# 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*

**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 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.

**This manuscript is a non-peer reviewed preprint that has been submitted for publication. Subsequent versions of this manuscript may have updated content. Feedback and comments are welcomed, feel free to contact the corresponding author:**
**Alexandre Wadoux**
alexandre.wadoux@sydney.edu.au

*Keywords:* Digital soil mapping, Machine learning, Geostatistics, Shapley, Partial dependence, H-statistic, Accumulated local effect, Surrogate modelling

*Corresponding author: Sydney Institute of Agriculture & School of Life and Environmental Sciences, The University of Sydney, New South Wales, Australia

*Email address:* alexandre.wadoux@sydney.edu.au (Alexandre M.J-C. Wadoux)

**Highlights**

- Describe a set of methods for the interpretation of complex mapping models.

- Methods are model-independent.

- Reveal global model functioning and local drivers of soil variation.

- Provide new insights into complex mapping models.

- Might assist formulating hypotheses on the mechanisms driving soil variation.

1. **1. Introduction**

2. Understanding the spatial variation of soil properties has become central to many sub-disciplines
3. of soil science. Digital soil mapping (DSM) techniques can be used for this purpose. Commonly
4. in DSM studies, statistical or non-statistical models are calibrated to exploit the quantitative re-
5. lationship between measured values of a soil property and a set of environmental covariates of
6. which maps are available, such as satellite imagery and terrain attributes. These models are used
7. to predict the soil property at unobserved locations and to identify and expose the importance
8. of environmental factors in the soil property spatial variation. Recent examples of studies using
9. this approach are Quist et al. (2019) for mapping soil nematodes and Heuvelink et al. (2021) for
10. mapping soil organic carbon in space and time.
11.

12. Since early soil mapping studies rooted in classical statistics and design-based inference in the
13. 70s, and based on geostatistics on the 80s (Heuvelink & Webster, 2001), the field has gradually
14. shifted attention towards more complex statistical and algorithmic tools from the field of machine
15. learning. Accuracy of such models is often higher than that of classical models. They are also
16. particularly useful in situation where the relationship between the soil property and environmental
17. covariates is too complex to be modelled mechanistically or with simple statistical models. How-
18. ever, popularization of complex models of soil variation was made at the expense of understanding
19. why the soil varies the way it does. Insight into the functioning and structure of the models are
20. difficult to obtain, so that these models are often referred to as "black boxes". Examples of such
21. models are random forest, support vector machines and neural networks. We refer to Hastie et al.
22. (2009) for an overview.
23.

24. In soil science, several attempts were made to obtain insights from complex models. The relative
25. effect of environmental covariates on model prediction is usually characterized by model-specific
26. variable importance statistics such as through the mean decrease in impurity for tree-based mod-
27. els (as is done in Vos et al., 2019, for example), or by calculating the partial dependence of the
28. prediction to environmental covariates (e.g. Zeng et al., 2017; Ottoy et al., 2017). For artificial
29. neural networks, the Garson's algorithm or the magnitude and direction of the connection weights
30. between neurons give indication on the variable importance (Olden & Jackson, 2002). An example
31. in soil mapping is Rivera & Bonilla (2020). While valid and useful to obtain insights into complex
32. models of soil variation, these methods are model-specific, i.e. they preclude comparison between
33. models (Wadoux et al., 2020a). A number of "model-agnostic" or model-independent interpre-
34. tation methods have recently been developed outside soil science, in the statistical and machine
35. learning literature. Model-independent means that these model interpretation methods are ap-
36. plicable to any model. It is worthwhile to introduce these recent developments, and to present a
37. strategy for the interpretation of complex soil mapping models. This was also recently highlighted
38. as one of the most pressing pedometric research topics (Wadoux et al., 2021b, Challenge 3).
39.

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. This 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 complex mapping models in a case study mapping topsoil organic carbon in France. Finally, we discuss 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 $\mathbf{s}$ 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 $\mathrm{E}\left[(Y - \hat{Y})^2\right]$.

Let $y(\mathbf{s}_i)$ be $n$ measurements of $Y$, $\mathbf{s}_i$ $(i = 1, \ldots, n; \mathbf{s}_i \in \mathcal{A})$ 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 $i^{\text{th}}$ row-vector and the $j^{\text{th}}$ column vector of $\mathbf{X}$, respectively, and $x_{i,j}$ 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. (2021a, 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 & 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 $j = 1, \ldots, p$:
   (a) Create modified (denoted by the asterisk $*$) covariate matrix $\mathbf{X}^*$ by permutation of the values in the $j^{\text{th}}$ 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 $j^{\text{th}}$ covariate by the ratio $\ell(\hat{f}(\mathbf{X}^*), \mathbf{y}) / \ell(\hat{f}(\mathbf{X}), \mathbf{y})$ or the difference $\ell(\hat{f}(\mathbf{X}^*), \mathbf{y}) - \ell(\hat{f}(\mathbf{X}), \mathbf{y})$.

3

Permutation of the covariate matrix involves randomness and is usually repeated to obtain a distribution of the importance metric. Figure 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 instead of a single covariate.
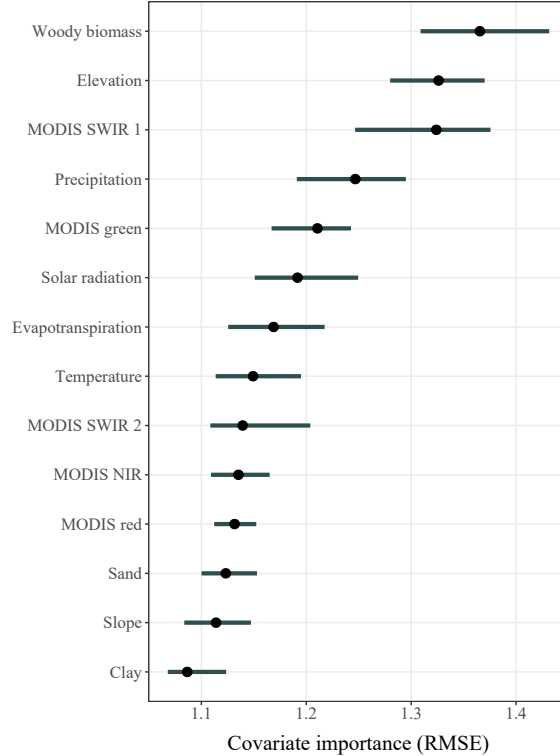


Figure 1: Example of a covariate importance with permutation assessed by the ratio of RMSE. The black dots represent the mean value of 100 permutations and the green lines the 90% confidence interval. Data and model from Wadoux et al. (2021a, Section 4.2).

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 permutation importance is given by Hooker & Mentch (2019). When covariates are correlated, alternatives to permutation importance are the conditional permutation strategy (Strobl et al., 2007; Watson & Wright, 2019; Fisher et al., 2019; Molnar et al., 2020b) and the dropped variables importance (Lei et al., 2018).

## 2.2. Individual conditional expectation

Individual conditional expectation (ICE, Goldstein et al., 2015) shows how the prediction at a location would change when the considered covariate would vary. Consider the subset of covariates $\mathbf{X}_{\mathcal{S}}$ of $\mathbf{X}$ composed of $l < p$ covariates, and $\mathbf{X}_{\mathcal{C}}$ its complement so that $f(\mathbf{X}) = f(\mathbf{X}_{\mathcal{S}}, \mathbf{X}_{\mathcal{C}})$. The subset $\mathbf{X}_{\mathcal{S}}$ usually contains one or two covariates (i.e. $l \approx 1, 2$). For any location in $\mathcal{A}$ with covariate values $(\mathbf{x}_{i,\mathcal{S}}, \mathbf{x}_{i,\mathcal{C}})$ and calibrated model $\hat{f}$, an ICE curve $\hat{f}_{\text{ICE}}$ shows model predictions

4

for a grid of $\mathbf{x}_{i,\mathcal{S}}$ while keeping fixed the values of $\mathbf{x}_{i,\mathcal{C}}$ (Fig. 2a).

When comparing ICE curves, it is convenient to center the individual ICE curves to a baseline value. 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:

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

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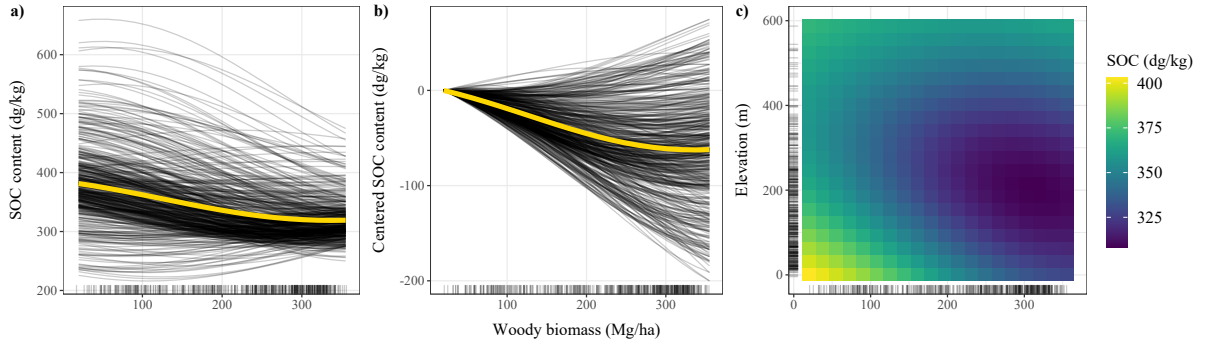


Figure 2: Examples of a) individual conditional expectation (ICE) curves (in black) for woody biomass against soil organic carbon (SOC) content. The yellow curve it the partial dependence plot (PDP). In b), ICE curves and the PDP are centered at the minimum of the covariate value (i.e. at a woody biomass value of 20). Plot c) shows the two-dimensional PDP of woody biomass against elevation. Data and model from Wadoux et al. (2021a, Section 4.2).

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 $f_{\text{PDP}}$ of $\hat{f}(\mathbf{X})$ on $\mathbf{X}_{\mathcal{S}}$ is formally expressed as the expected value of the model prediction over the distribution of the covariates in the subset $\mathcal{C}$ (Friedman, 2001):

$$f_{\text{PDP}}(\mathbf{x}_{\mathcal{S}}) = \mathrm{E}_{\mathbf{X}_{\mathcal{C}}}[\hat{f}(\mathbf{x}_{\mathcal{S}}, \mathbf{X}_{\mathcal{C}})]. \tag{2}$$

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

$$\hat{f}_{\text{PDP}}(\mathbf{x}_{\mathcal{S}}) = \frac{1}{n} \sum_{i=1}^{n} \hat{f}(\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{i\mathcal{C}}), \tag{3}$$

where $\mathbf{x}_{1\mathcal{C}}, \mathbf{x}_{2\mathcal{C}}, \ldots, \mathbf{x}_{n\mathcal{C}}$ are the row-vectors of $\mathbf{X}_{\mathcal{C}}$. 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 $\mathcal{S}$ of any size, only one or two covariates can reasonably be displayed. Note also that dependence among covariates in $\mathbf{X}_{\mathcal{S}}$ and $\mathbf{X}_{\mathcal{C}}$ 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 $\mathbf{X}_{\mathcal{S}}$ and $\mathbf{X}_{\mathcal{C}}$ are uncorrelated (Hastie et al., 2009, p. 370).

## 2.4. Accumulated local effect

An alternative to the PDP when covariates are dependent is the accumulated local effect (ALE, Apley & Zhu, 2020). The ALE shows the effect of changing the values of a covariate on the soil property. 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 $z_{0,\mathcal{S}}$.

$$f_{\mathrm{ALE}}(\mathbf{x}_{\mathcal{S}}) = \int_{z_{0,\mathcal{S}}}^{\mathbf{x}_{\mathcal{S}}} \mathrm{E}_{X_{\mathcal{C}}|X_{\mathcal{S}}} \left[ \hat{f}^{\mathcal{S}}(X_{\mathcal{S}}, X_{\mathcal{C}}) | X_{\mathcal{S}} = z_{\mathcal{S}} \right] \mathrm{d}z_{\mathcal{S}}, \qquad (4)$$

where $\hat{f}^{\mathcal{S}}(\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{C}}) = \frac{\delta \hat{f}(\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}})}{\delta \mathbf{x}_{\mathcal{S}}}$ is the derivative of the prediction function with respect to covariates $\mathbf{x}_S$. 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. Equation 4 can be approximated by a step function over the $K$ intervals:

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

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 smooth by linear interpolation. The ALE is centered at zero by:

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

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 & Zhu (2020, Eq. 13-16) for the equations describing the two-dimensional

165 ALE and see Molnar (2020, Chapter 5) for more details on the difference between PDP and ALE.
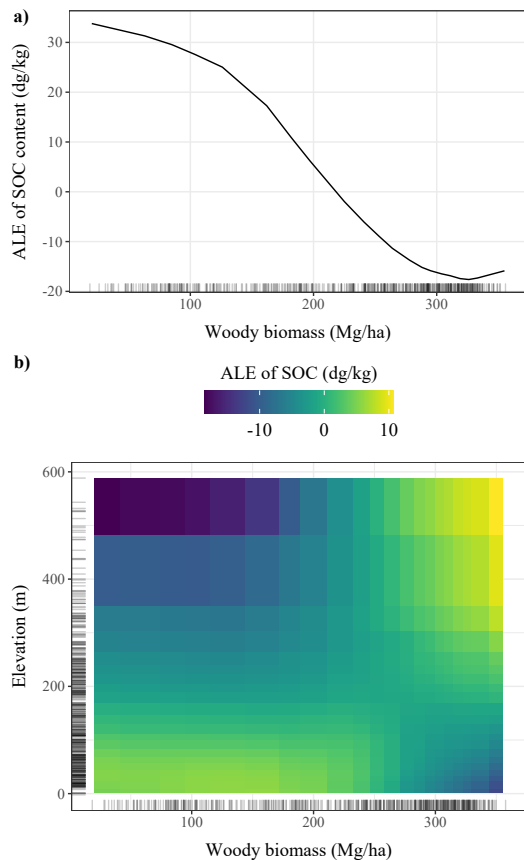166 An example of one and two-dimensional ALE plot is shown in Fig. 3.
167



Figure 3: Estimates of a one-dimensional accumulated local effect (ALE) plot of woody biomass on SOC content (a), and two-dimensional ALE of woddy biomass and elevation on SOC (b). Data and model from Wadoux et al. (2021a, Section 4.2).

168 Note that interpretation of the two-dimensional ALE plot is different from that of a PDP. ALE
169 is formally interpreted as being the centered difference in prediction (i.e. the effect) when the
170 observations within an interval are moved from one border of the interval to another other. Fig. 3a
171 shows the effect of woody biomass on SOC for a range of values of woody biomass, and compared
172 to the average prediction. Fig. 3b shows the pure interaction effect of woody biomass and elevation
173 on SOC compared to the average prediction. For example, the ALE estimate of woody biomass
174 in Fig. 3a illustrates that for large values of woody biomass (i.e. greater than 300 Mg/ha), the
175 predicted values of SOC are lower by nearly 20 dg/kg compared to the average prediction.
176
177 The estimates of ALE tend to be more robust than the PDP for correlated covariates, because of
178 averaging and accumulating the local effect over the conditional distribution. However, this comes
179 at the expense of determining of having a more localized interpretation (within intervals), and
180 possibly non-intuitive interpretations for some data-generating processes (Grömping, 2020).

7

### 2.5. Interaction between covariates

Interaction between covariates can be estimated with the H-statistic (Friedman & 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 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 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:

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

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

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

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.
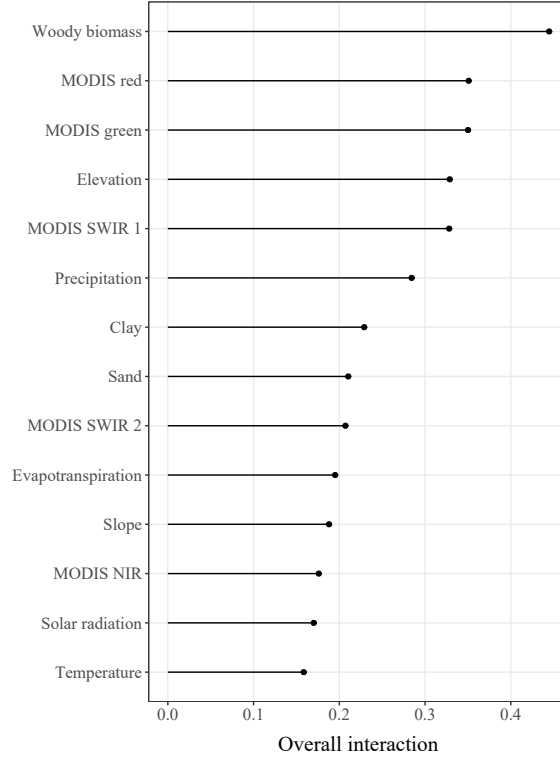
Figure 4: Estimate of the total interaction (Eq. 8) between 14 covariates used for prediction of SOC. Data and model from Wadoux et al. (2021a, Section 4.2).

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 $\mathbf{X}$ is approximated by calibrating a simple model $g$ on the $n$ prediction $\hat{y}(\mathbf{s}_i)$. 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 & Heuberger, 1995):

$$\text{MEC} = 1 - \frac{\sum_{i=1}^{n}(\hat{y}(\mathbf{s}_i) - \hat{y}^*(\mathbf{s}_i))^2}{\sum_{i=1}^{n}(\hat{y}(\mathbf{s}_i) - \overline{\hat{y}})^2}, \tag{9}$$

where $\hat{y}$ denote the predicted soil property at location $\mathbf{s}_i$ by model $\hat{f}$, and $\hat{y}^*$ 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.

9

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.

### 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 composed of $l < p$ covariates. $\mathcal{S} \subseteq \{1, \ldots, p\} \setminus \{j\}$ refers to any subset of covariates which excludes covariate $j$. The Shapley value $\phi$ for covariate $j$ for a data point $\mathbf{x}_0$ (not necessarily from the original data set) is given by:

$$\phi_{0,j} = \sum_{\mathcal{S} \subseteq \{1,\ldots,p\}-\{j\}} \frac{|\mathcal{S}|! \, (p - |\mathcal{S}| - 1)!}{p!} \left( \hat{f}^* \left( \mathbf{x}_{i,\mathcal{S} \cup \{j\}} \right) - \hat{f}^*(\mathbf{x}_{i,\mathcal{S}}) \right), \tag{10}$$

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 $\hat{f}^*(\mathbf{x}_{i,\mathcal{S}}) = \mathrm{E}_{X_C}[\hat{f}(\mathbf{x}_{i,\mathcal{S}}, X_C)]$ is the prediction function where covariates not contained in $\mathcal{S}$ are marginalized (similar for $\mathcal{S} \cup \{j\}$). Then $\hat{f}^* \left( \mathbf{x}_{i,\mathcal{S} \cup \{j\}} \right) - \hat{f}^* \left( \mathbf{x}_{i,\mathcal{S}} \right)$ 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 $\phi_{i,0}$.

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 & Kononenko (2014) and Lundberg & Lee (2017) proposed estimation methods to reduce the computational cost. Štrumbelj & 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 & 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 & Kononenko (2014).
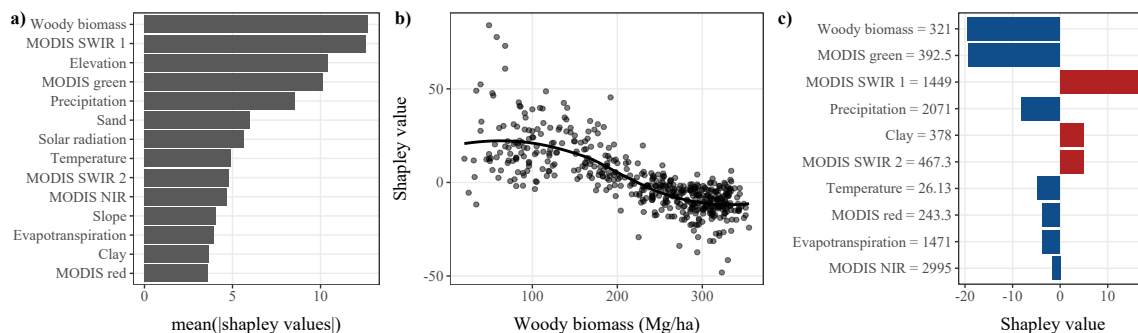
Figure 5: Average of the absolute Shapley values in the calibration dataset (a), dependence plot of SOC against woody biomass (b), and local interpretation of a single spatial location (c). Data and model from Wadoux et al. (2021a, Section 4.2).

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 (see also 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 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 seven factors influencing SOC spatial distribution: topography, vegetation, long-term average climatic conditions, climate seasonality, extreme climatic conditions and satellite imagery. 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 m × 250 m. The SOC data and their matching values of environmental covariates were then used to calibrate two mapping models.

The first model being used is 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. All other parameters where held to their default value. We used the R programming language (R Core Team, 2020) for the implementation and the R package ranger (Wright & Ziegler, 2017). The second model being used is 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 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.
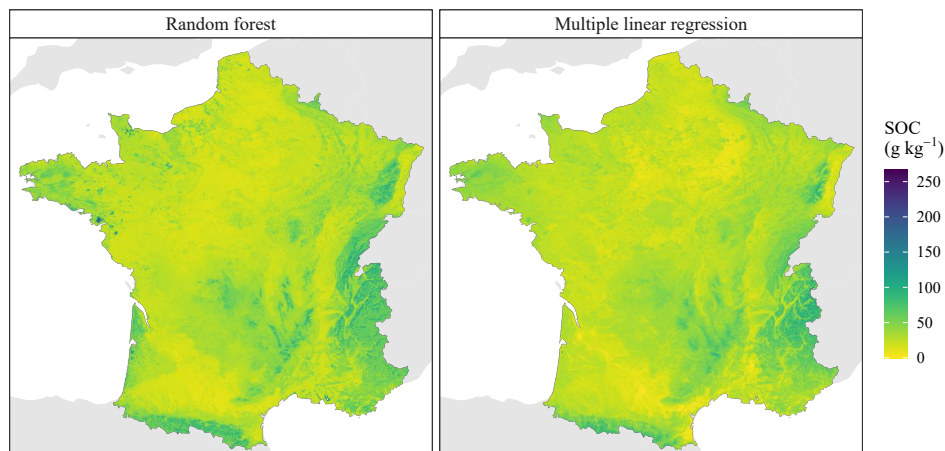


Figure 6: Spatial distribution of SOC (in g kg$^{-1}$) for Metropolitan France excluding Corsica. The SOC maps were made using random forest (left) and multiple linear regression (right).

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 two contrasting spatial locations and to a geographical area (Fig. 7). This allows us to understand how the importance of environmental covariates vary from one location to another and in space. The two spatial locations are denoted *Beauce* and *Landes*. Location *Beauce* is in a cropland-dominated region with fertile clay and/or silt-loam soils but relatively low carbon content due intensive agriculture whereas *Landes* is a coniferous forested area with sandy soils (i.e. Podsols), but with relatively high topsoil carbon content due to little interest in these soils for agricultural purposes (Meersmans et al., 2012). 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. Implementation of the interpretation methods was made with the R packages iml (Molnar et al., 2018) and fastshap (Greenwell, 2020).
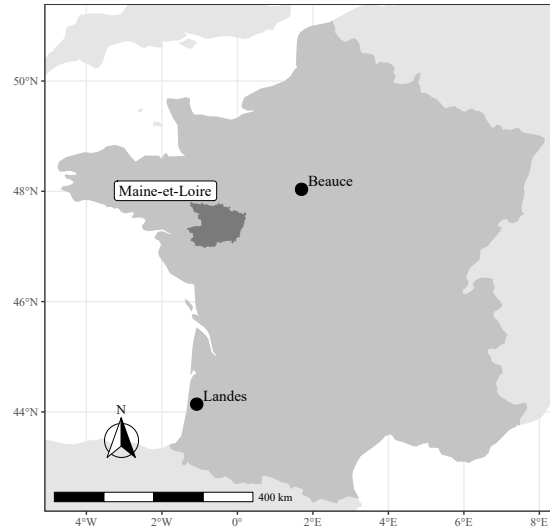
Figure 7: Location of the two spatial locations and the geographical area for the implementation of the local interpretation methods. The two black dots represent two spatial locations with contrasting SOC content. They are called *Beauce* and *Landes*. The dark grey area is called *Maine-et-Loire* and represents an administrative unit.

### 3.1. Global interpretation

*Which are the drivers of SOC spatial variation?*

Figure 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 satelite 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. Covariate representing soil water content for 1500kPa suction is, conversely, not essential to the RF model, because close to a ratio of RMSE value of 1 (i.e. removing covariate soil water content for 1500kPa does not affect model prediction accuracy). Figure 8b-c shows the covariate importance for group of covariates, for both RF (fig. 8b) and MLR (fig. 8c). All group of covariates are important in the RF model. Vegetation, and 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 group of covariates representing extreme and average climate conditions are the most important.
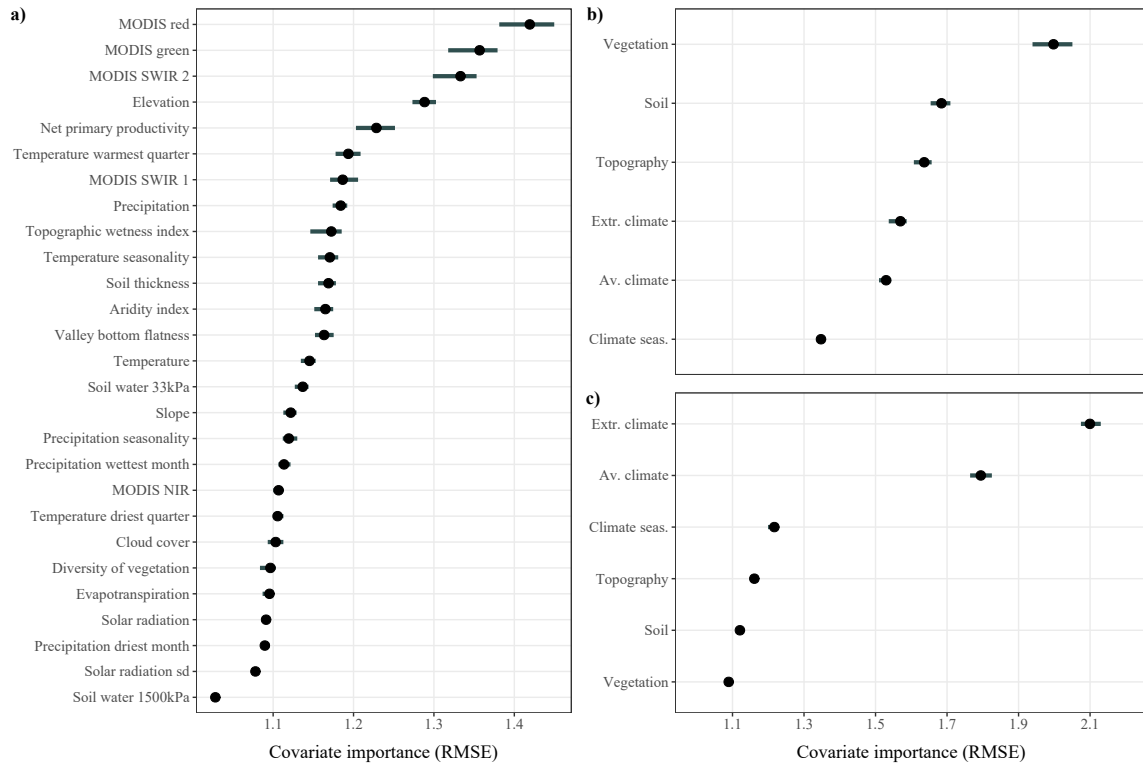
Figure 8: Mean and 90% confidence interval of the permutation-based covariate importance for a) all covariates of the random forest model, b) group of covariates for the random forest model and c) group of covariates in the multiple linear regression model. Covariate importance is assessed by the ratio of RMSE over 100 permutations. We refer to the Supplementary Material for information on the group of covariates.

Figure 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. Figure 9a indicates that the most important covariates are MODIS images and elevation. The overall ranking of covariate importance obtained by Shapley values is similar to that found with the permutation-based method. Figure 9b shows the covariate contribution to each individual location found in the calibration dataset. Most important covariates (e.g. MODIS red) have a large range of Shapley values (i.e. between -10 and 25), meaning that this covariate can have a relatively important contribution to the model prediction. Figure 9b also provides insight into the relationship between the relative covariate contriution 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.
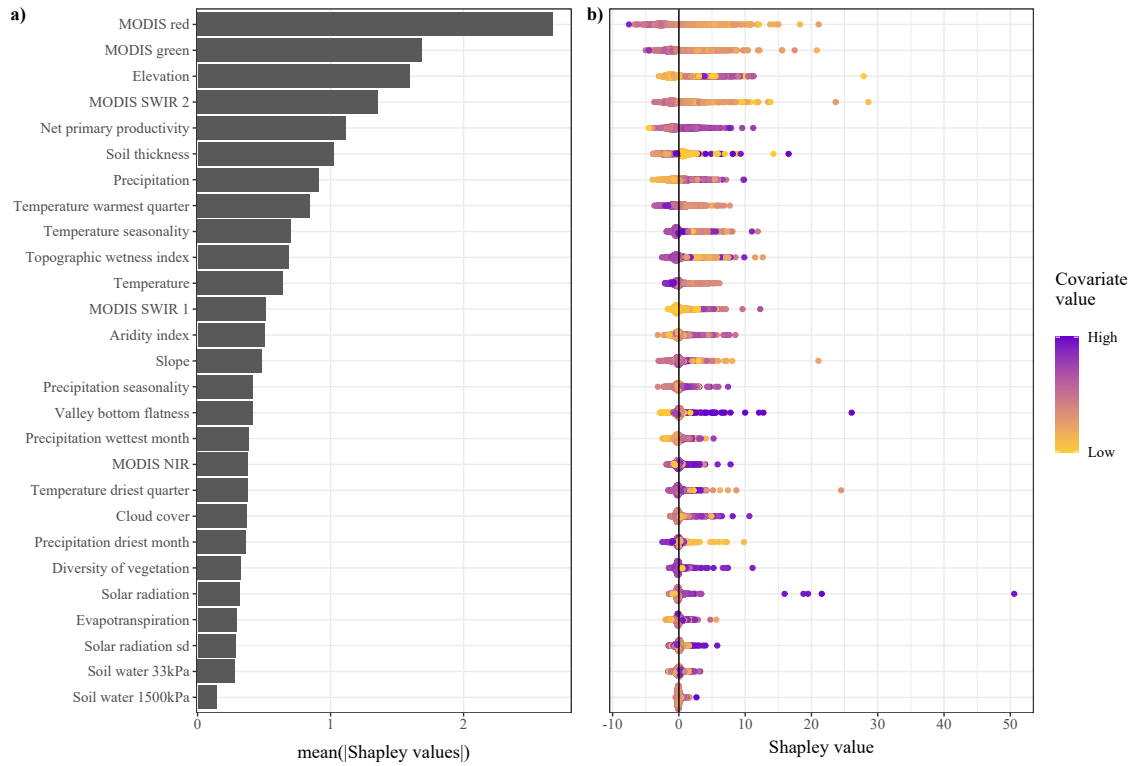
14

Figure 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. 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. Plot a) represents the averaged absolute values of plot b). The colour in b) represents the covariate value normalized in the range $(0, 1)$.

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

Figure 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 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 (the black dots).
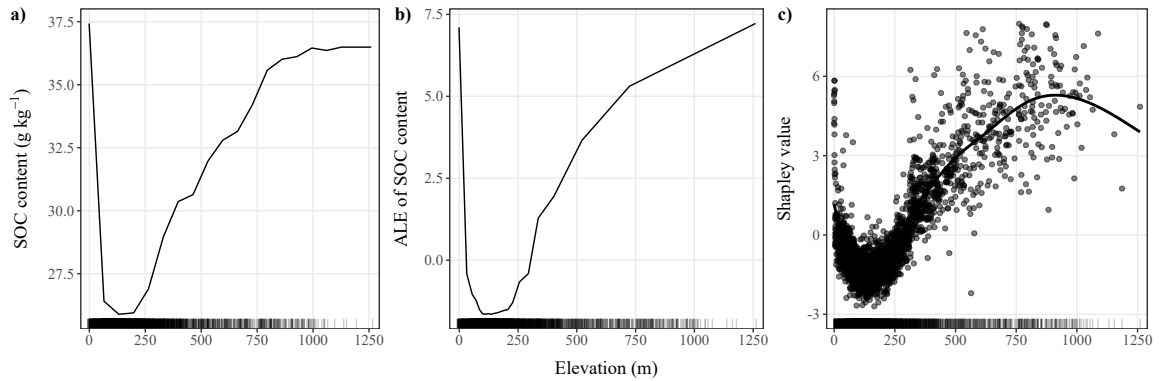
15

Figure 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.
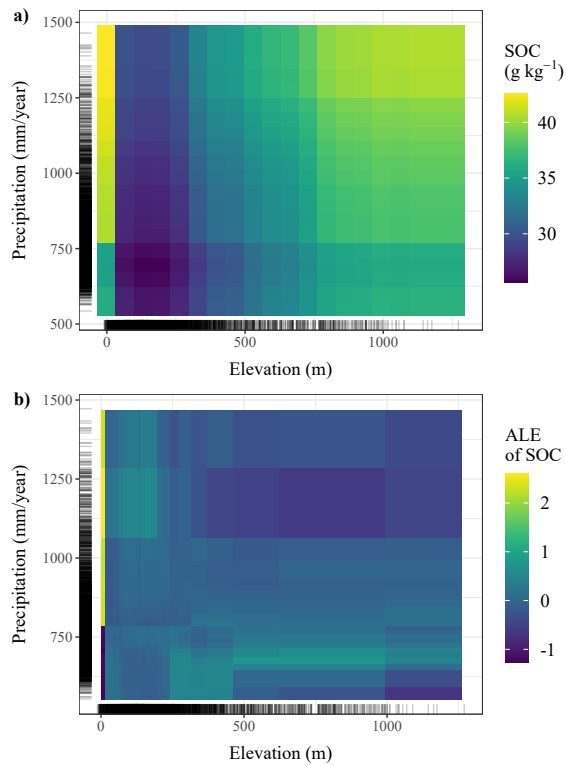


Figure 11: Two-dimensional partial dependence plot of the effect of elevation and temperature on SOC content (a), and accumulated local effect of elevation and temperature on SOC content (b).

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,

16

<sup>346</sup> while an opposite pattern is seen for elevation values larger than 250. In Fig. 11, both plots have
<sup>347</sup> a noticeable increasing pattern of SOC with higher precipitation, but only for low relief. Above an
<sup>348</sup> elevation of 1000 m, few SOC observations exist, which means that interpretations of effects for
<sup>349</sup> this elevation should be cautious.

<sup>350</sup> *How does SOC prediction depend on interactions among covariates?*

<sup>351</sup> Figure 12 shows the strength of the interaction between environmental covariates for the RF
<sup>352</sup> model. Note that the MLR is not expected to contain an interaction effect between covariates
<sup>353</sup> unless explicitly specified. Fig. 12a shows the presence of a strong overall interaction effect in
<sup>354</sup> the random forest model. Satellite imageries MODIS red, green and SWIR 2 are involved in
<sup>355</sup> interactions with other covariates. Elevation also substantially interacts with other covariates.
<sup>356</sup> Covariates standard deviation of monthly solar radiation and soil water content, conversely, have
<sup>357</sup> negligible interaction. Fig. 12b identifies how strong covariates interacts with elevation. Elevation is
<sup>358</sup> dominantly interacting with MODIS SWIR 2, precipitation seasonality and topographic covariates
<sup>359</sup> (e.g. wetness index). There is no strong interaction of elevation with soil water content, solar
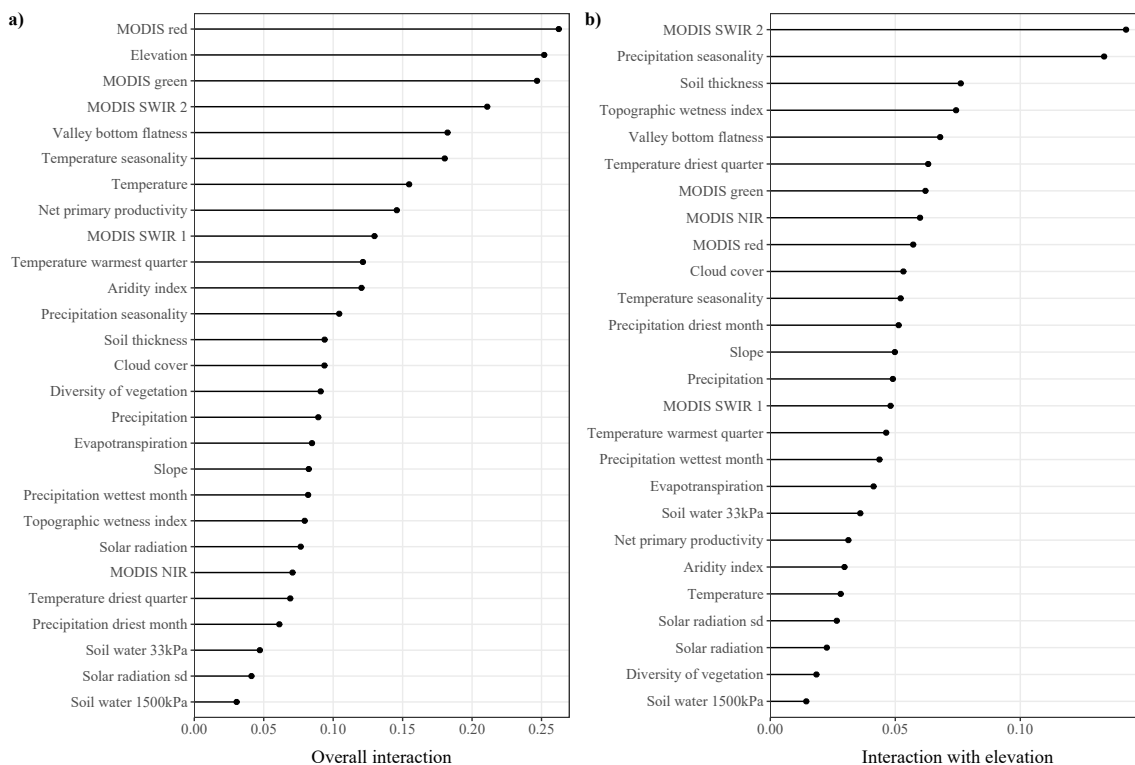<sup>360</sup> radiation and diversity of vegetation.



Figure 12: Estimate of the overall interaction (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.

<sup>361</sup> *How to summarize the model?*

<sup>362</sup> Figure 13 shows a surrogate model of the RF model. The surrogate model is a simple decision
<sup>363</sup> tree with a depth of three. It has a MEC of 0.3. The final nodes show the average predicted
<sup>364</sup> value and the percentage of data in the node. The colour of the final node is proportional to the
<sup>365</sup> 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, smallest predicted values of SOC (i.e. SOC <= 24) are found for locations with large values of the MODIS red band and low elevation (< 312 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 (i.e. less then 18 degrees) and precipitation of the driest month are relatively abundant (more than 67 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.
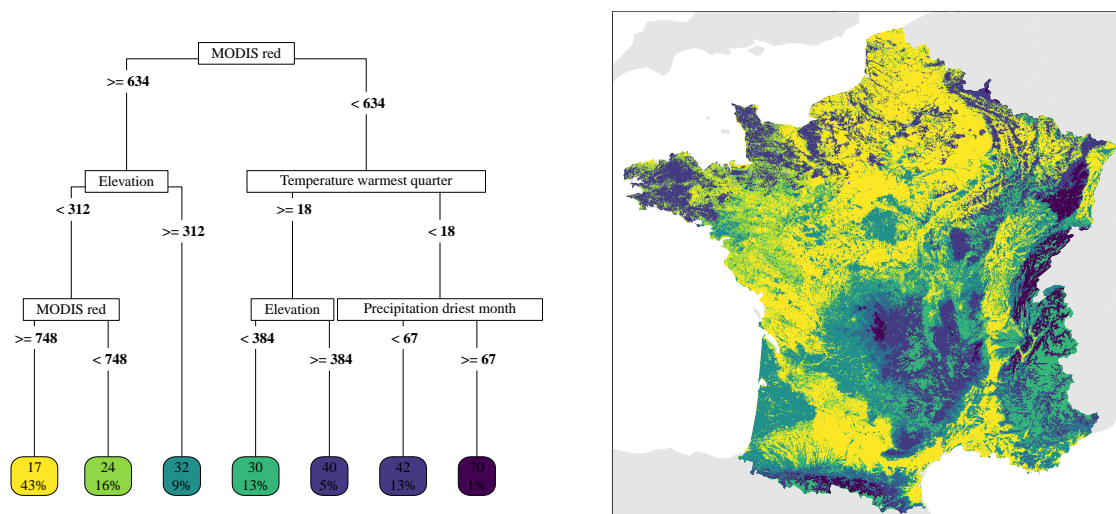


Figure 13: Surrogate model of the random forest model for prediction of SOC. The surrogate model (left) is a decision tree. The final node shows the average predicted value and the percentage of data in the node. The colour of the final node is proportional to the value in the node. The colour scheme is reported in France (right) using the rules of the decision tree.

### 3.2. Local interpretation

*What is the local functional form of the association between environmental covariates and SOC?*
Figure 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. $> 8$ 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 different from that observed on average for France, where elevation has a positive relationship with SOC content (see also Fig. 10a-b).
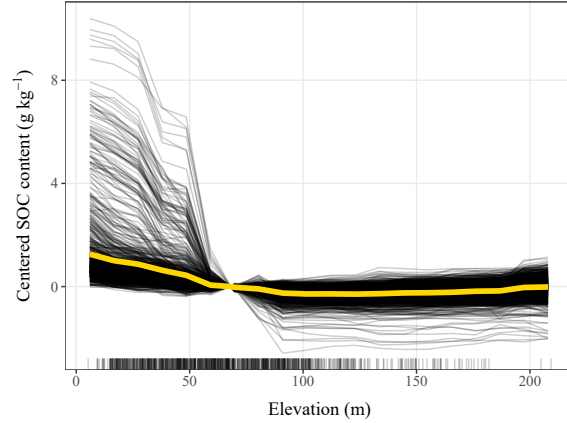
Figure 14: Centered effect of elevation on the SOC in the region *Maine-et-Loire*. The effect is centered at the average elevation value of the area ($x_0 = 68$). The black curves are the individual conditional expectation whereas the yellow curve is their average (i.e. their partial dependence function).

Figure 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 shown in Fig. 7. Figure 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, repectively. For random forest, the ICE curves have more variation: 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.
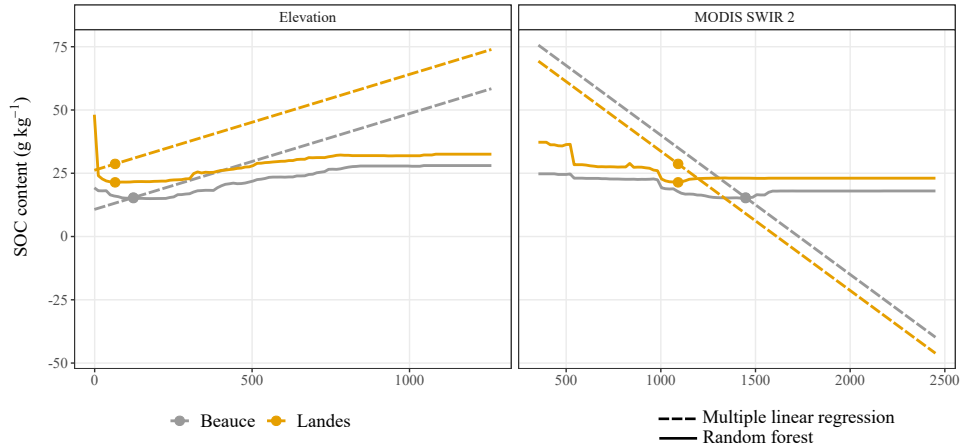


Figure 15: Individual expectation curves of the effect of elevation and MODIS SWIR 2 on SOC for the two locations of interest *Beauce* and *Landes* and the multiple linear regression and random forest models. The dots represent the SOC prediction made by the model at the locations.

19

*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 pattern. 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 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.
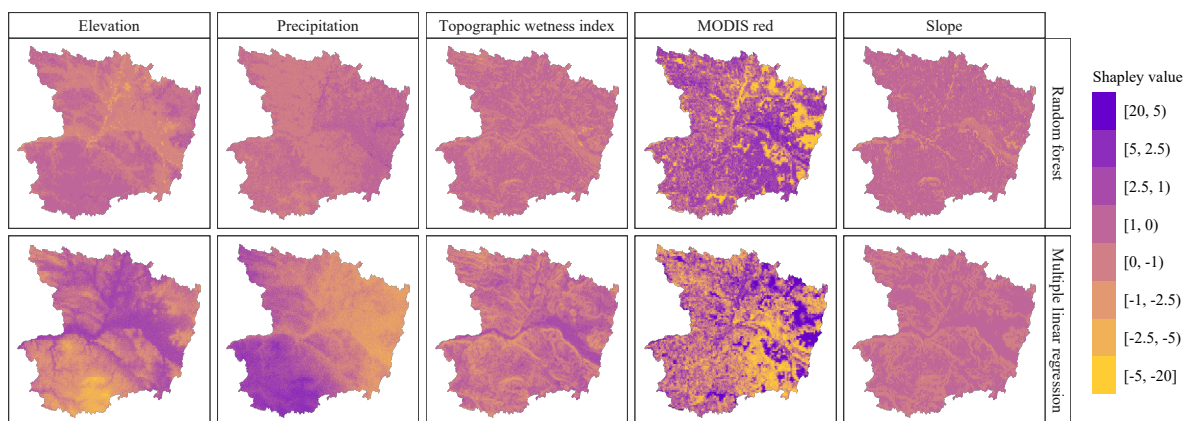


Figure 16: Spatial pattern of the Shapley values for five covariates and the two mapping models. Dark colour indicates that the covariate has a positive contribution to the SOC prediction while light colour indicates a negative contribution.

Figure 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 difference between the sum of Shapley values and the prediction is due to the approximation strategy. Fig. 17 shows that SOC prediction in the two spatial locations in 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, 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.
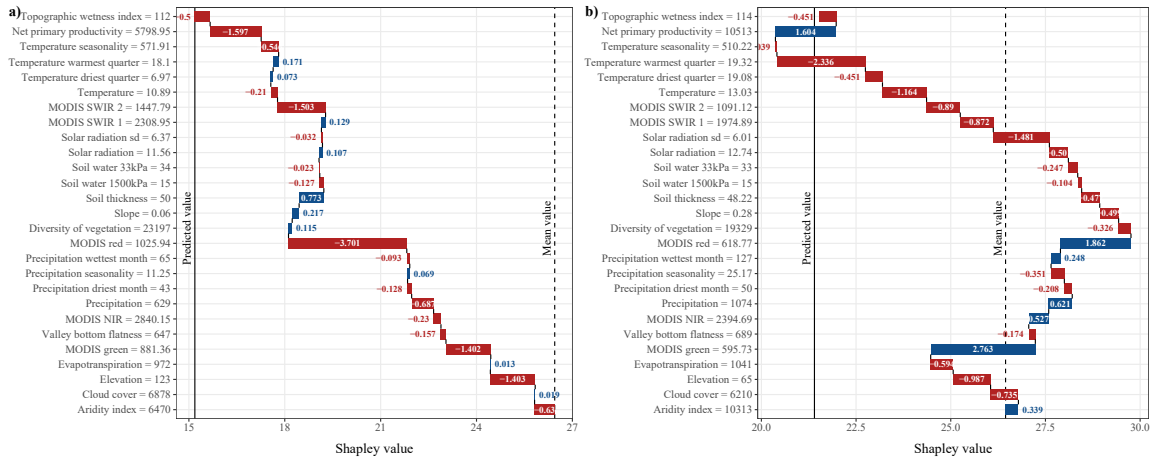
Figure 17: Contribution of the individual covariates to the prediction of SOC at location Beauce (a) and Landes (b). Contributions are estimated with Shapley values. The red colour indicates a negative contribution while a blue colour indicates a positive contribution. The y-axis indicates the value of the covariate at the prediction location.

## 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 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 in content lower than in a sandy soil covered by 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, 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

21

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. (2020c, 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), or by using variants that rely on the conditional distribution (e.g. conditional feature importance, Molnar et al., 2020b). 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.

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 & 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 & 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 & Kononenko (2014), Lundberg & Lee (2017) |

As mentioned in the Introduction the aim of this paper is 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 take out of the realm 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. Several recent studies have argued in this sense (e.g. Fourcade et al., 2018). Wadoux et al. (2020b), 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 diagnostic 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 for additive (Robnik-Šikonja & Kononenko, 2008) and non-additive (Gosiewska & Biecek, 2019) attribution, functional decomposition (Molnar et al., 2020a), or local interpretable model-agnostic explanations (LIME, Ribeiro et al., 2016). LIME is being 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 these 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 us insights into the nature of soil variation. Notably, geostatistical models are powerful for prediction and provision 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 & Burzykowski (2021, Section 1.5) for an overview and to Molnar et al. (2020b, 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 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 predictor and two complex 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 us formulating hypotheses on 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 evaluation of model prediction plausibility.

## Acknowledgment

## References

Apley, D. W., & Zhu, J. (2020). Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society: Series 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. *Machine Learning*, *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. *Journal of Machine Learning Research*, *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. *Global Ecology and Biogeography*, *27*, 245–256.

24

Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, *29*, 1189–1232.

Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, *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. *Journal of Computational and Graphical Statistics*, *24*, 44–65.

Gosiewska, A., & Biecek, P. (2019). IBreakDown: Uncertainty of model explanations for non-additive predictive models. arXiv:1903.11420.

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*. (2nd 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. *European Journal of Soil Science*, *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. (2019). Please stop permuting features: An explanation and alternatives. arXiv:1905.03151.

Janssen, P. H. M., & Heuberger, P. S. C. (1995). Calibration of process-oriented models. *Ecological Modelling*, *83*, 55–66.

Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R. J., & Wasserman, L. (2018). Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, *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., Lacarce, 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 for Sustainable Development*, *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. *Journal of Open Source Software*, *3*, 786.

25

Molnar, C., Casalicchio, G., & Bischl, B. (2020a). Quantifying model complexity via functional decomposition for better post-hoc interpretability. In P. Cellier, & K. Driessens (Eds.), *Machine Learning and Knowledge Discovery in Databases* (pp. 193–204). Springer International Publishing, New York.

Molnar, C., König, G., Bischl, B., & Casalicchio, G. (2020b). Model-agnostic feature importance and effects with dependent features–a conditional subgroup approach. arXiv:2006.04628.

Molnar, C., König, G., Herbinger, J., Freiesleben, T., Dandl, S., Scholbeck, C. A., Casalicchio, G., Grosse-Wentrup, M., & Bischl, B. (2020c). General pitfalls of model-agnostic interpretation methods for machine learning models. arXiv:2007.04131.

Olden, J. D., & Jackson, D. A. (2002). Illuminating the "black box": a randomization approach for understanding variable contributions in artificial neural networks. *Ecological Modelling*, *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. *European Journal of Soil Science*, *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. *Ecological Indicators*, *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 Biology and Biochemistry*, *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 Transactions on Knowledge and Data Engineering*, *20*, 589–600.

Shapley, L. S. (1953). A value for n-person games. In K. Harold William, & T. Albert William (Eds.), *Contributions to the Theory of Games* chapter 17. (pp. 31–40). Princeton University Press, Princeton volume 28 of *Annals of Mathematics Studies*.

Strobl, C., Boulesteix, A.-L., Zeileis, A., & Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics*, *8*, 1–21.

Štrumbelj, E., & Kononenko, I. (2014). Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, *41*, 647–665.

Vågen, T.-G., Winowiecki, L. A., Tondoh, J. E., Desta, L. T., & Gumbricht, 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. *European Journal of Soil Science*, *70*, 550–564.

Wadoux, A. M. J.-C., Dennis J J, W., & Brus, D. J. (2021a). 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. (2021b). Ten challenges for the future of pedometrics. *Geoderma*, *401*, 115155.

Wadoux, A. M. J.-C., Minasny, B., & McBratney, A. B. (2020a). Machine learning for digital soil mapping: applications, challenges and suggested solutions. *Earth-Science Reviews*, *210*, 103359.

Wadoux, A. M. J.-C., Samuel-Rosa, A., Poggio, L., & Mulder, V. L. (2020b). A note on knowledge discovery and machine learning in digital soil mapping. *European Journal of Soil Science*, *71*, 133–136.

Watson, D. S., & Wright, M. N. (2019). Testing conditional independence in supervised learning algorithms. arXiv:1901.09917.

Wright, M. N., & Ziegler, A. (2017). ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, *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 Science Society of America Journal*, *81*, 341–353.