

Beyond prediction: methods for interpreting complex models of soil variation

Alexandre M.J.-C. Wadoux^{a,*}, Christoph Molnar^b

^a Sydney Institute of Agriculture & School of Life and Environmental Sciences, The University of Sydney, Australia

^b Johner Institute GmbH, Konstanz, Germany

ARTICLE INFO

Keywords:

Machine learning
Shapley
Partial dependence
H-statistic
Accumulated local effect
Surrogate modelling

ABSTRACT

Understanding the spatial variation of soil properties is central to many sub-disciplines of soil science. Commonly in soil mapping studies, a soil map is constructed through prediction by a statistical or non-statistical model calibrated with measured values of the soil property and environmental covariates of which maps are available. In recent years, the field has gradually shifted attention towards more complex statistical and algorithmic tools from the field of machine learning. These models are particularly useful for their predictive capabilities and are often more accurate than classical models, but they lack interpretability and their functioning cannot be readily visualized. There is a need to understand how these models can be used for purposes other than making accurate prediction and whether it is possible to extract information on the relationships among variables found by the models. In this paper we describe and evaluate a set of methods for the interpretation of complex models of soil variation. An overview is presented of how model-independent methods can serve the purpose of interpreting and visualizing different aspects of the model. We illustrate the methods with the interpretation of two mapping models in a case study mapping topsoil organic carbon in France. We reveal the importance of each driver of soil variation, their interaction, as well as the functional form of the association between environmental covariate and the soil property. Interpretation is also conducted locally for an area and two spatial locations with distinct land use and climate. We show that in all cases important insights can be obtained, both into the overall model functioning and into the decision made by the model for a prediction at a location. This underpins the importance of going beyond accurate prediction in soil mapping studies. Interpretation of mapping models reveal how the predictions are made and can help us formulating hypotheses on the underlying soil processes and mechanisms driving soil variation.

1. Introduction

Understanding the spatial variation of soil properties has become central to many sub-disciplines of soil science. Digital soil mapping (DSM) techniques can be used for this purpose. Commonly in DSM studies, statistical or non-statistical models are calibrated to exploit the quantitative relationship between measured values of a soil property and a set of environmental covariates of which maps are available, such as satellite imagery and terrain attributes. These models are used to predict the soil property at unobserved locations and to identify and expose the importance of environmental factors in the soil property spatial variation. Recent examples of studies using this approach are [Quist et al. \(2019\)](#) for mapping soil nematodes and [Heuvelink et al. \(2021\)](#) for mapping soil organic carbon in space and time.

Since early soil mapping studies rooted in classical statistics and design-based inference in the 70s, and based on geostatistics in the 80s ([Heuvelink and Webster, 2001](#)), the field has gradually shifted attention towards more complex statistical and algorithmic tools from the field of machine learning. Accuracy of such models is often higher than that of classical models. They are also particularly useful in situation where the relationship between the soil property and environmental covariates is too complex to be modelled mechanistically or with simple statistical models. However, popularization of complex models of soil variation was made at the expense of understanding why the soil varies the way it does. Insight into the functioning and structure of the models are difficult to obtain, so that these models are often referred to as “black boxes”. Examples of such models are random forest, support vector machines and neural networks. We refer to [Hastie et al. \(2009\)](#) for an overview.

* Corresponding author at: Sydney Institute of Agriculture & School of Life and Environmental Sciences, The University of Sydney, New South Wales, Australia.
E-mail address: alexandre.wadoux@sydney.edu.au (A.M.J.-C. Wadoux).

In soil science, several attempts were made to obtain insights from complex models. The relative effect of environmental covariates on model prediction is usually characterized by model-specific variable importance statistics such as through the mean decrease in impurity for tree-based models (as is done in Vos et al., 2019 for example), or by calculating the partial dependence of the prediction to environmental covariates (e.g. Zeng et al., 2017; Ottoy et al., 2017). For artificial neural networks, the Garson's algorithm or the magnitude and direction of the connection weights between neurons give indication on the variable importance (Olden and Jackson, 2002). An example in soil mapping is Rivera and Bonilla (2020). While valid and useful to obtain insights into complex models of soil variation, these methods are model-specific, i.e. they preclude comparison between models (Wadoux et al., 2020). A number of "model-agnostic" or model-independent interpretation methods have recently been developed outside soil science, in the statistical and machine learning literature. This is evidenced by textbooks specifically addressing model interpretation (e.g. Molnar, 2020; Biecek and Burzykowski, 2021). Model-independent means that these model interpretation methods are applicable to any model. It is worthwhile to introduce these recent developments, and to present a strategy for the interpretation of complex soil mapping models. This was also recently highlighted as one of the most pressing pedometric research topics (Challenge 3, Wadoux et al., 2021).

At the higher level, one may distinguish between local and global model interpretation (Molnar, 2020). For mapping purpose, a local interpretation is appropriate when the objective is to evaluate how prediction to a single spatial location is made. It is indeed sensible to assume that the importance of certain environmental factors vary from one location to another, and between regions. A global interpretation, conversely, provides insights into the overall model functioning. Global methods expose the importance of each driver of spatial variation, their interaction, as well as the functional form of the association between environmental covariate and the soil property. In practice global and local methods are used jointly to interpret and visualize differentiable aspects of the model.

This paper is structured as follows. A first part introduces local and global interpretation methods for use in mapping studies. Such methods can be applied to any model (i.e. they are model-independent), although in practice it is not always sensible to apply them on simple models whose structure is readily interpretable (e.g. linear regression). The second part of the paper illustrates the methods for the interpretation of two models in a case study mapping topsoil organic carbon in France. Finally, we discuss in a third part the limitations of interpretation methods, possible alternatives, and summarize the utility of the methods as well as their pros and cons in a table.

2. Interpretation methods

Consider the soil property of interest Y modelled at any location s in the study area \mathcal{A} by $Y = f(X) + \varepsilon$, where f is the regression function that yields Y given values of one or more dependent variables X and $\varepsilon \in \mathbb{R}$ is a random error. Statistical regression techniques seek to estimate the form of the function f to make a prediction $\hat{Y} = \hat{f}(X)$ where the statistical model \hat{f} is estimated by minimizing the expected squared error term $E[(Y - \hat{Y})^2]$.

Let $y(s_i)$ be n measurements of Y , s_i ($i = 1, \dots, n$; $s_i \in \mathcal{A}$) and $\mathbf{X} \in \mathbb{R}^{n \times p}$ be the covariate matrix of size $n \times p$ where p is the number of environmental predictors. We denote \mathbf{x}_i and \mathbf{x}_j the i^{th} row-vector and the j^{th} column vector of \mathbf{X} , respectively, and x_{ij} a scalar value at row i and column j . We make no assumption on the functional form of \hat{f} and treat it as a "black-box". Hereafter, we describe methods to interpret the calibrated regression model \hat{f} and illustrate them with the data and support vector machine model from Wadoux et al. (2022), Section 4.2.

2.1. Covariate importance with permutation

Covariate importance obtained by permutation is a popular method to quantify the relative importance of an individual covariate or of a group of covariates on model prediction. A covariate is important if perturbing its values affects model prediction error: the larger the change in prediction error, the more important is the covariate. Prediction error is quantified by the error function $\ell(\hat{f}(\mathbf{X}), \mathbf{y})$, where \mathbf{y} is the n vector of observations. Error function $\ell(\hat{f}(\mathbf{X}), \mathbf{y})$ is usually the root mean square error (RMSE) or modelling efficiency coefficient (MEC, Janssen and Heuberger, 1995). Covariate importance is estimated with the following steps (Breiman, 2001; Fisher et al., 2019):

1. Estimate error function $\ell(\hat{f}(\mathbf{X}), \mathbf{y})$.
2. For each covariate $j = 1, \dots, p$:
 - (a) Create modified (denoted by the asterisk $*$) covariate matrix \mathbf{X}^* by permutation of the values in the j^{th} column.
 - (b) Estimate error function from prediction made with the permuted covariate matrix $\ell(\hat{f}(\mathbf{X}^*), \mathbf{y})$.
 - (c) Obtain covariate importance for the j^{th} covariate by the ratio $\ell(\hat{f}(\mathbf{X}^*), \mathbf{y}) / \ell(\hat{f}(\mathbf{X}), \mathbf{y})$ or the difference $\ell(\hat{f}(\mathbf{X}^*), \mathbf{y}) - \ell(\hat{f}(\mathbf{X}), \mathbf{y})$.

Permutation of the covariate matrix involves randomness and is usually repeated to obtain a distribution of the importance metric. Fig. 1 shows an example of permutation covariate importance using the ratio of RMSE. The technique can be extended to measure the importance of group of covariates, by permuting the group of covariate simultaneously

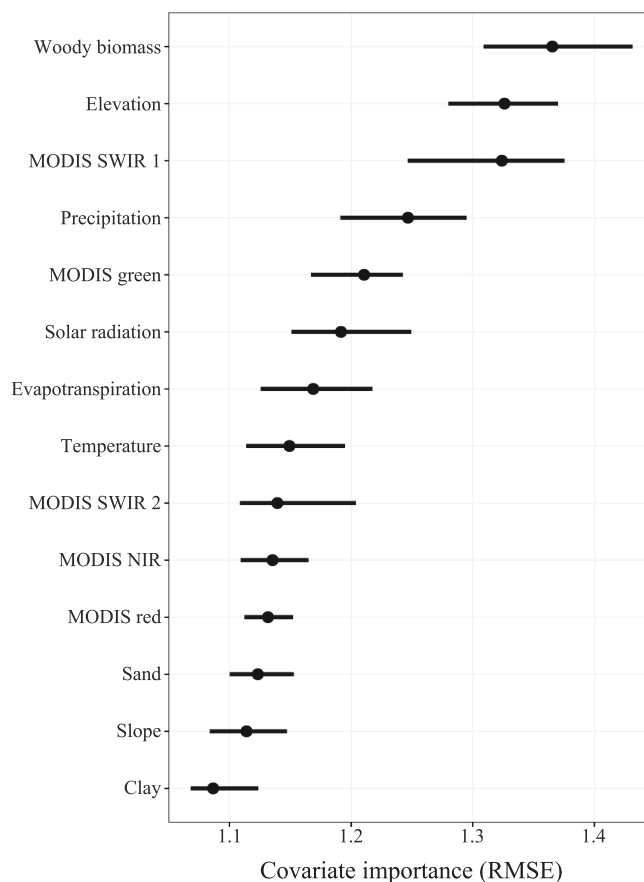


Fig. 1. Example of a covariate importance with permutation assessed by the ratio of RMSE. The black dots represent the mean value of 100 permutations and the lines the 90% confidence interval. Data and model from Wadoux et al. (2022), Section 4.2).

instead of a single covariate.

Calculation of covariate importance with permutation is computationally efficient as it does not require re-calibrating the model at each permutation. In case where covariates are dependent, however, the values obtained by permutation might be misleading and result in incorrect ranking of importance. In this situation, it is sensible to permute group of correlated covariates instead of each individual covariate. An extensive comparison of the impact that correlated covariates have on permutation importance is given by Hooker et al. (2021). When covariates are correlated, alternatives to permutation importance are the conditional permutation strategy (Strobl et al., 2007; Fisher et al., 2019; Watson and Wright, 2021) and the dropped variables importance (Lei et al., 2018).

2.2. Individual conditional expectation

Individual conditional expectation (ICE, Goldstein et al., 2015) plots shows one line per location indicating how the prediction for that location changes when a covariate changes. Each line in the ICE plot is computed by fixing the values of all covariates except the covariate of interest. The covariate of interest for that location is then iteratively replaced with different grid values along the range of the covariate of interest. Model predictions are computed for these newly created versions of the location. The results are then plotted as a line with covariate value on the x-axis and predicted value on the y-axis, representing the influence of that covariate on the prediction of the chose location. The ICE plot usually consists of many lines for different locations.

In statistical terms, consider the subset of covariates $\mathbf{X}_{\mathcal{J}}$ of \mathbf{X} composed of $l < p$ covariates, and $\mathbf{X}_{\mathcal{E}}$ its complement so that $f(\mathbf{X}) = f(\mathbf{X}_{\mathcal{J}}, \mathbf{X}_{\mathcal{E}})$. The subset $\mathbf{X}_{\mathcal{J}}$ usually contains one or two covariates (i.e. $l \approx 1, 2$). For any location in \mathcal{L} with covariate values $(\mathbf{x}_{i\mathcal{J}}, \mathbf{x}_{i\mathcal{E}})$ and calibrated model \hat{f} , an ICE curve \hat{f}_{ICE} shows model predictions for a grid of $\mathbf{x}_{i\mathcal{J}}$ while keeping fixed the values of $\mathbf{x}_{i\mathcal{E}}$ (Fig. 2a).

When comparing ICE curves, it is convenient to center the individual ICE curves to a baseline value. Without centering, it can be difficult to visually track differences in covariate effects. Centering makes the ICE curves comparable. The centered ICE curves show the partial dependence of the predicted value at a location to a covariate, expressed in terms of difference to the baseline value. The centered ICE curve is expressed as:

$$\text{centered } \hat{f}_{\text{ICE}} = \hat{f}_{\text{ICE}} - \hat{f}(x_0, \mathbf{x}_{i\mathcal{E}}), \quad (1)$$

where x_0 is the baseline value, usually the minimum, maximum or average of the values in the calibration dataset (Fig. 2b).

ICE curves are an intuitive way to explore the effect of covariates to individual spatial locations. ICE curves can further be computed for group of spatial locations within an area, and their average value (i.e.

their partial dependence plot, see also Section 2.3) compared to that of another area. This may provide insight into local or regional dependence to a covariate. However, ICE are also calculated from the marginal covariate distribution and are thus they are reliable only when covariates are independent. More information on this is provided in Section 2.3.

2.3. Partial dependence plots

Partial dependence plots (PDP) show how the model prediction behaves on average as a function of one or more covariates. This illustrates the effect of these covariates after averaging the effect of other covariates included in the model. The partial dependence function f_{PDP} of $\hat{f}(\mathbf{X})$ on $\mathbf{X}_{\mathcal{J}}$ for $\mathbf{x}_{\mathcal{J}}$ is formally expressed as the expected value of the model prediction over the distribution of the covariates in the subset \mathcal{E} (Friedman, 2001):

$$f_{\text{PDP}}(\mathbf{x}_{\mathcal{J}}) = \mathbb{E}_{\mathbf{X}_{\mathcal{E}}}[\hat{f}(\mathbf{x}_{\mathcal{J}}, \mathbf{X}_{\mathcal{E}})]. \quad (2)$$

In practice the numerical integration required to estimate the marginal distribution of $\mathbf{X}_{\mathcal{E}}$ is approximated by averaging over the n observation locations:

$$\hat{f}_{\text{PDP}}(\mathbf{x}_{\mathcal{J}}) = \frac{1}{n} \sum_{i=1}^n \hat{f}(\mathbf{x}_{\mathcal{J}}, \mathbf{x}_{i\mathcal{E}}), \quad (3)$$

where $\mathbf{x}_{1\mathcal{E}}, \mathbf{x}_{2\mathcal{E}}, \dots, \mathbf{x}_{n\mathcal{E}}$ are the row-vectors of $\mathbf{X}_{\mathcal{E}}$. Eq. 3 shows that the PDP of the calibration dataset is the average of the n ICE curves. Accordingly, Fig. 2a-b show the PDP of woody biomass on SOC as average of the n ICE curves. Fig. 2c is an example of two-dimensional PDP (i.e. for $l = 2$).

PDP are easy to implement and represent an intuitive way of interpreting a model. While PDP can be computed for subset \mathcal{J} of any size, only one or two covariates can reasonably be displayed. Note also that dependence among covariates in $\mathbf{X}_{\mathcal{J}}$ and $\mathbf{X}_{\mathcal{E}}$ can produce a PDP that is misleading. When covariates are dependent, taking the marginal expectation of one covariate leads to consider points that lie outside the multivariate joint distribution. We recommend testing independence using, for example, a combination of scatter plots and statistics such as the Spearman's rank correlation coefficient. The accumulated local effect (Section 2.4) is a sensible alternative to the PDP when covariates are dependent. Both marginal (Eq. 2) and conditional expectations are the same if covariates in $\mathbf{X}_{\mathcal{J}}$ and $\mathbf{X}_{\mathcal{E}}$ are uncorrelated (p. 370 in Hastie et al., 2009).

2.4. Accumulated local effect

An alternative to the PDP when covariates are dependent is the accumulated local effect (ALE, Apley and Zhu, 2020). The ALE shows the effect of changing the values of a covariate on the soil property.

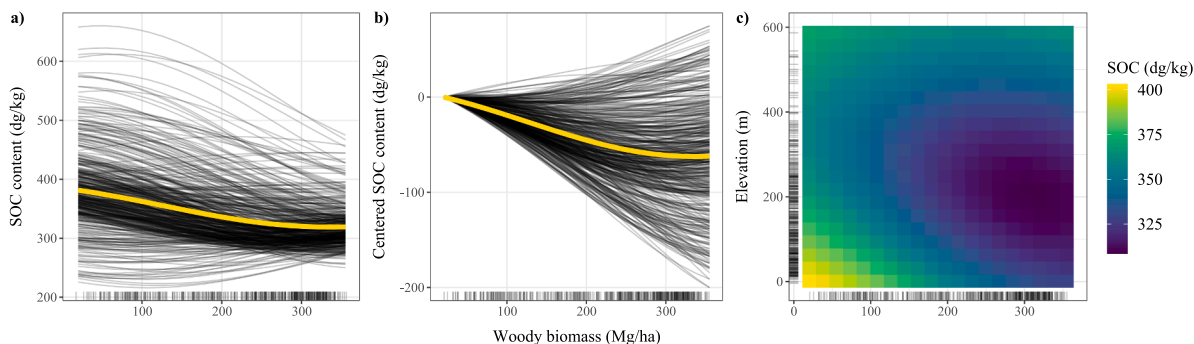


Fig. 2. Examples of a) individual conditional expectation (ICE) curves (in black) for woody biomass against soil organic carbon (SOC) content. The yellow curve is the partial dependence plot (PDP). In b), ICE curves and the PDP are centered at the minimum of the covariate value (i.e. at a woody biomass value of 20). Plot c) shows the two-dimensional PDP of woody biomass against elevation. Data and model from Wadoux et al. (2022), Section 4.2.

Formally, the ALE is defined as the accumulated derivative of the prediction function over the conditional distribution of the soil property, starting at the lower bound $z_{0,\mathcal{S}}$.

$$f_{ALE}(\mathbf{x}_{\mathcal{S}}) = \int_{z_{0,\mathcal{S}}}^{\mathbf{x}_{\mathcal{S}}} E_{X_{\mathcal{S}}|X_{\mathcal{S}}'}[\hat{f}^{\mathcal{S}}(X_{\mathcal{S}}, X_{\mathcal{S}}')|X_{\mathcal{S}} = z_{\mathcal{S}}] dz_{\mathcal{S}}, \quad (4)$$

where $\hat{f}^{\mathcal{S}}(\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}'}) = \frac{\partial \hat{f}(\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}'})}{\partial \mathbf{x}_{\mathcal{S}'}}$ is the derivative of the prediction function with respect to covariates $\mathbf{x}_{\mathcal{S}}$. For a single covariate (i.e. $\mathcal{S} = \{j\}$), the ALE is approximated as follows. Let the range of a covariate \mathbf{x}_j be partitioned into K intervals beginning with starting point $z_{0,j}$. $N_j(k)$ is the k -th interval with upper boundary $z_{k,j}$ and lower boundary $z_{k-1,j}$, i.e. $]z_{k-1,j}, z_{k,j}]$, and $n_j(k)$ is the total number of observations of \mathbf{x}_j within the interval. Scalar $x_{i,j}$ is the i -th observation of the p -vector \mathbf{x}_j and $\mathbf{x}_{i,-j}$ the values of the other covariates for this observation. Eq. 4 can be approximated by a step function over the K intervals:

$$\hat{f}_{ALE}(x_j) = \sum_{k=1}^{k_j(x_j)} \frac{1}{n_j(k)} \sum_{i: x_{i,j} \in N_j(k)} [\hat{f}(z_{k,j}, \mathbf{x}_{i,-j}) - \hat{f}(z_{k-1,j}, \mathbf{x}_{i,-j})], \quad (5)$$

where $k_j(x_j)$ is the interval that x_j falls into. The right-hand side of Eq. 5 is the difference in prediction computed over the range $]z_{k-1,j}, z_{k,j}]$, which quantifies the *effect* of the covariate for an individual observation within the interval. The sum of the individual effects within the interval is divided by the number of observation in the interval to obtain the *local* average difference of prediction. The left-hand sum of Eq. 5 defines the *accumulated* local effect over all intervals. The formula in Eq. 5 is a step function which can be smoothed by linear interpolation. The ALE is centered at zero by:

$$\text{centered } \hat{f}_{ALE}(x_j) = \hat{f}_{ALE}(x_j) - \frac{1}{n} \sum_{i=1}^n \hat{f}_{ALE}(x_{i,j}), \quad (6)$$

so that a point on the ALE curve is the difference to the average prediction of the model. For the estimation of two-dimensional ALE, the local effect is accumulated over rectangles instead of intervals. Refer to Apley and Zhu (2020), Eq. 13–16 for the equations describing the two-dimensional ALE and see Chapter 5 in Molnar (2020) for more details on the difference between PDP and ALE. An example of one and two-dimensional ALE plot is shown in Fig. 3.

Note that interpretation of the two-dimensional ALE plot is different from that of a PDP. ALE is formally interpreted as being the centered difference in prediction (i.e. the effect) when the observations within an interval are moved from one border of the interval to another other. Fig. 3a shows the effect of woody biomass on SOC for a range of values of woody biomass, and compared to the average prediction. Fig. 3b shows the pure interaction effect of woody biomass and elevation on SOC compared to the average prediction. For example, the ALE estimate of woody biomass in Fig. 3a illustrates that for large values of woody biomass (i.e. greater than 300 Mg/ha), the predicted values of SOC are lower by nearly 20 dg/kg compared to the average prediction.

The estimates of ALE tend to be more robust than the PDP for correlated covariates, because of averaging and accumulating the local effect over the conditional distribution. However, this comes at the expense of having a more localized interpretation (within intervals), and possibly non-intuitive interpretations for some data-generating processes (Grömping, 2020).

2.5. Interaction between covariates

Interaction between covariates can be estimated with the H-statistic (Friedman and Popescu, 2008). Interaction is the variation that remains unexplained after summing the individual effects of the covariates on the model prediction. In other words, there is interaction when the

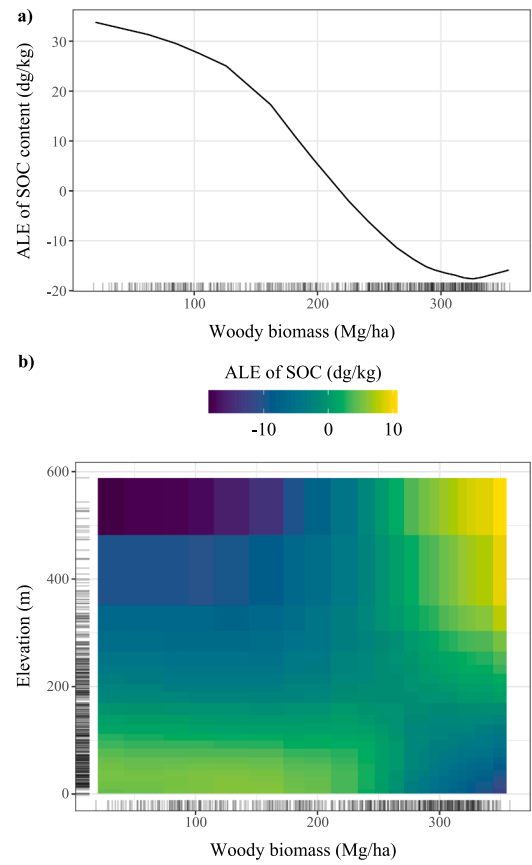


Fig. 3. Estimates of a one-dimensional accumulated local effect (ALE) plot of woody biomass on SOC content (a), and two-dimensional ALE of woody biomass and elevation on SOC (b). Data and model from Wadoux et al. (2022), Section 4.2. Note that the ALE are centered at zero, so that each point of the ALE curve is the difference to the average prediction.

combination of two covariates explains more of the data variance than the sum of these same two covariates taken separately. The H-statistics identifies the strength of the interaction, either between two covariates (*two-way* interaction) or between a covariate and all other combinations of covariates (*total* interaction). The individual covariate effect is measured by the PDP (Section 2). In a two-way interaction, the H-statistic measures the difference caused by the sum of the two individual covariates PDP, compared to the PDP of the combined two covariates. To measure the total interaction, the PDP of a single covariate is compared to that of the entire set of covariates. In each of the cases, the H-statistic is the amount of variance explained by the difference, and is an indication of the strength of the interaction. The interaction between two covariates ($\mathbf{x}_1, \mathbf{x}_2$), i.e. two-way interaction, is measured by the H-statistics as:

$$H_{12}^2 = \frac{\sum_{i=1}^n [\hat{f}_{PDP}(x_{i,1}, x_{i,2}) - \hat{f}_{PDP}(x_{i,1}) - \hat{f}_{PDP}(x_{i,2})]^2}{\sum_{i=1}^n \hat{f}_{PDP}^2(x_{i,1}, x_{i,2})}. \quad (7)$$

The interaction between a single covariate \mathbf{x}_j with all combinations of covariates is:

$$H_j^2 = \frac{\sum_{i=1}^n [\hat{f}(\mathbf{x}_i) - \hat{f}_{PDP}(x_{i,j}) - \hat{f}_{PDP}(\mathbf{x}_{i,-j})]^2}{\sum_{i=1}^n \hat{f}^2(\mathbf{x}_i)}. \quad (8)$$

The H-statistics is dimensionless and usually between 0 and 1, but can exceed one if the variance of the two-way interaction exceeds the

variance of the 2D-PDP (e.g. due to uncertainty in the estimation). A value close to 0 indicates no interaction, whereas a large value means that interaction between the covariates explains most of prediction variance. Fig. 4 shows an example visualization for the total interaction between a set of covariates.

The H-statistic has valid theoretical underpinning through the decomposition of the PDP, and can detect interaction between an arbitrary number of covariates. Further, it is dimensionless, which makes comparison possible between group of covariates and models. However, as for the PDP the H-statistic is sensitive to deviation from the assumption of independence between covariates, and is computationally expensive to estimate when the number of covariates is large.

2.6. Surrogate modelling

A surrogate model is a simple and interpretable model that is calibrated to approximate the prediction of a black-box model. In surrogate modelling, the prediction model \hat{f} which yields prediction of Y with X is approximated by calibrating a simple model g on the n prediction. Model g is interpretable, usually a linear model or a regression tree. The quality of the surrogate model g is evaluated by calculating validation statistics that compare the prediction made by the model \hat{f} to that made by model g , for example the modelling efficiency coefficient (Janssen and Heuberger, 1995):

$$MEC = 1 - \frac{\sum_{i=1}^n (\hat{y}(s_i) - \hat{y}^*(s_i))^2}{\sum_{i=1}^n (\hat{y}(s_i) - \bar{\hat{y}})^2}, \quad (9)$$

where \hat{y} denote the predicted soil property at location s_i by model \hat{f} , and

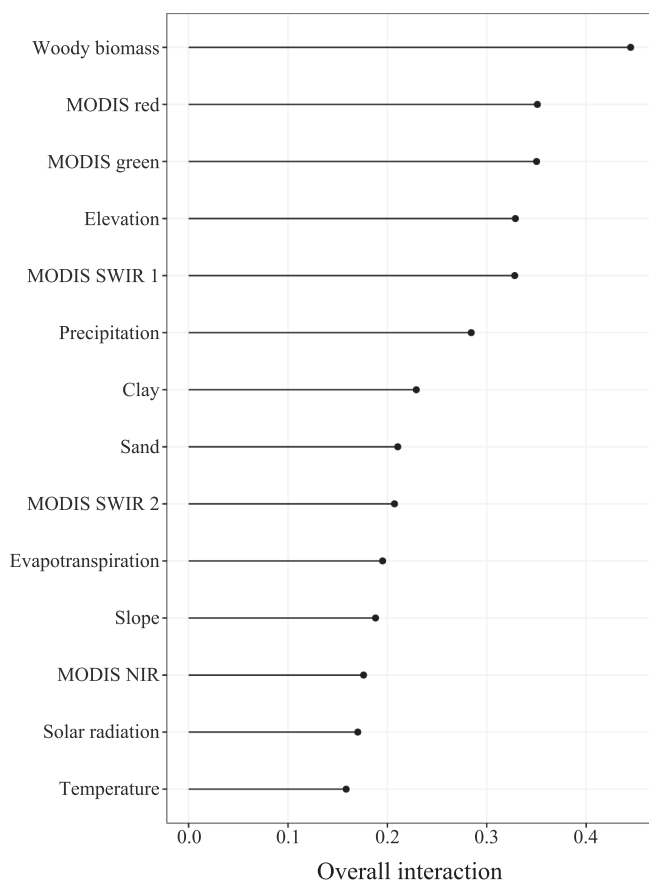


Fig. 4. Estimate of the total interaction (Eq. 8) between 14 covariates used for prediction of SOC. Data and model from Wadoux et al. (2022), Section 4.2.

\hat{y}^* is the predicted value of \hat{y} by model g at the same location. A MEC value of 1 indicates that the surrogate model is a perfect predictor of the values predicted by the black box model, whereas a value of 0 indicate that the surrogate model is as good predictor as the mean or the original predicted values. Note that the MEC can be negative, in which case the model is not a useful predictor of the soil property.

The main advantage of surrogate modelling lies in the intuitive interpretation of the model for non-specialists. There is also flexibility in the choice of surrogate model, usually a linear model or simple decision tree. Note that the surrogate model is an approximation of the predicted values, and thus interpretation should be made cautiously if the variance explained by the surrogate model (as indicated by the MEC) is insufficiently high. To date there is not clear cut-off value of the MEC for which we can be confident that the surrogate model is sufficiently close to the model it approximates.

2.7. Shapley values

Shapley values (Shapley, 1953) originate from coalitional game theory. In a game where a prediction is the “payout”, Shapley values aim to fairly distribute the payout among the covariates. Compared to the other methods, Shapley value is a local method, designed to explain individual predictions. However, Shapley values can be combined to create global interpretations. Recall that a covariate subset is \mathcal{S} , and is composed of $l < p$ covariates. $\mathcal{S} \subseteq \{1, \dots, p\} \setminus \{j\}$ refers to any subset of covariates which excludes covariate j . The Shapley value $\phi_{0,j}$ for covariate j for a data point \mathbf{x}_0 (not necessarily from the original data set) is given by:

$$\phi_{0,j} = \sum_{\mathcal{S} \subseteq \{1, \dots, p\} \setminus \{j\}} \frac{|\mathcal{S}|!(p - |\mathcal{S}| - 1)!}{p!} (\hat{f}^*(\mathbf{x}_0, \mathcal{S} \cup \{j\}) - \hat{f}^*(\mathbf{x}_0, \mathcal{S})), \quad (10)$$

where $|\mathcal{S}|$ is the size of the subset which excludes the j th covariate, $\mathcal{S} \cup \{j\}$ is the subset \mathcal{S} with the j th covariate added, and $\hat{f}^*(\mathbf{x}_0, \mathcal{S}) = E_{\mathbf{x}_{\mathcal{S}^c}}[\hat{f}(\mathbf{x}_0, \mathbf{x}_{\mathcal{S}^c})]$ is the prediction function where covariates not contained in \mathcal{S} are marginalized (similar for $\mathcal{S} \cup \{j\}$). Recall that p is the number of covariates. Then $\hat{f}^*(\mathbf{x}_0, \mathcal{S} \cup \{j\}) - \hat{f}^*(\mathbf{x}_0, \mathcal{S})$ can be interpreted as marginal contribution to the prediction when adding covariate j to the subset of covariates \mathcal{S} . The right hand-side of Eq. 10 is the marginal contribution for a subset of covariates, whereas the left hand-side is a weighted average, giving equal weight to each of marginal contributions of all possible subsets of covariates. The contribution of a covariate to the prediction of a single spatial location is then given by $\phi_{0,j}$.

The exact solution to Eq. 10 requires estimating the sum of the marginal contribution over $2^p - 1$ combinations of covariates, which is computationally inefficient if the number of covariates is large. Štrumbelj and Kononenko (2014) and Lundberg and Lee (2017) proposed estimation methods to reduce the computational cost. Štrumbelj and Kononenko (2014) introduced an approximation algorithm for Eq. 10 based on Monte-Carlo sampling. They further approximate the covariate effect on the prediction by integrating over the n observations of the calibration dataset. Lundberg and Lee (2017), reduce estimation of Shapley values as the optimal solution of a (local) weighted linear least squares regression (called KernelSHAP). Hereafter, Shapley values are estimated by the algorithm presented in Štrumbelj and Kononenko (2014).

A Shapley value is interpreted as the average contribution of a covariate to the prediction, in the unit of the soil property. Shapley values are commonly used to evaluate the individual contribution of each covariate to the prediction of the soil property at a particular location (i.e. local interpretation), compared to the average prediction of the calibration dataset (Fig. 5c). The absolute value of the Shapley values for individual observations in the calibration dataset can be summed to obtain an overall covariate importance, see also Section 2.1 and Fig. 5a

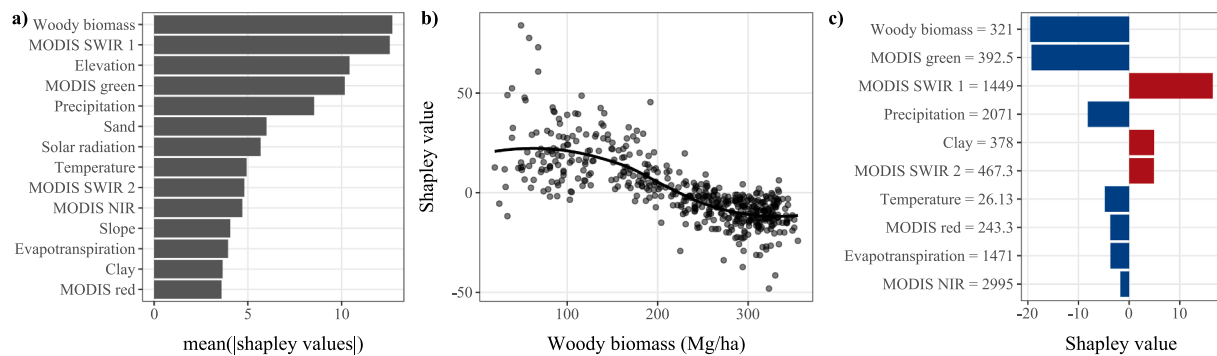


Fig. 5. Average of the absolute Shapley values in the calibration dataset (a), dependence plot of SOC against woody biomass (b), and local interpretation of a single spatial location (c). Data and model from Wadoux et al. (2022), Section 4.2.

for an example). Note, however, that overall covariate importance obtained by permutation is based on decrease in model accuracy whereas covariate importance based on Shapley values shows the overall contribution of the covariates to the prediction of the calibration dataset. Finally, the average of Shapley values in the calibration dataset for a covariate plotted against the covariate values is an indication of the partial dependence (Fig. 5b).

3. Illustration with soil data

We built and interpreted two models for mapping soil organic carbon content in France. We used as calibration sample ($n = 2947$) composed of topsoil (0–20 cm) values of organic carbon content (in g kg^{-1}) from the land use and cover area frame statistical survey (LUCAS, Orgiazzi et al., 2018) dataset. We collected a set of 29 environmental covariates covering France and representing six factors influencing SOC spatial distribution: topography, vegetation (including remote sensing imagery), long-term average climatic conditions, climate seasonality, extreme climatic conditions and soil. The list of covariates, their description and source is provided in the Supplementary Materials. All covariates were resampled using bilinear interpolation or aggregated to conform with a spatial resolution with grid cells of $250 \text{ m} \times 250 \text{ m}$. The SOC data and their matching values of environmental covariates were then used to calibrate two mapping models.

The first model used was random forest (RF, Breiman, 2001) which we calibrated using 250 trees and a *mtry* parameter fixed at the rounded down square root of the number of covariates. The *mtry* parameter is the size of the random partition from the set of covariates during the splitting of a tree. All other parameters were held to their default value. We

used the R programming language (R Core Team, 2020) for the implementation and the R package ranger (Wright and Ziegler, 2017). The second model used was a multiple linear regression (MLR, Hastie et al., 2009) fitted using ordinary least squares and the default implementation from the R package stats. Note that there is no fundamental objection to use interpretation methods on a MLR model, although this model structure is simple and can readily be interpreted. This allows us to compare the linear regression model with the random forest model and reveal the functioning of the interpretation methods. Both RF and MLR models were validated using random 10-fold cross-validation. The model predictions did not have a systematic over- or under-prediction (mean error close to zero) and had a RMSE value of 21.19 and 21.65 g kg^{-1} for random forest and linear regression, respectively. Finally, we used all the SOC data for model calibration and prediction. The resulting SOC maps are shown in Fig. 6.

We apply the local and global interpretation methods described in Section 2. We interpret the RF model and compare it with the MLR model when relevant. The global methods are applied on the models whereas the local methods are applied to a geographical area and to two contrasting spatial locations (Fig. 7). This allows us to understand how the importance of environmental covariates varies in space and from one location to another. The geographic region of study is called *Maine-et-Loire*, located in Western France in the Loire basin, and characterized by large variety of arable soils with overall relatively low carbon content. We interpreted the model for a region by dividing a geographical area into a fine grid and by treating each predicted pixel as an individual “observation”. The two spatial locations are denoted *Beauce* and *Landes* and can be referred to as individual pedons with the same support than the observations from LUCAS. Location *Beauce* is in a cropland-

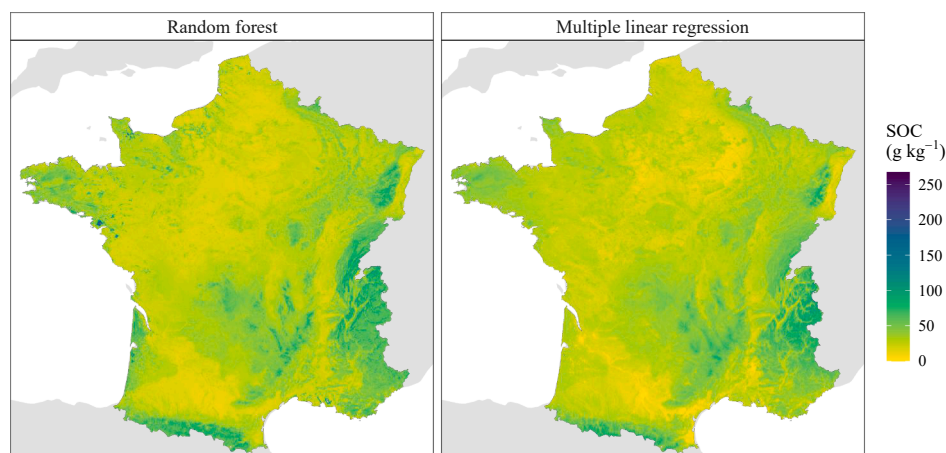


Fig. 6. Spatial distribution of SOC (in g kg^{-1}) for Metropolitan France excluding Corsica. The SOC maps were made using random forest (left) and multiple linear regression (right).

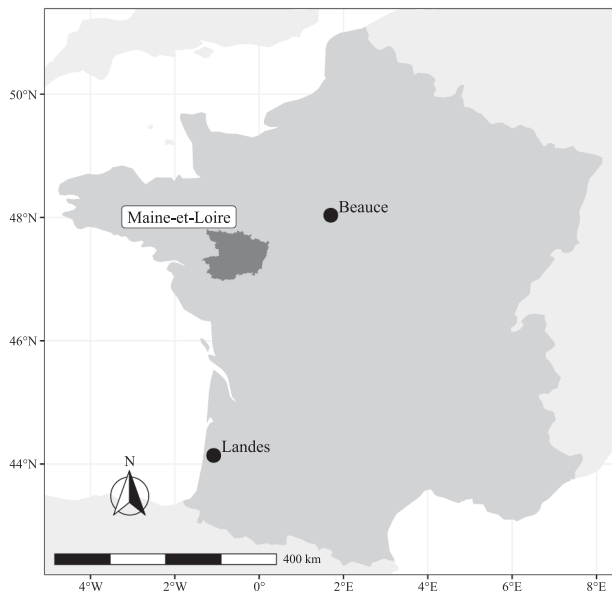


Fig. 7. Location of the two spatial locations and the geographical area for the implementation of the local interpretation methods. The two black dots represent two spatial locations with contrasting SOC content. They are called *Beauce* and *Landes*. The dark grey area is called *Maine-et-Loire* and represents an administrative unit.

dominated region with fertile clay and/or silt-loam soils and relatively low carbon content due to intensive agriculture whereas *Landes* is in a coniferous forested area with sandy soils (i.e. Podzols), but with

relatively high topsoil carbon content due to little interest in these soils for agricultural purposes (Meersmans et al., 2012). Implementation of the interpretation methods was made with the R packages *iml* (Molnar et al., 2018) and *fastshap* (Greenwell, 2020).

3.1. Global interpretation

3.1.1. Which are the drivers of SOC spatial variation?

Fig. 8a shows the covariate importance of the RF model (ratio of RMSE) obtained by 100 permutations. Nearly all covariates are important for the RF model. The figure indicates that three MODIS satellite imagery covariates (i.e. MODIS red, green and SWIR 2) are the most important. Removing them would decrease the RMSE by a factor of 1.33, 1.36 and 1.41 for the MODIS SWIR 2, green and red images, respectively. Elevation and net primary productivity are important covariates too. The covariate representing soil water content for 1500 kPa suction is, conversely, not essential to the RF model, because its importance value is close to a ratio of RMSE value of 1 (i.e. removing covariate soil water content for 1500 kPa does not affect model prediction accuracy). **Fig. 8b-c** shows the covariate importance for groups of covariates, for both RF (**Fig. 8b**) and MLR (**Fig. 8c**). All groups of covariates are important in the RF model. Vegetation, soil and topographic covariates are the most important. An opposite pattern is found in the MLR model, where these group of covariates appear the least important. For the MLR model, the two groups of covariates representing extreme and average climate conditions are the most important.

Fig. 9 shows an alternative interpretation of the RF covariate importance with Shapley values. Note that while **Fig. 8** shows the change in model RMSE, **Fig. 9** shows the magnitude of individual covariate contributions to the prediction of the SOC data used for calibration. **Fig. 9a** indicates that the most important covariates are MODIS images and elevation. The overall ranking of covariate importance

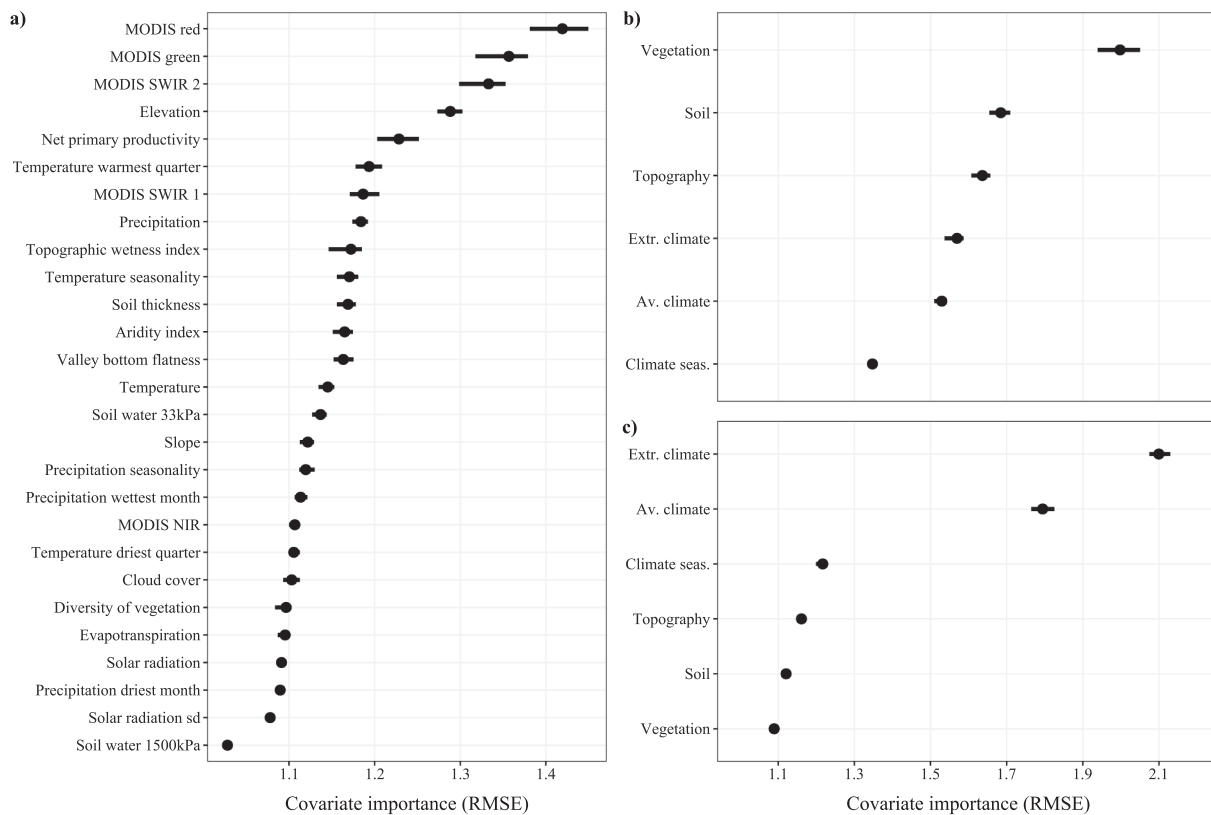


Fig. 8. Mean and 90% confidence interval of the permutation-based covariate importance for a) all covariates of the random forest model, b) group of covariates for the random forest model and c) group of covariates in the multiple linear regression model. Covariate importance is assessed by the ratio of RMSE over 100 permutations. We refer to the Supplementary Material for information on the group of covariates.

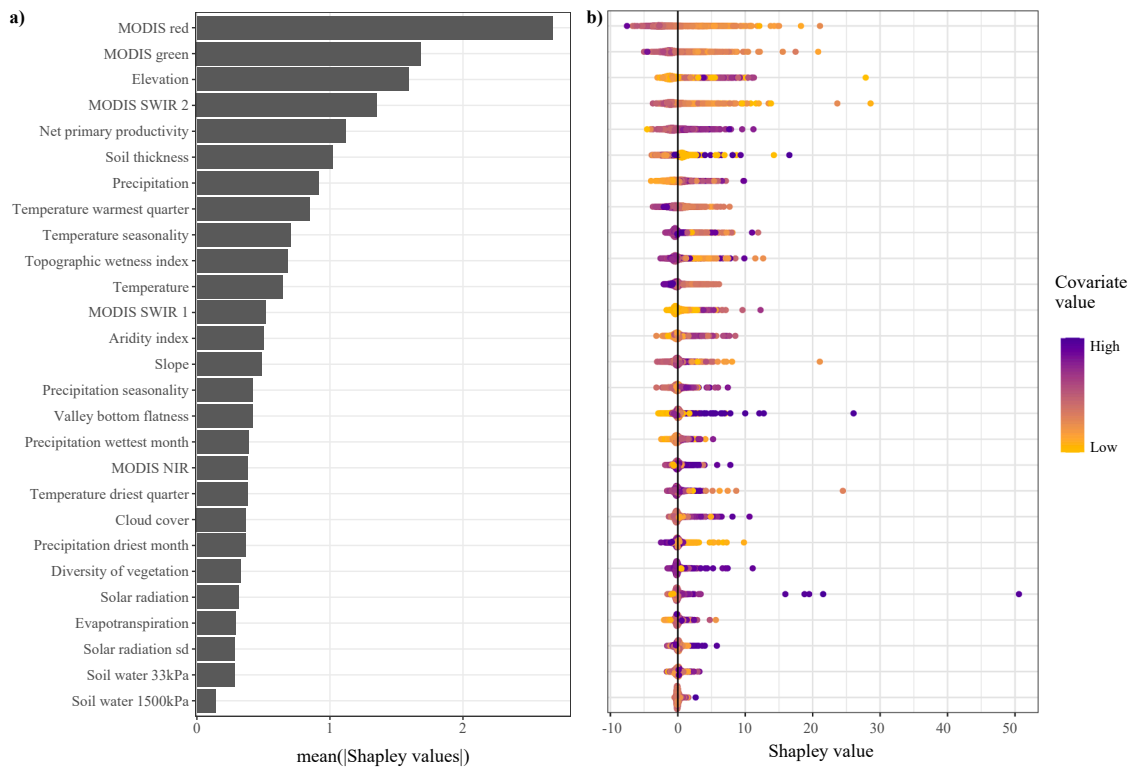


Fig. 9. Covariate importance estimated with Shapley values for the RF model. Plot a) shows the average covariate contribution to the prediction in the calibration dataset, that is, the averaged absolute values of plot b). Plot b) shows the individual Shapley values for each location of the calibration dataset, i.e. the contribution of the covariate to the prediction at this location. The colour in b) represents the covariate value normalized in the range (0–1).

obtained by Shapley values is similar to that found with the permutation-based method. Fig. 9b shows the covariate contribution to each individual location found in the calibration dataset. The most important covariates (e.g. MODIS red) have a large range of Shapley values (i.e. between -10 and 25 g kg^{-1}), meaning that this covariate can have a relatively important contribution to the model prediction. Fig. 9b also provides insight into the relationship between the relative covariate contribution to the prediction and the value of this covariate. For example, valley bottom flatness has, on average, a moderate impact in model prediction (Fig. 9a), but this is more subtle than that (Fig. 9b). For large values of valley bottom flatness, the covariate has a positive relationship with the SOC (i.e. it increases the SOC content), while it is the opposite for small values of valley bottom flatness.

3.1.2. What is the functional form of the association between environmental covariates and SOC?

Fig. 10 shows the effect of elevation on SOC, estimated with three difference methods (i.e. PDP in Section 2.3, ALE in Section 2.4 and Shapley values in Section 2.7). In each of the cases, SOC sharply decreases with elevation and then steadily increases for values of elevation larger than 250 m. With elevation values larger than 900 m, SOC levels off in the PDP, continues to increase in the ALE plot and decreases in the plot with shapley values. Note the different interpretation between the plots of Fig. 10. Fig. 10a (PDP) shows how the predicted SOC values change with elevation whereas Fig. 10b (ALE) shows the effect of elevation on SOC compared to the average prediction of SOC (i.e. centered at zero). Finally, Fig. 10c shows the relative contribution of elevation to the individual SOC observations of the calibration dataset

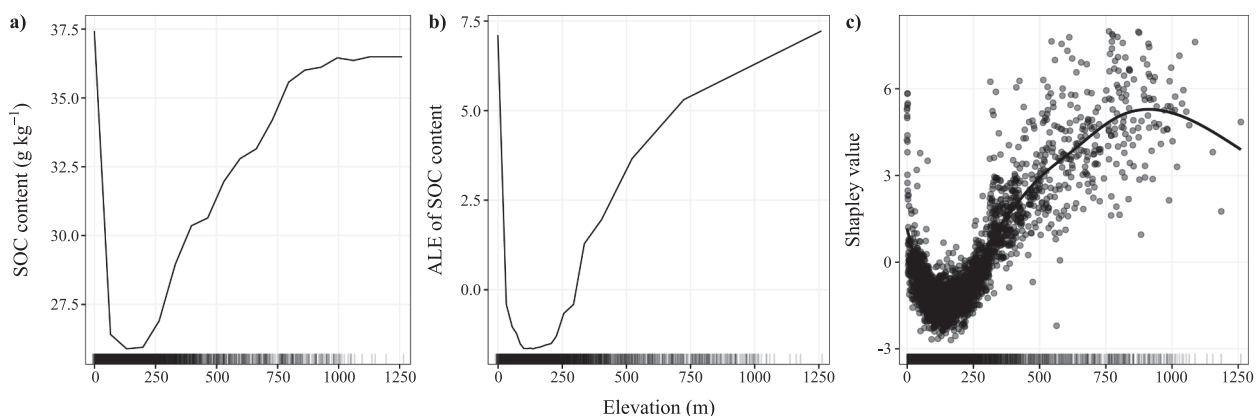


Fig. 10. Effect of elevation of SOC estimated with a) partial dependence, b) accumulated local effect and c) Shapley values. The x-axis shows the marginal distribution of elevation in the calibration dataset. In c) the black dots represent the individual Shapley values and the black curve is a smoothed line obtained over the Shapley values with a conditional mean function. Note that these results were obtained with the RF model.

(the black dots).

The two-dimensional relationship of SOC with elevation and precipitation seems more complex (Fig. 11) than the one-dimensional figures in Fig. 10. Fig. 11a shows that SOC content generally increases with higher elevation and more precipitation. However, the ALE plot in Fig. 11b has a different pattern: for elevation lower than 250 m, the SOC content increases with precipitation, while an opposite pattern is seen for elevation values larger than 250. In Fig. 11, both plots have a noticeable increasing pattern of SOC with higher precipitation, but only for low elevation. Above an elevation of 1000 m, few SOC observations exist, which means that interpretations of effects for this elevation should be cautious.

3.1.3. How does SOC prediction depend on interactions among covariates?

Fig. 12 shows the strength of the interaction between environmental covariates for the RF model. Note that the MLR is not expected to contain an interaction effect between covariates unless explicitly specified. Fig. 12a shows the presence of a strong overall interaction effect in the random forest model. Satellite imageries MODIS red, green and SWIR 2 are involved in interactions with other covariates. Elevation also substantially interacts with other covariates. Covariates standard deviation of monthly solar radiation and soil water content, conversely, have negligible interaction. Fig. 12b identifies how strong covariates interact with elevation. Elevation is dominantly interacting with MODIS SWIR 2, precipitation seasonality and topographic covariates (e.g. wetness index). There is no strong interaction of elevation with soil water content, solar radiation and diversity of vegetation.

3.1.4. How to summarize the model?

Fig. 13 shows a surrogate model of the RF model. The surrogate model is a simple decision tree with a depth of three. It has a MEC of 0.3. The final nodes show the average predicted value and the percentage of data in the node. The colour of the final node is proportional to the value

in the node. The colours associated to the rules are reported in the map of France. Fig. 13 shows that MODIS red band, elevation and climate seasonality covariates were selected by the surrogate model. Accordingly, the smallest predicted values of SOC (i.e. $\text{SOC} \leq 24 \text{ g kg}^{-1}$) are found for locations with large values of the MODIS red band and low elevation ($< 312 \text{ m}$). Large values of predicted SOC, conversely, are found for locations with relatively low values of MODIS red, when temperature of the warmest quarter are moderate (< 18 degrees) and precipitation of the driest month are relatively abundant ($> 67 \text{ mm}$). The pattern of the decision rules shown in the right-hand side of Fig. 13 shows regions where the RF model is likely to predict similar values of SOC. The map pattern shows that large SOC content is predicted in mountainous regions, and in a relatively large amount in Brittany and Normandy. Cropland and vineyard have low predicted carbon, whereas forested areas such as in the Landes have a high carbon content.

3.2. Local interpretation

3.2.1. What is the local functional form of the association between environmental covariates and SOC?

Fig. 14 shows the local association between SOC and elevation in the *Maine-et-Loire* area. The association is estimated for the RF model with ICE curves and their average value (i.e. their PDP), centered at the average value of elevation in the area (68 m). Each ICE curve is a location in the area. In *Maine-et-Loire*, SOC decreases with higher elevation, but this effect is relatively minor, as shown by the PDP curve that is nearly always close to zero. The ICE curves show a different association for individual locations. While most of the ICE curves are close to the PDP, for some locations the SOC content is relatively high (i.e. $> 8 \text{ g kg}^{-1}$ at 0 m) for low elevation and sharply decreases with higher elevation. Overall, there is more variability in the individual ICE curves for low elevation than for high elevation, which suggests that SOC content is higher and more variable with low elevation than it is with high elevation in *Maine-et-Loire*. The pattern of ICE curves observed in this area is thus similar from that observed on average for France for the elevation range 0–200 m, where elevation has a positive relationship with SOC content (see also Fig. 10a-b).

Fig. 15 shows the ICE curves of SOC with elevation and MODIS SWIR 2 band, for the MLR and RF models, and the two locations of interest, Beauce and Landes (Fig. 7). Fig. 15 shows that the two models predicted different values of SOC for Landes, but predicted similar values for Beauce. The predicted SOC of Beauce is also lower than that of Landes. The association between the SOC content and the two covariates (i.e. elevation and SWIR 2) is different between models. The linear model has ICE curves that increase and decrease linearly with elevation and MODIS SWIR 2, respectively. For random forest, the ICE curves have more variation and are not linear: in both locations SOC content slightly increases with elevation up to about 1000 m, after which SOC content levels off. At location *Landes*, a sharp decrease of SOC content is observed for increasing elevation in the first 20 m. Covariate MODIS SWIR 2 has negative relationship with SOC for the location in *Landes* up to values of about 1100, after which the SOC values are stable around 25 g kg^{-1} . For the location in *Beauce* SOC slightly decreases between 1000 and 1500, then remains constant.

3.2.2. How do environmental covariates contribute to the local prediction?

The spatial pattern of the Shapley values for the multiple linear regression and random forest models and five covariates is shown in Fig. 16. The figure shows clear differences in the contribution of covariates to the predictions and clear spatial patterns. The MODIS red band has large positive or negative Shapley values. This is also the case for elevation and precipitation. All covariates have a detailed spatial pattern of change in Shapley values with increasing distance from the Loire river. Substantial differences are also observed between the multiple linear regression and random forest models. The contribution of the MODIS red band to the SOC prediction made by the random forest model

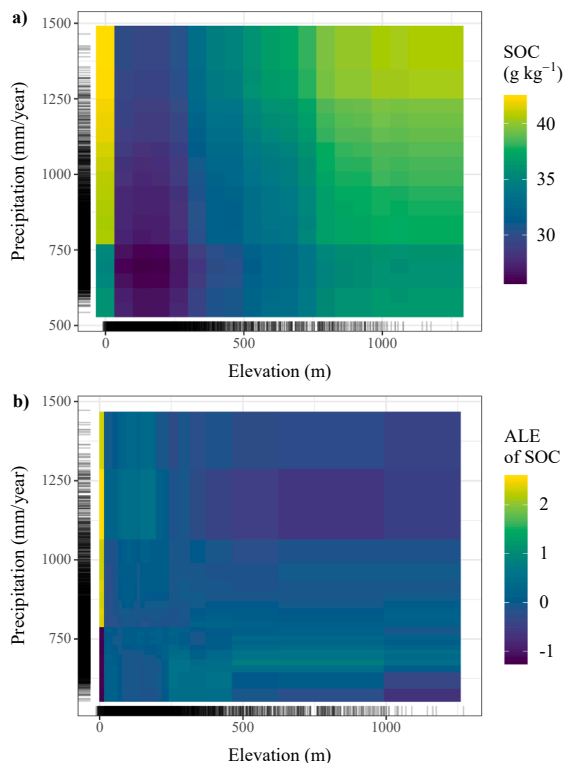


Fig. 11. Two-dimensional partial dependence plot of the effect of elevation and precipitation on SOC content (a), and accumulated local effect of elevation and precipitation on SOC content (b). Note that these results were obtained with the RF model.

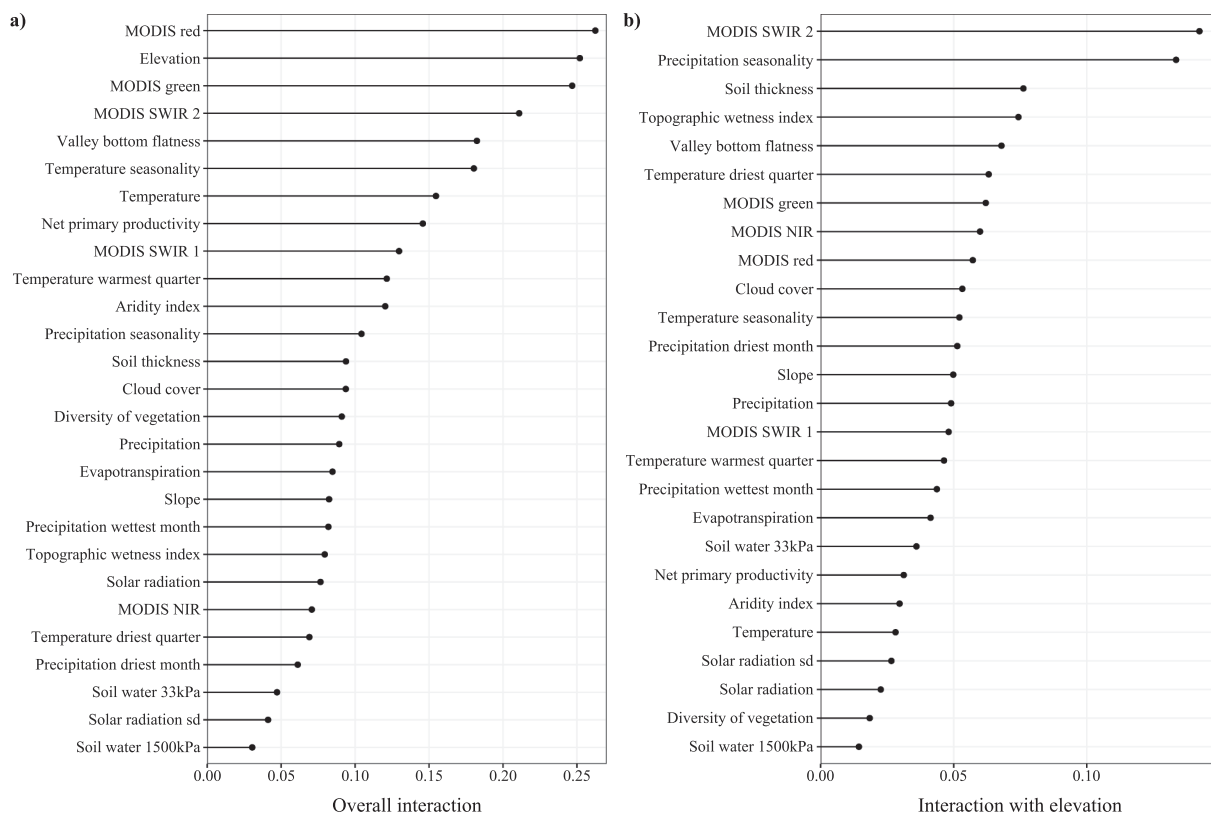


Fig. 12. Estimate of the overall interaction (i.e. the H-values calculated with Eq. 8) between the environmental covariates used in the random forest model (a) and estimate of the two-way interaction (Eq. 7) with elevation (b) for the RF model.

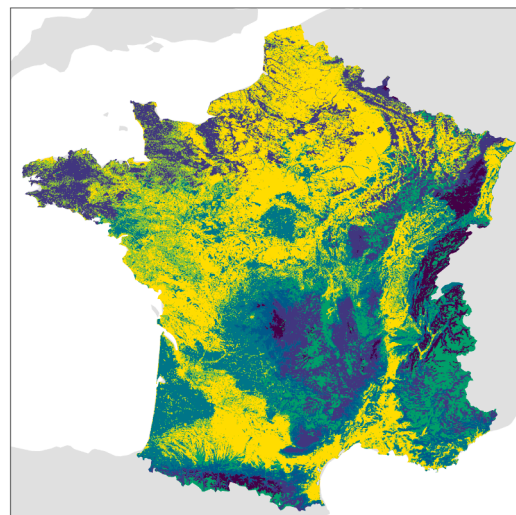
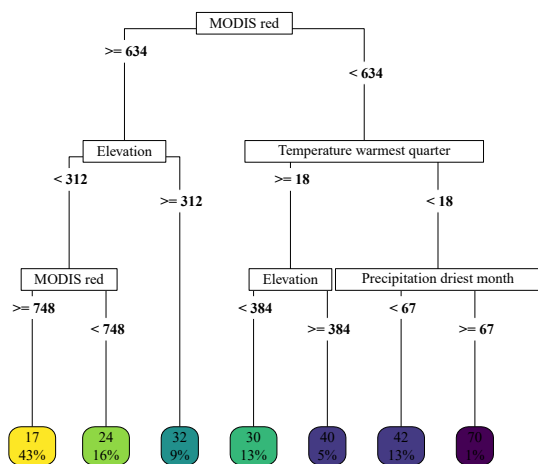


Fig. 13. Surrogate model of the random forest model for prediction of SOC. The surrogate model (left) is a decision tree. The final node shows the average predicted value and the percentage of data in the node. The colour of the final node is proportional to the value in the node. The colour scheme is reported in France (right) using the rules of the decision tree. Values of the final node range from 17 (43%) in yellow to 70 (1%) in dark purple.

is very different from that made by the multiple linear regression model. Also the pattern of Shapley values for precipitation and elevation is different between models. The linear regression model has a gradient of increasing Shapley values from North to South for the covariate precipitation. In the large floodplain of the river, elevation, topographic wetness index and slope have a negative contribution to the SOC prediction while it is the opposite for the linear model.

Fig. 17 shows the covariates contribution to the SOC prediction made

by RF at two spatial locations, in *Beauce* and *Landes*. The Shapley values of Fig. 17 show the positive or negative contribution to the prediction, in the unit of the SOC, using the average prediction from the calibration dataset as baseline. Slight differences between the sum of Shapley values and the predictions are due to the approximation strategy. Fig. 17 shows that SOC prediction in the two spatial locations are made in a very different way. The location in *Beauce* has low SOC content, and so contribution of covariates is mostly negative. MODIS red, green, SWIR 2,

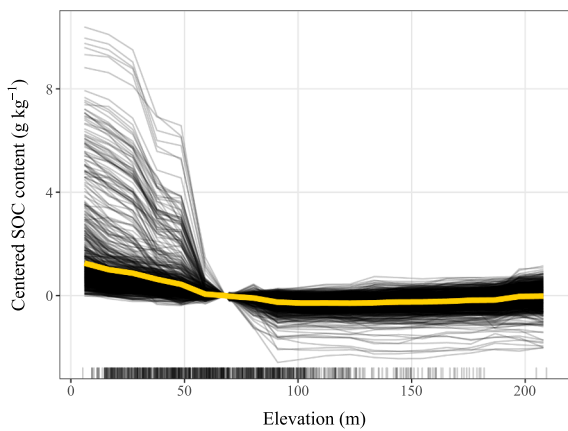


Fig. 14. Centered effect of elevation on the SOC in the region *Maine-et-Loire* and for the RF model. The effect is centered at the average elevation value of the area ($x_0 = 68$). The black curves are the individual conditional expectation whereas the yellow curve is their average (i.e. their partial dependence function).

net primary productivity and elevation had a large negative contribution, whereas a small positive contribution to the SOC prediction is made by the soil thickness. In the location in *Landes*, the SOC content is also lower than the average. Large positive contributions to the SOC predictions are made by the MODIS green and red bands, and by the net primary productivity. The temperature of the warmest quarter and standard deviation of the solar radiation show negative contributions to the SOC prediction.

4. Discussion

The methods tested for the interpretation of two mapping models provided valuable information on the drivers of SOC variation in France, their interaction, as well as on the functional form of the association between environmental covariates and SOC. This information was obtained either for a single spatial location or globally from the model as a whole. In our case study, for example, MODIS remote sensing images were on average the most important variables contributing to SOC prediction. The overall importance of MODIS images to predict SOC does not come as a surprise, because spectral characteristics of MODIS images correlate to biogeochemical properties relevant to explain the spatial distribution of SOC. MODIS red band strongly correlates with soil

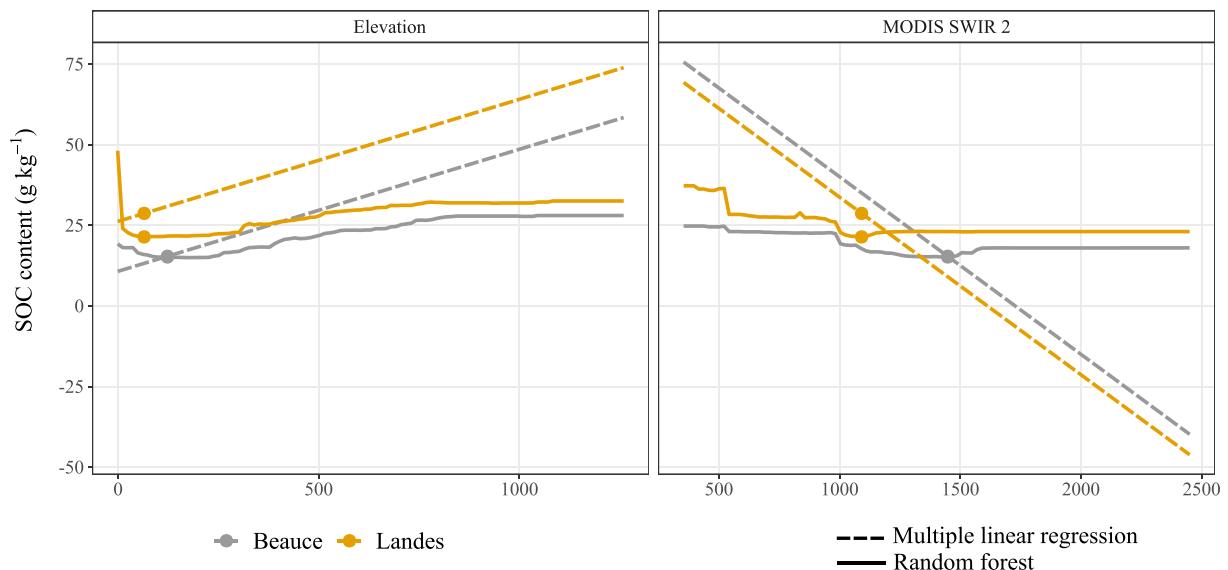


Fig. 15. Individual expectation curves of the effect of elevation and MODIS SWIR 2 on SOC for the two locations of interest *Beauce* and *Landes* and the multiple linear regression and random forest models. The dots represent the SOC prediction made by the model at the locations.

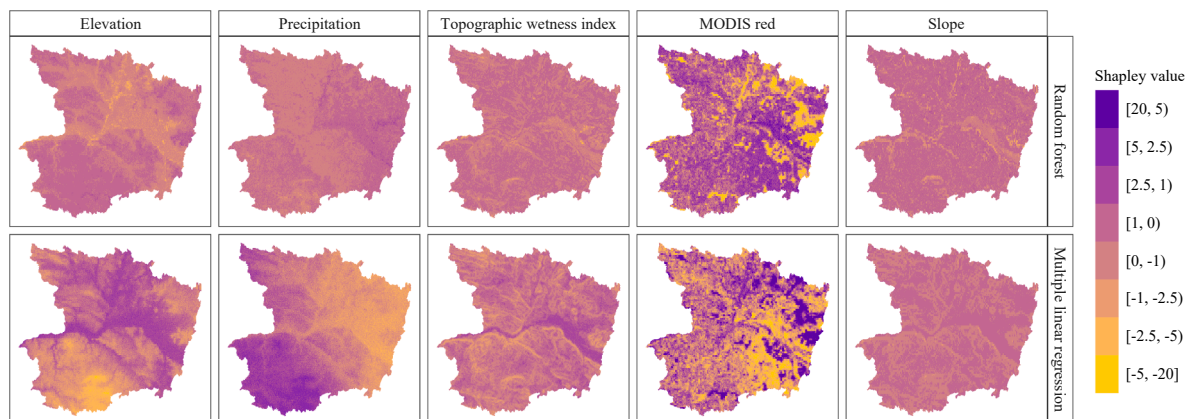


Fig. 16. Spatial pattern of the Shapley values for five covariates and the two mapping models. Dark colour indicates that the covariate has a positive contribution to the SOC prediction while light colour indicates a negative contribution.

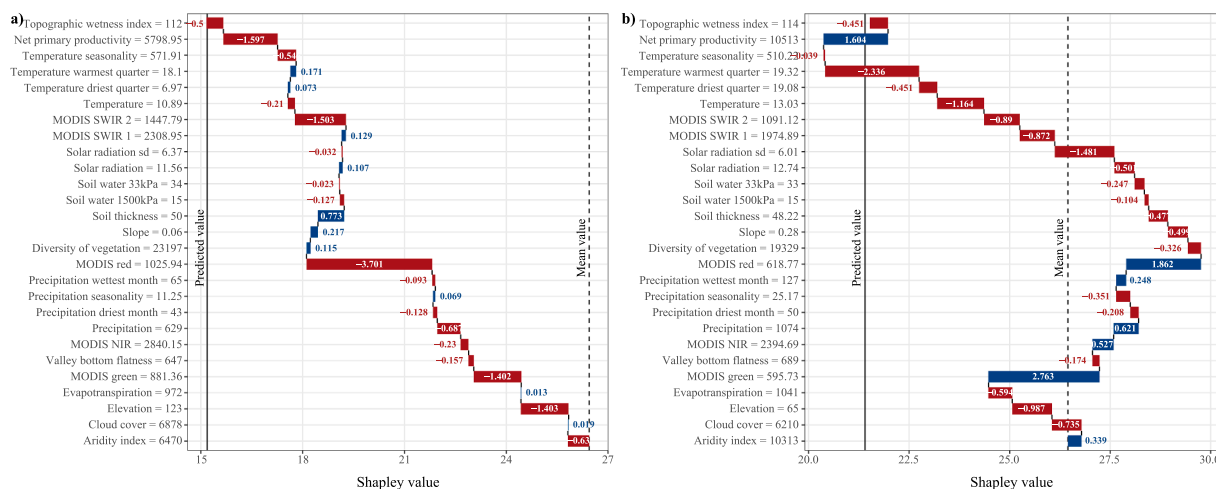


Fig. 17. Contribution of the individual covariates to the prediction made by the RF model of SOC at locations Beauce (a) and Landes (b). Contributions are estimated with Shapley values. The red colour indicates a negative contribution while a blue colour indicates a positive contribution. The y-axis indicates the value of the covariate at the prediction location.

organic matter (Dou et al., 2019). Vågen et al. (2016) used MODIS reflectance data only to predict SOC, pH, sand content, sum of exchangeable bases, as well as root-depth restrictions with high accuracy in Africa. Further, our results suggested that locally, elevation, precipitation, or valley bottom flatness could outweigh MODIS images. Admittedly, the functional form of the association between environmental covariates and SOC varies from one location to another. In a clayey agricultural soil, SOC was not only lower in content than in a sandy soil covered by a coniferous forest, but the environmental covariates contributed differently to the predictions. While these results should be interpreted with care, the average predicted value of SOC and the main covariates contributing to the prediction for these two locations appear realistic compared to existing studies (see for instance Meersmans et al., 2012).

The utility of the methods used in this paper, along with their are pros and cons are summarized in Table 1. We stress that in spite of apparent similarities between the methods (as illustrated, for example, in Fig. 10), the results actually differ in which aspects of the relationship between SOC and covariates they describe. Also, the representation of

the covariate importance obtained by permutation (Fig. 1) and Shapley values (Fig. 9) is seemingly similar, but covariate ranking in the two methods is made differently. Because of these similarities ample attention should be paid to the conclusions that can effectively be drawn with the interpretation methods. There is a risk that practitioners misinterpret the output of these methods. Apart from an understanding of which conclusions can potentially be drawn, a number of assumptions underlie the methods, the most important of which is that of independence between covariates. Permutation-based methods (e.g. covariate importance with permutation, PDP, Shapley values) might lead unrealistic results when covariates are dependent, because perturbation can produce data points that lead outside the multivariate covariate space. An illustration of this problem along with a simulated example is provided in Molnar et al. (2022), Section 5. It does not mean that permutation-based methods cannot be used when covariates are dependent, as is almost always the case in DSM studies, but that one must take care when interpreting the output of these methods.

Alternatively, methods that better account for dependence between covariates exist, such as when using the ALE instead of the PDP

Table 1
Summary table of the model-independent methods for global and local interpretation of mapping models.

Method	Level	Utility	Pros	Cons	Reference
Covariate importance with permutation	Global	Quantifies the importance of a covariate or group of covariates on model accuracy.	Intuitive interpretation. Takes into account interaction among covariates. Fast to compute.	Misleading when covariates are dependent.	Fisher et al. (2019)
Partial dependence plot	Global	Shows the association between covariates and soil property	Intuitive interpretation. Fast to estimate for small n .	One or two covariates can realistically be displayed in a single plot. Misleading when covariates are dependent.	Friedman (2001)
Accumulated local effect	Global	Shows the association between covariate and soil property.	Suited for dependent covariates. Fast to compute.	One or two covariates can realistically be displayed in a single plot. Cannot be estimated for a single location. Not available for categorical covariates.	Apley and Zhu (2020)
H-statistic	Global	Identifies the strength of the interaction between covariates.	Dimensionless. Has an underlying theory.	Slow to compute. Misleading when covariates are dependent.	Friedman and Popescu (2008)
Surrogate modelling	Global	Gives a summary of the model.	Intuitive interpretation. Flexibility in the choice of surrogate model.	Comes with the disadvantages of the surrogate model. Often difficult to approximate the black box model.	Molnar (2020)
Individual conditional expectation	Local	Shows the association between covariate and soil property at a single location.	Intuitive interpretation. Fast to estimate.	A single covariate can realistically be displayed in a plot. Misleading when covariates are dependent.	Goldstein et al. (2015)
Shapley values	Local/global	Quantifies the relative contribution of a covariate to a prediction	Has an underlying theory. Intuitive interpretation. Additive, and can be used for global interpretation.	Slow to compute. Misleading when covariates are dependent.	Shapley (1953), Štrumbelj and Kononenko (2014), Lundberg and Lee (2017)

(Table 1), or by using variants that rely on the conditional distribution (e.g. conditional feature importance). Note, however, that in each of the cases using a different method or a method that relies on the conditional distribution, might give results that are non-intuitive and more difficult to interpret.

As mentioned in the Introduction the aim of this paper was to show how insights can be obtained from complex empirical soil models, but interpretation of such models to explain the origin or causal mechanisms of the spatial distribution of soil properties should be made with care. Soil scientists are usually interested in obtaining insights into the data generation process by interpretation of the empirical relationships found by the model. While that is a worthy objective, empirical models do not aim to provide a diagnosis of causalities in the spatial pattern of soil properties, nor do they account for mechanisms derived from our knowledge of major soil processes. In our study, the strong dependence on MODIS satellite (spectral) imagery to produce the maps takes out of the realm direct assessment of causalities between soil forming factors and SOC, because satellite data are not intended to represent any pedological mechanism involved in the spatial distribution of SOC, although energy from MODIS images interact with photosynthetic vegetation and provide a proxy for vegetation and may also capture differences in geology/parent material in drylands. Any interpretation on mechanisms involved in SOC distribution should however made with care. Several recent studies have argued in this sense (e.g. Fourcade et al., 2018). Wadoux et al. (2020), for example, demonstrated that a complex empirical model is able predict accurately SOC, even when the covariates used to fit the model were meaningless and unrelated to known soil forming factors. They concluded that the pattern found by these complex models are not a reliable way to obtain new pedological knowledge. We recommend to use the interpretation methods described in this paper to obtain insights into the pattern found by the model, and then to translate the pattern into the formulation of hypotheses through connection of patterns to possible soil processes.

Another option, especially applicable when producing quantitative soil information (i.e. prediction) is the main objective, is to use interpretation methods to perform a diagnostic on the model. In many soil mapping studies issues of hypothesis generation are not present, so an assessment of potential causalities is not a priority. Often however, the modelling process is made of refining, possibly including manual selection of covariates and visual examination of some portions of the map. The overall model validation statistics might be acceptable, but the predicted pattern in some areas might not conform with expectations. Take, for instance, a model that predicts abnormally high SOC content in a sandy soil. Should we collect more data in this area or incorporate more relevant covariates? Model diagnostics further motivates the application of the methods described in this paper.

This study explored a complementary set of methods for the local and global interpretation of complex soil models. Within the framework of model-independent techniques we might also explore recent developments such as breakdown plots (Robnik-Šikonja and Kononenko, 2008) for additive and non-additive attribution, functional decomposition (Molnar et al., 2020), or local interpretable model-agnostic explanations LIME, Ribeiro et al., 2016). LIME is a popular local interpretation method potentially suited when the number of covariates (explanatory variable) is very large. However, this method also has disadvantages such as instability in the results and sensitivity to the local neighborhood size. Also here Shapley values might provide a computationally tractable alternative method for the interpretation of complex soil models. Thus, we did not present LIME in this study but we acknowledge that this might be a valuable approach too.

The alternative to model-independent methods is the use of prediction models that are not “black boxes” or interpretation methods that are specific to a model. In many instances sufficient insights into soil processes can be obtained through the rule sets generated by methods that rely on a statistical model. Geostatistical models of soil variation, for example, through the analysis of the variogram and kriging, can be

interpreted in terms of the estimated variogram parameters and plausibility of the assumptions, which all give insights into the nature of soil variation. Notably, geostatistical models are powerful for prediction and provisions to address complex non-stationary soil variation exist (e.g. through wavelet transform).

Finally, in the Introduction we presented a set of interpretation methods that are specific to a model. These methods are valid and useful for the interpretation of complex models. We refer to Biecek and Burzykowski (2021), Section 1.5 for an overview and to Molnar (2020), Section 10 for a summary of model-specific methods for interpreting artificial neural networks. Further investigations are needed to understand how these methods can be used for the interpretation of digital soil models.

5. Conclusion

We have presented methods to obtain insights into complex models of soil variation. These methods were reviewed and evaluated in a case study for mapping topsoil organic carbon in France using a large set of environmental covariates as predictors and two models. From the results and discussion we draw the following conclusions:

- The methods presented in this paper allows one to extract and visualize different aspects of a complex model.
- In a case study, we reveal i) the importance of each driver of soil variation, ii) their interaction and iii) the functional form of the association between environmental covariates and the soil property.
- Interpretation could also be performed locally, for an area or a spatial location of interest.
- The use of Shapley values for interpreting complex models of soil variation is a promising future line of research because it is versatile, enables both local and global interpretation, is easy to interpret and has an underlying theory.
- Different methods might produce seemingly similar results. Ample attention should be paid to the conclusions that can effectively be drawn with the interpretation methods.
- A number of assumptions underlie the use of the interpretation methods, the most common of which is that of independence between covariates. Deviation from this assumption does not preclude the use of the methods, but results should be interpreted with care.
- We presented a summary table as a guide for selecting the interpretation method, given the purpose of the study and the pros and cons of the method.

We stress the importance of going beyond prediction in the use of complex statistical or non-statistical models. Interpretation of models reveal how the predictions are made and can help formulate hypotheses about the underlying soil processes and mechanisms driving soil variation. Interpretation methods are also valuable when the production of quantitative soil information (i.e. prediction) is the main interest, to assist model refining and the evaluation of model prediction plausibility.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgment

The LUCAS topsoil dataset used in this work was made available by the European Commission through the European Soil Data Centre managed by the Joint Research Centre (JRC), <http://esdac.jrc.ec.europa.eu/>.

Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <https://doi.org/10.1016/j.geoderma.2022.115953>.

References

- Apley, D.W., Zhu, J., 2020. Visualizing the effects of predictor variables in black box supervised learning models. *J. R. Stat. Soc.: Ser. B (Statistical Methodology)* 82, 1059–1086.
- Biecek, P., Burzykowski, T., 2021. *Explanatory Model Analysis: Explore, Explain, and Examine Predictive Models*. CRC Press, Boca Raton.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32.
- Dou, X., Wang, X., Liu, H., Zhang, X., Meng, L., Pan, Y., Yu, Z., Cui, Y., 2019. Prediction of soil organic matter using multi-temporal satellite images in the Songnen Plain, China. *Geoderma* 356, 113896.
- Fisher, A., Rudin, C., Dominici, F., 2019. All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously. *J. Mach. Learn. Res.* 20, 1–81.
- Fourcade, Y., Besnard, A.G., Secondi, J., 2018. Paintings predict the distribution of species, or the challenge of selecting environmental predictors and evaluation statistics. *Glob. Ecol. Biogeogr.* 27, 245–256.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann. Stat.* 29, 1189–1232.
- Friedman, J.H., Popescu, B.E., 2008. Predictive learning via rule ensembles. *Ann. Appl. Stat.* 2, 916–954.
- Goldstein, A., Kapelner, A., Bleich, J., Pitkin, E., 2015. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *J. Computat. Graph. Stat.* 24, 44–65.
- Greenwell, B., 2020. Package "fastshap". url: <https://CRAN.R-project.org/package=fastshap> R package version 0.0.5 [Accessed 10.08.2021].
- Grömping, U., 2020. *Model-Agnostic Effects Plots for Interpreting Machine Learning Models*. Technical Report Mathematics, Physics and Chemistry, Department II, Beuth University of Applied Sciences Berlin.
- Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning*, second ed. Springer Series in Statistics, New York.
- Heuvelink, G.B.M., Angelini, M.E., Poggio, L., Bai, Z., Batjes, N.H., van den Bosch, H., Bossio, D., Estella, S., Lehmann, J., Olmedo, G.F., Sanderman, J., 2021. Machine learning in space and time for modelling soil organic carbon change. *Eur. J. Soil Sci.* 72, 1607–1623.
- Heuvelink, G.B.M., Webster, R., 2001. Modelling soil variation: past, present, and future. *Geoderma* 100, 269–301.
- Hooker, G., Mentch, L., Zhou, S., 2021. Unrestricted permutation forces extrapolation: variable importance requires at least one more model, or there is no free variable importance. *Stat. Comput.* 31, 1–16.
- Janssen, P.H.M., Heuberger, P.S.C., 1995. Calibration of process-oriented models. *Ecol. Model.* 83, 55–66.
- Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R.J., Wasserman, L., 2018. Distribution-free predictive inference for regression. *J. Am. Stat. Assoc.* 113, 1094–1111.
- Lundberg, S.M., Lee, S.-I., 2017. A unified approach to interpreting model predictions. In v. L. Ulrike, G. Isabelle, B. Samy, W. Hanna, & F. Rob (Eds.), *Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 4768–4777). Curran Associates Inc., Red Hook, New York.
- Meersmans, J., Martin, M.P., Lacerce, E., De Baets, S., Jolivet, C., Boulonne, L., Lehmann, S., Saby, N.P.A., Bispo, A., Arrouays, D., 2012. A high resolution map of French soil organic carbon. *Agronomy Sustain. Devel.* 32, 841–851.
- Molnar, C., 2020. *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*. Lulu Press, Raleigh.
- Molnar, C., Casalicchio, G., Bischl, B., 2018. iml: An R package for interpretable machine learning. *J. Open Source Software* 3, 786.
- Molnar, C., Casalicchio, G., Bischl, B., 2020. Quantifying model complexity via functional decomposition for better post-hoc interpretability. In: Cellier, P., Driessens, K. (Eds.), *Machine Learning and Knowledge Discovery in Databases*. Springer International Publishing, New York, pp. 193–204.
- Molnar, C., König, G., Herbringer, J., Freisleben, T., Dandl, S., Scholbeck, C.A., Casalicchio, G., Grosse-Wentrup, M., Bischl, B., 2022. General pitfalls of model-agnostic interpretation methods for machine learning models. In A. Holzinger, R. Goebel, R. Fong, T. Moon, K.-R. Müller, & W. Samek (Eds.), *xxAI – Beyond Explainable Artificial Intelligence*. Lecture Notes in Artificial Intelligence (pp. 55–84). Springer, Cham.
- Olden, J.D., Jackson, D.A., 2002. Illuminating the "black box": a randomization approach for understanding variable contributions in artificial neural networks. *Ecol. Model.* 154, 135–150.
- Orgiazzi, A., Ballabio, C., Panagos, P., Jones, A., Fernández-Ugalde, O., 2018. LUCAS Soil, the largest expandable soil dataset for Europe: a review. *Eur. J. Soil Sci.* 69, 140–153.
- Ottoy, S., De Vos, B., Sindayihebura, A., Hermy, M., Van Orshoven, J., 2017. Assessing soil organic carbon stocks under current and potential forest cover using digital soil mapping and spatial generalisation. *Ecol. Ind.* 77, 139–150.
- Quist, C.W., Gort, G., Mooijman, P., Brus, D.J., van den Elsen, S., Kostenko, O., Vervoort, M., Bakker, J., van der Putten, W.H., Helder, J., 2019. Spatial distribution of soil nematodes relates to soil organic matter and life strategy. *Soil Biol. Biochem.* 136, 107542.
- R Core Team, 2020. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing Vienna, Austria. url: <https://www.R-project.org/> [Accessed 10.08.2021].
- Ribeiro, M.T., Singh, S., Guestrin, C., 2016. Why should I trust you? Explaining the predictions of any classifier. In J. DeNero, M. Finlayson, & S. Reddy (Eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations* (pp. 1135–1144). Association for Computational Linguistics.
- Rivera, J.I., Bonilla, C.A., 2020. Predicting soil aggregate stability using readily available soil properties and machine learning techniques. *CATENA* 187, 104408.
- Robnik-Šikonja, M., Kononenko, I., 2008. Explaining classifications for individual instances. *IEEE Trans. Knowl. Data Eng.* 20, 589–600.
- Shapley, L.S., 1953. A value for n-person games. In: Harold William, K., Albert William, T. (Eds.), *Contributions to the Theory of Games* chapter, volume 28 of *Annals of Mathematics Studies*, 17. Princeton University Press, Princeton, pp. 31–40.
- Strobl, C., Boulesteix, A.-L., Zeileis, A., Hothorn, T., 2007. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinform.* 8, 1–21.
- Štrumbelj, E., Kononenko, I., 2014. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.* 41, 647–665.
- Vågen, T.-G., Winowiecki, L.A., Tondoh, J.E., Desta, L.T., Gumbrecht, T., 2016. Mapping of soil properties and land degradation risk in Africa using MODIS reflectance. *Geoderma* 263, 216–225.
- Vos, C., Don, A., Hobley, E.U., Prietz, R., Heidkamp, A., Freibauer, A., 2019. Factors controlling the variation in organic carbon stocks in agricultural soils of Germany. *Eur. J. Soil Sci.* 70, 550–564.
- Wadoux, A.M.J.-C., Dennis J J, W., Brus, D.J., 2022. An integrated approach for the evaluation of quantitative soil maps through Taylor and solar diagrams. *Geoderma*, 405, 115332.
- Wadoux, A.M.J.-C., Heuvelink, G.B.M., Lark, R.M., Lagacherie, P., Bouma, J., Mulder, V. L., Libohova, Z., Yang, L., McBratney, A.B., 2021. Ten challenges for the future of pedometrics. *Geoderma* 401, 115155.
- Wadoux, A.M.J.-C., Minasny, B., McBratney, A.B., 2020. Machine learning for digital soil mapping: applications, challenges and suggested solutions. *Earth Sci. Rev.* 210, 103359.
- Wadoux, A.M.J.-C., Samuel-Rosa, A., Poggio, L., Mulder, V.L., 2020. A note on knowledge discovery and machine learning in digital soil mapping. *Eur. J. Soil Sci.* 71, 133–136.
- Watson, D.S., Wright, M.N., 2021. Testing conditional independence in supervised learning algorithms. *Mach. Learn.* 110, 2107–2129.
- Wright, M.N., Ziegler, A., 2017. ranger: A fast implementation of random forests for high dimensional data in C++ and R. *J. Stat. Softw.* 77, 1–17.
- Zeng, C., Yang, L., Zhu, A.-X., 2017. Construction of membership functions for soil mapping using the partial dependence of soil on environmental covariates calculated by random forest. *Soil Sci. Soc. Am. J.* 81, 341–353.