

SoilGrids 2.0: producing ~~quality-assessed~~ soil information for the globe with quantified spatial uncertainty

Luis M. de Sousa¹, Laura Poggio¹, Niels H. Batjes¹, Gerard B. M. Heuvelink¹, Bas Kempen¹, Eloi Ribeiro¹, and David Rossiter¹

¹ISRIC - World Soil Information - Wageningen (NL)

Correspondence: Laura Poggio (laura.poggio@isric.org)

Abstract. SoilGrids produces maps of soil properties for the entire globe at medium spatial resolution (250 metres cell size) using state-of-the-art machine learning methods to generate the necessary models. It takes as inputs soil observations from about 240 000 locations worldwide and over 400 global environmental covariates describing vegetation, terrain morphology, climate, geology and hydrology. The aim of this work was the production of ~~quality-assessed~~ global maps of soil properties, with cross-validation, hyper-parameters selection and quantification of spatially explicit uncertainty, as implemented in the SoilGrids version 2.0 product incorporating state of the art practices and adapting them for global digital soil mapping with legacy data. The paper presents the evaluation of the global predictions produced for soil organic carbon content, total nitrogen, coarse fragments, pH(water), cation exchange capacity, bulk density and texture fractions at six standard depths (up to 200 cm). The quantitative evaluation showed metrics in line with previous global, continental and large regions studies. The qualitative evaluation showed that coarse scale patterns are well reproduced. The spatial uncertainty at global scale highlighted the need for more soil observations, especially in high latitude regions.

1 Introduction

Healthy soils provide important ecosystem services at the local, landscape and global level, and are important for the functioning of terrestrial ecosystems (Banwart et al., 2014; FAO and ITPS, 2015; UNEP, 2012). ~~Up-to-date information~~ Information on world soil resources, based on the currently 'best-available' (shared) soil profile data, at a scale level commensurate with user needs, is required to address a range of pressing global issues. These include avoiding and reducing soil erosion through land rehabilitation and development (Borrelli et al., 2017; WOCAT, 2007), mitigation and adaptation to climate change (Batjes, 2019; Harden et al., 2017; Sanderman et al., 2017; Yigini and Panagos, 2016; Smith et al., 2019), ensuring water security (Rockstroem et al., 2012), food production and food security (FAO et al., 2018; Soussana et al., 2017; Springmann et al., 2018), as well as the preservation of biodiversity (Barnes, 2015; IPBES, 2019; van der Esch et al., 2017) and human livelihood (Bouma, 2015).

~~Quality-assessed~~ The 'best available' soil data are required to support the Land Degradation Neutrality (LDN) (Cowie et al., 2018) initiative, achieve several of the Sustainable Development Goals, and provide input for e.g. earth system modelling by the IPCC (Dai et al., 2019; Luo et al., 2016; Todd-Brown et al., 2013) and crop modelling (Han et al., 2019; van Bussel

25 et al., 2015; van Ittersum et al., 2013), among many other applications. Such information can in turn help inform international conventions such as the United Nations Framework Convention on Climate Change (UNFCCC), the United Nation Convention to Combat Desertification (UNCCD), and the United Nations Convention on Biological Diversity (UNCBD).

Until the last decade, most global scale assessments requiring soil data used the Digital Soil Map of the World (DSMW) FAO (1995), an updated version of the original printed 1:5 M scale Soil Map of the World (SMW) (FAO-Unesco, 1971-1981).
30 The soil-geographical data from the DSMW provided the basis for generating a range of derived soil property databases that drew on a larger selection of soil profile data held in the WISE database (Batjes, 2012) and more sophisticated (taxotransfer) procedures for deriving various soil properties (Batjes et al., 2007). Subsequently, in a joint effort coordinated by the Food and Agriculture Organization of the United Nations (FAO), the best available (newer) soil information collated for central and southern Africa, China, Europe, northern Eurasia, and Latin America were combined into a new product known as the
35 Harmonised World Soil Database (HWSD) (FAO et al., 2012).

Until recently, the HWSD was the only digital map annex database available for global analyses. However, it has several limitations (GSP and FAO, 2016; Hengl et al., 2014; Ivushkin et al., 2019; Omuto et al., 2012). Some of these relate to the partly outdated soil-geographic data as well as the use of a two-layer model (0–30 and 30–100 cm) for deriving soil properties. Others concern the derived attribute data themselves, in particular their unquantified uncertainty, and the use of three different
40 versions of the FAO legend (i.e., FAO74, FAO85 and FAO90). These issues have been addressed to varying degrees in various new global soil datasets (Batjes, 2016; Shangguan et al., 2014; Stoorvogel et al., 2017) that still largely draw on a traditional, pedology-based mapping approach (Dai et al., 2019).

The last decade, Digital Soil Mapping (DSM) became a widely used approach to obtain maps of soil information (Minasny and McBratney, 2016). DSM consists primarily in building a qualitative-quantitative numerical model between soil obser-
45 vations and environmental information acting as proxies for the soil forming factors (McBratney et al., 2003; Minasny and McBratney, 2016). DSM can also integrate direct information as proxies for soil properties, for example proximal sensing measurements. The number of studies using DSM to produce maps of soil properties is ever growing. Numerous modelling approaches are considered, from linear model to geo-statistics, machine learning and artificial intelligence (e.g., deep learning). Keskin and Grunwald (2018) provide a recent review of methods and applications in the field of DSM. DSM tech-
50 niques have been applied at various spatial resolutions (e.g., 30m to 1000m) to support precision farming (~~Piikki et al., 2017~~) (e.g. Piikki et al., 2017) as well as applications at landscape (~~Ellili et al., 2019; Kempen et al., 2015~~), country (~~Mora-Vallejo et al., 2008; Ni~~), regional (~~Dorji et al., 2014; Moulatlet et al., 2017~~), continental (~~Grunwald et al., 2011; Guevara et al., 2018; Hengl et al., 2017a~~) (e.g. Ellili et al., 2019; Kempen et al., 2015), country (e.g. Mora-Vallejo et al., 2008; Nijbroek et al., 2018; Vitharana et al., 2019; Poggio a), regional (e.g. Dorji et al., 2014; Moulatlet et al., 2017), continental (e.g. Grunwald et al., 2011; Guevara et al., 2018; Hengl et al., 2017a),
55 , and global level (~~Hengl et al., 2014, 2017b; GSP and ITPS, 2018; Stockmann et al., 2015~~) (e.g. Hengl et al., 2014, 2017b; GSP and ITPS,

The aim of this paper is to present the development of new soil property maps for the world at 250 metres grid resolution with a process incorporating state-of-the-art practices and adapting them to the challenges of global Digital Soil Mapping with legacy data. It builds on previous global soil properties maps (SoilGrids250m) (Hengl et al., 2017b), integrating up-to-date

60 machine learning methods, the increased availability of standardised soil profile data for the world (Batjes et al., 2020) and environmental eo-variates-covariates (Nussbaum et al., 2018; Poggio et al., 2013; Reuter and Hengl, 2012). In particular, this paper addresses at global scale the following elements:

1. quality-assessed-Incorporation of soil profile data derived from ISRIC’s World Soil Information Service (WoSIS), with expanded number and spatial distribution of observations (Batjes et al., 2020);
- 65 2. a reproducible co-variant selection procedure, relying on Recursive Feature Elimination (Guyon et al., 2002);
3. improved cross-validation procedure, based on spatial stratification; and
4. quantification of prediction uncertainty using Quantile Regression Forest (Meinshausen, 2006).

2 Materials and Methods

This study uses Quantile Regression Forest (Meinshausen, 2006), a method with a limited number of parameters to be tuned and
70 that has proven an effective compromise between accuracy and feasibility for large datasets. Selected primary soil properties as defined and described in the GlobalSoilMap specifications (Arrouays et al., 2014) were modelled. The following sections describe each step of the workflow (Figure 1) in detail. These include:

1. Input soil data preparation
2. Covariates selection
- 75 3. Model tuning and cross-validation
4. Final model fitting for prediction
5. Predictions with uncertainty estimation.

2.1 Soil observation data

Soil property data for this study were derived from the ISRIC World Soil Information Service (WoSIS), which provides con-
80 sistent, standardised soil profile data for the world (Batjes et al., 2020). All soil data shared with ISRIC to support global mapping activities are first stored in the ISRIC Data Repository together with their metadata (including the name of the data owner and licence defining access rights). Subsequently, the source data are imported ‘as is’ into PostgreSQL, after which they are ingested into the WoSIS data model itself. Following data quality assessment and control -(including consistency checks on latitude-longitude and depth of horizon/layer; flagging of duplicate profiles; and providing measures for geographic and attribute accuracy, as well as time stamps), the descriptions for the soil analytical methods and the units of measurement
85 are standardised using consistent procedures, with additional checks for possible erroneous entries for the soil analytical data

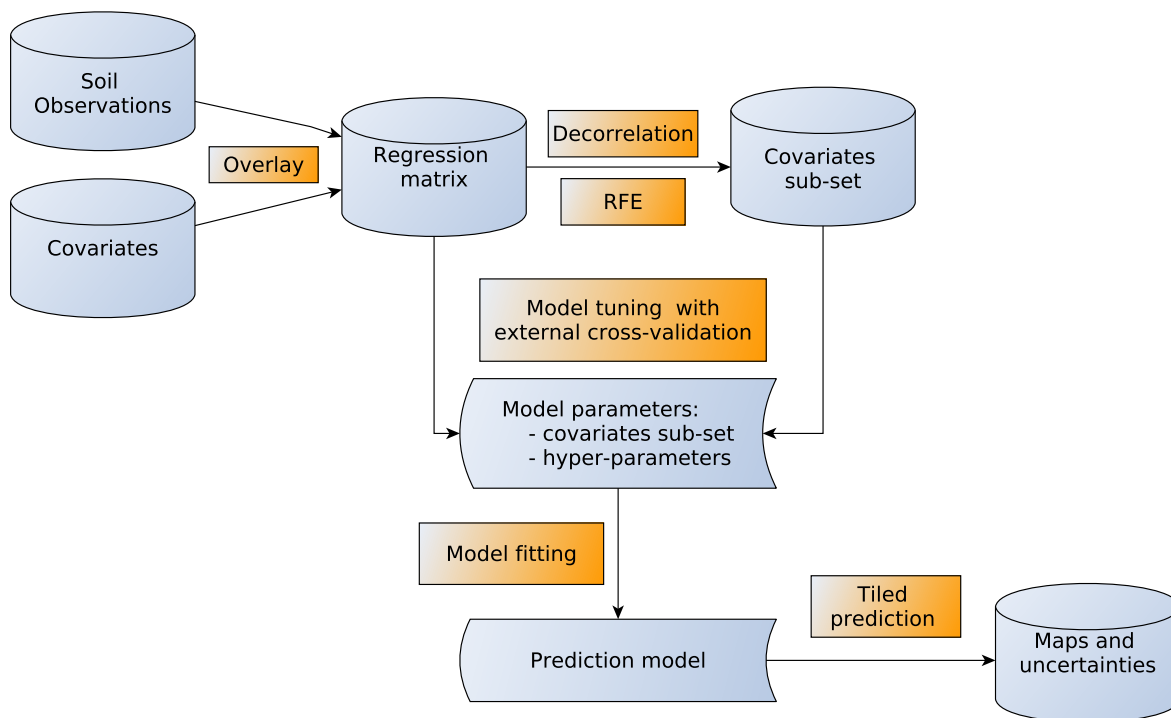


Figure 1. Workflow of the methodological approach.

[themselves](#) (Ribeiro et al., 2018). Ultimately, upon final consistency checks, the ~~quality-assessed-and~~ standardised data are made available via the ISRIC Soil Data Hub (<https://data.isric.org>) in accord with the licence specified by the data providers. As a result, not all data ~~standardized-standardised~~ in WoSIS are freely available to the international community. Hence, this study considers two 'sources' of point data.

First, the latest publicly available snapshot of WoSIS (Batjes et al., 2020). It contains, among others, data for chemical (organic carbon, total nitrogen, soil pH, cation exchange capacity) and physical properties (soil texture (sand, silt, and clay), coarse fragments). The snapshot comprises 196 498 geo-referenced profiles originating from 173 countries, representing over 832 000 soil layers (or horizons), in total over 5.8 million records. Generally, there are more observations for the superficial than the deeper layers. ~~Detailed information about the snapshot can be found in~~ [About 5% of the profiles were sampled before 1960, 14% between 1961-1980, 32% between 1981-2000, and 16% between 2001-2020; the date of sampling is unknown for 34% of the shared profiles](#) Batjes et al. (2020).

Second, in addition to the freely-shareable data, several soil [profile-observation](#) databases in our repository have licences stipulating that ISRIC may only use them for SoilGrids applications or visualizations, for example EU-LUCAS (Tóth et al., 2013) and soil data for the state of Victoria (Australia). The corresponding source data sets were screened and processed using the same procedures as used for the regular WoSIS workflow (some 42 000 profiles). As a result, some 240 000 profiles in total

Table 1. Soil properties description and units.

Soil property	Acronym	Units	Mapped Units	Description
Bulk density	BDOD	kg/dm ³	cg/cm ³	Bulk density of the fine earth fraction oven dry
Cation exchange capacity	CEC	cmol(c)/kg	mmol(c)/kg	Capacity of the fine earth fraction to hold exchangeable cations
Coarse fragments	CFVO	cm ³ /100cm ³ (volume %)	cm ³ /dm ³	Volumetric content of fragments larger than 2 mm in the whole soil
Nitrogen	N	g/kg	cg/kg	Sum of total nitrogen (ammonia, organic and reduced nitrogen) as measured by Kjeldahl digestion plus nitrate-nitrite
pH (water)	pH	-	*10	Negative common logarithm of the activity of hydronium ions (H ⁺) in water
Organic carbon concentration	SOC	g/kg	dg/kg	Gravimetric content of organic carbon in the fine earth fraction of the soil
Soil texture fractions	STF	%	g/kg	Gravimetric contents of sand, silt and clay in the fine earth fraction of the soil

were used as the data source for the present 2020 SoilGrids run, comprising more than 920 000 observed soil layers. During data processing some minor corrections were made to the merged input dataset, [for example further depth congruence checks](#).

2.1.1 Soil properties

105 For the purposes of SoilGrids, “soil” is up to 2m thick unconsolidated material at the Earth’s epidermis in direct contact with the atmosphere; thus sub-aqueous and tidally-exposed soils are not considered here, neither are materials deeper than 2m. This decision has consequences for computations of total stocks, in particular “soil” organic carbon.

Table 1 describes the soil properties that are considered in this version of SoilGrids: organic carbon content, total nitrogen content, soil pH (measured in water), cation exchange capacity, soil texture fractions, and proportion of coarse fragments.

110 These properties were modelled for the six standard depths intervals as defined in the GlobalSoilMap specifications (Arrouays et al., 2014): 0-5, 5-15, 15-30, 30-60, 60-100 and 100-200 cm.

‘Litter layers’ on top of mineral soils were excluded from further modelling using the following assumptions. Consistency in layer depth (e.g., sequential increase in the upper and lower depth reported for each layer down the profile) in WoSIS was checked using automated procedures. In accord with current internationally accepted conventions, such depth increments are given as “measured from the surface, including organic layers and mineral covers” (FAO, 2006; Schoeneberger et al., 2012). 115 Prior to 1993, however, the start (zero depth) of the profile was set at the top of the mineral surface (the solum proper), except

when “thick” organic layers as defined for peat soils (FAO-ISRIC, 1986) were present at the surface. Then the top of the peat layer was taken as the soil surface. Organic horizons were recorded as above and mineral horizons recorded as below, relative to the mineral surface (Schoeneberger et al., 2012) (pp. 2–6). Insofar as is possible, “superficial litter” on top of mineral layers was flagged as an auxiliary (Boolean) variable, also with reference to the original soil horizon designation when provided, so they can be filtered out during auxiliary computations of soil properties.

2.1.2 Transformation of texture data

A transformation was applied to the texture fractions, as follows. The relative percentage of sand, silt and clay can be treated as compositional variables, as the sum of the components always equals 100%. Therefore, these components were transformed using the Addictive Log-Ratio (ALR) transformation with the Gauss–Hermite quadrature (Aitchison, 1986). ALR has previously been applied to soil texture data (Lark and Bishop, 2007; Akpa et al., 2014; Ballabio et al., 2016; Poggio and Gimona, 2017a), and it has been shown (Lark and Bishop, 2007) that ALR-transformed variables preserve information on the spatial correlation and maintain the compositional integrity of the original components. In this study, clay was used as the denominator variable. Therefore the two ALR components that were interpolated can be defined as:

$$\begin{aligned} ALR1 &= \log\left(\frac{\text{sand}}{\text{clay}}\right) \\ ALR2 &= \log\left(\frac{\text{silt}}{\text{clay}}\right) \end{aligned} \quad (1)$$

2.1.3 Spatial stratification of observations

Random splitting of profile observations into n ~~validation~~ cross-validation folds is not suitable in this context, considering the high spatial variation in observation density as it would provide biased results (Brus, 2014). For regions like Europe and North America there are over 4 profiles per 10 km², whereas for large countries in Asia, such as Kazakhstan, India or Mongolia the number of available profiles is still quite limited (< 1 profile per 100 km²) (see Batjes et al. (2020) for further details).

Therefore, soil observations were spatially stratified in the geodetic domain to guarantee a balanced spatial distribution within each ~~validation~~ cross-validation fold. Spatial strata, in the form of hexagons, were created with an Icosahedral Snyder Equal-Area Grid (ISEAG) of aperture 3 and resolution 6, resulting in 7 292 strata ~~;(i.e., hexagonal cells),~~ each with an area around 70 000 km². This ISEAG was generated with the `dggridR` package for the R language (Barnes et al., 2016).

The profiles were assigned to one of ten folds, each equally represented in each stratum, i.e., each cell of the grid previously described. ~~The `caret` R package was used.~~ All observations (layers or horizons) belonging to a profile were always in the same fold for both model calibration and evaluation. The `caret` R package was used to subdivide the locations in the folds while maintaining the spatial distribution.

2.2 Environmental covariates

145 Over four hundred geographic layers were available as environmental covariates for this work. These were chosen for their presumed relation to the major soil forming factors, including long-term soil conditions, i.e., the “time” factor. Appendix A provides a list of the products used as covariates and their sources. The layers considered can be grouped as follow:

- **Climate:** temperature, precipitation, snowfall, cloud cover, solar radiation, wind speed.
- **Ecology:** bioclimatic zones and ecophysiological regions.
- 150 – **Geology:** soil and sedimentary thickness, rock types.
- **Land Use/Cover:** from sources such as the European Space Agency (ESA) and U.S. Geological Survey (USGS).
- **Elevation and terrain morphology:** including numerous morphology indexes and landform classes.
- **Vegetation Indexes:** such as the Normalized Difference Vegetation Index (NDVI), Enhanced Vegetation Index (EVI) and Net Primary Production (NPP).
- 155 – **Raw bands** from Landsat and Modis products.
- **Hydrography:** global water table, inundation and glaciers extent, and surface water change.

The ~~long-term~~ average and standard deviation of climatic variables and vegetation indices over 15 years (2001 - 2015) were computed from monthly data to capture their seasonal dynamics.

All covariates were projected to a common coordinate reference system (CRS), i.e., Goode’s Homolosine projection for land masses applied to the WGS84 *datum*. This projection was selected since among the equal-area projections supported by open source software it is the most effective minimising distortions over land (de Sousa et al., 2019). The projected covariates were imported to GRASS GIS in a normalised raster structure with cells of 250 m by 250 m. Covariates, and hence mapped areas, were restricted to land areas without built-up, water and glacier areas using a mask created from the ESA Land Cover layer for 2015 (Buchhorn et al., 2020). Thus properties of urban and subaqueous soils are not considered.

165 2.3 Covariates selection

Considering the large number of available environmental layers, a standardized and reproducible procedure to select covariates used for modelling was implemented to i) reduce redundancy between covariates, ii) obtain a more parsimonious and computationally-efficient model, iii) to decrease the risk of over-fitting (Gomes et al., 2019) and iv) to avoid a biased assessment of variable importance (Strobl et al., 2008). The covariates selection procedure consisted of two steps, de-correlation and Recursive Feature Elimination.

170

2.3.1 De-correlation analysis

De-correlation analysis was carried out as initial step to reduce the redundancy of information from more than 400 environmental layers. Only covariate layers that had a pairwise correlation coefficient ≤ 0.85 with all other covariates were included in the subsequent analyses. For each pair of covariates correlated above this threshold, only the first one in alphabetical order was selected for inclusion in the modelling phase. This step reduced the number of initial covariates to approximately 150 layers.

2.3.2 Recursive Feature Elimination

Recursive feature elimination (RFE) (Guyon et al., 2002) is a methodology that has proven effective to select an optimal set of covariates for regression trees models (Gomes et al., 2019; Hounkpatin et al., 2018). In this study, the RFE procedure implemented in the `caret` package for the R language (Kuhn, 2015) was used, as it offers a good compromise between accuracy and computation time. The algorithm starts by fitting a model using all covariates, assessing its performance and ranking covariate importance. The least important covariates are then removed from the pool, and again the model is fitted, assessed, and the least important covariates removed. The procedure repeats down to a pool between 0 and N covariates. This procedure is based on Out-Of-the-Bag (OOB) cross-validation and does not test all covariates combinations, but it is considered one of the most robust covariates selection approaches for models like Random Forests (Nussbaum et al., 2018).

The RFE procedure on the full set of observations and covariates would prove computationally prohibitive. To improve computational feasibility for large datasets, additional steps were developed. Four sets of observations were used for RFE, each obtained using three cross-validation folds (see Section 2.1.3 for further details): set1 contained folds 1 to 3, set2 folds 4 to 6, set3 folds 7 to 9 and set4 contained fold 10 and 2 other random selected folds. In a first step, the RFE procedure from `caret` was run independently on each set with default model hyper-parameters for [the Random Forests algorithm as implemented in the `ranger` package](#) (i.e., `ntree` as 500 and `mtry` as the rounded square root of the number variables). In each set the optimal number and combination of covariates was automatically selected when the model performances stopped increasing, i.e., when the loss function reached its minimum. In this study, the loss function was the OOB RMSE.

In the second step, the RFE procedure was applied with all observations and all covariates selected in at least one of the four sets used in the previous step. The final covariate set was the set minimising the loss function.

2.4 Hyper-parameter selection and cross-validation

Figure 2 summarises the approach used for the selection of the model hyper-parameters and the cross-validation. Further details are provided in the following sections.

2.4.1 Model tuning and [validation numeric evaluation](#)

Model tuning was performed with a 10-fold cross-validation procedure applied to multiple combinations of hyper-parameters.

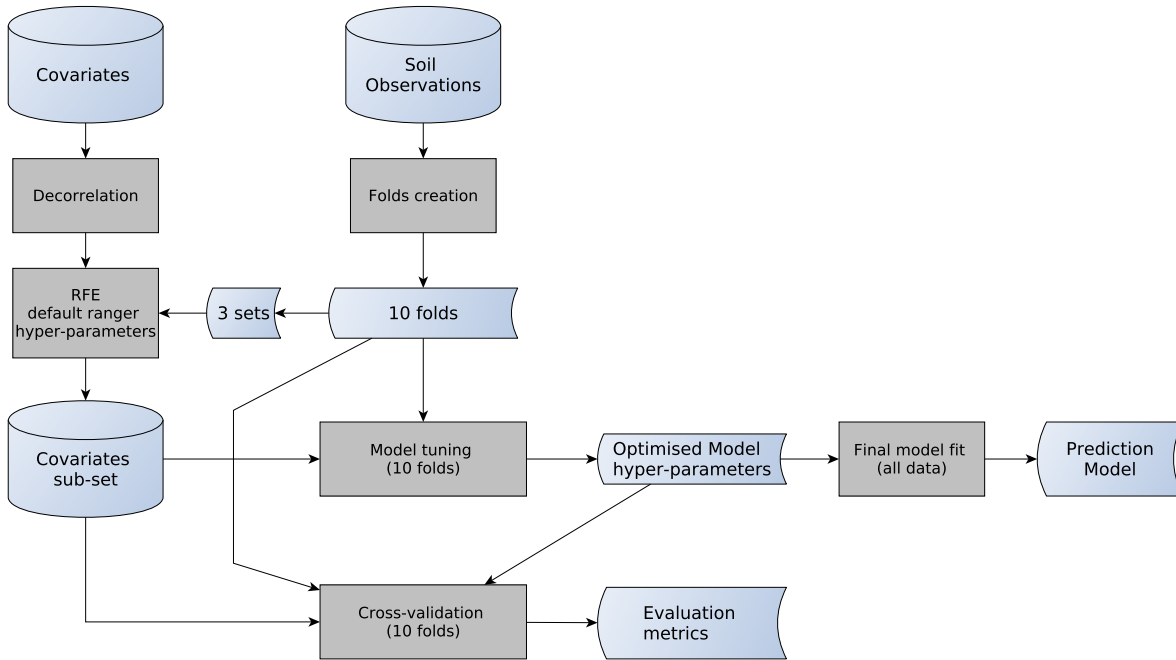


Figure 2. Detailed workflow for the Hyper-parameter selection and cross-validation.

Different numbers of decision trees (*n_{tree}* parameter) were combined with different numbers of covariates used in tree splits (*m_{try}* parameter). The number of trees were progressively increased with the following values: 100, 150, 200, 250 500, 750 and 1000. The different *m_{try}* values were multiples of the square root of the number of covariates. Four multipliers were tested, 1 (default in `ranger`), 1.5, 2 and 3. For example, if the RFE procedure identified a set of 50 covariates, the *m_{try}* values assessed were 7, 11, 14 and 21.

Each of the resulting combinations of *n_{tree}* and *m_{try}* parameters was used to train a different model with observations from nine folds. Predictions were then assessed on the remaining fold with classical performance measures, i.e., root mean squared error (RMSE) and model efficiency coefficient (MEC; Janssen and Heuberger, 1995). MEC is equal to the fraction of the explained variance based on the 1:1 line of predicted versus observed that is defined as 1 minus the ratio between residual sum of squares and total sum of squares. The final hyper-parameter selection was based on an optimisation of model performance and computational constraints, in this case memory consumption. For example an increase of the *n_{tree}* parameter above 200 provided a minor increment in the metrics (usually less than 0.1%, not reported here) while requiring considerably more memory and computation time.

The model evaluation was based on the performance metrics of the selected hyper-parameters combination. Predictions at the centre of the six standard depth intervals were compared with observations having the midpoint included within the considered interval.

2.5 Prediction and Uncertainty quantification

2.5.1 Model fit

The final model for each soil property was fitted with all available observations, the covariates and the hyperparameters selected in the previous steps. Observation depth was included in the model as a covariate. It was calculated at the mid-point of the sampled layer or horizon.

Models were obtained with the `ranger` package (Wright and Ziegler, 2017), with the option `quantreg` to build Quantile Random Forests (QRF; Meinshausen, 2006). With this option the prediction is not a single value, e.g., the average of predictions from the group of decision trees in the random forest, but rather a cumulative probability distribution of the soil property at each location and depth.

For each property (see Table 1) and standard depth from the GlobalSoilMap specification (0-5, 5-15, 15-30, 30-60, 60-100 and 100-200 cm) four different values were computed to characterise this distribution: median (0.50 quantile, $q_{0.50}$), mean, 0.05 quantile ($q_{0.05}$), and 0.95 quantile ($q_{0.95}$), i.e., the lower and upper limits of a 90% prediction interval. This uncertainty interval is as described in the GlobalSoilMap specifications (Arrouays et al., 2014). The predictions were computed for the mid-point of the depth interval and considered constant for the whole depth interval.

In order to compute the prediction uncertainty for soil texture, the back-transformation was applied at the level of individual tree predictions, and the quantiles of the tree prediction distributions obtained from the resulting values.

2.5.2 Uncertainty

The percentage of ~~validation~~-cross-validation observations contained in the 0.9 prediction interval was calculated (Prediction Interval Coverage probability, PICP) (Shrestha and Solomatine, 2006). Ideally the PICP is close to 0.9, indicating that the uncertainty was correctly assessed. A PICP substantially greater than 0.9 suggests that the uncertainty was underestimated, a substantially smaller PICP indicates that it was overestimated.

Furthermore, to visualise the uncertainty as a map, the following indicators were calculated:

1. 90th prediction interval (PI90)

$$PI90 = q_{0.95} - q_{0.05} \quad (2)$$

2. ratio of the interquartile range over the median (Prediction Interval Ratio, PIR):

$$PIR = \frac{q_{0.95} - q_{0.05}}{q_{0.50}} \quad (3)$$

2.6 Qualitative evaluation of spatial patterns

Expert judgement was used to evaluate the reasonableness of the maps, by comparing well-known spatial patterns at global, regional, and local scales with SoilGrids predictions (see subsection 3.4). Obviously these are not definitive evaluations, only indicative.

2.7 Software and computational framework

SoilGrids requires an intensive computational workflow, with numerous steps integrating different software. SoilGrids is entirely based on open source software, in particular: SLURM (Yoo et al., 2003) for job management, GRASS GIS (GRASS Development Team, 2020) for data and tiles management, and R statistical software (R Core Team, 2020) for model fitting and statistical analysis.

Predictions were computed in a high-performance computing cluster. A dynamic geographic tiling system was developed with GRASS GIS to maximise the use of memory for each job. Technical details on this parallelisation scheme are given in de Sousa et al. (2020).

The predictions were multiplied by a conversion factor of 10 or 100 to maintain the required precision while using integer type in the file geotiff to reduce space occupied on disk. Application of the conversion factor resulted in mapped layers with units differing from those of the input observations (see Table 1).

The total computation time with the selected covariates and hyper-parameters differed per property. On average, the complete computation of the 24 maps (mean and three quantiles for each of six standard depths) for a single property, including: (i) RFE, (ii) model training, and (iii) prediction, took approximately 1 500 CPU-hours. The prediction accounted for about two thirds of the total time.

3 Results and discussion

3.1 Input soil observations

Table 2 breaks down the distribution of the legacy soil observations for each soil property by depth interval. Table B1, in Appendix B, shows the number of observations by bio-climatic region.

Figures 3 and 4 show examples of observation density of the soil calibration data for two soil properties, pH_{water} and proportion of coarse fragments, that show a large difference in density.

As indicated, the number of observations for each property varies greatly with depth and bioclimatic region, with higher densities observed for North America and Europe (Batjes et al., 2020). Generally, there are more observations for agricultural areas. Further, the available profiles have been collated over several decades, some 62% of the data being from 1960-2020; the time of sampling is unknown for around 34% of the profiles. As indicated by Batjes et al. (2020), in principle, the age of the observations should be taken into account during the mapping process via covariate layers for time periods commensurate with the sampling dates, especially for soil properties that are readily affected by changes in land use or management practices. However, for these so-called 'dynamic' soil properties, such as pH and soil organic matter content, we consider that the spatial variation will be much greater than the temporal variation, so that not taking into account the age of observations will not greatly affect the map. In addition, it is difficult or impossible to find comparable covariates, in particular remote-sensing derived covariates, for each time period. Space-time relations should be considered in future assessments (Heuvelink et al., 2020).

Table 2. Number of observations per standard depth interval for each soil property. See Table 1 for abbreviations [and units of the considered soil properties](#).

Depth interval	BDOD	CEC	CFVO	N	pH	SOC	STF
0 - 5 cm	8122	20576	15541	27192	44049	48616	42983
5 - 15 cm	19817	49463	66833	82856	146677	148918	155302
15 - 30 cm	17819	40673	35254	39568	91326	91682	98659
30 - 60 cm	27146	63444	56755	48804	141812	122338	140353
60 - 100 cm	23130	58038	50912	36946	131172	102687	127073
100 - 200 cm	23396	66236	49995	28135	129373	92327	116847

This study considers standardised data for some 240 000 profiles, derived from WoSIS. This is over 60 000 more profiles than considered in the data compilation underpinning the preceding SoilGrids runs (Hengl et al., 2017b), thus providing substantial new information for calibration of the new global models. However, as indicated, there are still significant geographic gaps (e.g., arid regions, boreal regions, and 'forest' soils). Some of these are related to the physical remoteness or unaccessibility of some regions, while others are related to the fact that many soil datasets still are not or can not be shared for various reasons as described by Arrouays et al. (2017).

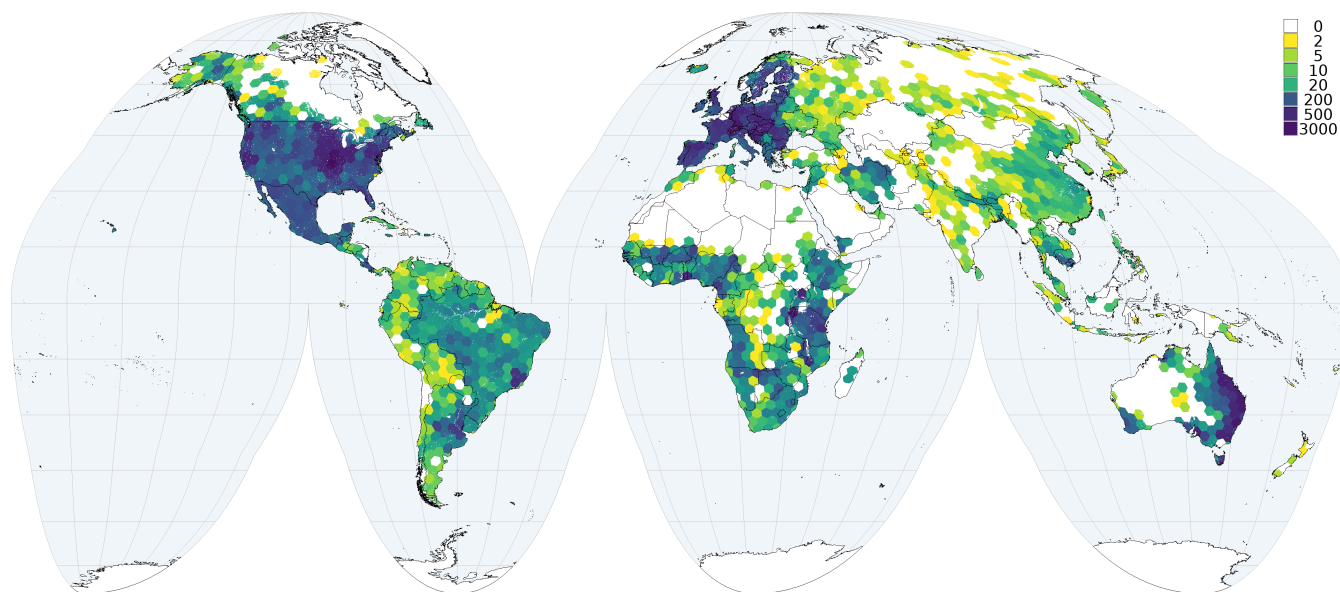


Figure 3. Number of observations per grid cell ($70\,000\text{ km}^2$) for soil pH_{water} .

In the previous version of SoilGrids (Hengl et al., 2017b), synthetic observations were randomly placed in regions with few or no observations, e.g., the Sahara and the Arabian Peninsula. This approach is worth further exploring, including information

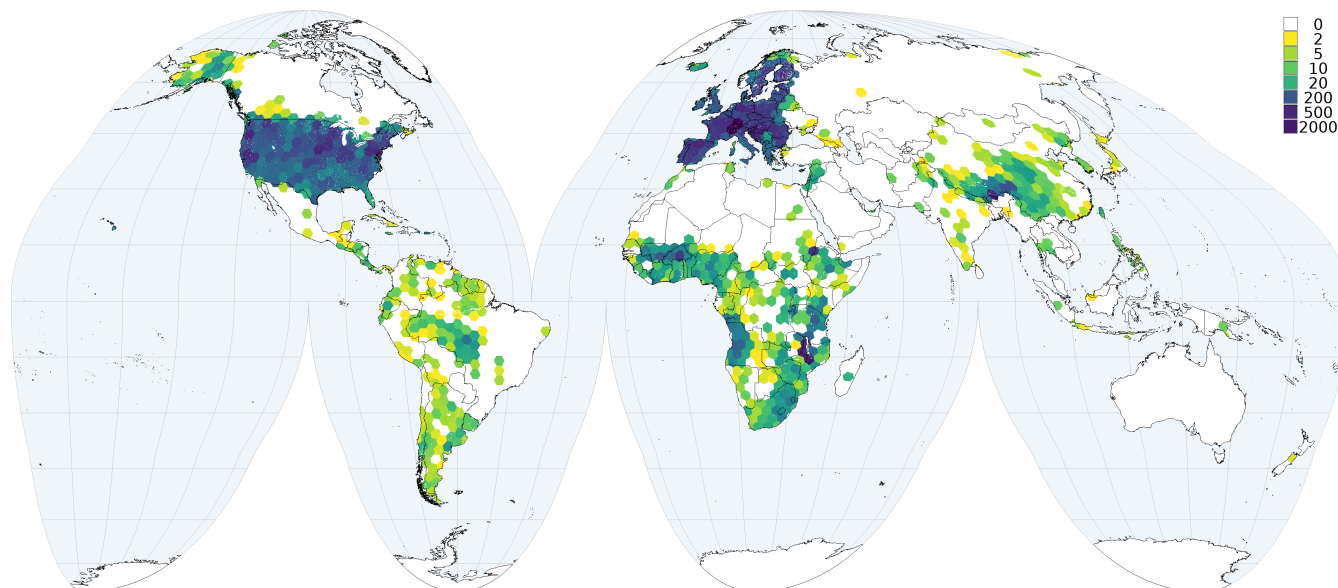


Figure 4. Number of of observations per grid cell (70 000 km^2) for coarse fragments.

derived from other regional datasets, expert opinion and by transfer learning from similar areas according to the *Homosoil* concept (Mallavan et al., 2010), which assumes similarity of soil-forming factors across regions. However, SoilGrids already implicitly incorporates the *Homosoil* concept, as long as there are sufficient observations in a given soil-forming environment anywhere in the world. Therefore, no synthetic observations ('pseudo-points') were included in this version of SoilGrids, also
 290 by a lack of confidence about the accuracy of the synthetic data.

In future studies, it will be relevant to identify beforehand areas of the world with a low observation density that are not yet represented by a high density of observations in other areas with similar soil-forming factors. A set of synthetic profiles could then be generated to describe these areas, by consulting soil scientists knowledgeable on the soils and soil properties of these areas.

295 **3.2 Model tuning and hyper-parameters selection**

Model hyper-parameters selected for each property are presented in Table 3.

The numbers of covariates selected using the two-step approach for covariates selection was fairly small in comparison with the full set (Table 3), resulting in more parsimonious models. Figure 6 shows two examples of the loss-function for RFE for two soil properties with different number and distributions of input observations. In both cases there is a clear improvement of
 300 performances while using 15 to 20 covariates. The curve reaches a minimum of the loss function and then stays on a plateau with a slight decline after the identified minimum.

All final models were trained with a maximum of 200 decision trees, a number beyond which performance gains did not noticeably increase.

The *mtry* parameter mainly depended on the number of covariates and was always between 1.5 and 2 times the square root of the number of covariates, which is the default provided by common random forest packages such as *ranger* (Wright and Ziegler, 2017). This confirms the need to determine optimum model hyper-parameters, especially when dealing with large numbers of input data (Nussbaum et al., 2018) as is the case here.

Table 3. Hyper-parameters for each considered soil property. See Table 1 for abbreviations [and units of the considered soil properties](#).

	Number of covariates	Number of trees	<i>mtry</i>
BDOD	40	200	12
CEC	25	200	10
CFVO	20	200	6
N	30	200	10
pH	32	200	9
SOC	40	200	12
Texture ALR I	25	150	10
Texture ALR II	27	150	10

3.3 Quantitative evaluation

Table 4. Global cross-validation results for both mean and median predictions. See Table 1 for abbreviations [and units of the considered soil properties](#).

Property	RMSE (median)	RMSE (mean)	MEC (median)	MEC (mean)
BDOD	0.19	0.19	0.73	0.74
CEC	11.01	10.69	0.40	0.43
CFVO	13.46	12.69	0.22	0.31
N	2.62	2.50	0.47	0.52
pH	0.78	0.77	0.67	0.68
SOC	39.67	36.48	0.37	0.47
Sand	0.19	0.18	0.51	0.54
Silt	0.13	0.13	0.60	0.62
Clay	0.13	0.13	0.42	0.43

Cross-validation results are summarised in Table 4, presenting the root mean squared error (RMSE) and Model Efficiency Coefficient (MEC). The MEC varies from a minimum of 0.31 for coarse fragments to a maximum of 0.74 for BDOD. Clay is less well modelled than the other two particle-size classes. This may be an effect of the chosen ALR transformation that had

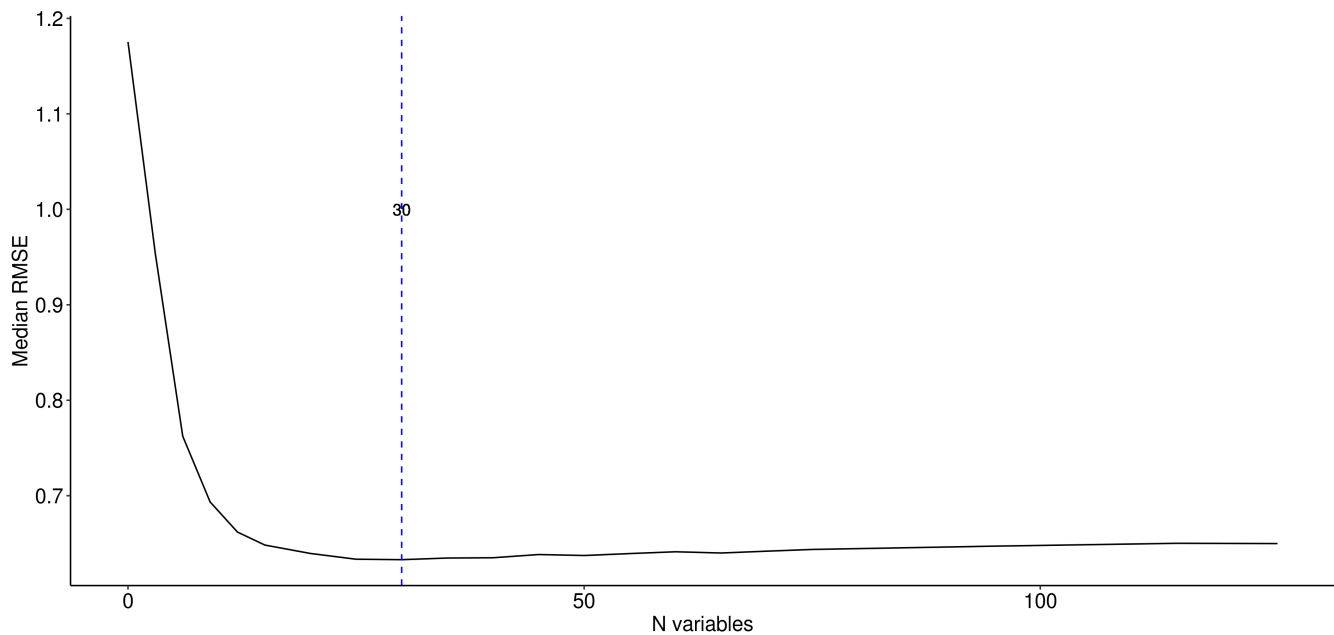


Figure 5. pH

—

Figure 6. ~~CFVO~~ Example of loss function (RMSE) used in the RFE step of covariates selection.

~~Loss function for RFE:~~

Table 5. MEC per depth layer for mean predictions. See Table 1 for abbreviations and units of the considered soil properties.

Depth Layer	BDOD	CEC	CFVO	N	pH	SOC	Sand	Silt	Clay
0-5cm	0.78	0.46	0.33	0.65	0.69	0.55	0.59	0.71	0.45
5-15cm	0.74	0.42	0.35	0.41	0.66	0.39	0.58	0.64	0.42
15-30cm	0.72	0.39	0.33	0.44	0.68	0.38	0.57	0.68	0.42
30-60cm	0.70	0.42	0.31	0.46	0.68	0.38	0.54	0.62	0.41
60-100cm	0.61	0.41	0.29	0.48	0.68	0.42	0.50	0.57	0.40
100-200cm	0.59	0.45	0.29	0.49	0.67	0.59	0.48	0.54	0.40

clay as denominator (Lark and Bishop, 2007). Metrics of the mean were always better than or equal to those for the median for all properties.

Overall, these metrics are in line with continental or large regions DSM studies (Keskin and Grunwald, 2018). However, they are slightly lower than those ~~preseneted~~ presented by Hengl et al. (2017b). The latter difference can be explained by the

315

Table 6. Prediction Interval Coverage Probability, global and by predicted depth interval. See Table 1 for abbreviations [and units of the considered soil properties](#).

Property	Global	[0, 5]	[5, 15]	[15, 30]	[30, 60]	[60, 100]	[100, 200]
BDOD	0.90	0.89	0.91	0.91	0.91	0.90	0.88
CEC	0.88	0.89	0.90	0.88	0.88	0.88	0.87
CFVO	0.95	0.96	0.95	0.95	0.95	0.94	0.94
N	0.92	0.91	0.92	0.93	0.92	0.92	0.92
pH	0.90	0.91	0.91	0.90	0.91	0.90	0.89
SOC	0.92	0.91	0.92	0.92	0.92	0.92	0.92
Sand	0.79	0.82	0.82	0.80	0.78	0.78	0.78
Silt	0.96	0.95	0.96	0.96	0.96	0.96	0.96
Clay	0.96	0.96	0.96	0.96	0.95	0.95	0.96

more prudent cross-validation approach now taken, with spatially balanced folds and all observations belonging to the same profile in the same fold. This prevents using data from the same profile both for calibration and [validation numerical evaluation](#).

Table 4 shows that the models with a higher number of [covariates-retained covariates \(Table 3\)](#) have better predictive performances. However, these models are also the models with the largest number of observations (Table 2). [The considered soil properties are also different. Therefore, no general conclusion can be drawn from this observation.](#)

Table 5 shows the MEC for mean predictions by depth interval. Performances decreased with depth, in line with many other DSM studies (Keskin and Grunwald, 2018). This pattern can be explained [in part by the reduced number of observations for deeper layers, but also by mainly by](#) weakened relationships between environmental layers and soil properties of the deeper [horizons layers](#).

In this study, the vertical dimension of soil variability was only taken into account by using the depth of the observation as a covariate. [Recent publications \(Ma et al., 2021; Nauman and Duniway, 2019\) indicate that such approach can be too simplistic or lead to problems with consistency over the predicted depth sequence. This may be true for local datasets where the short-range spatial variability is of a similar magnitude as the vertical variability. Further research is necessary to assess the effects of using depth as a covariate on global datasets and models.](#) Alternatives such as 3D smoothers (Poggio and Gimona, 2017b) or geostatistical models exploiting 3D spatial auto-correlation are worth exploring in further studies.

Table 6 summarises the PICPs, globally and by predicted depth intervals. Most of the values are between 0.88 and 0.92, indicating that the predictions intervals obtained with QRF are a realistic representation of the prediction uncertainty, as the expected value for a 90% prediction interval is 0.90. Exceptions are the models for coarse fragments with higher values around 0.95, indicating an overestimation of prediction uncertainty. The texture components have values with a larger spread, around 0.78 to 0.80 for sand and closer to 0.96 for silt and clay. These indicate a potential under-estimation of prediction intervals for sand and over-estimation for silt and clay. These results may be related with the range of these properties in the input

observations. The transformation method used to derive the prediction intervals for the texture components could also be a contributing factor. Further exploration of the causes is worthwhile.

A key ~~problem-issue~~ for DSM applications using legacy soil data is the evaluation (~~“validation”~~) ~~of the results~~ of the results.
340 The aspect of evaluation which compares actual with predicted values is numerical (or statistical) evaluation, often termed "validation" in the DSM literature; Oreskes (1998) and Rossiter (2017) explain why the term "evaluation" is preferred for the overall process of assessing the success of models, including DSM models. The best approach to numerical evaluation is to have an independent dataset obtained with probability sampling using a known sampling design (Brus et al., 2011; Brus, 2014). However, this is not feasible when only legacy data are available. In this case, a so-called "cross-validation" approach is often
345 used. ~~Cross-validation~~ This needs to be tuned to avoid over- or under-estimation of the numeric evaluation metrics, especially in case of large differences in observation density, i.e., clustered spatial observations. This is especially important at global scale, as the distribution of the soil observations is not uniform across the globe. It can not be guaranteed that the numeric evaluation metrics derived from cross-validation are unbiased estimates of the ~~true validation~~ true numeric evaluation metrics, i.e., those that would have been computed on a probability sample of the whole population. It is also not possible to quantify
350 how close the cross-validation metrics estimates are to the true ~~evaluation~~ metrics, as it is not possible to obtain confidence intervals (Brus et al., 2011). When using cross-validation it is important to prevent over- or under-optimistic estimates. For example, it is likely that prediction errors are smaller in areas where the sampling density is higher. Because of their high sampling density, such areas will be over-represented in the sample as the percentage of cross-validation points in clustered areas will be higher than the percentage of the total land area covered by those areas. Results of standard cross-validation will
355 be strongly influenced by the performances in clustered areas. Using spatial cross-validation as suggested by Meyer et al. (e.g. 2018), where it is ensured that calibration data are never too close to a ~~validation~~ cross-validation point, on the other hand could produce over-pessimistic results. In order to address some of these concerns, this study ~~followed~~ adopted a practical solution ~~where in which~~ the folds were created to guarantee a ~~spatially-balanceed distribution between validation folds, i.e., maintain spatially-balanced distribution between cross-validation folds, maintaining~~ the same densities of the input data in each fold so
360 that they represent approximately the same population.

Although the numerical evaluation procedure used in this work takes into account the spatial distribution of the observations and their density, further improvement is ~~necessary~~ possible in both model training and evaluation. For example, the ~~weight assigned to~~ weights assigned to observations in heavily sampled areas could be reduced. The USA and large regions of Europe and Australia have ~~very~~ high numbers of observations that could be ~~reduced to further~~ down-weighted to strengthen the spatial
365 robustness of the ~~validation~~ evaluation procedure. Declustering or debiasing techniques (Goovaerts, 1997; Deutsch and Journel, 1998) have been applied with success in other geo-statistics exercises and could be adapted to ~~the particular case of~~ global soil mapping. The creation of the folds could also be modified to take into account the density of the observations.

3.4 Qualitative evaluation of spatial patterns

At global scale well-known patterns are reproduced, and typical properties associated with many World Reference Base for
370 Soil Resources (WRB) (IUSS Working Group WRB, 2015) Reference Soil Groups can be recognised.

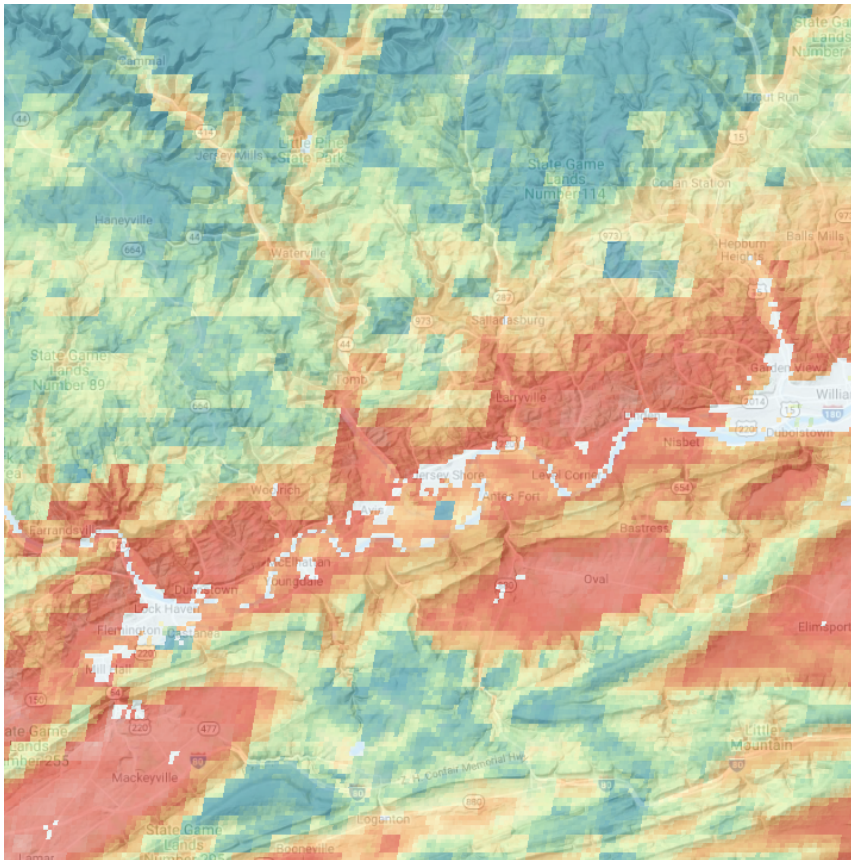
For example, the pH map identifies the large regions of alkaline soils (Solonetz, Solochak), highly-weathered soils (e.g., Acrisols, Alisols, Plithosols), acid forest soils (e.g., Podzols), and young soils from calcareous glacial deposits (e.g., Luvisols). The low pH of Andosols (e.g., Pacific Northwest USA, Japan, New Zealand) is also correctly represented. The texture components (PSC) maps correctly identify the siltier deltas (e.g., Yellow/Yangtze, Ganges/Brahmaputra), broad river plains (e.g., Po, Danube, Mississippi, Rio Plate, upper Amazon), the loess regions (e.g., midwestern USA, NW Europe, Ukraine), and the sandy North German/Polish plain. The cation exchange capacity (CEC) map clearly identifies large regions of highly-weathered clays (e.g., southeastern USA and China, central Brasil) and high-CEC 2:1 clays (e.g., “black cotton” Vertisols in the Deccan plateau and the Sudan). This map, together with the soil organic carbon (SOC) concentration maps, identify large regions of Histosols (e.g., northern Canada, Scotland, Siberia, Borneo). The SOC stock map identifies deep Histosols and cool, wet regions (e.g., Pacific Northwest North American coast, Ireland, southern Chile). The coarse fragment map identifies large areas of the Tibetan plateau and the principal mountain chains, as well as recently-glaciated soils on igneous bedrocks (e.g., Scandinavia, northern Quebec and Ontario).

Many regional patterns are also clear, for example the pH transition from Sahara through Sahel to the West Africa coast, and the PSC transitions from the Des Moines glacial lobe to the pro-glacial loess deposits in Iowa (USA) as well as the PSC transition from clayey marine sediments along the North/Baltic sea coasts through the sandy plains to the central German loess belt. The CEC map identifies contrasting areas of Vertisols (e.g., “black belts” in Alabama/Mississippi and Texas USA). The coarse fragments map shows the detailed pattern of the basin-and-range region of the western USA and the ridge-and-valley region of Appalachia.

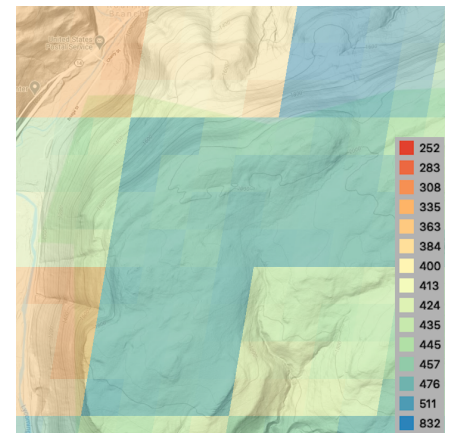
However, at the local scale, a preliminary assessment of SoilGrids in the USA, comparing with a gridded version of the national detailed gSSURGO (NRCS National Soil Survey Center, 2016) soil geographic database based on detailed field survey, reveals that SoilGrids may fail to account for local parent material transitions, e.g., sedimentary facies of coastal plain marine sediments, as well as glacial features such as proglacial lacustrine sediments and relic beach lines, so that the local PSC pattern is not accurate, sometimes on the order of 20-30% of a particle-size class.

For example, Figure 7 (left) shows predicted sand concentration of the 0-5 cm layer in an approximately 50 x 50 km area in central Pennsylvania (USA), ranging from approximately 25% (~~darker~~ darkest red) to ~~60%~~ (~~darker~~ 80% (darkest blue)). Important local differences are clear: the low sand concentrations of the clayey soils in the limestone valleys trending SW-NE and the high concentrations in soils from glacial till developed on sandstones in the north, as well as the residual soils on the resistant sandstone ridges of the Ridge and Valley province of Appalachia in the south. These do generally agree with the detailed soil survey.

At the detailed scale (250 m pixel) SoilGrids typically shows fine details that do not always appear to be related to obvious landscape or land use differences, when the map is viewed as a ground overlay in Google Earth. For example, Figure 7 (right) shows detail of predicted sand concentration of the 0-5 cm layer in an approximately 3 x 3 km area of the previous figure. The effect of some covariates being at 1 km resolution and others at 250 m is apparent, but the reason for the fine-scale differences is not. This area is of similar lithology, relief and land cover (second-growth dense forest) except the narrow valley at the northwest edge, yet the predictions are quite different.



(a) Overview; centre $\approx -77^{\circ} 14' E, 41^{\circ} 14' N$, near Jersey Shore, PA



(b) Detail; centre $\approx -76^{\circ} 56' E, 41^{\circ} 33' N$;
southeast of Roaring Branch, PA

Figure 7. Predicted sand concentration, %, 0-5 cm, ground overlay in © Google Earth

In this context, it should be realised that SoilGrids250m predictions are not meant for use at a detailed scale, i.e., at the sub-national or local level, as national data providers often have access to more detailed point datasets and covariate layers for their country than [have been provided to the point dataset on which SoilGrids250m can consider](#) (Chen et al., 2020; Roudier et al., 2020; Vitharana et al., 2020). [is based](#) (Chen et al., 2020; Roudier et al., 2020; Vitharana et al., 2019; Liu et al., 2020).

410 3.5 Prediction uncertainty

In general, the least sampled areas present the highest prediction uncertainties as expressed by the PICP. Figures 8 and 9 show an example for two properties and depths (Maps for all properties and depths can be accessed at <https://data.isric.org>). Figure 10 shows an example representing the quantiles for pH_{water} for the 60-100 cm layer. The north of Russia and the centre and north-west of Canada are large regions for which few soil observations are available, therefore prediction distributions are wider than in more densely sampled areas. However, these patterns are different for different properties. For example, arid areas actually have the narrowest prediction ranges of pH_{water} . The uncertainty range is often wide for properties and regions with a wider range of the property being modelled. This can be explained by the modelling approach performing more accurately within a limited range of options. These regions also have larger local spatial variation with more difficulties for predictions.

The communication of uncertainty is an open challenge (Arrouays et al., 2020). Uncertainty should provide information for policy makers and other stakeholders and not only scientists and modellers. The maps computed with eq. 3 are a first step in this direction, but their limitations must be understood. For properties that have values at or near zero, e.g., coarse fragments, they do not provide an entirely accurate uncertainty estimate. The use of uncertainty classes could be a further step to help domain stakeholders.

3.6 Limitations and outlook

425 4 Conclusions and future work

~~The aim of this work was the production of quality-assessed global maps of soil properties, with cross-validation, hyper-parameters selection and quantification of uncertainty as implemented in the SoilGrids 2.0 product.~~ This study represents a considerable effort to provide a globally-consistent product using the point dataset available to ISRIC, a large number of relevant covariates, and some optimisation of a well-established machine learning method, within the limits of practical computation. Yet it is clear that this product has some limitations, which will be considered in further work.

~~There are constant improvements in the DSM process, including the development and implementation of new methods and the use of new~~ First, there is an ever-expanding group of new covariates that can help explain and model the spatial variation of soil properties. Products derived from Earth Observation are particularly relevant in this regard and have considerably improved over the last decade. For example, the European Space Agency Sentinel missions (both optical and radar) provide high-resolution data that have been shown to improve DSM models performances.

Second, a fundamental problem is lack of well distributed point observations within the soil property geographic and features space. Additional soil data for so far under-represented regions, for example the northern boreal regions as being collated by the International Soil Carbon Network (Malhotra et al., 2019), will be sought for possible consideration in the WoSIS workflow that provides the point data underpinning the SoilGrids mapping effort. This effort would be aided by the provision by more data providers of at least a representative part of their point data to WoSIS, under suitable license. It is also important to consider the distribution of the observations in the covariates space to minimise the issues related to predictions into unknown regions of feature space (Meyer and Pebesma, 2020).

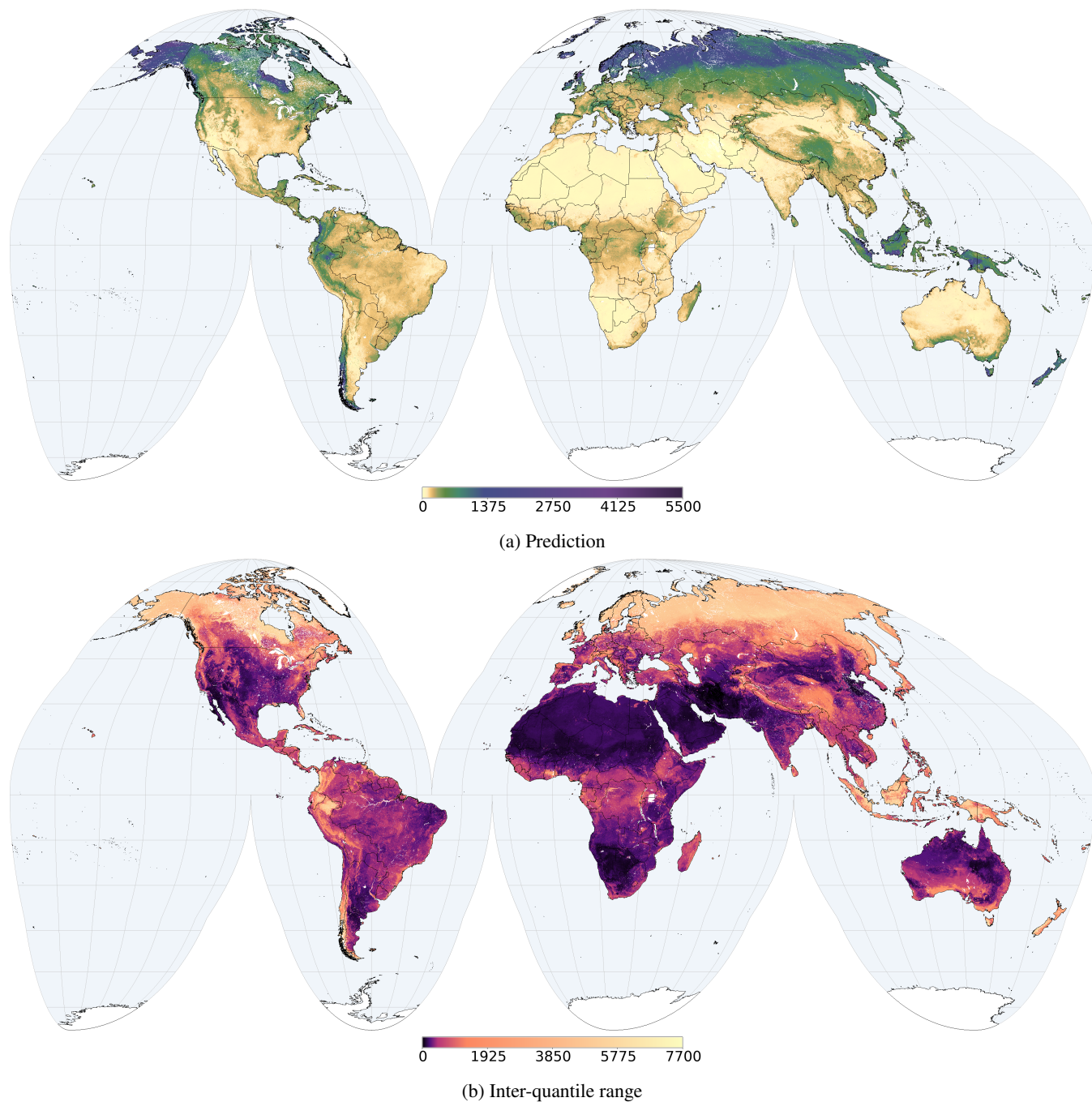


Figure 8. Mean soil organic carbon content (dg/kg) prediction and range between 5% and 95% quantiles in the 5 cm to 15 cm depth interval.

[Third, DSM methods are under active development, both new methods and improvements to established methods.](#) The use of decision tree-based models in DSM has become fairly common in recent years. Models such as Random Forests, XGBoost or

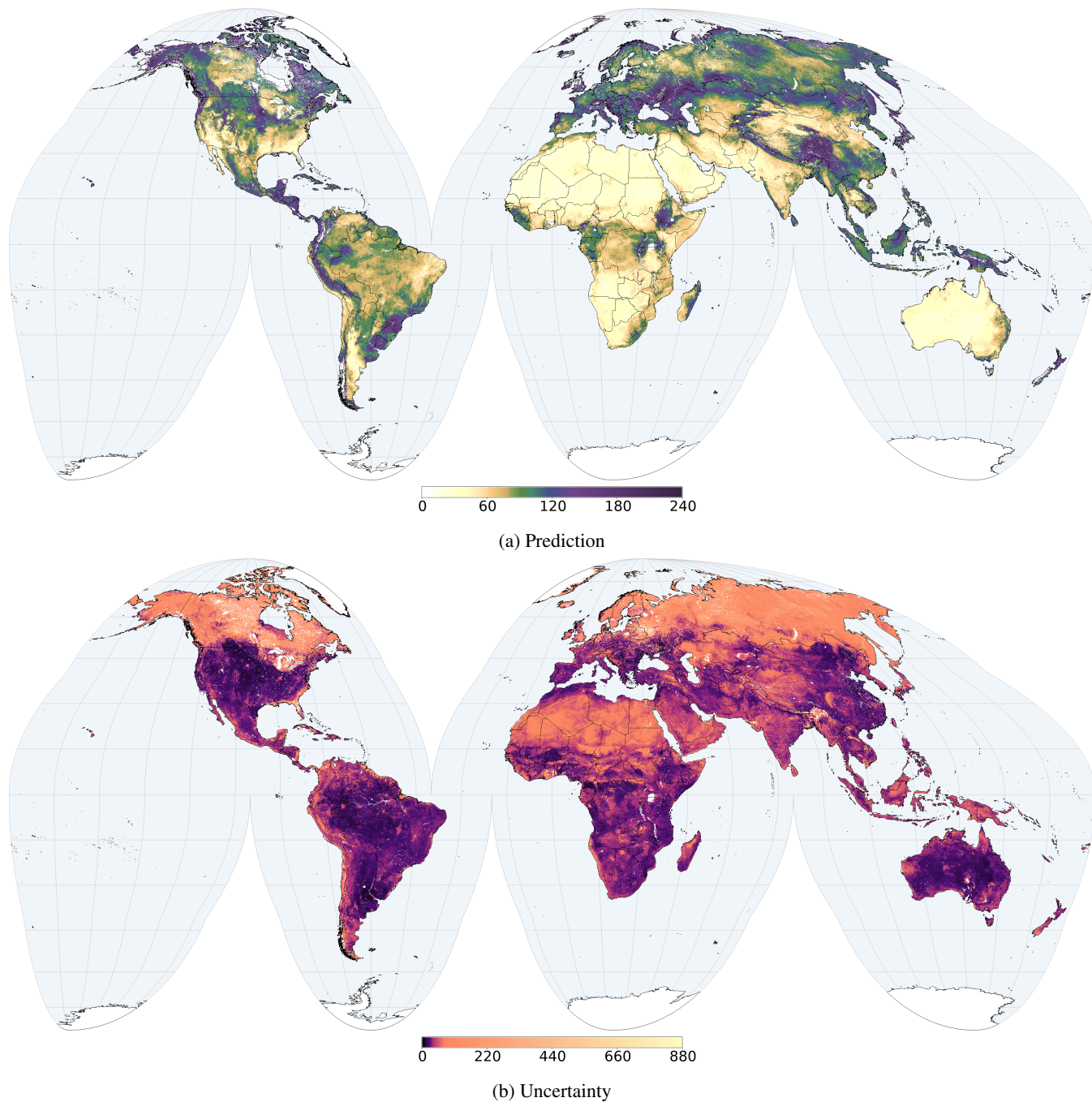


Figure 9. Median total Nitrogen prediction (cg/kg) and associated uncertainty for the 15 cm to 30 cm depth interval.

445 Cubist tend to provide better results than most multiple linear regression methods with reasonable computation costs (Khaledian and Miller, 2020). However, methods such as Artificial Neural Networks promise further improvements in model performances

if the amount and distribution of the data support these highly-complex models. This is the case in particular with convolutional or recursive neural networks (*deep learning*). However, these methods present computational challenges with the amount of

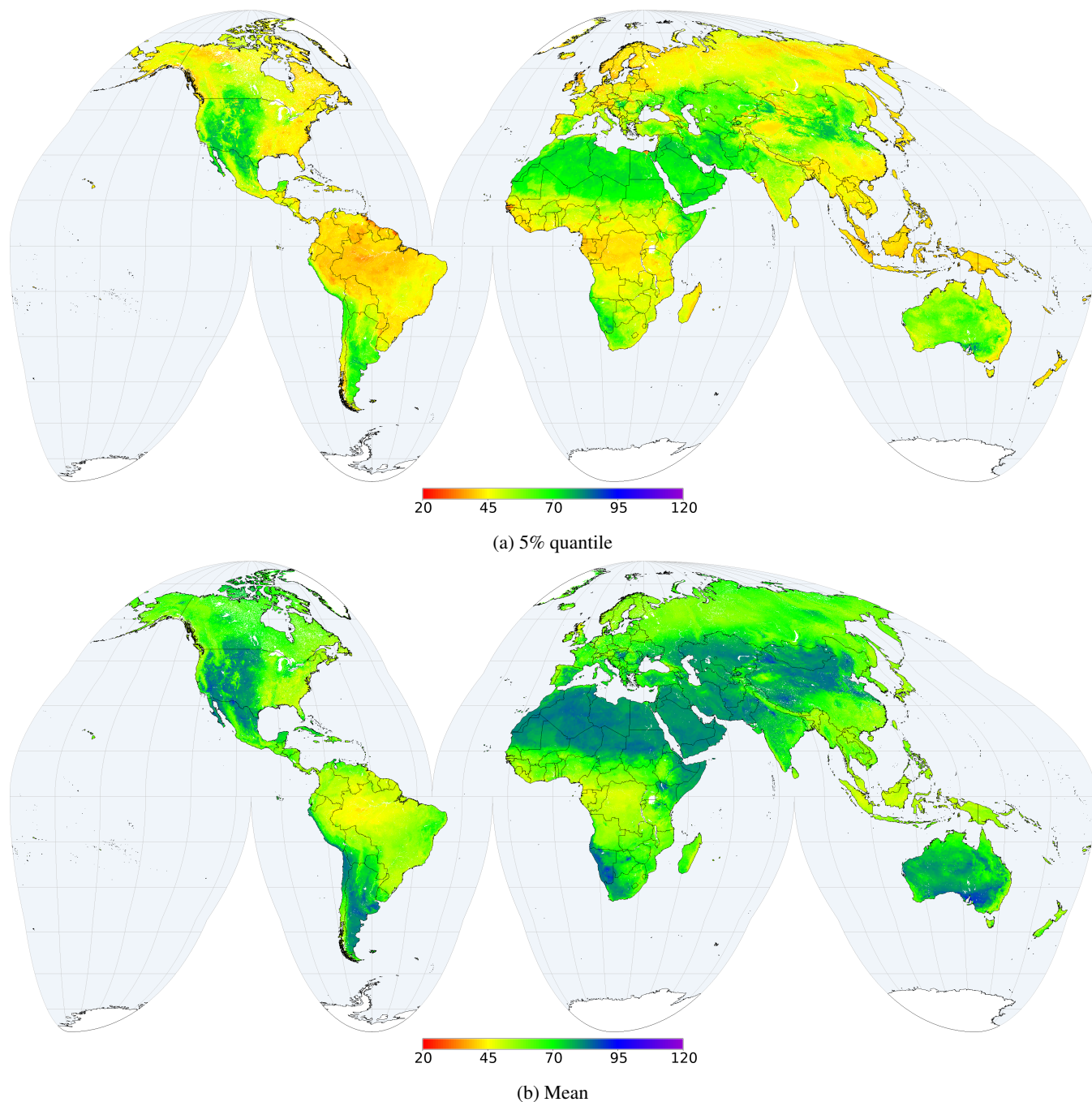
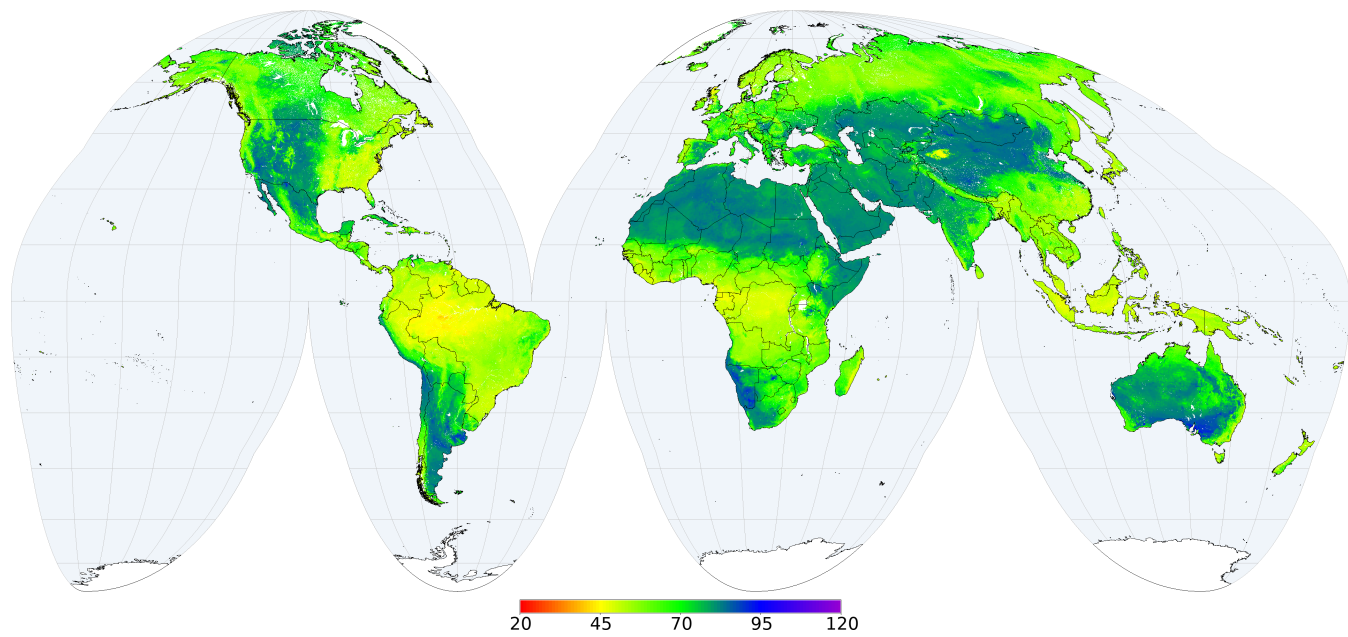
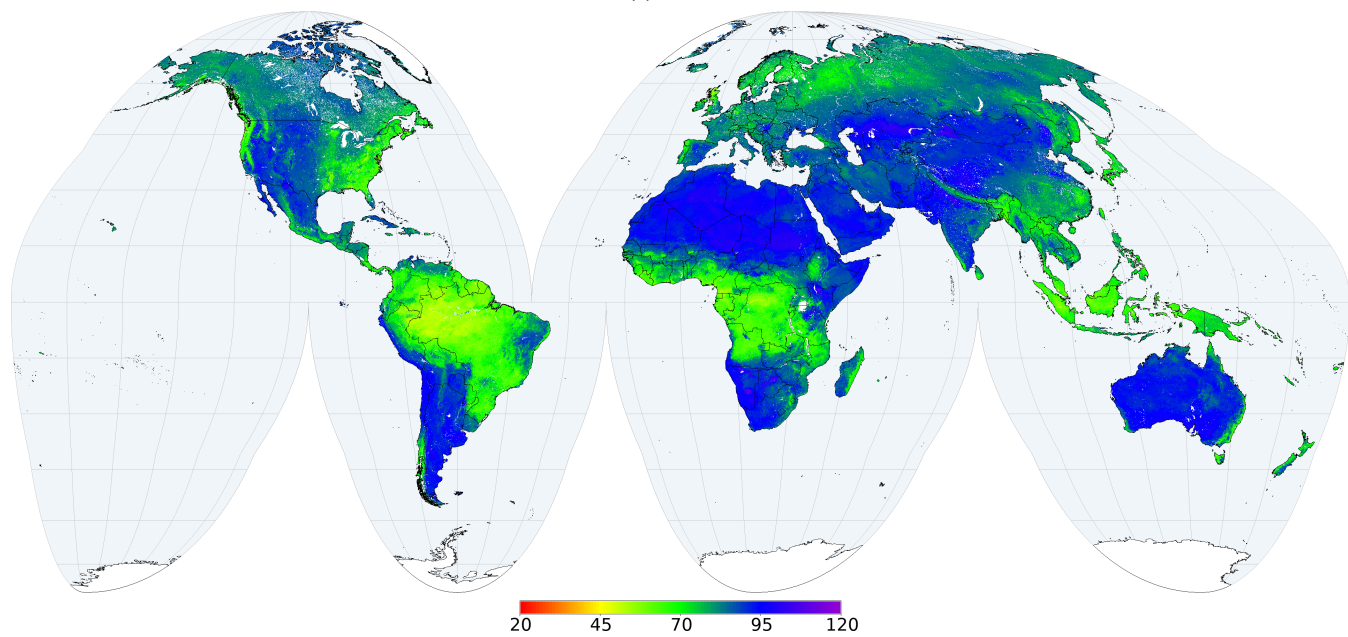


Figure 10. Prediction distribution for pH_{water} ($10 \cdot \text{pH}$) in the 60 cm to 100 cm depth interval.



(c) Median



(d) 95% quantile

Figure 10. Prediction distribution for pH_{water} ($10 \cdot \text{pH}$) in the 60 cm to 100 cm depth interval (cont.).

training data necessary for a sufficiently accurate DSM exercise, especially when working at global scale at medium to fine resolutions.

~~Cross-validation~~ Fourth, the proper method of cross-validation is another important aspect when considering how to assess and improve model performances. In particular, spatial cross-validation and declustering of the data need to be further explored.

~~This work described~~ Fifth, this research considered only the modelling of some ~~of the~~ primary soil properties, as defined and described in the GlobalSoilMap specifications. More work is necessary to obtain maps for soil thickness (either rooting zone, pedogenetic solum or regolith), soil properties derived with pedo-transfer functions, e.g. ~~hydrological soil properties~~, hydrologic soil properties such as saturated hydraulic conductivity (Pachepsky and Rawls, 2004) and complex properties that depend on multiple primary properties, e.g., carbon stocks. These layers are important inputs to model and map soil functions in the present and in the future as well as to support Earth System Modelling (Luo et al., 2016; Dai et al., 2019).

~~The~~ Sixth, the quantification of uncertainty is recommended and is becoming more common in DSM studies; ~~this~~. This work introduced it at global scale for the first time to our knowledge. While the provision of quantiles is mentioned in the GlobalSoilMap specifications, the representation and communication of uncertainty to end users and stakeholders remain an important research field to be further explored. The appropriate uncertainty intervals, both in terms of user acceptance and modelling feasibility, also need to be investigated.

Finally, the integration of highly automatised workflows with expert opinion should be further explored. DSM products use statistical models to describe soils and it is important to take into account the expertise and experience of pedologists, at least in an evaluation loop if not as part of the modelling itself. We made a first attempt at this in the "Qualitative Evaluation" section, above, but do not have a method to effectively incorporate expert observations into a workflow.

4 Conclusions

This study presents and discusses the production of global maps of soil properties as implemented in the SoilGrids 2.0 product, with cross-validation, hyper-parameters selection and quantification of uncertainty, using the 'best available' (shared) soil profile data for the world. In particular, the study describes a robust and reproducible DSM work-flow addressing the challenges of global data modelling: 1. non-homogeneous spatial distribution of input soil observations; 2. robust quantitative evaluation with a cross-validation procedure balancing accuracy and performances; 3. qualitative evaluation of the spatial patterns of the maps to include information about matching with well recognised pedo-landscape features; 4. quantification and mapping of the spatial uncertainty to provide users with a measure for and warning for users of the products.

As such, it describes a next step into global modelling and mapping of soil properties, explicitly highlighting the importance of quantitative and qualitative evaluation and uncertainty communication. The actual use of SoilGrids 2.0 in global and wide-area regional applications, where soil properties are important model inputs, will be the real test of its applicability and usefulness.

480 *Code and data availability.* The code underpinning the SoilGrids 2.0 workflow is available under the GPL3 license at the SoilGrids git repository.

SoilGrids predictions themselves are available to the public under the Creative Commons CC-BY 4.0 licence, facilitating their widespread use. They may be obtained as world mosaics in the Virtual Raster Tile (VRT) format from a WebDAV service. A suite of Web Coverage Services (WCS) facilitates automated access, e.g., from computer programmes or modelling frameworks. A set of notebooks was developed
485 with examples for the use of the WCS. A new web based portal was also developed with this release, providing users with a light and swift means to visualise and explore the new predictions, making the best of state-of-the-art technologies for the web. A ReST API in beta stage is also available <https://rest.soilgrids.org/>.

Appendix A: Environmental covariates

Over four hundred geographic layers were available as environmental covariates for this work. These are chosen for their presumed relation to the major soil forming factors.

Table A1: Covariates sets and sources

Weather and Climate
<ul style="list-style-type: none">• Temperature and precipitation from the Climatologies at high resolution for the earth's land surface areas (CHELSEA) dataset (Karger et al., 2016).• Snowfall from ESA's CCI Land Cover dataset (Bontemps et al., 2013).• Cloud cover by EarthEnv (Wilson and Jetz, 2016).• Temperature and water vapour from NASA's MODIS products (Wan et al., 2006).• Precipitation, solar radiation, temperature, water vapour, wind speed plus various indexes from the WorldClim version 2 climate data series (Fick and Hijmans, 2017).
Ecology and ecosystems
<ul style="list-style-type: none">• Bioclimatic zones in the Global Ecophysiology product by the USGS Geosciences and Environmental Change Science Center (GECSC) (Dinerstein et al., 2017).
Geology
<ul style="list-style-type: none">• Average soil and sedimentary-deposit thickness by the Distributed Active Archive Centre (DAAC) (Pelletier et al., 2016).• Rock types by the USGS Geosciences and Environmental Change Science Center (GECSC), based on the Global Lithological Map database v1.1 (Hartmann and Moosdorf, 2012).
Land Use and Land Cover
<ul style="list-style-type: none">• 2010 land cover classes from ESA's land cover CCI (Bontemps et al., 2013).• Bare ground and tree cover from the USGU Global Land Cover dataset (Hansen et al., 2013).• 2010 land cover classes from the NGCC's GLOBE Land30 product (Chen et al., 2015).
Elevation and morphology
<ul style="list-style-type: none">• Land surface elevation from the EarthEnv-DEM90 dataset (Robinson et al., 2014).• Land surface elevation and various morphology indexes from the WorldGrids dataset (Reuter and Hengl, 2016).• Land form classes in the Global Ecophysiology product by the USGS Geosciences and Environmental Change Science Center (GECSC) (Sayre et al., 2014).
Core satellites outputs

- Bands 3, 4, 5 and 7 from Landsat (Zanter, 2019)
 - Middle and near infra-red bands from MODIS (Savtchenko et al., 2004).
-

Vegetation Indexes

- NDVI from Landsat (Zanter, 2019).
 - EVI and NPP from MODIS (Savtchenko et al., 2004).
-

Hydrography

- Global Inundation Extent from Multi-Satellites (GIEMS) dataset by Estellus (Fluet-Chouinard et al., 2015).
 - Extent of glaciers, surface water change and occurrence probability by the JRC (Pekel et al., 2016).
 - Global water table depth (Fan et al., 2013).
-

Appendix B: Bio-climatic regions

Table B1 summarises the number of observations per property for each bio-climatic region. An interactive map of the regions is available on-line.

Table B1. Number of observations per property for each bio-climatic region. See Table 1 for abbreviations.

Biome	CEC	CFVO	N	pH	SOC	STP
Tropical & Subtropical Moist Broadleaf Forests	4185	2117	8378	12872	11901	11651
Tropical & Subtropical Dry Broadleaf Forests	558	205	1370	3264	2724	3051
Tropical & Subtropical Coniferous Forests	59	30	54	1336	878	1331
Temperate Broadleaf & Mixed Forests	12585	29708	24711	56569	49727	61822
Temperate Conifer Forests	6058	6417	5812	7597	9490	9834
Boreal Forests/Taiga	1443	3210	4834	4140	6819	5358
Tropical & Subtropical Grasslands, Savannas & Shrublands	8391	8259	20181	27633	24951	23135
Temperate Grasslands, Savannas & Shrublands	13442	9885	9812	23654	24421	25416
Flooded Grasslands & Savannas	246	124	503	754	818	798
Montane Grasslands & Shrublands	479	1865	1073	1386	3994	3568
Tundra	312	199	466	548	807	695
Mediterranean Forests, Woodlands & Scrub	1747	5951	8034	9126	12532	11428
Deserts & Xeric Shrublands	3342	3412	3224	8994	8163	9862
Mangroves	88	26	165	264	437	250

Author contributions. LP and LdS conceived and executed the research and wrote the paper. NB, GH, BK gave suggestions about the approach and contributed extensively to the paper. DR designed and executed the qualitative evaluation and wrote the corresponding sections in the paper. NB and ER designed and created the database of soil observations. All authors reviewed the paper.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This work was funded from ISRIC core funding, with additional support from the EU-H2020 CIRCASA project. ISRIC – World Soil Information, legally registered as the International Soil Reference and Information Centre, receives core funding from the Netherlands Government. We thank Rik van den Bosch (ISRIC Director) for the internal support to the project. We specially thank the organizations and experts that provided soil point data for consideration in WoSIS and SoilGrids. LP is member of a consortium supported by LE STUDIUM Loire Valley Institute for Advanced Studies through its LE STUDIUM Research Consortium Programme.

References

- Aitchison, J.: The statistical analysis of compositional data., Chapman & Hall, London., 1986.
- 505 Akpa, S. I. C., Odeh, I. O. A., Bishop, T. F. A., and Hartemink, A. E.: Digital Mapping of Soil Particle-Size Fractions for Nigeria, *Soil Science*, 78, 1953–1966, <https://doi.org/10.2136/sssaj2014.05.0202>, 2014.
- Arrouays, D., Grundy, M. G., Hartemink, A. E., Hempel, J. W., Heuvelink, G. B., Hong, S. Y., Lagacherie, P., Lelyk, G., McBratney, A. B., McKenzie, N. J., d.L. Mendonca-Santos, M., Minasny, B., Montanarella, L., Odeh, I. O., Sanchez, P. A., Thompson, J. A., and Zhang, G.-L.: Chapter Three - GlobalSoilMap: Toward a Fine-Resolution Global Grid of Soil Properties, vol. 125 of *Advances in Agronomy*, pp. 93 – 134, Academic Press, <https://doi.org/http://dx.doi.org/10.1016/B978-0-12-800137-0.00003-0>, 2014.
- 510 Arrouays, D., Leenaars, J. G. B., Richer-de Forges, A. C., Adhikari, K., Ballabio, C., Greve, M., Grundy, M., Guerrero, E., Hempel, J., Hengl, T., Heuvelink, G., Batjes, N., Carvalho, E., Hartemink, A., Hewitt, A., Hong, S.-Y., Krasilnikov, P., Lagacherie, P., Lelyk, G., Libohova, Z., Lilly, A., McBratney, A., McKenzie, N., Vasquez, G. M., Mulder, V. L., Minasny, B., Montanarella, L., Odeh, I., Padarian, J., Poggio, L., Roudier, P., Saby, N., Savin, I., Searle, R., Solbovoy, V., Thompson, J., Smith, S., Sulaeman, Y., Vintila, R., Rossel, R. V., Wilson, P., Zhang, G.-L., Swerts, M., Oorts, K., Karklins, A., Feng, L., Ibelle Navarro, A. R., Levin, A., Laktionova, T., Dell'Acqua, M., Suvannang, N., Ruam, W., Prasad, J., Patil, N., Husnjak, S., Pásztor, L., Okx, J., Hallet, S., Keay, C., Farewell, T., Lilja, H., Juilleret, J., Marx, S., Takata, Y., Kazuyuki, Y., Mansuy, N., Panagos, P., Van Liedekerke, M., Skalsky, R., Sobocka, J., Kobza, J., Eftekhari, K., Alavipanah, S. K., Moussadek, R., Badraoui, M., Da Silva, M., Paterson, G., Gonçalves, M. d. C., Theocharopoulos, S., Yemefack, M., Tedou, S., Vrscaj, B., Grob, U., Kozák, J., Boruvka, L., Dobos, E., Taboada, M., Moretti, L., and Rodriguez, D.: Soil legacy data rescue via GlobalSoilMap and other international and national initiatives, *GeoResJ*, 14, 1–19, <https://doi.org/10.1016/j.grj.2017.06.001>, 2017.
- 520 Arrouays, D., McBratney, A., Bouma, J., Libohova, Z., de Forges, A. C. R., Morgan, C. L., Roudier, P., Poggio, L., and Mulder, V. L.: Impressions of digital soil maps: The good, the not so good, and making them ever better, *Geoderma Regional*, 20, e00255, <https://doi.org/https://doi.org/10.1016/j.geodrs.2020.e00255>, <http://www.sciencedirect.com/science/article/pii/S2352009420300043>, 2020.
- 525 Ballabio, C., Panagos, P., and Monatanarella, L.: Mapping topsoil physical properties at European scale using the LUCAS database., *Geoderma*, 261, 110–123, 2016.
- Banwart, S., Black, H., Cai, Z., Gicheru, P., Joosten, H., Victoria, R., Milne, E., Noellemeyer, E., Pascual, U., Nziguheba, G., Vargas, R., Bationo, A., Buschiazzi, D., de Brogniez, D., Melillo, J., Richter, D., Termansen, M., van Noordwijk, M., Goverse, T., Ballabio, C., Bhattacharyya, T., Goldhaber, M., Nikolaidis, N., Zhao, Y., Funk, R., Duffy, C., Pan, G., la Scala, N., Gottschalk, P., Batjes, N., Six, J., van Wesemael, B., Stocking, M., Bampa, F., Bernoux, M., Feller, C., Lemanceau, P., and Montanarella, L.: Benefits of soil carbon: report on the outcomes of an international scientific committee on problems of the environment rapid assessment workshop, *Carbon Management*, 5, 185–192, <http://dx.doi.org/10.1080/17583004.2014.913380>, 2014.
- Barnes, M.: Aichi targets: Protect biodiversity, not just area, *Nature*, 526, 195–195, <http://dx.doi.org/10.1038/526195e>, 2015.
- Barnes, R., Sahr, K., Evenden, G., Johnson, A., and Warmerdam, F.: dggridR: discrete global grids for R, <https://github.com/r-barnes/dggridR>, 2016.
- 535 Batjes, N.: ISRIC-WISE derived soil properties on a 5 by 5 arc-minutes global grid (ver. 1.2), Report Report 2012/01, ISRIC - World Soil Information, http://www.isric.org/sites/default/files/isric_report_2012_01.pdf, 2012.

- Batjes, N.: Harmonised soil property values for broad-scale modelling (WISE30sec) with estimates of global soil carbon stocks, *Geoderma*, 269, 61–68, <http://dx.doi.org/10.1016/j.geoderma.2016.01.034>(withsupplementalinformation:http://www.isric.org/sites/default/files/isric_report_2015_01.pdf), 2016.
- 540 Batjes, N., Al-Adamat, R., Bhattacharyya, T., Bernoux, M., Cerri, C., Gicheru, P., Kamoni, P., Milne, E., Pal, D., and Rawajfih, Z.: Preparation of consistent soil data sets for SOC modelling purposes: secondary SOTER data sets for four case study areas, *Agriculture, Ecosystems and Environment*, 122, 26–34, <http://dx.doi.org/10.1016/j.agee.2007.01.005>, 2007.
- Batjes, N. H.: Technologically achievable soil organic carbon sequestration in world croplands and grasslands, *Land Degradation & Development*, 30, 25–32, <https://onlinelibrary.wiley.com/doi/abs/10.1002/ldr.3209>, 2019.
- 545 Batjes, N. H., Ribeiro, E., and van Oostrum, A.: Standardised soil profile data to support global mapping and modelling (WoSIS snapshot 2019), *Earth System Science Data*, 2020, 299–320, <https://doi.org/10.5194/essd-12-299-2020>, <https://doi.org/10.5194/essd-12-299-2020>, 2020.
- Bontemps, S., Defourny, P., Radoux, J., Van Bogaert, E., Lamarche, C., Achard, F., Mayaux, P., Boettcher, M., Brockmann, C., Kirches, G., et al.: Consistent global land cover maps for climate modelling communities: Current achievements of the ESA’s land cover CCI, in: *Proceedings of the ESA Living Planet Symposium*, Edimburgh, pp. 9–13, 2013.
- 550 Borrelli, P., Robinson, D. A., Fleischer, L. R., Lugato, E., Ballabio, C., Alewell, C., Meusburger, K., Modugno, S., Schütt, B., Ferro, V., Bagarello, V., Oost, K. V., Montanarella, L., and Panagos, P.: An assessment of the global impact of 21st century land use change on soil erosion, *Nature Communications*, 8, 2013, <https://doi.org/10.1038/s41467-017-02142-7>, 2017.
- 555 Bouma, J.: Engaging Soil Science in Transdisciplinary Research Facing “Wicked” Problems in the Information Society, *Soil Sci. Soc. Am. J.*, 79, 454–458, <https://dx.doi.org/10.2136/sssaj2014.11.0470>, 2015.
- Brus, D., Kempen, B., and Heuvelink, G.: Sampling for validation of digital soil maps, *European Journal of Soil Science*, 62, 394–407, <https://doi.org/10.1111/j.1365-2389.2011.01364.x>, <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1365-2389.2011.01364.x>, 2011.
- Brus, D. J.: Statistical sampling approaches for soil monitoring, *European Journal of Soil Science*, 65, 779–791, <https://doi.org/10.1111/ejss.12176>, <https://onlinelibrary.wiley.com/doi/abs/10.1111/ejss.12176>, 2014.
- 560 Buchhorn, M., Lesiv, M., Tsendbazar, N. E., Herold, M., Bertels, L., and Smets, B.: Copernicus Global Land Cover Layers—Collection 2, *Remote Sensing*, 108, 1044, <https://doi.org/10.3390/rs12061044>, 2020.
- Chen, J., Chen, J., Liao, A., Cao, X., Chen, L., Chen, X., He, C., Han, G., Peng, S., Lu, M., et al.: Global land cover mapping at 30 m resolution: A POK-based operational approach, *ISPRS Journal of Photogrammetry and Remote Sensing*, 103, 7–27, 2015.
- 565 Chen, S., Mulder, V. L., Heuvelink, G. B. M., Poggio, L., Caubet, M., Román Dobarco, M., Walter, C., and Arrouays, D.: Model averaging for mapping topsoil organic carbon in France, *Geoderma*, 366, 114–127, <https://doi.org/https://doi.org/10.1016/j.geoderma.2020.114237>, <http://www.sciencedirect.com/science/article/pii/S0016706119316519>, 2020.
- Cowie, A. L., Orr, B. J., Castillo Sanchez, V. M., Chasek, P., Crossman, N. D., Erlewein, A., Louwagie, G., Maron, M., Metternicht, G. I., Minelli, S., Tengberg, A. E., Walter, S., and Welton, S.: Land in balance: The scientific conceptual framework for Land Degradation Neutrality, *Environmental Science & Policy*, 79, 25–35, <https://doi.org/10.1016/j.envsci.2017.10.011>, 2018.
- 570 Dai, Y., Shangguan, W., Wang, D., Wei, N., Xin, Q., Yuan, H., Zhang, S., Liu, S., and Yan, F.: A review on the global soil datasets for earth system modeling, *SOIL*, 5, 137–158, <https://doi.org/10.5194/soil-5-137-2019>, 2019.
- de Sousa, L. M., Poggio, L., and Kempen, B.: Comparison of FOSS4G Supported Equal-Area Projections Using Discrete Distortion Indicators, *ISPRS International Journal of Geo-Information*, 8, 351, 2019.

- 575 de Sousa, L. M., Poggio, L., Dawes, G., Kempen, B., and Van Den Bosch, R.: Computational Infrastructure of SoilGrids 2.0, in: International Symposium on Environmental Software Systems, pp. 24–31, Springer, 2020.
- Deutsch, C. and Journel, A.: GSLIB: Geostatistical Software Library and User's Guide, Oxford University Press, New York, second edition edn., 1998.
- Dinerstein, E., Olson, D., Joshi, A., Vynne, C., Burgess, N. D., Wikramanayake, E., Hahn, N., Palminteri, S., Hedao, P., Noss, R., Hansen, 580 M., Locke, H., Ellis, E. C., Jones, B., Barber, C. V., Hayes, R., Kormos, C., Martin, V., Crist, E., Sechrest, W., Price, L., Baillie, J. E. M., Weeden, D., Suckling, K., Davis, C., Sizer, N., Moore, R., Thau, D., Birch, T., Potapov, P., Turubanova, S., Tyukavina, A., de Souza, N., Pintea, L., Brito, J. C., Llewellyn, O. A., Miller, A. G., Patzelt, A., Ghazanfar, S. A., Timberlake, J., Klöser, H., Shennan-Farpón, Y., Kindt, R., Lillesø, J.-P. B., van Breugel, P., Graudal, L., Voge, M., Al-Shammari, K. F., and Saleem, M.: An Ecoregion-Based Approach to Protecting Half the Terrestrial Realm, *BioScience*, 67, 534–545, <https://doi.org/10.1093/biosci/bix014>, <https://doi.org/10.1093/biosci/bix014>, 2017.
- 585 Dorji, T., Odeh, I. O. A., Field, D. J., and Baillie, I. C.: Digital soil mapping of soil organic carbon stocks under different land use and land cover types in montane ecosystems, Eastern Himalayas, *Forest Ecology and Management*, 318, 91–102, <https://doi.org/http://dx.doi.org/10.1016/j.foreco.2014.01.003>, <http://www.sciencedirect.com/science/article/pii/S037811271400005X>, 2014.
- 590 Ellili, Y., Walter, C., Michot, D., Pichelin, P., and Lemerrier, B.: Mapping soil organic carbon stock change by soil monitoring and digital soil mapping at the landscape scale, *Geoderma*, 351, 1–8, <https://doi.org/10.1016/j.geoderma.2019.03.005>, 2019.
- Fan, Y., Li, H., and Miguez-Macho, G.: Global Patterns of Groundwater Table Depth, *Science*, 339, 940–943, <https://doi.org/10.1126/science.1229881>, <https://science.sciencemag.org/content/339/6122/940>, 2013.
- FAO: Digital Soil Map of the World and derived properties (ver. 3.5), Report FAO Land and Water Digital Media Series 1, FAO, <http://www.fao.org/geonetwork/srv/en/metadata.show?id=14116>, 1995.
- 595 FAO: Guidelines for soil description (Fourth ed.), Report, FAO, <http://www.fao.org/docrep/019/a0541e/a0541e.pdf>, 2006.
- FAO and ITPS: Status of the world's soil resources (SWSR) - Main report, Report, Food and Agriculture Organization of the United Nations and Intergovernmental Technical Panel on Soils, <http://www.fao.org/3/a-i5199e.pdf>, 2015.
- FAO, IIASA, ISRIC, ISSCAS, and JRC: Harmonized World Soil Database (version 1.2), Report, Prepared by Nachtergaele, FO and van 600 Velthuisen, H and Verelst, L and Wiberg, D and Batjes, NH and Dijkshoorn, JA, and van Engelen, VWP, and Fischer, G, and Jones, A, and Montanarella, L., and Petri, M, and Prieler, S, and Teixeira, E and Xuezheng, Shi. Food and Agriculture Organization of the United Nations (FAO), International Institute for Applied Systems Analysis (IIASA), ISRIC - World Soil Information, Institute of Soil Science - Chinese Academy of Sciences (ISSCAS), Joint Research Centre of the European Commission (JRC), http://webarchive.iiasa.ac.at/Research/LUC/External-World-soil-database/HWSD_Documentation.pdf, 2012.
- 605 FAO, IFAD, UNICEF, WFP, and WHO: The State of Food Security and Nutrition in the World 2018. Building climate resilience for food security and nutrition, Report, FAO, <http://www.fao.org/3/I9553EN/i9553en.pdf>, 2018.
- FAO-ISRIC: Guidelines for soil description (3rd Edition, Rev.), FAO, Rome, 1986.
- FAO-Unesco: FAO-Unesco Soil Map of the World, 1:5,000,000 (Vol. 1 to 10), United Nations Educational, Scientific, and Cultural Organization, Paris, 1971-1981.
- 610 Fick, S. E. and Hijmans, R. J.: WorldClim 2: new 1-km spatial resolution climate surfaces for global land areas, *International journal of climatology*, 37, 4302–4315, 2017.

- Fluet-Chouinard, E., Lehner, B., Rebelo, L.-M., Papa, F., and Hamilton, S. K.: Development of a global inundation map at high spatial resolution from topographic downscaling of coarse-scale remote sensing data, *Remote Sensing of Environment*, 158, 348–361, 2015.
- 615 Gomes, L. C., Faria, R. M., de Souza, E., Veloso, G. V., Schaefer, C. E. G., and Fernandes Filho, E. I.: Modelling and mapping soil organic carbon stocks in Brazil, *Geoderma*, 340, 337–350, 2019.
- Goovaerts, P.: *Geostatistics for Natural Resources Evaluation*, Oxford University Press, 1997.
- GRASS Development Team: *Geographic Resources Analysis Support System (GRASS GIS) Software*, version 7.8.0, <http://www.grass.osgeo.org>, 2020.
- Grunwald, S., Thompson, J. A., and Boettinger, J. L.: Digital soil mapping and modeling at continental scales: Finding solutions for global issues, *Soil Science Society of America Journal*, 75, 1201–1213, <https://dx.doi.org/10.2136/sssaj2011.0025>, 2011.
- 620 GSP and FAO: *Pillar 4 Implementation Plan – Towards a Global Soil Information System*, Report, Global Soil Partnership, <http://www.fao.org/3/a-bl102e.pdf>, 2016.
- GSP and ITPS: *Global soil organic carbon map (GSOCmap)*, Report Technical Report, Global Soil Partnership (GSP) and International Panel on Soils (ITPS), <http://www.fao.org/3/i8891en/I8891EN.pdf>, 2018.
- 625 Guevara, M., Olmedo, G. F., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C., Arévalo, G. E., Arroyo-Cruz, C. E., Bolivar, A., Bunning, S., Bustamante Cañas, N., Cruz-Gaistardo, C. O., Davila, F., Dell Acqua, M., Encina, A., Figueredo Tacona, H., Fontes, F., Hernández Herrera, J. A., Ibelle Navarro, A. R., Loayza, V., Manueles, A. M., Mendoza Jara, F., Olivera, C., Osorio Hermosilla, R., Pereira, G., Prieto, P., Alexis Ramos, I., Rey Brina, J. C., Rivera, R., Rodríguez-Rodríguez, J., Roopnarine, R., Rosales Ibarra, A., Rosales Riveiro, K. A., Schulz, G. A., Spence, A., Vasques, G. M., Vargas, R. R., and Vargas, R.: No Silver Bullet for Digital Soil Mapping: Country-specific Soil Organic Carbon Estimates across Latin America, *SOIL*, 2018, 173–193, <https://doi.org/10.5194/soil-4-173-2018>, 2018.
- 630 Guyon, I., Weston, J., Barnhill, S., and Vapnik, V.: Gene selection for cancer classification using support vector machines, *Machine learning*, 46, 389–422, 2002.
- Han, E., Ines, A. V. M., and Koo, J.: Development of a 10-km resolution global soil profile dataset for crop modeling applications, *Environmental Modelling & Software*, <https://doi.org/10.1016/j.envsoft.2019.05.012>, 2019.
- 635 Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S., Tyukavina, A., Thau, D., Stehman, S., Goetz, S., Loveland, T. R., et al.: High-resolution global maps of 21st-century forest cover change, *science*, 342, 850–853, 2013.
- Harden, J. W., Hugelius, G., Ahlström, A., Blankinship, J. C., Bond-Lamberty, B., Lawrence, C. R., Loisel, J., Malhotra, A., Jackson, R. B., Ogle, S., Phillips, C., Ryals, R., Todd-Brown, K., Vargas, R., Vergara, S. E., Cotrufo, M. F., Keiluweit, M., Heckman, K. A., Crow, S. E., Silver, W. L., DeLonge, M., and Nave, L. E.: Networking our science to characterize the state, vulnerabilities, and management opportunities of soil organic matter, *Global Change Biology*, 24, e705–e718, <https://doi.org/10.1111/gcb.13896>, 2017.
- 640 Hartmann, J. and Moosdorf, N.: The new global lithological map database GLiM: A representation of rock properties at the Earth surface, *Geochemistry, Geophysics, Geosystems*, 13, 2012.
- Hengl, T., Mendes de Jesus, J., MacMillan, R. A., Batjes, N. H., Heuvelink, G. B., Ribeiro, E. C., Samuel-Rosa, A., Kempen, B., Leenaars, J. G., Walsh, M. G., and Gonzalez, M. R.: SoilGrids1km— global soil information based on automated mapping, *PLoS ONE*, 9, <https://doi.org/10.1371/journal.pone.0105992>, 2014.
- 645 Hengl, T., Leenaars, J. G. B., Shepherd, K. D., Walsh, M. G., Heuvelink, G. B. M., Mamo, T., Tilahun, H., Berkhout, E., Cooper, M., Fegraus, E., Wheeler, I., and Kwabena, N. A.: Soil nutrient maps of Sub-Saharan Africa: assessment of soil nutrient content at 250 m spatial resolution using machine learning, *Nutrient Cycling in Agroecosystems*, <https://doi.org/10.1007/s10705-017-9870-x>, 2017a.

- 650 Hengl, T., Mendes de Jesus, J., Heuvelink, G. B. M., Ruiperez Gonzalez, M., Kilibarda, M., Blagotić, A., Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas, R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G. B., Ribeiro, E., Wheeler, I., Mantel, S., and Kempen, B.: SoilGrids250m: Global gridded soil information based on machine learning, *PLOS ONE*, 12, e0169748, <http://dx.doi.org/10.1371/journal.pone.0169748>, 2017b.
- Heuvelink, G. B. M., Angelini, M. E., Poggio, L., Bai, Z., Batjes, N. H., van den Bosch, R., Bossio, D., Estella, S., Lehmann, J., Olmedo, 655 G. F., and Sanderman, J.: Machine learning in space and time for modelling soil organic carbon change, *European Journal of Soil Science*, n/a, <https://doi.org/https://doi.org/10.1111/ejss.12998>, <https://doi.org/10.1111/ejss.12998>, 2020.
- Hounkpatin, O. K., de Hipt, F. O., Bossa, A. Y., Welp, G., and Amelung, W.: Soil organic carbon stocks and their determining factors in the Dano catchment (Southwest Burkina Faso), *Catena*, 166, 298–309, 2018.
- IPBES: Global assessment report on biodiversity and ecosystem services of the Intergovernmental Science- Policy Platform on Biodi- 660 versity and Ecosystem Services. E. S. Brondizio, J. Settele, S. Díaz, and H. T. Ngo (editors), Report, IPBES, <https://www.ipbes.net/global-assessment-report-biodiversity-ecosystem-services>, 2019.
- IUSS Working Group WRB: World Reference Base for Soil Resources 2014; Update 2015. International Soil Classification System for Naming Soils and Creating Legends for Soil Maps, no. 106 in *World Soil Resources Reports*, FAO, Rome, <http://www.fao.org/3/i3794en/I3794en.pdf>, 2015.
- 665 Ivushkin, K., Bartholomeus, H., Bregt, A. K., Pulatov, A., Kempen, B., and de Sousa, L.: Global mapping of soil salinity change, *Remote Sensing of Environment*, 231, <https://doi.org/10.1016/j.rse.2019.111260>, <https://doi.org/10.1016/j.rse.2019.111260>, 2019.
- Janssen, P. H. M. and Heuberger, P. S. C.: Calibration of process-oriented models, *Ecological Modelling*, 83, 55–66, 1995.
- Karger, D., Conrad, O., Böhrer, J., Kawohl, T., Kreft, H., Soria-Auza, R., Zimmermann, N., Linder, H., and Kessler, M.: CHELSA climatologies at high resolution for the earth's land surface areas (Version 1.1), *World Data Center for Climate*, 2016.
- 670 Kempen, B., Brus, D. J., and de Vries, F.: Operationalizing digital soil mapping for nationwide updating of the 1:50,000 soil map of the Netherlands, *Geoderma*, 241–242, 313 – 329, <https://doi.org/https://doi.org/10.1016/j.geoderma.2014.11.030>, <http://www.sciencedirect.com/science/article/pii/S0016706114004285>, 2015.
- Kempen, B., Dalsgaard, S., Kaaya, A. K., Chamuya, N., Ruipérez-González, M., Pekkarinen, A., and Walsh, M. G.: Mapping topsoil organic carbon concentrations and stocks for Tanzania, *Geoderma*, 337, 164 – 180, <https://doi.org/https://doi.org/10.1016/j.geoderma.2018.09.011>, 675 <http://www.sciencedirect.com/science/article/pii/S0016706118305111>, 2019.
- Keskin, H. and Grunwald, S.: Regression kriging as a workhorse in the digital soil mapper's toolbox, *Geoderma*, 326, 22 – 41, <https://doi.org/https://doi.org/10.1016/j.geoderma.2018.04.004>, 2018.
- Khaledian, Y. and Miller, B. A.: Selecting appropriate machine learning methods for digital soil mapping, *Applied Mathematical Modelling*, 81, 401–418, 2020.
- 680 Kuhn, M.: A Short Introduction to the caret Package, *R Found Stat Comput*, pp. 1–10, 2015.
- Lark, R. and Bishop, T.: Cokriging particle size fractions of the soil., *European Journal of Soil Science*, 58, 763–774, 2007.
- Liu, F., Zhang, G.-L., Song, X., Li, D., Zhao, Y., Yang, J., Wu, H., and Yang, F.: High-Resolution and Three-Dimensional Mapping of Soil Texture of China, *Geoderma*, 361, 114–121, <https://doi.org/10.1016/j.geoderma.2019.114061>, 2020.
- Luo, Y., Ahlström, A., Allison, S., Batjes, N., Brovkin, V., Carvalhais, N., Chappell, A., Ciais, P., Davidson, E., Finzi, A., Georgiou, K., 685 Guenet, B., Hararuk, O., Harden, J., He, Y., Hopkins, F. M., Jiang, L., Koven, C., Jackson, R., Jones, C., Lara, M., Liang, J., McGuire, A., PARTON, W. J., Peng, C., Randerson, J., Salazar, A., Sierra, C., Smoth, M., Tian, H., Todd-Brown, K., Torn, M., van Groeningen,

- K., Wang, Y. P., Westm, O., Wei, Y., Wieder, W., Xia, J., Xia, X., Xu, X., and Zhu, T.: Towards more realistic projections of soil carbon dynamics by Earth System Models, *Global Biogeochem. Cycles*, 30, 40–56, <http://dx.doi.org/10.1002/2015GB005239>, 2016.
- 690 Ma, Y., Minasny, B., McBratney, A., Poggio, L., and Fajardo, M.: Predicting soil properties in 3D: Should depth be a covariate?, *Geoderma*, 383, 114–124, <https://doi.org/https://doi.org/10.1016/j.geoderma.2020.114794>, <http://www.sciencedirect.com/science/article/pii/S0016706120325490>, 2021.
- Malhotra, A., Todd-Brown, K., Nave, L. E., Batjes, N. H., Holmquist, J. R., Hoyt, A. M., Iversen, C. M., Jackson, R. B., Lajtha, K., Lawrence, C., Vindušková, O., Wieder, W., Williams, M., Hugelius, G., and Harden, J.: The landscape of soil carbon data: emerging questions, synergies and databases, *Progress in Physical Geography: Earth and Environment*, 43(5), 707–719, <https://doi.org/10.1177/0309133319873309>,
695 <http://dx.doi.org/10.1177/0309133319873309>, 2019.
- Mallavan, B., Minasny, B., and Mcbratney, A.: Homosoil, a Methodology for Quantitative Extrapolation of Soil Information Across the Globe, vol. 2, pp. 137–150, https://doi.org/10.1007/978-90-481-8863-5_12, 2010.
- McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping, *Geoderma*, 117, 3–52, [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.
- 700 Meinshausen, N.: Quantile regression forests, *Journal of Machine Learning Research*, 7, 983–999, 2006.
- Meyer, H. and Pebesma, E.: Predicting into unknown space? Estimating the area of applicability of spatial prediction models, *ArXiv*, <abs/2005.07939>, 2020.
- Meyer, H., Reudenbach, C., Hengl, T., Katurji, M., and Nauss, T.: Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation, *Environmental Modelling & Software*, 101, 1
705 – 9, <https://doi.org/https://doi.org/10.1016/j.envsoft.2017.12.001>, <http://www.sciencedirect.com/science/article/pii/S1364815217310976>, 2018.
- Minasny, B. and McBratney, A. B.: Digital soil mapping: A brief history and some lessons, *Geoderma*, 264, 301–311, <https://doi.org/10.1016/j.geoderma.2015.07.017>, 2016.
- Mora-Vallejo, A., Claessens, L., Stoorvogel, J., and Heuvelink, G. B. M.: Small scale digital soil mapping in Southeastern Kenya, *CATENA*,
710 76, 44–53, <http://dx.doi.org/10.1016/j.catena.2008.09.008>, 2008.
- Moulatlet, G. M., Zuquim, G., Figueiredo, F. O. G., Lehtonen, S., Emilio, T., Ruokolainen, K., and Tuomisto, H.: Using digital soil maps to infer edaphic affinities of plant species in Amazonia: Problems and prospects, *Ecology and Evolution*, pp. n/a–n/a, <http://dx.doi.org/10.1002/ece3.3242>, 2017.
- Nauman, T. W. and Duniway, M. C.: Relative prediction intervals reveal larger uncertainty in 3D approaches to predictive digital soil
715 mapping of soil properties with legacy data, *Geoderma*, 347, 170 – 184, <https://doi.org/https://doi.org/10.1016/j.geoderma.2019.03.037>, <http://www.sciencedirect.com/science/article/pii/S0016706118319347>, 2019.
- Nijbroek, R., Piikki, K., Söderström, M., Kempen, B., Turner, K., Hengari, S., and Mutua, J.: Soil Organic Carbon Baselines for Land Degradation Neutrality: Map Accuracy and Cost Tradeoffs with Respect to Complexity in Otjozondjupa, Namibia, *Sustainability*, 10,
1610, <https://dx.doi.org/10.3390/su10051610>, 2018.
- 720 NRCS National Soil Survey Center: Gridded Soil Survey Geographic (gSSURGO) Database User Guide, Version 2.2, Tech. rep., NRCS, http://www.nrcs.usda.gov/wps/PA_NRCSCconsumption/download?cid=nrcs142p2_051847&ext=pdf, 2016.
- Nussbaum, M., Spiess, K., Baltensweiler, A., Grob, U., Keller, A., Greiner, L., Schaeppman, M. E., and Papritz, A.: Evaluation of digital soil mapping approaches with large sets of environmental covariates, *SOIL*, 4, 1–22, <https://dx.doi.org/10.5194/soil-4-1-2018>, 2018.

- Omuto, C., Nachtergaele, F., and Vargas Rojas, R.: State of the Art Report on Global and Regional Soil Information: Where are we? Where to go?, Report, FAO, <http://www.fao.org/docrep/017/i3161e/i3161e.pdf>, 2012.
- 725 Oreskes, N.: Evaluation (Not Validation) of Quantitative Models., *Environmental Health Perspectives*, 106, 1453–1460, <https://doi.org/10.1289/ehp.98106s61453>, 1998.
- Pachepsky, Y. and Rawls, W. J., eds.: Development of Pedotransfer Functions in Soil Hydrology, no. 30 in *Developments in Soil Science*, Elsevier, 2004.
- 730 Pekel, J.-F., Cottam, A., Gorelick, N., and Belward, A. S.: High-resolution mapping of global surface water and its long-term changes, *Nature*, 540, 418, 2016.
- Pelletier, J., Broxton, P., Hazenberg, P., Zeng, X., Troch, P., Niu, G., Williams, Z., Brunke, M., and Gochis, D.: Global 1-km Gridded Thickness of Soil, Regolith, and Sedimentary Deposit Layers, <https://doi.org/10.3334/ORNLDAAAC/1304>, http://daac.ornl.gov/cgi-bin/dsvviewer.pl?ds_id=1304, 2016.
- 735 Piikki, K., Söderström, M., and Stadig, H.: Local adaptation of a national digital soil map for use in precision agriculture, *Advances in Animal Biosciences*, 8, 430–432, <https://dx.doi.org/10.1017/s2040470017000966>, 2017.
- Poggio, L. and Gimona, A.: 3D mapping of soil texture in Scotland, *Geoderma Regional*, 9, 5 – 16, <https://doi.org/https://doi.org/10.1016/j.geodrs.2016.11.003>, <http://www.sciencedirect.com/science/article/pii/S2352009416301286>, digital soil mapping across the globe, 2017a.
- 740 Poggio, L. and Gimona, A.: Assimilation of optical and radar remote sensing data in 3D mapping of soil properties over large areas, *Science of The Total Environment*, 579, 1094 – 1110, <https://doi.org/https://doi.org/10.1016/j.scitotenv.2016.11.078>, <http://www.sciencedirect.com/science/article/pii/S0048969716325177>, 2017b.
- Poggio, L., Gimona, A., and Brewer, M. J.: Regional scale mapping of soil properties and their uncertainty with a large number of satellite-derived covariates, *Geoderma*, 209–210, 1–14, <http://dx.doi.org/10.1016/j.geoderma.2013.05.029>, 2013.
- 745 R Core Team: R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, <http://www.R-project.org/>, ISBN 3-900051-07-0, 2020.
- Reuter, H. and Hengl, T.: Worldgrids—a public repository of global soil covariates, *Computing Ethics: A Multicultural Approach*, p. 287, 2016.
- Reuter, H. I. and Hengl, T.: Worldgrids-a public repository of global soil covariates, in: *Digital Soil Assessments and Beyond - Proceedings of the Fifth Global Workshop on Digital Soil Mapping, 5th Global Workshop on Digital Soil Mapping*, pp. 287–292, sg2019, <https://tinyurl.com/yyx7vd9s>, 2012.
- 750 Ribeiro, E., Batjes, N., and Van Oostrum, A.: World Soil Information Service (WoSIS) - Towards the standardization and harmonization of world soil data. *Procedures Manual 2018*, Report ISRIC Report 2018/01, ISRIC - World Soil Information, <http://dx.doi.org/10.17027/isric-wdcsoils.20180001>, 2018.
- 755 Robinson, N., Regetz, J., and Guralnick, R. P.: EarthEnv-DEM90: A nearly-global, void-free, multi-scale smoothed, 90m digital elevation model from fused ASTER and SRTM data, *ISPRS Journal of Photogrammetry and Remote Sensing*, 87, 57–67, 2014.
- Rockstroem, J., Falkenmark, M., Lannerstad, M., and Karlberg, L.: The planetary water drama: Dual task of feeding humanity and curbing climate change, *Geophys. Res. Lett.*, 39, L15 401, <http://dx.doi.org/10.1029/2012gl051688>, 2012.
- Rossiter, D. G.: Maps and Models Are Never Valid, but They Can Be Evaluated (with Responses), *Pedometron*, 41, 19–28, <http://www.pedometrics.org/Pedometron/Pedometron41.pdf>, 2017.
- 760

- Roudier, P., Burge, O. R., Richardson, S. J., McCarthy, J. K., Grealish, G., and Ausseil, A.-G.: National Scale 3D Mapping of Soil pH Using a Data Augmentation Approach, *Remote Sensing*, 12, 2872, <https://www.mdpi.com/2072-4292/12/18/2872>, 2020.
- Sanderman, J., Hengl, T., and Fiske, G. J.: Soil carbon debt of 12,000 years of human land use, *Proceedings of the National Academy of Sciences*, <http://dx.doi.org/10.1073/pnas.1706103114>, 2017.
- 765 Savtchenko, A., Ouzounov, D., Ahmad, S., Acker, J., Leptoukh, G., Koziana, J., and Nickless, D.: Terra and Aqua MODIS products available from NASA GES DAAC, *Advances in Space Research*, 34, 710–714, 2004.
- Sayre, R., Dangermond, J., Frye, C., Vaughan, R., Aniello, P., Breyer, S., Cribbs, D., Hopkins, D., Nauman, R., Derrenbacher, W., et al.: A new map of global ecological land units—an ecophysiological stratification approach, Washington, DC: Association of American Geographers, 2014.
- 770 Schoeneberger, P., Wysicki, D., Benham, E., and Broderson, W.: Field book for describing and sampling soils (ver. 3.0), Natural Resources Conservation Service, National Soil Survey Center, Lincoln NE, http://www.nrcs.usda.gov/Internet/FSE_DOCUMENTS/nrcs142p2_052523.pdf, 2012.
- Shangguan, W., Dai, Y., Duan, Q., Liu, B., and Yuan, H.: A global soil data set for earth system modeling, *Journal of Advances in Modeling Earth Systems*, 6, 249–263, <http://dx.doi.org/10.1002/2013MS000293>, 2014.
- 775 Shrestha, D. L. and Solomatine, D. P.: Machine learning approaches for estimation of prediction interval for the model output, *Neural Networks*, 19, 225 – 235, <https://doi.org/https://doi.org/10.1016/j.neunet.2006.01.012>, *earth Sciences and Environmental Applications of Computational Intelligence*, 2006.
- Smith, P., Soussana, J.-F., Angers, D., Schipper, L., Chenu, C., Rasse, D., Batjes, N. H., van Egmond, F., McNeill, S., Kuhnert, M., Arias-Navaro, C., Olesen, J. E., Chirinda, N., Fornara, D., Joosse, P., Wollenberg, L., Alvaro-Fuentes, J., and Cobena, A.: How to measure, report and verify soil carbon change to realise the potential of soil carbon sequestration for atmospheric greenhouse gas removal, *Global Change Biology*, pp. 1–23, <https://doi.org/10.1111/gcb.14815>, <https://doi.org/10.1111/gcb.14815>, 2019.
- 780 Soussana, J.-F., Lutfalla, S., Ehrhardt, F., Rosenstock, T., Lamanna, C., Havlík, P., Richards, M., Wollenberg, E., Chotte, J.-L., Torquebiau, E., Ciais, P., Smith, P., and Lal, R.: Matching policy and science: Rationale for the ‘4 per 1000 - soils for food security and climate’ initiative, *Soil and Tillage Research*, <https://doi.org/10.1016/j.still.2017.12.002>, 2017.
- 785 Springmann, M., Clark, M., Mason-D’Croz, D., Wiebe, K., Bodirsky, B. L., Lassaletta, L., de Vries, W., Vermeulen, S. J., Herrero, M., Carlson, K. M., Jonell, M., Troell, M., DeClerck, F., Gordon, L. J., Zurayk, R., Scarborough, P., Rayner, M., Loken, B., Fanzo, J., Godfray, H. C. J., Tilman, D., Rockström, J., and Willett, W.: Options for keeping the food system within environmental limits, *Nature*, <https://doi.org/10.1038/s41586-018-0594-0>, 2018.
- Stockmann, U., Padarian, J., McBratney, A., Minasny, B., de Brogniez, D., Montanarella, L., Hong, S. Y., Rawlins, B. G., and Field, D. J.: Global soil organic carbon assessment, *Global Food Security*, 6, 9–16, <https://doi.org/http://dx.doi.org/10.1016/j.gfs.2015.07.001>, <http://www.sciencedirect.com/science/article/pii/S2211912415000231>, 2015.
- 790 Stoorvogel, J. J., Bakkenes, M., Temme, A. J. A. M., Batjes, N. H., and ten Brink, B.: S-World: a Global Soil Map for Environmental Modelling, *Land Degradation and Development*, 28, 22–33, <http://dx.doi.org/10.1002/ldr.2656>, 2017.
- Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., and Zeileis, A.: Conditional variable importance for random forests, *BMC Bioinformatics*, 9, 307, <https://doi.org/10.1186/1471-2105-9-307>, 2008.
- 795 Todd-Brown, K. E. O., Randerson, J. T., Post, W. M., Hoffman, F. M., Tarnocai, C., Schuur, E. A. G., and Allison, S. D.: Causes of variation in soil carbon simulations from CMIP5 Earth system models and comparison with observations, *Biogeosciences*, 10, 1717–1736, <https://dx.doi.org/10.5194/bg-10-1717-2013>, 2013.

- Tóth, G., Jones, A., and Montanarella, L.: LUCAS Topsoil survey: methodology, data and results, Report, Land Resource Management Unit
800 - Soil Action, European Commission Joint Research Centre Institute for Environment and Sustainability, https://esdac.jrc.ec.europa.eu/ESDB_Archive/eusoils_docs/other/EUR26102EN.pdf, 2013.
- UNEP: The benefits of soil carbon - managing soils for multiple, economic, societal and environmental benefits, pp. 19–33, United Nations Environmental Programme, Nairobi, <http://www.unep.org/yearbook/2012>, 2012.
- van Bussel, L. G. J., Grassini, P., Van Wart, J., Wolf, J., Claessens, L., Yang, H., Boogaard, H., de Groot, H., Saito, K., Cassman, K. G.,
805 and van Ittersum, M. K.: From field to atlas: Upscaling of location-specific yield gap estimates, *Field Crops Research*, 177, 98–108, <http://dx.doi.org/10.1016/j.fcr.2015.03.005>, 2015.
- van der Esch, S., Brink, B. t., Stehfest, E., Bakkenes, M., Sewell, A., Bouwman, A., Meijer, J., Westhoek, H., and van den Berg, M.: Exploring future changes in land use and land condition and the impacts on food, water, climate change and biodiversity: Scenarios for the UNCCD Global Land Outlook, Report, UNCCD, <https://tinyurl.com/yagvs9vu>, 2017.
- 810 van Ittersum, M. K., Cassman, K. G., Grassini, P., Wolf, J., Tiftonell, P., and Hochman, Z.: Yield gap analysis with local to global relevance—A review, *Field Crops Research*, 143, 4–17, <http://dx.doi.org/10.1016/j.fcr.2012.09.009>, 2013.
- Vitharana, U. W. A., Mishra, U., and Mapa, R. B.: National soil organic carbon estimates can improve global estimates, *Geoderma*, 337, 55–64, <https://doi.org/https://doi.org/10.1016/j.geoderma.2018.09.005>, <http://www.sciencedirect.com/science/article/pii/S0016706118305925>, 2019.
- 815 Wan, Z. et al.: MODIS land surface temperature products users’ guide, Institute for Computational Earth System Science, University of California: Santa Barbara, CA, USA, 2006.
- Wilson, A. M. and Jetz, W.: Remotely sensed high-resolution global cloud dynamics for predicting ecosystem and biodiversity distributions, *PLoS biology*, 14, e1002415, 2016.
- WOCAT: Where the land is greener: Case studies and analysis of soil and water conservation initiatives worldwide, CTA, UNEP, FAO and
820 CDE, Berne, 2007.
- Wright, M. N. and Ziegler, A.: ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R, *Journal of Statistical Software*, 77, 1–17, <https://doi.org/10.18637/jss.v077.i01>, 2017.
- Yigini, Y. and Panagos, P.: Assessment of soil organic carbon stocks under future climate and land cover changes in Europe, *Science of The Total Environment*, 557–558, 838–850, <http://dx.doi.org/10.1016/j.scitotenv.2016.03.085>, 2016.
- 825 Yoo, A. B., Jette, M. A., and Grondona, M.: Slurm: Simple linux utility for resource management, in: *Workshop on Job Scheduling Strategies for Parallel Processing*, pp. 44–60, Springer, 2003.
- Zanter, K.: Landsat 4-7 Surface Reflectance Code (LEDAPS) Product Guide, <https://www.usgs.gov/land-resources/nli/landsat/landsat-normalized-difference-vegetation-index>, 2019.