

A scalable framework for soil property mapping tested across a highly diverse tropical data-scarce region

Rodrigo de Q. Miranda^{a,1}, Rodolfo L.B. Nóbrega^{b,c,d,*,1} , Anne Verhoef^e,
 Estevão L.R. da Silva^{a,f}, Jadson F. da Silva^a, José C. de Araújo Filho^f ,
 Magna S.B. de Moura^{g,h} , Alexandre H.C. Barros^f, Alzira G.S.S. Souzaⁱ, Wanhong Yang^j ,
 Hui Shao^j, Raghavan Srinivasan^{k,l}, Feras Ziadat^m, Suzana M.G.L. Montenegroⁿ,
 Maria do S.B. Araújo^o , Josiclêda D. Galvêncio^a

^a PRODEMA, Universidade Federal de Pernambuco, Recife, Brazil

^b University of Bristol, School of Geographical Sciences, University Road, Bristol BS8 1SS, UK

^c Cabot Institute for the Environment, University of Bristol, Bristol, UK

^d Imperial College London, Georgina Mace Centre for the Living Planet, Department of Life Sciences, Silwood Park Campus, Buckhurst Road, Ascot SL5 7PY, UK

^e The University of Reading, Department of Geography and Environmental Science, Reading, UK

^f Brazilian Agricultural Research Corporation – Embrapa Soils, Recife, Brazil

^g Brazilian Agricultural Research Corporation – Embrapa Semi-arid Region, Petrolina, Brazil

^h Brazilian Agricultural Research Corporation - Embrapa Tropical Agroindustry, Fortaleza, Brazil

ⁱ Instituto Federal Baiano, Uruçuca, Bahia 45680-000, Brazil

^j University of Guelph, Department of Geography, Guelph, Ontario N1G 2W1, Canada

^k Spatial Sciences Laboratory, Texas A&M University, College Station, USA

^l Blackland Research and Extension Center, Agrilife Research, Temple, USA

^m Food and Agriculture Organization of the United Nations (FAO), Rome 00153, Italy

ⁿ Departamento de Engenharia Civil, Universidade Federal de Pernambuco, Recife, Brazil

^o Departamento de Ciências Geográficas, Universidade Federal de Pernambuco, Recife, Brazil

ARTICLE INFO

Dataset link: [Model outputs from the study "A scalable framework for soil property mapping tested across a highly diverse tropical data-scarce region"](#)

Keywords:

Digital Soil Mapping
 Tropical Soil Properties
 Gradient Boosting Model
 SLEEP
 Pernambuco
 Northeast region
 Brazil

ABSTRACT

Reliable soil property maps are essential for environmental modeling, yet conventional mapping methods remain costly and time-consuming. We developed a machine learning framework that integrates the Soil-Landscape Estimation and Evaluation Program (SLEEP) with gradient boosting to predict soil properties at regional scales and multiple depths. Our approach addresses multicollinearity through a recursive feature selection algorithm. We applied this framework to a tropical region characterized by a ~700-km longitudinal gradient of contrasting topography, climate, and vegetation (~98,000 km²; NE Brazil), where scarce soil physicochemical data limit environmental modeling. We used six topographical, ten climate, and two vegetation covariates, along with data from 223 soil profiles (~1 profile per 440 km²). Training and testing of our framework demonstrated strong spatial performance ($r^2 = 0.79$ – 0.98 and percent bias = -1.39 – 1.14 %). Topographic and climatic factors held greater weight than other variables in predicting soil layers, texture, and sum of bases. Moreover, we used our soil parameters combined with multiple pedotransfer functions (PTFs) to derive soil hydraulic properties. Our PTFs-derived estimates of hydraulic conductivity were considerably lower than high-resolution global predictions available for our study area due to differences in clay fraction and mineralogy. Therefore, we recommend the use of region-specific PTFs for hydraulic properties based on multi-covariate soil property maps. This cost-effective framework accurately integrates diverse environmental covariates, adapts to varying soil data availability, and scales across spatial resolutions, making it highly transferable to other data-scarce regions.

* Corresponding author at: University of Bristol, School of Geographical Sciences, University Road, Bristol BS8 1SS, UK.

E-mail address: r.nobrega@bristol.ac.uk (R.L.B. Nóbrega).

¹ These authors contributed equally to this work.

1. Introduction

Soils are a key component in many landscape models that focus on providing solutions to global environmental issues such as food and water scarcity, unsustainable energy production, and biodiversity losses (Bouma and McBratney, 2013). For a more comprehensive understanding of the role of soils in addressing these global challenges, as well as their interactions with other environmental factors, it is necessary to map the spatial distribution of soil properties robustly. Soil mapping is complex and highly resource-intensive (Li and Heap, 2014; Mendonça-Santos and dos Santos, 2006), and the majority of the existing maps were produced using conventional soil survey protocols (Hartemink et al., 2012), which remains the primary approach to capture soil spatial variability. However, this surveying approach has been criticized for being heuristically dependent on the practical knowledge of pedologists, and for deriving interpretations using sometimes insufficient or incomplete datasets (Scully et al., 2003).

Digital Soil Mapping (DSM) is a quantitative approach to mapping soil properties using statistical relationships between soil observations and environmental variables. It was formalized with the SCORPAN model, which considers factors such as soil properties, climate, vegetation, topography, and spatial position to guide the selection of covariates in DSM (McBratney et al., 2003) to produce models capable of interpolating and extrapolating data with high resolution (Scully et al., 2003). DSM reduces survey costs and improves access to soil data by leveraging advances in remote sensing, geospatial analysis, and machine learning (ML) (Kempen et al., 2012; Lagacherie and McBratney, 2006). It has been widely applied to map soil attributes such as texture, organic carbon, and pH at regional to continental scales (e.g., Ballabio et al., 2016; Guevara et al., 2018).

DSM has been widely used across the world to reduce soil mapping costs over large areas (e.g., Tóth et al., 2017; Guevara et al., 2018; Padarian et al., 2017; Teng et al., 2018). The methodological core of DSM includes mathematical models capable of performing both interpolations and extrapolations of soil properties across multiple scales (Barros et al., 2013; Laurent et al., 2017; Saxton and Rawls, 2006; Tomasella et al., 2000; Wang et al., 2018; Zeraatpisheh et al., 2019). These models can predict the distribution of a given soil property horizontally, e.g., over the topsoil of a landscape, or vertically, i.e., along soil profiles. In soil science, spatial extrapolations are usually made by (i) applying a conceptual model to the survey area to simulate the distribution of soil patches (Scully et al., 2003), (ii) using geostatistical interpolations (Li and Heap, 2014), (iii) delimiting geographical subdivisions where environmental processes follow a relatively homogeneous pattern, such as the facets, described by Ziadat et al. (2015), or (iv) by applying pedotransfer functions (PTFs) to basic properties available for each soil location. PTFs are predictive statistical models, typically regression equations, that use basic soil information to estimate soil properties that are costly to measure, such as water retention characteristics and bulk density (Barros and de Jong van Lier, 2014).

There is an ever-growing need for soil data, e.g., for research and applications related to environmental solutions, especially in the tropics where soil data are scarce and soils exhibit the highest global diversity (Minasny and Hartemink, 2011; Scharlemann et al., 2014; Orgiazzi et al., 2016). The hydro-thermal behavior of tropical soils is quite different compared to temperate soils, often due to their distinct mineralogies and soil-forming processes (Ito and Wagai, 2017; Nóbrega et al., 2020). In Brazil, various polynomial PTFs have been calibrated at both national (Tomasella et al., 2000) and sub-national scales (Barros et al., 2013; Oliveira et al., 2002) for estimating soil properties such as hydraulic conductivity, water retention characteristics and bulk density. However, high uncertainties are expected when conducting both horizontal and vertical soil properties extrapolations, especially for vertical extrapolations because data on soil profiles across extensive terrain extents are rarely available (Yost and Hartemink, 2020).

ML techniques have been increasingly applied as an approach to

circumvent issues typical of conventional soil mapping methods and those issues that are due to the complexity caused by modeling the soil with ever-increasing amounts of information stored in databases on soil parameters and covariates (Wadoux et al., 2020). If trained properly, ML techniques allow for more accurate predictions of soil parameters, whereas other approaches with underlying assumptions on statistical distributions may not be applicable or even fail to produce sensible values (Taghizadeh-Mehrjardi et al., 2016). However, many ML studies used for soil mapping do not predict soil properties at different depths (e.g., van der Westhuizen et al., 2023; Bao et al., 2024; Hateffard et al., 2024; Qu et al., 2024; Sun et al., 2024). When depth predictions are made, it is common to follow standardized output specifications, such as those defined by GlobalSoilMap (Ballabio et al., 2016; Rahmati et al., 2018), which uses six fixed depth intervals within the 0–200 cm soil depth. However, this approach is inconsistent with established soil classification systems, consequently limiting the pedological interpretation of the results (Wadoux et al., 2020).

ML approaches in digital soil mapping (DSM) offer improved estimates of soil parameters, with the accuracy strongly influenced by the choice of soil maps and pedotransfer functions (PTFs) (Montzka et al., 2017). For instance, Gupta et al. (2021) demonstrated that a ML approach involving various soil and environmental covariates improved predictions of saturated hydraulic conductivity compared to traditional PTF-based methods. They generated a final dataset with a spatial resolution of 1 km by using a random forest algorithm and data from 821 sites distributed around the world; however, with only ~12 % of these data from the tropics. Indeed, soil maps for the tropics often exhibit a coarse exaggeration of soil properties. This occurs because the common statistical techniques applied to perform extrapolations are heavily dependent on how dense the collection of soil profiles is, and this is generally sparse due to financial and time limitations.

The possibility of using high-resolution environmental covariates offers new opportunities for adding local information into soil property modeling. In hydrology, for example, the Soil and Water Assessment Tool (SWAT; Arnold et al., 1998) employs the Soil–Landscape Estimation and Evaluation Program (SLEEP; Ziadat et al., 2015), which goes beyond a simple point-by-point approach by aggregating pixels into more homogeneous areas according to topographic features. This subdivision reduces noise from abrupt terrain changes and captures the influence of landscape context on soil formation more effectively. However, relying on these covariates alone, i.e., without ML, often involves simple regressions that struggle to account for both gradual and abrupt soil variability (Wadoux et al., 2020). The use of ML techniques, such as random forest (RF) or gradient boosting models (GBMs), has improved the prediction accuracy of soil organic matter and total N when compared to geostatistical methods, and further gains have been achieved when these approaches are combined (Auzas et al., 2024; Nozari et al., 2024; Tziachris et al., 2019). While geostatistics uses spatial autocorrelation to refine local estimates, ML captures complex interactions among environmental variables, thereby improving overall model robustness and predictive performance.

In this study, we address the growing need for improved soil models that capture the spatial variability of physical and chemical properties in the tropics by developing a bespoke machine learning framework. Applied across a ~700-km longitudinal gradient in Brazil with contrasting topography, climate, and vegetation, our approach targets a long-standing gap in tropical soil observations within global soil databases. We hypothesize that our framework can accurately capture both vertical and horizontal variability in soil properties in a large tropical region with highly contrasting environmental conditions and land use. It combines SLEEP with calibrated GBMs to produce high-resolution (30 m) predictions across multiple depths. The framework was developed to enable the generation of soil maps that support: (1) assimilation of legacy soil data in their native format; (2) fine-scale prediction of key soil properties; (3) identification of environmental drivers for each pedological feature, and; (4) generation of soil datasets for

environmental modeling.

2. Materials and methods

2.1. Methodology workflow

We developed and applied our modeling framework by integrating SLEEP and a calibrated GBM, which we tested for a 700-km longitudinal gradient in Northeast Brazil (see [Section 2.2](#)). The stage-wise additive trees of GBMs can capture higher-order interactions between soil properties and climate, vegetation, and topographic predictors without the need for additional feature engineering (e.g., transformations). GBMs also adapt to depth-dependent heteroscedasticity while maintaining linear scalability for 30 m resolution predictions across large datasets, such as the 100 million pixels used in this study. Our methodology comprises a three-step process that starts with the collection and pre-processing of six topographical, ten climate, and two vegetation parameters acquired from different data sources ranging from remotely sensed datasets to meteorological stations (see [Section 2.3](#)). These independent variables are correlated with soil physical and chemical properties, referred to as basic soil properties, as described in [Table 1](#) and [Section 2.4](#), to allow for their subsequent horizontal and vertical predictions.

We used SLEEP to create a non-distributed grid formed by facets, which, in this study, are treated as the smallest spatial units representing homogeneous conditions where soil formation factors may produce similar soil types. To define these facets, SLEEP first creates preliminary versions of these facets by delineating watersheds. Each watershed is divided into multiple catchments, and then the facets are defined by the division of the catchments into two parts, i.e., each side of their main drainage stream ([Ziadat et al., 2015](#)). The size of the catchments is determined by a user-defined threshold assigned during stream definition. The smaller this threshold, the denser the stream network, resulting in a greater number of delineated catchments and facets. Once the facets are created, SLEEP aggregates them based on their slope similarity in a process called facet classification, which ultimately creates contiguous patches, which are clusters of facets that share similar slope characteristics and are treated as unified mapping units. The patches allow SLEEP to reduce the number of facets by grouping them into a single mapping unit. This approach reduces the processing time when working with large areas and avoids the ‘salt-and-pepper’ noise in the mapping process. Next, we estimated the ten basic soil properties (indicated in [Table 1](#)) in each patch at multiple depths by calibrating one model for each basic soil property using ML instead of traditional SLEEP multiple regressions because they can capture a wider range of data distributions (see [Section 2.5](#)). The calibration mechanism is composed of a recursive feature selector and a randomized searcher, which were configured to perform a 2-fold cross-validation (see [Section 2.6](#)). At the end of this step, all patches are turned into virtual soil profiles, i.e., simulated soil patches with their own depth-dependent simulated physical and chemical properties, and the uncertainty was calculated for each estimated soil property (see [Section 2.7](#)). Finally, in the third step, we used the dataset composed of virtual profiles to generate PTF-estimated soil parameters (see [Section 2.8](#)).

2.2. Study area

The study area is in Northeast Brazil; it covers an area of approx. 98,000 km², and closely follows the domain of the state of Pernambuco ([Fig. 1](#)). This region exhibits a longitudinal gradient of contrasting topography, climate and vegetation. The elevation ranges from approx. 0 to over 1150 m a.s.l. in a variable gradient from East to West. This region is influenced by three meteorological phenomena, namely Frontal Systems (FS), Upper Tropospheric Cyclonic Vortices (UTCV), and the Intertropical Convergence Zone (ITC) ([Salgueiro et al., 2016](#)). There are three predominant climate types (Köppen’s classification) in

Table 1

Summary of variables and parameters with their corresponding descriptions and units.

Variable	Type	Description	Unit
AAT	T	Prefix used to denote accumulated variables	-
ASPECT	T	Downslope direction at each cell	°
CTI	T	Compound Topographic Index	-
CURV	T	Surface curvature at each cell	-
DEM	T	Digital elevation model	m
PCTSLP	T	Surface slope at each cell	%
LST	V	Land surface temperature	K
NDVI	V	Normalized difference vegetation index	-
RHAV	C	Mean air relative humidity	fraction (0–1)
PCPMM	C	Mean total monthly precipitation	mm
PCPSKW	C	Skew coefficient for daily precipitation in month	mm
PCPSTD	C	Standard deviation for daily precipitation in month	mm
SOLARAV	C	Mean daily solar radiation for month	MJ m ⁻² day ⁻¹
TMPMN	C	Mean daily minimum air temperature	°C
TMPMX	C	Mean daily maximum air temperature	°C
TMPSTDMN	C	Standard deviation for daily minimum air temperature	°C
TMPSTDMX	C	Standard deviation for daily maximum air temperature	°C
WNDVAV	C	Mean daily wind speed in month	m s ⁻¹
CS	B	Coarse sand content	%
FS	B	Fine sand content	%
L_MAX	B	Number of soil layers	-
SB	B	Sum of bases (Ca ²⁺ , Mg ²⁺ , K ⁺ and Na ⁺)	cmol _c kg ⁻¹
SOL_CBN	B	Organic carbon content	%
SOL_CLAY	B	Clay content	%
SOL_ROCK	B	Rock fragments content	%
SOL_SAND	B	Sand content	%
SOL_SILT	B	Silt content	%
SOL_Z	B	Depth from soil surface to bottom of the soil layer	mm
R _v	P	Volume fraction of gravel	cm ³ cm ⁻³
R _w	P	Weight fraction of gravel	g g ⁻¹
θ ₁₅₀₀	P	Water content at −1500 kPa	m ³ m ⁻³
θ ₃₃	P	Water content at −33 kPa	m ³ m ⁻³
θ _s	P	Saturated water content	m ³ m ⁻³
θ _r	P	Residual water content	m ³ m ⁻³
ρ _N	P	Normal density	g cm ⁻³
ρ _R	P	Gravel density	g cm ⁻³
OM	P	Organic matter	%
SN1	P	Non-sand content	fraction
SOL_AWC	P	Available water capacity of the soil layer	mm mm ⁻¹
SOL_BD	P	Moist bulk soil density	g cm ⁻³
SOL_K	P	Saturated hydraulic conductivity	mm hr ⁻¹
USLE_K	P	USLE equation soil erodibility (K) factor	-
ψ	P	Matric potential	kPa
α	P	Parameter of van Genuchten (1980) usually expressing inverse length (pressure head)	m ⁻¹
n and m	P	Shape-fitting parameters of van Genuchten (1980)	-

In column 2: T = topography, V = vegetation, C = climate, B = basic property, and P = pedotransfer function parameter.

the study area: hot semi-arid (steppe) climate (BSh; 61.4 % of the area), tropical with dry summer (As; 32.7 %) and tropical monsoon (Am; 4.9 %); the remaining 1 % is composed of areas with a tropical climate with dry winter (Aw; 0.1 %), and humid subtropical with dry winter and hot summer (Cwa; 0.3 %), temperate summer (Cwb; 0.3 %), or dry and hot summer (Csa; 0.3 %) ([Alvares et al., 2013](#)). Precipitation has a high spatial variability ([Souza et al., 2021](#)) with the annual mean precipitation rates reaching approx. 2000 mm in the East and decreasing westwards to less than 400 mm. As for the vegetation, near the coast, the predominant land-uses are Atlantic rain forest and rainfed croplands (a mosaic of sugarcane plantations and fruticulture) ([Souza Jr et al., 2020](#)). Approaching the middle transition, around longitude 36° 47', high altitudes contribute to microclimatic conditions that favor rainfed corn and

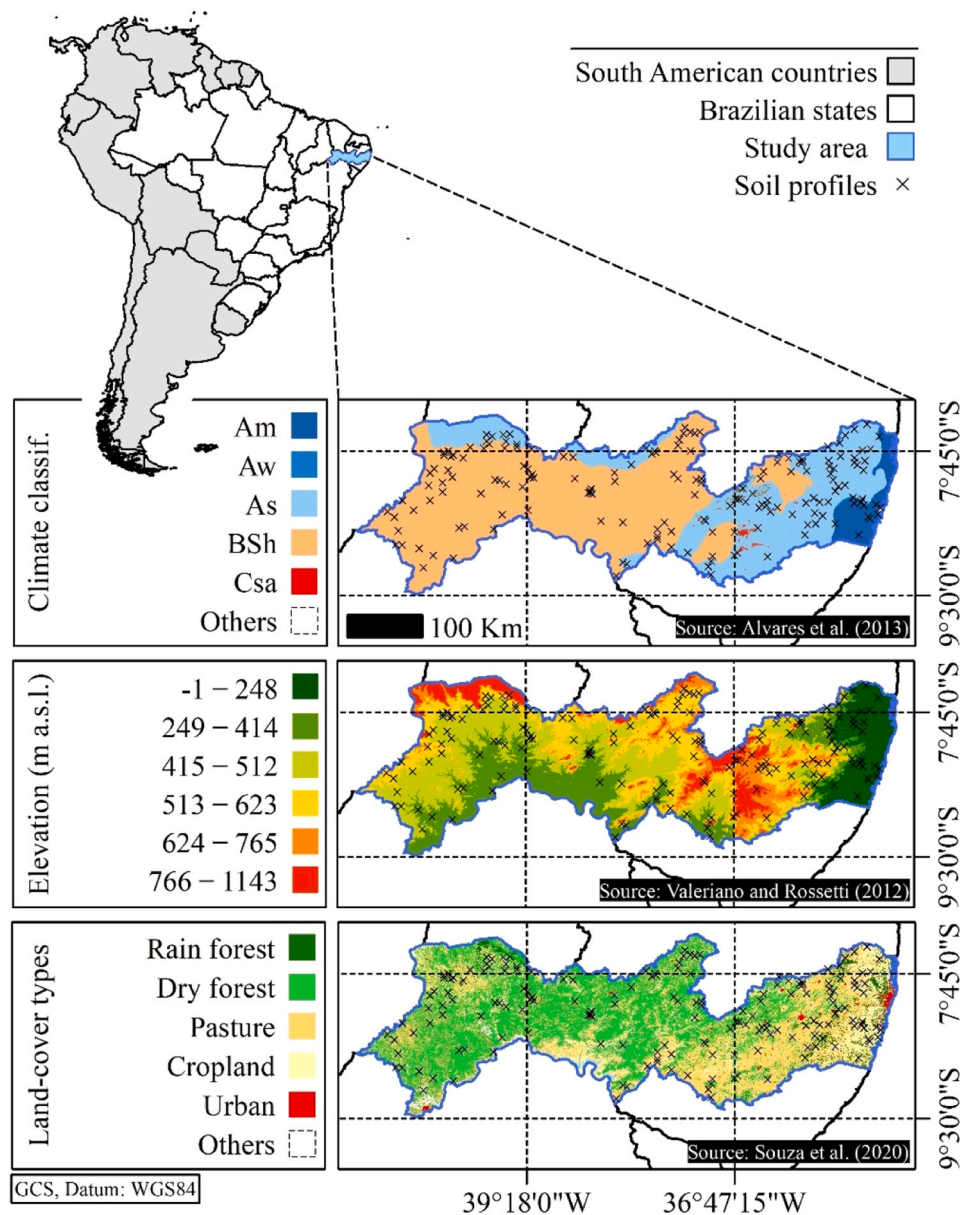


Fig. 1. Spatial distribution of the surveyed soil profiles across a longitudinal gradient of environmental conditions over the study area.

bean cultivation, and mixed natural vegetation formations. With rainfall decreasing, the vegetation changes to a seasonally dry tropical forest, i. e., the Brazilian Caatinga. Pastures become a common land-use activity, and the soil gets shallower and rocky (Souza Jr et al., 2020). According to the Brazilian and FAO system of soil classification, the dominant soils are, respectively, *Argissolos*, i.e., Acrisols and Lixisols (25 % of the area), *Neossolos*, i.e., Leptosols, Arenosols, Regosols, or Fluvisols (32 %) and *Planossolos*, i.e., Planosols and Solonetz (16 %), *Latosolos*, i.e., Ferralsols (9 %) and *Luvissolos*, i.e., Luvisols (9 %) (Araújo Filho et al., 2014). The geology maps for the state of Pernambuco show predominantly (90 %) pre-Cambrian rocks belonging to the São Francisco Craton and the Borborema Province, and the remaining area is mainly composed of Paleomesozoic sedimentary basins and Mesocenezic coastal basins (Torres and Pfaltzgraff, 2014).

2.3. Input data collection

We selected the input parameters based on their widely known role on soil formation. **Elevation data:** we collected data from the

TOPODATA database (<http://www.dsr.inpe.br/topodata>), which is a bias-corrected version of the data produced by the NASA SRTM (Shuttle Radar Topography Mission) for the Brazilian territory made by the National Institute for Spatial Research (INPE) at 1 arc-second (approx. 30 m) (de de Morisson Valeriano and de Fátima Rossetti, 2012).

Soil data: we digitized georeferenced data regarding morphological (number and depth of soil horizons), physical (particle size distribution), and chemical (Ca^{2+} , Mg^{2+} , K^{+} , Na^{+} and C) soil properties, acquired from the ZAPE (Agroecological Zoning of the state of Pernambuco) project of the Brazilian Agricultural Research Corporation (EMBRAPA) (Silva et al., 2001). This legacy soil database comprises 223 soil profiles distributed over the study area (Fig. 1).

Meteorological data: we obtained data for air temperature ($^{\circ}\text{C}$), air relative humidity (%), solar radiation ($\text{MJ m}^{-2} \text{ day}^{-1}$), wind speed (m s^{-1}), and precipitation (mm) from the 1961–2016 period through two open-access databases: daily precipitation data from the Water and Climate Agency of Pernambuco (APAC; <http://www.apac.pe.gov.br/meteorologia/monitoramento-pluvio.php>), and the other meteorological parameters from the National Water Agency of Brazil (ANA; <https://www.ana.gov.br/>).

<http://www.snirh.gov.br/hidroweb/>). The preprocessing of these data is detailed in the [Supplementary Material](#) (Section 1 of the [Supplementary Material](#)).

Remotely sensed data: we obtained data regarding NDVI (Normalized Difference Vegetation Index) from MOD13A3 (monthly composition and 1 km spatial resolution) (Didan, 2015), and LST (Land Surface Temperature) from MOD11A2 (8-day composition and 1 km spatial resolution) (Wan et al., 2015) from <https://earthdata.nasa.gov/> (Greenbelt, 2019).

2.4. Soil survey data description

Our soil dataset includes the total number of soil horizons (L_MAX), but for modeling purposes in this study we will refer to it as the number of soil layers since we did not validate the model's efficacy in distinguishing horizons through further field experiments. Thus, a soil layer here refers to a vertical depth interval used to represent distinct soil properties within the soil profile. The database also contains each soil layer's depth from the land surface (SOL_Z; mm), soil clay content (≤ 0.002 mm; SOL_CLAY; %), silt (> 0.002 and ≤ 0.05 mm; SOL_SILT; %),

sand (> 0.05 and ≤ 2 mm; SOL_SAND; %), rock fragments (> 2 mm; SOL_ROCK; %), organic carbon (SOL_CBN; %), and sum of bases (sum of Ca^{2+} , Mg^{2+} , K^{+} and Na^{+} ; SB; $\text{cmol}_c \text{ kg}^{-1}$). In this study, we define the rock parameter as the proportion of rock fragments greater than 2 mm (ABNT, 1995; FAO, 2006). The sand fraction was divided into fine (> 0.05 and ≤ 0.2 mm; FS) and coarse sand (> 0.2 and ≤ 2 mm; CS) (Table 1). All particle classification followed the Brazilian technical standards described in ABNT (1995), and physical and chemical analyses were performed as described in Embrapa (1997).

Soil profiles exhibit an average depth of 1228 ± 613 mm, ranging from 120 to 2550 mm. The number of soil layers varies from one to seven. Rock fragments (> 2 mm) exhibit 4.4 ± 11 % of total content. If we only consider particles ≤ 2 mm, the average soil texture has the following composition: sand (55 ± 19 %), clay (27 ± 14 %), and silt (18 ± 9 %) (Fig. S1 in the [Supplementary Material](#)).

2.5. Inputs for the preprocessing workflow

The core of our modeling framework combines SLEEP and a calibrated GBM. Soil data were modeled in SLEEP by creating facets (see

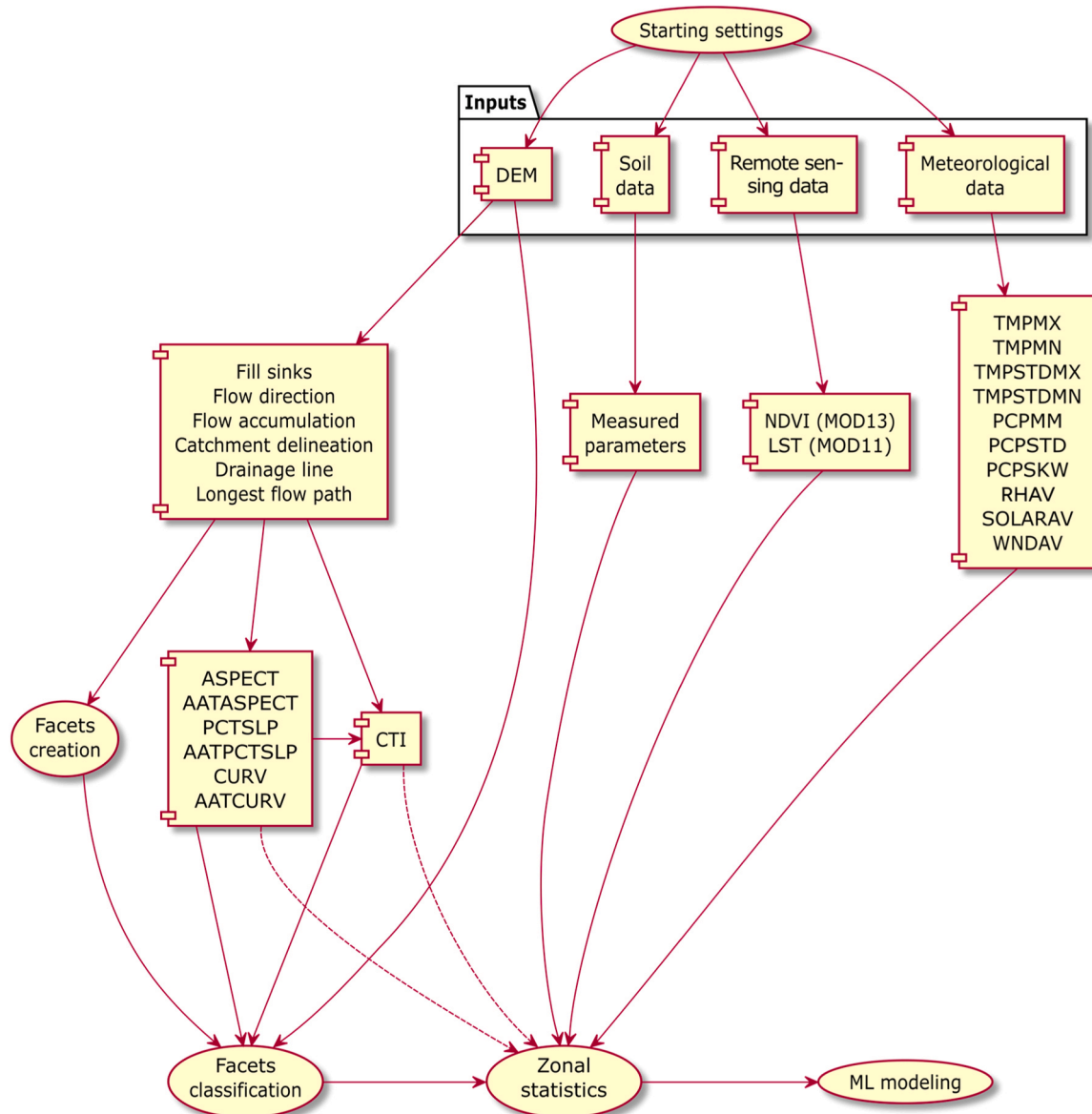


Fig. 2. Processing scheme of the integration of the SLEEP algorithm and the Gradient Boosting Models. The description of the parameters can be found in [Table 1](#).

Section 2.1), for which basic soil properties, i.e., L_MAX, SOL_Z, SOL_CLAY, SOL_SILT, SOL_SAND, CS, FS, SOL_ROCK, SOL_CBN, and SB, were calculated.

SLEEP requires three inputs: (i) a digital elevation model (DEM), (ii) a shapefile containing the data observed for each soil profile, and (iii) the auxiliary data including meteorological and vegetation data in raster format (Fig. 2). In this algorithm, we extracted the drainage network following Tarboton et al. (1991) by setting the size of the catchments to 0.001 % of the total study area, i.e., on average 1803 pixels per catchment, which was obtained based on a visual evaluation of different thresholds with a focus on providing a balance between satisfactory spatial resolution and processing efficiency. We aggregated the facets based on their slope similarity using the clustering technique IsoCluster (Richards, 2013) to create patches.

Finally, we modified the way the basic properties were modeled, replacing the original SLEEP algorithm's simple multiple linear regression with GBMs. GBM is an ensemble learner that consists of a set of decision trees composed of weak predictive models (WPM) often prone to overfitting, but, when combined, produce highly accurate outputs (Friedman, 2001). Each of these trees is a rule-based system, whose terminal nodes can either be a WPM, i.e., leaf node, or an if-then-else rule, i.e., regular node, applied to an input variable. The trees are created through an iterative sequence of improvements of WPMs using boosting, while simultaneously optimizing, via minimization of a loss function using gradient-based optimization (Natekin and Knoll, 2013).

For GBM processing, two datasets were produced: (i) one composed of only the information from the patches that overlie the observed data for each profile to be used as the dataset for fitting, and (ii) consisting of all available input information for every patch in the study area to be used as the dataset for prediction. The dataset for fitting was split using

the Holdout method at 20 %, e.g., Whitney (1971), creating two sub-datasets, where 80 % of the records were used for model calibration (training dataset), and the remaining 20 % for model verification (verification dataset) (Fig. S2 in the Supplementary Material).

The sampling technique used in this process is a variation of the k-fold cross-validation (Wong, 2015), which ensures stratified folds with a balanced distribution of each target class. For continuous dependent variables without predefined classes, a quantile-based discretization function (*qcut* function in Python; The pandas development team, 2024) was applied to discretize these variables into equal-sized groups based on sample quantiles, allowing the entire data distribution to be sampled.

The GBMs had four basic parameters derived from the DEM (Table 1) as input features, namely the downslope direction (ASPECT), the Compound Topographic Index (CTI), the surface curvature (CURV) and slope (PCTSLP), as well as 12 auxiliary data series from remote sensing (NDVI, LST) and meteorological stations (see Table 1). As targets, they had eight basic soil properties (labeled as Type B in Table 1, see 'ML outputs' in the upper half of Fig. 3). GBM was used as a multiclass classifier to simulate the number of soil layers, i.e., L_MAX, and a regressor for the other targets. In the GMB model, SOL_ROCK was not directly estimated but was computed as a residual component of sand, silt and clay, which were not rescaled to sum to 100 % as inputs. Coarse sand (CS) and fine sand (FS) were normalized to sum up to 100 %.

2.6. Model calibration and validation

To calibrate the hyperparameters, we submitted all our GBMs to a Recursive Feature Selector (RFS; Guyon et al., 2002) followed by randomized 2-fold cross-validation to optimize hyperparameter selection. The RFS here is an input feature selection algorithm that fits a model and

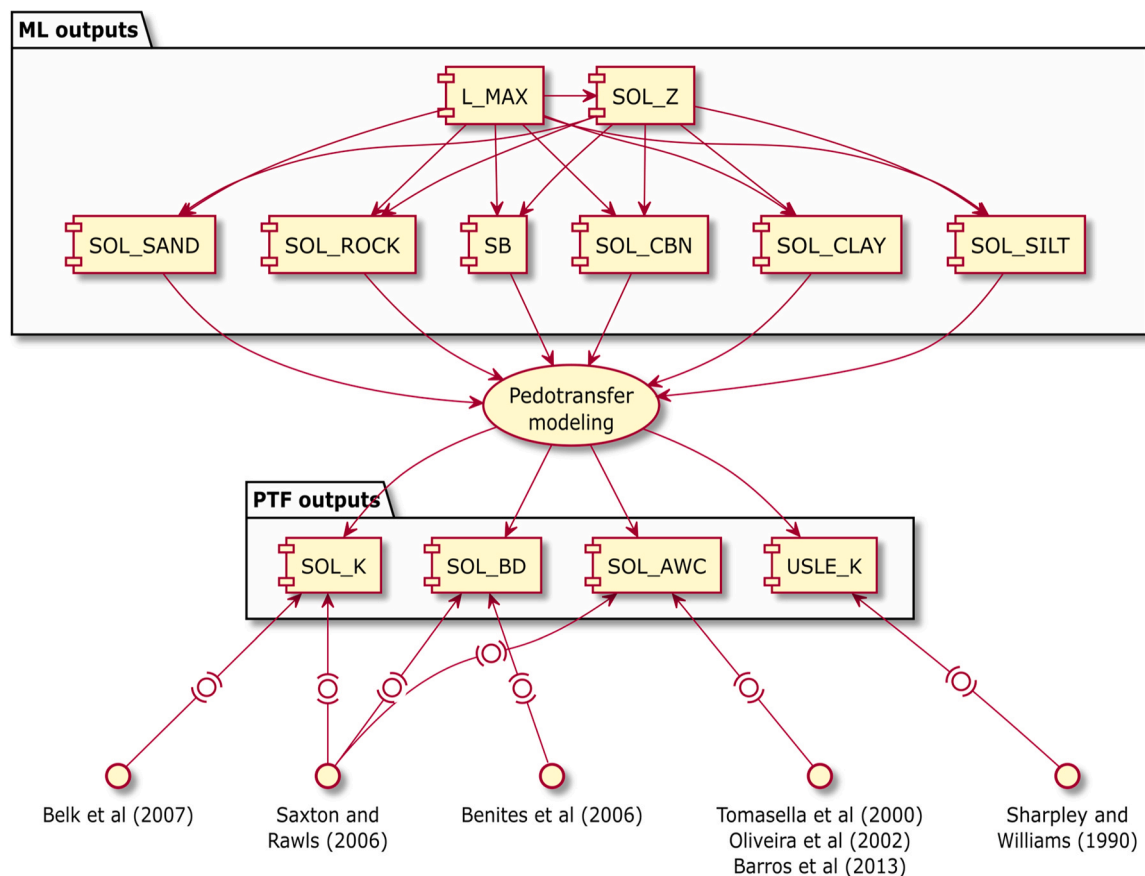


Fig. 3. Processing workflow of all model outputs. The top half of this figure explains the machine learning processing of the basic soil characteristics, whereas the bottom half summarizes the PTF-derived products. The description of the parameters can be found in Table 1.

eliminates the weakest ranked inputs recursively, considering each iteration a smaller set of features until the best combination is found. We determined the optimal cross-validation splitting strategy for our model's calibration by performing a small-scale test using all data and one variable, i.e., L_MAX, with different fractions of data splits for validation (10, 15, 20, 25, and 30 %) combined in a factorial design with different levels of data slicing for cross-validation (2, 3, 4, 5, and 10 folds). All tested data splits, and cross-validation configurations for both RFS and hyperparameters calibration resulted in accuracy between 0.96 and 0.97, with 20 % data split and 2-fold cross-validation yielding an accuracy of 0.97 (Eq. 1). Therefore, we used the 2-fold calibration to reduce computing demand. This means that 50 % of the calibration data were used to test each hyperparameter combination's impact. With this configuration, the full simulation ran for 232 h (~10 days) on a super-computer with 120 cores distributed across 10 Intel i7 processors (3.2–3.33 GHz), 80 GB DDR3 RAM (1333 MHz), 10 TB HDD storage, and 20 Gigabit network cards. The modeling algorithm is freely available at GitHub and is compatible with Python 2.7.15 and 3.6.9. For details, see Miranda et al. (2022).

The performance indices used in all calibrations were the accuracy (Eq. 1) for the classifier, i.e., for L_MAX, and the coefficient of determination (r^2) (Eq. 2) for the regressors. For model verification, the most efficient models were evaluated using the testing dataset, and the same performance indices plus the Root Mean Square Error (RMSE) (Eq. 3) and Percent Bias (PBIAS) (Eq. 4) were applied. This final verification allowed us to evaluate the potential of the best models to perform extrapolations.

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + FP + FN + TN)} \quad (1)$$

$$r^2 = \frac{\sum (obs - \overline{obs}) \times (sim - \overline{sim})}{\sqrt{\sum (obs - \overline{obs})^2} \times \sqrt{\sum (sim - \overline{sim})^2}} \quad (2)$$

$$RMSE = \sqrt{\frac{\sum (obs - sim)^2}{n}} \quad (3)$$

$$PBIAS = \frac{\sum (obs - sim)}{\sum (obs)} \times 100 \quad (4)$$

TP, FP, FN, and TN in Eq. 2 represent True Positives, False Positives, False Negatives, and True Negatives, respectively, in a contingency table. The variable *obs* in Eqs. 2–4 refers to the observed parameter value for a given soil layer, while *sim* represents the simulated value, with the overbar indicating their average values.

In this study, the classification problem involves distinguishing between soil properties based on observed and simulated values. However, due to an imbalance in class representation, where certain soil conditions, e.g., a specific texture class or rock presence are underrepresented, the model may become biased toward the dominant class, leading to poor detection of minority cases. To mitigate this issue, we applied the Synthetic Minority Oversampling Technique (SMOTE) to balance the class distribution. SMOTE generates synthetic samples for the under-represented soil properties, ensuring they contribute more effectively to the model training process. This technique promotes balanced learning and improves the detection of minority soil conditions. Details of this technique can be found in Chawla et al. (2002). To calibrate the hyperparameters, we created a set of possible values for each parameter. Details for this procedure can be found in Section 3 of the Supplementary Material. The calibrated models were applied to predict basic properties for each patch, creating 64,415 virtual soil profiles. The entire predicted dataset was converted to a raster format, and each raster is a different soil attribute. All outputs are available from Miranda et al. (2025).

2.7. Sensitivity and uncertainty analysis

The model sensitivity to input data was calculated as the importance, i.e., a weighted factor of each selected property for the most accurate GBMs. The importance (w) ranges from 0 to 1, where 1 reflects the highest weight a given input can receive in a model, and 0 the lowest. The sum of all weights is 1 for each model. More specifically, w values reflect indirectly how much the performance metric changes every time a given input is used to split a node in the whole model (Natekin and Knoll, 2013).

For the uncertainty analysis of the modeled variables, the selected inputs for each model and patch used in the predictions were classified into two categories (e), i.e., whether they extrapolated the calibration range of values (1) or not (0), as summarized in the following equation:

$$u_f = \sum_{i=0} (e_i \times w_i), \quad (5)$$

where u_f is the uncertainty of each model; patch, e_i is the binary category that reflects the extrapolation and w_i is its importance in the model (weight) of a given selected input i . As u_f gets close to 1, extrapolation is greater indicating higher associated uncertainty. The opposite occurs when it approaches 0, which means that all inputs used for a given prediction were in the range of values used for calibration.

2.8. Application and comparison of pedotransfer functions

All data from the virtual soil profiles were submitted to a series of pre-established PTFs (see bottom-half of Fig. 5) to generate four soil properties: SOL_K (saturated hydraulic conductivity; mm hr⁻¹), SOL_BD (moist bulk density; g cm⁻³), SOL_AWC (available water capacity; mm mm⁻¹), and USLE_K (factor K from the USLE equation; unitless). SOL_K was modeled using the equations described in Saxton and Rawls (2006) and Belk et al. (2007), and USLE_K using Sharpley et al. (1993) (equation groups S1–S3 described in Table S2 in the Supplementary Material). SOL_AWC was calculated with the equations from Saxton and Rawls (2006), Tomasella et al. (2000), Oliveira et al. (2002) and Barros et al. (2013) as described in equation groups S4–S9 in Table S3 in the Supplementary Material. Saxton and Rawls (2006) produced PTFs using a soil dataset from extensive soil sampling across the entire United States. Tomasella et al. (2000) used a similar database for Brazil, while Barros et al. (2013) used data for the Northeast region of Brazil only. Finally, Oliveira et al. (2002) created PTFs with data that originated strictly from the state of Pernambuco.

All SOL_AWC models require SOL_BD as an input. Thus, SOL_BD derived from Saxton and Rawls (2006) was coupled with their corresponding SOL_AWC model, while SOL_BD from Benites et al. (2007) was used in the models of Tomasella et al. (2000), Oliveira et al. (2002) and Barros et al. (2013). To distinguish between PTF sources, subscripts were assigned to variables as follows: BK for Belk et al. (2007), BR for Barros et al. (2013), OL for Oliveira et al. (2002), SR for Saxton and Rawls (2006), and TM for Tomasella et al. (2000). Additionally, SOL_K_{SR/BR} and SOL_K_{SR/TM} refer to SOL_K estimated using Saxton and Rawls (2006)'s PTF, where θ_5 , θ_{33} , and θ_{1500} were derived from Barros et al. (2013) and Tomasella et al. (2000), respectively.

We compared our SOL_K results derived from Saxton and Rawls (2006) to the dataset generated by Gupta et al. (2021), who generated high-resolution, i.e., 1 km, global SOL_K values using a ML framework. We chose Saxton and Rawls (2006) because it is a widely used PTF. That way we avoided bias caused by comparing Gupta et al. (2021)'s results to SOL_K estimates derived from PTFs that were specific to our area of study, such as from Barros et al. (2013) and Oliveira et al. (2002). Nevertheless, we made available all results of all PTFs and their combinations, e.g., using the SOL_K model from Saxton and Rawls (2006) using the field capacity model from Barros et al. (2013), at <https://zenodo.org/deposit/5918544> (Miranda et al., 2025). To enable the

SOL_K comparison, we cropped the dataset from Gupta et al. (2021) to our spatial extent and resampled our dataset to Gupta et al. (2021)'s spatial resolution. We also compared the clay fraction obtained in this study with the one used by Gupta et al. (2021), provided by Hengl (2018), because this is an important component of many SOL_K models, including the one by Saxton and Rawls (2006) (Table S2 in the Supplementary Material). We calculated mean SOL_K and clay fraction as a weighted mean for each grid cell for Gupta et al. (2021)'s SOL_K and respective soil depth since our SOL_K values are representative for the entire soil layer. For the SOL_K dataset from Gupta et al. (2021) and clay fraction from Hengl (2018), we calculated the vertical value mean using the trapezoidal rule suggested by Hengl et al. (2017). This approach was chosen because the SOL_K values were predicted at discrete soil depths rather than being representative of the midpoint of the predefined depth intervals.

3. Results and discussion

3.1. Model performance

The spatial modeling produced 64,415 patches with an average area of $1.35 \pm 4.54 \text{ km}^2$, and an average density of 0.75 patches per km^2 . Each one of these was considered as a virtual soil profile for which GBM outputs were calculated. In this study, the models demonstrated a consistent ability to perform such extrapolations, as the performance of the models during the verification was similar to that found by the calibration algorithm (Table 2). The r^2 and PBIAS values varied from 0.79 to 0.98, and from -1.39 to 1.14 , respectively. Among all models for the prediction of percentages of each soil parameter, the lowest r^2 value was found for the modeled SOL_SILT at 0.79 (Table 2). We believe that the large number of predictors, each with similar importance, for the SOL_SILT model (Table 3) may have caused prediction redundancies and probably degraded the model strength by increasing its variance, even though we applied a RFS algorithm for feature selection.

When comparing the simulated and observed reference datasets (Table S4 in the Supplementary Material), some differences are expected because the soil survey data used as observed dataset (Section 2.4) was not systematically sampled. Therefore, there will be locations with simulated interpolated soil properties exhibiting values that exceed those in the observed dataset. The largest relative differences between simulated and observed values were for SOL_ROCK (44.4 %), SB (53.1 %), CS (103.3 %), and FS (31.9 %). Despite the lack of systematic sampling, these differences would be expected to be modest, as the observed dataset covers the entire study area and diverse environments (Fig. 1). We attribute these large differences in SOL_ROCK to the fact that this parameter was calculated as the residual of all soil separates (see Fig. S4 in the Supplementary Material). That is, it was the only parameter that was not directly modeled from independent covariates. As for CS and FS, they were directly modeled but had to be resampled to sum to 100 %. Rather than applying the same approach to texture parameters, we opted to sacrifice SOL_ROCK's prediction accuracy. Its

spatial variance produced a high number of zeros (38.5 % of total values) compared to other parameters (<0.01 %), resulting in insufficient variance for accurate modeling. Although 21.98 % of SB predictions ranged between 0.1 and $3.84 \text{ cmol}_c \text{ kg}^{-1}$ and no zeros, they exhibited a higher concentration near zero, similar to SOL_ROCK. Finally, 51.49 % of the 135,934 virtual profiles exhibited some degree of uncertainty. Most uncertainty values were below 15 %, while the highest values (50–60 %) were observed for L_MAX, SOL_SAND, and SB (Fig. 4). We would like to highlight that our approach to estimate uncertainty relies on identifying extrapolations beyond the calibration range and does not fully account for model structural uncertainty or the propagation of cumulative errors.

The models developed in this study used a dataset of *in situ* observations from a range of different climate types, vegetation covers and topographical characteristics. The diversity in this dataset ensured sufficient variance for the GBM, as evidenced by the model metrics (Table 2), and was a key factor in the successful application of the framework. These results show that our framework is highly transferable to other tropical regions with similar environmental modulators. Furthermore, it can be adapted for regions with different characteristics, provided that multiple variations of a single parameter are used without violating the assumption of multicollinearity.

3.2. Environmental modulators

Results showed that simulated soil properties the most influential environmental modulators were climate, topography, and vegetation (Fig. 5). This consistently reflects broader soil-forming processes, including climate-driven weathering, erosion, and vegetation–soil feedback. A better understanding of how these environmental factors affect physical and chemical soil properties can help manage their changes in response to future climate conditions or land use modifications, such as deforestation (Badía et al., 2016). In our study area, the properties related to topographic and climatic conditions were dominant predictors for all soil properties, whereas the weights for covariates related to vegetation were slightly greater for soil property estimates related to sand, i.e., SOL_SAND, CS, and FS. Topography is consistently included as an input variable in our models (Fig. 5) because it is a key factor in soil formation in Northeast Brazil (Oliveira et al., 2018). The topographic conditions (see Table 1) comprise slope, which may affect the quantity of soil deposition or erosion; aspect, which drives the direction of surface and subsurface runoff, and relative exposure of soils to sunlight; and finally curvature, which changes water flow velocity, controlling erosion and deposition processes (Barbieri et al., 2009; Patton et al., 2018).

The model weights for the L_MAX model were largest for NDVI (18 %) and terrain elevation (DEM, 13 %) as its main inputs. Elevation is well related to climate conditions (Badía et al., 2016), which impact the speed at which parent materials weather and erode, and hence the rate of soil development, e.g., via accumulation of organic matter on top of the soil. As for NDVI, it most likely indirectly reflects the vertical

Table 2

Calibrated values for the hyperparameters $n_{\text{estimators}}$ (NE), max_depth (MD), min_samples_split (MSS) and min_samples_leaf (MSL) of the Gradient Boosting Models (GBM), for each estimated soil property and their corresponding calibration performance. The description of the variables can be found in Table 1.

Output variable	Calibrated hyperparameters				Calibration Accuracy ^(a) or $r^{2(b)}$	Verification		
	NE	MD	MSS	MSL		Accuracy ^(a) or $r^{2(b)}$	RMSE	PBIAS
L_MAX	1325	23	41	70	0.91 ^(a)	0.96 ^(a)	-	-
SOL_Z (mm)	4445	3	36	7	0.92 ^(b)	0.98 ^(b)	73.19	0.02
SOL_SAND (%)	2521	87	73	6	0.77 ^(b)	0.91 ^(b)	6.27	1.14
SOL_CLAY (%)	1518	38	85	12	0.78 ^(b)	0.93 ^(b)	4.48	0.29
SOL_SILT (%)	1624	85	15	3	0.76 ^(b)	0.79 ^(b)	4.77	-1.36
SOL_CBN (%)	1265	27	17	43	0.78 ^(b)	0.91 ^(b)	0.14	-3.39
SB ($\text{cmol}_c \text{ kg}^{-1}$)	1026	46	23	2	0.82 ^(b)	0.95 ^(b)	1.79	2.97
CS (%)	2893	38	40	63	0.92 ^(b)	0.98 ^(b)	2.46	1.04
FS (%)	2282	3	7	13	0.89 ^(b)	0.97 ^(b)	2.03	-0.03

Table 3
List of input parameters used for calibrating the Gradient Boosting Models of basic soil properties. The weights (w) calculated for each input in the models are between parentheses. The description of the variables and parameters can be found in Table 1.

Output variable	Inputs (in fractions)
L_MAX	NDVI (0.18), DEM (0.13), ASPECT (0.07), PCPMM (0.07), WNDV (0.07), AAT_ASPECT (0.05), CUR (0.05), TMPSTD (0.03), CTI (0.03), SPR (0.03), PCPSTD (0.03), TMPMN (0.03), TMPSTD (0.03), ATT_SPR_F (0.02), LST (0.02), PCPSKW (0.02), RHAV (0.02), SOLARAV (0.02).
SOL_Z	LAYER (0.83), AAT_ASPECT (0.02), CUR (0.02), NDVI (0.02), TMPMN (0.02), L_MAX (0.02), CTI (0.01), PCPSKW (0.01), PCPMM (0.01), SOLARAV (0.01), WNDV (0.01), TMPSTD (0.01).
SOL_SAND	NDVI (0.09), WNDV (0.09), CTI (0.08), LST (0.08), ASPECT (0.07), CUR (0.07), TMPMN (0.07), PCPSKW (0.06), DEM (0.06), LAYER (0.06), ATT_CUR (0.05), TMPSTD (0.05), L_MAX (0.05).
SOL_CLAY	AAT_ASPECT (0.08), PCPMM (0.08), LST (0.07), ASPECT (0.06), CUR (0.06), WNDV (0.06), DEM (0.05), CTI (0.04), NDVI (0.04), PCPSTD (0.04), ATT_CUR (0.03), RHAV (0.02), SOLARAV (0.02), TMPSTD (0.02), TMPMN (0.02), TMPSTD (0.02), ATT_SPR_F (0.01), SPR (0.01), PCPSKW (0.01), TMPMN (0.01).
SOL_SILT	TMPMN (0.11), SOL_Z (0.1), DEM (0.09), ASPECT (0.07), PCPMM (0.07), CTI (0.05), CUR (0.05), RHAV (0.05), L_MAX (0.05), AAT_ASPECT (0.04), ATT_SPR_F (0.04), SOLARAV (0.03), TMPSTD (0.03), TMPSTD (0.03), LAYER (0.03), SPR (0.02), WNDV (0.02), TMPMN (0.02), PCPSKW (0.01), PCPSTD (0.01).
SOL_CBN	LAYER (0.24), SOL_Z (0.2), ATT_CUR (0.07), NDVI (0.06), CUR (0.04), WNDV (0.04), AAT_ASPECT (0.03), CTI (0.03), SPR (0.03), PCPSKW (0.03), PCPSTD (0.03), PCP_MM (0.03), DEM (0.03), ASPECT (0.02), ATT_SPR_F (0.02), LST (0.02), SOLARAV (0.02), TMPSTD (0.02), TMPSTD (0.02), L_MAX (0.02), RHAV (0.01), TMPSTD (0.01).
SB	RHAV (0.19), WNDV (0.14), PCPSTD (0.08), DEM (0.07), SOL_Z (0.07), TMPMN (0.06), LST (0.05), TMPSTD (0.05), ASPECT (0.04), CUR (0.04), PCPMM (0.04), L_MAX (0.04), AAT_ASPECT (0.03), TMPSTD (0.03), NDVI (0.02), LAYER (0.02), ATT_CUR (0.01), SOLARAV (0.01), TMPMN (0.01).
CS	SOL_SAND (0.65), TMPSTD (0.06), DEM (0.05), TMPMX (0.05), SPR (0.04), LST (0.04), NDVI (0.04), SOLARAV (0.03), WNDV (0.03), PCPSTD (0.02).
FS	SOL_SAND (0.4), SOLARAV (0.09), NDVI (0.07), ATT_CUR (0.05), SPR (0.05), DEM (0.05), TMPMX (0.05), LST (0.04), PCPMM (0.04), RHAV (0.03), TMPSTD (0.03), SOL_Z (0.03), WNDV (0.02).

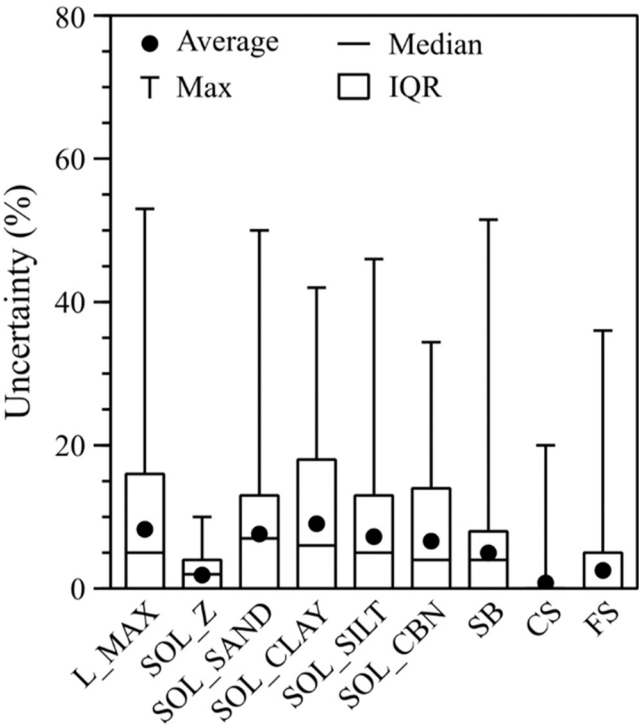


Fig. 4. Uncertainty analysis of the Gradient Boosting Models (GBM) for basic soil parameters. IQR stands for interquartile range, and variable descriptions can be found in Table 1.

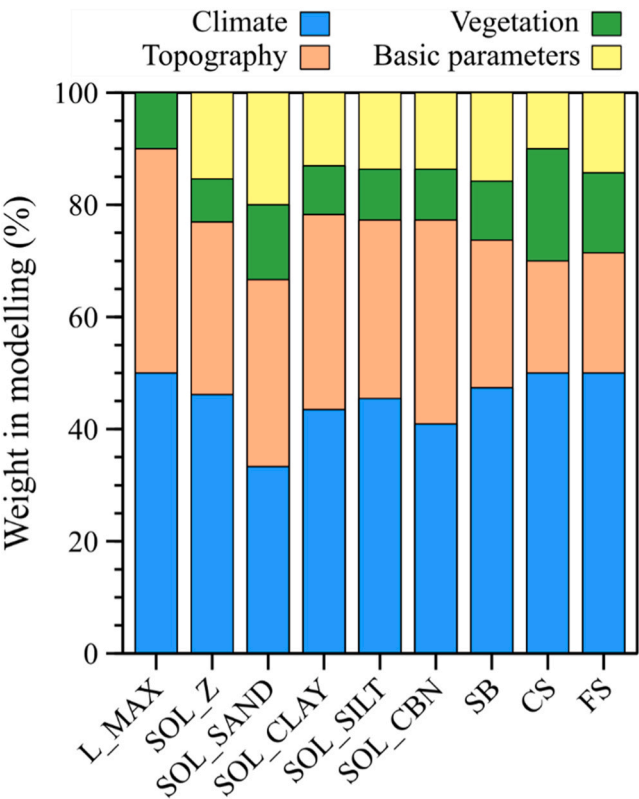


Fig. 5. Proportional weights (w, as in Eq. 5) of the different input variables for modeling each basic soil parameter. The weights for 'basic parameters' represent the influence of other basic soil parameters on the predicted parameter. The description of the variables can be found in Table 1.

variability of soil properties, as soils formed under forests tend to be weathered to greater depth. This occurs because forests grow in higher rainfall areas (Bonan, 2008) and have deeper rooting systems that often create biopores, facilitating internal drainage.

Our model for SB was mainly influenced by relative humidity (19 %) and wind speed (14 %). These variables are known for controlling the intensity of biochemical reactions, and wind erosion (Ravi et al., 2004), respectively. Wind erosion can remove and redistribute topsoil nutrients (Zobeck et al., 1989), affecting local soil nutrient levels, especially in arid and semi-arid regions, as seen in the western region of our study area, where soils are dry and covered by sparse vegetation (Miranda et al., 2020; Ravi et al., 2004). Regarding precipitation, although it may be an important climate factor for soil formation in other regions (e.g., Dixon et al., 2016), its characteristics, i.e., PCPSTD and PCPMM, together weighted only 12 % of the variance in SB in our model.

Regarding the overall importance of the model inputs, key parameters are CTI, L_{MAX}, SOL_Z, and SOL_{SAND} (Table 3). The key role of CTI can be explained by its ability to encapsulate the terrain structure (Gessler et al., 1995; Moore et al., 1993). The influence of SOL_Z on SOL_{SAND} and SOL_{SILT} was relatively strong, suggesting that soil depth plays a critical role in determining sand and silt distribution. The prevalence of sand in surface layers is well-documented, particularly in soils prone to erosion due to their lower structural stability (Valentin and Bresson, 1992). Furthermore, vegetation cover, represented by NDVI, emerged as a key predictor of SOL_{SAND}. High vegetation density often indicates advanced soil weathering or lower sand content, as soils beneath dense forests in high-rainfall regions tend to be more leached and clay-rich (Souza et al., 2016), a pattern observed in the eastern part of our study area.

3.3. Hydraulic parameters predictions via PTFs

The bulk density estimates SOL_{BD_{SR}} (Saxton and Rawls, 2006) and SOL_{BD_{OL}} (Benites et al., 2007) were similar, with a mean difference of only 0.09 g cm⁻³ (Table 4). While both models produced an acceptable range of values, SOL_{BD_{SR}} yielded a small percentage of very high estimates, with 0.85 % of SOL_{BD_{SR}} values exceeding 1.8 g cm⁻³ when considered as a weighted average across all soil layers. Although Benites et al. (2007) reported SOL_{BD} values as high as 2.25 g cm⁻³ in Brazil, we

Table 4

Descriptive statistics of all calculated pedotransfer functions (PTF) data using basic soil properties derived from Gradient Boosting Models. Table 1 contains the description of acronyms that represent the soil hydraulic properties in column 1.

PTF outputs	Mean (SD)		Minimum	Maximum	Invalid values (%)
SOL _{BD_{SR}} (g cm ⁻³)	1.54	(0.09)	1.01	2.60	0
SOL _{BD_{OL}} (g cm ⁻³)	1.45	(0.07)	1.12	1.76	0
SOL _{AWC_{SR}} (mm mm ⁻¹)	0.11	(0.01)	0.01	0.18	0
SOL _{AWC_{BR}} (mm mm ⁻¹)	0.05	(0.03)	0.001	0.17	0.75
SOL _{AWC_{TM}} (mm mm ⁻¹)	0.03	(0.01)	0.001	0.13	5.01
SOL _{AWC_{OL}} (mm mm ⁻¹)	0.07	(0.01)	0.01	0.16	0
SOL _{K_{SR}} (mm hr ⁻¹)	11.17	(14.24)	0.003	932.54	0
SOL _{K_{SR/BR}} (mm hr ⁻¹)	1101.28	(350.5)	10.41	1900.21	0
SOL _{K_{SR/TM}} (mm hr ⁻¹)	26.72	(26.58)	0.001	219.47	12.07
SOL _{K_{BR}} (mm hr ⁻¹)	63.85	(333.9)	8.85	12112	0
USLE _K (unitless)	0.22	(0.03)	0.01	0.41	0

recommend caution when interpreting values above ~2 g cm⁻³. With regards to SOL_{AWC}, the equation by Oliveira et al. (2002), SOL_{AWC_{OL}}, which was calibrated strictly using data from our study area, was the only equation that did not ‘saturate’ when PTFs were applied. Since we evaluate and map soils in a region similar to that of Oliveira et al. (2002), our results highlight the common tendency of PTFs to exhibit overfitting, becoming over-adjusted to the specific datasets that are used for their calibration (De Vos et al., 2005).

Two of the four SOL_K estimates were derived from variations of Saxton and Rawls (2006) (Tables S1 and S2 in the Supplementary Material). The difference between them depends on the calculation of the inputs θ_s , θ_{33} and θ_{1500} , which differ from the approaches originally proposed by Saxton and Rawls (2006), SOL_{K_{SR}}, i.e. those by Barros et al. (2013), SOL_{K_{SR/BR}}, and the one by Tomasella et al. (2000), SOL_{K_{SR/TM}}. Maximum values ranged from 219.47 (SOL_{K_{SR/TM}}) to 1900.21 mm h⁻¹ (SOL_{K_{SR/BR}}). The approach that generates SOL_{K_{BR}} is the simplest; it only uses SOL_Z as input, and therefore it does not exhibit differences for soils with different textures and the same depths. A small number of invalid values was found only for SOL_{AWC_{BR}}, SOL_{AWC_{TM}}, and SOL_{K_{SR/TM}} due to inaccurate extrapolations, i.e., out of the a priori parameter range expected or acceptable for these parameters or PTFs, of θ_r and n . For USLE_K the applied model expects values varying from 0.1 to 0.5 (Sharpley et al., 1993). However, we found values below this range because our simulated dataset included soils with high coarse-sand content.

The SOL_K dataset from Gupta et al. (2021) predominantly exhibited higher values than our SOL_K estimates using the PTF from Saxton and Rawls (2006) (Fig. 6A). Differences in SOL_K exceeded 100 mm h⁻¹ (as indicated by red dashed rectangles in Fig. 6A), and the highest concentration of differences is approximately fivefold (Fig. 6B). For the region with the most humid climate (Am climate in Fig. 1, dashed rectangle 4 in Fig. 6A), we also found a higher clay content (up to 50 %) in our dataset (Fig. 6C) when compared to the data from Hengl (2018) used as an input by Gupta et al. (2021), which we identify as one of the reasons for the SOL_K differences between the datasets for this specific area, despite a lack of overall apparent correlation between clay fraction differences and differences in SOL_K for the entire study region (Fig. 6D). The semi-arid areas with some of the highest differences in SOL_K (Fig. 6A, rectangles 1–3) also exhibit some of the shallowest soils (Fig. 6E). Although we cannot draw a direct relationship between the SOL_K differences and soil depth, it is important to note that deeper soils in this region hold greater clay fractions (Fig. 6F). The dataset by Gupta et al. (2021) follows a standardized soil layer protocol with a total depth of 200 cm for all grid cells, whereas our results were produced following a methodology designed to provide pedological meaning with a more realistic number of soil layers and respective soil profile depths. The impact of these differences goes beyond the disparities in saturated hydraulic values, which themselves carry high uncertainties (Zhang and Schaap, 2019). Estimates of hydraulic properties, even when in a realistic range, can be highly misleading if the soil layers and depth are being assumed spatially homogeneous (Dai, Shangguan, et al., 2019). A better representation of soil profile characteristics in models, such as soil profile depth (Brunke et al., 2016), will lead to more realistic soil maps, as we have shown here, and consequently improve the performance of land surface models (Dy and Fung, 2016; Kearney and Maino, 2018), for example.

We note that only 12 % of the measurements used to train the ML algorithm that generated Gupta et al. (2021)’s dataset were located in the tropics and none in our study area, and that the soil datasets used in their methodology are likely to be substantially different from the one we generated in our study, particularly regarding clay fraction. Also, our comparison of SOL_K values was based on the prediction of SOL_K using the PTF from Saxton and Rawls (2006), which predicted the lowest SOL_K values among the PTFs used in this study (Table 4). This set of PTFs was developed using data from North America, which can lead to high errors and uncertainty when used in other regions (Vereecken

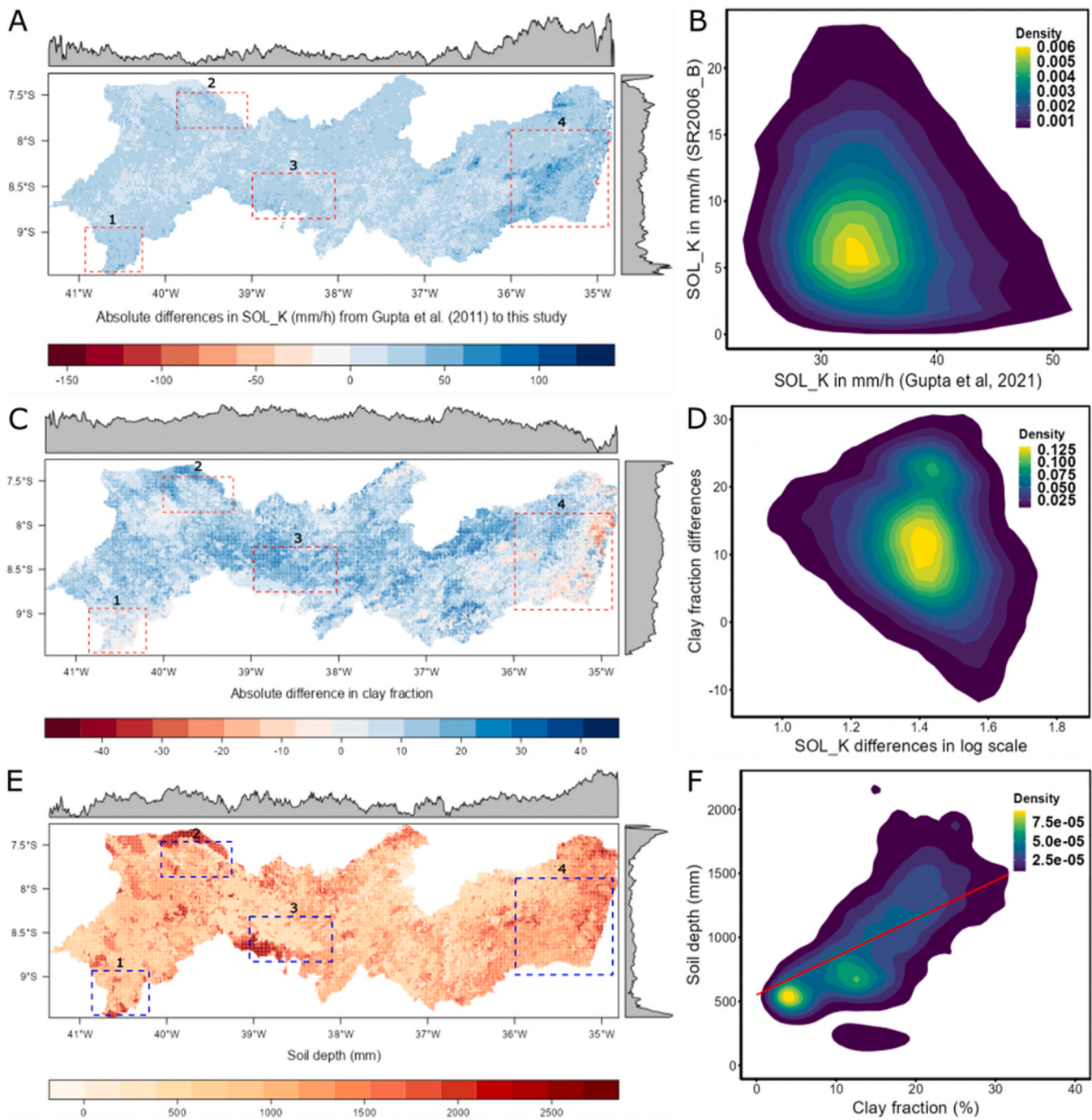


Fig. 6. Differences in saturated hydraulic conductivity (SOL_K) and clay fraction between the data generated and used by Gupta et al. (2021) and results in this study, and total soil depth from our study. The maps (panels A, C, and E) highlight some areas (within dashed rectangles) where the SOL_K differences were the greatest, and the top and right margins exhibit the distribution of the latitudinal and longitudinal means, respectively. The density estimates in panels B, D, and F were calculated using the kde2d function available in the MASS package (Venables and Ripley, 2003) in the R language (R Core Team, 2017).

et al., 2016). Nevertheless, our ML framework was able to generate a soil map with high accuracy (mean $r^2 > 0.9$, Table 2) and low mean uncertainty ($< 10\%$, Fig. 4), thus capturing the variability of basic soil properties that drive most common PTFs. Note that Lehmann et al. (2021) showed that tropical soils can have a higher SOL_K than soils from temperate climates due to the predominance of kaolinite clays over illite clays, for example, in many tropical regions. From a soil hydraulic point of view, kaolinite clays behave more like sandy soils than clay soils. However, based on the dominant clay type data provided by Ito and Wagai (2017); see also Lehmann et al., (2021) in Pernambuco the

prevalence of low activity clays, such kaolinite, is relatively low. This sets this area apart from other South American tropical regions such as the Amazon rainforest. Lehmann et al. (2021) point out that clay mineral-informed pedotransfer functions and machine learning algorithms trained with datasets including different clay types and soil structure formation processes may improve soil hydraulic properties prediction. In that case it is important to consider that not all tropical clay types are necessarily kaolinite.

4. Conclusions

In this study, we produced robust soil property maps using a data-driven ML framework based on integration of a covariance model (SLEEP) with decision trees (GBM), for a tropical region with highly variable topography, climate, and vegetation characteristics that is not well represented in global soil property datasets. Good model performance is reflected in our models' statistics that present r^2 and PBIAS values varying from 0.79 to 0.98, and from -1.39 to 1.14 , respectively. Decision tree methods are highly advantageous because they are free of strict assumptions and can simultaneously handle diverse variables, scales, distributions, and relationships. We explored this characteristic in detail in this study, by employing multiple freely available datasets with an extensive array of data types (e.g., number of soil layers and chemical composition) to improve the soil information in our study area. GBM models can be considered semi-black-box models due to the complexity introduced by combining multiple individual trees, which often limits their direct interpretability. We addressed this challenge by incorporating a feature selector during calibration, which enabled us to perform uncertainty analyses and identify the primary environmental modulators of various soil properties.

Our results are especially important for soil management in response to climate change, land-use changes, and environmental degradation, such as deforestation and desertification, at multiple spatial scales. Our machine learning framework offers enhanced flexibility, enables regular short-term map updates, and supports the integration of future economic and environmental modelling (e.g., <https://super.hawqs.tamu.edu/>), while drastically reducing capital investments compared to in situ surveys and mapping. We believe that these promising findings will enhance all modelling efforts that require detailed soil information and encourage the development of new frameworks and datasets for soil sciences. Our new dataset can be further used to create a new portfolio of applications, such as agricultural zoning and environmental management strategies.

CRediT authorship contribution statement

Anne Verhoef: Writing – review & editing, Writing – original draft, Visualization, Validation, Methodology, Investigation, Formal analysis, Conceptualization. **Montenegro Suzana:** Writing – review & editing, Funding acquisition. **Nóbrega Rodolfo L. B.:** Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Resources, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Feras Ziadat:** Writing – review & editing, Validation, Software. **Raghavan Srinivasan:** Writing – review & editing, Writing – original draft, Methodology. **Miranda Rodrigo:** Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Hui Shao:** Writing – review & editing. **Wanhong Yang:** Writing – review & editing. **Souza Alzira:** Writing – review & editing, Methodology, Data curation, Conceptualization. **Barros Alexandre:** Writing – review & editing, Methodology, Data curation, Conceptualization. **Moura Magna:** Writing – review & editing, Writing – original draft, Validation, Methodology, Conceptualization. **Araújo Filho José:** Writing – review & editing, Validation, Methodology, Data curation, Conceptualization. **Silva Jadson:** Writing – review & editing, Data curation. **Galvêncio Josicléda:** Writing – review & editing, Writing – original draft, Supervision, Funding acquisition, Conceptualization. **Silva Estevão:** Writing – review & editing, Data curation. **Araújo Maria:** Writing – review & editing, Funding acquisition.

Declaration of Competing Interest

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

Acknowledgments

We thank Budiman Minasny and two anonymous reviewers for their valuable inputs, and Surya Gupta for clarifying aspects on [Gupta et al. \(2021\)](#) that allowed us to use and compare it to our results. Soil data access was provided by the Brazilian Agricultural Research Corporation (EMBRAPA) through the Agroecological Zoning of the state of Pernambuco (ZAPE) project. The authors also acknowledge the following funding sources: The Brazilian Coordination for the Improvement of Higher Level Personnel (CAPES 88887.371850/2019-00) for AGSSS and JDG; The Fundação de Amparo a Ciência e Tecnologia do Estado de Pernambuco (Project FACEPE APQ, 0646-9.25/16) for RQM and JDG; The National Council for Scientific and Technological Development of Brazil (CNPq) through the projects MCTIC/CNPq 28/2018 (431980/2018-7), PEGASUS MCTI/CNPq N° 19/2017 (441305/2017-2), CNPq/MCTIC/BRICS 29/2017 (442335/2017-2), INCT Mudanças Climáticas II, INCT ONSEADAPTA (406919/2022-4), and productivity grants (448236/2014-1 and 313469/2020-2) for SMGLM and MSBA, and CNPq/MCTI/FNDCT N° 21/2024 447433/2024-5) for RLBN; and the UK Natural Environment Research Council (NE/N012526/1 ICL and NE/N012488/1 UoR) and FAPESP (The São Paulo Research Foundation) (FAPESP 2015/50488-5) for the UK/Brazil Nordeste project for AV, RLBN and MSBM. For the purpose of open access, the authors have applied a Creative Commons Attribution (CC BY) license to any Author Accepted Manuscript version arising from this submission.

Appendix A. Supporting information

Supplementary information and analyses associated with this article can be found in the online version of the Supplementary Material at [doi:10.1016/j.soilad.2025.100064](https://doi.org/10.1016/j.soilad.2025.100064).

Data Availability

The authors do not have permission to share data.

Model outputs from the study "A scalable framework for soil property mapping tested across a highly diverse tropical data-scarce region" (ZENODO)

References

- ABNT, 1995. *Rochas e Solo* (No. NBR 6502). Associação Brasileira de Normas Técnicas, Rio de Janeiro.
- Alvares, C.A., Stape, J.L., Sentelhas, P.C., de Moraes Gonçalves, J.L., Sparovek, G., 2013. Köppen's climate classification map for Brazil. *Meteorol. Z.* 22 (6), 711–728. <https://doi.org/10.1127/0941-2948/2013/0507>.
- Araújo Filho, J.C. de, Araújo, M. do S.B. de, Marques, F.A., Lopes, H.L., 2014. Solos. In: Torres, F.S. de M., Pfaltzgraff, P.A. dos S. (Eds.), *Geodiversidade do estado de Pernambuco*. MINISTÉRIO DE MINAS E ENERGIA.
- Arnold, J.G., Srinivasan, R., Muttiah, R.S., Williams, J.R., 1998. Large area hydrologic modeling and assessment part i: Model development. *J. Am. Water Resour. Assoc.* 34 (1), 73–89. <https://doi.org/10.1111/j.1752-1688.1998.tb05961.x>.
- Auzzas, A., Capra, G.F., Jani, A.D., Ganga, A., 2024. An improved digital soil mapping approach to predict total N by combining machine learning algorithms and open environmental data. *Model. Earth Syst. Environ.* 10, 6519–6538. <https://doi.org/10.1007/s40808-024-02127-8>.
- Badía, D., Ruiz, A., Girona, A., Martí, C., Casanova, J., Ibarra, P., Zufiaurre, R., 2016. The influence of elevation on soil properties and forest litter in the Siliceous Moncayo Massif, SW Europe. *J. Mt. Sci.* 13 (12), 2155–2169. <https://doi.org/10.1007/s11629-015-3773-6>.
- Ballabio, C., Panagos, P., Monatanarella, L., 2016. Mapping topsoil physical properties at European scale using the LUCAS database. *Geoderma* 261, 110–123. <https://doi.org/10.1016/j.geoderma.2015.07.006>.
- Bao, Y., Yao, F., Meng, X., Wang, J., Liu, H., Wang, Y., Liu, Q., Zhang, J., Mouazen, A.M., 2024. A fine digital soil mapping by integrating remote sensing-based process model and deep learning method in Northeast China. *Soil Tillage Res.* 238, 106010. <https://doi.org/10.1016/j.still.2024.106010>.
- Barbieri, D.M., Marques Júnior, J., Alleoni, L.R.F., Garbuio, F.J., Camargo, L.A., 2009. Hillslope curvature, clay mineralogy, and phosphorus adsorption in an Alfisol cultivated with sugarcane. *Sci. Agric.* 66 (6), 819–826. <https://doi.org/10.1590/s0103-90162009000600015>.

- Barros, A.H.C., de Jong van Lier, Q., 2014. Pedotransfer functions for Brazilian soils. Application of Soil Physics in Environmental Analyses. Springer International Publishing, Cham, pp. 131–162. https://doi.org/10.1007/978-3-319-06013-2_6.
- Barros, A.H.C., van Lier, Q., de J., Maia, A., de H.N., Scarpere, F.V., 2013. Pedotransfer functions to estimate water retention parameters of soils in northeastern Brazil. *Rev. Bras. De. Cienc. Do Solo* 37 (2), 379–391. <https://doi.org/10.1590/s0100-06832013000200009>.
- Belk, E.L., Markewitz, D., Rasmussen, T.C., Carvalho, E.J.M., Nepstad, D.C., Davidson, E. A., 2007. Modeling the effects of throughfall reduction on soil water content in a Brazilian Oxisol under a moist tropical forest. *Water Resour. Res.* 43 (8). <https://doi.org/10.1029/2006wr005493>.
- Benites, V.M., Machado, P.L.O.A., Fidalgo, E.C.C., Coelho, M.R., Madari, B.E., 2007. Pedotransfer functions for estimating soil bulk density from existing soil survey reports in Brazil. *Geoderma* 139 (1–2), 90–97. <https://doi.org/10.1016/j.geoderma.2007.01.005>.
- Bonan, G.B., 2008. Forests and climate change: forcings, feedbacks, and the climate benefits of forests. *Science* 320 (June), 1444–1450. <https://doi.org/10.1126/science.1155121>.
- Bouma, J., McBratney, A., 2013. Framing soils as an actor when dealing with wicked environmental problems. *Geoderma* 200–201, 130–139. <https://doi.org/10.1016/j.geoderma.2013.02.011>.
- Brunke, M.A., Broxton, P., Pelletier, J., Gochis, D., Hazenberg, P., Lawrence, D.M., et al., 2016. Implementing and evaluating variable soil thickness in the Community Land Model, version 4.5 (CLM4.5). *J. Clim.* 29 (9), 3441–3461. <https://doi.org/10.1175/jcli-d-15-0307.1>.
- Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. SMOTE: synthetic minority over-sampling technique. *J. Artif. Intell. Res.* 16, 321–357. <https://doi.org/10.1613/jair.953>.
- Core Team, R., 2017. R: A Language and Environment for Statistical Computing (Version 3.3.3). R Foundation for Statistical Computing, Vienna, Austria. (<https://www.r-project.org/>).
- Dai, Y., Shangguan, W., Wei, N., Xin, Q., Yuan, H., Zhang, S., et al., 2019. A review of the global soil property maps for Earth system models. *SOIL* 5 (2), 137–158. <https://doi.org/10.5194/soil-5-137-2019>.
- de Morisson Valeriano, M., de Fátima Rossetti, D., 2012. Topodata: Brazilian full coverage refinement of SRTM data. *Appl. Geogr. (Seven. Engl.)* 32 (2), 300–309. <https://doi.org/10.1016/j.apgeog.2011.05.004>.
- De Vos, B., Van Meirvenne, M., Quataert, P., Deckers, J., Muys, B., 2005. Predictive quality of pedotransfer functions for estimating bulk density of forest soils. *Soil Sci. Soc. Am. J. Soil Sci. Soc. Am.* 69 (2), 500–510. <https://doi.org/10.2136/sssaj2005.0500>.
- Didan, K., 2015. MOD13A3 MODIS/terra vegetation indices monthly L3 global 1km SIN grid V006 [Data set]. NASA EOSDIS Land Process. DAAC. <https://doi.org/10.5067/MODIS/MOD13A3.006>.
- Dixon, J.L., Chadwick, O.A., Vitousek, P.M., 2016. Climate-driven thresholds for chemical weathering in postglacial soils of New Zealand. *J. Geophys. Res. Earth Surf.* 121 (9), 1619–1634. <https://doi.org/10.1002/2016jf003864>.
- Dy, C.Y., Fung, J.C.H., 2016. Updated global soil map for the Weather Research and Forecasting model and soil moisture initialization for the Noah land surface model. *J. Geophys. Res. Atmos.* 121 (15), 8777–8800. <https://doi.org/10.1002/2015jd024558>.
- Embrapa, 1997. *Manual de Métodos de Análise de Solo*, 2nd ed. EMBRAPA-CNPq, Rio de Janeiro, p. 212.
- Food and Agriculture Organization (FAO), 2006. *Guidelines for soil description*. Food & Agriculture Organization of the United Nations, 4th ed. FAO, Rome, Italy.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann. Stat.* 29. <https://doi.org/10.1214/aos/1013203451>.
- Gessler, P.E., Moore, I.D., McKENZIE, N.J., Ryan, P.J., 1995. Soil-landscape modelling and spatial prediction of soil attributes. *Int. J. Geogr. Inf. Syst.* 9 (4), 421–432. <https://doi.org/10.1080/02693799508902047>.
- Greenbelt. (2019). Earthdata Search. Earth Science Data and Information System (ESDIS) Project, Earth Science Projects Division (ESPD), Flight Projects Directorate, Goddard Space Flight Center (GSFC) National Aeronautics and Space Administration (NASA). Retrieved April 11, 2021, from (<https://search.earthdata.nasa.gov/>).
- 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., Bustamante Cañas, N., Cruz-Gaistardo, C.O., Davila, F., Dell Acqua, M., Encina, A., Figueredo Tacona, H., Fontes, F., Hernández Herrera, Ibelles Navarro, Loayza, V., Manueles, A. M., Mendoza Jara, F., Olivera, C., Osorio Hermsilla, R., Pereira, G., Prieto, P., Ramos, I.A., Rey Brina, Rivera, R., Rodríguez-Rodríguez, J., Roopnarine, R., Rosales Ibarra, A., Rosales Riveiro, Schulz, G.A., Spence, A., Vasques, G.M., Vargas, R.R.R.R. R., Vargas, R.R.R.R.R., Federico Olmedo, G., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C., Arévalo, G.E., Eduardo Arroyo-Cruz, C., Bolívar, A., Bunning, S., Bustamante Cañas, N., Omar Cruz-Gaistardo, C., Davila, F., Dell Acqua, M., Encina, A., Tacona, H.F., Fontes, F., Herrera, J.A.H., Roberto Ibelles Navarro, A., Loayza, V., Manueles, A.M., Mendoza Jara, F., Olivera, C., Osorio Hermsilla, R., Pereira, G., Prieto, P., Ramos, I.A., Carlos Rey Brina, J., Rivera, R., Rodríguez-Rodríguez, J., Roopnarine, R., Ibarra, A.R., Amaury Rosales Riveiro, K., Andrés Schulz, G., Spence, A., Vasques, G.M., Vargas, R.R.R.R.R., Vargas, 2018. No silver bullet for digital soil mapping: Country-specific soil organic carbon estimates across Latin America 4, 173–193. <https://doi.org/10.5194/soil-4-173-2018>.
- Gupta, S., Lehmann, P., Bonetti, S., Papritz, A., Or, D., 2021. Global prediction of soil saturated hydraulic conductivity using random forest in a covariate-based GeoTransfer function (CoGTF) framework. *J. Adv. Model. Earth Syst.* 13 (4). <https://doi.org/10.1029/2020ms002242>.
- Guyon, I., Weston, J., Barnhill, S., Vapnik, V., 2002. *Mach. Learn.* 46 (1/3), 389–422. <https://doi.org/10.1023/a:1012487302797>.
- Hartemink, A.E., Lowery, B., Wacker, C., 2012. Soil maps of Wisconsin. *Geoderma* 189–190, 451–461. <https://doi.org/10.1016/j.geoderma.2012.05.025>.
- Hateffard, F., Steinbuch, L., Heuvelink, G.B.M., 2024. Evaluating the extrapolation potential of random forest digital soil mapping. *Geoderma* 441, 116740. <https://doi.org/10.1016/j.geoderma.2023.116740>.
- Hengl, T., 2018. Clay content in % (kg / kg) at 6 standard depths (0, 10, 30, 60, 100 and 200 cm) at 250 m resolution [Data set]. Zenodo. <https://doi.org/10.5281/ZENODO.2525663>.
- Hengl, T., Mendes de Jesus, J., Heuvelink, G.B.M., Ruiperez Gonzalez, M., Kilibarda, M., Blagotić, A., et al., 2017. SoilGrids250m: Global gridded soil information based on machine learning. *PloS One* 12 (2), e0169748. <https://doi.org/10.1371/journal.pone.0169748>.
- Ito, A., Wagai, R., 2017. Global distribution of clay-size minerals on land surface for biogeochemical and climatological studies. *Sci. Data* 4, 170103. <https://doi.org/10.1038/sdata.2017.103>.
- Kearney, M.R., Maino, J.L., 2018. Can next-generation soil data products improve soil moisture modelling at the continental scale? An assessment using a new microclimate package for the R programming environment. *J. Hydrol.* 561, 662–673. <https://doi.org/10.1016/j.jhydrol.2018.04.040>.
- Kempen, B., Brus, D.J., Stoorvogel, J.J., Heuvelink, G.B.M., de Vries, F., 2012. Efficiency comparison of conventional and digital soil mapping for updating soil maps. *Soil Sci. Soc. Am. J. Soil Sci. Soc. Am.* 76 (6), 2097–2115. <https://doi.org/10.2136/sssaj2011.0424>.
- Lagacherie, P., McBratney, A.B., 2006. Chapter 1 spatial soil information systems and spatial soil inference systems: perspectives for digital soil mapping. *Developments in Soil Science*. Elsevier, pp. 3–22. [https://doi.org/10.1016/s0166-2481\(06\)31001-x](https://doi.org/10.1016/s0166-2481(06)31001-x).
- Laurent, F., Pocard-Chapuis, R., Plassin, S., Pimentel Martinez, G., 2017. Soil texture derived from topography in North-eastern Amazonia. *J. Maps* 13 (2), 109–115. <https://doi.org/10.1080/17445647.2016.1266524>.
- Lehmann, P., Leshchinsky, B., Gupta, S., Mirus, B.B., Bickel, S., Lu, N., Or, D., 2021. Clays are not created equal: how clay mineral type affects soil parameterization. *Geophys. Res. Lett.* 48, e2021GL095311. <https://doi.org/10.1029/2021GL095311>.
- Li, J., Heap, A.D., 2014. Spatial interpolation methods applied in the environmental sciences: a review. *Environ. Model. Softw. Environ. Data N.* 53, 173–189. <https://doi.org/10.1016/j.envsoft.2013.12.008>.
- McBratney, A.B., Mendonça Santos, M.L., Minasny, B., 2003. On digital soil mapping. *Geoderma* 117 (1–2), 3–52. [https://doi.org/10.1016/s0016-7061\(03\)00223-4](https://doi.org/10.1016/s0016-7061(03)00223-4).
- Mendonça-Santos, M.L., dos Santos, H.G., 2006. Chapter 3 the state of the art of Brazilian soil mapping and prospects for digital soil mapping. *Developments in Soil Science*. Elsevier, pp. 39–601. [https://doi.org/10.1016/s0166-2481\(06\)31003-3](https://doi.org/10.1016/s0166-2481(06)31003-3).
- Minasny, B., Hartemink, A.E., 2011. Predicting soil properties in the tropics. *Earth Sci. Rev.* 106 (1–2), 52–62. <https://doi.org/10.1016/j.earscirev.2011.01.005>.
- Miranda, R. de Q., Nóbrega, R.L.B., Galvão, J.D., 2022. SLEEPy - an implementation of the soil-landscape estimation and evaluation program using machine learning modeling (Version 1.1) [Python]. Github. Retrieved from. (<https://github.com/razeayres/sleepy>).
- Miranda, R. de Q., Nóbrega, R.L.B., da Silva, E.L.R., da Silva, J.F., de Araújo Filho, J.C., de Moura, M.S.B., et al., 2025. Model outputs from the study "A scalable framework for soil property mapping tested across a highly diverse tropical data-scarce region" [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.15603168>.
- Miranda, R. de Q., Nóbrega, R.L.B., Moura, M.S.B. de, Raghavan, S., Galvão, J.D., 2020. Realistic and simplified models of plant and leaf area indices for a seasonally dry tropical forest. *Int. J. Appl. Earth Obs. Geoinf.* 85, 101992. <https://doi.org/10.1016/j.jag.2019.101992>.
- Montzka, C., Herbst, M., Weihermüller, L., Verhoef, A., Vereecken, H., 2017. A global data set of soil hydraulic properties and sub-grid variability of soil water retention and hydraulic conductivity curves. *Earth Syst. Sci. Data* 9 (2), 529–543. <https://doi.org/10.5194/essd-9-529-2017>.
- Moore, I.D., Gessler, P.E., Nielsen, G.A., Peterson, G.A., 1993. Soil attribute prediction using terrain analysis. *Soil Sci. Soc. Am. J. Soil Sci. Soc. Am.* 57 (2), 443–452. <https://doi.org/10.2136/sssaj1993.036159950057000200026x>.
- Natekin, A., Knoll, A., 2013. Gradient boosting machines, a tutorial. *Front. Neuroinformatics* 7. <https://doi.org/10.3389/fnbot.2013.00021>.
- Nóbrega, R.L.B., Ziembowicz, T., Torres, G.N., Guzha, A.C., Amorim, R.S.S., Cardoso, D., Johnson, M.S., Santos, T.G., Couto, E., Gerold, G., 2020. Ecosystem services of a functionally diverse riparian zone in the Amazon-Cerrado agricultural frontier. *Global Ecology and Conservation* 21, e00819. <https://doi.org/10.1016/j.gecco.2019.e00819>.
- Nozari, S., Pahlavan-Rad, M.R., Brungard, C., Heung, B., Borůvka, L., 2024. Digital soil mapping using machine learning-based methods to predict soil organic carbon in two different districts in the Czech Republic. *Soil Water Res.* 19, 32–49. <https://doi.org/10.17221/119/2023-SWR>.
- Oliveira, L.B., Ribeiro, M.R., Jacomine, P.K.T., Rodrigues, J.J.V., Marques, F.A., 2002. Funções de pedotransferência para predição da umidade retida a potenciais específicos em solos do estado de Pernambuco. *Rev. Bras. De. Cienc. Do Solo* 26 (2), 315–323. <https://doi.org/10.1590/s0100-06832002000200004>.
- Oliveira, D.P., Sartor, L.R., Souza Júnior, V.S., Corrêa, M.M., Romero, R.E., Andrade, G. R.P., Ferreira, T.O., 2018. Weathering and clay formation in semi-arid calcareous soils from Northeastern Brazil. *Catena* 162, 325–332. <https://doi.org/10.1016/j.catena.2017.10.030>.
- Orgiazzi, A., Bardgett, R.D., Barrios, E., Behan-Pelletier, V., Briones, M.J.I., Chotte, J.-L., et al. (Eds.), 2016. *Global soil biodiversity atlas*. European Commission, Publications Office of the European Union, Luxembourg. Retrieved from. (<https://data.europa.eu/doi/10.2788/2613>).

- Padarian, J., Minasny, B., McBratney, A.B., 2017. Chile and the Chilean soil grid: A contribution to. *GlobalSoilMap* 9, 17–28. <https://doi.org/10.1016/j.geodrs.2016.12.001>.
- Patton, N.R., Lohse, K.A., Godsey, S.E., Crosby, B.T., Seyfried, M.S., 2018. Predicting soil thickness on soil mantled hillslopes. *Nat. Commun.* 9 (1), 3329. <https://doi.org/10.1038/s41467-018-05743-y>.
- Qu, L., Lu, H., Tian, Z., Schoorl, J.M., Huang, B., Liang, Yonghong, Qiu, D., Liang, Yin, 2024. Spatial prediction of soil sand content at various sampling density based on geostatistical and machine learning algorithms in plain areas. *Catena* 234, 107572. <https://doi.org/10.1016/j.catena.2023.107572>.
- Rahmati, M., Weihermüller, L., Vanderborght, J., Pachepsky, Y.A., Mao, L., Sadeghi, S. H., et al., 2018. Development and analysis of the soil water infiltration global database. *Earth Syst. Sci. Data* 10 (3), 1237–1263. <https://doi.org/10.5194/essd-10-1237-2018>.
- Ravi, S., D'Odorico, P., Over, T.M., Zobeck, T.M., 2004. On the effect of air humidity on soil susceptibility to wind erosion: the case of air-dry soils. *Geophys. Res. Lett.* 31 (9). <https://doi.org/10.1029/2004gl019485>.
- Richards, J.A., 2013. *Remote Sensing Digital Image Analysis*. Springer Berlin Heidelberg, Berlin, Heidelberg. <https://doi.org/10.1007/978-3-642-30062-2>.
- Salgueiro, J.H.P. de B., Montenegro, S.M.G.L., Pinto, E.J. de A., Silva, B.B. da, Souza, W. M. de, Oliveira, L.M.M. de, 2016. Influence of oceanic-atmospheric interactions on extreme events of daily rainfall in the Sub-basin 39 located in Northeastern Brazil. *RBRH* 21 (4), 685–693. <https://doi.org/10.1590/2318-0331.011616023>.
- Saxton, K.E., Rawls, W.J., 2006. Soil water characteristic estimates by texture and organic matter for hydrologic solutions. *Soil Sci. Soc. Am. J. Soil Sci. Soc. Am.* 70 (5), 1569. <https://doi.org/10.2136/sssaj2005.0117>.
- Scharlemann, J.P.W., Tanner, E.V.J., Hiederer, R., Kapos, V., 2014. Global soil carbon: understanding and managing the largest terrestrial carbon pool. *Carbon manag* 5, 81–91. <https://doi.org/10.4155/cmt.13.77>.
- Scull, P., Franklin, J., Chadwick, O.A., McArthur, D., 2003. Predictive soil mapping: a review. *Prog. Phys. Geogr.* 27 (2), 171–197. <https://doi.org/10.1191/0309133303pp366ra>.
- Sharpley, A.N., Williams, J.R., United States, & Agricultural Research Service, 1993. EPIC Eros. /Product. Impact Calc. 1 Model Doc. Retrieved from (<https://handle.nal.usda.gov/10113/CAT10698097>).
- Silva, F.B.R., Santos, J.C.P., Silva, A.B., Calvacanti, A.C., Silva, F.H.B.B., Burgos, N., Parahyba, R.B.V., Oliveira Neto, Souza Neto, Araújo Filho, Lopes, O.F., Luz, L.R.Q.P., Leite, A.P., Souza, L.G.M.C., Silva, C.P., Vazão-Silva, M.A., Barros, 2001. *Zonamento agroecológico do Estado de Pernambuco*.
- Souza, R., Feng, X., Antonino, A., Montenegro, S., Souza, E., Porporato, A., 2016. Vegetation response to rainfall seasonality and interannual variability in tropical dry forests. *Hydrol. Process.* 30 (20), 3583–3595. <https://doi.org/10.1002/hyp.10953>.
- Souza, C.M., Jr, Z. Shimbo, J., Rosa, M.R., Parente, L.L., A. Alencar, A., Rudorff, B.F.T., et al., 2020. Reconstructing three decades of land use and land cover changes in Brazilian biomes with Landsat archive and earth engine. *Remote Sens.* 12 (17), 2735. <https://doi.org/10.3390/rs12172735>.
- Souza, A.G.S.S., Ribeiro Neto, A., Souza, de, L.L., 2021. Soil moisture-based index for agricultural drought assessment: SMADI application in Pernambuco. *State-Brazil* 252, 112124. <https://doi.org/10.1016/j.rse.2020.112124>.
- Sun, L., Liu, F., Zhu, X., Zhang, G., 2024. High-resolution digital mapping of soil erodibility in China. *Geoderma* 444, 116853. <https://doi.org/10.1016/j.geoderma.2024.116853>.
- Taghizadeh-Mehrjardi, R., Nabiollahi, K., Kerry, R., 2016. Digital mapping of soil organic carbon at multiple depths using different data mining techniques in Baneh region, Iran. *Geoderma* 266, 98–110. <https://doi.org/10.1016/j.geoderma.2015.12.003>.
- Tarboton, D.G., Bras, R.L., Rodriguez-Iturbe, I., 1991. On the extraction of channel networks from digital elevation data. *Hydrol. Process.* 5 (1), 81–100. <https://doi.org/10.1002/hyp.3360050107>.
- Teng, H.T., Viscarra Rossel, Shi, Z., Behrens, T., 2018. Updating a national soil classification with spectroscopic predictions. and digital soil mapping 164, 125–134. <https://doi.org/10.1016/j.catena.2018.01.015>.
- The pandas development team, 2024. pandas-dev/pandas: Pandas. <https://doi.org/10.5281/ZENODO.3509134>.
- Tomasella, J., Hodnett, M.G., Rossato, L., 2000. Pedotransfer functions for the estimation of soil water retention in Brazilian soils. *Soil Sci. Soc. Am. J. Soil Sci. Soc. Am.* 64 (1), 327–338. <https://doi.org/10.2136/sssaj2000.641327x>.
- Torres, F.S. de M., Pfaltzgraff, P.A. dos S. (Eds.), 2014. *Geodiversidade do estado de Pernambuco*. CPRM. Retrieved from. (<http://rigeo.cprm.gov.br/handle/doc/16771>).
- Tóth, B., Weynants, M., Pásztor, L., Hengl, T. 3D soil hydraulic database of Europe at 250 m resolution, 31, 2662–2666. <https://doi.org/10.1002/hyp.11203>.
- Tziachris, P., Aschonitis, V., Chatzistathis, T., Papadopolou, M., 2019. Assessment of spatial hybrid methods for predicting soil organic matter using DEM derivatives and soil parameters. *Catena* 174, 206–216. <https://doi.org/10.1016/j.catena.2018.11.010>.
- Valentin, C., Bresson, L.-M., 1992. Morphology, genesis and classification of surface crusts in loamy and sandy soils. *Geoderma* 55 (3–4), 225–245. [https://doi.org/10.1016/0016-7061\(92\)90085-1](https://doi.org/10.1016/0016-7061(92)90085-1).
- van der Westhuizen, S., Heuvelink, G.B.M., Hofmeyr, D.P., 2023. Multivariate random forest for digital soil mapping. *Geoderma* 431, 116365. <https://doi.org/10.1016/j.geoderma.2023.116365>.
- van Genuchten, M.T., 1980. A closed-form equation for predicting the hydraulic conductivity of unsaturated soils. *Soil Sci. Soc. Am. J.* 44 (5), 892–898. <https://doi.org/10.2136/sssaj1980.03615995004400050002x>.
- Venables, W.N., Ripley, B.D., 2003. *Modern Applied Statistics with S*. Springer Science & Business Media. Retrieved from. (<https://play.google.com/store/books/details?id=974c4vKurNkC>).
- Vereecken, H., Schnepf, A., Hopmans, J.W., Javaux, M., Or, D., Roose, T., et al., 2016. Modeling soil processes: review, key challenges and new perspectives. *Vadose Zone J.* <https://doi.org/10.2136/vzj2015.09.0131>.
- Wadoux, A.M.J.-C., Minasny, B., McBratney, A.B., 2020. Machine learning for digital soil mapping: Applications, challenges and suggested solutions. *EarthSci. Rev.* 210, 103359. <https://doi.org/10.1016/j.earscirev.2020.103359>.
- Wan, Z., Hook, S., Hulley, G., 2015. MOD11A2 MODIS/Terra Land Surface Temperature/Emissivity 8-Day L3 Global 1km SIN Grid V006 [Data set]. NASA EOSDIS Land Process. DAAC. <https://doi.org/10.5067/MODIS/MOD11A2.006>.
- Wang, Q., Wu, B., Stein, A., Zhu, L., Zeng, Y., 2018. Soil depth spatial prediction by fuzzy soil-landscape model. *J. Soils Sediment.* 18 (3), 1041–1051. <https://doi.org/10.1007/s11368-017-1779-0>.
- Whitney, A.W., 1971. A direct method of nonparametric measurement selection. *IEEE Trans. Comput. Inst. Electr. Electron. Eng. C.* 20 (9), 1100–1103. <https://doi.org/10.1109/t-c.1971.223410>.
- Yost, J.L., Hartemink, A.E., 2020. How deep is the soil studied – an analysis of four soil science journals. *Plant Soil* 452 (1), 5–18. <https://doi.org/10.1007/s11104-020-04550-z>.
- 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. <https://doi.org/10.1016/j.geoderma.2018.09.006>.
- Zhang, Y., Schaap, M.G., 2019. Estimation of saturated hydraulic conductivity with pedotransfer functions: a review. *J. Hydrol.* 575, 1011–1030. <https://doi.org/10.1016/j.jhydrol.2019.05.058>.
- Ziadat, F.M., Yegantham, D., Shoemate, D., Srinivasan, R., Narasimhan, B., Tech, J., 2015. Soil-Landscape Estimation and Evaluation Program (SLEEP) to predict spatial distribution of soil attributes for environmental modeling. *Int. J. Agric. Biol. Eng.* 8 (3), 158–172. <https://doi.org/10.25165/ijabe.v8i3.1270>.
- Zobeck, Fryrear, D.W., Pettit, R.D., 1989. Management effects on wind-eroded sediment and plant nutrients. *J. Soil Water Conserv.* 44 (2), 160. Retrieved from. (<http://www.jswconline.org/content/44/2/160.abstract>).