

Review Article

Machine learning for digital soil mapping: Applications, challenges and suggested solutions



Alexandre M.J.-C. Wadoux*, Budiman Minasny, Alex B. McBratney

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

ARTICLE INFO

Keywords:

Soil science
Pedometrics
Data mining
Spatial data
Geostatistics
Random forest

ABSTRACT

The uptake of machine learning (ML) algorithms in digital soil mapping (DSM) is transforming the way soil scientists produce their maps. Within the past two decades, soil scientists have applied ML to a wide range of scenarios, by mapping soil properties or classes with various ML algorithms, on spatial scale from the local to the global, and with depth. The wide adoption of ML for soil mapping was made possible by the increase in data availability, the ease of accessing environmental spatial data, and the development of software solutions aided by computational tools to analyse them. In this article, we review the current use of ML in DSM, identify the key challenges and suggest solutions from the existing literature. There is a growing interest in the use of ML in DSM. Most studies emphasize prediction and accuracy of the predicted maps for applications, such as baseline production of quantitative soil information. Few studies account for existing soil knowledge in the modelling process or quantify the uncertainty of the predicted maps. Further, we discuss the challenges related to the application of ML for soil mapping and suggest solutions from existing studies in the natural sciences. The challenges are: sampling, resampling, accounting for the spatial information, multivariate mapping, uncertainty analysis, validation, integration of pedological knowledge and interpretation of the models. Overall, the current literature shows few attempts in understanding the underlying soil structure or process using the predicted maps and the ML model, for example by generating hypotheses on mechanistic relationships among variables. In this regard, several additional challenging aspects need to be considered, such as the inclusion of pedological knowledge in the ML algorithm or the interpretability of the calibrated ML model. Tackling these challenges is critical for ML to gain credibility and scientific consistency in soil science. We conclude that for future developments, ML could incorporate three core elements: *plausibility*, *interpretability*, and *explainability*, which will trigger soil scientists to couple model prediction with pedological explanation and understanding of the underlying soil processes.

1. Introduction

In recent years, soil science has witnessed a considerable increase in digital soil mapping activities. This is caused by the convergence of several timely factors which are, among others, a huge demand for quantitative and spatial soil information, the accumulation of databases of measured or inferred soil properties coupled with exhaustively known environmental variables, and the development of numerical models combined with computer resources to mine these stores of soil data. The digital soil mapping (DSM) framework was formalized by the publication of McBratney et al. (2003) which builds on Jenny's $S = clorpt$ model (Jenny, 1941) of soil formation, where S is the soil and the acronym *clorpt* stands for climate, organisms, relief, parent material and time, respectively. In short, *clorpt* is a list of variables which, if they are

known without error, are likely to explain the soil over a region. McBratney et al. (2003) supplemented Jenny's formulation with n , which stands for spatial position, and advocated the *scorpan* model for soil mapping. This updated equation provides a spatial model to express quantitatively the relationship between a soil property or class and environmental variables, for a given spatial location.

Conventionally, spatial prediction of soil has been embedded in the geostatistical framework (Heuvelink and Webster, 2001) in which a sample of a soil property is modelled as a sum of a linear combination of environmental covariates and a spatially autocorrelated (stochastic) residual, and prediction at unobserved locations is made by kriging. Geostatistical models are often used in soil mapping because they have several advantages (Oliver, 1987). First, a statistically sound model is assumed for spatial variation. This enables interpretation of the

* Corresponding author.

E-mail address: alexandre.wadoux@sydney.edu.au (A.M.J.-C. Wadoux).<https://doi.org/10.1016/j.earscirev.2020.103359>

Received 23 January 2020; Received in revised form 2 September 2020; Accepted 7 September 2020

Available online 11 September 2020

0012-8252/ © 2020 Elsevier B.V. All rights reserved.

underlying physical processes conveyed by the model. Secondly, spatial autocorrelation is explicitly modelled. This is relevant for environmental variables such as soil which vary from place to place, but exhibit correlation between places. Thirdly, an explicit measure of the uncertainty is associated with the prediction. In many circumstances such as in a decision making process, the prediction is not the only interest and uncertainty maps are required.

Geostatistical mapping of soil has, conversely, several limitations which have only partially been resolved in the current literature. To begin, the residuals are assumed normally distributed, stationary (with constant mean and unit variance) and isotropic. Next, modelling the non-linear relation between a soil property or class and numerous cross-correlated covariates is not straightforward and introduces additional challenges (e.g. many parameters have to be estimated). In heterogeneous areas, the model of spatial variation (i.e. the variogram) also fails to capture both gradual and abrupt changes in soil variation. Finally, geostatistical models are computationally demanding if the sample size and/or the number of prediction locations are large (Cressie and Johannesson, 2008). Often in practice, the soil data and the user-demand make conformity to these assumptions a challenging task.

As an alternative, machine learning (ML) has emerged since the 1990s as a tool for DSM (Lagacherie, 2008). ML techniques refer to a large class of non-linear data-driven algorithms employed primarily for data mining and pattern recognition purposes, and now frequently used for regression and classification tasks in all fields of science. ML algorithms do not make an assumption of the observations' distribution, unlike geostatistical methods where transformation of the original observations is often required to satisfy the assumptions. ML algorithms can also handle a large number of cross-correlated covariates as predictor.

In parallel to this emergence, there has been a tremendous increase in the production and availability of regional and global soil databases. For example, the Soil and Terrain Digital Database (SOTER, Oldeman and Van Engelen (1993)) made by FAO-UNESCO compiled quantitative information on soil and terrain for different parts of the world while WoSIS is a harmonised database of more than 6 million geo-referenced soil records (Batjes et al., 2017). Additionally, numerous spatially exhaustive *scorpan* covariates are available at global scale for climate (Fick and Hijmans, 2017), elevation (Yamazaki et al., 2017), and parent material (Hartmann and Moosdorf, 2012). Further potential covariates are provided by remote sensing such as by the MODIS (Mira et al., 2015) satellite or Sentinel-2A hyperspectral sensor (Gascon et al., 2017). Corollary to the increase of soil databases size, soil mappers face an increasing complexity in modelling the soil with soil data and covariates. Conventional regression techniques, such as linear regression, seem to some extent outdated to accommodate this complexity. This supports the growing use of ML algorithms for DSM.

An essential distinction between (geo-)statistical models and ML algorithms applied in DSM is their purpose: ML algorithms mostly emphasize prediction accuracy whereas statistical models attempt to infer the process which generated the data through a pre-defined model of spatial variation. In the latter case, any interpretation is made in light of the model functions and the value of the covariates or input data. In ML, a predictive model is constructed to predict output values from input values using an error-minimization procedure. No explicit assumptions on the functional form of the relationship between independent and dependent variables are made. Since ML algorithms are not conditioned to follow any statistical assumptions, their predictions often appear more accurate than those made by conventional models, but models lack interpretability, i.e. the model structure is very complex and cannot be readily visualized.

In DSM, as in statistical modelling (Breiman, 2001), the objective is twofold: 1. getting better understanding of mechanistic soil processes, 2. making accurate predictions. The current literature on ML for DSM has focused almost exclusively on the latter. In many situations, the soil

map accuracy is crucial, for example in applications such as policy making or baseline production of quantitative soil information. Following this argument, there has been an increasing number of publications where prediction (viz. mapping) of a soil property or class was the main interest. Many “easy-to-follow” software implementations have supported this increase. When the interest is not only in mapping accuracy, but also in understanding the underlying soil structure and to obtain scientific insights on the soil processes, additional aspects need to be considered, in particular the inclusion of pedological knowledge in the ML algorithm or the need to increase the interpretability of the (often very complex) calibrated ML model structure.

This article aims to review the development of ML applied to DSM by identifying key challenges and opportunities. In this review, we define ML as the computer-assisted practice of using data-driven (and mostly non-linear) statistical models which resort to a large amount of input data to learn a pattern and make a prediction. We start by reviewing and summarizing the current use of ML in DSM. We discuss supervised ML algorithms, since the DSM literature almost exclusively used supervised algorithms for mapping. Based on this summary, we identify gaps in the knowledge and define areas in which adaptation in the use of ML algorithms would be beneficial to increase our understanding of the underlying soil structure and processes. We suggest potential solutions based on the literature from different fields of the natural sciences. Finally, we define three core elements that could trigger soil scientists to couple model prediction with hypothesis generation and explanation of soil processes.

2. A summary of applications

Table 1 summarizes some recent case studies of digital soil maps that have been produced using a ML algorithm. The example studies presented in Table 1 are by no means an exhaustive representation of the available literature but are deemed representative of the diversity of applications and modelling procedures found among DSM with ML studies. The papers were chosen from a search in Web of Science (WoS) and Scopus databases. From this search, a representative sample of more than 150 papers was selected for this review, 75 of which are summarized in Table 1.

2.1. What is mapped?

2.1.1. Quantitative variables

ML algorithms have been successfully applied for quantitative mapping of various soil properties such as soil organic carbon concentration (Henderson et al., 2005; Bui et al., 2009; Kheir et al., 2010b; Dai et al., 2014; Siewert, 2018; Pouladi et al., 2019) and associated stocks (Grimm et al., 2008; Adhikari et al., 2014; Ließ et al., 2016; Wang et al., 2017; McNicol et al., 2019), to map soil texture (viz. clay, silt and sand content) (Ließ et al., 2012; Akpa et al., 2014; Vaysse and Lagacherie, 2015; da Silva Chagas et al., 2016), pH (Dharumarajan et al., 2017), or cation exchange capacity (Forkuor et al., 2017).

ML algorithms have also been applied to make maps of soil nutrients such as nitrogen (Viscarra-Rossel et al., 2015; Forkuor et al., 2017), phosphorus (Viscarra-Rossel et al., 2015; Hengl et al., 2017b; Song et al., 2018), potassium, calcium or magnesium (Hengl et al., 2017b), among others.

A number of studies have also predicted soil attributes and conditions with ML such as bulk density (Viscarra-Rossel et al., 2015) or soil pollutants (Kheir et al., 2010a). Wu et al. (2016) mapped soil background concentrations of arsenic in the Jiangxi Province in China. Taghizadeh-Mehrjardi et al. (2016a) mapped soil salinity in Iran. Tajik et al. (2019) mapped abundance and diversity of soil invertebrates using environmental covariates in a deciduous forest ecosystem in northern Iran while Malone et al. (2009) mapped carbon storage and available water capacity in an area in eastern Australia. Van Den

Table 1
Non-exhaustive list with summary of case studies in which machine learning algorithms are used for digital soil mapping.

Spatial extent ^a	Sample size	Sampling design	Number of covariates	Machine learning model ^b	Covariate selection	Parameter tuning	Map quality indices ^c	Uncertainty quantification	Reference
Quantitative maps									
Plot	285	grid-based	19	cubist, RF	no	no	R ² , RMSE	no	Pouladi et al. (2019)
Local	47	stratified random	41	RF	yes	yes	RMSE, IQR	yes	Blanco et al. (2018)
Local	70	clHS	19	cubist	no	no	MAE, RMSE, R ² , CCC	yes	Lacoste et al. (2014)
Local	75	grid-based	9	ANN	no	no	R ² , MSE	no	Kalambukattu et al. (2018)
Local	98	varied sources	173	RF	yes	no	RMSE, R ²	no	Shi et al. (2018)
Local	116	simple random	20	RF	no	no	R ² , RMSE, CCC	no	Dharumaran et al. (2017)
Local	117	not specified	412	GBM	yes	yes	R ² , RMSE, MAE	yes	Hamezpour et al. (2019)
Local	120	stratified random	not specified	cubist	yes	no	ME, MAE, R ² , R _{adj} ²	no	Miller et al. (2015b)
Local	120	stratified random	22	RF	no	no	ME, RMSE, R ² , MSE	no	Wiesmeier et al. (2011)
Local	137	systematic random	20	ANN, BRT	yes	yes	R ² , RMSE, ME	no	Mosleh et al. (2016)
Local	138	not specified	15	ANN, GEP	yes	yes	RMSE, R ² , MBE	no	Mahmoudabadi et al. (2017)
Local	150	grid-based	not specified	RF	yes	no	RMSE, R ² , CCC	no	Zhu et al. (2019)
Local	151	not specified	not specified	ANN	no	yes	Willmott's index of agreement, R ² , RMSE, correlation coefficient	no	Sergeev et al. (2019)
Local	153	grid-based	26	RF	yes	no	R ² , NRMSE	no	Kovačević et al. (2010)
Local	159/34	not specified	37	RF, cubist, QRF, NN, avNNet, ctree, evtree, GBM, k-NN, RT, SVM	yes	no	RMSE, R ²	no	Tajik et al. (2019)
Local	165	stratified random	18	RF	no	yes	R ² , RMSE, MAE, MARE	yes	Rudiyanto et al. (2018)
Local	173 profiles	clHS	19	RF	no	no	MSE, NMSE	no	Grimm et al. (2008)
Local	188 profiles	clHS	16	ANN, SVR, k-NN, RF, RT	no	yes	ME, RMSE, R ²	no	Taghizadeh-Mehrjardi et al. (2014)
Local	234	not specified	410	cubist	yes	no	RMSE, CCC	no	Taghizadeh-Mehrjardi et al. (2016b)
Local	330 profiles	not specified	12	BRT, ANN, least-square SVM	no	yes	MAE, R ²	yes	Miller et al. (2015a)
Local	330	simple random	10	RF, GBM	no	yes	R ² , R _{adj} ² , RMSE, relative RMSE	no	Ottay et al. (2017)
Local	334	clHS	16	cubist, RF, RT	yes	no	ME, MAE, RMSE, R ²	no	Tziachris et al. (2019)
Local	342/321	-	14	MARS, SVR, RF, Cubist, NN	-	yes	R ² , RMSE	no	Zeraatpisheh et al. (2019)
Local	399	not specified	12	RF	no	no	R ² , RMSE	no	Behrens et al. (2018b)
Local	440	varied sources	19	RF, SVM, ANN	no	yes	RMSE, ME	no	da Silva Chagas et al. (2016)
Local	460	grid-based	21	RF	no	yes	ME, MAE, RMSE	no	Were et al. (2015)
Local	568	simple random	26	QRF	no	no	R ² , RMSE, range-normalized RMSE, Moran's I	yes	Pahlavan-Rad and Akbarimoghaddam (2018)
Local	1104	expert	29	RF, SVM, SGB	no	yes	R ² , RMSE, range-normalized RMSE, bias, RMSE, SS, R ²	no	Kirkwood et al. (2016)
Local	≤ 1052/ 2050/ 2379	varied sources	300-500	BRT, RF	yes	yes	RMSE, sMAPE	no	Forkuor et al. (2017)
Local	2388	varied sources	3	CNN, RF	no	yes	bias, RMSE, SS, R ²	no	Nussbaum et al. (2018)
Regional	not specified	not specified	20	cubist	no	yes	ME, RMSE, R ² , CCC	no	Wadoux et al. (2019b)
Regional	125 profiles	purposive	12	BRT, RF	no	no	R ² , RMSE, bias, CCC	yes	Mulder et al. (2016)
Regional	244	grid-based	4	ANN	no	yes	MAE, RMSE, R ² , CCC	no	Yang et al. (2016)
Regional	339/961	varied sources	40	QRF	no	yes	ME, MAE, RMSE, CCC	no	Dai et al. (2014)
Regional	485 profiles	not specified	5	CNN	no	yes	R ² , RMSE	yes	Nauman and Duniway (2019)
Regional	500	not specified	12	RF, BRT	yes	no	R ² , RMSE	yes	Padarian et al. (2019)
							R ² , RMSE	no	Beguin et al. (2017)

(continued on next page)

Table 1 (continued)

Spatial extent ^a	Sample size	Sampling design	Number of covariates	Machine learning model ^b	Covariate selection	Parameter tuning	Map quality indices ^c	Uncertainty quantification	Reference
Regional	528	subset from a systematic grid	18	k-NN	yes	no	RMSE, R ² , Bias, coefficient of variance	no	Mansuy et al. (2014)
Regional	705	simple random	16	RF, BRT, SVM	yes	yes	R ² , MAE, RMSE, CCC	yes	Wang et al. (2018)
Regional	978 profiles	not specified	24	RF	no	no	R ² , ME, RMSE, CCC	no	Akpa et al. (2014)
Regional	1,014	stratified random	327	CART, BRT, BRT, RF, SVM	yes	no	R ² , RMSD, RPD, RPIQ	no	Keskin et al. (2019)
Regional	1,134	not specified	81	NN	no	no	R ² , ME, MAE, RMSE	no	Aitkenhead and Coull (2016)
Regional	1,300	not specified	6	RF	no	no	CCC, RMSE	yes	McNicol et al. (2019)
Regional	profiles	1,626	40	SVM	no	yes	R ² , MSE	no	Wu et al. (2016)
Regional	2,024	legacy data	16	QRF	no	no	ME, RMSE, R ² , accuracy plot	yes	Vayse and Lagacherie (2017)
Regional	profiles	2,024	16	legacy data	no	yes	MSE, R ²	no	Vayse and Lagacherie (2015)
Regional	2,943	two-stage systematic	37	CNN, RF	no	yes	ME, RMSE, R ² , CCC	yes	Wadoux (2019)
Regional	4,859	not specified	26	QRF	no	no	ME, RMSE, accuracy plot	yes	Szabó et al. (2019)
Regional	4,859	not specified	32	QRF	no	no	ME, RMSE, accuracy plot	yes	Szabó and Pásztor (2019)
Regional	5,386	varied sources	6	cubist, SVM	no	no	R ² , MSE, CCC	yes	Somarathna et al. (2016)
Regional	13,000	not specified	18	RF	no	no	R ²	yes	Koch et al. (2019)
Regional	19,790	two-stage systematic	197	RF	no	no	ME	no	Wadoux et al. (2019a)
Regional	37,693	legacy soil data	74	RF, Cubist, SVM	yes	yes	R ² , RMSE, MAE	yes	Gomes et al. (2019)
Regional - Global	2,268-27,262	varied sources	34	cubist	no	yes	CCC, RMSE, SDE, ME	yes	Viscarra-Rossel et al. (2015)
Regional - Global	366,034	varied sources	> 200	RF, GBM	no	yes	R ² , ME, RMSE, MAE	yes	Ramcharan et al. (2018)
Global	11,268	legacy soil data	118	SVM, kernel weighted NN, RF	yes	no	EC, RMSE, R ²	yes	Guevara et al. (2018)
Global	150,000	legacy soil data	> 200	RF, GBM	no	yes	R ²	no	Hengl et al. (2017a)
Categorical maps	-	not specified	125	ANN	no	no	Accuracy, recall, precision	no	Behrens et al. (2005)
Local	33 profiles	not specified	16	RF, J48	no	no	not specified	no	Massawe et al. (2018)
Local	103/297/ 57	cLHS	130	k-NN, NSC, CT, BCT, RF, linear SVM, radial-basis SVM, NN, ANN	yes	yes	Kappa analysis, Brier scores, visual inspection, confusion index	no	Brungard et al. (2015)
Local	125 profiles	cLHS	17	RF	no	no	map purity, Cohen's kappa, Shannon entropy index, relative purity, relative diversity	no	Zeraatpisheh et al. (2017)
Local	151	not specified	not specified	SVM	no	no	NRMSD, micro averaged F1 measure, kappa statistics	no	Kovačević et al. (2010)
Local	175, 63 profiles	varied sources	27	k-NN, SVM, DT, RF	no	no	OA, PA, UA, kappa coefficient, AUROC	no	Vermeulen and Van Niekerk (2017)
Local	452 profiles	regular grid	6	DT, RF	yes	no	OA, UA, PA, Kappa coefficient of agreement	no	Shariff et al. (2019)
Local	917	grid-based	33	RF	yes	no	Kappa index	no	Houkpatin et al. (2018)
Local	3,121	by-polygon, equal-class, area-weighted, and area-weighted with random over sampling	20	CART, CART with bagging, RF, k-NN, NSC, ANN, LMT, SVM	no	yes	overall agreement, quantity disagreement, allocation disagreement, total disagreement	no	Heung et al. (2016)

(continued on next page)

Table 1 (continued)

Spatial extent ^a	Sample size	Sampling design	Number of covariates	Machine learning model ^b	Covariate selection	Parameter tuning	Map quality indices ^c	Uncertainty quantification	Reference
Regional	89,323	random sampling	26	k-NN, RF	yes	no	recall, accuracy	no	Subburayalu and Slater (2013)
Regional	366,034	varied sources	> 200	RF, GBM	no	yes	OA, regional dataset	yes	Ramcharan et al. (2018)
Regional	7,664	varied sources	110	DT, RF, EGB, SVM, k-NN	yes	no	OA, precision, recall, F-score, K-index	no	Taghizadeh-Mehrjardi et al. (2019b)
Regional	9,924	not specified	23	RF	yes	no	error matrix	no	Haring et al. (2012)
Global	150,000	legacy data	> 200	RF, GBM	no	yes	map purity, weighted kappa metrics, AUC, true positive rate, scaled Shannon's entropy index	no	Hengl et al. (2017a)

^a Plot: 0–1 km²; Local: > 1 km²–10⁴ km²; Regional: > 10⁴ km²–10⁷ km²; Global: > 10⁷ km².

^b RF: random forest; ANN: artificial neural networks; CNN: convolutional neural networks; GBM: gradient boosting machine; BRT: boosted regression tree; GEP: gene expression programming; QRF: quantile regression forest; avNNet: neural networks using model averaging; ctree: conditional inference trees; evtree: evolutionary algorithm for classification and regression tree; NN: neural networks; GBM: generalized boosted regression; k-NN: k-nearest neighbours; RT: regression tree; SVM: support vector machine; MARS: multivariate adaptive regression splines; SGB: stochastic gradient boosting; CART: classification and regression tree; NSC: nearest shrunken centroids; CT: classification tree; BCT: bagged classification tree; DT: decision tree; LMT: logistic model tree; EGB: extreme gradient boosting.

^c R² coefficient of determination; R²adj: adjusted coefficient of determination; RMSE: root mean square error; IQR: interquartile range; MAE: mean absolute error; CCC: Lin's concordance correlation coefficient; MSE: mean square error; ME: mean error; MBE: mean bias error; RPIQ: ratio of performance to interquartile distance; NRMSE: normalized root mean squared deviation; MARE: median absolute relative error; NMSE: normalized mean square error; SMAPE: symmetric mean absolute percentage error; SS: skill score; RMSD: minimum root mean square deviation; RPD: residual prediction deviation; SDE: standard deviation of the error; EC: overall ratio; OA: overall accuracy; UA: user accuracy; AUROC: area under receiver operating characteristic curve; AUC: area under the curve.

Hoogen et al. (2019) mapped the soil nematodes density at the global scale.

2.1.2. Categorical variables

Compared with continuous soil property mapping, fewer studies have applied ML to categorical variables. Digital mapping of soil classes using ML started in the 90s. Probably the first of its kind, Lagacherie and Holmes (1997) predicted soil classes in a regional area while Cialella et al. (1997) predicted soil drainage classes using remote sensing and elevation covariates. Behrens et al. (2005) mapped soil units in a 600 km² area of Western Germany. These studies have recently been completed by a number of publications comparing the maps predicted by a ML model to conventional soil maps (e.g. Zeraatpisheh et al., 2017). Scull et al. (2005); Brungard et al. (2015); Heung et al. (2016); Hounkpatin et al. (2018) employed ML to classify soil taxonomic units. Vermeulen and Van Niekerk (2017) mapped salt-affected areas in irrigation schemes in South Africa. Table 1 provides an additional summary of case studies.

A special case of categorical mapping occurs when the map of soil class already exists but needs to be disaggregated. Disaggregation involves producing new soil maps by increasing the spatial resolution of the original soil map (downscaling), updating the original map using new information or harmonizing the map with other maps. Bui et al. (1999) and Moran and Bui (2002) used a decision tree (DT) to disaggregate an existing map and obtain a realization of the disaggregated soil class distribution. With multiple realizations, the most probable soil class is obtained for a given location. This was further investigated by Hansen et al. (2009) to disaggregate a reconnaissance soil map using a binary DT. A similar approach with DT was used in Häring et al. (2012) to downscale soil types within existing map unit boundaries. More recently, Odgers et al. (2014) used ML to model and disaggregate soil classes and reported the probability associated to each soil class at a given location in the area of interest. A growing number of publications exploits this approach (e.g. Holmes et al., 2014; Vincent et al., 2018; Ellili et al., 2019).

2.2. Extent, resolution, depths

There is a large range of case studies mapping soil properties or classes from the plot (< 1 km²) to the global (> 10⁷ km²) scale. Most studies in our literature review predict at a local to regional scale. The mean extent of the study area was 3,900 km², but most studies (90%) considered a study area smaller than 650,000 km² (equivalent to the size of metropolitan France). Few studies mapped at plot or global scales. For example, Pouladi et al. (2019) made a quantitative map over a 10 ha (0.1 km²) field in Denmark while Hengl et al. (2017a) produced quantitative and categorical maps for the whole world.

We found a clear correlation between the spatial extent of the study area and the grid spacing (i.e. the spacing between point predictions) at which the soil property or class is mapped: the larger the study area, the coarser the resolution. The resolution spans between 2 m x 2 m (Lacoste et al., 2014) to 1 km x 1 km for large, regional or continental study areas (e.g. Hengl et al., 2014). Most studies, however, mapped at a standard spatial resolution of 30, 90 or 250 m.

While most of the studies (70%) predicted a soil property or class for a single depth (topsoil), a number of studies accounted for the soil variation at multiple depths. Viscarra-Rossel et al. (2015) followed the GlobalSoilMap project specifications (Arrouays et al., 2014) to produce a quantitative three-dimensional map of several soil properties for six depths intervals, namely 0–0.05 m, 0.05–0.15 m, 0.15–0.30 m, 0.30–0.60 m, 0.60–1.00 m and 1.00–2.00 m. Similar depth intervals were used in Mulder et al. (2016) and Adhikari et al. (2014) for soil organic prediction in France or Denmark, respectively. Several other studies (e.g. Grimm et al., 2008; Lacoste et al., 2014) used standard depth intervals for prediction, based on national mapping requirements or suitable for their specific case study.

2.3. Sampling design and density, sample size

The sampling design is the position in the two-dimensional geographic space of the units used to calibrate or validate the ML algorithm. Most studies do not specify the sampling design of their study. It is speculated that the sample originates from multiple sources, such as legacy data, expert-based designs, and combination of several surveys, each of which had a different sampling design. When specified, non-probability sampling such as grid-based sampling designs are by far the most used for calibration (e.g. by Pahlavan-Rad and Akbarimoghaddam, 2018; Sergeev et al., 2019; Shariffar et al., 2019). Another non-probability sampling design is conditioned Latin Hypercube (cLHS), used to collect a sample in Lacoste et al. (2014) and Brungard et al. (2015). Probability sampling is used in about one fourth of the studies. For example, simple random sampling was used in Tziachris et al. (2019), while a sample was collected based on stratified random sampling using land-use and topography as stratifying variables in Wiesmeier et al. (2011).

In our literature review, we found that the sample size varies considerably between studies. While the average sample is composed of 1,000 units, about one third of the studies use a sample with less than 150 units, mostly for local or small-scale regional areas. For example, Blanco et al. (2018) used a sample of size 47 for mapping soil water retention in a 93 km² area while Massawe et al. (2018) observed 33 soil profiles to calibrate a ML algorithm and to predict soil taxa over a 11,600 km² area. As expected, global studies have very large sample sizes. Hengl et al. (2017a) and Ramcharan et al. (2018) used a sample composed of more than 150,000 units to make soil property or class maps of the whole world, and of the United States, respectively.

When the sample size is associated to the extent of the study area, our review shows that large-scale studies have a very coarse sampling density. While the average sampling density in our literature survey is 0.24 units/km², studies by Beguin et al. (2017) and Wang et al. (2017) had both a sampling density smaller than 3 units/10,000 km² for mapping soil properties in the rangelands of eastern Australia or in the Canadian boreal forests. Small-scale studies have, conversely, high sampling density. All studies with area size less than 50 km² have a sampling density larger than 7 units/km².

2.4. Covariates

Environmental covariates are used as predictors in ML algorithms. They are supposed to explain part of the physical and chemical processes governing soil spatial variation. Most studies used about 20 covariates. Only a few used less than five (e.g. Dai et al., 2014; Padarian et al., 2019) while others used more than 100 (e.g. Hengl et al., 2017a; Ramcharan et al., 2018). Since the covariates represent soil forming factors, numerous studies (e.g. Viscarra-Rossel and Chen, 2011; Wang et al., 2018; Gomes et al., 2019; Szatmári and Pásztor, 2019) logically selected the covariates to represent the key factors of the *scorpan* model of soil spatial variation. The most common ones are existing soil property or class maps, (long-term) average annual precipitation and temperature, remote sensing images (e.g. SPOT satellite images or vegetation indices derived from satellite images), elevation, terrain attributes (e.g. slope, local curvature, topographic wetness index) and existing geological maps.

Covariates representing *scorpan* factor of soil variation might not be available or easily obtainable in all case studies. In some cases, covariates are selected based on expert knowledge. A number of studies therefore calibrated ML algorithms using sets of climatic variables, remote sensing images or terrain attributes only, or a combination of them. For example, Mansuy et al. (2014) used a set of eight climatic and eight terrain attribute variables to map soil carbon, nitrogen and texture in a large area in Canada. Shariffar et al. (2019) used six terrain attributes as covariates. These were chosen from a large set of environmental covariates using knowledge on the expected relationship

between the covariate and the soil property to be mapped. We note that a few studies (e.g. Hengl et al., 2018; Miller et al., 2015a) have considered a very large (> 100) number of covariates for calibration. This large amount of covariates relies mostly on remote sensing images, such as MODIS land products (long-term averages, several near- or mid-infrared bands) or Landsat products (near-, short-wave near-infrared, or γ radiometric bands, bare ground images). Using temporal covariates in a ML algorithm is also a way to map the temporal dynamic of soil properties. For examples Heuvelink et al. (2020) used time series of MODIS products to map the temporal dynamic of organic carbon in Argentina between 1982 and 2007.

A few studies account for the multi-scale variation of the environmental covariates. In other words, terrain derivatives may well be aggregated to account for physical processes in soil that are not visible on a finer scale. Examples of studies using multi-scale covariates for mapping with ML algorithms are Behrens et al. (2010), Miller et al. (2015b) and more recently Behrens et al. (2018a). Miller et al. (2015b), for example, used a total of 412 covariates, several of which are derived from the aggregation of terrain attributes from a fine (i.e. a grid cell size of 2 m x 2 m) elevation map.

A growing number of studies have advocated the use of spatial surrogate covariates as an indicator of spatial position in the *scorpan* model of soil variation. The most common surrogate is the use of geographical coordinates (easting and northing) as covariates in the model. Maps of distances from observation locations, or group of locations, have been used by Hengl et al. (2018). They are categorized into Euclidean, downslopes or “resistance” distances. More recently, Behrens et al. (2018b) used Euclidean distance fields, which are maps of distance from reference locations in the study area such as a corner or a center.

2.5. Covariate selection

Covariate (*aka* feature) selection aims at reducing the number of covariates used to calibrate the ML models. There are several reasons for selecting a subset of covariates to calibrate the model. Some of them are: (i) to calibrate the ML model faster, (ii) to reduce complexity, (ii) to increase the prediction accuracy, (iv) to avoid multicollinearity, or (v) to prevent over-fitting of the ML model, i.e. to prevent poor prediction accuracy on unseen data. In our literature review, about one third of the studies applied covariate selection. Two main categories of covariate selection techniques are found. The first applies the covariate selection as a pre-processing step, i.e. before calibrating the ML model. This is the case in Zhu et al. (2019); Hamzehpour et al. (2019) and Zeraatpisheh et al. (2019). Hamzehpour et al. (2019) selected the covariates to be used in calibration by computing the Pearson's *r* correlation coefficient between the covariates, and by discarding the ones that were highly correlated, while Mosleh et al. (2016) selected the covariates based on the Pearson's *r* correlation coefficient between the soil property values and the covariates, and selected a subset of covariates which are strongly correlated with the property. The second type of covariate selection is called “wrapper” methods and relies on the inference made by a calibrated ML model to determine whether covariates are important. By re-calibrating a ML model several times, each time removing the least important covariate, one may expect to reduce considerably the overall number of covariates with little or no decrease in model prediction accuracy. Examples on the use of “wrapper” methods are found in Taghizadeh-mehrjardi et al. (2016); Shi et al. (2018); Rudiyanto and Minasny (2018); Tajik et al. (2019) or Gomes et al. (2019). The most used of “wrapper” methods is an optimization algorithm called recursive feature elimination.

2.6. Machine learning algorithms

A large number of ML algorithms and their variants have been used in the DSM literature. For quantitative mapping, tree-based algorithms

are the most popular ones, the simplest version of which is the classification and regression trees (CART), used for example by Taghizadeh-Mehrjardi et al. (2016b). CART is known to be sensitive to the quality and size of the calibration sample. To solve this problem, the bagging (bootstrap and aggregating) procedure (Breiman, 2017) has been introduced in random forest (RF). Our literature review shows that RF is currently the most popular ML algorithm for regression purposes. Example of case studies using RF for mapping are Tziachris et al. (2019); Vaysse and Lagacherie (2015); Forkuor et al. (2017); Dharumarajan et al. (2017); Liu et al. (2019). More recently, Vaysse and Lagacherie (2017) used a variant of random forest (RF), called quantile regression forest, as a method to map the uncertainty associated with the prediction of the soil property. Another tree-based method is cubist, employed in about 10% of the reviewed literature (e.g. by Mulder et al., 2016; Viscarra-Rossel et al., 2015; Miller et al., 2015a). A few studies (less than five) used boosted regression tree (Yang et al., 2016; Beguin et al., 2017). In addition, a number of studies used artificial (Aitkenhead and Coull, 2016; Guevara et al., 2018; Lamichhane et al., 2019) or convolutional neural networks (CNN, Wadoux, 2019). A relatively small number of studies used alternative algorithms such as support vector machines (Guevara et al., 2018), *k*-nearest neighbours (Mansuy et al., 2014) or generalized boosted regression (Tziachris et al., 2019; Gomes et al., 2019).

For classification purposes, tree-based algorithms are also the most popular ones. About 80% of the case studies used at least one tree-based algorithm such as regression tree (e.g. Taghizadeh-Mehrjardi et al., 2019b; Heung et al., 2016), RF (e.g. Häring et al., 2012) or boosted regression tree (e.g. Lorenzetti et al., 2015). Alternatively, gradient boosting was used by Hengl et al. (2017a), *k*-nearest neighbours by Vermeulen and Van Niekerk (2017) and compared to support vector machines. The latter algorithm was also used in Taghizadeh-Mehrjardi et al. (2019b). Artificial neural networks (ANN) is also popular and was used in Behrens et al. (2005) and Heung et al. (2016).

Recent studies have proposed to use model ensemble modelling, in which predictions from several individual models are aggregated. Taghizadeh-Mehrjardi et al. (2019a) combined seven ML model predictions for soil class mapping in a case study in Iran while Song et al. (2020) implemented a weighted ensemble learning model to map soil organic carbon in consideration of pedoclimatic zones in China. Ensembles were also considered in Hengl et al. (2017a) for global soil mapping.

2.7. Parameter tuning

The performance of a ML model is affected by the values of its parameters. Almost half of the studies perform a search to find optimal values. Padarian et al. (2019) manually selected the number of ANN neurons in each layer of the model. This search is automated by a so-called grid-search process. This is by far the most used technique for parameter tuning. In a grid-search process, a number of parameter values are evaluated based on the model prediction error. The process is computationally intensive (i.e. the ML model must be calibrated for each parameter set proposal). Examples of studies using a grid-search to find ML parameters are Ottoy et al. (2017), Taghizadeh-Mehrjardi et al. (2016a), Pahlavan-Rad and Akbarimoghaddam (2018), Sergeev et al. (2019), Forkuor et al. (2017) and Ramcharan et al. (2018). An alternative to the grid search is to apply an optimization algorithm, such as the particle swarm method, to find optimal parameter values. For example, Wu et al. (2016) compared two genetic algorithms and a grid search process to find the ML parameters. Recently, Wadoux et al. (2019b) used Bayesian optimization to optimize the number of layers, the neuron number, the learning rate and the batch size of an CNN algorithm for mapping soil organic carbon.

2.8. Validation and uncertainty quantification

In our literature review, all studies computed at least one map quality index. A list of map quality indices is provided in Table 1. About 30% of the studies obtained the map quality index through cross-validation, while 30% through data-splitting. The remaining studies either repeated data-splitting several times, validated through visual examination or used a grid-based sampling design. Only two studies provided model-free estimates of the map quality index with the collection of an additional probability sample for validation (Subburayalu and Slater, 2013; Lacoste et al., 2014).

In addition to the map quality indices, about 30% of the studies quantified the uncertainty associated with the prediction. These studies reported confidence interval, obtained by bootstrapping the original set of observations (e.g. Chen et al., 2019; Padarian et al., 2019; Hamzehpour et al., 2019). A few studies used the kriging variance computed on the residuals of a trend obtained by predicting with a ML algorithm (e.g. Koch et al., 2019), or a combination of bootstrap and kriging variance (e.g. Viscarra-Rossel et al., 2015). In three studies, prediction intervals were obtained through the quantile regression forest. Wadoux (2019) obtained the prediction intervals following a two-step procedure called mean plus variance estimate (MVE) for mapping several soil properties using CNN.

3. Challenges and opportunities

Based on the review, we identify some knowledge gaps and challenges in the current use of ML algorithms for DSM. We will also outline some opportunities for research.

3.1. Sampling

Despite abundant evidence that the sampling design and sample size play a key role in the resulting map accuracy (De Gruijter et al., 2006), sampling designs suitable for mapping with ML are yet to be discovered. The effect of the sample size for mapping with ML was discussed in Somarathna et al. (2017) where the efficiency of several ML algorithms were compared for the spatial prediction of soil carbon. The study showed that having a sufficiently large sample size is more important than choosing a sophisticated ML algorithm, and that when the sample size is small, it is best to use simple models. About sampling design, Brus (2019) speculated that ML algorithms would benefit from a spread of the sampling units in the feature (covariate) space, and suggested the use of feature space coverage sampling (FSCS) using *k*-means clustering or conditioned Latin Hypercube sampling (cLHS). Both sampling designs aim at covering the space spanned by the covariates, but in different ways. Experimental results are provided by Wadoux et al. (2019a) in a study comparing five sampling designs (viz. simple random sampling, cLHS, spatial coverage sampling (SCS), FSCS and a design optimized in terms of mean square error) for soil property mapping with RF. The results show large differences in mapping accuracy between the designs, and that a FSCS design optimized in the most important covariates of the RF model had the closest match to an optimized design. By performing further diagnostics, the study concludes that RF does not benefit from an uniform spread of the units in the geographic/feature space, nor from reproducing the marginal distribution of the covariates (as it is done in cLHS). These results were confirmed by Ma et al. (2020) for soil classes mapping with RF, and similar conclusions were drawn by Lagacherie et al. (2020) in a numerical experiment mapping a continuous soil property with quantile regression forest. These results apply for RF but there is a need to further investigate sampling designs for other ML algorithms. While most studies in our literature review (Table 1) used a grid-based sampling or

cLHS, there are now clear indications that most conventional sampling designs (e.g. spatial coverage sampling) are not effective for the purpose of mapping with ML.

To discover what makes a good design for mapping with ML, one should ideally derive optimal designs. More importantly, one should investigate the characteristics of these designs, so that future research can generate simple designs that resemble the optimal ones (Wadoux, 2019a). It is likely that optimal designs differ between ML algorithms. We speculate that a somewhat uniform spread in the feature (i.e. covariate) space remains important for all ML algorithms since they all link the covariates and the sample values in a non-linear way, but that additional considerations might outweigh or overtake this uniform spread. An example of optimal design is given by the studies of Pozdnoukhov and Kanevski (2006) and Tuia et al. (2013) where the sampling configurations are optimized with active learning for mapping with support vector machines. In the first study, the selected sampling units were the most beneficial for the algorithm, avoiding misclassification between temperature below or above 20°C (categorical mapping) by becoming support vectors. In Tuia et al. (2013), a similar methodology was adopted and tested in three case studies to subsample an existing sample for quantitative mapping, to add new sampling units optimally in a continuous map or to define suitable areas for sampling. In all case studies, the authors obtained a design optimal for the purpose of mapping with support vectors machine. They conclude that while a sampling design can be representative of the geographical space, the sampling design can be judged unrepresentative if other dimensions are considered. These results encourage the use of new methods for sampling design optimization such as active learning. Active learning is a model-based sequential re-design algorithm. In active learning, the objective function (e.g. the spatially averaged prediction uncertainty) is explicitly quantified and used to define the additional sampling units that are the most beneficial for the model (e.g. the boundary between two classes). In this sense, active learning is similar to optimization with spatial simulated annealing (Van Groenigen, 1997) routinely used in geostatistical sampling design optimization. Besides the optimization algorithms, a set of objective functions needs to be tested. MacKay (1992) defined an objective function that searches for the optimal units in the space spanned by the predictors (i.e. covariates) for prediction using a ANN algorithm. Taking these considerations and testing active learning for sampling design optimization would certainly make a valuable contribution to DSM research.

3.2. Resampling

Regional and global scale studies almost invariably use legacy soil data (Stumpf et al., 2016). Legacy soil samples provide valuable information on soil classes and properties but are often highly clustered in areas of specific interest. In modelling with ML, it is assumed that the sample is composed of independent and identically distributed sampling units whereas soil observations within an area typically exhibit spatial autocorrelation (i.e. close observations are more similar than remote ones) This has important implications in terms of sampling, resampling of the observations and validation of the models. A ML algorithm calibrated with a spatially clustered sample may lead to biased predictions over the area because of the over-representation in the calibration process of regions of high sampling density. Despite being critical, this has yet been disregarded in DSM studies. In geostatistics, spatial declustering has been applied to reduce the effect of clustered data in the calculation of experimental variogram (Marchant et al., 2013). One form, called cell declustering involves overlaying a grid over the area and assigning a weight to the sampling units based on the inverse of the number of units in the cell. In ecology, a first attempt was made by Bel et al. (2005) and later Bel et al. (2009) to decluster the sampling units used in the calibration of a CART model. In Bel et al. (2005), the sampling units are given weights which are obtained from kriging the spatial mean. Bel et al. (2009) elaborated a more complex

procedure in which all quantities involved in the CART algorithm (e.g. the proportion of leaves) have a spatial estimate. This has been further considered by Stojanova et al. (2013) for both categorical and quantitative mapping of ecological variables. Illés et al. (2019) applied polygonal declustering technique to spatially clustered samples by assigning weights on the units based on Voronoi's area proportion.

We point out that the clustering may also occur in the feature (i.e. covariates) space and speculate that this may also affect the prediction if most units are clustered at some specific areas of the feature space. For example, a model trained to predict organic carbon in a mountainous area will exhibit biased prediction if most sampling units originate from valley, and that elevation is used as covariates. Similar to Bel et al. (2005), weights can be assigned to the units to down-weight the importance of over-sampled areas in the feature space. An example method is provided by Carré et al. (2007). The authors assumed that a good sample has an uniform spread in the feature space and thus covers all strata of a hypercube based on the covariates. A weight is assigned to each unit in the sample based on the density of the units in each stratum. The larger the density within the stratum, the smaller the weight assigned to a single unit.

The nature of the legacy soil data in categorical mapping also poses additional challenges. ML algorithms for categorical mapping rely on balanced sets of units. In other words, all classes shall comprise a comparable number of sampling units. Legacy soil samples are considered imbalanced in that all classes are not represented equally. Most ML algorithms are calibrated by maximizing the average (classification) accuracy on an independent validation sample. This often results in very low predictive accuracy for under-sampled classes and in models biased towards the over-sampled classes (He and Garcia, 2009). In the ML literature, several approaches have been developed to handle class imbalanced samples. At the higher level, one may distinguish between cost function and resampling based approaches. In the first approach, the model is penalized for misclassification of under-represented classes. This stems from the calibration of ML algorithms, which minimize a loss function to find optimal parameter values (e.g. in ANN). In the second approach, resampling of the sample is performed by either adding units in the under-sampled class, removing units from the over-sampled class, or a mix of the two. The second approach has been recently applied in soil mapping studies, in particular by Heung et al. (2016) and Shariffar et al. (2019). Taghizadeh-Mehrjardi et al. (2019b) tested eight resampling approaches and their effect on the prediction accuracy of five ML algorithms in two large-scale case studies. However, to date resampling techniques are applied the same way as in other disciplines while soil data often present spatial autocorrelation which may affect the resampling strategies. This has not yet been investigated in the literature.

3.3. Accounting for spatial information

ML algorithms do not account for spatial autocorrelation contained in the raw soil data, unless explicitly specified. Sinha et al. (2019) have tested RF for different scenarios of spatial autocorrelation in the observations and confirmed that the presence of spatial autocorrelation leads to high variance of the residuals. ML algorithms accounting for autocorrelated observations have recently been formulated, such as geographical RF (Georganos et al., 2019), or spatial ensemble techniques (Jiang et al., 2017). The two methods boil down to geographically weighted regression by fitting spatially local sub-models using only neighbouring observations. Jiang et al. (2017) decomposed the area into geographic disjoint sub-areas, and fitted a local model in each sub-area. Georganos et al. (2019) fitted a sub-model to each observation using RF, accounting for both non-stationarity and spatial autocorrelation. Finally, random forest spatial interpolation, developed by Sekulić et al. (2020), is another mapping procedure with RF in which the local context is included by adding the observations at the nearest locations of the prediction location as covariates.

Applying a non-spatial model for DSM is not a problem in itself. This is corroborated by the definition of DSM given in [Lagacherie and McBratney \(2006\)](#), which gives provision for mapping using “non-spatial soil inference systems”. In theory, if one includes all relevant environmental variables to model the soil property or class, there should be no spatial autocorrelation in the residuals of the fitted models. If there is spatial autocorrelation in the residuals, some important covariates are likely to be missing. More importantly, this also means that predictions made by the ML algorithm might be biased or the model underfitted because this is a violation of the assumption of independence between data points that is implicitly assumed ([Bel et al., 2005](#)). [Kühn and Dormann \(2012\)](#) recommend mapping the spatial distribution of the residual autocorrelation to facilitate the identification of a missing spatial process. In some cases, a map of residuals exhibits a clear pattern (e.g. increasing residuals with distance from the river) and might help to generate a new hypothesis or to refine the existing model.

Despite the availability of datasets and care taken during modelling, residual autocorrelation is still likely to occur. Several authors have advocated the use of spatial surrogate covariates (i.e. pseudo-covariates) as an indicator of spatial position in the *scorpan* model of soil variation or to account for spatial autocorrelation contained in the data. The most common surrogate is the use of geographical coordinates (easting and northing) as covariate in the model. This has led to maps with visible artefacts, in particular when used in combination with tree-based algorithms. Alternatively, maps of distances from observation locations, or a group of locations, have been proposed by [Hengl et al. \(2018\)](#). They are categorized into Euclidean, downslopes or “resistance” distances. Maps of distance to observation locations generally have no direct meaning in terms of soil process over an area (e.g. distance from the river). [Behrens et al. \(2018b\)](#) proposed to use Euclidean distance fields, which are maps of distance from reference locations in the study area such as the corner or the centre. The studies using distance maps as covariates have shown for several case studies an important reduction of the residual autocorrelation, when compared to a model without distance maps in the set of covariates.

In the context of DSM, we hypothesize that the current use of distance maps is not entirely satisfactory for several reasons. While the goal of modelling is to account for the data variability, and thus to reduce the residuals to zero, in practice no model is entirely confirmed by the observations and there is always a discrepancy between observed and predicted values. Analysis of the residuals may provide valuable insights on the model limitations. Including pseudo-covariates with the set of pedologically relevant covariates, however, can be harmful because it precludes analysis of the residuals and the generation of new hypotheses from these residuals ([Hawkins, 2012](#)). It also hampers the interpretation of the most important covariates ([Meyer et al., 2019](#)), which is key in several studies on soil mapping. Finally, pseudo-covariates of distance may well integrate over several pedologically relevant covariates, making them better covariates or masking the effect of pedologically relevant covariates. In spatial ecology, alternatives to distance maps are found in the use of spatial eigenvector maps, spatial filters or trend-surface regression computed on, or optimized for, the residuals of a model calibrated using ecologically relevant covariates ([Kühn et al., 2009](#)).

3.4. Multivariate mapping

Several authors (e.g. [Hengl et al., 2018](#); [Wadoux, 2019](#); [Wadoux et al., 2019b](#); [Padarian et al., 2019](#)) have shown that it is possible to calibrate a single ML model to predict either multiple soil properties or a single soil property at multiple depths. This reduces the risk of overfitting, computational resources that would be otherwise required to calibrate several disjoint models ([Wadoux, 2019](#)), and increases prediction accuracy if there is correlation between the variables to predict. [Padarian et al. \(2019\)](#) use a multivariate CNN model to predict

SOC at multiple soil depths and report a significant increase of prediction accuracy for the deeper soil depths, compared to predictions made for each depth separately by a cubist model. [Wadoux \(2019\)](#) have shown that for a CNN model, it was feasible to constrain the prediction to avoid inconsistent prediction between compositional soil properties, in particular soil texture. It was done by adding an additional layer to the model, but we speculate that this could also be realized by modifying the objective function used to calibrate the model. Despite a few recent studies, there has been little interest in multivariate soil mapping using ML algorithms. In the ML literature, it appears that almost all conventional ML algorithms have a multivariate counterpart. Multivariate NNs have already been tested in soil mapping studies. An adaptation of the RF algorithm for multivariate mapping is proposed by [Hengl et al. \(2018\)](#) but has several limitations. For example, the calibrated model size increases dramatically when the number of soil properties to predict also increases and it does not allow to separate the contribution of the covariates to each predicted property separately. A theoretical framework for multivariate RF is described by [Segal and Xiao \(2011\)](#) and was further implemented in the R language by [Rahman et al. \(2017\)](#). For support vector machines, a multivariate extension is described in [Xu et al. \(2013\)](#).

One objective when mapping soil properties or classes is to learn from the calibrated model. A calibrated multivariate model can provide insights on the soil property and horizon interrelations. Unfortunately, in a multivariate ML model, the correlation between soil properties or depths is not modelled explicitly (e.g. using a cross-covariance matrix between soil properties). As a result, the correlation between properties or depths cannot be assessed internally and no pedological interpretation can be derived from the calibrated model. More research is needed to discover whether the correlation between original and predicted soil properties (or depths) is preserved in a multivariate ML model. To model the correlation between properties explicitly, two solutions exist. The first is to calibrate additional stochastic parameters together with the ML parameters (e.g. in a neural network algorithm). This can take the form of an auto-regressive model between the predictions ([Uria et al., 2016](#)). Another straightforward solution is to calibrate the model with a criterion related to the absolute difference in correlation between the measured properties and predicted properties. While this is easy to implement in ML calibration based on an objective function (e.g. neural network), this is not straightforward for models such as RF. Overall, including correlation between properties or depths when predicting with a ML algorithm requires further investigation so as to build pedologically realistic and interpretable models.

3.5. Uncertainty analysis

Uncertainty analysis in DSM is crucial to decide whether the predicted soil map is reliable to be used for agricultural production systems or decision making. Uncertainty analysis is also about acknowledging the limits of the models and is therefore one step towards model interpretability. At the higher level, the ML literature distinguishes two sources of uncertainty: aleatoric and epistemic uncertainties ([Fig. 1](#)). Aleatoric uncertainty is the data noise variance (in other terms, the data error), and arises from noise in the data and measurement error. Epistemic uncertainty refers to model and model parameter uncertainty and represents our ignorance about a true model that generated the data. While epistemic uncertainty is easy to reduce (e.g. by collecting more data at areas of low sampling density), aleatoric uncertainty is rather difficult to assess (one must repeat the measurement several times) and to reduce. Methods to quantify epistemic uncertainty are bootstrapping, or Bayesian modelling. Quantifying epistemic uncertainty enables to obtain confidence intervals of the prediction. Aleatoric uncertainty is mainly quantified by quantile regression methods, but Monte-Carlo simulation from the probability distribution of the observations might also be a possible approach. The quantification of both aleatoric and epistemic uncertainty provides prediction

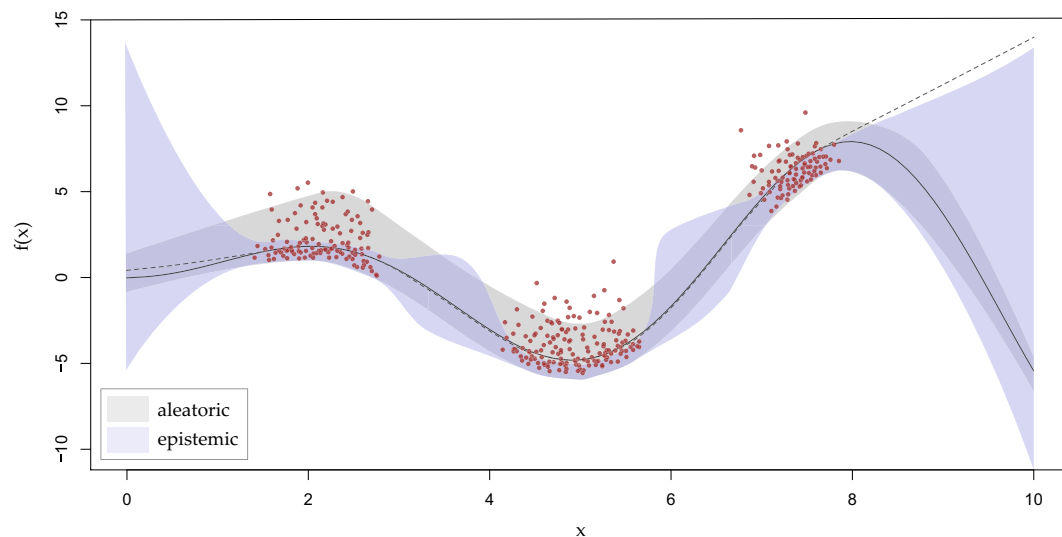


Fig. 1. Transect with location of the sampling units in red, the true (solid line) and predicted (dash line) value of the variable of interest, the aleatoric uncertainty (grey shade) and epistemic uncertainty (blue shape). When no observations are present, the epistemic uncertainty increases. The aleatoric uncertainty remains somewhat constant across the transect. Adapted from [Tagasovska and Lopez-Paz \(2019\)](#). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

intervals with methods such as quantile regression forest (QRF), the Delta or Bayesian methods and the MVE for ANN algorithms.

The recent development of conditional generative adversarial networks (cGAN) ([Mirza and Osindero, 2014](#)) to generate possible realizations of the observations with specific conditions or characteristics seem to be of particular interest to account for measurement error in DSM. Accounting for measurement error is considered by [Wadoux et al. \(2019b\)](#) for mapping soil organic carbon using uncertain measurement of the soil property. However, the authors do not propose a method to quantify the uncertainty of the measurements, nor propagate the measurement error to the predicted map. With cGAN, a probability distribution of the observations is built, which might be used for Monte Carlo simulations. Each Monte Carlo sample is used as input in the ML algorithm, and the final map is the integration of all these simulations. This would effectively tackle the aleatoric uncertainty of the ML model. More importantly, this would also quantify the uncertainty in the measurements, which is currently one of the most important challenges in DSM.

Most studies to date do not provide estimates of the uncertainty ([Table 1](#)). Successful attempts have been made by [Vaysse and Lagacherie \(2017\)](#) and [Wadoux \(2019\)](#) to report prediction intervals for random forest and neural networks models, respectively. Confidence intervals were reported in several studies (e.g. [Hamzehpour et al., 2019](#); [Gomes et al., 2019](#)) and were obtained by training multiple disjoint models using bootstrapped samples of the original data. In a few studies, the variance obtained by bootstrapping is averaged with the variance obtained by kriging of the residuals ([Viscarra-Rossel et al., 2015](#)). From [Fig. 1](#) it follows that if sampling units are selected from a small area in the feature or geographic space, then there will be little uncertainty in this area. Likewise the uncertainty dramatically increases when areas of the feature space are under-sampled, or even worse, ignored. When sampling units are clustered, (spatial) cross-validation might not be sufficient to define realistic prediction accuracy measures because the sampling units used for validation are taken from similar regions of the feature space while the model is biased towards these same regions ([Gahegan, 2000](#)). While the (spatial) cross-validation results might show strong agreement between predicted and measured soil property or class and therefore validate a ML model with very high predictive abilities, an uncertainty quantification would show unrealistic predictions characterized by a large uncertainty (see right-

hand side of [Fig. 1](#)). This results from ML algorithms being very poor for extrapolating to areas of the covariate space that are not comprised in the calibration sample. Uncertainty quantification that separates out data and model uncertainties is thus recommended to complete the evaluation of the predicted maps.

3.6. Validation

Studies by [Roberts et al. \(2017\)](#) and [Ruß and Brenning \(2010\)](#) found that the estimated performance of the ML algorithms applied to spatial data depends on the validation strategy. In DSM, model performance is usually assessed using random k -fold cross-validation (CV) or single random split of a sample into calibration and validation and/or test subsamples. These strategies can give over-optimistic validation statistics estimates in case of clustered data because of the presence of autocorrelation in the observations ([Micheletti et al., 2014](#); [Gasch et al., 2015](#); [Meyer et al., 2018](#)). Validation statistics estimated from a random split of the master sample assess the ability of the model to reproduce the calibration sample but fail to assess the model performance in terms of spatial mapping ([Meyer et al., 2019](#)). As an alternative, several methods ([Brenning, 2012](#); [Le Rest et al., 2014](#); [Pohjankukka et al., 2017](#); [Meyer et al., 2019](#)) for spatial CV are proposed to account for spatial autocorrelation of the observations. Two main strategies are adopted. [Roberts et al. \(2017\)](#); [Brenning \(2012\)](#); [Meyer et al. \(2019\)](#) used a spatial block approach for k -fold CV where the master sample is divided into k spatially disjoint subsamples using clustering algorithms on the coordinates or by dividing the spatial domain based on k cells. In [Le Rest et al. \(2014\)](#) and [Pohjankukka et al. \(2017\)](#), observations from the calibration subsample that are within a given geographic distance of the validation subsample were omitted from the calibration subsample, after which the model was fitted using the remaining observations from the calibration subsample. While these two approaches account for spatial autocorrelation of the observations during validation, further research is required to provide guidelines to select the realistic distance from which a validation data point is statistically independent from the calibration sample so as to avoid the opposite effect, i.e. assessment of extrapolation and subsequent underoptimistic validation statistics estimates. The major limitation of the spatial CV techniques is that they preclude a calibration point to be close to a validation point, while in practice the prediction is made at all locations in an area, including

points that are close and far-away from the calibration data. In this sense, validation statistics estimates using spatial CV may well be underoptimistic.

Research on spatial CV has drawn attention to the role of autocorrelation on the calibration on the ML algorithms. [Schratz et al. \(2019\)](#) show that hyperparameter tuning is also affected by spatial autocorrelation, and that overoptimistic results are reported when the same data are used for performance assessment and parameter tuning. They proposed a nested (block) CV approach for hyperparameter tuning ([Schratz et al., 2019](#)) where spatial blocks are split a second time into spatially disjoint geographic subsamples used to optimize the hyperparameters. The major disadvantage of this method is the dramatic increase in computing time, which can be solved by distributed (parallel) computing solutions. Similarly to the hyperparameter tuning using nested spatial CV, [Meyer et al. \(2018\)](#) showed that autocorrelated covariates lead to overfitting and visible artefacts in the predicted map. The study proposes an iterative procedure for variable selection where a group of two variables is first selected based on the error computed with spatial CV, and new variables are iteratively added only if these increase the model performance. The study of [Meyer et al. \(2018\)](#) gave another argument against the use of covariates describing the spatial dependency as these lead to misinterpretation of the model's important contributors and impossibility for the model to generalize.

[Meyer et al. \(2019\)](#) emphasized the value of visual examination of the predicted maps in addition to the statistical validation. In [Meyer et al. \(2019\)](#), two maps with similar map validation accuracy statistics have a different spatial pattern. This phenomena is due to different sets of selected covariates, some having strong spatial autocorrelation, leading to visible artefacts in the predicted map. This highlights the need for research on the evaluation of predicted maps in terms of spatial pattern. [Poggio et al. \(2019\)](#) compared the spatial structure of predicted versus observed values by computing the area under the curve of variograms fitted on the validation locations for both predicted and observed probabilities of having a peat soil. This relies, however, on the assumption that the variogram of the validation locations is representative of the mapped area. More research in this direction will be valuable for future DSM studies. To date, visual assessment of the map to detect artefacts, and in consideration of our knowledge of soil forming processes, is the best option.

3.7. Machine learning and pedological knowledge

Accounting for existing expert soil knowledge in DSM with ML is a challenging exercise ([Ma et al., 2019](#)). ML algorithms do not build on any existing a priori conceptual model of the soil processes and only processes that are conveyed by the input data are represented in the map ([Coveney et al., 2016](#); [Koch et al., 2019](#)). To prevent extrapolation, [Hengl et al. \(2014\)](#) did not provide soil maps in some under-sampled areas of the globe such as deserts and glaciers for global mapping of several soil properties. This stems from incomplete datasets of soil observations for these areas, despite the fact that extensive expert knowledge exists. In [Hengl et al. \(2017a\)](#) this was solved by integrating the expert knowledge in the form of expert-based pseudo-points to guide the ML model in areas of evident extrapolation. In [Koch et al. \(2019\)](#), 600 pseudo-points were also added in under-represented areas of the geographic space. The study stressed the importance of consulting an expert when building a ML model. In the same study, meaningful covariates were selected based on existing knowledge on the soil process, and plausibility of the predicted soil map was made in consideration of the knowledge of soil forming process. On many occasions, meaningful covariates are selected for mapping soil properties or classes. For example, [Brungard et al. \(2015\)](#) used a set of covariates selected a priori by an expert on the area under study. In [Viscarra-Rossel and Chen \(2011\)](#) a set of *scorpan* covariates was selected for mapping soil properties in Australia. These examples show that in the literature, adding expert-based pseudo-points and selecting meaningful

covariates are, to date, two straightforward options to include existing knowledge into a ML algorithm for DSM.

The above shows that little is known on how to account for existing knowledge in ML. Unfortunately, the complexity of the models increase on the same order as our understanding of the model functioning decreases. Improvement in this situation is made by ensuring that the map made by the ML algorithm matches the existing knowledge, for example by reflecting or corroborating the hypothesized soil pattern for an area. Pedological knowledge can be integrated to enforce results consistent with the existing scientific principles. This can be done during model building, calibration and validation. One can incorporate pedological knowledge by selecting appropriate covariates or by adding pseudo-points. In model building, incorporating knowledge takes the form of a hybrid model, a specific model architecture or objective function (in ANN models) constraining the calibration process according to specific knowledge. For example, [Wadoux \(2019\)](#) added the constraint that the prediction of topsoil clay, silt and sand must sum to 100% in a CNN model. Finally, pedological knowledge is used to make *post-hoc* checks on the plausibility of the calibrated model and predicted maps.

[Gahegan \(2019\)](#) stressed that since ML models (the author used the term “predictive process model” in the sense in which “machine learning model” is used in this article) have no connection to established theory, one can never be sure that the outcome is realistic given the real-world processes involved. The problem is that non-valid models are thus difficult to recognize and to reject since they are often not interpretable by a human. To ensure that models fit the existing pedological knowledge, they must be opened and understood in their functioning. Opening the “black box” is thus necessary but not straightforward (see next section on interpretability), and is often reduced to the analysis of which environmental covariates are the most often used by the model to make a prediction (e.g. [Mahmoudabadi et al. \(2017\)](#) or [McNicol et al. \(2019\)](#)) by using metrics based on variable importance.

Several authors, however, have warned against the use of these metrics for pedological interpretation (e.g. [Fourcade et al. \(2018\)](#) or [Wadoux et al. \(2020\)](#)). [Wadoux et al. \(2020\)](#) used meaningless, pseudo-covariates to map soil organic carbon over a hypothetical area. The authors obtained an accurate map, and concluded that the ML algorithm should not be used to interpret causal relationships. [Wadoux et al. \(2020\)](#) suggested using calibrated ML models as a “hypothesis discovery” tool, in which the mechanisms conveyed by the calibrated ML model are supplied to the researcher for possible explanations of the soil process, which can then be confronted to results of controlled experiments and principles of soil genesis. The challenge that then arises, noticed by [Gahegan \(2019\)](#) is the conversion of the mechanisms of the ML model (the model “knowledge”) from a data language to a human one. The data language is typically parameters or metrics such as the “mean decrease of purity” or “Gini importance index” of a covariate to assess its importance in the prediction of a soil property or class. Such metrics are not interpretable in terms of human explanation and they do not relate to soil processes. Translating the data language to the domain (the human language) requires some attention and further research. More discussion on this issue is found in [Gahegan et al. \(2001\)](#).

3.8. Interpretation of the models

When prediction is not the only interest, soil scientists rely on ML algorithms to gain insights into the modelled processes, for example to generate hypotheses. Despite providing higher prediction accuracy than other conventional models, most ML models are considered as a black box because their structure is complex and contains many adjustable parameters ([Breiman, 2001](#)). Broadly speaking, we do not learn from the model how the input covariates are related to the output soil property or classe, and what are the underlying mechanisms behind the prediction. This is unfortunate for soil science because in many cases the model itself is considered as a source of knowledge in addition to

the collected soil data. Scientific findings remain hidden when the model only gives a prediction without explanations. Interpretability of the model could warrant the extraction of the knowledge captured by the calibrated model. Miller (2019) defined interpretability as the degree to which human can understand the cause of a decision. In general, the need for interpretability of a ML algorithm stems from a deficiency in problem formalization (Doshi-Velez and Kim, 2017; Molnar, 2019). This means that for a given task (e.g. mapping the spatial distribution of soil organic carbon), the prediction itself does not fully solve the original problem if the interest is in understanding the mechanisms. We suggest three reasons that drive the demand for interpretability in DSM (adapted from Doshi-Velez and Kim (2017)). The first and most obvious reason is to increase our scientific understanding of the soil system by extracting knowledge from the mechanisms captured by the model. Scientists wish to know what are the drivers of a soil process and, more importantly, whether the mechanisms captured by the model confirm our scientific understanding of the system (see Section 3.7). The second reason is to audit the calibrated ML algorithm. Is the ML algorithm predicting for the right reasons? If a scientist makes a model for mapping the topsoil nitrogen content of a field, the interpretation might reveal that the model is actually predicting soil clay, that is, a proxy of the initial objective. The third reason is to avoid financial loss or to prevent a safety issue. Take the example of the remediation of the soil due to radioactive fallout after the Fukushima nuclear accident. A map of contaminated soils made by a ML algorithm would typically predict the dominant soil type characteristics, i.e. forest soil (about 75% of the area), for classification into contaminated or not contaminated areas. Interpretation of the model might then reveal that the important features learned by the model are unrealistic for agricultural landscapes and residential areas whose remediation is yet critical to safely move back the population (Evrard et al., 2019).

A straightforward way to increase interpretability is to decrease model complexity, for example by building a single DT instead of a RF composed of several thousand trees. A simple model enables visualization of the important mechanisms of the model and resultant explanations. For DT algorithms, it is possible to map the predicted values for specific rules (if the model is sufficiently simple). Decreasing complexity, however, is done at the expense of model prediction accuracy. For more complex ML algorithms, built-in features allow the user to retrieve the variable importance. In DT-like algorithms, the variable importance is derived from the thresholds used for the splits. For neural networks, the output weights associated with the input layer neurons provides an indication of the important features (Gahegan, 2000). One drawback of these techniques is their inability to provide information on whether the covariates have a causal link to the modelled soil property or class, which leads several authors to warn against their use for knowledge discovery (e.g. Fourcade et al., 2018). More importantly, these variable importance metrics are summary statistics that are not always meaningful and that are model-specific, i.e. they preclude comparison between models or parts of the predicted map.

Molnar (2019) reviewed techniques to interpret ML algorithms and defined two main categories of interpretation techniques. The first are model-specific ones, and routinely used in DSM activities (e.g. RF variable importance). The second category is model-free or model-agnostic techniques (Molnar, 2019). They enable the users to interpret any model, thus not restraining them to simple models or models with embedded features of interpretation. A summary of how model-agnostic techniques are employed is shown in Fig. 2. Examples of model-agnostic techniques are the partial dependence plot (Friedman, 2001) if the number of covariates is small (two maximum), individual conditional expectation (Goldstein et al., 2015), global or local model-agnostic explanation (LIME, Ribeiro et al., 2016), or Shapley values (as described in Molnar, 2019). Finally, sensitivity analysis is also a straightforward means of *post-hoc* interpretation of how the model output depends upon the different covariates.

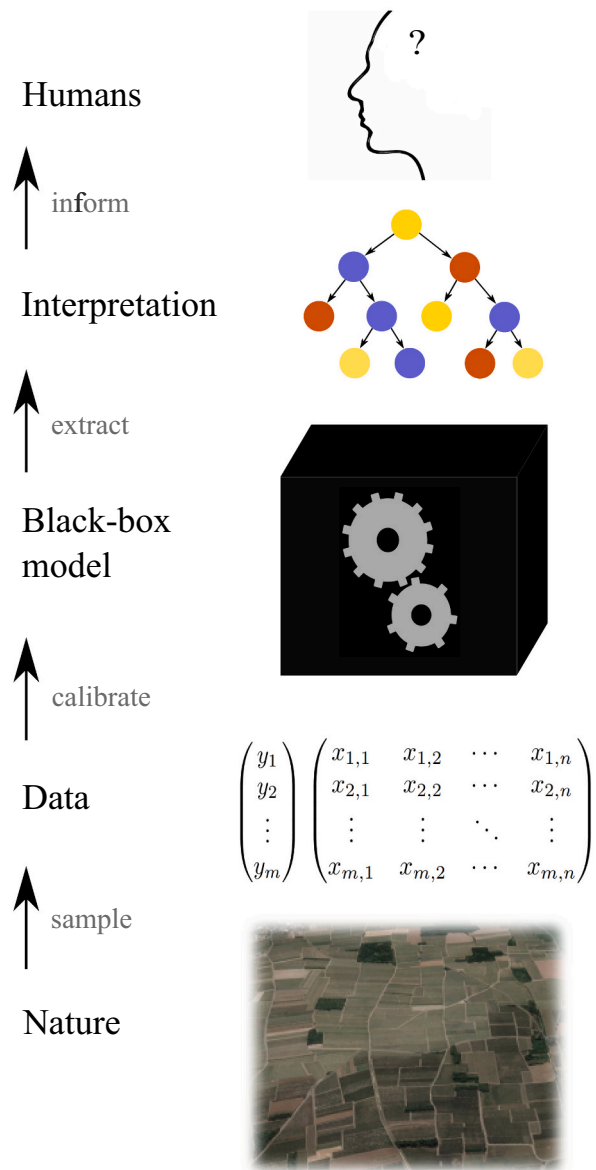


Fig. 2. Summary framework for model-agnostic interpretable ML, adapted from Molnar (2019). The lowest level is the reality, the unknown real-world soil that one wants to predict. The second level is the dataset that is extracted from the reality. We collect a fraction of the reality, a sample, and link it to exhaustively known environmental covariates. The relationships between the covariates and the sample is learned by a black-box machine learning model (level 3), on top of which comes the interpretation level to extract some knowledge from the structure of the calibrated model. The structure of the model is converted to human understandable knowledge.

4. The way forward

ML algorithms are now extensively used in soil mapping for regression and classification, much solely for prediction purposes. There is no doubt that prediction accuracy benefits from these data-driven models because ML algorithms are not constrained by a pre-defined model of the soil spatial variation, in comparison to geostatistical or mechanistic models. Recall that the goal of modelling is to both predict and obtain information on the soil processes, i.e. to increase our scientific understanding of the soil. This section provides a conceptual framework which can help ML-users to increase the scientific

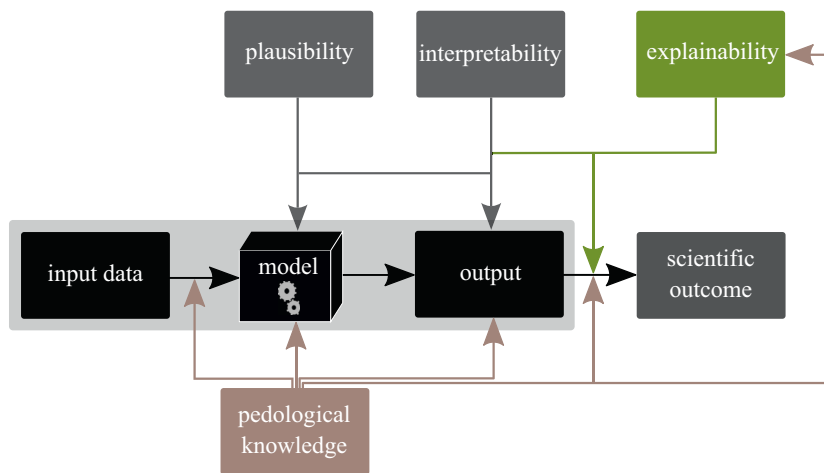


Fig. 3. Conceptual framework for the derivation of a scientific outcome from a ML model, adapted from Roscher et al. (2020). The light grey box represents the conventional use of ML algorithms in digital soil mapping, in which an output is derived from a calibrated ML model given a set of input data. A scientific outcome is obtained by explaining the output of a model using pedological knowledge, but also by ensuring scientific consistency at each link of the chain. Alternatively, a plausible and interpretable model can be explained using pedological knowledge.

understanding of the soil by building plausible models that can guide the soil scientist to obtain explanation and generate hypotheses. To increase the amount of information that we obtain from the modelling, future research on soil mapping with ML could incorporate the three core elements proposed by Roscher et al. (2020) and Lipton (2018) which we adapted, as follows:

Plausibility: Models should not only be accurate but also valid in light of the current knowledge and scientific theories. The plausibility is the solution path taken by the ML algorithm to link the input to the output, and does not depend directly on the data (Lipton, 2018). In practical terms, it starts with the model building step, by feeding the model with credible covariates and by accounting for the spatial particularities of soil data. Spatial or temporal correlation among data should be modelled, either by using a specific model (e.g. a CNN), or by using a model architecture that accounts for this particularity (see Section 3.3). Plausibility also takes the form of model constraints, to avoid the prediction of unrealistic proportions or ratios. The plausibility can be further tested in terms of model simulatability (Lipton, 2018). Since ML algorithms can model arbitrary patterns, there should be some attempts to test the model with synthetic data (e.g. as it is done in Lagacherie et al. (2020)) or data from a calibrated mechanistic model representing a large range of dynamics (Reichstein et al., 2019). Increasing model plausibility will facilitate the acceptance of ML to a large range of scenarios in soil science.

Interpretability: Interpretability is the translation of an abstract model or model output into terms understandable by humans (Montavon et al., 2018). Model interpretability pairs with model plausibility and hypothesis discovery. Complex and arbitrary patterns extracted from the data by an algorithm can be understood only through the transparency of the model. Interpretation is obtained by model-specific and model-agnostic methods, described in Section 3.8. Visual examination of the maps is also a means of interpretability. While complex ML models are potentially harmful because they often do not model any real-world process, there is an opportunity to challenge existing knowledge by *post-hoc* comparison of existing maps produced by expert knowledge with the maps predicted by a ML model, and through analysis of the striking differences. This is possible only if the model is interpretable by humans and the physical relationships between variables are realistic (i.e. if the model is *plausible*). Model interpretation is also an opportunity to generate new hypotheses, by interpreting the relationships found by the ML algorithm in the stores of soil data. The new hypotheses derived by these interpretations may challenge existing knowledge on the soil spatial variation and genesis.

Explainability: Modellers should shy away from mindless model fitting and prediction, and intensify research on models that both predict and explain. Explanations aim to answer the three questions: *what* is the modelled process?, *how* has it been modelled?, and *why* has this process

been modelled? (Miller, 2019). In this sense, explaining a process is an interpretation of a ML model plus expert knowledge and contextual information. For example, a different explanation is warranted when one wants to explain the pattern of a predicted soil map or the reason for two close predicted soil classes to be different. To explain, the modeller uses the data, the plausibility of the model and its interpretation using expert knowledge. Explainability is helped by model structure providing algorithmic explanations in the form of graphs or equations.

An example of model structure providing algorithmic explanations in DSM is found with the use of Bayesian belief networks (BBN, Cooper, 1990) in Mayr et al. (2010) and later Taalab et al. (2015). BBN is a probabilistic graphical model predicting the likely value of a soil property or class given conditional dependencies between covariates. Recent advances in ML have made a step further by discovering the graph structure directly from the data. However, while BBN is an interpretable ML model of conditional dependence between variables, the process that generated these dependencies remains hidden. To discover new processes from data, automated model discovery process is perhaps the way forward, in which equations describing a process are inductively assembled (i.e. using the data) into a single predictive model by heuristic search methods (Bridewell et al., 2008; Gahegan, 2019). The calibrated model is a set of equations constrained by existing, verified equations (e.g. differential equations of the water flow) representing known mechanisms. The model can be refined using expert knowledge and additional data (Dale et al., 1989). More importantly, these models produce explanations, which can be refuted or approved in light of scientific principles.

Fig. 3 illustrates the central role played by the three elements *plausibility*, *interpretability* and *explainability* in obtaining a scientific outcome from ML. Fig. 3 shows that the three core elements are conditioned to the use of pedological knowledge at each link of the chain. Enforcing pedological knowledge during modelling restricts the solution space to scientifically consistent results and may decrease the overall prediction accuracy. For DSM purposes, it is not obvious whether an increase of predictive accuracy is worth a substantial decrease in model consistency. For this reason, recent studies (e.g. Bennett et al., 2013; Lapuschkin et al., 2019) advocated the use of other criteria to measure the overall performance, such as model complexity or consistency (Karpadne et al., 2017). Including other criteria to assess the overall performance of a ML model would certainly make one step towards “conscious” DSM, and contribute to the uptake ML for knowledge discovery in soil science.

5. Conclusion

In this contribution, we have reviewed the current use of ML algorithms for DSM, and identified key challenges for which we provided

partial solutions. We draw the following conclusions:

- There has been a large number of studies mapping soil properties or classes using a ML algorithm. A wide range of soil properties, attributes and types have been predicted. Likewise, an increasing number of ML algorithms have been tested. Case studies are dominated by the use of legacy samples for local and regional scale (about 10^4 km²) areas. Ensembles of different algorithms are gaining more attention to improve prediction. Few studies reported the uncertainty associated with predictions.
- Attempts have been made to model soil variation in space and time, but we did not find any study developing a spatio-temporal ML algorithm for soil mapping.
- The configuration of a good sampling design for mapping with ML is largely unknown. The effect of the sampling design on model calibration and prediction has generally been disregarded. More research is needed in this direction.
- A large number of studies has focused exclusively on achieving a high mapping accuracy, thus are goal-oriented (producing a map) and favour prediction rather than inference, i.e. obtaining new knowledge on the underlying soil structure and process.
- Corollary to the previous item, comparisons between models and studies are made based on map quality indices, disregarding model complexity or consistency with respect to the existing pedological knowledge. The mapping accuracy seem to be the sole standard by which most maps quality is measured.
- More research is needed to account for pedological knowledge in the modelling process and to increase the interpretability of the ML models. Increasing interpretability may serve several purposes, among which obtaining a better understanding of mechanistic soil processes.

Overall, our review of the literature suggested that in recent studies inference is relegated to the background with the emergence of the mapping accuracy as the sole standard by which progress is measured. Most studies applied ML to case studies, and thus focused on high mapping accuracy. In this sense, any prediction can become a soil map, whether it contains soil knowledge or not. Several additional challenging aspects could help the soil scientist to obtain insights from the model, such as interpretability and the accounting of pedological knowledge.

We suggest that there is opportunity to include pedological knowledge at each step of the modelling chain, to improve or correct the existing dataset, to design the model architecture, to constrain the model calibration, or to analyse the output using *post-hoc* checks on the predicted soil maps. Future studies on DSM should use *plausible*, *interpretable* and *explainable* ML models to extract important scientific insights from soil data. One step towards achieving this goal is to integrate model consistency in addition to model prediction accuracy to evaluate the overall performance of the mapping approach. This will ensure that future studies use models that are not only accurate but also valid in light of the current knowledge and scientific theories, and that such models will serve as a source of information to generate hypotheses.

Declaration of Competing Interest

None

Acknowledgement

Budiman Minasny is member of a consortium supported by LE STUDIUM Loire Valley Institute for Advanced Studies through its LE STUDIUM Research Consortium Programme.

References

- Adhikari, K., Hartemink, A.E., Minasny, B., Kheir, R.B., Greve, M.B., Greve, M.H., 2014. Digital mapping of soil organic carbon contents and stocks in denmark. *PLoS One* 9, e105519.
- Aitkenhead, M.J., Coull, M.C., 2016. Mapping soil carbon stocks across Scotland using a neural network model. *Geoderma* 262, 187–198.
- Akpa, S.I.C., Odeh, I.O.A., Bishop, T.F.A., Hartemink, A.E., 2014. Digital mapping of soil particle-size fractions for Nigeria. *Soil Sci. Soc. Am. J.* 78, 1953–1966.
- Arrouays, D., McKenzie, N., Hempel, J., de Forges, A.R., McBratney, A.B., 2014. *GlobalSoilMap: Basis of the Global Spatial Soil Information System*. CRC Press, Boca Raton, USA.
- Batjes, N.H., Ribeiro, E., van Oostrum, A., Leenaars, J., Hengl, T., Mendes de Jesus, J., 2017. WoSIS: providing standardised soil profile data for the world. *Earth System Science Data* 9, 1–14.
- Beguín, J., Fuglistad, G.-A., Mansuy, N., Pare, D., 2017. Predicting soil properties in the Canadian boreal forest with limited data: Comparison of spatial and non-spatial statistical approaches. *Geoderma* 306, 195–205.
- Behrens, T., Förster, H., Scholten, T., Steinrücken, U., Spies, E.-D., Goldschmidt, M., 2005. Digital soil mapping using artificial neural networks. *J. Plant Nutr. Soil Sci.* 168, 21–33.
- Behrens, T., Zhu, A.-X., Schmidt, K., Scholten, T., 2010. Multi-scale digital terrain analysis and feature selection for digital soil mapping. *Geoderma* 155, 175–185.
- Behrens, T., Schmidt, K., MacMillan, R.A., Viscarra-Rossel, R.A., 2018a. Multi-scale digital soil mapping with deep learning. *Sci. Rep.* 8, 15244.
- Behrens, T., Schmidt, K., Viscarra-Rossel, R.A., Gries, P., Scholten, T., MacMillan, R.A., 2018b. Spatial modelling with Euclidean distance fields and machine learning. *Eur. J. Soil Sci.* 69, 757–770.
- Bel, L., Laurent, J.M., Bar-Hen, A., Allard, D., Cheddadi, R., 2005. A spatial extension of CART: application to classification of ecological data. In: Renard, P., Demougeot-Renard, H., Froidevaux, R. (Eds.), *Geostatistics for Environmental Applications*. Springer, Berlin, Heidelberg, pp. 99–109.
- Bel, L., Allard, D., Laurent, J.M., Cheddadi, R., Bar-Hen, A., 2009. CART algorithm for spatial data: Application to environmental and ecological data. *Computat. Stat. Data Anal.* 53, 3082–3093.
- Bennett, N.D., Croke, B.F.W., Guariso, G., Guillaume, J.H.A., Hamilton, S.H., Jakeman, A.J., Marsili-Libelli, S., Newham, L.T.H., Norton, J.P., Perrin, C., et al., 2013. Characterising performance of environmental models. *Environ. Model. Softw.* 40, 1–20.
- Blanco, C.M.G., Gomez, V.M.B., Crespo, P., Lief, M., 2018. Spatial prediction of soil water retention in a Páramo landscape: Methodological insight into machine learning using random forest. *Geoderma* 316, 100–114.
- Breiman, L., 2001. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Stat. Sci.* 16, 199–231.
- Breiman, L., 2017. *Classification and Regression Trees*. Routledge, New York, USA.
- Brenning, A., 2012. Spatial cross-validation and bootstrap for the assessment of prediction rules in remote sensing: The R package *sprr*. In: 2012 International Geoscience and Remote Sensing Symposium. IEEE, pp. 5372–5375.
- Bridewell, W., Langley, P., Todorovski, L., Džeroski, S., 2008. Inductive process modeling. *Mach. Learn.* 71, 1–32.
- Brungard, C.W., Boettinger, J.L., Duniway, M.C., Wills, S.A., Edwards Jr., T.C., 2015. Machine learning for predicting soil classes in three semi-arid landscapes. *Geoderma* 239, 68–83.
- Brus, D.J., 2019. Sampling for digital soil mapping: A tutorial supported by R scripts. *Geoderma* 338, 464–480.
- Bui, E.N., Loughhead, A., Corner, R., 1999. Extracting soil-landscape rules from previous soil surveys. *Soil Res.* 37, 495–508.
- Bui, E., Henderson, B., Viergever, K., 2009. Using knowledge discovery with data mining from the Australian Soil Resource Information System database to inform soil carbon mapping in Australia. *Glob. Biogeochem. Cycles* 23, GB4033.
- Carré, F., McBratney, A.B., Minasny, B., 2007. Estimation and potential improvement of the quality of legacy soil samples for digital soil mapping. *Geoderma* 141, 1–14.
- Chen, S., Mulder, V.L., Martin, M.P., Walter, C., Lacoste, M., Richer-de Forges, A.C., Saby, N.P.A., Loiseau, T., Hu, B., Arrouays, D., 2019. Probability mapping of soil thickness by random survival forest at a national scale. *Geoderma* 344, 184–194.
- Cialella, A.T., Dubayah, R., Lawrence, W., Levine, E., 1997. Predicting soil drainage class using remotely sensed and digital elevation data. *Photogramm. Eng. Remote. Sens.* 63, 171–177.
- Cooper, G.F., 1990. The computational complexity of probabilistic inference using bayesian belief networks. *Artif. Intell.* 42, 393–405.
- Coveney, P.V., Dougherty, E.R., Highfield, R.R., 2016. Big data need big theory too. *Philos. Trans. R. Soc. A Math. Phys. Eng. Sci.* 374, 20160153.
- Cressie, N., Johannesson, G., 2008. Fixed rank kriging for very large spatial data sets. *J. Roy. Stat. Soc.* 70, 209–226.
- Dai, F., Zhou, Q., Lv, Z., Wang, X., Liu, G., 2014. Spatial prediction of soil organic matter content integrating artificial neural network and ordinary kriging in Tibetan Plateau. *Ecol. Indic.* 45, 184–194.
- Dale, M.B., McBratney, A.B., Russell, J.S., 1989. On the role of expert systems and numerical taxonomy in soil classification. *J. Soil Sci.* 40, 223–234.
- De Grujter, J.J., Brus, D.J., Bierkens, M.F.P., Kotters, M., 2006. *Sampling for Natural Resource Monitoring*. Springer Science & Business Media, Dordrecht, NL.
- Dharumarajan, S., Hegde, R., Singh, S.K., 2017. Spatial prediction of major soil properties using Random Forest techniques-A case study in semi-arid tropics of South India. *Geoderma Regional* 10, 154–162.
- Doshi-Velez, F., Kim, B., 2017. Towards a rigorous science of interpretable machine

- learning. arXiv:1702.08608.
- Ellili, Y., Malone, B.P., Michot, D., Minasny, B., Vincent, S., Walter, C., Lemerrier, B., 2019. Comparing three approaches of spatial disaggregation of legacy soil maps based on DSMART algorithm. *SOIL Discussions*.
- Evraud, O., Lacey, J.P., Nakao, A., 2019. Effectiveness of landscape decontamination following the Fukushima nuclear accident: a review. *SOIL* 5, 333–350.
- Fick, S.E., Hijmans, R.J., 2017. Worldclim 2: new 1-km spatial resolution climate surfaces for global land areas. *Int. J. Climatol.* 37, 4302–4315.
- Forkuor, G., Hounkpatin, O.K.L., Welp, G., Thiel, M., 2017. High resolution mapping of soil properties using remote sensing variables in south-western Burkina Faso: a comparison of machine learning and multiple linear regression models. *PLoS One* 12, e0170478.
- 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.* 25, 1189–1232.
- Gahegan, M., 2000. On the application of inductive machine learning tools to geographical analysis. *Geogr. Anal.* 32, 113–139.
- Gahegan, M., 2019. Fourth paradigm GIScience? prospects for automated discovery and explanation from data. *Int. J. Geogr. Inf. Sci.* 34, 1–21.
- Gahegan, M., Wachowicz, M., Harrower, M., Rhyne, T.-M., 2001. The integration of geographic visualization with knowledge discovery in databases and geocomputation. *Cartogr. Geogr. Inf. Sci.* 28, 29–44.
- Gasch, C.K., Hengl, T., Gräler, B., Meyer, H., Magney, T.S., Brown, D.J., 2015. Spatio-temporal interpolation of soil water, temperature, and electrical conductivity in 3D + T: The Cook Agronomy Farm data set. *Spatial Statistics* 14, 70–90.
- Gascon, F., Bouzinac, C., Thépaut, O., Jung, M., Francesconi, B., Louis, J., Lonjou, V., Lafrance, B., Massera, S., Gaudel-Vacaresse, A., et al., 2017. Copernicus Sentinel-2A calibration and products validation status. *Remote Sens.* 9, 584.
- Georganos, S., Grippa, T., Gadiaga, A.N., Linard, C., Lennert, M., Vanhuyse, S., Mboga, N.O., Wolff, E., Kalogirou, S., 2019. Geographical random forests: A spatial extension of the random forest algorithm to address spatial heterogeneity in remote sensing and population modelling. *Geocarto International* 1, 1–12.
- Goldstein, A., Kapelner, A., Bleich, J., Pitkin, E., 2015. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *J. Comput. Graph. Stat.* 24, 44–65.
- Gomes, L.C., Faria, R.M., de Souza, E., Veloso, G.V., Schaefer, C.E.G.R., Fernandes Filho, E.I., 2019. Modelling and mapping soil organic carbon stocks in Brazil. *Geoderma* 340, 337–350.
- Grimm, R., Behrens, T., Märker, M., Elsenbeer, H., 2008. Soil organic carbon concentrations and stocks on Barro Colorado Island—Digital soil mapping using Random Forests analysis. *Geoderma* 146, 102–113.
- Guevara, M., Olmedo, G.F., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C., Arévalo, G.E., Arroyo-Cruz, C.E., Bolívar, A., Bunning, S., et al., 2018. No silver bullet for digital soil mapping: country-specific soil organic carbon estimates across Latin America. *SOIL* 4, 173–193.
- Hamzehpour, N., Shafizadeh-Moghadam, H., Valavi, R., 2019. Exploring the driving forces and digital mapping of soil organic carbon using remote sensing and soil texture. *CATENA* 182, 104141.
- Hansen, M.K., Brown, D.J., Dennison, P.E., Graves, S.A., Bricklemeyer, R.S., 2009. Inductively mapping expert-derived soil-landscape units within dambo wetland catenae using multispectral and topographic data. *Geoderma* 150, 72–84.
- Häring, T., Dietz, E., Osenstetter, S., Koschitzki, T., Schröder, B., 2012. Spatial disaggregation of complex soil map units: A decision-tree based approach in Bavarian forest soils. *Geoderma* 185, 37–47.
- Hartmann, J., Moosdorf, N., 2012. The new global lithological map database GLiM: A representation of rock properties at the Earth surface. *Geochim. Geophys. Geosyst.* 13, 1–37.
- Hawkins, B.A., 2012. Eight (and a half) deadly sins of spatial analysis. *J. Biogeogr.* 39, 1–9.
- He, H., Garcia, E.A., 2009. Learning from imbalanced data. *IEEE Trans. Knowl. Data Eng.* 21, 1263–1284.
- Henderson, B.L., Bui, E.N., Moran, C.J., Simon, D.A.P., 2005. Australia-wide predictions of soil properties using decision trees. *Geoderma* 124, 383–398.
- Hengl, T., de Jesus, J.M., MacMillan, R.A., Batjes, N.H., Heuvelink, G.B.M., Ribeiro, E., Samuel-Rosa, A., Kempen, B., Leenaars, J.G., Walsh, M.G., et al., 2014. SoilGrids1km—global soil information based on automated mapping. *PLoS One* 9, e105992.
- Hengl, T., de Jesus, J.M., Heuvelink, G.B.M., Gonzalez, M.R., Kilibarda, M., Blagotić, A., Shanguan, W., Wright, M.N., Geng, X., Bauer-Marschallinger, B., et al., 2017a. SoilGrids250m: Global gridded soil information based on machine learning. *PLoS One* 12, e0169748.
- Hengl, T., Leenaars, J.G.B., Shepherd, K.D., Walsh, M.G., Heuvelink, G.B.M., Mamo, T., Tilahun, H., Berkhout, E., Cooper, M., Fegraus, E., et al., 2017b. Soil nutrient maps of Sub-Saharan Africa: assessment of soil nutrient content at 250 m spatial resolution using machine learning. *Nutr. Cycl. Agroecosyst.* 109, 77–102.
- Hengl, T., Nussbaum, M., Wright, M.N., Heuvelink, G.B.M., Gräler, B., 2018. Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables. *PeerJ* 6, e5518.
- Heung, B., Ho, H.C., Zhang, J., Knudby, A., Bulmer, C.E., Schmidt, M.G., 2016. An overview and comparison of machine-learning techniques for classification purposes in digital soil mapping. *Geoderma* 265, 62–77.
- Heuvelink, G.B.M., Webster, R., 2001. Modelling soil variation: past, present, and future. *Geoderma* 100, 269–301.
- 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., 2020. Machine learning in space and time for modelling soil organic carbon change. *Eur. J. Soil Sci.* 1–17.
- Holmes, K.W., Odgers, N.P., Griffin, E.A., van Gool, D., 2014. Spatial disaggregation of conventional soil mapping across Western Australia using DSMART. In: Arruays, D., McKenzie, N., Hempel, J., Richer de Forges, A., McBratney, A.B. (Eds.), *GlobalSoilMap: Basis of the Global Spatial Soil Information System*. Taylor & Francis, London, UK, pp. 273–279.
- Hounkpatin, O.K.L., Schmidt, K., Stumpf, F., Forkuor, G., Behrens, T., Scholten, T., Amelung, W., Welp, G., 2018. Predicting reference soil groups using legacy data: A data pruning and Random Forest approach for tropical environment (Dano catchment, Burkina Faso). *Sci. Rep.* 8, 9959.
- Illés, G., Sutikno, S., Szatmári, G., Sandhyavriti, A., Pásztor, L., Kristijono, A., Molnár, G., Yusa, M., Székely, B., 2019. Facing the peat CO₂ threat: digital mapping of Indonesian peatlands—a proposed methodology and its application. *J. Soils Sediments* 1–16.
- Jenny, H., 1941. *Factors of Soil Formation: A System of Quantitative Pedology*. McGrawHill, New York, USA.
- Jiang, Z., Li, Y., Shekhar, S., Rampi, L., Knight, J., 2017. Spatial ensemble learning for heterogeneous geographic data with class ambiguity: A summary of results. In: *Proceedings of the 25th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*. 23. ACM, pp. 1–10.
- Kalambukattu, J.G., Kumar, S., Raj, R.A., 2018. Digital soil mapping in a Himalayan watershed using remote sensing and terrain parameters employing artificial neural network model. *Environ. Earth Sci.* 77, 203.
- Karpatne, A., Atluri, G., Faghmous, J.H., Steinbach, M., Banerjee, A., Ganguly, A., Shekhar, S., Samatova, N., Kumar, V., 2017. Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Trans. Knowl. Data Eng.* 29, 2318–2331.
- Keskin, H., Grunwald, S., Harris, W.G., 2019. Digital mapping of soil carbon fractions with machine learning. *Geoderma* 339, 40–58.
- Khair, R.B., Greve, M.H., Abdallah, C., Dalgaard, T., 2010a. Spatial soil zinc content distribution from terrain parameters: A GIS-based decision-tree model in Lebanon. *Environ. Pollut.* 158, 520–528.
- Khair, R.B., Greve, M.H., Bøcher, P.K., Greve, M.B., Larsen, R., McCloy, K., 2010b. Predictive mapping of soil organic carbon in wet cultivated lands using classification-tree based models: The case study of Denmark. *J. Environ. Manag.* 91, 1150–1160.
- Kirkwood, C., Cave, M., Beamish, D., Grebb, S., Ferreira, A., 2016. A machine learning approach to geochemical mapping. *J. Geochem. Explor.* 167, 49–61.
- Koch, J., Stisen, S., Refsgaard, J.C., Erntsen, V., Jakobsen, P.R., Højberg, A.L., 2019. Modeling depth of the redox interface at high resolution at national scale using random forest and residual Gaussian simulation. *Water Resour. Res.* 55, 1451–1469.
- Kovačević, M., Bajat, B., Gajić, B., 2010. Soil type classification and estimation of soil properties using support vector machines. *Geoderma* 154, 340–347.
- Kühn, I., Dormann, C.F., 2012. Less than eight (and a half) misconceptions of spatial analysis. *J. Biogeogr.* 39, 995–998.
- Kühn, I., Nobis, M.P., Durka, W., 2009. Combining spatial and phylogenetic eigenvector filtering in trait analysis. *Glob. Ecol. Biogeogr.* 18, 745–758.
- Lacoste, M., Minasny, B., McBratney, A., Michot, D., Viaud, V., Walter, C., 2014. High resolution 3d mapping of soil organic carbon in a heterogeneous agricultural landscape. *Geoderma* 213, 296–311.
- Lagacherie, P., 2008. Digital soil mapping: a state of the art. In: Hartemink, A.E., McBratney, A., de Lourdes Mendonça-Santos, M. (Eds.), *Digital Soil Mapping with Limited Data*. Springer, Dordrecht, Netherlands, pp. 3–14.
- Lagacherie, P., Holmes, S., 1997. Addressing geographical data errors in a classification tree for soil unit prediction. *Int. J. Geogr. Inf. Sci.* 11, 183–198.
- Lagacherie, P., McBratney, A., 2006. Spatial soil information systems and spatial soil inference systems: perspectives for digital soil mapping. *Dev. Soil Sci.* 31, 3–22.
- Lagacherie, P., Arruays, D., Bourenane, H., Gomez, C., Nkuba-Kasanda, L., 2020. Analysing the impact of soil spatial sampling on the performances of digital soil mapping models and their evaluation: A numerical experiment on quantile random forest using clay contents obtained from vis-nir-swir hyperspectral imagery. *Geoderma* 375, 114503.
- Lamichhane, S., Kumar, L., Wilson, B., 2019. Digital soil mapping algorithms and covariates for soil organic carbon mapping and their implications: A review. *Geoderma* 395–413.
- Lapushkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W., Müller, K.-R., 2019. Unmasking Clever Hans predictors and assessing what machines really learn. *Nat. Commun.* 10, 1096.
- Le Rest, K., Pinaud, D., Monestiez, P., Chadoeuf, J., Bretagnolle, V., 2014. Spatial leave-one-out cross-validation for variable selection in the presence of spatial autocorrelation. *Glob. Ecol. Biogeogr.* 23, 811–820.
- Ließ, M., Glaser, B., Huwe, B., 2012. Uncertainty in the spatial prediction of soil texture: comparison of regression tree and random forest models. *Geoderma* 170, 70–79.
- Ließ, M., Schmidt, J., Glaser, B., 2016. Improving the spatial prediction of soil organic carbon stocks in a complex tropical mountain landscape by methodological specifications in machine learning approaches. *PLoS One* 11, e0153673.
- Lipton, Z.C., 2018. The mythos of model interpretability. *Queue* 16, 31–57.
- Liu, F., Zhang, G.-L., Song, X., Li, D., Zhao, Y., Yang, J., Wu, H., Yang, F., 2019. High-resolution and three-dimensional mapping of soil texture of China. *Geoderma* 361, 114061.
- Lorenzetti, R., Barbetti, R., Fantappiè, M., L'Abate, G., Costantini, E.A., 2015. Comparing data mining and deterministic pedology to assess the frequency of wrb reference soil groups in the legend of small scale maps. *Geoderma* 237, 237–245.
- Ma, Y., Minasny, B., Malone, B.P., Mcbratney, A.B., 2019. Pedology and digital soil mapping (DSM). *Eur. J. Soil Sci.* 70, 216–235.
- Ma, T., Brus, D.J., Zhu, A.-X., Zhang, L., Scholten, T., 2020. Comparison of conditioned Latin hypercube and feature space coverage sampling for predicting soil classes using

- simulation from soil maps. *Geoderma* 370, 114366.
- MacKay, D.J.C., 1992. Information-based objective functions for active data selection. *Neural Comput.* 4, 590–604.
- Mahmoudabadi, E., Karimi, A., Haghnia, G.H., Sepehr, A., 2017. Digital soil mapping using remote sensing indices, terrain attributes, and vegetation features in the rangelands of northeastern Iran. *Environ. Monit. Assess.* 189, 500.
- Malone, B.P., McBratney, A.B., Minasny, B., Laslett, G.M., 2009. Mapping continuous depth functions of soil carbon storage and available water capacity. *Geoderma* 154, 138–152.
- Mansuy, N., Thiffault, E., Paré, D., Bernier, P., Guindon, L., Villemaire, P., Poirier, V., Beaudoin, A., 2014. Digital mapping of soil properties in Canadian managed forests at 250 m of resolution using the k-nearest neighbor method. *Geoderma* 235, 59–73.
- Marchant, B.P., Viscarra Rossel, R.A., Webster, R., 2013. Fluctuations in method-of-moments variograms caused by clustered sampling and their elimination by declustering and residual maximum likelihood estimation. *Eur. J. Soil Sci.* 64, 401–409.
- Massawe, B.H.J., Subburayalu, S.K., Kaaya, A.K., Winowiecki, L., Slater, B.K., 2018. Mapping numerically classified soil taxa in Kilombero valley, Tanzania using machine learning. *Geoderma* 311, 143–148.
- Mayr, T., Rivas-Casado, M., Bellamy, P., Palmer, R., Zawadzka, J., Corstanje, R., 2010. Two methods for using legacy data in digital soil mapping. In: Boettinger, J., Howell, D., Moore, A., Hartemink, A., Kienast-Brown, S. (Eds.), *Digital Soil Mapping*. Springer, Dordrecht, NL, pp. 191–202.
- McBratney, A.B., Santos, M.M., Minasny, B., 2003. On digital soil mapping. *Geoderma* 117, 3–52.
- McNicol, G., Bulmer, C., D'Amore, D., Sanborn, P., Saunders, S., Giesbrecht, I., Arriola, S.G., Bidlack, A., Butman, D., Buma, B., 2019. Large, climate-sensitive soil carbon stocks mapped with pedology-informed machine learning in the North Pacific coastal temperate rainforest. *Environ. Res. Lett.* 14, 014004.
- Meyer, H., Reudenbach, C., Hengl, T., Katurji, M., Naus, T., 2018. Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation. *Environ. Model. Softw.* 101, 1–9.
- Meyer, H., Reudenbach, C., Wöllauer, S., Naus, T., 2019. Importance of spatial predictor variable selection in machine learning applications—moving from data reproduction to spatial prediction. *Ecol. Model.* 411, 108815.
- Micheletti, N., Foresti, L., Robert, S., Leuenberger, M., Pedrazzini, A., Jaboyedoff, M., Kanevski, M., 2014. Machine learning feature selection methods for landslide susceptibility mapping. *Math. Geosci.* 46, 33–57.
- Miller, T., 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artif. Intell.* 267, 1–38.
- Miller, B.A., Koszinski, S., Wehrhan, M., Sommer, M., 2015a. Comparison of spatial association approaches for landscape mapping of soil organic carbon stocks. *Soil* 1, 217–233.
- Miller, B.A., Koszinski, S., Wehrhan, M., Sommer, M., 2015b. Impact of multi-scale predictor selection for modeling soil properties. *Geoderma* 239, 97–106.
- Mira, M., Weiss, M., Baret, F., Courault, D., Hagolle, O., Gallego-Elvira, B., Olioso, A., 2015. The MODIS (collection V006) BRDF/albedo product MCD43D: Temporal course evaluated over agricultural landscape. *Remote Sens. Environ.* 170, 216–228.
- Mirza, M., Osindero, S., 2014. Conditional generative adversarial nets. arXiv:1411.1784.
- Molnar, C., 2019. Interpretable machine learning. Lulu, Morrisville, USA.
- Montavon, G., Samek, W., Müller, K.-R., 2018. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing* 73, 1–15.
- Moran, C.J., Bui, E.N., 2002. Spatial data mining for enhanced soil map modelling. *Int. J. Geogr. Inf. Sci.* 16, 533–549.
- Mosleh, Z., Salehi, M.H., Jafari, A., Borujeni, I.E., Mehnatkesh, A., 2016. The effectiveness of digital soil mapping to predict soil properties over low-relief areas. *Environ. Monit. Assess.* 188, 195.
- Mulder, V.L., Lacoste, M., Richer-de Forges, A.C., Martin, M.P., Arrouays, D., 2016. National versus global modelling the 3D distribution of soil organic carbon in mainland France. *Geoderma* 263, 16–34.
- Nauman, T.W., Duniway, M.C., 2019. Relative prediction intervals reveal larger uncertainty in 3D approaches to predictive digital soil mapping of soil properties with legacy data. *Geoderma* 347, 170–184.
- Nussbaum, M., Spiess, K., Baltensweiler, A., Grob, U., Keller, A., Greiner, L., Schaepman, M.E., Papritz, A., 2018. Evaluation of digital soil mapping approaches with large sets of environmental covariates. *Soil* 4, 1–22.
- Odgers, N.P., Sun, W., McBratney, A.B., Minasny, B., Clifford, D., 2014. Disaggregating and harmonising soil map units through resampled classification trees. *Geoderma* 214, 91–100.
- Oldeman, L.R., Van Engelen, V.W.P., 1993. A world soils and terrain digital database (SOTER)—An improved assessment of land resources. *Geoderma* 60, 309–325.
- Oliver, M.A., 1987. Geostatistics and its application to soil science. *Soil Use Manag.* 3, 8–20.
- Otto, 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. Indic.* 77, 139–150.
- Padarian, J., Minasny, B., McBratney, A.B., 2019. Using deep learning for digital soil mapping. *Soil* 5, 79–89.
- Pahlavan-Rad, M.R., Akbarimoghaddam, A., 2018. Spatial variability of soil texture fractions and pH in a flood plain (case study from eastern Iran). *CATENA* 160, 275–281.
- Poggio, L., Lassauce, A., Gimona, A., 2019. Modelling the extent of northern peat soil and its uncertainty with sentinel: Scotland as example of highly cloudy region. *Geoderma* 346, 63–74.
- Pohjankukka, J., Pahikkala, T., Nevalainen, P., Heikkonen, J., 2017. Estimating the prediction performance of spatial models via spatial k-fold cross validation. *Int. J. Geogr. Inf. Sci.* 31, 2001–2019.
- Pouladi, N., Möller, A.B., Tabatabai, S., Greve, M.H., 2019. Mapping soil organic matter contents at field level with cubist, random forest and kriging. *Geoderma* 342, 85–92.
- Pozdnoukhov, A., Kanevski, M., 2006. Monitoring network optimisation for spatial data classification using support vector machines. *Int. J. Environ. Pollut.* 28, 465–484.
- Rahman, R., Otridge, J., Pal, R., 2017. IntegratedMRF: random forest-based framework for integrating prediction from different data types. *Bioinformatics* 33, 1407–1410.
- Ramcharan, A., Hengl, T., Nauman, T., Brungard, C., Waltman, S., Wills, S., Thompson, J., 2018. Soil property and class maps of the conterminous United States at 100-meter spatial resolution. *Soil Sci. Soc. Am. J.* 82, 186–201.
- Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., et al., 2019. Deep learning and process understanding for data-driven earth system science. *Nature* 566, 195–204.
- Ribeiro, M.T., Singh, S., Guestrin, C., 2016. Why should I trust you?: Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, pp. 1135–1144.
- Roberts, D.R., Bahn, V., Ciuti, S., Boyce, M.S., Elith, J., Guilleri-Arroita, G., Hauenstein, S., Lahoz-Monfort, J.J., Schröder, B., Thuiller, W., et al., 2017. Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure. *Ecography* 40, 913–929.
- Roscher, R., Bohn, B., Duarte, M.F., Garcke, J., 2020. Explainable machine learning for scientific insights and discoveries. *IEEE Access* 8, 42200–42216.
- Rudiyanto, B., Minasny, Setiawan, Saptomo, S.K., McBratney, A.B., et al., 2018. Open digital mapping as a cost-effective method for mapping peat thickness and assessing the carbon stock of tropical peatlands. *Geoderma* 313, 25–40.
- Ruß, G., Brenning, A., 2010. Data mining in precision agriculture: management of spatial information. In: *International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*. Springer, pp. 350–359.
- Schratz, P., Muenchow, J., Iturriza, E., Richter, J., Brenning, A., 2019. Hyperparameter tuning and performance assessment of statistical and machine-learning algorithms using spatial data. *Ecol. Model.* 406, 109–120.
- Scully, P., Franklin, J., Chadwick, O.A., 2005. The application of classification tree analysis to soil type prediction in a desert landscape. *Ecol. Model.* 181, 1–15.
- Segal, M., Xiao, Y., 2011. Multivariate random forests. *Wiley Interdisc. Rev.* 1, 80–87.
- Sekulić, A., Kilibarda, M., Heuvelink, G.B.M., Nikolić, M., Bajat, B., 2020. Random forest spatial interpolation. *Remote Sens.* 12, 1687.
- Sergeev, A., Buevich, A., Baglaeva, E., Shichkin, A., 2019. Combining spatial autocorrelation with machine learning increases prediction accuracy of soil heavy metals. *CATENA* 174, 425–435.
- Shariffar, A., Sarmadian, F., Malone, B.P., Minasny, B., 2019. Addressing the issue of digital mapping of soil classes with imbalanced class observations. *Geoderma* 350, 84–92.
- Shi, J., Yang, L., Zhu, A., Qin, C., Liang, P., Zeng, C., Pei, T., et al., 2018. Machine-learning variables at different scales vs. knowledge-based variables for mapping multiple soil properties. *Soil Sci. Soc. Am. J.* 82, 645–656.
- Siewert, M.B., 2018. High-resolution digital mapping of soil organic carbon in permafrost terrain using machine learning: a case study in a sub-Arctic peatland environment. *Biogeosciences* 15, 1663–1682.
- da Silva Chagas, C., de Carvalho Junior, W., Bhering, S.B., Calderano Filho, B., 2016. Spatial prediction of soil surface texture in a semiarid region using random forest and multiple linear regressions. *CATENA* 139, 232–240.
- Sinha, P., Gaughan, A.E., Stevens, F.R., Nieves, J.J., Sorichetta, A., Tatem, A.J., 2019. Assessing the spatial sensitivity of a random forest model: Application in gridded population modeling. *Comput. Environ. Urban. Syst.* 75, 132–145.
- Somarathna, P.D.S.N., Malone, B.P., Minasny, B., 2016. Mapping soil organic carbon content over New South Wales, Australia using local regression kriging. *Geoderma Regional* 7, 38–48.
- Somarathna, P.D.S.N., Minasny, B., Malone, B.P., 2017. More data or a better model? figuring out what matters most for the spatial prediction of soil carbon. *Soil Sci. Soc. Am. J.* 81, 1413–1426.
- Song, Y.-Q., Zhao, X., Su, H.-Y., Li, B., Hu, Y.-M., Cui, X.-S., 2018. Predicting spatial variations in soil nutrients with hyperspectral remote sensing at regional scale. *Sensors* 18, 3086.
- Song, X.-D., Wu, H.-Y., Ju, B., Liu, F., Yang, F., Li, D.-C., Zhao, Y.-G., Yang, J.-L., Zhang, G.-L., 2020. Pedoclimatic zone-based three-dimensional soil organic carbon mapping in China. *Geoderma* 363, 114145.
- Stojanova, D., Ceci, M., Appice, A., Malerba, D., Džeroski, S., 2013. Dealing with spatial autocorrelation when learning predictive clustering trees. *Ecological Informatics* 13, 22–39.
- Stumpf, F., Schmidt, K., Behrens, T., Schönbrodt-Stitt, S., Buzzo, G., Dumperth, C., Wadoux, A., Xiang, W., Scholten, T., 2016. Incorporating limited field operability and legacy soil samples in a hypercube sampling design for digital soil mapping. *J. Plant Nutr.* *Soil Sci.* 179, 499–509.
- Subburayalu, S.K., Slater, B.K., 2013. Soil series mapping by knowledge discovery from an Ohio county soil map. *Soil Sci. Soc. Am. J.* 77, 1254–1268.
- Szathmári, G., Pásztor, L., 2019. Comparison of various uncertainty modelling approaches based on geostatistics and machine learning algorithms. *Geoderma* 337, 1329–1340.
- Szathmári, G., Pirkó, B., Koós, S., Laborci, A., Bakacsi, Z., Szabó, J., Pásztor, L., 2019. Spatio-temporal assessment of topsoil organic carbon stock change in Hungary. *Soil Tillage Res.* 195, 104410.
- Taalab, K., Corstanje, R., Mayr, T., Whelan, M., Creamer, R., 2015. The application of expert knowledge in bayesian networks to predict soil bulk density at the landscape scale. *Eur. J. Soil Sci.* 66, 930–941.
- Tagasovska, N., Lopez-Paz, D., 2019. Single-model uncertainties for deep learning. In: *Advances in Neural Information Processing Systems*, pp. 6417–6428.
- Taghizadeh-Mehrjardi, R., Minasny, B., Sarmadian, F., Malone, B.P., 2014. Digital mapping of soil salinity in Ardakan region, central Iran. *Geoderma* 213, 15–28.

- Taghizadeh-Mehrjardi, R., Nabiollahi, K., Kerry, R., 2016a. Digital mapping of soil organic carbon at multiple depths using different data mining techniques in Baneh region, Iran. *Geoderma* 266, 98–110.
- Taghizadeh-mehrjardi, R., Toomanian, N., Khavaninzadeh, A.R., Jafari, A., Triantafyllis, J., 2016b. Predicting and mapping of soil particle-size fractions with adaptive neuro-fuzzy inference and ant colony optimization in central Iran. *Eur. J. Soil Sci.* 67, 707–725.
- Taghizadeh-Mehrjardi, R., Minasny, B., Toomanian, N., Zeraatpisheh, M., Amirian-Chakan, A., Triantafyllis, J., 2019a. Digital mapping of soil classes using ensemble of models in Isfahan region, Iran. *Soil Systems* 3, 37.
- Taghizadeh-Mehrjardi, R., Schmidt, K., Eftekhari, K., Behrens, T., Jamshidi, M., Davatgaar, N., Toomanian, N., Scholten, T., 2019b. Synthetic resampling strategies and machine learning for digital soil mapping in Iran. *Eur. J. Soil Sci.* 71 (3), 352–368.
- Tajik, S., Ayoubi, S., Shirani, H., Zeraatpisheh, M., 2019. Digital mapping of soil invertebrates using environmental attributes in a deciduous forest ecosystem. *Geoderma* 353, 252–263.
- Tuia, D., Pozdnoukhov, A., Foresti, L., Kanevski, M., 2013. Active learning for monitoring network optimization. In: Mateu, J., Müller, W.G. (Eds.), *Spatio-Temporal Design: Advances in Efficient Data Acquisition*. Wiley Online Library, Chichester, UK, pp. 285–318.
- Tziachris, P., Aschonitis, V., Chatzistathis, T., Papadopoulou, M., 2019. Assessment of spatial hybrid methods for predicting soil organic matter using DEM derivatives and soil parameters. *CATENA* 174, 206–216.
- Uria, B., Côté, M.-A., Gregor, K., Murray, I., Larochelle, H., 2016. Neural autoregressive distribution estimation. *J. Machine Learning Res.* 17, 7184–7220.
- Van Den Hoogen, J., Geisen, S., Routh, D., Ferris, H., Traunspurger, W., Wardle, D.A., De Goede, R.G.M., Adams, B.J., Ahmad, W., Andriuzzi, W.S., et al., 2019. Soil nematode abundance and functional group composition at a global scale. *Nature* 572, 194–198.
- Van Groenigen, J.W., 1997. Spatial simulated annealing for optimizing sampling. In: *geoENV I—Geostatistics for Environmental Applications*. Springer, Berlin, Germany, pp. 351–361.
- Vaysse, K., Lagacherie, P., 2015. Evaluating digital soil mapping approaches for mapping GlobalSoilMap soil properties from legacy data in Languedoc-Roussillon (France). *Geoderma Regional* 4, 20–30.
- Vaysse, K., Lagacherie, P., 2017. Using quantile regression forest to estimate uncertainty of digital soil mapping products. *Geoderma* 291, 55–64.
- Vermeulen, D., Van Niekerk, A., 2017. Machine learning performance for predicting soil salinity using different combinations of geomorphometric covariates. *Geoderma* 299, 1–12.
- Vincent, S., Lemerrier, B., Berthier, L., Walter, C., 2018. Spatial disaggregation of complex soil map units at the regional scale based on soil-landscape relationships. *Geoderma* 311, 130–142.
- Viscarra-Rossel, R.A., Chen, C., 2011. Digitally mapping the information content of visible-near infrared spectra of surficial Australian soils. *Remote Sens. Environ.* 115, 1443–1455.
- Viscarra-Rossel, R.A., Chen, C., Grundy, M.J., Searle, R., Clifford, D., Campbell, P.H., 2015. The Australian three-dimensional soil grid: Australia's contribution to the GlobalSoilMap project. *Soil Research* 53, 845–864.
- Wadoux, A.M.J.-C., 2019a. Sampling Design Optimization for Geostatistical Modelling and Prediction. Ph.D. thesis. Wageningen University & Research.
- Wadoux, A.M.J.-C., 2019. Using deep learning for multivariate mapping of soil with quantified uncertainty. *Geoderma* 351, 59–70.
- Wadoux, A.M.J.-C., Brus, D.J., Heuvelink, G.B.M., 2019a. Sampling design optimization for soil mapping with random forest. *Geoderma* 355, 113913.
- Wadoux, A.M.J.-C., Padarian, J., Minasny, B., 2019b. Multi-source data integration for soil mapping using deep learning. *Soil* 5, 107–119.
- 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.
- Wang, S., Zhuang, Q., Wang, Q., Jin, X., Han, C., 2017. Mapping stocks of soil organic carbon and soil total nitrogen in Liaoning Province of China. *Geoderma* 305, 250–263.
- Wang, B., Waters, C., Orgill, S., Gray, J., Cowie, A., Clark, A., Li Liu, D., 2018. High resolution mapping of soil organic carbon stocks using remote sensing variables in the semi-arid rangelands of eastern Australia. *Sci. Total Environ.* 630, 367–378.
- Were, K., Bui, D.T., Dick, Ø.B., Singh, B.R., 2015. A comparative assessment of support vector regression, artificial neural networks, and random forests for predicting and mapping soil organic carbon stocks across an Afrotropical landscape. *Ecol. Indic.* 52, 394–403.
- Wiesmeier, M., Barthold, F., Blank, B., Kögel-Knabner, I., 2011. Digital mapping of soil organic matter stocks using random forest modeling in a semi-arid steppe ecosystem. *Plant Soil* 340, 7–24.
- Wu, J., Teng, Y., Chen, H., Li, J., 2016. Machine-learning models for on-site estimation of background concentrations of arsenic in soils using soil formation factors. *J. Soils Sediments* 16, 1787–1797.
- Xu, S., An, X., Qiao, X., Zhu, L., Li, L., 2013. Multi-output least-squares support vector regression machines. *Pattern Recogn. Lett.* 34, 1078–1084.
- Yamazaki, D., Ikeshima, D., Tawatari, R., Yamaguchi, T., O'Loughlin, F., Neal, J.C., Sampson, C.C., Kanae, S., Bates, P.D., 2017. A high-accuracy map of global terrain elevations. *Geophys. Res. Lett.* 44, 5844–5853.
- Yang, R.-M., Zhang, G.-L., Liu, F., Lu, Y.-Y., Yang, F., Yang, F., Yang, M., Zhao, Y.-G., Li, D.-C., 2016. Comparison of boosted regression tree and random forest models for mapping topsoil organic carbon concentration in an alpine ecosystem. *Ecol. Indic.* 60, 870–878.
- Zeraatpisheh, M., Ayoubi, S., Jafari, A., Finke, P., 2017. Comparing the efficiency of digital and conventional soil mapping to predict soil types in a semi-arid region in Iran. *Geomorphology* 285, 186–204.
- Zeraatpisheh, M., Ayoubi, S., Jafari, A., Tajik, S., Finke, P., 2019. Digital mapping of soil properties using multiple machine learning in a semi-arid region, central Iran. *Geoderma* 338, 445–452.
- Zhu, M., Feng, Q., Zhang, M., Liu, W., Deo, R.C., Zhang, C., Yang, L., 2019. Soil organic carbon in semiarid alpine regions: the spatial distribution, stock estimation, and environmental controls. *J. Soils Sediments* 19, 1–15.