Firestone et al., 1979). Thus, N_2O
663 fluxes are characterized by a very high spatial variability and is much
664 harder to predict than nitrogen leaching (Britz and Leip, 2009; Leip
665 et al., 2011a).
- 666 • The best results are obtained for the largest training dataset. Mostly,
667 for all methods, the performance increases with the size of the learn-
668 ing dataset despite some exceptions: sometimes, using a larger dataset
669 makes the training process harder and can slightly deteriorate the per-
670 formance (e.g., for MLP, large datasets leads to harder local minima
671 problems in the optimization procedure: for this method, the best pre-
672 diction of N leaching estimates is not obtained from the largest training
673 set).
- 674 • In a similar way, the variability of the errors tends to decrease with
675 the size of the training dataset but some methods behave differently
676 (see, e.g., Acosso whose variability strictly increases with the size of
677 the training dataset for N leaching prediction).
- 678 • In most cases, the most accurate predictions (according to MSE or
679 R^2 values) are also the predictions that have the smallest variability
680 either from the standard deviation point of view or from the smallest
681 maximum point of view.

682 Looking deeper into the methods themselves, the following conclusions
683 can also be derived:

- 684 • LM1 gives poor performance because the plain linear model is probably
685 too simple to catch the complexity of the modeled phenomenon.
- 686 • LM2 performs very badly for small training datasets since it is over-
687 specified (the number of parameters to be estimated is close to the
688 size of the dataset; R^2 are negative which means that the model is
689 less accurate than the trivial model predicting any observation by the
690 mean value of the outputs). But for large training datasets, it behaved
691 correctly. Additionally, the number of variables selected during the step
692 AIC, in function of the training dataset size, is given in Table 6. The
693 number of selected variables for N leaching prediction is higher than
694 the number of selected variables for N₂O prediction but it also tends
695 to be more stable regarding the dataset size. Also note that, in any
696 case, the number of selected variables is high compared to the original
697 number of variables (120): this means that the underlying model under
698 study is certainly not plain linear and this explains why LM1 fails to
699 approximate it accurately.

700 [Table 5 about here.]

- 701 • Splines and kriging based methods have the best performance for small
702 and medium training datasets (especially for N leaching prediction) but
703 they can not be run for large training datasets (up to 2 000 observa-
704 tions) due to the calculation costs. The Dace and SDR models have
705 the best performance. Additionally, the number of selected variables
706 for ACOSSO and SDR are given in Table 7. The number of compo-
707 nents effectively included in the model tend decrease with the training
708 set size, especially for N₂O prediction. Comparing this table with Ta-
709 ble 6, the number of components is also quite small, even smaller than
710 the number of original variables for some cases.

711 [Table 6 about here.]

- 712 • Machine learning methods (MLP, SVM and RF) behave correctly for
713 medium training datasets and obtain the best performance for large
714 training datasets. SVM and RF have the best results with a very good

715 overall accuracy, as, for these methods, R^2 are greater than 90% and
716 95%, respectively for N₂O and N leaching predictions.

717 Moreover, Wilcoxon paired tests on the residuals (absolute value) were
718 computed to understand if the differences in accuracy between the best meth-
719 ods were significant: for N₂O prediction, the difference between the best per-
720 formance (RF) and the second one (SVM) is significant (p-value equal to
721 0.16%) whereas, for N leaching prediction, the difference between the best
722 performance (SVM) and the second one (RF) is not significant. This test con-
723 firms the differences between the best performance of metamodells obtained
724 with different dataset sizes: for example, the difference between SVM trained
725 with about 15 000 observations and Dace trained with 2 000 observations is
726 significant (p-value smaller than $2.2 \cdot 10^{-16}$).

727 Finally, we took into account the time needed to train the metamodel
728 and subsequently to use it for prediction. The time for training is not so
729 important as it is spent only once during the calibration step. The time for
730 prediction is a key point for CCAT project and so it played a leading role in
731 choosing the best metamodel; it must be quite limited to allow fast multi-
732 scenario simulations or sensitivity analysis. Table 8 provides the approximate
733 time spent to train and use each method with large datasets (respectively,
734 about 15 000 observations for the training step and about 19 000 observations
735 for the prediction one) on a desktop computer.

736 [Table 7 about here.]

737 4.4.2. Computational time

738 The training time for LM1 was the best one but the corresponding per-
739 formance is very poor. RF had a low training time since it does not require
740 any parameter tuning and it is not very sensitive to the size of dataset thanks
741 to the bootstrapping procedure. The prediction time is really low for all the
742 methods compared to the DNDC-EUROPE runs which had demanded about
743 1 day to simulate the same outputs on a high performance computer cluster.
744 Even though RF was not the fastest approach it provides the best compro-
745 mise between speed and accuracy. SVM spent more time in prediction since
746 it required the calculation of the kernel matrix whose size is proportional
747 (and thus much more sensitive) to the number of new predictions to make.
748 The same issues applies to splines approach, where the kernel matrix has to
749 be re-computed for every ANOVA term in the decomposition, as well as for
750 kriging, thus explaining the larger computational cost. The highest cost for

751 SDR predictions are linked to the more detailed decomposition, which im-
752 plies a larger number of reproducing kernels. To compute the large amount
753 of 19 000 model outputs, the time required for predictions does not exceed a
754 few minutes in any cases.

755 4.4.3. *Metamodeling interpretation*

756 To give an indication of which variables are important in the prediction
757 of both inputs, an “importance” measure was calculated for each variable of
758 the best final model (i.e., random forest trained with the full training dataset
759 for N₂O prediction and SVM trained with the full dataset for N leaching
760 prediction). For random forests, the importance is quite common: for a
761 given input variable, the values of out-of-sample observations are randomly
762 permuted; the mean squared error is then calculated based on all out-of-
763 sample sets for all trees in the forest. The increase in the mean squared
764 error compared to the out-of-sample mean squared error calculated with the
765 true values of the predictor is called the importance of the predictor (see
766 (Genuer et al., 2010) for a complete study of this quantity in the framework
767 of variable selection problems). Unfortunately, MLPs and SVMs are not
768 based on bootstrapping so out-of-sample observations do not exist for these
769 methods. Hence, importance cannot be defined or directly compared to
770 the one given for random forests. Nevertheless, a close definition can be
771 introduced by using the validation set selected for the tuning process and by
772 comparing the mean squared error of permuted inputs to the true squared
773 error on this validation set.

774 Figure 4 illustrates the values of the importance measure in both cases. It
775 can be seen that the two metamodels are very different: that (RF) which aims
776 at estimating N₂O fluxes (left) is mainly based on two important variables
777 (SOC and pH) whereas SVM, used to estimate N leaching, has a less strict
778 behavior: at least four variables are important in that last modeling, N_MR,
779 N_FR, pH and Nres.

780 [Figure 4 about here.]

781 N₂O fluxes are mainly related to denitrification processes, which require
782 anaerobic conditions and organic material as substrate (Firestone et al.,
783 1979). Anaerobic conditions form if diffusion of oxygen is blocked in wet
784 soils, or in denitrification “hotspots” around organic matter promoting very
785 high oxygen consumption rates (Parkin, 1987). It is therefore not surprising

786 that the soil organic carbon content (SOC) was found to be the most impor-
787 tant for the prediction of N₂O fluxes. Soil pH is also an important parameter,
788 influencing both the reduction of nitrate (total denitrification) but also the
789 reduction of N₂O to N₂ (Granli and Bøckman, 1994). For nitrogen leaching,
790 on the other hand, we found that the most important factor was the most
791 important factor was nitrogen input as manure amendment, mineral fertilizer
792 spreading, and N from crop residue incorporation in the soil before sowing
793 (these are indeed even more important than pH). To a large degree, nitrogen
794 leaching is determined by soil texture which controls the percolation rate of
795 water through the soil profile and precipitation. As a consequence, it is not
796 surprising to find the top-factors determining nitrogen leaching in a relatively
797 narrow range, as compared to N₂O fluxes.

798 4.5. Conclusion about the comparison of metamodeling strategies

799 The experiments described in the following subsections enlighten several
800 facts: first, metamodeling strategies were able to approximate accurately
801 N₂O and N leaching predictions at a low computational cost. Even with
802 small dataset sizes (100 HSMUs to train the data), the overall accuracy rate,
803 measured by R^2 , is greater than 80% for at least one metamodel. In this case
804 study, N₂O was harder to predict than N leaching. Then, increasing the size
805 of the training dataset is time consuming but also leads to a better accuracy
806 in the prediction for (almost) all the methods. Hence, the selection of a
807 metamodeling approach has to be based on a careful compromise between
808 computational costs and accuracy. This choice strictly depends on the size
809 of available training data and on the project's target. We pointed out that
810 splines and kriging based methods should be chosen when the number of
811 training data is smaller than 2 000 since they provided the most accurate
812 solution with a reasonable running time. With large datasets, random forests
813 were able to handle the training step and to calculate accurate predictions
814 with low computational costs (more than 15 000 observations were trained in
815 about 15 minutes and only several seconds were needed in predicting 19 000
816 new values).

817 Finally, we pointed out, in Section 4.4.3, that combining metamodeling
818 with an importance measure can also be used to provide a simplified insight
819 on the important processes and on the main input variables involved in the
820 prediction of N₂O fluxes and N leaching. This can help to find strategies to
821 control nitrogen surplus or to perform a fast sensitivity analysis. This last
822 issue is currently under investigation.

823 5. Conclusion

824 This article provides a full comparison of several metamodel approaches
825 for the prediction of N₂O fluxes and N leaching from European farmlands.
826 The conclusions of the meta-model comparison are general enough to be
827 extended to other similar case studies. A more valuable and detailed impact
828 assessment of CC standards at European or country level is possible only
829 by simulating all the 207000 HSMUs that cover the EU27. This approach
830 demands the collection of enormous amounts of data and their storage into
831 large datasets. From our work, random forest proved to be a reliable and
832 effective tool for elaborating large datasets with low computational costs and
833 an acceptable accuracy. For these reasons it has been chosen to be integrated
834 into the CCAT platform to estimate the N₂O fluxes and N leaching from the
835 EU27 farmlands.

836 References

- 837 Akaike, H., 1974. A new look at the statistical model identification. *IEEE*
838 *Transactions on Automatic Control* 19, 716–723.
- 839 Alliance Environnement, 2007. Evaluation of the application of cross com-
840 pliance as foreseen under regulation 1782/2003. Part I: descriptive report.
841 Technical Report. D.G. Agriculture. Institute for European Environmental
842 Policy (IEEP), London, UK and Orade-Brche Sarl, Auzeville, France.
- 843 Amit, Y., Geman, D., 1997. Shape quantization and recognition with random
844 trees. *Neural Computation* 9, 1545–1588.
- 845 Anderson, M., Kustas, W., Norman, J., 2003. Upscaling and downscaling - a
846 regional view of the soil-plant-atmosphere continuum. *Agronomy Journal*
847 95, 1408–1423.
- 848 Askegaard, M., Olesen, J., Kristensen, K., 2005. Nitrate leaching from or-
849 ganic arable crop rotations: effects of location, manure and catch crop.
850 *Soil Use and Management* 21, 181–188.
- 851 Berlinet, A., Thomas-Agnan, C., 2004. *Reproducing Kernel Hilbert Spaces in*
852 *Probability and Statistics*. Kluwer Academic Publisher, Boston, Norwell,
853 MA, USA / Dordrecht, The Netherlands.

- 854 Bishop, C., 1995. Neural Networks for Pattern Recognition. Oxford Univer-
855 sity Press, New York, USA.
- 856 Blanning, R., 1975. The construction and implementation of metamodels.
857 Simulation 24, 177–184.
- 858 Boser, B., Guyon, I., Vapnik, V., 1992. A training algorithm for optimal
859 margin classifiers, in: 5th annual ACM Workshop on COLT, D. Haussler
860 Editor, ACM Press. pp. 144–152.
- 861 Bouzaher, A., Lakshminarayan, P., Cabe, R., Carriquiry, A., Gassman, P.,
862 Shogren, J., 1993. Metamodels and nonpoint pollution policy in agricul-
863 ture. Water Resources Research 29, 1579–1587.
- 864 Breiman, L., 2001. Random forests. Machine Learning 45, 5–32.
- 865 Breiman, L., Friedman, J., Olsen, R., Stone, C., 1984. Classification and
866 Regression Trees. Chapman and Hall, New York, USA.
- 867 Britz, W., 2008. Automated model linkages: the example of CAPRI. Agrar-
868 witschaft 57, 363–367.
- 869 Britz, W., Leip, A., 2009. Development of marginal emission factors for N
870 losses from agricultural soils with the DNDC-CAPRI metamodel. Agriculture,
871 Ecosystems and Environment 133, 267–279.
- 872 Britz, W., Witzke, P., 2008. CAPRI model documentation 2008. CAPRI
873 project, Institute for Food and Resource Economics. Bonn, Germany. On-
874 line at: <http://www.capri-model.org>.
- 875 Chadwick, D., 2005. Emissions of ammonia, nitrous oxide and methane from
876 cattle manure heaps: effect of compaction and covering. Atmospheric
877 Environment 39, 787–799.
- 878 Chang, C., Lin, C., 2001. LIBSVM: a library for support vector machines.
879 Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 880 Christmann, A., Steinwart, I., 2007. Consistency and robustness of kernel-
881 based regression in convex risk minimization. Bernoulli 13, 799–819.
- 882 Cressie, N., 1990. The origins of kriging. Mathematical Geology 22, 239–352.

- 883 De Vries, W., Kros, H., Velthof, G., Oudendag, D., Leip, A., van der Velde,
884 M., Kempen, M., 2008. Development and application methodology of en-
885 vironmental impact tool to evaluate cross compliance measures. Technical
886 Report Deliverable 4.2.1 of EU STREP 44423. CCAT Project.
- 887 Del Grosso, S., Parton, W., Mosier, A., Walsh, M., D.S., O., Thornton, P.,
888 2006. DAYCENT national-scale simulation of nitrous oxide emissions from
889 cropped soils in the United States. *Journal of Environment Quality* 35,
890 1451–1460.
- 891 Dimopoulus, D., Fermantziz, I., Vlahos, G., 2007. The responsiveness of cross
892 compliance standards to environmental pressure. Technical Report Deliv-
893 erable 12 of the CC Network Project, SSPE-CT-2005-022727. Agricultural
894 University of Athens. Athens, Greece.
- 895 Elbersen, B., Jongeneel, R., Kempen, M., Klein-Lankhorst, R., De Vries, W.,
896 Lesschen, J., Onate, J., Alonso, M., Kasperczyk, N., Schramek, J., Mik,
897 M., Peepson, A., Bouma, F., Staritsky, I., Kros, H., 2010. Final report
898 of CCAT project results. Technical Report Deliverable 2.8 of EU STREP
899 44423-CCAT. CCAT project.
- 900 European Commission, 2002. Implementation of Council Directive
901 91/676/EEC concerning the protection of waters against pollution caused
902 by nitrates from agricultural sources. Technical Report. D.G. Environ-
903 ment, European Commission. Luxembourg.
- 904 European Council, 1991. Nitrate Directive 91/676/EEC. Technical Report.
905 European Union. Brussels, Belgium.
- 906 European Council, 2000. Water Framework Directive 2000/60/EC. Technical
907 Report. European Union. Brussels, Belgium.
- 908 European Council, 2003. Council regulation (EC). Technical Report
909 1782/2003. European Union. Brussels, Belgium.
- 910 European Council, 2009. Council regulation (EC) No 73/2009. Technical
911 Report. European Council. Brussels, Belgium.
- 912 European Environment Agency, 1995. Europe's environment, the dobriss
913 assessment. Technical Report. European Environment Agency.

- 914 European Environment Agency, 2010. Annual European community green-
915 house inventory 1980-2008 and inventory report 2010. Submission to the
916 UNFCCC secretariat. Technical Report. European Environment Agency.
917 Copenhagen, Denmark.
- 918 European Union Commission, 2004. Commission Regulation No 796/2004.
919 Technical Report. European Commission. Brussels, Belgium.
- 920 European Union Commission, 2009. Commission Regulation No 1122/2009.
921 Technical Report. European Commission. Brussels, Belgium.
- 922 European Union Commission, 2010. Agriculture in the EU - Statistical and
923 Economic Information - Report 2009. Technical Report. European Com-
924 mission. Brussels, Belgium.
- 925 FAO, 2005. Key to drought-resistant soil and sustained food production.
926 The importance of soil organic matter. Technical Report Soils bulletin 80.
927 Food Agricultural Organization. Rome, Italy.
- 928 FAO, 2007. Payer les agriculteurs pour les services environnementaux. La
929 situation mondiale de l'alimentation et de l'agriculture. Technical Report.
930 Food Agricultural Organization. Rome, Italy.
- 931 Firestone, M., Smith, M., Firestone, R., Tiedje, J., 1979. The influence of
932 nitrate, nitrite, and oxygen on the composition of the gaseous products of
933 denitrification in soil. Soil Science Society of America Journal 43, 1140–
934 1144.
- 935 Follador, M., Leip, A., 2009. Derivation of DNDC meta-models to evaluate
936 the impact of cross compliance measures on nitrogen N surplus, N leaching,
937 N₂O emissions at European scale. Technical Report Deliverable 4.2.3 of EU
938 STREP 44423-CCA. European Commission, D.G Joint Research Centre,
939 Institute for Environment and Sustainability, Climate Change Unit. Ispra,
940 Italy.
- 941 Follador, M., Leip, A., Orlandini, L., 2011. Assessing the impact of cross
942 compliance measures on nitrogen fluxes from european farmlands with
943 DNDC-EUROPE. Environmental Pollution In Press.
- 944 Forrester, A., Keane, A., 2009. Recent advances in surrogate-based optimiza-
945 tion. Progress in Aerospace Sciences 45, 50–79.

- 946 Friedman, J., 1991. Multivariate adaptive regression splines. *Annals of Statistics*
947 19, 1–67.
- 948 Genuer, R., Poggi, J., Tuleau-Malot, C., 2010. Variable selection using ran-
949 dom forests. *Pattern Recognition Letters* 31, 2225–2236.
- 950 van Gighc, J., 1991. System design modeling and metamodeling. Springer,
951 New York, USA.
- 952 Granli, T., Bøckman, O., 1994. Nitrous oxide from agriculture. *Norwegian*
953 *Journal of Agricultural Sciences Supplement No. 12*. Norsk Hydro Research
954 Centre: Porsgrunn, Norway.
- 955 van Grinsven, H., Ward, M., Benjamin, N., De Kok, T., 2006. Does the
956 evidence about health risks associated with nitrate ingestion warrant an
957 increase of the nitrate standard for drinking water? *Environmental Health:*
958 *A Global Access Science Source* 5.
- 959 Gu, C., 2002. Smoothing Spline ANOVA Models. Springer-Verlag, New
960 York, USA.
- 961 Haberlandt, U., Krysanova, V., Bardossy, A., 2002. Assessment of nitrogen
962 leaching from arable land in large river basins. Part II: regionalisation using
963 fuzzy rule based modelling. *Ecological Modelling* 150, 277–294.
- 964 Ho, T., 1998. The random subspace method for constructing decision forests.
965 *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20, 832–
966 844.
- 967 Hornik, K., 1991. Approximation capabilities of multilayer feedforward net-
968 works. *Neural Networks* 4, 251–257.
- 969 Intergovernmental Panel on Climate Change, 2007. Fourth Assessment Re-
970 port of the Intergovernmental Panel on Climate Change. Working Group
971 I - The Physical Science Basis. Cambridge University Press, Cambridge,
972 United Kingdom / New York, NY, USA,.
- 973 Jongeneel, R., Elbersen, B., Klein-Lankhorst, R., De Vries, W., Kros, H.,
974 Velthof, G., Kempen, M., Annen, D., Onate, J., van der Velde, M., Leip,
975 A., 2008. Report describing the operationalisation of the first selection of
976 indicators into impacts of Cross Compliance for the implementation in the

- 977 first prototype of the analytical tool. Technical Report Deliverable 2.3 of
978 EU STREP 44423-CCAT.
- 979 Jongeneel, R., Elbersen, B., de Vries, V., Klein-lankhorst, J., Schramek, J.,
980 Rudloff, B., Heckelei, T., Kempen, M., Annen, D., van der Velde, M., Leip,
981 A., Redman, M., Mikk, M., Onate, J., Slangen, L., 2007. General approach
982 to the assessment of the impacts of CC in the EU and list of indicators.
983 Technical Report Deliverables 2.1 and 2.2 of EU STREP 44423. CCAT
984 Project.
- 985 Kalman, R., 1960. A new approach to linear filtering and prediction prob-
986 lems. *Transactions of the ASME, Journal of Basic Engineering* 82D, 35–45.
- 987 Keerthi, S., Shevade, S., Bhattacharyya, C., Murty, K., 2001. Improvements
988 to platt’s SMO algorithm for SVM classifier design. *Neural Computation*
989 13, 637–649.
- 990 Kennedy, M., O’Hagan, A., 2001. Bayesian calibration of computer models
991 (with discussion). *Journal of the Royal Statistical Society, series B* 63,
992 425–464.
- 993 Kirchmann, H. and Esala, M., Morken, J., Ferm, M., Bussink, W., Gus-
994 tavsson, J., Jakobsson, C., 1998. Ammonia emissions from agriculture.
995 *Nutrient Cycling in Agroecosystems* 51, 1–3.
- 996 Kleijnen, J., 1975. A comment on Blanning’s “Metamodel for sensitivity
997 analysis: the regression metamodel in simulation”. *Interfaces* 5, 21–23.
- 998 Kleijnen, J., 2009. Kriging metamodelin in simulation: a review. *European*
999 *Journal of Operational Research* 1992, 707–716.
- 1000 Kleijnen, J., Sargent, R., 2000. A methodology for fitting and validating
1001 metamodels in simulation. *European Journal of Operational Research* 120,
1002 14–29.
- 1003 Krogh, A., Hertz, J., 1992. A simple weight decay can improve generalization,
1004 in: *Advances in Neural Information Processing Systems*, Kaufmann, M.,
1005 pp. 950–957.
- 1006 Krysanova, V., Haberlandt, U., 2002. Assessment of nitrogen leaching from
1007 arable land in large river basins. Part I: simulation experiments using a
1008 process-based model. *Ecological Modelling* 150, 255–275.

- 1009 Leip, A., Achermann, B. and Billen, G., Bleeker, A., Bouwman, L., de Vries,
1010 W., Dragosits, U., Döring, U., Fernall, D., Geupel, M., Johnes, P., Le Gall,
1011 A.C., Monni, S., Nevečerál, R., Orlandini, L., Prud'homme, M., Reuter,
1012 H., Simpson, D., Seufert, G., Spranger, T., Sutton, M., van Aardenne, J.,
1013 Voss, M., Winiwarter, W., 2011a. Integrating nitrogen fluxes at the Euro-
1014 pean scale, in: Sutton, M., Howard, C., Erisman, J., Billen, G., Bleeker,
1015 A., van Grinsven, H., Grennfelt, P., Grizzetti, B. (Eds.), European Nitro-
1016 gen Assessment. Cambridge University Press, New York, USA. chapter 16,
1017 pp. 345–376.
- 1018 Leip, A., Britz, W., De Vries, W., Weiss, F., 2011b. Farm, land, and soil
1019 nitrogen budgets for agriculture in Europe. Environmental Pollution In
1020 Press.
- 1021 Leip, A., Busto, M., Winiwarter, W., 2011c. Developing stratified N₂O
1022 emission factors for Europe. Environmental Pollution In Press.
- 1023 Leip, A., Dämmgen, U., Kuikman, P., van Amstel, A., 2005. The quality of
1024 European (EU-15) greenhouse gas inventories from agriculture. Environ-
1025 mental Sciences 2, 177–192.
- 1026 Leip, A., Marchi, G., Koebler, R., Kempen, M., Britz, W., Li, C., 2008.
1027 Linking an economic model for european agriculture with a mechanistic
1028 model to estimate nitrogen and carbon losses from arable soils in europe.
1029 Biogeosciences 5, 73–94.
- 1030 Li, C., 2000. Modeling trace gas emissions from agricultural ecosystems.
1031 Nutrient Cycling in Agroecosystems 58, 259–273.
- 1032 Li, C., Frohling, S., Frohling, T., 1992. Model of nitrous oxide evolution from
1033 soil driven by rainfall events: 1, model structure and sensitivity. Journal
1034 of Geophysical Research 97, 9759–9776.
- 1035 Lin, Y., Zhang, H., 2006. Component selection and smoothing in smoothing
1036 spline analysis of variance models. Annals of Statistics 34, 2272–2297.
- 1037 Lophaven, S., Nielsen, H., Sondergaard, J., 2002. DACE - A Matlab kriging
1038 toolbox, Version 2.0. Technical Report IMM-TR-2002-12. Informatics and
1039 Mathematical Modelling, Technical University of Denmark. DK-2800 Kgs.
1040 Lyngby, Denmark.

- 1041 Matson, P., Parton, W., Power, A., Swift, M., 1997. Agricultural intensifi-
1042 cation and ecosystem properties. *Science* 277, 504–509.
- 1043 Mattera, D., Haykin, S., 1998. Support vector machines for dynamic re-
1044 construction of a chaotic system, in: Schölkopf, B., Burges, C., Smola,
1045 A. (Eds.), *Advances in Kernel Methods: Support Vector Learning*. MIT
1046 Press, Cambridge, MA, USA, pp. 209–241.
- 1047 Meckesheimer, M., Booker, A., Barton, R., Simpson, T., 2002. Computa-
1048 tionally inexpensive metamodel assessment strategies. *AIAA Journal* 40,
1049 2053–2060.
- 1050 Montès, P., 1994. Smoothing noisy data by kriging with nugget effects, in:
1051 Laurent, P., Le Méhauté, A., Shumaker, L., Peters, A. (Eds.), *Wavelets,*
1052 *Images and Surface Fitting*. Wellesley, MA, USA. chapter AK Peters, pp.
1053 371–378.
- 1054 Ng, C., Young, P., 1990. Recursive estimation and forecasting of non-
1055 stationary time series. *Journal of Forecasting* 9, 173–204.
- 1056 Oenema, O., Oudendag, D., Velthof, G., 2007. Nutrient losses from manure
1057 management in the European Union. *Livestock Science* 112, 261–272.
- 1058 Oenema, O., Witzke, H., Klimont, Z., Lesschen, J., Velthof, G., 2009. Inte-
1059 grated assessment of promising measures to decrease nitrogen losses from
1060 agriculture in EU-27. *Agriculture, Ecosystems and Environment* 133, 280–
1061 288.
- 1062 Parkin, T., 1987. Soil microsites as a source of denitrification variability. *Soil*
1063 *Science Society of America Journal* 51, 1194–1199.
- 1064 Pineros Garcet, J., Ordonez, A., Roosen, J., Vanclouster, M., 2006. Meta-
1065 modeling: theory, concepts and application to nitrate leaching modelling.
1066 *Ecological Modelling* 193, 629–644.
- 1067 Power, A., 2010. Ecosystem services and agriculture: tradeoffs and synergies.
1068 *Philosophical Transactions of the Royal Society B* 365, 2959–2971.
- 1069 Ratto, M., Pagano, A., 2010. Recursive algorithms for efficient identification
1070 of smoothing spline anova models. *Advances in Statistical Analysis* 94,
1071 367–388.

- 1072 Ratto, M., Pagano, A., Young, P.C., 2007. State dependent parameter meta-
1073 modelling and sensitivity analysis. *Computer Physics Communications*
1074 177, 863–876.
- 1075 Ripley, B., 1994. Neural networks and related methods for classification.
1076 *Journal of the Royal Statistical Society, Series B* 56, 409–456.
- 1077 Sacks, J., Welch, W., Mitchell, T., Wynn, H., 1989. Design and analysis of
1078 computer experiments. *Statistical Science* 4, 409–435.
- 1079 Santner, T., Williams, B., Notz, W., 2003. *The Design and Analysis of*
1080 *Computer Experiments*. Springer, New York, NY, USA.
- 1081 Scherr, S., Sthapit, S., 2009. Farming and land use to cool the planet, in: In-
1082 stitute, T.W. (Ed.), *State of World 2009*. W.W. Norton & Co., Washington
1083 DC, USA. chapter 3, pp. 30–49.
- 1084 Schwarz, G., 1978. Estimating the dimension of a model. *Annals of Statistics*
1085 6, 461–464.
- 1086 Schweppe, F., 1965. Evaluation of likelihood functions for Gaussian signals.
1087 *IEEE Transactions on Information Theory* 11, 61–70.
- 1088 Shaffer, M., Ma, L., 2001. Carbon and nitrogen dynamics in upland soils,
1089 in: Shaffer, M., Ma, L., Hansen, S. (Eds.), *Modeling Carbon and Nitrogen*
1090 *Dynamics for Soil Management*. CRC Press LLC, Bacon Raton, Florida,
1091 USA. chapter 2, pp. 11–27.
- 1092 Simpson, T., Peplinski, J., Koch, P., Allen, J., 2001. Metamodels for
1093 computer-based engineering design: survey and recommendations. *En-*
1094 *gineering with Computers* 17, 129–150.
- 1095 Singh, R., 2000. Environmental consequences of agricultural development: a
1096 case study from the green revolution state of Haryana, India. *Agriculture,*
1097 *Ecosystems and Environment* 82, 97–103.
- 1098 Steinwart, I., Christman, A., 2008. *Support Vector Machines*. *Information*
1099 *Science and Statistics*.
- 1100 Steinwart, I., Christmann, A., 2008. *Support Vector Machines*. *Information*
1101 *Science and Statistics*, Springer.

- 1102 Storlie, C., Bondell, H., Reich, B., Zhang, H., 2011. Surface estimation,
1103 variable selection, and the nonparametric oracle property. *Statistica Sinica*
1104 21, 679–705.
- 1105 Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal*
1106 *of the Royal Statistical Society, series B* 58, 267–288.
- 1107 Tilman, D., Cassman, K., Matson, P., Naylor, R., Polasky, S., 2002. Agri-
1108 cultural sustainability and intensive production practices. *Nature* 418,
1109 671–677.
- 1110 Vanclooster, M., Viaene, P., Christiaens, K., Ducheyne, S., 1996. Wave: a
1111 mathematical model for simulating water and agrochemicals in the soil and
1112 vadose environment, reference and user’s manual (release 2.1). Technical
1113 Report. Institute for Land and Water Management, Katholieke Univer-
1114 siteit Leuven. Leuven, Belgium.
- 1115 Vapnik, V., 1995. *The Nature of Statistical Learning Theory*. Springer Verlag,
1116 New York, USA.
- 1117 Vapnik, V., 1998. *Statistical Learning Theory*. Wiley, New York, USA.
- 1118 van der Velde, M., Bouraoui, F., de Vries, W., 2009. Derivation of EPIC
1119 meta-models to evaluate the impact of cross compliance measures on leach-
1120 ing and runoff of nitrogen and soil erosion at European scale. CCAT De-
1121 liverable (report) 4.2.3.2.. CCAT project.
- 1122 Velthof, G., Oudendag, D., Witzke, H., Asman, W., Klimont, Z., Oenema,
1123 O., 2009. Integrated assessment of nitrogen emissions from agriculture in
1124 EU-27 using MITERRA-EUROPE. *Journal of Environmental Quality* 38,
1125 1–16.
- 1126 Venables, W., Ripley, B., 2002. *Modern Applied Statistics with S-PLUS*.
1127 Springer, New York, USA.
- 1128 Wahba, G., 1990. *Spline Models for Observational Data*. Society for Indus-
1129 trial and Applied Mathematics, Philadelphia, Pennsylvania, USA.
- 1130 Wang, G., Shan, S., 2007. Review of metamodeling techniques in support
1131 of engineering design optimization. *Journal of Mechanical Design* 129,
1132 370–380.

- 1133 Webb, J., Menzi, H., Pain, B., Misselbrook, T., Dämmgen, U., Hendriks, H.,
1134 Döhler, H., 2005. Managing ammonia emissions from livestock production
1135 in Europe. *Environmental Pollution* 135, 399–406.
- 1136 Werbos, P., 1974. Beyond regression: New tools for prediction and analysis
1137 in behavior sciences. Ph.D. thesis. Committee on Applied Mathematics,
1138 Harvard University. Cambridge, MA, USA.
- 1139 Young, P., 1999. Nonstationary time series analysis and forecasting. *Progress*
1140 *in Environmental Science* 1, 3–48.
- 1141 Young, P., 2000. Stochastic, dynamic modelling and signal processing: time
1142 variable and state dependent parameter estimation, in: Fitzgerald, W.,
1143 Walden, A., Smith, R., Young, P. (Eds.), *Nonlinear and Nonstationary*
1144 *Signal Processing*. Cambridge University Press, Cambridge, USA, pp. 74–
1145 114.
- 1146 Young, P., 2001. The identification and estimation of nonlinear stochastic sys-
1147 tems, in: Mees, A. (Ed.), *Nonlinear Dynamics and Statistics*. Birkhauser,
1148 USA, Boston, USA, pp. 127–166.
- 1149 Young, P., Ng, C., 1989. Variance intervention. *Journal of Forecasting* 8,
1150 399–416.
- 1151 Zhang, Y., Li, C., Zhou, X., Moore, B., 2002. A simulation model linking crop
1152 growth and soil biogeochemistry for sustainable agriculture. *Ecological*
1153 *Modelling* 151, 75–108.

1154 **List of Figures**

1155	1	Flow of data through the DNDC-EUROPE model (M) and	
1156		relationship with the metamodel's one (MM). The input	
1157		variables of the metamodels were selected from the original	
1158		DNDC-EUROPE dataset (screening). The estimated (*) out-	
1159		put were compared with the detailed model's output during	
1160		the training and test phases to improve the metamodel and to	
1161		evaluate the goodness of the approximation.	47
1162	2	R^2 evolution in function of the size of the train set (log scale)	
1163		for N_2O prediction	48
1164	3	R^2 evolution in function of the size of the train set (log scale)	
1165		for N leaching prediction	49
1166	4	Importance measure for each variable in the case of (left) N_2O	
1167		prediction with the full training dataset and random forest and	
1168		of (right) N leaching prediction with the full training dataset	
1169		and SVM (For the meaning of the acronyms, please refer to	
1170		Section 4.2)	50

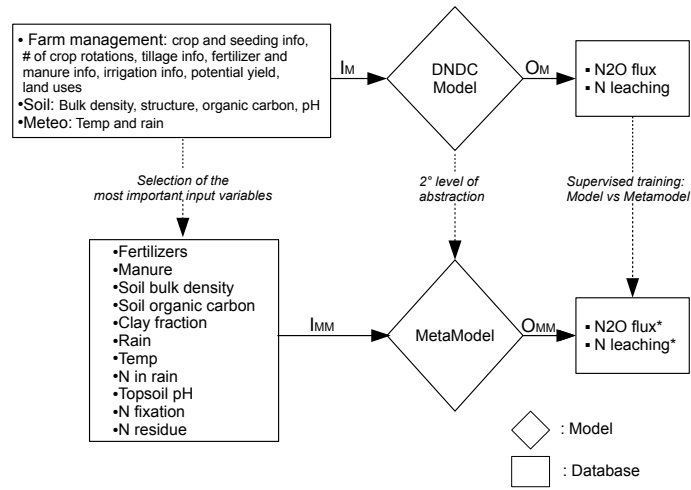


Figure 1: Flow of data through the DNDC-EUROPE model (M) and relationship with the metamodel's one (MM). The input variables of the metamodels were selected from the original DNDC-EUROPE dataset (screening). The estimated (*) output were compared with the detailed model's output during the training and test phases to improve the metamodel and to evaluate the goodness of the approximation.

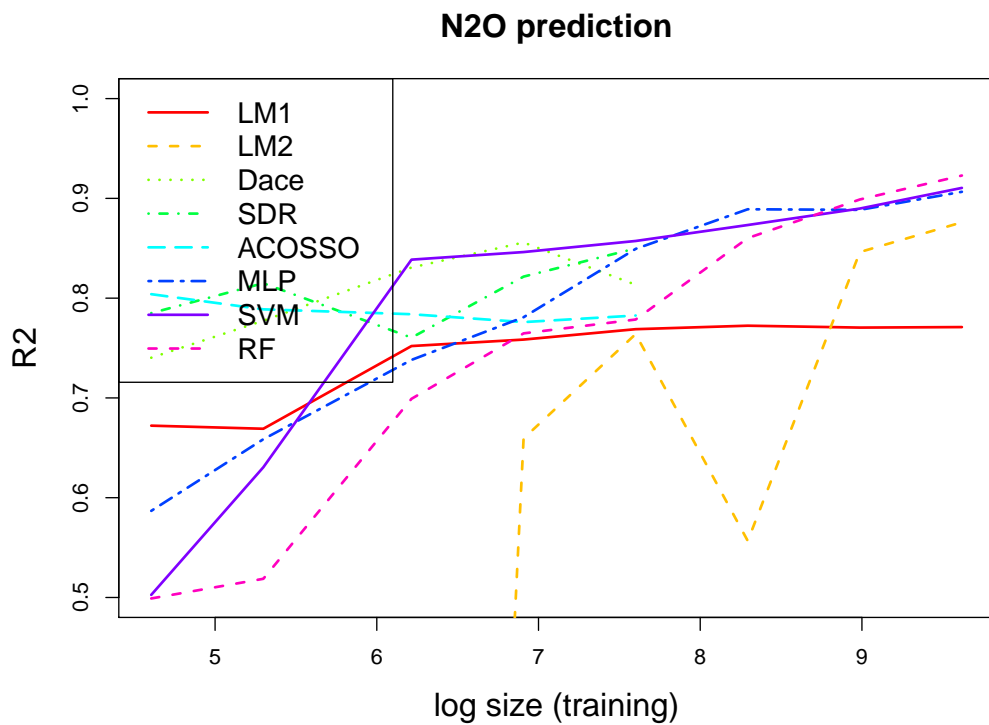


Figure 2: R^2 evolution in function of the size of the train set (log scale) for N_2O prediction

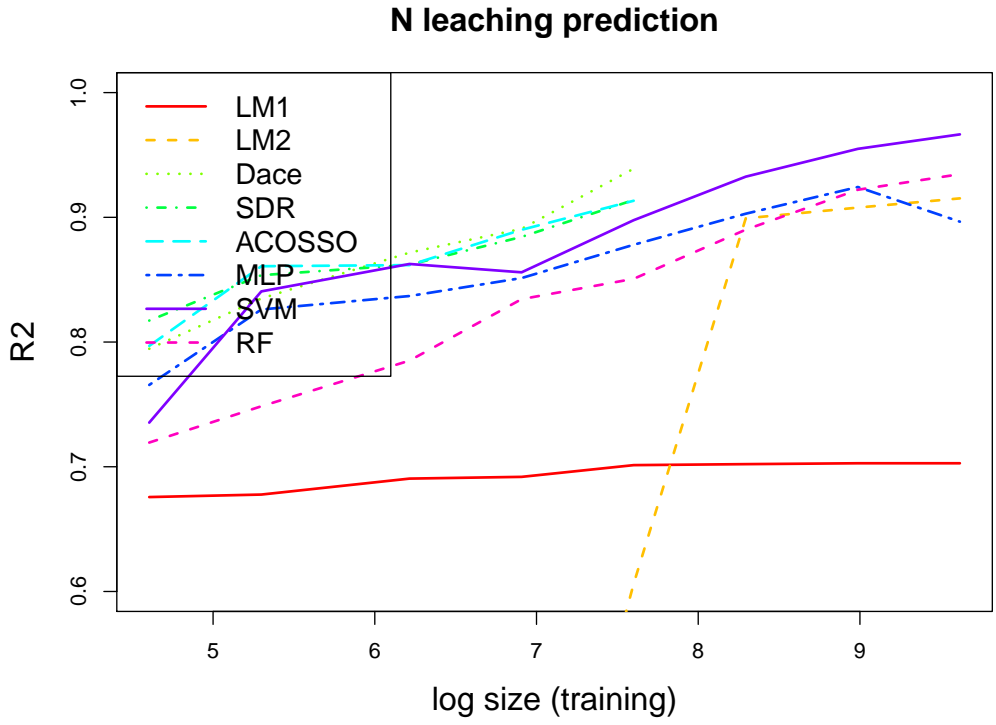


Figure 3: R^2 evolution in function of the size of the train set (log scale) for N leaching prediction

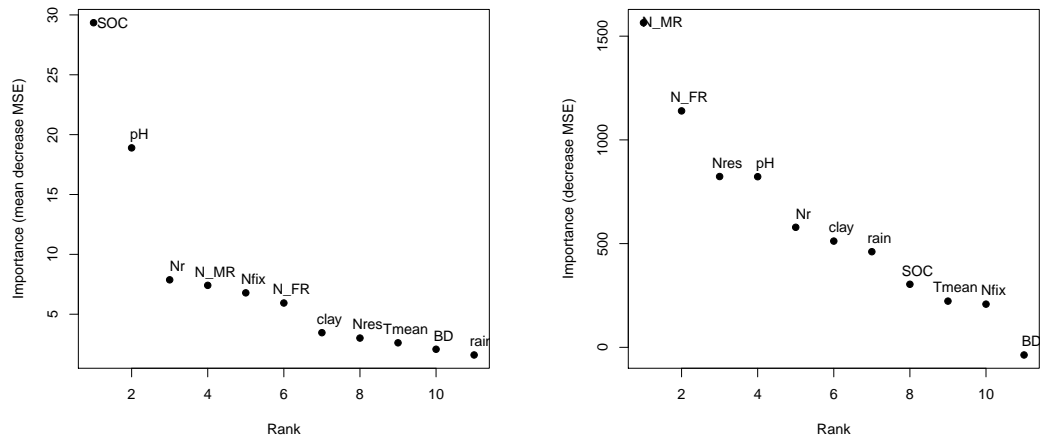


Figure 4: Importance measure for each variable in the case of (left) N₂O prediction with the full training dataset and random forest and of (right) N leaching prediction with the full training dataset and SVM (For the meaning of the acronyms, please refer to Section 4.2)

1171 **List of Tables**

1172 1 Summary of the main features for training, validation (hyper-
1173 parameters tuning) and test steps of each method. 25
1174 2 R^2 (first line) and MSE (second line) on the test set for each
1175 method and various sizes of the training dataset for N₂O pre-
1176 diction. For each size, the best R^2 is in bold. ✓ corresponds
1177 to cases impossible to train, either because the model is over-
1178 specified (more parameters to estimate than the number of
1179 observations: LM2) or because the training size is too large
1180 for the method to be used (Dace/SDR/Acosso) 53
1181 3 R^2 (first line) and MSE (second line) on the test set for each
1182 method and various sizes of the training dataset for N leaching
1183 prediction. For each size, the best R^2 is in bold. ✓ corresponds
1184 to cases impossible to train, either because the model is over-
1185 specified (more parameters to estimate than the number of
1186 observations: LM2) or because the training size is too large
1187 for the method to be used (Dace/SDR/Acosso) 54
1188 4 Standard deviation (first line) and maximum (second line) of
1189 the squared errors on the test set for each method and vari-
1190 ous sizes of the training dataset for N₂O prediction. For each
1191 size, the minimal standard deviation and the minimal value
1192 of the maxima are in bold. ✓ corresponds to cases impossible
1193 to train, either because the model is over-specified (more pa-
1194 rameters to estimate than the number of observations: LM2)
1195 or because the training size is too large for the method to be
1196 used (Dace/SDR/Acosso) 55
1197 5 Standard deviation (first line $\times 10^3$) and maximum (second
1198 line $\times 10^3$) of the squared errors on the test set for each method
1199 and various sizes of the training dataset for N leaching predic-
1200 tion. For each size, the minimal standard deviation and the
1201 minimal value of the maxima are in bold. ✓ corresponds to
1202 cases impossible to train, either because the model is over-
1203 specified (more parameters to estimate than the number of
1204 observations: LM2) or because the training size is too large
1205 for the method to be used (Dace/SDR/Acosso) 56

1206	6	Number of variables selected by AIC stepwise procedure in	
1207		LM2 for N ₂ O prediction and N leaching prediction in function	
1208		of the training dataset size	57
1209	7	Number of ANOVA components selected by the COSSO	
1210		penalty in ACOSSE and SDR for N ₂ O prediction and N leach-	
1211		ing prediction as a function of the training dataset size.	58
1212	8	Approximative time for training from about 15 000 observa-	
1213		tions (first line) and for predicting about 19 000 observations	
1214		(second line) on a desktop computer (Processor 2GHz, 1.5GO	
1215		RAM). In the case of SDR, ACOSSE and DACE we report the	
1216		time for training using samples with 2 000 model runs because	
1217		the method can not be used for largest training datasets.	59

Size of the dataset	LM1	LM2	Dace	SDR	Acosso	MLP	SVM	RF
100	67.22% 11.50	✓ ✓	74.03% 9.11	78.50% 7.54	80.40% 6.88	58.68% 14.50	50.26% 17.45	49.90% 17.57
200	66.91% 11.61	-13 093% 4 626	77.74% 7.81	81.50% 6.49	78.88% 7.41	65.86% 11.98	63.05% 12.96	51.87% 16.89
500	75.20% 8.70	-163% 92.35	83.07% 5.94	76.04% 8.41	78.39% 7.58	73.81% 9.19	83.86% 5.66	69.91% 10.56
1 000	76.85% 8.47	65.94% 11.95	85.58% 5.06	82.16% 6.26	77.60% 7.86	78.81% 7.69	84.62% 5.40	76.47% 8.25
2 000	76.89% 8.11	76.40% 8.28	81.34% 6.55	84.16% 5.27	78.26% 7.63	84.94% 5.28	85.73% 5.01	77.86% 7.77
4 000	77.24% 7.99	55.67% 15.55	✓ ✓	✓ ✓	✓ ✓	88.91% 3.89	87.33% 4.45	86.01% 4.90
8 000	77.05% 8.05	84.62% 5.40	✓ ✓	✓ ✓	✓ ✓	88.85% 3.91	88.98% 3.86	89.89% 3.55
≈ 15 000	77.10% 8.03	87.60% 3.28	✓ ✓	✓ ✓	✓ ✓	90.66% 3.28	91.05% 3.14	92.29% 2.71

Table 2: R^2 (first line) and MSE (second line) on the test set for each method and various sizes of the training dataset for N_2O prediction. For each size, the best R^2 is in bold. ✓ corresponds to cases impossible to train, either because the model is over-specified (more parameters to estimate than the number of observations: LM2) or because the training size is too large for the method to be used (Dace/SDR/Acosso)

Size of the dataset	LM1	LM2	Dace	SDR	Acosso	MLP	SVM	RF
100	67.57% 1 742	✓ ✓	79.46% 1 103	81.72% 982	79.69% 1 091	76.56% 1 259	73.54% 1 421	71.94% 1 507
200	67.77% 1 731	-2 086% > 10 ⁶	83.49% 887	85.36% 786	86.08% 747	82.61% 934	84.06% 856	74.85% 1 351
500	69.05% 1 662	36.92% 3 388	87.17% 689	86.20% 741	86.17% 743	83.69% 876	86.26% 738	78.51% 1 154
1 000	69.19% 1 655	27.24% 3 908	89.08% 587	88.43% 621	89.00% 591	85.13% 799	85.59% 774	83.44% 889
2 000	70.13% 1 604	60.62% 2 115	93.90% 328	91.39% 462	91.33% 466	84.94% 655	89.77% 549	85.07% 802
4 000	70.21% 1 600	89.92% 541	✓ ✓	✓ ✓	✓ ✓	93.26% 521	87.33% 362	89.01% 590
8 000	70.28% 1 596	90.78% 495	✓ ✓	✓ ✓	✓ ✓	92.43% 406	95.49% 242	92.21% 418
≈ 15 000	70.28% 1 596	91.52% 455	✓ ✓	✓ ✓	✓ ✓	89.65% 556	96.65% 180	93.46% 351

Table 3: R^2 (first line) and MSE (second line) on the test set for each method and various sizes of the training dataset for N leaching prediction. For each size, the best R^2 is in bold. ✓ corresponds to cases impossible to train, either because the model is over-specified (more parameters to estimate than the number of observations: LM2) or because the training size is too large for the method to be used (Dace/SDR/Acosso)

Size of the dataset	LM1	LM2	Dace	SDR	Acosso	MLP	SVM	RF
100	80.4	✓	72.7	52.4	50.2	125.5	159.6	150.0
	2 400	✓	2 319	1 845	1 597	2 911	3 816	3 538
200	84.5	$> 10^5$	68.1	52.3	64.6	100.3	113.7	145.4
	2 461	$> 10^6$	2 207	1 915	2 098	2 534	2 636	3 352
500	59.3	1 472.9	49.6	74.0	60.2	84.9	42.5	99.1
	2 027	48 769	1 928	2 589	2 303	2 172	1 753	2718
1 000	56.9	203.5	48.6	51.0	63.4	53.9	48.5	77.7
	1 980	8 384	1 643	1 633	2 065	1 888	1 874	2 348
2 000	50.3	81.5	66.7	37.8	62.9	38.4	41.6	70.4
	1 826	2 890	2 456	1 212	3 000	1 039	1 663	2 421
4 000	46.1	539.2	✓	✓	✓	33.0	37.6	52.8
	1 711	32 290	✓	✓	✓	1 110	1 519	2 040
8 000	42.2	60.9	✓	✓	✓	31.0	43.2	38.3
	1 564	2 846	✓	✓	✓	1 072	1 773	1 645
$\simeq 15\ 000$	42.2	29.0	✓	✓	✓	29.0	35.7	25.6
	1 568	1 339	✓	✓	✓	1 339	1 833	807

Table 4: Standard deviation (first line) and maximum (second line) of the squared errors on the test set for each method and various sizes of the training dataset for N₂O prediction. For each size, the minimal standard deviation and the minimal value of the maxima are in bold. ✓ corresponds to cases impossible to train, either because the model is over-specified (more parameters to estimate than the number of observations: LM2) or because the training size is too large for the method to be used (Dace/SDR/Acosso)

Size of the dataset	LM1	LM2	Dace	SDR	Acosso	MLP	SVM	RF
100	6.11	✓	5.83	6.45	8.14	6.72	9.79	7.99
	173.1	✓	180.5	177.0	241.1	238.5	367.7	275.7
200	6.26	$> 10^4$	5.24	7.15	8.61	6.28	5.36	7.95
	184.7	$> 10^5$	152.4	290.6	279.5	213.8	146.8	284.6
500	6.99	45.7	7.34	7.38	8.62	6.89	6.77	7.83
	204.1	1 427.7	238.2	280.0	280.9	213.8	302.8	290.7
1 000	7.37	82.3	7.64	7.10	8.90	7.72	10.24	7.47
	220.9	4 090.4	270.6	239.8	255.3	289.0	358.1	291.1
2 000	5.91	71.9	2.66	3.15	9.13	5.74	6.63	5.53
	177.4	4 309.1	96.6	113.3	128.7	225.9	320.6	212.7
4 000	5.71	4.94	✓	✓	✓	3.50	3.61	4.51
	167.0	213.5	✓	✓	✓	134.5	123.1	218.2
8 000	5.59	4.31	✓	✓	✓	2.80	2.38	2.60
	162.0	161.8	✓	✓	✓	77.8	77.4	70.4
$\simeq 15\ 000$	5.53	2.54	✓	✓	✓	4.74	1.35	3.00
	157.2	72.1	✓	✓	✓	147.0	36.1	128.7

Table 5: Standard deviation (first line $\times 10^3$) and maximum (second line $\times 10^3$) of the squared errors on the test set for each method and various sizes of the training dataset for N leaching prediction. For each size, the minimal standard deviation and the minimal value of the maxima are in bold. ✓ corresponds to cases impossible to train, either because the model is over-specified (more parameters to estimate than the number of observations: LM2) or because the training size is too large for the method to be used (Dace/SDR/Acosso)

Training dataset size	Number of selected variables (N ₂ O prediction)	Number of selected variables (N leaching prediction)
200	79	95
500	74	84
1 000	75	89
2 000	79	94
4 000	94	95
8 000	98	97
≈ 15 000	96	100

Table 6: Number of variables selected by AIC stepwise procedure in LM2 for N₂O prediction and N leaching prediction in function of the training dataset size

Training dataset size	Number of selected variables (N ₂ O prediction)		Number of selected variables (N leaching prediction)	
	ACOSSO	SDR	ACOSSO	SDR
100	30	23	24	39
200	13	17	31	26
500	17	32	18	19
1 000	7	9	28	31
2 000	9	10	29	30

Table 7: Number of ANOVA components selected by the COSSO penalty in ACOSSO and SDR for N₂O prediction and N leaching prediction as a function of the training dataset size.

Use	LM1	LM2	Dace	SDR	Acosso	MLP	SVM	RF
Train	<1 s.	50 min	80 min	4 hours	65 min	2.5 hours	5 hours	15 min
Prediction	<1 s.	<1 s.	90 s.	14 min	4 min.	1 s.	20 s.	5 s.

Table 8: Approximative time for training from about 15 000 observations (first line) and for predicting about 19 000 observations (second line) on a desktop computer (Processor 2GHz, 1.5GO RAM). In the case of SDR, ACOSSO and DACE we report the time for training using samples with 2 000 model runs because the method can not be used for largest training datasets.