Review article

# What can artificial intelligence do for soil health in agriculture?

Stefan Schweng [a], Luca Bernardini [b,c], Katharina Keiblinger [b], Hans-Peter Kaul [c], Iztok Fister Jr. [d], Niko Lukač [d], Javier Del Ser [e,f], Andreas Holzinger [a] [iD],*

[a] Human-Centered AI Lab, Institute of Forest Engineering, Department of Ecosystem Management, Climate and Biodiversity, BOKU University, Vienna, Austria
[b] Institute of Soil Research, Department of Ecosystem Management, Climate and Biodiversity, BOKU University, Vienna, Austria
[c] Institute of Agronomy, Department of Agricultural Sciences, BOKU University, Vienna, Austria
[d] University of Maribor, Faculty of Electrical Engineering and Computer Science, Slovenia
[e] TECNALIA, Basque Research & Technology Alliance (BRTA), Spain
[f] Department of Mathematics, University of the Basque Country (UPV/EHU), Spain

## ARTICLE INFO

## ABSTRACT

The integration of artificial intelligence (AI) into soil research presents significant opportunities to advance the understanding, management, and conservation of soil ecosystems. This paper reviews the diverse applications of AI in soil health assessment, predictive modeling of soil properties, and the development of pedotransfer functions within the context of agriculture, emphasizing AI's advantages over traditional analytical methods. We identify soil organic matter decline, compaction, and biodiversity loss as the most frequently addressed forms of soil degradation. Strong trends include the creation of digital soil maps, particularly for soil organic carbon and chemical properties using remote sensing or easily measurable proxies, as well as the development of decision support systems for crop rotation planning and IoT-based monitoring of soil health and crop performance. While random forest models dominate, support vector machines and neural networks are also widely applied for soil parameter modeling. Our analysis of datasets reveals clear regional biases, with tropical, arid, mild continental, and polar tundra climates remaining underrepresented despite their agricultural relevance. We also highlight gaps in predictor–response combinations for soil property modeling, pointing to promising research avenues such as estimating heavy metal content from soil mineral nitrogen content, microbial biomass, or earthworm abundance. Finally, we provide practical guidelines on data preparation, feature extraction, and model selection. Overall, this study synthesizes recent advances, identifies methodological limitations, and outlines a roadmap for future research, underscoring AI's transformative potential in soil science.

## Contents

---

## 1. Introduction and motivation

Agriculture faces a dual challenge: meeting the rising global demand for food, feed, energy, and fiber while ensuring the sustainability of natural resources. The increasing demand for agricultural products is driven by population growth, higher per capita calorie consumption, and shifts in dietary preferences towards animal-based products, particularly in developing countries [1,2]. These trends exert significant pressure on agricultural systems to intensify production while minimizing environmental impacts.

A key strategy to address this challenge is sustainable intensification, which aims to increase productivity without disrupting natural nutrient cycles [3], degrading soil organic carbon (SOC) content [4] or impairing other soil properties [5]. Soil health is particularly critical, as an estimated 98.8% of the calories consumed by humans globally originate from soil, with only 1.2% derived from aquatic sources [6]. Furthermore, [2] predicts that global crop calorie production will rise by 47% from 2011 to 2050, driven by population growth, rising incomes, and evolving dietary preferences. However, intensification practices to meet growing demands often degrade soil quality, threatening long-term agricultural productivity and environmental stability [7, 8]. These factors underscore the necessity of innovative solutions to sustainably manage agricultural resources.

The European Green Deal [9], launched in 2019, represents the European Union's overarching strategy to transform the EU into a climate-neutral, resource-efficient, and competitive economy by 2050. A central pillar of this agenda is the EU Biodiversity Strategy for 2030 [10], which emphasizes restoring degraded ecosystems, improving soil fertility, and ensuring sustainable land use. Recognizing soil as a non-renewable resource critical to food security, biodiversity, and climate resilience, the EU has introduced the Soil Monitoring and Resilience Directive to ensure that all EU soils are in a healthy condition by 2050 [11]. This directive requires the establishment of harmonized monitoring systems across Member States, defines soil health indicators, and obliges national governments to develop plans for restoring degraded soils. These initiatives highlight the EU's determination to address soil degradation and foster sustainable land use, offering a large-scale example for similar efforts in a global context.

Achieving these political and environmental goals depends critically on the ability to monitor and model soil parameters with sufficient accuracy and spatial coverage. Traditionally, this has required a combination of labor-intensive field sampling, laboratory analysis, mechanistic modeling, and geostatistical methods. A foundational contribution in this field was made by [12], who proposed a conceptual model relating soil properties, such as organic carbon content, pH, and soil profile classes, to climatic, topographic, lithological, and temporal factors, along with biological influences from vegetation and human activity. This work was the basis for many quantitative mechanistic models,

as summarized by [13]. To address limitations in spatial resolution, spatial interpolation techniques such as *trend-surface* analysis [14] and geostatistical approaches like *kriging* [15–17] were developed. Variants such as *co-kriging* showed that some soil properties could be inferred from others [18,19]. Building on these ideas, [20] extended the model presented in [12] by incorporating both soil properties and spatial attributes into a unified predictive framework known as the *Scorpan* model.

In recent years, the increasing availability of high-resolution spatial, temporal, and spectral data has opened new avenues for more efficient and scalable soil health monitoring. This data abundance enables the development of data-driven modeling approaches that can capture complex, nonlinear interactions between soil properties, environmental factors, and management practices. In this context, Artificial Intelligence (AI), including machine learning (ML) and deep learning (DL) techniques, has emerged as a powerful tool for extracting actionable insights from large, heterogeneous datasets [21]. AI-driven models can enhance the prediction of soil parameters [22], improve the resolution of digital soil maps [23], and support decision-making processes in precision agriculture by offering data-driven, localized recommendations [24]. By complementing or extending traditional modeling methods, AI holds the potential to play a transformative role in achieving both the productivity and sustainability targets defined by recent agricultural and environmental policy frameworks.

A common challenge in data-driven soil property modeling is the uneven geographic distribution of available data. Dense coverage in some regions and sparse coverage in others can bias predictions, particularly when transferring models to new areas or scaling to broader regions. For instance, [25] found that drier regions are underrepresented in studies of soil microbial responses, which is problematic given the high sensitivity of microbial communities to moisture changes. Similarly, [26] showed that expanding the global soil respiration database [27] with measurements from previously underrepresented regions improved global representativeness but also increased model uncertainty by revealing greater variability, especially in tropical and southern hemisphere regions. These examples illustrate how geographic biases in datasets can shape both local soil property inferences and global-scale predictions, underscoring the importance of regionally balanced data collection.

Previous review papers have explored the use of AI in agriculture from various perspectives, including some that touch upon soil health. For example, [28] focus on ML-based estimation of soil indicators using remote sensing, while [29,30] review spectral and AI-based methods for predicting soil nutrients and diagnosing nutrient deficiencies. Other works address specific challenges such as soil pollution apportionment [31], input use efficiency [32], or soil carbon pools [33]. Broader technological perspectives include reviews on computer vision for food security [34], robotics in sustainable farming [35], and Internet of Things (IoT) applications for monitoring soil quality [36].

**Table 1**

Research questions for structured literature review.

| Research questions |
| --- |
| RQ1: What soil health challenges have been addressed using AI? |
| RQ2: Which AI methodologies and modeling techniques have been used to address these challenges? |
| RQ3: Does AI research in agriculture exhibit regional bias? |
| RQ4: What are the key research gaps and future opportunities for AI in improving soil health in agriculture? |

Several reviews offer regional or systemic overviews, such as digital transformation in agriculture [37], smart agriculture in developing countries [38], and Indian agricultural contexts [39]. Finally, [40] compare ML and DL approaches to conventional techniques for soil quality assessment. While these works offer valuable insights, they tend to focus on specific technologies, target narrow aspects of soil health, or provide high-level overviews without synthesizing the connection between AI methods, data availability, and soil-related challenges in a structured and actionable way.

In this work, we conduct a structured literature review (SLR) to synthesize the development of soil parameter modeling practices, with a particular focus on recent advances in AI. Distinct from prior reviews, we aim to bridge methodological and practical perspectives by making key AI techniques more accessible to non-experts and offering guidance on designing robust data-processing pipelines. Our review also contextualizes soil parameter modeling within the broader goals of sustainable agriculture, examining how AI is applied to address prevalent soil-related threats. Furthermore, we analyze global data availability for soil modeling and relate it to the degree of agricultural intensification across regions in order to highlight underrepresented areas. Finally, we summarize the general limitations of AI applications for soil health in agriculture and identify opportunities for future research. The specific research questions (RQs) addressed in this work are summarized in Table 1.

The remainder of this paper is structured as follows. Section 2 outlines the methodology, provides an overview of how AI is applied for soil health in agriculture and presents the results of our SLR, addressing RQ1 to RQ3. Section 3 addresses RQ4 by identifying general research gaps and outlining future research directions. Finally, Section 4 summarizes key findings and contributions of this work.

## 2. Artificial intelligence for soil health in agriculture: A structured literature review

In order to gather the state-of-the-art of AI applications for soil health in agriculture we performed an SLR and identified 115 articles to be reviewed. These articles were later distributed among six experts from the fields of Computer Science and Soil Science to be full-text reviewed. This section presents the SLR methodology (Section 2.1) as well as an overview of the review findings, addressing RQ1 (Section 2.3.1), RQ2 (Section 2.3.2) and RQ3 (Section 2.3.3).

### 2.1. Methodology

Inspired by the SLR processes described in [41,42] the SLR conducted in this paper consists of the following steps: (1) The definition of RQs (see Table 1), (2) identification of relevant articles by selecting appropriate search keywords and synonyms, followed by (3) article filtering based on abstract screening and predefined inclusion criteria (ICs) and exclusion criteria (ECs). Last but not least (4) the extraction of data from the remaining articles using data extraction questions (DEQs). The output of step (4) will eventually be used to answer the RQs defined in step (1).

**Table 2**

Keywords and synonyms for article search.

| Keyword | Synonyms |
| --- | --- |
| Artificial intelligence | AI, Machine learning, ML |
| Agriculture | Field, Crop, Farm |
| Soil health | Soil quality |

**Table 3**

Inclusion criteria for article search and filtering.

| Criterion | Condition |
| --- | --- |
| IC1 | Article was published in peer-reviewed journal or other types of reputable sources. |
| IC2 | Article is written in English. |
| IC3 | Article presents an AI approach to tackle a problem in agriculture related to soil health. |

**Table 4**

Exclusion criteria for article search and filtering.

| Criterion | Condition |
| --- | --- |
| EC1 | Article is a review, survey or meta-analysis. |
| EC2 | Article is not accessible (last attempt of access: January 21st, 2025). |
| EC3 | Article is marked with expression(s) of concern from the journal's editorial board. |

#### 2.1.1. Definition of research questions

Planning a SLR involves defining the RQs to be answered. In this work the SLR aims at answering the questions defined in Table 1. All following steps will be guided by these questions.

#### 2.1.2. Identification of relevant articles

The search engine used for identifying relevant publications is Web Of Science (WoS).[1] The next step in the SLR is defining keywords and respective synonyms to be used for the article search. Table 2 shows the list of selected keywords and synonyms.

The search string used in WoS was: TS=((''artificial intelligence'' OR AI OR ''machine learning'' OR ML) AND (agri* OR crop* OR field* OR farm*) AND (''soil health'' OR ''soil quality'')). This search string queries articles by analyzing the articles' title, abstract, keywords and KeyWords Plus, where the latter is a set of keywords derived algorithmically by WoS from the titles of the article's cited references. The search terms in the string represent the search keywords and synonyms defined in Table 2.

The WoS article search was carried out on December 2nd, 2024 and resulted in 334 articles matching the search string specified above. Note that there was no date range specified for the publication date. In order to reproduce the search results the index date range 1900-01-01 to 2024-12-01 must be specified in the WoS advanced search feature. Fig. 1 shows the number of matched articles per publication year.

#### 2.1.3. Article filtering

Articles included in the final set of papers for full-text review must meet all ICs outlined in Table 3 and must not meet any of the ECs detailed in Table 4.

IC1 is inherently fulfilled by using WoS as a search engine. Every article indexed in the *Web of Science Core Collection* must be published in a peer-reviewed journal, with just a few exceptions of journals that are *generally long-established*.[2] IC2 and IC3 are ensured by title and abstract screening of the initial set of 334 articles matching the

---

[1] Web Of Science: https://www.webofscience.com, accessed May 7th, 2025.
[2] WoS Core Collection: https://support.clarivate.com/ScientificandAcademicResearch/s/article/Web-of-Science-Core-Collection-Explanation-of-peer-reviewed-journals, accessed January 7th, 2025.
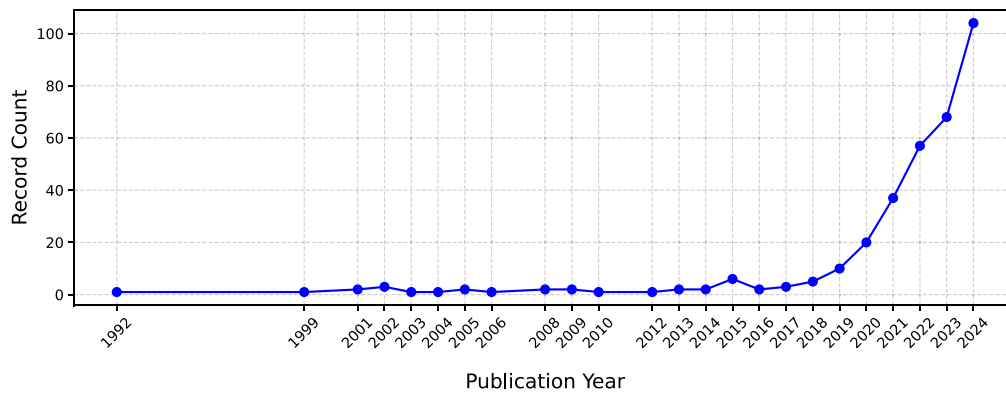
**Fig. 1.** Number of articles per publication year matching the search string in WoS. A total number of 334 articles was found, where the first publication was published in 1992. The WoS search was performed on December 2nd, 2024.
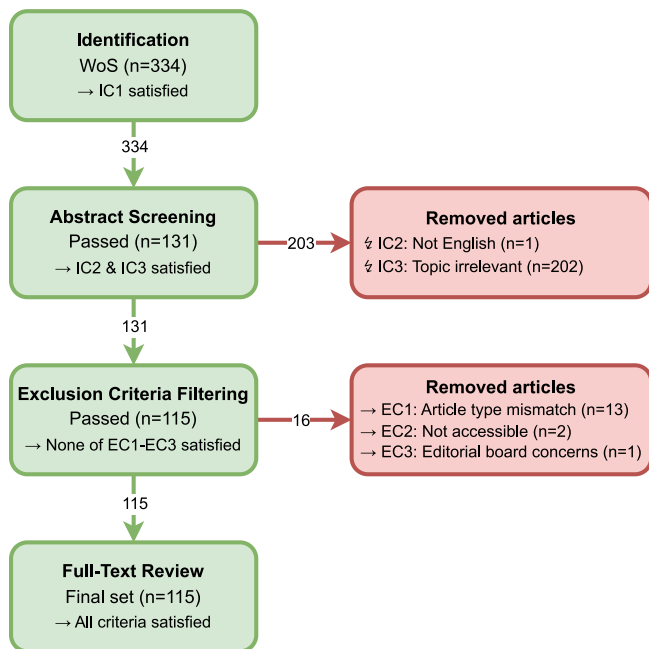


**Fig. 2.** PRISMA-inspired flowchart [44] illustrating the selection process of articles for the full-text review of the SLR. The initial search in the WoS database identified 334 articles. During title and abstract screening using ASReview, 131 articles met the ICs outlined in Table 3. Finally, the ECs defined in Table 4 were applied, which lead to the exclusion of 16 articles. This process yielded a final set of 115 articles for full-text review.

WoS search string. Three experts from the fields of Computer Science and Soil Science participated in this initial screening process using ASReview [43]. All 334 titles and abstracts have been reviewed by one expert and filtered based on IC2 and IC3, which resulted in a set of 131 articles.

In accordance with the ECs, thirteen articles have been excluded because they were identified as review or survey papers (EC1), two articles were excluded because they were not accessible (EC2) and one article was excluded due to an expression of concern regarding the integrity of the paper (EC3). This results in a final article set for full-text reviews containing 115 articles. The full selection and filtering process is visualized as a flowchart diagram in Fig. 2 adapted from the PRISMA method [44].

### 2.1.4. Data extraction

To conduct the full-text reviews, the article set of 115 papers was distributed among six experts from the fields of Computer Science and

Soil Science. To ensure consistency in the extracted information, a set of DEQs was defined to guide the review process, helping reviewers focus on a common set of aspects. Each DEQ corresponds to a specific RQ. The complete list of DEQs and their related RQs is provided in Table 5.

While reviewing an article, experts were presented with each DEQ in a Google Form, and each article was reviewed by a single expert. All DEQs were optional, allowing experts the flexibility to skip questions that were not applicable or could not be confidently answered. The collected responses were exported to an Excel file. Since eight out of nine questions included free-text answer fields, a post-processing step was performed using a Python script to categorize and group related responses. The resulting Excel file, containing the post-processed DEQ responses for each of the 115 articles, is provided in Appendix A.

### 2.2. How to apply artificial intelligence for soil health challenges?

The majority of articles start by describing the response variables (to be predicted) and the predictor variables used in the analysis. After this, the authors typically follow these steps: (1) data collection, (2) preprocessing, (3) modeling, and (4) prediction. To provide a schematic overview, Fig. 3 illustrates this general workflow.

#### 2.2.1. Data collection

Data collection usually involves integrating multiple data sources. The most common sources include public databases, in-situ measurements (e.g., soil parameters in the study area), agricultural management data and remote sensing data. The choice of data sources depends on the specific task, regional scope, and the desired resolution in time, space, or both.

While publicly available satellite data (e.g., Sentinel or Landsat) and climatic data (e.g., WorldClim) are easily accessible for global-scale applications, obtaining ground truth data for soil parameters is more challenging. For regional (i.e., country- or continent-level) or global studies, soil parameter databases such as the Land Use and Coverage Area Frame Survey (LUCAS) or the International Soil Reference and Information Centre (ISRIC) provide sparse point-scale soil data. However, downscaling soil parameters to high-resolution temporal or spatial grids remains difficult and is often addressed by incorporating data from local weather stations [45] or by using high-resolution soil maps of other *Scorpan* covariates [23].

At the sub-national scale, studies often rely on regional databases such as the Losan database [46] or the South Dakota Geological Survey.[3] An additional example, not mentioned in the reviewed literature,

---

[3] South Dakota Geological Survey: https://www.sdgs.usd.edu, accessed February 4th, 2025.

**Table 5**
DEQs used to extract relevant data from each article in the SLR and to answer the specified RQs.

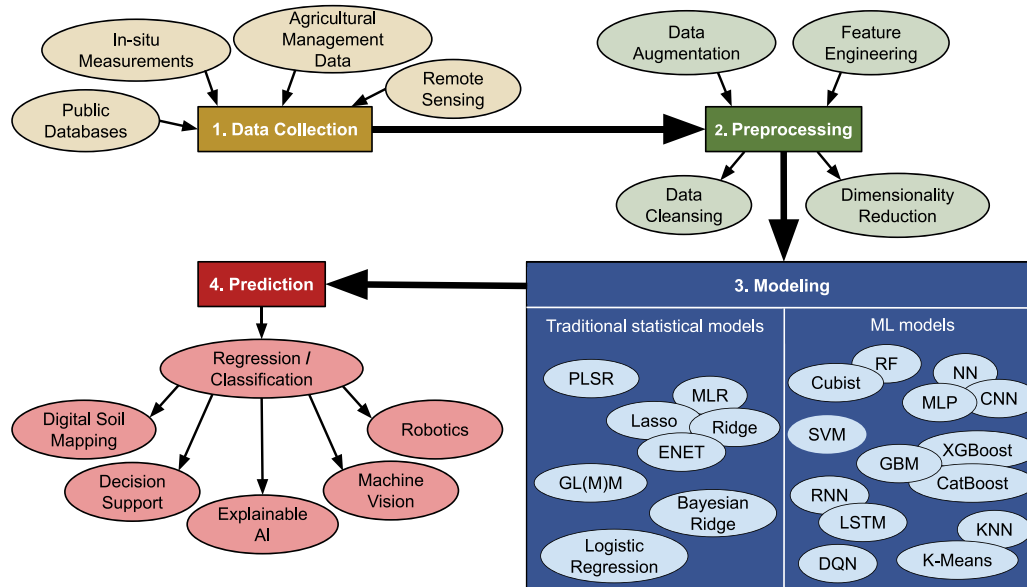| DEQ | Question | Related RQ |
|---|---|---|
| DEQ1 | Which soil threats were addressed by the presented research? | RQ1 |
| DEQ2 | Which soil quality indicators were assessed? | RQ1 |
| DEQ3 | Which AI algorithms and/or modeling techniques were used in the presented approach? | RQ2 |
| DEQ4 | Which type of input data was used for the presented method? | RQ2 |
| DEQ5 | Which type of ground truth data was used? | RQ2 |
| DEQ6 | What is the sample size of the data used in this article? | RQ2 |
| DEQ7 | What is the name of the database serving the method's input data? | RQ2 |
| DEQ8 | What is the regional scope of the presented method? | RQ3 |
| DEQ9 | What are the geographic origins of the assessed data points? | RQ3 |



**Fig. 3.** Basic workflow of the data pipeline commonly presented in the investigated articles.

is the Austrian digital soil map *eBod*,[4] which summarizes soil types, parent material, and key hydrological, physical, and chemical properties, along with agricultural value. While these datasets enable detailed local modeling, transferring such models to other regions remains difficult due to substantial regional variability [47]. Section 3.1.2 discusses this challenge in greater depth and outlines potential strategies to address it.

On a local scale, studies often rely on long-term crop experiments conducted by research institutions. These experiments provide records of weather conditions, agricultural management practices, harvested yields, and in-situ soil measurements. Typical examples of weather parameters are daily maximum and minimum air temperature, precipitation, and solar radiation. Management data often include the amounts of nitrogen fertilizer or manure applied per hectare. Soil parameters commonly measured in-situ are bulk density, soil texture, and soil mineral nitrogen content.

More details on the types of data collected across different regional scales in the reviewed studies are discussed in Section 2.3.1.

*2.2.2. Preprocessing*
Common steps in data preprocessing include data augmentation (i.e., generating additional artificial samples) or feature engineering (i.e., creating new based on existing features), data cleansing (e.g., standardization, duplicate/outlier removal), and dimensionality reduction

---

[4] eBod soil map: https://www.bodenkarte.at, accessed September 9th, 2025.

(i.e., feature selection to improve predictive performance and training efficiency while mitigating overfitting).

Incomplete and unbalanced datasets present a major challenge in data driven AI because they compromise the reliability, generalizability, and interpretability of models. Missing data can lead to biased parameter estimation, reduced statistical power, and unreliable predictions [48]. Preprocessing is necessary to address these issues through techniques such as imputation of missing values, normalization, class balancing (e.g., oversampling or SMOTE [49]), and feature selection. This ensures data quality, reduces model bias, and enhances the robustness of downstream predictive and interpretative tasks, ultimately enabling more accurate and ecologically valid insights.

Examples of feature engineering include the study by [50], which uses Sentinel-2 spectral bands to compute six vegetation indices that serve as additional predictors for mapping the regional distribution of cover crops. Another example is the work of [51], which simulates potential evapotranspiration of sorghum–sudangrass from weather variables, soil properties, and plant characteristics, thereby creating input features for biomass yield prediction. On the other hand, a typical example of data augmentation in machine vision involves expanding training datasets through image transformations such as random cropping, rotation, or adding artificial noise [52].

Dimensionality reduction techniques are also widely applied. For example, [53] use Principal Coordinate Analysis (PCoA) and the t-SNE algorithm to reduce the dimensionality of soil microbiome DNA sequencing data, while [54] apply Principal Component Analysis (PCA) to spectral data. In addition, feature selection approaches such as Random Forest Recursive Feature Elimination (RF-RFE) and Guided

Regularized Random Forest (GRRF) are used to identify compact, informative subsets of features by leveraging the importance scores derived from random forest (RF) models [55,56].

Data cleansing is a key step in preprocessing that enhances data quality and model performance. Standardization, which adjusts features to have zero mean and unit variance, helps ensure that each feature contributes equally during model training and also improves the convergence dynamics of many ML algorithms by enabling faster and more stable optimization. Outlier filtering is especially important for models sensitive to extreme values, such as linear regression, where outliers can disproportionately influence results. Imputation techniques are also used to address missing values, with methods ranging from simple mean substitution to more sophisticated model-based approaches. The relevance of outlier filtering and imputation depends on the model type, as some models are more robust to such issues than others.

Handling categorical data is a critical preprocessing step that directly affects how models learn. Distinguishing between nominal and ordinal variables ensures appropriate encoding. Nominal features such as crop or land use types lack inherent order and are best represented using one-hot encoding to avoid implying hierarchy. In contrast, ordinal features, like erosion risk levels or soil quality ratings, have a natural order and require encoding strategies that preserve this structure. Improper handling can lead to models misinterpreting category relationships or missing important ordinal information.

### 2.2.3. Modeling

Modeling refers to the process of learning representations of the relationship between one or more response variables and a set of predictor variables, often referred to as *covariates* or *features*. For example, bulk density (the response variable) can be modeled using spectral features [57] or elevation and mean annual precipitation among other covariates [58].

In general, the models used for modeling this process can be categorized into traditional statistical and ML models, as illustrated in Fig. 3. The models listed in the box of the third step (*Modeling*) represent the models most commonly applied in the reviewed studies. Overlapping ellipses indicate models that are related or of a similar type.

Traditional statistical models are based on predefined assumptions about data distributions and the relationships between predictor and response variables. They estimate parameters through statistical techniques such as maximum likelihood estimation and follow a parametric approach with fixed functional forms. The reviewed studies employed methods such as partial least squares regression (PLSR), multiple linear regression (MLR) and its regularized variants – lasso, ridge regression, and elastic net (ENET) – as well as generalized linear (mixed) models (GL(M)Ms). Bayesian ridge regression was also applied for regression tasks, while logistic regression was used for classification tasks.

In contrast, ML models do not require explicit distributional assumptions and rely on optimization techniques rather than closed-form parameter estimation. Whereas traditional models emphasize interpretability, ML approaches prioritize predictive performance by directly learning patterns from data. In the reviewed literature, commonly used methods included tree-based models such as RF and Cubist, neural networks (NNs) – often implemented as multilayer perceptrons (MLPs) for regression tasks or convolutional neural networks (CNNs) for computer vision tasks – support vector machines (SVMs), and gradient boosting machines (GBMs) such as XGBoost and CatBoost. Other approaches included recurrent neural networks (RNNs) such as Long Short-Term Memory (LSTM) networks, clustering algorithms (k-nearest neighbors and k-means), and reinforcement learning techniques such as Deep Q-Networks (DQNs).

The choice of model typically depends on factors such as dataset size, feature dimensionality, and the complexity of the relationship between predictors and response variables. Further guidance on selecting models for specific tasks is provided in Section 3.3.2.

### 2.2.4. Prediction

The final step of the schematic workflow, *Prediction*, usually involves regression, classification, or hybrid approaches. Regression is used in cases where the response variable is numerical, while classification typically refers to predicting categorical values. Hybrid approaches arise when categorization is based on numerical predictions. Predictions from regression and classification models are commonly applied in digital soil mapping (DSM), decision support systems (DSS), and fields such as robotics and machine vision [59].

When examining the specific tasks addressed in recent literature, DSM appeared prominently, likely due to the complexity and cost of traditional soil attribute measurement methods. Remote sensing for SOC mapping has emerged as a particularly strong trend [60–63]. Nevertheless, researchers have also applied DSM techniques to characterize other soil properties, such as salinity [64,65], pH value [66,67], soil texture [67], available water capacity [45] or soil erodibility [68].

DSS were also widely discussed, with applications ranging from crop rotation planning [24] to IoT-based platforms for monitoring soil health and crop performance [69]. Other systems typically focus on optimizing input use efficiency for fertilization, irrigation, and weed management [70]. For instance, [71] developed a method to delineate soil management zones using RF and soil properties, thereby enabling site-specific fertilizer recommendations.

In addition, explainable AI (XAI) [72,73] has gained importance for enhancing the transparency and interpretability of model predictions. For example, [74] combined RF with Shapley Additive Explanations (SHAP) [75] to analyze pesticide effects on earthworm lethality, while [51] applied SHAP values to evaluate how environmental variables and management practices influence biomass yields.

Less frequently covered applications among the reviewed articles include machine vision tasks based solely on RGB images and robotics. For example, [52] employed the Inception V3 model to detect potato plant diseases from RGB leaf images, while [76] developed a prototype of a cotton harvesting robot.

### 2.3. Literature study and responses to the research questions

The following sections address RQ1–RQ3 in detail. These responses form the foundation for identifying research gaps (RQ4), which are discussed in Section 3. A summary of the data extraction results is shown in Fig. 4.

#### 2.3.1. Addressed soil threats and assessed quality indicators

DEQ1 focused on the soil threats addressed in each article. The corresponding question offered eight predefined options, based on the definitions by [77]. The most frequently cited threat was soil organic matter (SOM) decline, which appeared in a significant majority of responses. Other common threats included compaction, biodiversity loss, and erosion. Contamination was also frequently noted, typically referring to heavy metals (e.g., copper, lead, cadmium), metalloids (arsenic), and water pollution due to eutrophication. Less frequently mentioned threats included salinization, landslides and floods.

DEQ2 examined which soil quality indicators were considered. Respondents could select from sixteen predefined indicators, also derived from [77]. The most frequently mentioned indicator was organic matter, reflecting its importance in evaluating soil health. Organic matter includes both stable and labile fractions. The labile fraction is particularly responsive to conservation farming practices [78], making it especially relevant in the context of this work. Since SOM is typically not measured directly, it is commonly estimated from measured SOC using the Van Bemmelen factor. Despite some debate about its accuracy [79], this conversion remains widely applied in practice.

Other commonly cited indicators included chemical properties, particularly pH, electrical conductivity, and cation exchange capacity, as well as macronutrients, with nitrogen, phosphorus, and potassium most frequently mentioned. Soil structure also featured prominently, with
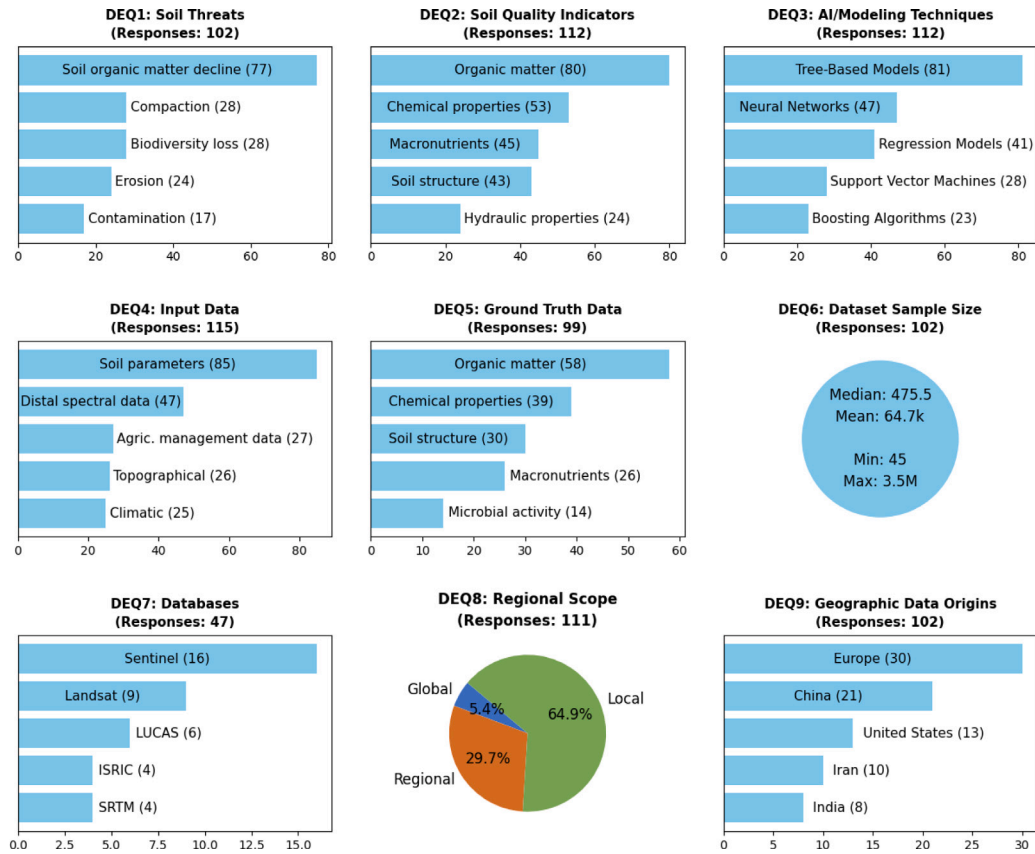
**Fig. 4.** Overview of the data extraction results. The total number of responses for each DEQ is indicated in parentheses below the chart titles. Bar charts present the five most frequent response categories, with the number of mentions shown in parentheses. DEQ8 was posed as a single-choice question; its response distribution is illustrated using a pie chart. For DEQ9, statistics on the size of the utilized datasets are provided.

texture and bulk density as key attributes. Hydraulic properties, such as water storage capacity and soil moisture, were also often reported.

Additional indicators included microbial parameters (e.g., microbial biomass, soil respiration), contaminants (mainly heavy metals), micronutrients (e.g., iron, copper, zinc, boron), and inorganic carbon (calcium carbonate, bicarbonate). Less frequently mentioned were erosion risk and salinity, and only a few responses mentioned earthworms as a soil quality indicator.

### 2.3.2. Modeling techniques and dataset characteristics

This section summarizes key trends in AI techniques and dataset characteristics. Insights are drawn from evaluation of DEQ3–DEQ7, covering aspects such as modeling approaches, input and ground truth data types, dataset sizes, and referenced data sources.

In what refers to DEQ3, a wide range of modeling techniques was reported. Tree-based models (e.g., RFs) stood out as the most frequently used, followed by NNs, regression models (e.g., MLR or PLSR), SVMs and boosting algorithms. Less frequently mentioned methods were clustering, nature-inspired optimization techniques, and RL approaches such as Q-learning.

Regarding input data types (DEQ4), soil parameters (e.g., SOC, pH, bulk density) were the most commonly utilized, followed by distal spectral data collected via satellites or drones. Agricultural management data, topographic features, climatic conditions, and proximal spectral data from handheld sensors also played important roles. Pure RGB imagery, land use information, and RGB-D data were rarely cited, the latter being mentioned only once in a robotics application for cotton harvesting [76]. Regarding DEQ5, the ground truth data categories reflected a strong emphasis on organic matter, particularly SOC, as well as chemical properties like pH, electrical conductivity, and cation exchange capacity.

Soil structure features such as texture and bulk density, macronutrients (e.g., nitrogen, phosphorus, potassium), microbial activity, and hydraulic properties were also common. Other categories such as micronutrients, inorganic carbon, agricultural management outputs (e.g., crop rotations), contaminants, salinity, and erosion risk were cited less frequently.

Figs. 5 and 6 highlight the interplay between modeling techniques, input data types, and prediction targets (i.e., ground truth data). Regression techniques (i.e., traditional statistical models), NNs, SVMs and especially tree-based models were frequently applied to predict organic matter, chemical properties, soil structure, and macronutrients. Similarly, spectral data, agricultural management data, and especially soil parameters emerged as key inputs for these predictions. Notably, the use of soil parameters to infer other soil properties, such as organic matter or soil structure, emphasizes the relevance of pedotransfer functions and process-based modeling approaches.

RF was frequently chosen as the modeling technique due to several advantages: its simplicity of parameterization [80], ability to reduce bias towards dominant predictors [81], robustness against overfitting [82], and built-in support for feature selection [83]. The latter is particularly important for the common use case of predicting SOC based on distal spectral data, where the challenge lies in selecting relevant covariates from a large number of spectral bands [84]. Moreover, RF is often employed to analyze the relative importance of individual predictors, making it valuable for model interpretation [60].

Regarding DEQ6, dataset sizes varied widely across studies. Although the median dataset size was below 500 samples, the mean was substantially higher due to a few studies using very large datasets. Most models, including SVMs, regression models, tree-based models, and NNs, were typically applied to datasets ranging from a few hundred to a few thousand samples. In contrast, boosting algorithms were, on
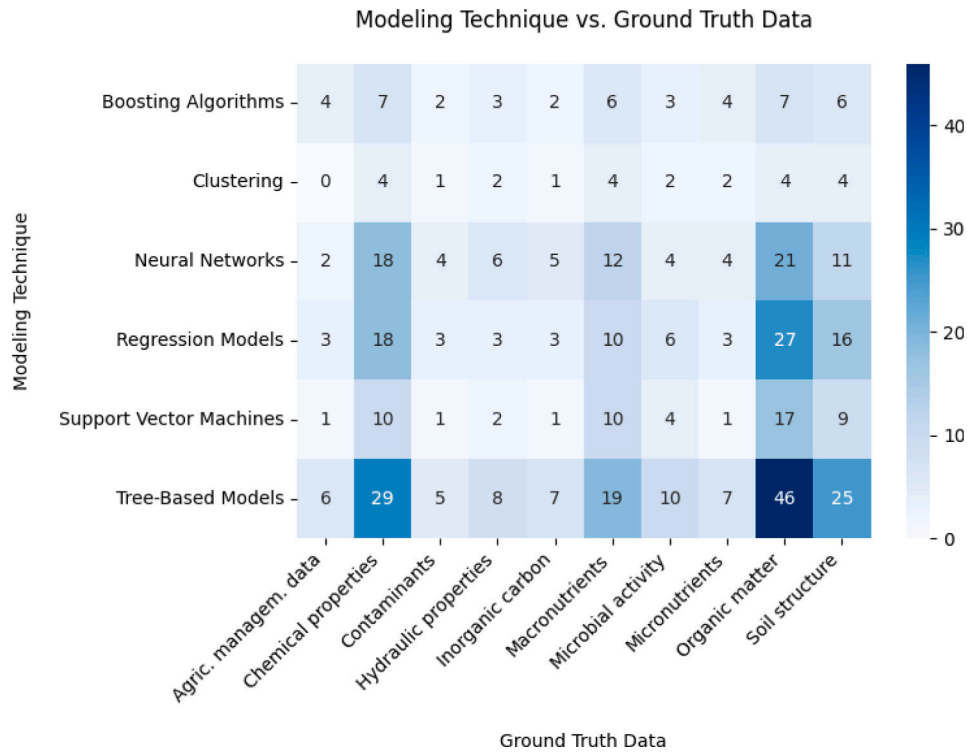
**Fig. 5.** Heat map showing the AI algorithms/modeling techniques referenced in DEQ responses in combination with specific ground truth data types.
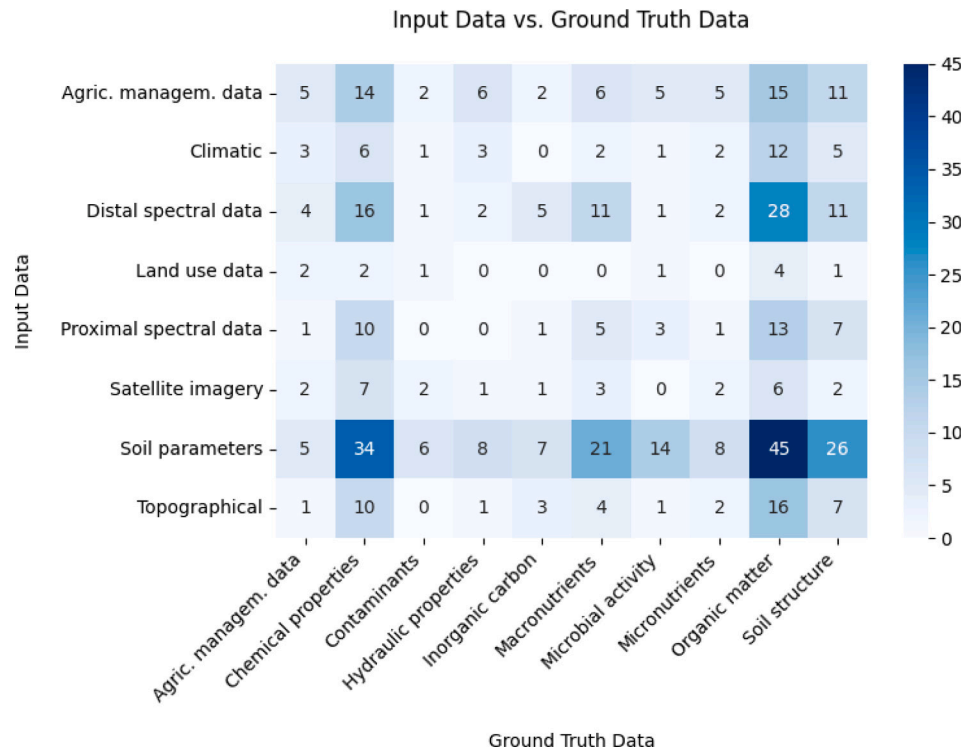


**Fig. 6.** Heat map showing input data types referenced in DEQ responses in combination with specific ground truth data types.

average, associated with slightly larger datasets, as shown in Fig. 7(a). It is important to note that most reported boosting algorithms are in fact tree-based, since the most commonly used methods, such as XGBoost [85] and CatBoost [86], are based on sequentially combining decision trees. Despite this, the observed trends in model choice might be incidental rather than technically motivated, even though 5000 bootstrap resamples were used to reduce sampling variability. For instance, NNs are generally expected to perform better with larger datasets [87], yet were not predominantly applied in such cases. One possible explanation for the relatively more frequent use of boosting algorithms on larger datasets may lie in their built-in support for predictor importance analysis, which improves model interpretability [83].
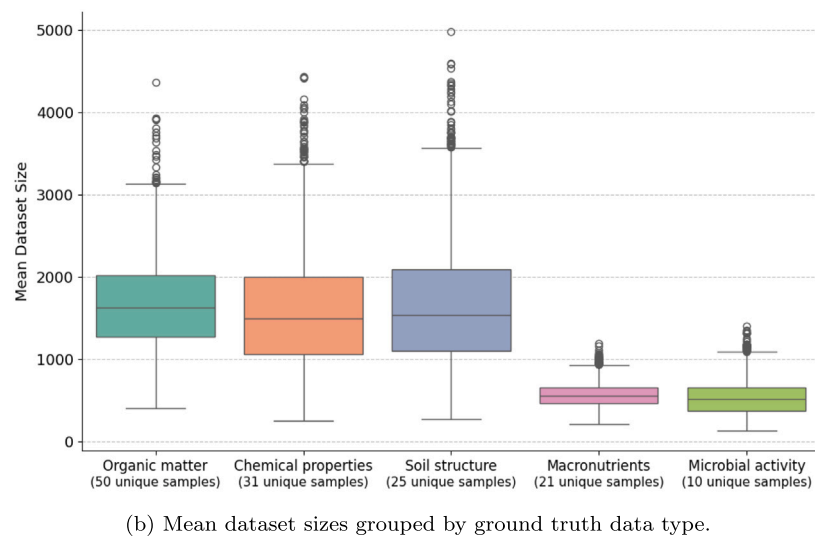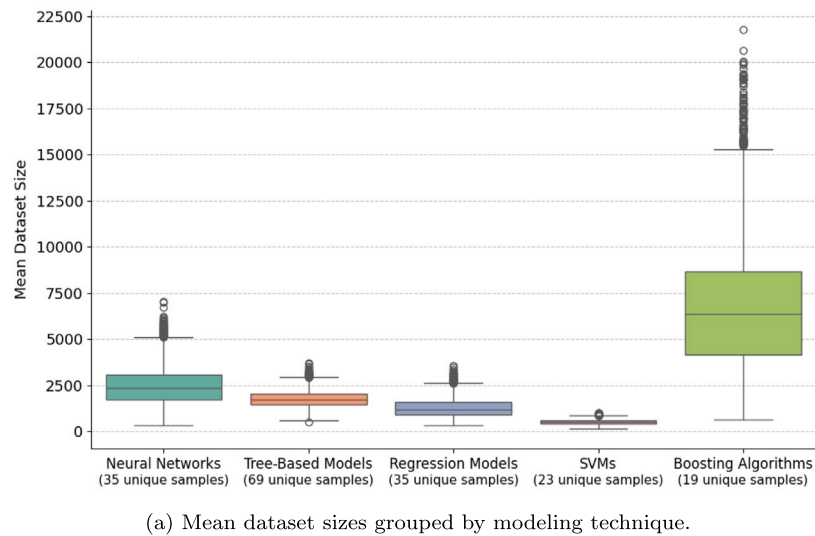
(a) Mean dataset sizes grouped by modeling technique.



(b) Mean dataset sizes grouped by ground truth data type.

**Fig. 7.** Box plots of mean dataset sizes estimated via bootstrapped resampling (5000 resamples). Prior to resampling, original dataset sizes were restricted to the 90th percentile to reduce the influence of extreme outliers, which is common strategy in exploratory data analysis [89]. The number of unique samples used for bootstrapping is given in parentheses below each modeling technique or ground truth data type label.

Furthermore, boosting algorithms are often considered easier to tune than NNs, particularly in terms of hyperparameter optimization [88], making them especially attractive for applied research settings.

A similar analysis was conducted for the various types of ground truth data. Fig. 7(b) illustrates the relationship between dataset size and the targeted soil property, based on the reviewed studies. The trends indicate that research on organic matter, chemical properties, and soil structure typically employed moderately sized datasets with a mean of approximately 1500 samples, while studies focusing on macronutrients or microbial activity more frequently relied on smaller datasets with a mean of about 500 samples.

In terms of data availability (DEQ7), several commonly used and publicly available data sources were identified. Sentinel[5] and Landsat[6] were the most frequently cited, followed by LUCAS,[7] ISRIC,[8] and the

Shuttle Radar Topography Mission (SRTM).[9] Although Sentinel and Landsat are technically satellite missions, they were frequently referenced as primary sources for remote sensing data. Additional mentions included Corine Land Cover[10] and WorldClim.[11]

### 2.3.3. Regional scope and geographic distribution of datasets

The regional scope of the studies was examined through DEQ8, which asked respondents to classify the spatial extent of their research as local, regional, or global. The majority (65%) indicated a local scope, referring to sub-national or site-specific studies. A further 30% reported a regional scope, representing research conducted at the national or continental level. Only 5% of the studies addressed issues at a global scale.

In DEQ9, the geographic origins of the datasets used in the literature were investigated. As shown in Fig. 4, the top five regions in terms of

---

[5] Sentinel: https://sentinel.esa.int/web/sentinel/missions/sentinel-2, accessed May 8th, 2025.

[6] Landsat: https://www.usgs.gov/landsat-missions, accessed May 8th, 2025.

[7] LUCAS: https://ec.europa.eu/eurostat/web/lucas, accessed May 8th, 2025.

[8] ISRIC Data Hub: https://data.isric.org/, accessed May 8th, 2025.

[9] SRTM: https://www.earthdata.nasa.gov/data/instruments/srtm, accessed May 8th, 2025.

[10] Corine Land Cover: https://land.copernicus.eu/en/products/corine-land-cover, accessed May 8th, 2025.

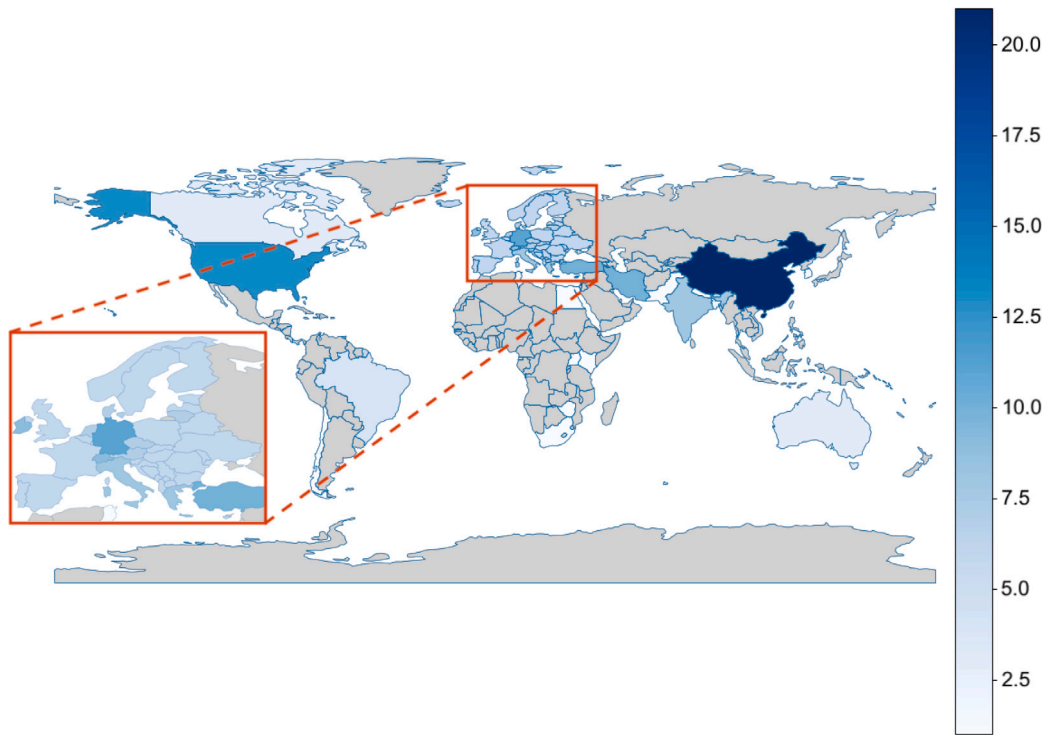[11] WorldClim: https://worldclim.org, accessed May 8th, 2025.

**Fig. 8.** A world map showing the number of dataset origins per country. In contrast to DEQ9 in Fig. 4, European countries are considered individually.

dataset origin are Europe, China, the United States, Iran, and India. Notably, Europe as a region aggregates all mentions of individual European countries as well as general references to the continent.

Fig. 8 provides a global visualization of dataset origins mentioned in the DEQ responses. Well-represented regions include Canada, the United States, Brazil, Europe, Iran, India, Bangladesh, China, and Australia. The extensive coverage of European countries is primarily attributable to the LUCAS database.

To identify underrepresented regions in terms of data availability, we utilize 1-km Köppen–Geiger climate classification maps [90]. This approach is based on the premise that similar climate zones tend to exhibit similar ecosystems and, consequently, similar representation in existing datasets. Since this study focuses on soil parameter modeling in an agricultural context, we emphasize regions of high agricultural relevance, as estimated by harvested area of primary crops (FAOSTAT, 2023).

Fig. 9 compares the proportion of global harvested area per climate zone (blue bars) with the proportion of dataset origins from the literature (orange bars). Climate zones with less than 1% of global harvested area are excluded from the analysis as they are considered less relevant for this context. We define underrepresented climate zones as those where the proportion of harvested area exceeds the proportion of literature-referenced datasets.

Fig. 10 visualizes the underrepresented climate zones. These zones span large areas of Central and South America, Africa, Central and South Asia, and Australia. They are primarily characterized by tropical, arid (dry), and mild continental climates, as well as regions of polar tundra.

## 3. Limitations and research opportunities

This section identifies gaps in the current literature by first outlining general limitations of state-of-the-art research (Section 3.1) and then exploring specific research opportunities (Section 3.2) and workflow suggestions (Section 3.3) for potential future studies.

### 3.1. General research gaps for soil health applications in agriculture

Our analysis highlights several research gaps and opportunities for advancing AI applications in soil health.

#### 3.1.1. Limited use of advanced artificial intelligence techniques

In terms of AI methodologies, RF dominated the field, mentioned in 68% of responses. In contrast, more advanced approaches, such as RL and Graph Neural Networks (GNNs), were rarely applied. RL was cited in only three DEQ responses, yet it offers promising opportunities for DSS, especially in areas like crop rotation planning and fertilizer application. When guided by rule-based constraints (or guardrails [91]), RL agents can be aligned with specific soil health requirements, enabling more controlled and sustainable decision-making. However, the limited use of RL in current literature may stem from the lack of realistic simulators and the delayed, often sparse, reward signals in agricultural environments. Moreover, the risk associated with trial-and-error learning on real-world farms makes RL challenging to apply in practice.

Similarly, GNNs were referenced in just one study, despite their strong potential for modeling spatiotemporal autocorrelation. This is a highly relevant aspect for soil parameter models, as demonstrated by [92]. The limited uptake of GNNs in current literature suggests an underexplored opportunity to better capture complex spatial and temporal dependencies in agricultural and soil health data. This underutilization may be due to the additional effort required to construct graph structures from spatial data, combined with the technical challenges and limited interpretability associated with GNNs in practice.

Given that spatiotemporal autocorrelation leads to overlapping patterns in soil parameter modeling [21], multi-task learning (MTL) presents a promising yet largely overlooked approach. MTL addresses this by incorporating shared sub-models that capture common structures across related prediction tasks. Despite its potential to enhance both model efficiency and accuracy, only one study explicitly applied an MTL approach [54]. By enabling joint learning of multiple
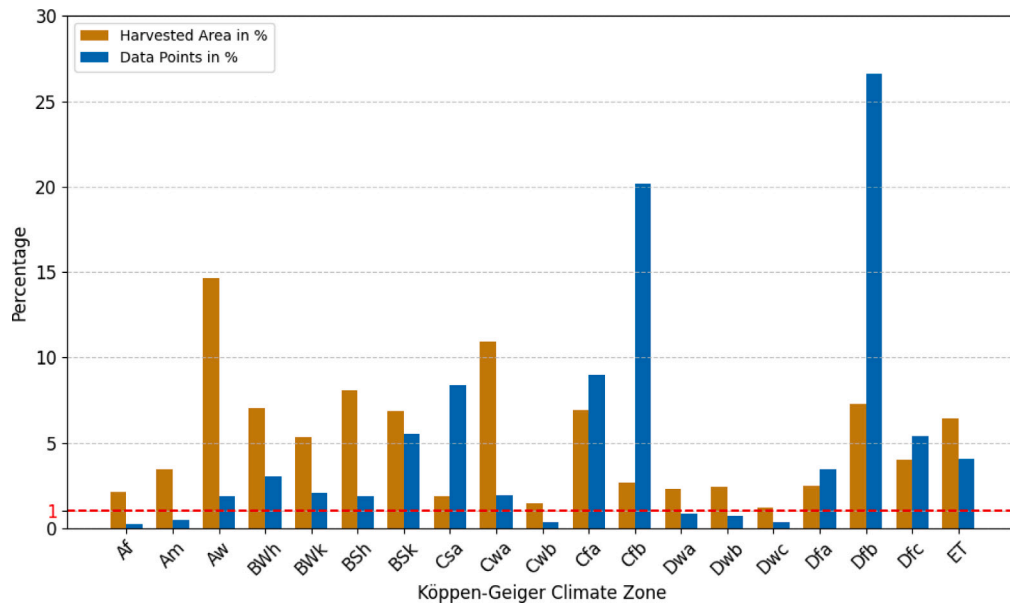
**Fig. 9.** Comparison of primary crop harvested area (FAOSTAT, 2023) and dataset origins across Köppen–Geiger climate zones. Climate zones representing less than 1% of the total harvested area were excluded. For a description of the climate zones shown, refer to the legend in Fig. 10.
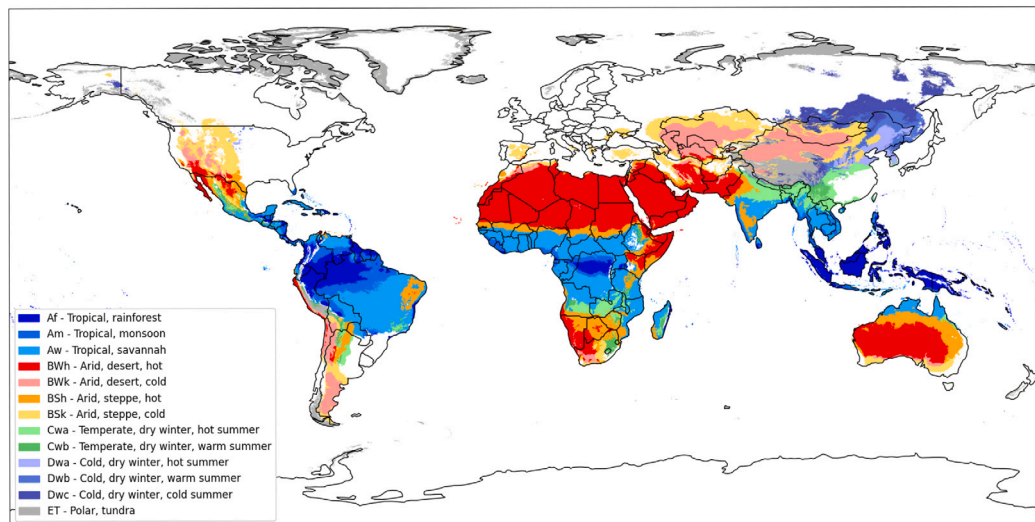


**Fig. 10.** Köppen–Geiger climate zones with at least 1% of global primary crop harvested area (FAOSTAT, 2023) and underrepresented in existing datasets from the literature (i.e., harvested area exceeds data point share; see Fig. 9). Based on 1991–2020 climate data from [90].

soil-related targets, MTL could support the development of more holistic and adaptable models, better aligned with real-world agricultural decision-making needs. The limited adoption of MTL may be due to the challenge of obtaining well and consistently labeled datasets for multiple related tasks, as well as the risk of negative transfer when task relationships are poorly understood or imbalanced.

Another noteworthy finding is the potential of hybrid modeling techniques that integrate data-driven approaches, such as RFs or DL, with process-based models like RothC [93] and C-Tool [94] for SOC modeling. There are various ways of implementing hybrid models: through upstream process-based models followed by data-driven techniques aiming to minimize residuals [95], by incorporating physical laws into the neural network's loss function [96], through physics-based augmentation of datasets [97,98], or by embedding physics into the neural network architecture itself [99]. These approaches offer a promising avenues for improving both the accuracy and interpretability of predictions. Yet, hybrid modeling remains complex, often requiring deep interdisciplinary expertise and significant effort to couple legacy

simulation models with modern ML frameworks. The lack of standardized tools or workflows further hinders widespread adoption in the soil health domain.

In addition, time-series modeling remains underutilized in soil health applications, despite its clear relevance. These models are designed to process sequential data and make temporally-aware predictions, which is an essential capability given the time-dependent nature of many agricultural and soil processes. While a few studies employed LSTMs [100] for tasks such as disease and yield prediction [69] or forecasting soil moisture and temperature [101], the overall uptake remains modest. Other time-series approaches, including Facebook Prophet [102,103], echo state networks (ESNs) [104], state–space models (SSMs) [105], and transformer-based models like Informer [106] have seen minimal use, though they offer powerful capabilities for capturing complex temporal dynamics. This may be partly due to the limited temporal granularity and irregularity of soil datasets, which restrict the applicability of time-series techniques.

Additionally, missing data and noisy measurements complicate model training and validation.

Furthermore, another largely underexplored field is the intersection of machine vision, robotics, and soil health, such as efforts to reduce soil compaction through the use of lightweight robots instead of heavy harvesting machinery. This represents a notable research gap, as AI-driven automation holds considerable potential to enhance the precision and efficiency of soil monitoring and management. One possible reason is that key soil attributes, such as microbial activity or nutrient cycling, are not easily detectable through visual cues. Moreover, acquiring labeled image data in diverse field conditions remains a time-consuming and costly challenge.

As automation and AI increasingly influence decision-making and field operations, ensuring transparency and trust in these systems becomes crucial. In the area of XAI, researchers commonly investigate feature importance (FI) and distinguish between global (GFI) and local (LFI) perspectives. GFI refers to the importance of individual covariates across the entire dataset, while LFI captures the contribution of covariates to a single prediction.

Traditional statistical models, such as linear or logistic regression, typically require no additional FI analysis, as their fitted coefficients directly indicate covariate importance. In contrast, FI is generally more difficult to assess for ML models. Tree-based models, such as RF, are a notable exception in terms of GFI, offering built-in methods like *Mean Decrease in Impurity* and *Permutation Importance*. However, for LFI, particularly in the context of black-box models such as NNs and DL, covariate importance must typically be analyzed using dedicated techniques. SHAP values and Local Interpretable Model-agnostic Explanations (LIME) [107] are the most widely used methods for LFI analysis [108], with SHAP also applicable to GFI. Despite their utility, few studies have applied FI techniques such as SHAP or LIME. In general, XAI methods like these should be adopted more frequently to enhance the plausibility, interpretability, and explainability of ML models. Such improvements are essential for enhancing the trustworthiness of ML systems [21,109]. Their limited use may reflect a general lack of awareness or technical familiarity with XAI tools among agricultural researchers, as well as the additional computational burden and complexity they introduce to already resource-intensive workflows.

Last but not least, there is still little research on the use of large language models (LLMs) for soil health in agriculture, despite their significant potential in DSS. Common AI methods in agriculture, such as DL, often require large, annotated datasets, which are costly and time-consuming to obtain. In contrast, LLMs, such as OpenAI's ChatGPT models, have demonstrated high accuracy in agricultural text classification tasks [110], even without fine-tuning. This capability could be leveraged to inform DSS by extracting insights from unstructured data sources, such as news about natural hazards, pest outbreaks, or market trends. Moreover, LLMs offer unique potential in translating complex analytical or scientific information into formats that are easily understandable across multiple languages and varying levels of education [111].

### 3.1.2. Regional bias and global applicability

AI research in soil health predominantly focuses on sub-national applications, which may limit its relevance for policymakers who require broader regional or global perspectives. Some approaches aim to define soil quality indices using PCA, with the objective of identifying a minimum dataset from a larger set of soil properties. These selected properties are then transformed to a uniform scale and combined into a final index [112]. Such indices enable the assessment and comparison of soil health across nearby regions with similar land use types [113]. However, developing a globally applicable soil quality index continues to be a major challenge due to the wide variation in climatic and soil conditions worldwide [114].

Future research should explore methods to develop robust, scalable indices that can guide AI-driven soil health assessments on a larger scale. Key AI research areas that can support this scaling include: (i) transfer learning [115,116], to adapt models trained on data-rich regions to data-scarce areas; (ii) federated ML [117–119], to enable collaborative model development across countries without sharing sensitive local data; (iii) self-supervised learning [120], to leverage vast amounts of unlabeled soil and environmental data for feature extraction; and (iv) multi-modal learning [121], to integrate heterogeneous data sources such as remote sensing, field measurements, and environmental simulations for more comprehensive models [122,123]. These approaches can help bridge the gap between local insights and global policy needs in soil health management.

The global applicability of soil health monitoring also depends on the availability of accurate on-field measurement methods. While laboratory analyses remain the gold standard for precision, the reliability of field-based techniques for key indicators such as SOC and microbial activity has been demonstrated [124]. Promoting easy-to-use, on-site measurement tools that can be applied by farmers is essential to enable regular and widespread soil data collection.

Another approach for scaling from local to regional or even global levels is the establishment of *lighthouse project networks*. These lighthouse projects aim to inform policymakers about soil health targets under comparable conditions across different regions of a study area. Additionally, locally relevant solutions for soil health challenges are regularly evaluated and reported. The EU mission *A Soil Deal for Europe* has the target of establishing a network of 100 living labs and lighthouses to co-create knowledge, test solutions and demonstrate their value in real-life conditions [11,125]. By implementing standardized measurement protocols, this network has the potential to enable large-scale and harmonized assessment of soil data across Europe, potentially inspiring similar projects on a global scale.

Furthermore, our analysis of data origins shows clear regional biases. As extensively discussed in Section 2.3.3, entire climate zones, such as tropical, arid, mild continental, and polar tundra regions, are significantly underrepresented in AI research for soil health in agriculture. These zones span vast areas of Central and South America, Africa, Central and South Asia, and Australia. This lack of data from key climatic regions raises concerns about whether AI models trained on existing datasets can generalize effectively across diverse environmental conditions. To ensure the global applicability of AI-driven solutions, future research should prioritize the creation of high-quality, standardized datasets from these underrepresented areas. In this context, methods for bias detection and mitigation in datasets and AI models [126] will be surely essential to address existing imbalances and support the development of robust and generalizable AI-based soil health assessments.

In summary, while AI has made significant contributions to soil health research, several areas remain underexplored. Future research should diversify AI methodologies, expand data availability in neglected regions and develop more generalizable soil quality indices. By addressing these gaps, AI can play a more effective role in promoting sustainable soil management worldwide.

### 3.2. Research opportunities in modeling soil response variables

In the literature review, we analyzed the data types used to model various soil response variables. Fig. 11 highlights the most relevant response variable associated with each soil threat, as identified in the review. For each response variable, the figure also presents potential input data types (e.g., spectral or topographical data) that have been underrepresented in existing studies but remain promising avenues for further investigation. It is important to note that these response–predictor connections are not intended to suggest modeling a response variable using only the indicated predictor types, but rather to encourage their inclusion as additional inputs in future modeling efforts.
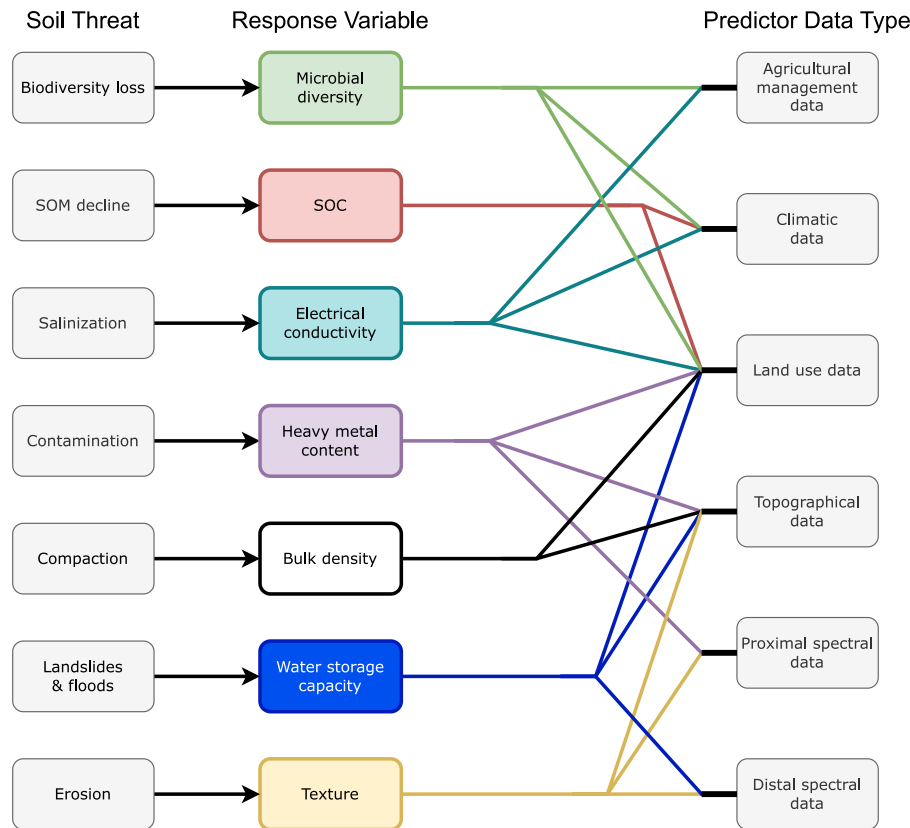
**Fig. 11.** Future Research: Mapping potential input data types to the most relevant soil response variables for each soil threat, based on insights from the literature review. These mappings highlight feasible but underrepresented combinations of response variables and predictor data types in the existing literature.

An example for the soil threat *compaction* is predicting bulk density based on topographical data, which is a plausible but underutilized approach. Terrain attributes such as slope and curvature affect soil compaction by influencing machinery movement, water accumulation, and erosion. Flatter areas often experience higher traffic and water retention, leading to increased bulk density. AI models using topographic inputs, such as those derived from digital elevation models, could help estimate bulk density where direct measurements are unavailable, offering a scalable solution for compaction risk assessment [127].

Predicting electrical conductivity using agricultural management data is another promising but underutilized approach. Management practices such as fertilizer application and crop rotation directly influence salt accumulation in the soil. In addition, excessive irrigation without proper drainage can lead to elevated electrical conductivity. Incorporating such management data into the modeling process could improve the detection and monitoring of salinity risks, particularly in intensively farmed or irrigated areas [128].

Beyond topographical and management data, other underused input types also hold potential for soil modeling. Climatic data (e.g., mean annual temperature and precipitation) influence key processes like organic matter turnover and microbial activity [129]. Land use data, such as EUNIS habitat classifications,[12] reflect human impact and ecological context. Proximal spectral data from hand-held sensors provide detailed surface measurements, while distal spectral data from sensors [130] or satellites [131] enable large-scale monitoring of variables like SOC and soil moisture.

Similar to Figs. 11, 12 provides a more detailed view of potential soil parameters that could serve as predictor variables. The figure displays the same set of soil threats and response variables on the left, now paired with relevant soil parameters on the right. As before, the mappings highlight combinations that are technically feasible but have been rarely explored in the existing literature.

For example, predicting SOC using soil mineral nitrogen data is both feasible and underexplored. Soil mineral nitrogen content, reflecting microbial activity and organic matter decomposition, can improve SOC predictions by capturing key biological processes, especially in contexts where direct measurements are limited or costly [132]. However, especially for SOC prediction it is important to incorporate additional soil parameters, such as soil texture, in order to account for SOC accumulation potential [133].

Another example involves predicting heavy metal content for the soil threat contamination using earthworm data. Earthworms are well-established bioindicators of soil pollution, as their abundance, biomass, and tissue composition often correlate with concentrations of metals like lead, cadmium, and arsenic. Changes in earthworm populations can signal contamination earlier than traditional chemical analyses, making them valuable for early detection. Integrating such biological indicators into AI models could enhance the assessment of soil contamination, particularly in regions where detailed chemical data is lacking [134].

### 3.3. Workflow suggestions for soil parameter modeling

The primary use cases of soil parameter modeling in soil research encompass various types of prediction tasks. In general, this involves estimating soil response parameters based on soil predictor variables. These estimates can predict the current state using past observations or forecast future conditions. Often, these predictions are used to generate maps (DSM), where estimated soil parameters are assigned to geographic areas lacking measurements of the response variable, or where measurements are unavailable at the desired spatial resolution.

---

[12] EUNIS habitat classification: https://eunis.eea.europa.eu/habitats-code-browser.jsp, accessed April 3rd, 2025.
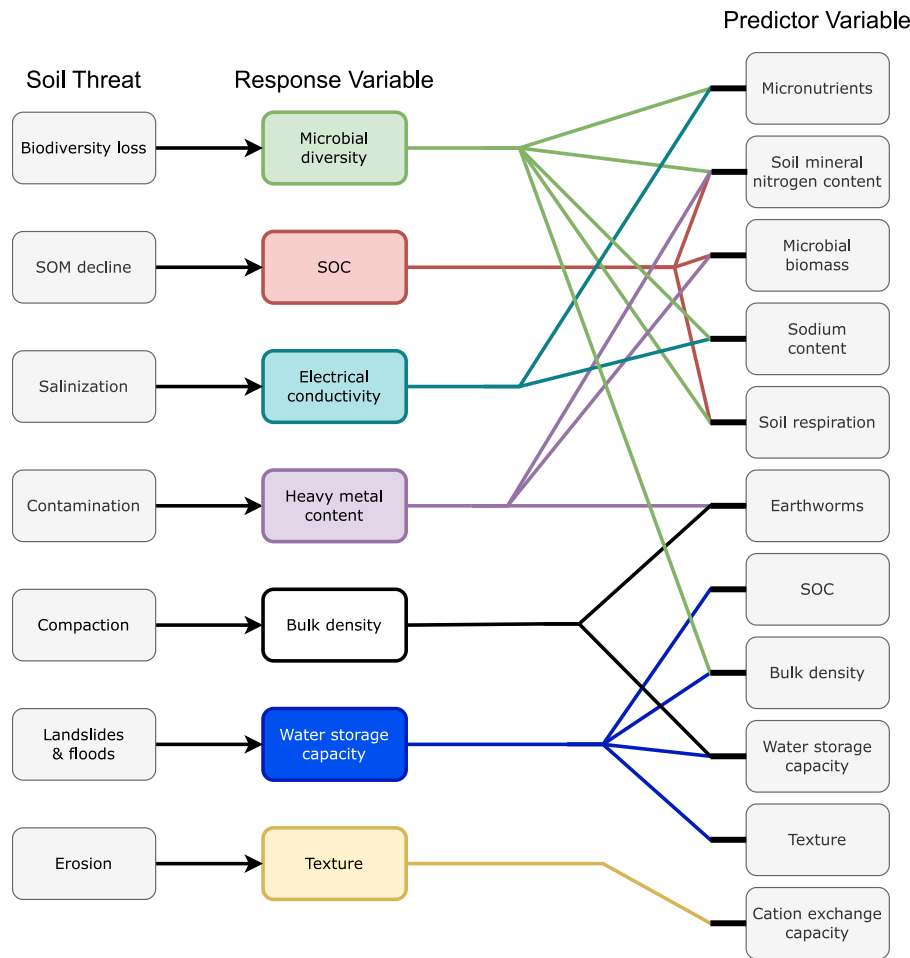
**Fig. 12.** Future Research: Mapping potential input soil parameters (predictor variables) to the most relevant soil response variables for each soil threat, based on insights from the literature review. These mappings highlight feasible but underrepresented combinations of response variables and predictors in the existing literature.

The following sections present recommendations for soil parameter prediction. In particular, Section 3.3.1 addresses feature extraction and data preparation in agricultural contexts, while Section 3.3.2 focuses on selecting appropriate modeling techniques based on dataset characteristics.

### 3.3.1. Feature extraction and data preparation

When data collection is finished, the typical workflow for soil parameter prediction proceeds with feature extraction and preprocessing. If the data includes RGB imagery or spectral data, it is often necessary to extract features from these input types, as shown in Fig. 13. For RGB imagery, extracted features may include morphological plant parameters or stress indicators, such as leaf rolling [135]. In contrast, common spectral features include crop residue coverage estimates [136,137] or vegetation indices like the normalized difference vegetation index (NDVI) [62,63], eventually used to estimate yield or crop quality [138, 139]

However, feature extraction is not limited to RGB and spectral data. From low-level inputs such as elevation or land use data, more specific features such as terrain flatness or distances to water bodies can be derived [140]. Additionally, pedotransfer functions and process-based models can be used to estimate unknown soil parameters based on available data, thereby generating new features and augmenting the dataset [141].

In most cases, the objective is to obtain a set of tabular features suitable for the subsequent prediction task. Once all input data has been transformed into tabular format, optional steps in the data pipeline include feature reduction, data cleansing and outlier filtering, as described in Section 2.2.2. Based on the resulting dataset, an appropriate modeling technique can then be selected to match the structure and characteristics of the data.

### 3.3.2. Choosing appropriate modeling techniques

One of the most challenging aspects of soil parameter modeling is selecting an appropriate modeling technique. The choice largely depends on the characteristics of the available dataset. Fig. 14 supports soil researchers in choosing suitable approaches based on key dataset properties, such as sample size, feature dimensionality, and data linearity. For each scenario, the flowchart suggests applicable models, with a legend indicating whether models are robust to noise, support missing inputs, or provide uncertainty estimates. The primary focus remains on regression models, as they dominate soil parameter prediction tasks.

Ensemble methods such as RF and XGBoost have gained significant popularity due to their ability to efficiently extract patterns from large and complex datasets. This success is largely driven by their combination of multiple weak learners and the widespread availability of well-optimized software implementations [142]. Among these, RF is considered a versatile and reliable default. It is inherently robust to noise, can handle missing data through surrogate splits, and provides intuitive metrics for assessing FI. Extensions like Quantile Regression Forests (QRF) add the ability to quantify prediction uncertainty, making them especially useful when uncertainty estimates are required for better informed decision making. Other modeling approaches [143] also offer algorithmic means to provide confidence estimates for their
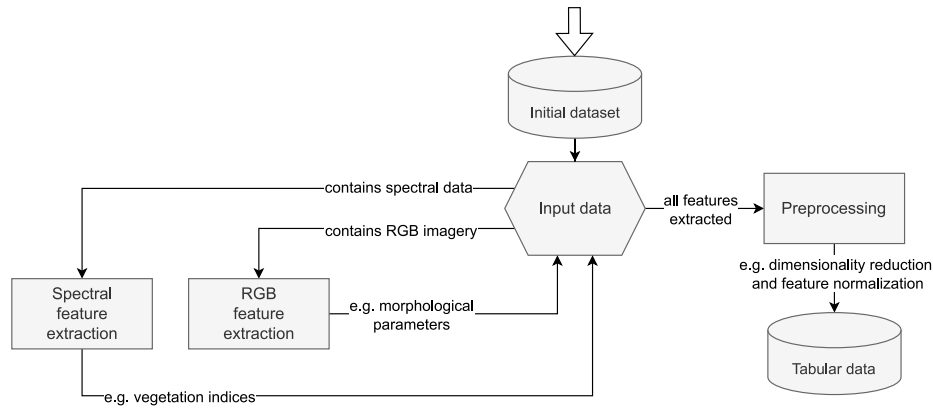
**Fig. 13.** Feature extraction workflow based on an example with input data in soil research.
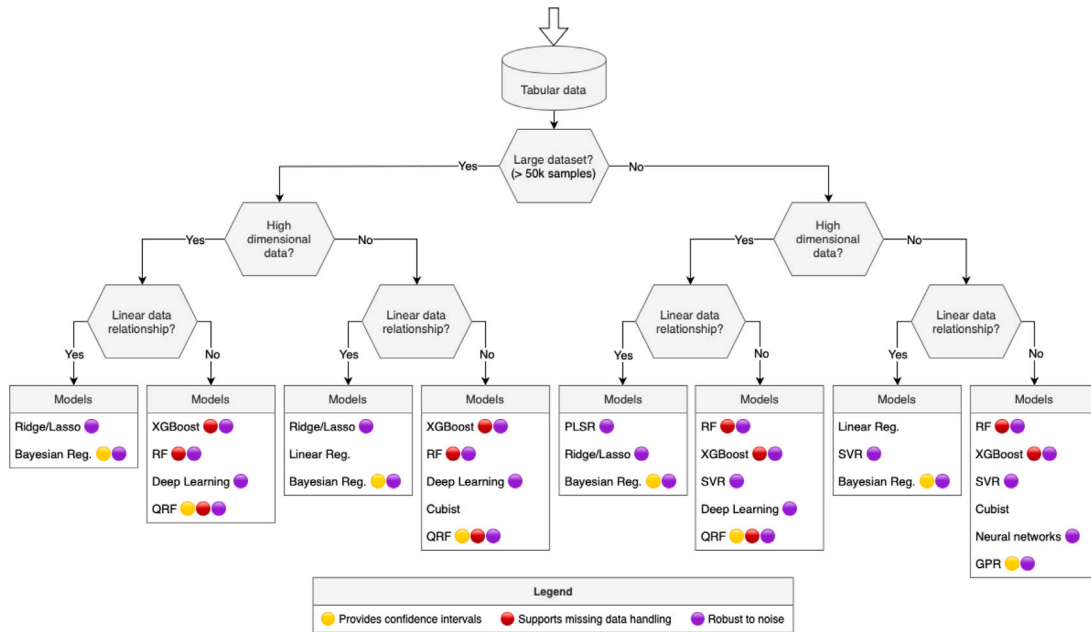


**Fig. 14.** Recommended regression models based on dataset size, dimensionality and linearity of data.

predictions. Combined with model-agnostic uncertainty quantification techniques such as conformal prediction [144], these methods present valuable opportunities for estimating confidence in practical soil health scenarios, which are often subject to high epistemic uncertainty or inherent measurement noise.

Cubist, a hybrid rule-based model, also offers a balance between interpretability and predictive power. By combining decision trees with multivariate linear regressions at the leaf level, Cubist captures nonlinear patterns while maintaining transparency. Like RF, it handles missing values via surrogate features. However, it may be more sensitive to outliers due to its reliance on linear regression components.

When working with high-dimensional data, overfitting becomes a key concern. Regularization techniques such as Ridge and Lasso regression help mitigate this by penalizing complex models. While Ridge is effective with correlated predictors, Lasso can shrink irrelevant features to zero, effectively performing feature selection. PLSR offers an alternative by projecting data into a lower-dimensional latent space. Though suitable for small sample sizes, PLSR becomes computationally intensive for large datasets, where tree-based or DL models may perform better.

Model interpretability is often critical, particularly in agricultural decision-making contexts. While DL models, such as convolutional NNs (CNNs) or Transformers, excel at handling unstructured or high-dimensional data (e.g., imagery), they are often perceived as black boxes. Techniques like SHAP and LIME can help explain predictions even in complex architectures. For instance, SHAP has been used in SOC modeling to reveal the importance of variables like crop residue coverage, temperature, and clay content [136].

Data linearity can be assessed through visual inspection (e.g., scatter or residual plots), or by fitting a simple linear model and evaluating its residuals. When data is nonlinear, models like RF, XGBoost, or Support Vector Regression (SVR) with nonlinear kernels (e.g., radial basis function or polynomial kernels) become more appropriate. SVR performs particularly well on small to medium sized datasets, provided its regularization parameters are carefully tuned [145].

In noisy or uncertain environments, Bayesian Regression, Gaussian Process Regression (GPR), and QRF offer the advantage of uncertainty quantification. Bayesian models incorporate prior distributions over parameters, which act as a form of regularization and provide confidence bounds around predictions [146]. Similarly, GPR directly models the distribution over outputs for any given input, and QRF yields prediction intervals by analyzing the spread of tree predictions.

Multi-modal data, such as combinations of imagery, sensor data, and farm management records, can be addressed using a serial modeling approach: first applying modality-specific feature extraction
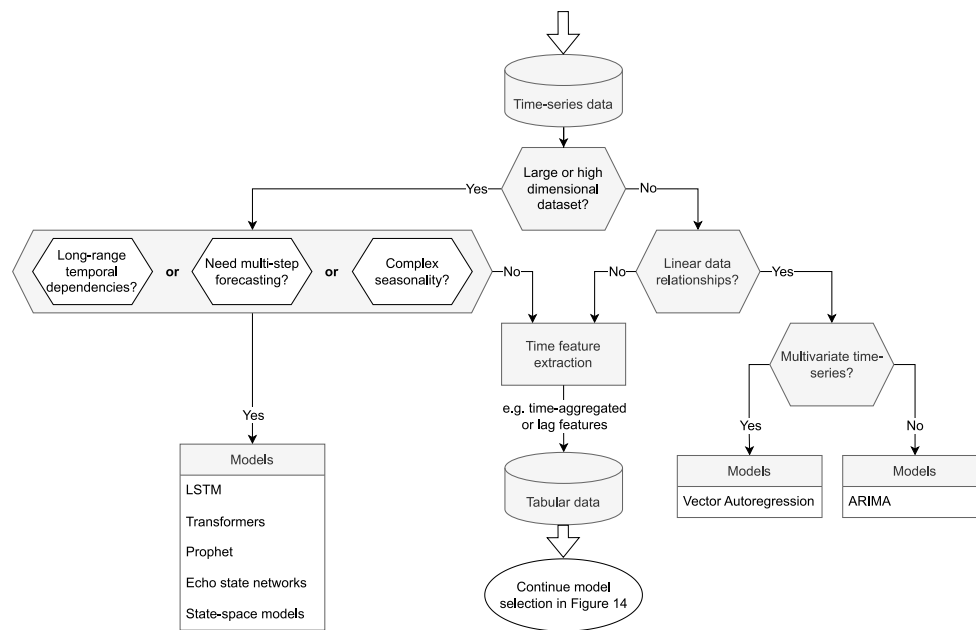
**Fig. 15.** Recommended models for time-series data and time-series prediction.

(e.g., CNNs for RGB images, spectral indices for hyperspectral data) and then aggregating extracted features into a unified tabular format. This allows the final regression step to use any classical tabular model, such as RF or XGBoost, depending on the dataset properties (Fig. 14).

Time-series data introduces additional complexity due to temporal dependencies. For linear trends, traditional models like ARIMA or Vector Autoregression are effective. For nonlinear temporal patterns, classical ML models (e.g., RF, SVR, gradient boosting) can be used if time is encoded as features (e.g., lags or timestamps). For larger or more complex datasets, advanced models such as LSTMs, Transformers, Facebook Prophet, ESNs or SSMs are better suited to capture seasonality and long-term trends. Fig. 15 outlines recommended models for time-series data.

As discussed in Section 3.1.1, a promising direction in AI-based modeling for soil health assessment is the application of hybrid modeling techniques. One such approach is physics-informed neural networks (PINNs) [147,148], which show strong potential for soil health monitoring by integrating physical models of soil processes, such as moisture retention, organic matter decomposition, and nutrient cycling, into AI predictions. Unlike purely data-driven models, PINNs incorporate these known relationships as constraints during training, enabling more accurate and physically consistent predictions, even with limited or noisy data. This is particularly valuable for estimating hard-to-measure soil health indicators like microbial activity or carbon fluxes, where direct observations are sparse. By grounding AI predictions in well-understood soil dynamics, PINNs can enhance the reliability and interpretability of soil health assessments across diverse environmental conditions.

Finally, when spatial dependencies are not captured directly in input features, geostatistical methods such as residual kriging can be applied in a post-processing step. This hybrid approach has shown improved accuracy in DSM applications, for instance, in the spatial prediction of SOC [140].

In summary, the model selection process in soil parameter prediction is highly context-dependent. Factors such as data size, structure, noise level, interpretability needs, and modality should guide the choice.

## 4. Conclusions and outlook

The integration of AI in soil health research has shown transformative potential in advancing our understanding, monitoring, and management of soil ecosystems, particularly in agriculture. From predictive modeling and DSM to DSS and XAI, this work has highlighted the wide range of AI-driven solutions currently addressing soil health challenges. Table 6 summarizes the key insights of this work as well as identified research gaps and future research directions.

AI is especially effective in analyzing large-scale, heterogeneous datasets, enabling precision agriculture, and detecting early indicators of soil degradation. Despite these advancements, several challenges and opportunities remain. One major limitation is the lack of standardized, globally applicable soil health benchmarks and ground truth datasets. Most AI models today rely on region-specific data, limiting their ability to generalize across diverse pedo-climatic conditions. To address this, future research should focus on developing foundational, universal AI models capable of adapting to varying environments while maintaining high accuracy. Techniques such as MTL, transfer learning, and federated learning offer promising pathways for training models across distributed datasets while preserving data privacy.

Moreover, several modeling approaches remain underutilized in current literature. These include hybrid or physics-informed models, which combine data-driven and process-based approaches, as well as dedicated time-series models, and RL for dynamic decision-making in soil management.

Another promising direction lies in the integration of multi-modal AI systems, which combine remote sensing, in-situ sensor data, and genomic analysis of soil microbiomes. While ML has been widely applied to physical and chemical soil properties, biological aspects, such as microbial community structures, remain underexplored. DL applications in microbiome analysis could offer new insights by linking microbial diversity to soil functions and ecosystem services.

Ethical considerations must also be addressed as AI becomes more embedded in agricultural decision-making. Issues such as fairness, inclusivity, and data ownership are critical. Future AI tools should be co-developed with local stakeholders to ensure equitable access and avoid reinforcing regional inequalities.

In addition, explainability remains essential for transparency, traceability, and trust in AI-driven soil health assessments. Many current models operate as black or gray boxes, limiting their interpretability for farmers, policymakers, and soil scientists. Future work should prioritize interpretable ML approaches that provide actionable insights without sacrificing accuracy. Models that integrate domain-specific soil science

**Table 6**

Summary of key insights, research gaps, and future directions on AI for soil health in agriculture.

| Key insights | Research gaps | Future directions |
|---|---|---|
| AI enables soil health assessment and sustainable management | Soil data collection remains inconsistent and biased towards a few regions | Establish standardized monitoring protocols; expand data from underrepresented zones using transfer learning and federated ML |
| Novel predictor–response combinations can improve modeling (cf. Figs. 11, 12) | Many feasible combinations (e.g., topography → bulk density, earthworms → heavy metals) remain unexplored | Systematically integrate diverse input data types and biological indicators into soil models |
| Advanced AI approaches (RL, GNNs, PINNs, MTL, time-series models) show high potential | Underutilized due to technical complexity and data demands | Apply RL with guardrails for DSS, GNNs for spatial data, MTL for related tasks, and hybrid physics–ML models for SOC and other properties with known underlying physical dynamics |
| Interpretability and trust are essential for adoption | Few studies apply SHAP, LIME, or domain-informed XAI | Prioritize interpretable ML and XAI to ensure transparency and stakeholder trust, particularly those with low or even null expertise in AI |
| ML has mainly targeted physical and chemical soil properties | Biological aspects (e.g., microbial communities) are underexplored | Use DL and multi-modal AI (remote sensing, in-situ sensors, microbiome data) to link soil biology with functions and services |
| Emerging technologies offer new opportunities | Limited research on robotics for reduced soil compaction and LLMs for DSS | Explore lightweight robotics for sustainable field operations and LLMs for unstructured data, multilingual knowledge transfer, causal inference, and decision support. |

knowledge with AI techniques could further bridge the gap between expertise and automation.

In summary, while AI has already begun reshaping soil health research and agricultural practices, significant white spaces remain for innovation beyond frontiers. Addressing these open challenges will require interdisciplinary collaboration between AI researchers, soil scientists, agronomists, and policymakers. By advancing towards transparent, robust, and globally scalable AI solutions, we can harness AI's full potential to safeguard soil health and ensure the sustainability of agricultural ecosystems in the face of climate change and growing global food demands.

## Abbreviations

AI ... Artificial Intelligence
CNN ... Convolutional Neural Networks
DEQ ... Data Extraction Question
DL ... Deep Learning
DQN ... Deep Q-Network
DSM ... Digital Soil Mapping
DSS ... Decision Support Systems
EC ... Exclusion Criterion
ENET ... Elastic Net
ESN ... Echo State Network
FI ... Feature Importance
GBM ... Gradient Boosting Machine
GFI ... Global Feature Importance
GL(M)M ... Generalized Linear (Mixed) Model
GNN ... Graph Neural Network
GPR ... Gaussian Process Regression
GRRF ... Guided Regularized Random Forest
IC ... Inclusion Criterion
ISRIC ... International Soil Reference and Information Centre
LFI ... Local Feature Importance
LIME ... Local Interpretable Model-agnostic Explanations
LSTM ... Long Short-Term Memory
LUCAS ... Land Use and Coverage Area Frame Survey
ML ... Machine Learning
MLR ... Multiple Linear Regression
MTL ... Multi-task Learning
NDVI ... Normalized Difference Vegetation Index
NN ... Neural Network
PCA ... Principal Component Analysis
PCoA ... Principal Coordinate Analysis

PLSR ... Partial Least Squares Regression
QRF ... Quantile Regression Forest
RF ... Random Forest
RFE ... Recursive Feature Elimination
RL ... Reinforcement Learning
RNN ... Recurrent Neural Network
RQ ... Research Question
SHAP ... Shapley Additive Explanations
SLR ... Structured Literature Review
SOC ... Soil Organic Carbon
SOM ... Soil Organic Matter
SSM ... State-space Model
SVR ... Support Vector Regression
SVM ... Support Vector Machine
SRTM ... Shuttle Radar Topography Mission
WoS ... Web of Science

## CRediT authorship contribution statement

**Stefan Schweng:** Writing – original draft, Data curation, Software, Methodology, Investigation, Writing – review & editing. **Luca Bernardini:** Writing – original draft, Data curation, Methodology, Investigation, Validation, Writing – review & editing. **Katharina Keiblinger:** Validation, Supervision, Writing – review & editing. **Hans-Peter Kaul:** Supervision, Data curation, Writing – review & editing. **Iztok Fister Jr.:** Supervision, Data curation, Writing – review & editing. **Niko Lukač:** Supervision, Data curation, Writing – review & editing. **Javier Del Ser:** Supervision, Writing – review & editing. **Andreas Holzinger:** Conceptualization, Data curation, Writing – original draft, Supervision, Funding acquisition, Writing – review & editing.

## Declaration of Generative AI and AI-assisted technologies in the writing process

During the preparation of this work the authors used language tools such as Writefull, DeepL and ChatGPT in order to improve the article's readability. After using these tools, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

## Funding

and of the Government of Lower Austria, Project GFF NÖ FTI-22-I-004 "Infrastructure for the realistic testing of AI-supported robot systems in demanding environments without direct energy connection, for multiple use cases (e.g., monitoring/maintenance of forest roads) - human–robot teaming". Javier Del Ser acknowledges funding support from the Basque Government through the consolidated research group MATHMODE (IT1456-22).

## 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.

## Appendix A. Supplementary data

Supplementary material related to this article can be found online at https://doi.org/10.1016/j.cosrev.2025.100832. Supplementary data contains an Excel file listing all 115 articles used for the SLR and the extracted DEQ data.

## Data availability

Attached as Excel File.

## References

[1] P. Pingali, J. Boiteau, A. Choudhry, A. Hall, Making meat and milk from plants: A review of plant-based food for human and planetary health, World Dev. 170 (2023) http://dx.doi.org/10.1016/j.worlddev.2023.106316.

[2] R.D. Sands, B. Meade, J.L. Seale, S. Robinson, R. Seeger, Scenarios of global food consumption: Implications for agriculture, Econ. Res. Rep. 323 (2023) URL https://www.ers.usda.gov/publications/pub-details/?pubid=107473. USDA Economic Research Service.

[3] C. Rosinger, K.M. Keiblinger, L.G. Bernardini, D. Horn, G. Heller, H. Eigner, H.P. Kaul, S. Huber, O. Sae-Tun, G. Bodner, Significant advances in plant-available and replenishable macro- and micronutrients with soil health-oriented conservation farming: Novel insights from a multi-site on-farm evaluation, Geoderma 457 (2025) http://dx.doi.org/10.1016/j.geoderma.2025.117275.

[4] S. Huber, C. Rosinger, G. Bodner, L.G. Bernardini, M. Bieber, A. Mentler, O. Sae-Tun, B. Scharf, K.M. Keiblinger, Highway to health: Microbial pathways of soil organic carbon accrual in conservation farming systems, Geoderma 452 (2024) http://dx.doi.org/10.1016/j.geoderma.2024.117115.

[5] J. Rockström, J. Williams, G. Daily, A. Noble, N. Matthews, L. Gordon, H. Wetterstrand, F. DeClerck, M. Shah, P. Steduto, C. de Fraiture, N. Hatibu, O. Unver, J. Bird, L. Sibanda, J. Smith, Sustainable intensification of agriculture for human prosperity and global sustainability, Ambio 46 (2017) 4–17, http://dx.doi.org/10.1007/s13280-016-0793-6.

[6] P.M. Kopittke, N.W. Menzies, P. Wang, B.A. McKenna, E. Lombi, Soil and the intensification of agriculture for global food security, Environ. Int. 132 (2019) http://dx.doi.org/10.1016/j.envint.2019.105078.

[7] H.I. Bedolla-Rivera, M.d.l.L.X. Negrete-Rodríguez, F.P. Gámez-Vázquez, D. Álvarez-Bernal, E. Conde-Barajas, Analyzing the impact of intensive agriculture on soil quality: A systematic review and global meta-analysis of quality indexes, Agronomy 13 (8) (2023) http://dx.doi.org/10.3390/agronomy13082166, URL https://www.mdpi.com/2073-4395/13/8/2166.

[8] A.A.R. Sousa, J. Muñoz-Rojas, C. Brígido, S.A. Prats, Impacts of agricultural intensification on soil erosion and sustainability of olive groves in Alentejo (Portugal), Landsc. Ecol. 38 (2023) 3479–3498, http://dx.doi.org/10.1007/s10980-023-01682-2.

[9] European University Institute, T. Jevnaker, F. Agostini, E. Beckstedde, S. Nicolai, R. Belmans, I. Conti, A. Ferrari, L. Hancher, L. Heinrich, T. Iliopoulos, L. Iozzelli, J. Kneebone, E. Marro, L. Meeus, E. Menegatti, M. Münchmeyer, A. Nouicer, M. Olczak, A. Piebalgs, A. Pototschnig, V. Reif, N. Rossetto, M. Salvetti, M. del Carmen Sandoval, T. Schittekatte, P. Schlosser, D. Stampatori, The EU Green Deal – 2024 Edition, European University Institute, 2025, http://dx.doi.org/10.2870/1445509.

[10] European Commission and Directorate-General for Environment, EU Biodiversity Strategy for 2030 – Bringing Nature Back into Our Lives, Publications Office of the European Union, 2021, http://dx.doi.org/10.2779/677548.

[11] E. Commission, D.-G. for Research, Innovation, EU Mission, Soil Deal for Europe, Publications Office of the European Union, 2022, http://dx.doi.org/10.2777/706627.

[12] H. Jenny, Factors of Soil Formation: A System of Quantitative Pedology, McGraw-Hill Book Company, New York, 1941.

[13] D. Yaalon, Conceptual models in pedogenesis: Can soil-forming functions be solved? Geoderma 14 (3) (1975) 189–205, http://dx.doi.org/10.1016/0016-7061(75)90001-4.

[14] R. Gittins, Trend-surface analysis of ecological data, J. Ecol. 56 (3) (1968) 845–869, http://dx.doi.org/10.2307/2258110.

[15] T.M. Burgess, R. Webster, Optimal interpolation and isarithmic mapping of soil properties: I. The semi-variogram and punctual kriging, Eur. J. Soil Sci. 70 (1) (2019) 11–19, http://dx.doi.org/10.1111/ejss.12784.

[16] T.M. Burgess, R. Webster, Optimal interpolation and isarithmic mapping of soil properties: Ii. block kriging, J. Soil Sci. 31 (2) (1980) 333–341, http://dx.doi.org/10.1111/j.1365-2389.1980.tb02084.x.

[17] R. Webster, T.M. Burgess, Optimal interpolation and isarithmic mapping of soil properties: Iii. changing drift and universal kriging, J. Soil Sci. 31 (3) (1980) 505–524, http://dx.doi.org/10.1111/j.1365-2389.1980.tb02100.x.

[18] A.B. McBratney, R. Webster, Optimal interpolation and isarithmic mapping of soil properties, J. Soil Sci. 34 (1) (1983) 137–162, http://dx.doi.org/10.1111/j.1365-2389.1983.tb00820.x.

[19] M. Vauclin, S.R. Vieira, G. Vachaud, D.R. Nielsen, The use of cokriging with limited field soil observations, Soil Sci. Am. J. 47 (2) (1983) 175–184, http://dx.doi.org/10.2136/sssaj1983.03615995004700020001x.

[20] A.B. McBratney, M.L.M. Santos, B. Minasny, On digital soil mapping, Geoderma 117 (2003) 3–52, http://dx.doi.org/10.1016/S0016-7061(03)00223-4.

[21] A.M. Wadoux, B. Minasny, A.B. McBratney, Machine learning for digital soil mapping: Applications, challenges and suggested solutions, Earth-Sci. Rev. 210 (2020) http://dx.doi.org/10.1016/j.earscirev.2020.103359.

[22] L.G. Bernardini, C. Rosinger, G. Bodner, K.M. Keiblinger, E. Izquierdo-Verdiguier, H. Spiegel, C.O. Retzlaff, A. Holzinger, Learning vs. understanding: When does artificial intelligence outperform process-based models in soil organic carbon prediction? New Biotechnol. 81 (7) (2024) 20–31, http://dx.doi.org/10.1016/j.nbt.2024.03.001.

[23] R. Wang, W. Chen, H. Chen, X. Qin, Finer soil properties mapping framework for broad-scale area: A case study of Hubei Province, China, Geoderma 449 (2024) http://dx.doi.org/10.1016/j.geoderma.2024.117023.

[24] S. Fenz, T. Neubauer, J.K. Friedel, M.L. Wohlmuth, AI- and data-driven crop rotation planning, Comput. Electron. Agric. 212 (2023) http://dx.doi.org/10.1016/j.compag.2023.108160.

[25] E.J. Díaz-Vallejo, M. Seeley, A.P. Smith, E. Marín-Spiotta, A meta-analysis of tropical land-use change effects on the soil microbiome: Emerging patterns and knowledge gaps, Biotropica 53 (2021) 738–752, http://dx.doi.org/10.1111/btp.12931.

[26] E. Stell, D. Warner, J. Jian, B. Bond-Lamberty, R. Vargas, Spatial biases of information influence global estimates of soil respiration: How can we improve global predictions? Global Change Biol. 27 (2021) 3923–3938, http://dx.doi.org/10.1111/gcb.15666.

[27] J. Jian, R. Vargas, K. Anderson-Teixeira, E. Stell, V. Herrmann, M. Horn, N. Kholod, J. Manzon, R. Marchesi, D. Paredes, B. Bond-Lamberty, A restructured and updated global soil respiration database (SRDB-V5), Earth Syst. Sci. Data 13 (2) (2021) 255–267, http://dx.doi.org/10.5194/essd-13-255-2021, URL https://essd.copernicus.org/articles/13/255/2021/.

[28] F.A. Diaz-Gonzalez, J. Vuelvas, C.A. Correa, V.E. Vallejo, D. Patino, Machine learning and remote sensing techniques applied to estimate soil indicators – review, Ecol. Indic. 135 (2022) http://dx.doi.org/10.1016/j.ecolind.2021.108517.

[29] S. Jain, D. Sethia, K.C. Tiwari, A critical systematic review on spectral-based soil nutrient prediction using machine learning, Environ. Monit. Assess. 196 (8) (2024) 699, http://dx.doi.org/10.1007/s10661-024-12817-6.

[30] S. Barathkumar, K. Sellamuthu, K. Sathyabama, P. Malathi, R. Kumaraperumal, P. Devagi, Advancements in soil quality assessment: A comprehensive review of machine learning and AI-driven approaches for nutrient deficiency analysis, Commun. Soil Sci. Plant Anal. 56 (2) (2025) 257–282, http://dx.doi.org/10.1080/00103624.2024.2406484.

[31] B. Shi, J. Meng, T. Wang, Q. Li, Q. Zhang, G. Su, The main strategies for soil pollution apportionment: A review of the numerical methods, J. Environ. Sci. 136 (2024) 95–109, http://dx.doi.org/10.1016/j.jes.2022.09.027.

[32] S. Sow, S. Ranjan, M.F. Seleiman, H.M. Alkharabsheh, M. Kumar, N. Kumar, S.R. Padhan, D.K. Roy, D. Nath, H. Gitari, Artificial intelligence for maximizing agricultural input use efficiency: Exploring nutrient, water and weed management strategies, Phyton - Int. J. Exp. Bot. 93 (7) (2024) 1–30, http://dx.doi.org/10.32604/phyton.2024.052241.

[33] S. Raza, A. Irshad, A. Margenot, N. Li, S. Ullah, K. Mehmood, M.A. Khan, N. Siddique, J. Zhou, Inorganic carbon is overlooked in global soil carbon research: A bibliometric analysis, Geoderma 443 (2024) 116831, http://dx.doi.org/10.1016/j.geoderma.2024.116831.

[34] S. Sood, H. Singh, Computer vision and machine learning based approaches for food security: A review, Multimedia Tools Appl. 80 (18) (2021) 27973–27999, http://dx.doi.org/10.1007/s11042-021-11036-2.

[35] A.M. Husaini, M. Sohail, Robotics-assisted, organic agricultural-biotechnology based environment-friendly healthy food option: Beyond the binary of GM versus organic crops, J. Biotech. 361 (2023) 41–48, http://dx.doi.org/10.1016/j.jbiotec.2022.11.018.

[36] S.L. Ullo, G.R. Sinha, Advances in IoT and smart sensors for remote sensing and agriculture applications, Remote. Sens. 13 (13) (2021) 2585, http://dx.doi.org/10.3390/rs13132585.

[37] M. Khanna, S.S. Atallah, S. Kar, B. Sharma, L. Wu, C. Yu, G. Chowdhary, C. Soman, K. Guan, Digital transformation for a sustainable agriculture in the United States: Opportunities and challenges, Agricult. Econ. 53 (6) (2022) 924–937, http://dx.doi.org/10.1111/agec.12733.

[38] R.K. Goel, C.S. Yadav, S. Vishnoi, R. Rastogi, Smart agriculture–urgent need of the day in developing countries, Sustain. Comput.: Inform. Syst. 30 (2021) 100512, http://dx.doi.org/10.1016/j.suscom.2021.100512.

[39] N. Alex, C. Sobin, J. Ali, A comprehensive study on smart agriculture applications in India, Wirel. Pers. Commun. 129 (4) (2023) 2345–2385, http://dx.doi.org/10.1007/s11277-023-10234-5.

[40] R.A. El Behairy, H.M. El Arwash, A.A. El Baroudy, M.M. Ibrahim, E.S. Mohamed, D.E. Kucher, M.S. Shokr, How can soil quality be accurately and quickly studied? A review, Agronomy 14 (8) (2024) 1682, http://dx.doi.org/10.3390/agronomy14081682.

[41] B. Kitchenham, S. Charters, Guidelines for Performing Systematic Literature Reviews in Software Engineering, Tech. Rep., EBSE Technical Report EBSE-2007-01, Keele University and University of Durham, 2007, URL https://www.elsevier.com/__data/promis_misc/525444systematicreviewsguide.pdf. (Accessed 6 November 2024).

[42] T. Dybå, T. Dingsøyr, G. Hanssen, Applying systematic reviews to diverse study types: An experience report, in: First International Symposium on Empirical Software Engineering and Measurement, ESEM 2007, IEEE, 2007, pp. 225–234, http://dx.doi.org/10.1109/ESEM.2007.27.

[43] R. van de Schoot, J. de Bruin, R. Schram, P. Zahedi, J. de Boer, F. Weijdema, B. Kramer, M. Huijts, M. Hoogerwerf, G. Ferdinands, A. Harkema, J. Willemsen, Y. Ma, Q. Fang, S. Hindriks, L. Tummers, D.L. Oberski, An open source machine learning framework for efficient and transparent systematic reviews, Nat. Mach. Intell. 3 (2021) 125–133, http://dx.doi.org/10.1038/s42256-020-00287-7.

[44] M.J. Page, J.E. McKenzie, P.M. Bossuyt, I. Boutron, T.C. Hoffmann, C.D. Mulrow, L. Shamseer, J.M. Tetzlaff, E.A. Akl, S.E. Brennan, R. Chou, J. Glanville, J.M. Grimshaw, A. Hróbjartsson, M.M. Lalu, T. Li, E.W. Loder, E. Mayo-Wilson, S. McDonald, L.A. McGuinness, L.A. Stewart, J. Thomas, A.C. Tricco, V.A. Welch, P. Whiting, D. Moher, The PRISMA 2020 statement: An updated guideline for reporting systematic reviews, Int. J. Surg. 88 (2021) http://dx.doi.org/10.1016/j.ijsu.2021.105906.

[45] A.M. Gómez, Q. de Jong van Lier, N.E. Silvero, L. Inforsato, M.L.A. de Melo, H.S. Rodríguez-Albarracín, N.A. Rosin, J.T.F. Rosas, R. Rizzo, J.A. Demattê, Digital mapping of the soil available water capacity: tool for the resilience of agricultural systems to climate change, Sci. Total Environ. 882 (2023) http://dx.doi.org/10.1016/j.scitotenv.2023.163572.

[46] ESRAF, Losan Database, Tech. Rep., Ente Regionale per i Servizi all' Agricoltura e alle Foreste, 2008, URL https://www.cursa.it/wp-content/uploads/2021/09/pasSAGGI-Anno-3-n.8-2017.pdf. (Accessed 4 February 2025).

[47] M.E. Angelini, B. Kempen, G.B. Heuvelink, A.J. Temme, M.D. Ransom, Extrapolation of a structural equation model for digital soil mapping, Geoderma 367 (2020) http://dx.doi.org/10.1016/j.geoderma.2020.114226.

[48] A.P. Hassler, E. Menasalvas, F.J. Garcia-Garcia, L. Rodriguez-Manas, A. Holzinger, Importance of medical data preprocessing in predictive modeling and risk factor discovery for the frailty syndrome, Springer/Nature BMC Med. Inform. Decis. Mak. 19 (1) (2019) 1–17, http://dx.doi.org/10.1186/s12911-019-0747-6.

[49] N.V. Chawla, K.W. Bowyer, L.O. Hall, W.P. Kegelmeyer, SMOTE: Synthetic minority over-sampling technique, J. Artificial Intelligence Res. 16 (2002) 321–357, http://dx.doi.org/10.1613/jair.953.

[50] K. Jain, R. John, N. Torbick, V. Kolluru, S. Saraf, A. Chandel, G.M. Henebry, M. Jarchow, Monitoring the spatial distribution of cover crops and tillage practices using machine learning and environmental drivers across Eastern South Dakota, Environ. Manag. 74 (2024) 742–756, http://dx.doi.org/10.1007/s00267-024-02021-0.

[51] R. Sarkar, R. Sarkar, C. Long, B. Northup, Ex-ante analyses using machine learning to understand the interactive influences of environmental and agromanagement variables for target-oriented management practice selection, Eur. J. Agron. 162 (2025) http://dx.doi.org/10.1016/j.eja.2024.127432.

[52] M. Dutta, D. Gupta, Y. Gulzar, M.S. Mir, C.W. Onn, A.B. Soomro, Leveraging inception V3 for precise early and late blight disease classification in potato crops, Trait. Du Signal 41 (2024) 705–715, http://dx.doi.org/10.18280/ts.410213.

[53] Y. Mo, R. Bier, X. Li, M. Daniels, A. Smith, L. Yu, J. Kan, Agricultural practices influence soil microbiome assembly and interactions at different depths identified by machine learning, Commun. Biol. 7 (2024) 1349, http://dx.doi.org/10.1038/s42003-024-07059-8, URL https://www.nature.com/articles/s42003-024-07059-8.

[54] Z. Zhai, F. Chen, H. Yu, J. Hu, X. Zhou, H. Xu, PS-MTL-LUCAS: A partially shared multi-task learning model for simultaneously predicting multiple soil properties, Ecol. Inform. 82 (2024) http://dx.doi.org/10.1016/j.ecoinf.2024.102784.

[55] L. Lotfollahi, M.A. Delavar, A. Biswas, M. Jamshidi, R. Taghizadeh-Mehrjardi, Modeling the spatial variation of calcium carbonate equivalent to depth using machine learning techniques, Environ. Monit. Assess. 195 (2023) http://dx.doi.org/10.1007/s10661-023-11126-8.

[56] E. Izquierdo-Verdiguier, R. Zurita-Milla, An evaluation of guided regularized random forest for classification and regression tasks in remote sensing, Int. J. Appl. Earth Obs. Geoinf. 88 (2020) http://dx.doi.org/10.1016/j.jag.2020.102051.

[57] L. Shi, S. O'Rourke, F.B. de Santana, K. Daly, Prediction of soil bulk density in agricultural soils using mid-infrared spectroscopy, Geoderma 434 (2023) http://dx.doi.org/10.1016/j.geoderma.2023.116487.

[58] S. Chen, Z. Chen, X. Zhang, Z. Luo, C. Schillaci, D. Arrouays, A.C. Richer-De-Forges, Z. Shi, European topsoil bulk density and organic carbon stock database (0-20 cm) using machine-learning-based pedotransfer functions, Earth Syst. Sci. Data 16 (2024) 2367–2383, http://dx.doi.org/10.5194/essd-16-2367-2024.

[59] A. Holzinger, K. Keiblinger, P. Holub, K. Zatloukal, H. Müller, AI for life: Trends in artificial intelligence for biotechnology, New Biotechnol. 74 (2023) 16–24, http://dx.doi.org/10.1016/j.nbt.2023.02.001.

[60] M. Hosseinpour-Zarnaq, F. Moshiri, M. Jamshidi, R. Taghizadeh-Mehrjardi, M.M. Tehrani, F.E. Meymand, Monitoring changes in soil organic carbon using satellite-based variables and machine learning algorithms in arid and semi-arid regions, Environ. Earth Sci. 83 (2024) http://dx.doi.org/10.1007/s12665-024-11876-9.

[61] T. Zhang, Y. Li, M. Wang, Remote sensing-based prediction of organic carbon in agricultural and natural soils influenced by salt and sand mining using machine learning, J. Environ. Manag. 352 (2024) http://dx.doi.org/10.1016/j.jenvman.2024.120107.

[62] J. Ou, Z. Wu, Q. Yan, X. Feng, Z. Zhao, Improving soil organic carbon mapping in farmlands using machine learning models and complex cropping system information, Environ. Sci. Eur. 36 (2024) http://dx.doi.org/10.1186/s12302-024-00912-x.

[63] N. Samarinas, N.L. Tsakiridis, S. Kokkas, E. Kalopesa, G.C. Zalidis, Soil data cube and artificial intelligence techniques for generating national-scale topsoil thematic maps: A case study in Lithuanian croplands, Remote. Sens. 15 (2023) http://dx.doi.org/10.3390/rs15225304.

[64] S. Naimi, S. Ayoubi, M. Zeraatpisheh, J.A.M. Dematte, Ground observations and environmental covariates integration for mapping of soil salinity: A machine learning-based approach, Remote. Sens. 13 (2021) http://dx.doi.org/10.3390/rs13234825.

[65] N. Wang, J. Peng, J. Xue, X. Zhang, J. Huang, A. Biswas, Y. He, Z. Shi, A framework for determining the total salt content of soil profiles using time-series sentinel-2 images and a random forest-temporal convolution network, Geoderma 409 (2022) http://dx.doi.org/10.1016/j.geoderma.2021.115656.

[66] B. Hu, M. Xie, Z. Shi, H. Li, S. Chen, Z. Wang, Y. Zhou, H. Ni, Y. Geng, Q. Zhu, X. Zhang, Fine-resolution mapping of cropland topsoil pH of southern China and its environmental application, Geoderma 442 (2024) http://dx.doi.org/10.1016/j.geoderma.2024.116798.

[67] O. Yüzügüllü, N. Fajraoui, F. Liebisch, Soil texture and pH mapping using remote sensing and support sampling, IEEE J. Sel. Top. Appl. Earth Obs. Remote. Sens. 17 (2024) 12685–12705, http://dx.doi.org/10.1109/JSTARS.2024.3422494.

[68] A.D. Raj, S. Kumar, K.R. Sooryamol, K.J. George, Soil erodibility mapping using remote sensing and in situ soil data with random forest model in a mountainous catchment of Indian Himalayas, Environ. Monit. Assess. 196 (2024) http://dx.doi.org/10.1007/s10661-024-13173-1.

[69] B. Padmavathi, A. BhagyaLakshmi, G. Vishnupriya, K. Datchanamoorthy, IoT-based prediction and classification framework for smart farming using adaptive multi-scale deep networks, Expert Syst. Appl. 254 (2024) http://dx.doi.org/10.1016/j.eswa.2024.124318.

[70] S. Sow, S. Ranjan, M.F. Seleiman, H.M. Alkharabsheh, M. Kumar, N. Kumar, S.R. Padhan, D.K. Roy, D. Nath, H. Gitari, D.O. Wasonga, Artificial intelligence for maximizing agricultural input use efficiency: Exploring nutrient, water and weed management strategies, Phyton-Int. J. Exp. Bot. 93 (2024) 1569–1598, http://dx.doi.org/10.32604/phyton.2024.052241.

[71] S. Maleki, A. Karimi, A. Mousavi, R. Kerry, R. Taghizadeh-Mehrjardi, Delineation of soil management zone maps at the regional scale using machine learning, Agronomy 13 (2023) http://dx.doi.org/10.3390/agronomy13020445.

[72] A.B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, R. Chatila, F. Herrera, Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI, Inf. Fusion 58 (2020) 82–115, http://dx.doi.org/10.1016/j.inffus.2019.12.012.

[73] A. Bennetot, I. Donadello, A. El Qadi El Haouari, M. Dragoni, T. Frossard, B. Wagner, A. Sarranti, M. Trocan, R. Chatila, A. Holzinger, A.D. Garcez, N. Diaz-Rodriguez, A practical tutorial on explainable AI techniques, ACM Comput. Surv. 57 (2) (2025) 1–44, http://dx.doi.org/10.1145/3670685.

[74] M. Kotli, G. Piir, U. Maran, Pesticide effect on earthworm lethality via interpretable machine learning, J. Hazard. Mater. 461 (2024) http://dx.doi.org/10.1016/j.jhazmat.2023.132577.

[75] S. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, 2017, URL http://arxiv.org/abs/1705.07874.

[76] H. Gharakhani, J.A. Thomasson, Y. Lu, Integration and preliminary evaluation of a robotic cotton harvester prototype, Comput. Electron. Agric. 211 (2023) http://dx.doi.org/10.1016/j.compag.2023.107943.

[77] E.K. Bünemann, G. Bongiorno, Z. Bai, R.E. Creamer, G.D. Deyn, R. de Goede, L. Fleskens, V. Geissen, T.W. Kuyper, P. Mäder, M. Pulleman, W. Sukkel, J.W. van Groenigen, L. Brussaard, Soil quality – a critical review, Soil Biol. Biochem. 120 (2018) 105–125, http://dx.doi.org/10.1016/j.soilbio.2018.01.030.

[78] S. Wieser, K.M. Keiblinger, A. Mentler, C. Rosinger, K. Wriessnig, N. Bruhn, L.G. Bernardini, M. Bieber, S. Huber, G. Bodner, Labile not stable SOC fractions constitute the manageable drivers of soil health advances in carbon farming, Geoderma 449 (2024) http://dx.doi.org/10.1016/j.geoderma.2024.116991.

[79] D.W. Pribyl, A critical review of the conventional SOC to SOM conversion factor, Geoderma 156 (2010) 75–83, http://dx.doi.org/10.1016/j.geoderma.2010.02.003.

[80] M. Pal, Random forest classifier for remote sensing classification, Int. J. Remote Sens. 26 (2005) 217–222, http://dx.doi.org/10.1080/01431160412331269698.

[81] J.S. Evans, M.A. Murphy, Z.A. Holden, S.A. Cushman, Modeling species distribution and change using random forest, in: Predictive Species and Habitat Modeling in Landscape Ecology: Concepts and Applications, Springer New York, New York, NY, 2011, pp. 139–159, http://dx.doi.org/10.1007/978-1-4419-7390-0_8.

[82] S.K. Salas, Modified shape index for object-based random forest image classification of agricultural systems using airborne hyperspectral datasets, PLoS One 14 (3) (2019) 1–22, http://dx.doi.org/10.1371/journal.pone.0213356.

[83] R. Genuer, J.M. Poggi, C. Tuleau-Malot, Variable selection using random forests, Pattern Recognit. Lett. 31 (2010) 2225–2236, http://dx.doi.org/10.1016/j.patrec.2010.03.014.

[84] E.A.L. Salas, S.S. Kumaran, Perimeter-area soil carbon index (PASCI): modeling and estimating soil organic carbon using relevant explicatory waveband variables in machine learning environment, Geo-Spatial Inf. Sci. (2023) http://dx.doi.org/10.1080/10095020.2023.2211612.

[85] T. Chen, C. Guestrin, XGBoost: A scalable tree boosting system, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, Association for Computing Machinery, New York, NY, USA, 2016, pp. 785–794, http://dx.doi.org/10.1145/2939672.2939785.

[86] L. Prokhorenkova, G. Gusev, A. Vorobev, A.V. Dorogush, A. Gulin, CatBoost: unbiased boosting with categorical features, in: Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS '18, Curran Associates Inc., Red Hook, NY, USA, 2018, pp. 6639–6649.

[87] A.B. Cengiz, A.S. McGough, How much data do I need? A case study on medical data, in: 2023 IEEE International Conference on Big Data, BigData, IEEE, 2023, pp. 3688–3697, http://dx.doi.org/10.1109/BigData59044.2023.10386440.

[88] D. McElfresh, S. Khandagale, J. Valverde, V.P. C., B. Feuer, C. Hegde, G. Ramakrishnan, M. Goldblum, C. White, When do neural nets outperform boosted trees on tabular data?, 2024, arXiv:2305.02997. URL https://arxiv.org/abs/2305.02997.

[89] P.J. Huber, Robust Statistics, in: Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, New York, 1981, http://dx.doi.org/10.1002/0471725250.

[90] H.E. Beck, T.R. McVicar, N. Vergopolan, A. Berg, N.J. Lutsko, A. Dufour, Z. Zeng, X. Jiang, A.I. van Dijk, D.G. Miralles, High-resolution (1 km) Köppen-Geiger maps for 1901–2099 based on constrained CMIP6 projections, Sci. Data 10 (2023) http://dx.doi.org/10.1038/s41597-023-02549-6.

[91] J. Keane, S. Keyser, J. Kedziora, Strategy masking: A method for guardrails in value-based reinforcement learning agents, 2025, http://dx.doi.org/10.48550/arXiv.2501.05501, arXiv:2501.05501.

[92] Y. Zha, Y. Yang, Innovative graph neural network approach for predicting soil heavy metal pollution in the Pearl River Basin, China, Sci. Rep. 14 (2024) http://dx.doi.org/10.1038/s41598-024-67175-7.

[93] K. Coleman, D.S. Jenkinson, RothC-26.3-A model for the turnover of carbon in soil, in: Evaluation of Soil Organic Matter Models: Using Existing Long-Term Datasets, Springer, 1996, pp. 237–246, http://dx.doi.org/10.1007/978-3-642-61094-3_17.

[94] A. Taghizadeh-Toosi, B.T. Christensen, N.J. Hutchings, J. Vejlin, T. Kätterer, M. Glendining, J.E. Olesen, C-TOOL: A simple model for simulating whole-profile carbon storage in temperate agricultural soils, Ecol. Model. 292 (2014) 11–25, http://dx.doi.org/10.1016/j.ecolmodel.2014.08.016.

[95] G. Hu, F. You, An AI framework integrating physics-informed neural network with predictive control for energy-efficient food production in the built environment, 2023.

[96] S. Cuomo, M.D. Rosa, F. Piccialli, L. Pompameo, V. Vocca, A numerical approach for soil microbiota growth prediction through physics-informed neural network, Appl. Numer. Math. 207 (2025) 97–110, http://dx.doi.org/10.1016/j.apnum.2024.08.025.

[97] L.J. Koppensteiner, H.P. Kaul, S. Raubitzek, P. Weihs, P. Euteneuer, J. Bernas, G. Moitzi, T. Neubauer, A. Klimek-Kopyra, N. Barta, R.W. Neugschwandt-ner, Estimating wheat traits using artificial neural network-based radiative transfer model inversion, Remote. Sens. 17 (2025) http://dx.doi.org/10.3390/rs17111904.

[98] Y. Liu, A. Wang, B. Li, J. Šimůnek, R. Liao, Combining mathematical models and machine learning algorithms to predict the future regional-scale actual transpiration by maize, Agricult. Water. Manag. 303 (2024) http://dx.doi.org/10.1016/j.agwat.2024.109056.

[99] Z. Liu, P. Ma, Y. Wang, W. Matusik, M. Tegmark, KAN 2.0: Kolmogorov-Arnold networks meet science, 2024, URL http://arxiv.org/abs/2408.10205.

[100] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural Comput. 9 (8) (1997) 1735–1780, http://dx.doi.org/10.1162/neco.1997.9.8.1735.

[101] Y. Wu, Z. Yang, Y. Liu, Internet-of-things-based multiple-sensor monitoring system for soil information diagnosis using a smartphone, Micromachines 14 (2023) http://dx.doi.org/10.3390/mi14071395.

[102] Y. Ning, H. Kazemi, P. Tahmasebi, A comparative machine learning study for time series oil production forecasting: ARIMA, LSTM, and prophet, Comput. Geosci. 164 (2022) http://dx.doi.org/10.1016/j.cageo.2022.105126.

[103] C.B.A. Satrio, W. Darmawan, B.U. Nadia, N. Hanafiah, Time series analysis and forecasting of coronavirus disease in Indonesia using ARIMA model and PROPHET, in: Procedia Computer Science, 179, Elsevier B.V., 2021, pp. 524–532, http://dx.doi.org/10.1016/j.procs.2021.01.036.

[104] H. Jaeger, H. Haas, Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication, Science 304 (5667) (2004) 78–80.

[105] M. Aoki, State space modeling of time series, Springer Science and Business Media, Heidelberg, 2013, http://dx.doi.org/10.1007/978-3-642-75883-6.

[106] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, W. Zhang, Informer: Beyond efficient transformer for long sequence time-series forecasting, 2020, http://dx.doi.org/10.48550/arXiv.2012.07436.

[107] M.T. Ribeiro, S. Singh, C. Guestrin, Why should i trust you?: Explaining the predictions of any classifier, in: 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2016, ACM, 2016, pp. 1135–1144, http://dx.doi.org/10.1145/2939672.2939778.

[108] E. Doumard, J. Aligon, E. Escriva, J.-B. Excoffier, P. Monsarrat, C. Soulé-Dupuy, A quantitative approach for the comparison of additive local explanation methods, Inf. Syst. 114 (2023) 102162, http://dx.doi.org/10.1016/j.is.2022.102162.

[109] A. Holzinger, I. Fister Jr., I. Fister, H.-P. Kaul, S. Asseng, Human-centered AI in smart farming: Towards agriculture 5.0, IEEE Access 12 (2024) 62199–62214, http://dx.doi.org/10.1109/ACCESS.2024.3395532.

[110] B. Zhao, W. Jin, J.D. Ser, G. Yang, ChatAgri: Exploring potentials of ChatGPT on cross-linguistic agricultural text classification, Neurocomputing 557 (2023) http://dx.doi.org/10.1016/j.neucom.2023.126708.

[111] A. Tzachor, M. Devare, C. Richards, P. Pypers, A. Ghosh, J. Koo, S. Johal, B. King, Large language models and agricultural extension services, Nat. Food 4 (2023) 941–948, http://dx.doi.org/10.1038/s43016-023-00867-x.

[112] S.S. Andrews, D.L. Karlen, J.P. Mitchell, A comparison of soil quality indexing methods for vegetable production systems in Northern California, Agricult. Ecosys. Environ. 90 (2002) 25–45.

[113] Y. Zhou, H. Ma, Y. Xie, X. Jia, T. Su, J. Li, Y. Shen, Assessment of soil quality indexes for different land use types in typical steppe in the Loess Hilly Area, China, Ecol. Indic. 118 (2020) http://dx.doi.org/10.1016/j.ecolind.2020.106743.

[114] C.J. Feeney, L. Bentley, D.D. Rosa, P. Panagos, B.A. Emmett, A. Thomas, D.A. Robinson, Benchmarking soil organic carbon (SOC) concentration provides more robust soil health assessment than the SOC/clay ratio at European scale, Sci. Total Environ. 951 (2024) http://dx.doi.org/10.1016/j.scitotenv.2024.175642.

[115] A. Hosna, E. Merry, J. Gyalmo, Z. Alom, Z. Aung, M.A. Azim, Transfer learning: a friendly introduction, J. Big Data 9 (1) (2022) 102.

[116] K. Weiss, T.M. Khoshgoftaar, D. Wang, A survey of transfer learning, J. Big Data 3 (2016) 1–40, http://dx.doi.org/10.1186/s40537-016-0043-6.

[117] C. Zhang, Y. Xie, H. Bai, B. Yu, W. Li, Y. Gao, A survey on federated learning, Knowl.-Based Syst. 216 (2021) 106775, http://dx.doi.org/10.1016/j.knosys.2021.106775.

[118] M.V. Luzón, N. Rodríguez-Barroso, A. Argente-Garrido, D. Jiménez-López, J.M. Moyano, J. Del Ser, W. Ding, F. Herrera, A tutorial on federated learning from theory to practice: Foundations, software frameworks, exemplary use cases, and selected trends, IEEE/CAA J. Autom. Sin. 11 (4) (2024) 824–850, http://dx.doi.org/10.1109/JAS.2024.124215.

[119] M. Huelser, H. Mueller, N. Díaz-Rodríguez, A. Holzinger, On the disagreement problem in human-in-the-loop federated machine learning, J. Ind. Inf. Integr. 45 (5) (2025) 100827, http://dx.doi.org/10.1016/j.jii.2025.100827.

[120] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, J. Tang, Self-supervised learning: Generative or contrastive, IEEE Trans. Knowl. Data Eng. 35 (1) (2021) 857–876, http://dx.doi.org/10.1109/TKDE.2021.3090866.

[121] T. Baltrušaitis, C. Ahuja, L.-P. Morency, Multimodal machine learning: A survey and taxonomy, IEEE Trans. Pattern Anal. Mach. Intell. 41 (2) (2018) 423–443, http://dx.doi.org/10.1109/TPAMI.2018.2798607.

[122] A. Holzinger, M. Dehmer, F. Emmert-Streib, R. Cucchiara, I. Augenstein, J.D. Ser, W. Samek, I. Jurisica, N. Díaz-Rodríguez, Information fusion as an integrative cross-cutting enabler to achieve robust, explainable, and trustworthy medical artificial intelligence, Inf. Fusion 79 (2022) 263–278, http://dx.doi.org/10.1016/j.inffus.2021.10.007.

[123] A. Diez-Olivan, J.D. Ser, D. Galar, B. Sierra, Data fusion and machine learning for industrial prognosis: Trends and perspectives towards industry 4.0, Inf. Fusion 50 (2019) 92–111, http://dx.doi.org/10.1016/j.inffus.2018.10.005.

[124] S. Huber, L.G. Bernardini, A. Bennett, J. Fohrafellner, K. Dohnke, M. Bieber, F. Vuolo, A. Mentler, G. Bodner, K. Keiblinger, Suitability of microbial and organic matter indicators for on-farm soil health monitoring, Soil Use Manag. 40 (2024) http://dx.doi.org/10.1111/sum.12993.

[125] P. Panagos, N. Broothaerts, C. Ballabio, A. Orgiazzi, D. De Rosa, P. Borrelli, L. Liakos, D. Vieira, E. Van Eynde, C. Arias Navarro, How the EU soil observatory is providing solid science for healthy soils, Eur. J. Soil Sci. 75 (3) (2024) e13507, http://dx.doi.org/10.1111/ejss.13507.

[126] T.P. Pagano, et al., Bias and unfairness in machine learning models: a systematic review on datasets, tools, fairness metrics, and identification and mitigation methods, Big Data Cogn. Comput. 7 (1) (2023) 15.

[127] C. Schillaci, A. Perego, E. Valkama, M. Märker, S. Saia, F. Veronesi, A. Lipani, L. Lombardo, T. Tadiello, H.A. Gamper, New pedotransfer approaches to predict soil bulk density using WoSIS soil data and environmental covariates in Mediterranean agro-ecosystems, Sci. Total Environ. 780 (2021) 146609, http://dx.doi.org/10.1016/j.scitotenv.2021.146609.

[128] H. Wang, C. Zheng, S. Ning, C. Cao, K. Li, H. Dang, Y. Wu, J. Zhang, Impacts of long-term saline water irrigation on soil properties and crop yields under maize-wheat crop rotation, Agricult. Water. Manag. 286 (2023) 108383, http://dx.doi.org/10.1016/j.agwat.2023.108383.

[129] E.A. Davidson, I.A. Janssens, Temperature sensitivity of soil carbon decomposition and feedbacks to climate change, Nature 440 (2006) 165–173, http://dx.doi.org/10.1038/nature04514.

[130] F. Ehrlich-Sommer, F. Hoenigsberger, C. Gollob, A. Nothdurft, K. Stampfer, A. Holzinger, Sensors for digital transformation in smart forestry, Sensors 24 (3) (2024) 798, http://dx.doi.org/10.3390/s24030798.

[131] S. Das, D. Panday, Chapter 29: Soil health assessment and spatial characterization using remote sensing, in: S. Dharumarajan (Ed.), Remote Sensing of Soils, Elsevier, Amsterdam, 2024, pp. 455–467, http://dx.doi.org/10.1016/B978-0-443-18773-5.00034-X.

[132] Y. Liang, E.F. Leifheit, A. Lehmann, M.C. Rillig, Soil organic carbon stabilization is influenced by microbial diversity and temperature, Sci. Rep. 15 (1) (2025) 13990, http://dx.doi.org/10.1038/s41598-025-98009-9.

[133] C. Rosinger, K. Keiblinger, M. Bieber, L.G. Bernardini, S. Huber, A. Mentler, O. Sae-Tun, B. Scharf, G. Bodner, On-farm soil organic carbon sequestration potentials are dominated by site effects, not by management practices, Geoderma 433 (2023) http://dx.doi.org/10.1016/j.geoderma.2023.116466.

[134] R. Yadav, R. Kumar, R.K. Gupta, T. Kaur, A. Kour, S. Kaur, A. Rajput, Heavy metal toxicity in earthworms and its environmental implications: A review, Environ. Adv. 12 (2023) 100374, http://dx.doi.org/10.1016/j.envadv.2023.100374.

[135] Z. Jiang, H. Tu, B. Bai, C. Yang, B. Zhao, Z. Guo, Q. Liu, H. Zhao, W. Yang, L. Xiong, J. Zhang, Combining UAV-RGB high-throughput field phenotyping and genome-wide association study to reveal genetic variation of rice germplasms in dynamic response to drought stress, New Phytol. 232 (2021) 440–455, http://dx.doi.org/10.1111/nph.17580.

[136] Y. Dong, X. Wang, S. Wang, B. Li, J. Liu, J. Huang, X. Li, Y. Zeng, W. Su, Enhancing soil organic carbon prediction by unraveling the role of crop residue coverage using interpretable machine learning, Geoderma 455 (2025) http://dx.doi.org/10.1016/j.geoderma.2025.117225.

[137] M.W. Zhang, X.L. Sun, M.N. Zhang, H.X. Yang, H.J. Liu, H.X. Li, Improved soil organic matter monitoring by using cumulative crop residue indices derived from time-series remote sensing images in the central black soil region of China, Soil Tillage Res. 246 (2025) http://dx.doi.org/10.1016/j.still.2024.106357.

[138] A. Ali, H.P. Kaul, Monitoring yield and quality of forages and grassland in the view of precision agriculture applications—A review, 17, 2025, http://dx.doi.org/10.3390/rs17020279,

[139] A. Ali, M.U. Hassan, H.P. Kaul, Broad scope of site-specific crop management and specific role of remote sensing technologies within it—A review, J. Agron. Crop Sci. 210 (2024) http://dx.doi.org/10.1111/jac.12732.

[140] O.D. Adeniyi, A. Brenning, M. Maerker, Spatial prediction of soil organic carbon: Combining machine learning with residual kriging in an agricultural lowland area (Lombardy region, Italy), Geoderma 448 (2024) http://dx.doi.org/10.1016/j.geoderma.2024.116953.

[141] K.V. Looy, J. Bouma, M. Herbst, J. Koestel, B. Minasny, U. Mishra, C. Montzka, A. Nemes, Y.A. Pachepsky, J. Padarian, M.G. Schaap, B. Tóth, A. Verhoef, J. Vanderborght, M.J. van der Ploeg, L. Weihermüller, S. Zacharias, Y. Zhang, H. Vereecken, Pedotransfer functions in earth system science: Challenges and perspectives, Rev. Geophys. 55 (2017) 1199–1256, http://dx.doi.org/10.1002/2017RG000581.

[142] S. González, S. García, J.D. Ser, L. Rokach, F. Herrera, A practical tutorial on bagging and boosting based ensembles for machine learning: Algorithms, software tools, performance study, practical perspectives and opportunities, Inf. Fusion 64 (2020) 205–237, http://dx.doi.org/10.1016/j.inffus.2020.07.007.

[143] M. Abdar, F. Pourpanah, S. Hussain, D. Rezazadegan, L. Liu, M. Ghavamzadeh, P. Fieguth, X. Cao, A. Khosravi, U.R. Acharya, A review of uncertainty quantification in deep learning: Techniques, applications and challenges, Inf. Fusion 76 (2021) 243–297, http://dx.doi.org/10.1016/j.inffus.2021.05.008.

[144] A.N. Angelopoulos, S. Bates, Conformal prediction: A gentle introduction, Found. Trends® Mach. Learn. 16 (4) (2023) 494–591, http://dx.doi.org/10.1561/2200000101.

[145] M. Awad, R. Khanna, M. Awad, R. Khanna, Support Vector Regression, A Press, Berkeley (CA), 2015, pp. 67–80, http://dx.doi.org/10.1007/978-1-4302-5990-9_4.

[146] Z. Cao, Y. Wang, Bayesian model comparison and selection of spatial correlation functions for soil parameters, Struct. Saf. 49 (2014) 10–17, http://dx.doi.org/10.1016/j.strusafe.2013.06.003.

[147] M. Raissi, P. Perdikaris, G. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, J. Comput. Phys. 378 (2019) 686–707, http://dx.doi.org/10.1016/j.jcp.2018.10.045, URL https://www.sciencedirect.com/science/article/pii/S0021999118307125.

[148] G.E. Karniadakis, I.G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, L. Yang, Physics-informed machine learning, Nat. Rev. Phys. 3 (6) (2021) 422–440, http://dx.doi.org/10.1038/s42254-021-00314-5.

**Stefan Schweng** is an Assistant Researcher at Human-Centered AI Lab Vienna, where he contributes to advancing research in explainable artificial intelligence, with applications in smart farming and forestry. He is currently pursuing a doctoral degree in natural resources and life sciences, focusing on smart agriculture. His research interests include multi-objective optimization, machine vision, 3D point cloud analysis, decision support systems, and dynamic systems modeling.
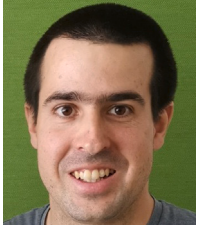
**Luca Bernardini** is an Assistant Researcher at the Institute of Soil Research and the Institute of Agronomy at BOKU University in Vienna. He is currently completing a doctoral degree in natural resources and life sciences, with a focus on modeling using process-based and ML models. His research interests include leveraging geospatial data and models of all kinds to inform decision-making for relevant stakeholders, meta-modeling, remote sensing applications and on-farm research. His overarching goal is to combine pedoclimatic data with crop information to support monitoring and decision-making towards a sustainable transition in Agriculture.

**Katharina Keiblinger** is a Full Professor at the Institute of Soil Research, BOKU University (Vienna), where she leads a working group in soil microbiology and microbial ecology. Her research focuses on the role of soil microbial communities in nutrient cycling, soil health, and ecosystem functioning. She currently leads several national and EU-funded projects, including the "living lab" initiative, which investigates how regenerative agricultural practices can restore degraded soils and enhance long-term soil health. Her scientific interests include soil biodiversity, microbial-driven nutrient dynamics, soil quality assessment, and the integration of mechanistic understanding with AI-based modeling approaches.
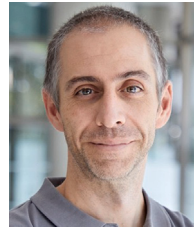
**Hans-Peter Kaul** received Dipl.Ing. (equivalent to M.Sc.) and Ph.D. degrees from the University of Bonn, Germany, and the Habilitation degree from Hohenheim University, Germany, in 1998. He is currently a recognized expert in the field of advanced cropping systems, precision farming, hyperspectral remote sensing, and crop modeling. He has published high quality research articles regarding sustainable-based cropping and crop modeling. Since 2001, he has been a Full Professor in agronomy and grassland management with BOKU University, Austria, where he acts as the Head of the Institute of Agronomy.

**Iztok Fister Jr.** received a B.Sc., M.Sc., and Ph.D. in Computer Science from the University of Maribor, Slovenia. He is currently an Associate Professor at the University of Maribor. He has published over 200 research articles in refereed journals, conferences, and book chapters. His research interests include Data Mining, Pervasive Computing, Optimization, and Sports Science. He has acted as a Program Committee member of more than 30 international conferences.

**Niko Lukač** is Professor of Computer Science at the University of Maribor, Slovenia. His research in 3D point cloud processing, GeoAI, environmental modeling, and parallel computing has led numerous national and international R&D projects with publications in high ranking journals, and he has assumed editorial roles in various scientific journals. From 2019 to 2022, he contributed as an executive committee member to the prominent European Umbrella Organisation for Geographic Information (EUROGI). He is expert advisor to the EC and in 2024, the University of Maribor honored him for his exceptional research achievement.

**Javier Del Ser** holds a telecom engineering degree from the University of the Basque Country (2003), a Ph.D. in Control Engineering and Industrial Electronics from the University of Navarra (2006, Cum Laude), and a second Ph.D. in Information and Communication Technologies from the University of Alcalá de Henares (2013, Cum Laude and recipient of the Extraordinary Ph.D. Award). He is a Research Professor at TECNALIA and a Distinguished Researcher at the Department of Mathematics of the University of the Basque Country (UPV/EHU). He is Visiting Professor at the University of Granada (Spain) and BOKU University, Vienna (Austria). His research focuses on AI and Machine Learning, with applications in practical modeling and optimization problems across diverse sectors, including industry, healthcare, transportation, energy, and mobility.

**Andreas Holzinger** advocates a synergistic approach of Human-Centered AI aligning AI with human values, ethical principles, and legal issues to ensure secure, safe and controllable AI. For his achievements he was elected 2019 to Academia Europaea, ELLIS 2020, and ifip Fellow 2021. He obtained his Ph.D. in Cognitive Science in 1998, and his second Ph.D. in Computer Science from TU Graz in 2003. He was Visiting Professor for Machine Learning & Knowledge Extraction in Verona (Italy), RWTH Aachen (Germany), UCL London (UK) and at the Alberta Machine Intelligence Institute in Edmonton (Canada). Since 2022 he is full professor for digital transformation at BOKU University, Vienna, Austria.