



# Developing the Swiss soil spectral library for local estimation and monitoring

Philipp Baumann<sup>1</sup>, Anatol Helfenstein<sup>1, 2</sup>, Andreas Gubler<sup>3</sup>, Armin Keller<sup>4</sup>, Reto Giulio Meuli<sup>3</sup>, Daniel Wächter<sup>3</sup>, Juhwan Lee<sup>5</sup>, Raphael Viscarra Rossel<sup>6</sup>, and Johan Six<sup>1</sup>

<sup>1</sup>Institute of Agricultural Sciences, Department of Environmental Systems Science (D-USYS), ETH Zürich, Switzerland

<sup>2</sup>Soil Geography and Landscape Group, Wageningen University, PO Box 47, 6700 AA Wageningen, The Netherlands

<sup>3</sup>Swiss Soil Monitoring Network (NABO), Agroscope, Reckenholzstrasse 191, 8046 Zürich, Switzerland

<sup>4</sup>Swiss Competence Center for Soils (KOBO), School of Agricultural, Forest and Food Sciences HAFL, Bern University of Applied Sciences BFH, Bern, Switzerland

<sup>5</sup>Department of Agronomy and Medicinal Plant Resources, Gyeongnam National University of Science and Technology, Jinju, Republic of Korea

<sup>6</sup>Soil and Landscape Science, School of Molecular and Life Sciences, Curtin University, GPO Box U1987, Perth WA 6845, Australia

**Correspondence:** Philipp Baumann (baumann-philipp@protonmail.com), Current Address: Swiss Competence Center for Soils (KOBO), School of Agricultural, Forest and Food Sciences HAFL, Bern University of Applied Sciences BFH, Bern, Switzerland

## Abstract.

Information on soils' composition and physical, chemical and biological properties is paramount to elucidate agroecosystem functioning in space and over time. For this purposes we developed a national Swiss soil spectral library (SSL;  $n = 4374$ ) in the mid-infrared (mid-IR), calibrating 17 properties from legacy measurements on soils from the Swiss biodiversity monitoring program ( $n = 3778$ ; 1094 sites) and the Swiss long-term monitoring network ( $n = 596$ ; 71 sites). General models were trained with the interpretable rule-based learner CUBIST, testing combinations of {5, 10, 20, 50, 100} committees of rules and {2, 5, 7, 9} neighbors to localize predictions with repeated by location grouped ten-fold cross-validation. To evaluate the information in spectra to facilitate long-term soil monitoring at a plot-level, we conducted 71 model transfers for the NABO sites to induce locally relevant information from the SSL, using the data-driven sample selection method RS-LOCAL. Eleven soil properties were estimated with discrimination capacity suitable for screening ( $R^2 > 0.6$ ), out of which total carbon (C), organic C (OC), total N, organic matter content, pH, and clay showed accuracy eligible for accurate diagnostics ( $R^2 > 0.8$ ). CUBIST and the spectra estimated total C accurately with RMSE = 0.84 % while the measured range was 0.1–58.3 %, and OC with RMSE = 1.20 % (measured range 0.0–27.3 %). Compared to general estimates of properties from CUBIST, local modeling on average reduced the root mean square error of total C per site fourfold. We found that the selected SSL subsets were highly dissimilar in terms of both their spectral input space and the measured values. This suggests that data-driven selection with RS-LOCAL leverages chemical diversity in composition rather than similarity. Our results suggest that mid-IR soil estimates were sufficiently accurate to support many soil applications that require a large volume of input data, such as precision agriculture, soil C accounting and monitoring, and digital soil mapping. This SSL can be updated continuously, for example with samples



from deeper profiles and organic soils, so that the measurement of key soil properties becomes even more accurate and efficient in the near future.

## 1 Introduction

Soils provide a manifold of functions within terrestrial ecosystems, many of which are vital for humankind. To quantify these functions from the soils' composition and properties, one typically relies on physical, chemical and biological laboratory analytical measurements. Doing this consumes both financial resources and time. For example, repeated measurements are needed to describe soil functioning in changing environments, for example in response to agronomic management. Soil visible (vis) and infrared (IR) spectroscopic measurements and modeling have become indispensable tools to gather quick, relatively accurate, and inexpensive soil information, both on the field and in the laboratory (Nocita et al., 2015; Viscarra Rossel et al., 2016, 2017). When soil chemical and physical properties are calibrated to the spectra, a single mid-IR ( $4000 - 500 \text{ cm}^{-1}$ ;  $2500 - 25000 \text{ nm}$ ) or vis-NIR ( $25000 - 4000 \text{ cm}^{-1}$ ,  $400 - 2500 \text{ nm}$ ) measurement can be used to estimate multiple soil properties. Soil is a complex matrix with many organic and mineral components. This yields spectra with absorptions that overlap and contain many and often highly correlated variables. Hence, to successfully develop calibrations and make predictions for attributes related to soil composition, multivariate statistical methods are needed to disentangle relationships between these variables and measured attributes. Further, it is important to consider that the diversity in spectral characteristics typically reflects soils' chemical and physical composition. Since the soil composition is influenced by the soil forming factors — soil parent material, climate, topography, organisms, and age of soils (Dokuchaev, 1899; Jenny, 1941) — these factors provide further means of causally interpreting and judging the applicability of the method for a particular set of soils. Compared to the NIR, mid-IR offers a more accurate characterization of soils' chemistry since this region contains the fundamental vibrations with more defined peaks (Janik et al., 1998).

A soil spectral library (SSL) can be defined as a well-ordered and harmonized collection of soil samples, their spectra, analytical reference measurements, contextual information, and additional metadata on samples and methods used. A central question behind the development of large SSLs is how to achieve accurate local predictions based on established collections of soil information — for example within a new landscape, ecosystem, farm, field, or plot in a new region — where reference data of only a few observations are available. More recently, SSLs that span large geographical extents are being developed (Sila et al., 2016; Viscarra Rossel et al., 2016; Padarian et al., 2019b; England and Viscarra Rossel, 2018; Briedis et al., 2020; Angelopoulou et al., 2020; Dangal et al., 2019). These efforts are motivated by the prospect that soil spectroscopy can supplement many conventional methods of soil analysis. A range of predictive modeling strategies and algorithms have been tested for soil spectral analysis, among others involving tools from chemometrics (e.g., partial least squares (PLS) regression) (Janik and Skjemstad, 1995), traditional machine learning (e.g., regression tree methods) (Viscarra Rossel and Webster, 2012), to convolutional neural networks (CNNs) (Padarian et al., 2019a, b; Tsakiridis et al., 2020).

There are two main strategies for estimating properties of new soils using spectra. The first one is to calibrate one general or global model that is applied to predict new samples, and the other is to derive local calibrations by conditioning on a specific set



of observations and features of the SSL to new data based on soil knowledge and/or algorithms. However, empirical evaluation of local and global methods are needed in different contexts where data on soil attributes is needed (i.e. soil study, soil mapping project). Such studies or applications should consider the "no-free-lunch" theorems for machine learning and optimization (Wolpert, 1996; Wolpert and Macready, 1997): there is no single algorithm or combinations of them that work best under all situations or applications.

General statistical learning makes use of all available training data to construct one parametric model. In contrast, local learning methods combine different learning methods, supervised and/or unsupervised, and together with domain knowledge produce more modular forms of learning (Solomatine, 2008). The available training set can be a subset and algorithmic sub-models can thereby be optimized to more accurately predict new single observations or groups of them. Local learning has also been termed *transfer learning*. Transfer learning is a general expression for adapting previous knowledge gained from existing data (i.e., model representation) for a new prediction case (Pratt et al., 1993; Pratt and Thrun, 1997; Thrun and Pratt, 1998). It has been defined as a transfer from knowledge in (a) source task(s) or domain(s)—here, an SSL—to a target domain (Pan and Yang, 2010), and thus comprising soils from new locations in this case.

The soil spectroscopy community has suggested several approaches to achieve local calibrations based on an established SSL and its feature space. One example is augmenting (spiking) SSLs with a few weighted (Guerrero et al., 2016) or unweighted (Seidel et al., 2019) local samples. Other studies calibrated separate models on partitions of training data that were derived from applying certain criteria (i.e., geographical region, terrain attributes, parent material, soil type, land use, spectra-based clustering) (Sila et al., 2016; Ogen et al., 2019). Still others used memory-based learning based on spectral similarity, extracting useful information from compositional relatedness of soils (Ramirez-Lopez et al., 2013; Clairotte et al., 2016; Hong et al., 2019; Dangal et al., 2019) or additionally geographic proximity (Tziolas et al., 2019). These all produce individual models for each sample to be predicted. Memory-based learning combines both *lazy learning*, where a subset of stored samples are only retrieved and modeled when new samples are predicted, and *local learning* principles, where modeled subsets define points within a local neighborhood (Dietterich et al., 1993). The spectrum-based learner developed by Ramirez-Lopez et al. (2013) is a prominent memory-based method for which each new prediction sample forms its own target domain. The selection of source instances is governed by spectral similarity. **These properties make it as well a transfer learning method.** Another approach used by Padarian et al. (2019a) was re-training weights within specific layers of a deep convolutional neural network using local (target) sets, which were spectral soil data sets per country (parameter-transfer approach). Finally, the selection of matching SSL samples using the resampling-based selection RS-LOCAL algorithm has also been used (Lobsey et al., 2017). Lobsey et al. (2017) showed that this data-driven transfer approach outperforms most other current methods for deriving local estimates. Still, despite these promising learners, transferring the useful information contained within large and diverse SSLs, and their resulting calibrations onto new, local target areas with unique soil characteristics remains challenging due to soil complexity.

**RS-LOCAL obtains locally-relevant information by selecting specific rows (instances) from the training set and transfer them to the prediction set. RS-LOCAL is an example of an instance or sample transfer approach. It heavily relies on sampling and performance-driven reduction of the library, yielding a subset of samples that can accurately estimate the properties of soils in the local target task. We wanted to investigate this promising new method for local soil estimation and monitoring in**



Switzerland because it makes no prior assumptions on which samples from the library could be useful. This makes it potentially more accurate and as well more flexible to new local soil contexts than when creating constraints with similarity measures. A further advantage for large SSLs is that it removes samples that might be spectrally similar but cause inaccurate calibrations (i.e., erroneous measurements or spectra with confounding effects). Such a local approach however requires an well-established

5 and sufficiently diverse SSL in order to extract useful soils that are locally relevant.

Thus, our first goal was to develop a national mid-IR SSL with reference measurements for Switzerland to deliver 17 key chemical and physical soil proxies. This SSL includes soils and their analysis data from the long-term Swiss soil monitoring network (NABO; 71 agricultural sites with times series measurements,  $n = 596$ ) and the Swiss biodiversity monitoring (BDM) network (1094 grid-locations,  $n = 3778$ ). This is the first operational SSL for Switzerland in the mid-IR that allows means  
10 for spectral estimation with sufficient existing soil diversity. The second goal was to develop general rule-based models for all available soils properties using the CUBIST algorithm. Further, we wanted to infer important spectral regions in the models and their chemical associations, which we illustrated with the estimation of total carbon (C).

For soil monitoring, it is crucial to obtain locally unbiased spectral estimates of key soil properties such as organic C, from high soil variability of large SSLs and over time. This was our motivation to design a predictive transfer workflow that was  
15 adaptive to soils' composition and properties for each long-term monitoring site. Hence, our third goal was to leverage the SSL with its spatial and temporal variability of soils to derive local calibrations by transfer learning with RS-LOCAL. Specifically, we aimed for reproducing time-series measurements (starting from 1985) of soil C at the Swiss agricultural long-term monitoring sites based on spectral analyses and two calibration samples per site. To the best of our knowledge, there is no study yet that has evaluated the usefulness of a large and diverse SSL for systematic soil monitoring. We therefore wanted to design a local  
20 calibration strategy using transfer learning, that would be effective in reducing (conditional) bias at monitoring plots compared to the general rules derived in the first aim. Furthermore, we had a strong interest in identifying the mechanisms, considering both soil knowledge and data distributions, of how such a local transfer would induce locally-adaptive soil estimation.

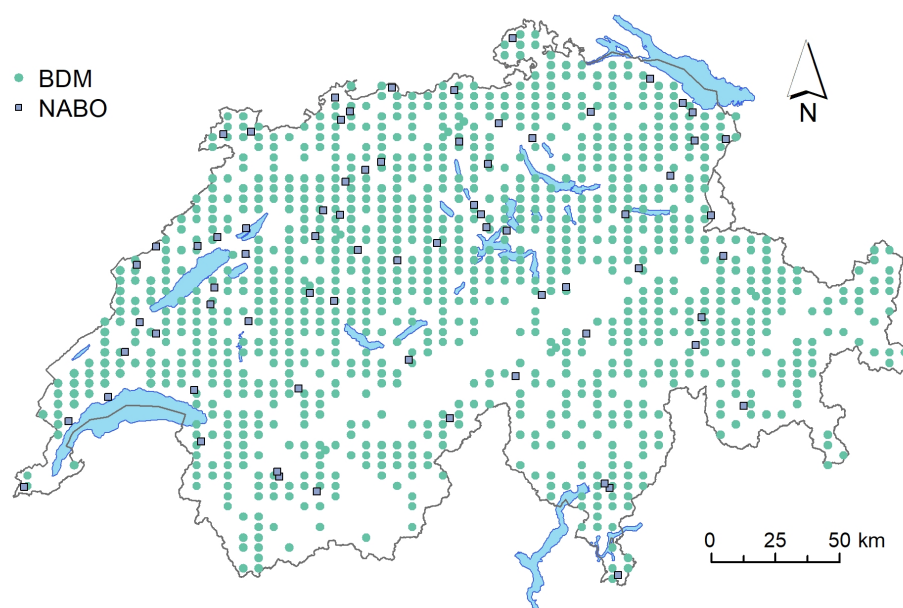




## 2 Material and Methods

### 2.1 Soils and data sets

To establish the Swiss SSL, we obtained soil samples and reference data from two different sources: 1) the Swiss soil monitoring network (NABO), and 2) the Swiss Biodiversity Monitoring (BDM) program (Bundesamt für Umwelt (BAFU), 2014) (Figure 1). The NABO currently consists of 108 sites where soils are being continuously measured every five years since 1985 for long-term soil monitoring. Out of the 108 sites, we selected 71 sites under agricultural management—comprising of arable land (33 sites), permanent grassland (26 sites), and special crops (11 sites; horticulture)—and one protected area. For the mid-infrared SSL, we used 596 NABO soil samples from 6 campaigns conducted between 1985 and 2015.



**Figure 1.** Swiss map with sampling locations of mid-infrared spectral library including the sites of the Biodiversity Monitoring Program (BDM;  $6 \times 4$  km;  $n = 1094$ ) and the National Soil Monitoring Network (NABO;  $n = 71$ ). 71 NABO sites ( $10 \text{ m} \times 10 \text{ m}$ ) were sampled with a grid-based stratified design. 1094 BDM samples were obtained from single sampling events. The NABO sites have been continuously sampled and measured in five-year intervals since 1985.

The plots at the NABO sites covered  $10 \text{ m} \times 10 \text{ m}$  each. These were repeatedly sampled for 0–20 cm soil depth. Four replicate samples were taken by stratified random sampling and bulking four times 25 cores from 100 sub-areas of  $1 \text{ m}^2$  to account for small-scale soil variability. Desaules et al. (2010) and Gubler et al. (2019) detailed the sample collection and data harmonization process of the measurements. The soils of the BDM were sampled at 0–0.2 m depth from positions on a regular grid of  $6 \text{ km} \times 4 \text{ km}$  laid over Switzerland (a total of 1094 locations). Each sampled location included four sub-samples that were taken at the intersection of the four cardinal directions from the center point and the circumference of a circle with a radius of



3 m to 3.5 m (Meuli et al., 2017). Due to its design covering all major geographic regions in Switzerland—the Jura mountain range, the Central Plateau, and the Alps—the BDM sampling campaign comprises a major part of the biogeochemical diversity of soils and predominant land use types in Switzerland. The wide coverage of soil conditions are an important source of soil chemical variability.

## 5 2.2 Chemical reference analysis

Data on chemical and physical soil properties were previously measured and provided by the NABO group. All laboratory soil analyses for the 17 properties were based on the protocols of the Swiss Standard Method (Agroscope, 1996). Mineral elements were determined by extraction with 1:10 ammonium acetate-EDTA solution (AAE10; method Agroscope 1996). The measured properties were total C, organic C (OC), total nitrogen (N), pH ( $\text{CaCl}_2$ ),  $\text{CaCO}_3$ , clay, silt, sand,  $\text{CEC}_{\text{pot}}$ , P(AAE),  
 10 K(AAE), Ca(AAE), Mg(AAE), Cu(AAE), Zn(AAE), and Fe(AAE). All chemical analyses of NABO soils were done on four bulked replicates per site and sampling event. For BDM locations, four spatial replicates were measured each.

## 2.3 Measuring and processing spectra

All milled soil samples from the NABO and the BDM archive ( $n = 4374$ ; with a particle size below  $100\mu\text{m}$ ) were measured with the Vertex 70 Fourier-transform spectrometer from Bruker (Bruker Optik GmbH, Ettlingen, Germany) at ETH Zurich,  
 15 using a high-throughput accessory (HTS-XT) and custom 24-well plates tailored to diffuse reflectance measurements. The mid-IR spectrometer was equipped with a KBr beamsplitter and a Mercury Cadmium Telluride (MCT) detector, which was permanently cooled with liquid nitrogen during the measurements. The reflectance spectra were acquired between  $7500\text{ cm}^{-1}$  ( $1333.3\text{ nm}$ ) and  $600\text{ cm}^{-1}$  ( $16666.7\text{ nm}$ ) at an effective resolution of  $2\text{ cm}^{-1}$ .

Each soil sample was measured twice. The soil surface was flattened evenly and without compression by the thin round  
 20 middle part of the spatula. The first measurement position of the 24-well plate contained a gold (Au) reference surface, which produced a single reflectance spectrum for normalizing the reflectance of the 23 following soil measurements. The "atmospheric compensation" routine implemented in the Bruker OPUS software was used to eliminate unwanted absorptions of  $\text{H}_2\text{O}$  vapour continuum and  $\text{CO}_2$  gas in the measurement chamber, based on the single channel reference spectrum measured once on each plate. All single channel reflectance spectra were obtained by averaging 32 internal measurements. The spectra were recorded  
 25 as  $\log(1/R)$ . Then, an average spectrum per sample was produced by calculating the mean of all spectral variables for the measured replicates. Finally, the spectrum offset and further scatter effects were reduced and the features were transformed with a Savitzky-Golay (Savitzky and Golay, 1964) first derivative smoother using a window size of 35 variables ( $70\text{ cm}^{-1}$ ) and third order polynomial fit. Finally, we selected every 8th spectral variable to reduce redundancy in the spectra (collinearity) and produce more parsimonious spectral estimates of soil properties. This resulted in 209 variables between  $634\text{ cm}^{-1}$  and  
 30  $3962\text{ cm}^{-1}$ , which formed the predictors for the subsequent general and local transfer modeling.



## 2.4 Data processing and statistical computing

All spectral and reference data were processed and modeled with the R software environment for statistical computing and graphics version 3.6.0 (R Core Team, 2019). We used the caret (Kuhn, 2020) R package to streamline the statistical learning process. Basic data transformations such as data preparation and aggregation were done using the tidyverse (Wickham, 2019) set of packages and data.table (Dowle and Srinivasan, 2019). The spectral data were handled and processed with the simplerspec (Baumann, 2019) and prospectr (Stevens and Ramirez-Lopez, 2013) packages.

## 2.5 General soil estimation: Rules for the entire SSL

The general soil estimation was done by rules trained with the CUBIST (Quinlan, 1993) learner, separately developed for each analytical soil measure. We chose this algorithm since it has shown excellent performance for SSLs with rather large soil variability and multicollinear spectral variables, and because its interpretation is mechanistically more intuitive as it is a form of data partitioning (simple conditions and linear equations). CUBIST first forms model trees using basic mechanisms of M5 (Quinlan, 1992). Wang and Witten (1996) outlined detailed principles behind the construction of the model trees and derivation of rules, and Viscarra Rossel and Webster (2012) described it for soil spectroscopic modeling.

### 2.5.1 Model development and validation

We tested a full-factorial combination of  $\{5, 10, 20, 50, 100\}$  committees of rules and  $\{2, 5, 7, 9\}$  neighbors to tune the CUBIST models. To get realistic estimates of the models' general performance, we defined a grouped ten-fold cross-validation scheme that treated the entire site (e.g., for total Carbon: NABO: 71 sites; BDM: 1079 sites) as independent in the modeling data sets. This made all observations from a site the unit of prediction, making the procedure equivalent to external cross-validation.

To reduce the bias-variance trade-off in the assessment, we repeated five times the grouped ten-fold cross-validation (CV) procedure (Friedman et al., 2008; Kuhn and Johnson, 2013). This was done with predefined random seeds to reproduce the division into training and validation proportions. We considered this site grouping factor as prior information when cross-validation segments were created, so that a particular site was only present within one segment (fold) of a cross-validation split. This grouped assignment prevented that the relationships were trained on the model fitting sets and prevented a particular site from leaking into the testing segments, yielding reliable generalization errors.

We tested the correspondence of mid-IR and model-derived predictions ( $x_i$ ) and measured standard reference measurements ( $\hat{x}_i$ ) with common regression metrics. We cross-validated the inaccuracy of the models with the root mean square error (RMSE). The mean squared error (MSE) was further decomposed into mean error (ME) or bias and the standard deviation of the error (SDE) or imprecision, so that  $RMSE^2 = ME^2 + SDE^2$  (Viscarra Rossel and McBratney, 1998). To describe the linear dependency between measurements and modeled values and give a relative goodness of fit, the coefficient of determination ( $R^2$ ) from linear regression is also reported. All metrics were aggregated from five estimates from independent resampling repeats. We reported mean values and standard deviations to provide uncertainties of the estimates.



## 2.5.2 Deriving important spectral variables

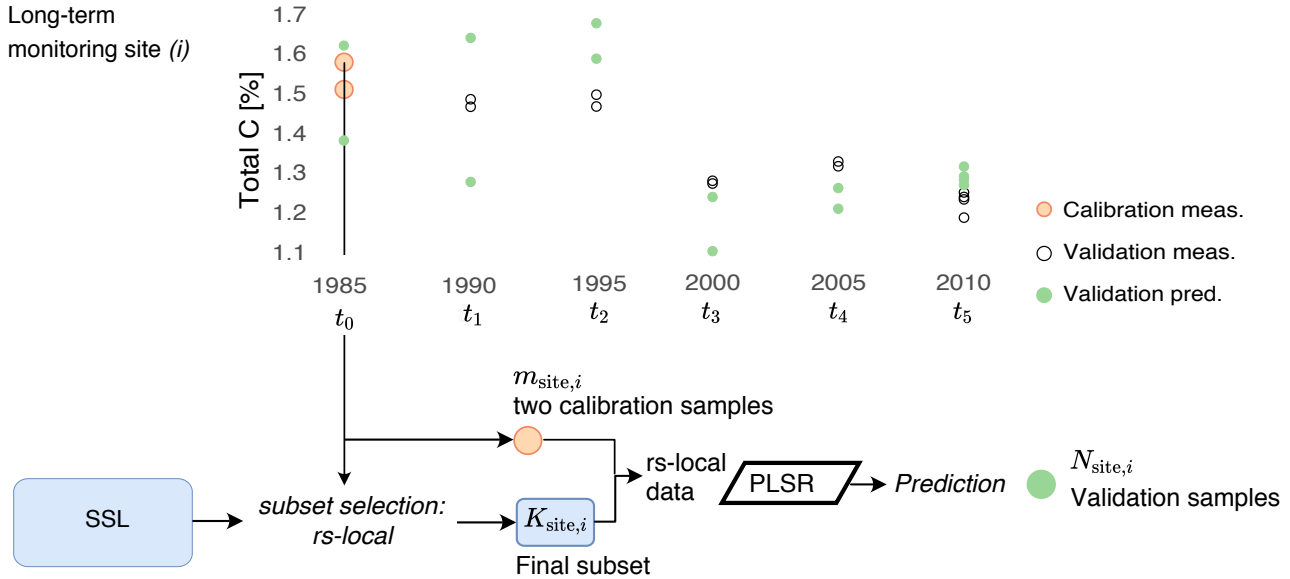
We used the recursive feature elimination (RFE) method, a backwards variable-selection algorithm described by Guyon et al. (2002), to test whether the modeling can be simplified and to find most important spectral features. Soil reflectance spectra typically contain many correlated and potentially redundant variables. Therefore, constraining them to relevant subsets that feed into the modeling can further improve predictive accuracy and reduce computation time and storage for model updates. We recursively eliminated subsets of variables with low CUBIST variable importance, calculated with equally weighted usage of a particular variable in split conditions and regressions. This step-wise variable reduction was based on the following predefined subset sizes  $S_i$ , starting with the full set at  $i = 1$  and ending with the most important predictor at  $i = 30$ :

$$S_i = \{209, 150, 120, 105, 90, 75, 60, 50, 40, 35, 30, 25, 20, 17, 14, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1\} \quad (1)$$

The dropped variables at each specific reduction step received identical importance ranks, from 30 (least important variables) to 1 (most important variable). Importance ranks were determined with step-wise variable reduction because model-based importance of a given input variable can substantially change when some correlated variables occur more frequently than others. Otherwise, using CUBIST importance measure on the entire spectrum would confound the importance of relevant but highly correlated variables. Since RFE is a wrapper method of variable selection, external test sets (resampling) was needed to exclude selection bias in estimating subset performance (RMSE) (Kuhn and Johnson, 2013). For this purpose, we nested another inner layer of resampling for RFE within the five times repeated 10-fold CV scheme. Importance ranks of variables and outer test RMSEs were averaged from the 50 CV folds. To decrease computation time, we conducted the RFE with 5 CUBIST committees. The RFE procedure and the resampling setup is explained further in the appendix 3.2.1.

## 2.6 Local soil estimation for plot-level monitoring

We defined a local soil estimation scenario where a new long-term monitoring site was initiated at time zero ( $t_0$ ). Each one of the 71 NABO sites was assumed to be novel while the remaining ones were established with spectral and reference data records. We therefore conducted 71 separate model transfers to test spectral-based soil monitoring using this SSL. We calibrated models at each site using two local samples per given site and **a relevant subset of the remaining Swiss SSL**. The two local samples were chosen from pooled samples at  $t_0$  (first two out of maximum four replicates), or in addition at  $t_1$  if there was only one sample in  $t_0$ . Figure 2 illustrates the local modeling workflow. All other samples per given site besides the two chosen during calibration (in other words the successive time-series measurements at a monitoring plot) were used as local validation samples ( $N_{\text{site},i}$ ). The respective samples from the remaining SSL included spectra and reference measurements from all BDM samples and NABO samples, excluding the ones from the respective target site. We used only two calibration samples per NABO site to capture the predictive mechanisms at site-level because we wanted to avoid overoptimistic local assessment; both local calibration and validation samples were repeated soil measurements, and are otherwise — if not adequately handled in the calibration sampling strategy — at risk of information leakage when soils' composition and relevant properties show constant trends over time.



**Figure 2.** Scheme illustrating the transfer of the soil spectral library (SSL) to a long-term monitoring site using the RS-LOCAL approach. The local calibration samples and a subset of the SSL are used to calibrate a partial least squares regression (PLSR) model, which predicts the local validation samples.

For each of the 71 sites, the spectral relevant samples from the remaining Swiss SSL were selected using the RS-LOCAL algorithm (see Lobsey et al. (2017)). The site-specific samples ( $m_{\text{site},i}$ ) denote local calibration samples from a NABO plot. The recursive reductions of the initial training data, which determined the finally yielded subsets ( $K_{\text{site},i}$ ) were driven by model performance (RMSE) for the two local calibration samples. For each NABO site  $i$ , the corresponding  $K_{\text{site},i}$  set was

5 spiked with the two local calibration samples. On this combined  $m_{\text{site},i} + K_{\text{site},i}$  data set, a final PLSR model, locally adapted for the monitoring plot by optimization on the calibration samples, was developed using 10-fold cross-validation. Finally, the local validation spectra ( $N_{\text{site},i}$ ) were predicted using the most accurate calibration model.

The search algorithm RS-LOCAL has three empirical parameters to control the samples that are selected for the local transfer from the SSL (Lobsey et al., 2017). Parameter  $k$  is both the number of samples drawn from the original and reduced library

10 without replacement and the number of samples of the returned SSL subset. Parameter  $b$  is the number of times  $k$  samples are randomly drawn from the remaining data at iteration  $i$  of the performance-driven library reduction. Parameter  $r$  is the proportion of samples, which are consistently in weakest models, that are removed at it each reduction step. The configuration of the RS-LOCAL search was optimized for each NABO site. For each site, we ran separate RS-LOCAL runs, testing a full-factorial combination of empirical parameter sets  $k = \{30, 50, 150\}$ ,  $b = \{10, 20, 50\}$ ,  $r = \{0.05, 0.1, 0.2\}$ . The finally selected

15 optimal subset per site yielded the smallest RMSE on the two local calibration samples, and was therefore used to infer predictions for the local validation samples.



### 2.6.1 Uncertainty of spectral monitoring uncertainty: CUBIST vs. RS-LOCAL transfer

To compare performance of the CUBIST and RS-LOCAL transfer, errors and concordance of both methods were conditionally assessed per individual NABO ( $n = 71$ ) site. For CUBIST, grouped cross-validation hold-outs were used. Thereby, the two local validation samples  $N_{\text{site}, i}$  were excluded, so that the test configuration was identical to the local transfer scenario. In addition to the mentioned assessment statistics, the ratio of performance to interquartile distance (Bellon-Maurel et al., 2010) (RPIQ; 75th and 25th percentiles) was used for relative comparisons between the local transfer and rule-based model because it is robust to non-normal and skewed distributions of measured values.

### 2.6.2 Evaluating the predictive mechanisms behind the local transfer

For each of the 71 statistical transfers at a plot-level, we quantified the similarity between the selected data sources  $K_i$  (from SSL) and the respective local target domain  $\{m_i; N_i\}$  (local calibration and validation) by multivariate distances across the spectral input variables. The distance of single observations within  $\{K_i; m_i\}$  was referenced to the center of all data, which lead to two respective distributions of distance measures for these sets of points and per site. This procedure involved two steps: 1) compress the input data to reduce the "curse of dimensionality" (Bellman, 1961) and be able to discriminate similarity with spectra (with many dimensions, distance to nearest neighbor becomes similar to distance to farthest neighbor); 2) calculate Mahalanobis distance using a robust method (see below; Varmuza and Filzmoser (2016)), so that the location and scatter were influenced by the main data rather than by atypical observations.

To condense the spectral information over the entire SSL, Savitzky-Golay preprocessed spectra that included all observations with carbon elemental measurements were mean-centered, scaled, and then transformed by PCA using singular value decomposition. Dimensionality reduction was necessary to avoid computationally singular values during the subsequent calculation of the covariance matrix (for the Mahalanobis distance). The first ten principal components that explained 86.5 % of the variation in preprocessed spectra were kept for distance calculations. Finally, the Mahalanobis distance of all the observations to their center was computed with robust estimates for both the center and the covariance matrix of the selected PCA scores, using the Minimum Covariance Determinant (MCD) estimator (Rousseeuw, 1984; Hubert and Debruyne, 2010).





### 3 Results

#### 3.1 Summary of reference measurements

The samples from the Swiss soil monitoring network (NABO) exhibited the highest variability across samples for total C and OC ( $n = 592$ ; Table 1). Organic C ranged from 1.1 % to 27.3 %. The texture of the soils varied considerably. The pH had values between 3.5 and 7.6 and the soils were slightly acidic overall with a median of 5.8. Compared to the NABO data set, the soils from the BDM program covered a wider set ( $n = 3723$  for total C) and range of measured soil properties. The measured range of total C for BDM (0.1 – 58.3 %) extended further than that for the NABO. The distribution of pH values was similar in the NABO and BDM sets. The BDM data included also the available cations extracted by AAE (see Table 1). The median  $\text{CEC}_{\text{pot}}$  was almost equivalent to the value of the NABO sites (24 vs. 23  $\text{cmol}(+) \text{kg}^{-1}$ ). Exchangeable Ca showed the largest coefficient of variation ( $\text{CV} = 1.56$ ) among the measured properties of the BDM set.

#### 3.2 General soil estimation with CUBIST modeling

For most of the properties, minimal cross-validated errors were achieved with 100 committees and 9 neighbors. The rule-based models explained a large proportion of the variation ( $R^2 > 0.9$ ) in properties that typically have a strong link to total C (humus, organic C, N) (Table 3; Figure 3). Clay was accurately estimated ( $\text{RMSE} = 5\%$ ; range = 0–60 %), whereas sand and silt were less accurately estimated. The pH was accurately estimated ( $\text{RMSE} = 0.3$ ). Our models discriminated a large proportion in the measured variation of Ca and Mg (ammonium acetate-EDTA) in the mid-IR ( $R^2 = 0.97$  and  $0.78$ ). Reference values of potential cation exchange capacity ranged from 0–136  $\text{cmol}(+) \text{kg}^{-1}$  and were modeled with an  $\text{RMSE}$  of 7  $\text{cmol}(+) \text{kg}^{-1}$  ( $R^2 = 0.70$ ). However, the extractable nutrients P, K, Cu and Zn were insufficiently explained by mid-infrared spectral rules. However, the rules achieved nearly unbiased property estimates over all observations. We found marginal local bias at the largest values, mostly for variables with positively skewed distributions such as total C.

Overall, out of the 17 available soil properties, total C, total N, total  $\text{CaCO}_3$ , Ca and Mg (ammonium acetate-EDTA), Organic Matter, OC,  $\text{CEC}_{\text{pot}}$ , pH, sand and clay (11) were modeled with relatively good discrimination capacity in the measured ranges (Figure 3).

##### 3.2.1 Model interpretation and filtering with variable importance

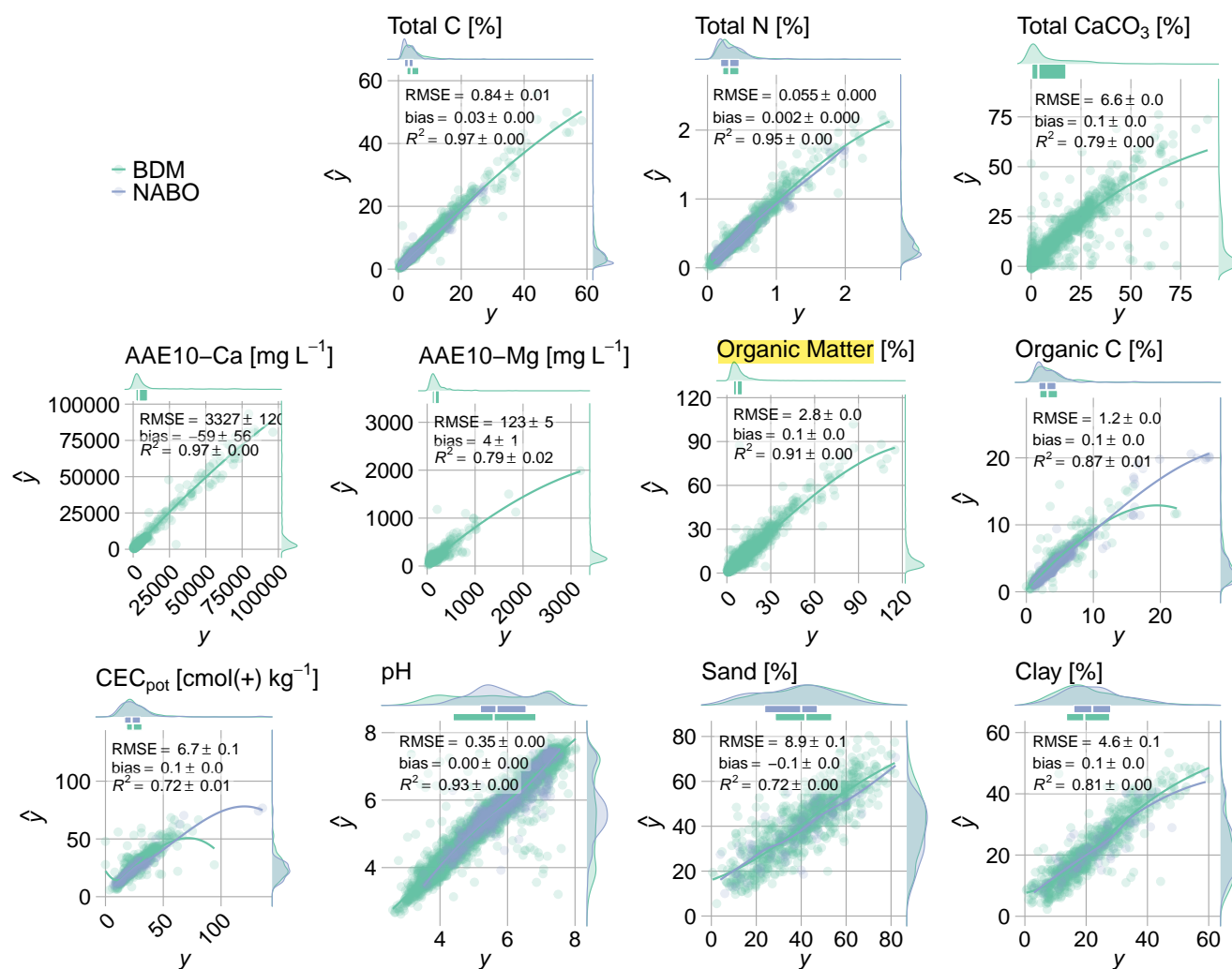
Figure 4 shows that the test  $\text{RMSE}$  overall increased with less spectral variables using CUBIST. The lowest error ( $\text{RMSE}_{\text{test}} = 0.834\%$  for total C) of spectroscopic estimation was achieved with the initial modeling spectra (209 variables). The predictive capacity for total C across all external cross-validations even increased when we reduced the number of spectral input variables further to 105 and 90 variables ( $\text{RMSE}_{\text{test}} = 0.810\%$  and  $0.820\%$  for total C, respectively). For the subsequent variable reduction steps, model performance steadily dropped until one wavenumber was left ( $\text{RMSE}_{\text{test}} = 1.88\%$  C).

The spectral feature between  $1786 \text{ cm}^{-1}$  and  $1754 \text{ cm}^{-1}$  was the most important one for the estimation of total C with CUBIST (Figure 4). The twelve spectral variables with the best importance ranks across all RFE iterations and test sets de-

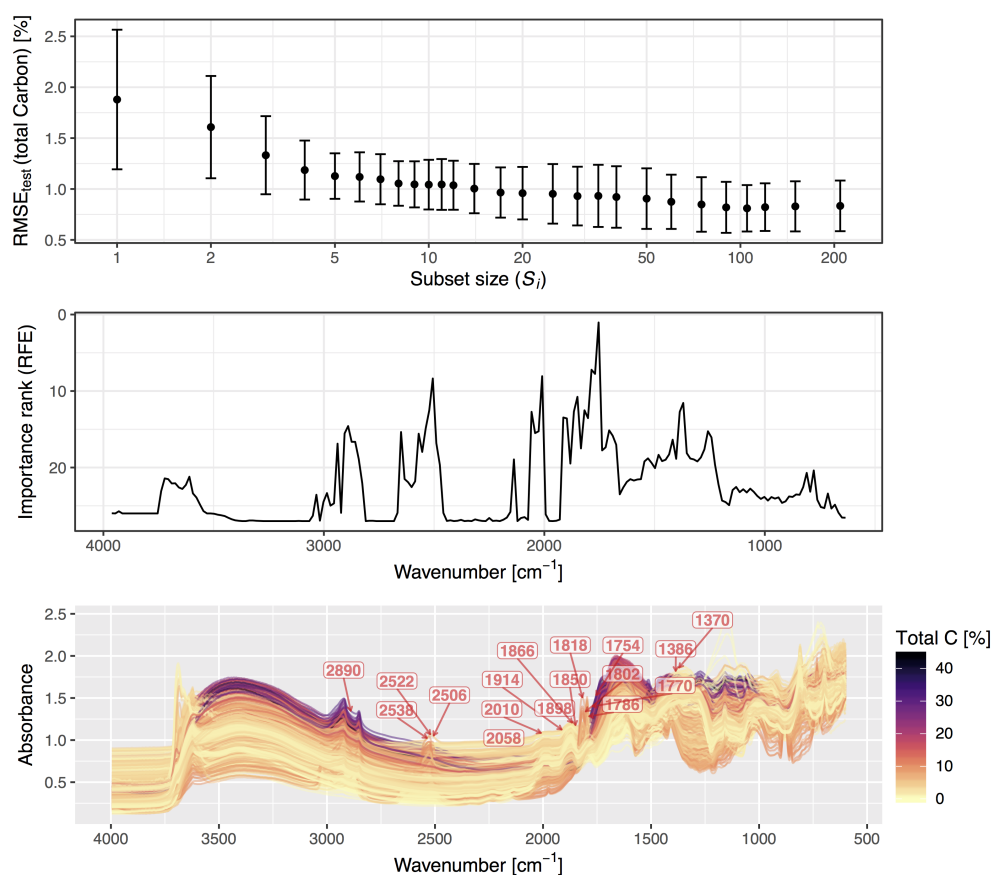


**Table 1.** Summary statistics of the measured soil properties of the Swiss soil spectral library derived from the sample archive of the Swiss soil monitoring network (NABO) and the Swiss Biodiversity Monitoring (BDM) program. C = total carbon, OC = organic carbon, CEC<sub>pot</sub> = potential cation exchange capacity. Extractable element contents were based on 1:10 ammonium acetate-EDTA (AAE) solution. All measured values were referenced to dry weight by water correction after drying at 105 °C.

	<i>n</i>	#Locations	Min	Max	Median	Mean	SD	CV
<b>NABO</b>								
C [%]	572	71	1.1	27.3	3.3	4.0	3.5	0.88
OC [%]	592	71	1.1	27.3	3.0	3.7	3.4	0.92
N [%]	572	71	0.11	1.99	0.32	0.36	0.26	0.71
pH	574	71	3.5	7.6	5.7	5.8	0.9	0.16
Clay [%]	80	55	3	59	22	23	10	0.45
Silt [%]	81	55	15	80	38	40	13	0.32
Sand [%]	80	55	4	82	40	37	17	0.45
CEC <sub>pot</sub> [cmol(+) kg <sup>-1</sup> ]	121	58	7	136	23	26	17	0.68
<b>BDM</b>								
C [%]	3723	1079	0.1	58.3	4.1	5.5	4.9	0.89
OC [%]	498	472	0.0	22.5	3.2	3.9	2.8	0.72
Humus [%]	3664	1073	0	115	7	9	9	1.01
N [%]	3724	1079	−0.00	2.64	0.32	0.38	0.25	0.67
pH	3765	1094	2.6	8.0	5.6	5.6	1.3	0.24
CaCO <sub>3</sub> [%]	1565	455	−0.2	88.5	3.6	10.7	14.4	1.35
Clay [%]	787	785	0	60	19	21	11	0.51
Silt [%]	787	785	10	71	30	31	10	0.31
Sand [%]	787	785	1	82	42	41	17	0.41
CEC <sub>pot</sub> [cmol(+) kg <sup>-1</sup> ]	674	190	0	94	24	26	12	0.44
P (AAE) [mg kg <sup>-1</sup> ]	417	417	1	1047	19	40	77	1.92
K (AAE) [mg kg <sup>-1</sup> ]	417	417	22	1255	106	136	115	0.84
Ca (AAE) [mg kg <sup>-1</sup> ]	417	417	141	96250	3927	12226	19127	1.56
Mg (AAE) [mg kg <sup>-1</sup> ]	417	417	16	3196	161	232	259	1.11
Cu (AAE) [mg kg <sup>-1</sup> ]	417	417	2	73	6	8	5	0.73
Zn (AAE) [mg kg <sup>-1</sup> ]	417	417	1	131	4	6	9	1.43
Fe (AAE) [mg kg <sup>-1</sup> ]	417	417	84	1640	342	387	194	0.50
DNA [μg g <sup>-1</sup> ]	718	225	2	155	15	20	17	0.87



**Figure 3.** Agreement between measured and mid-infrared predicted values that was obtained from CUBIST models. Models' performance was assessed by site grouped cross-validation holdouts (five times repeated ten-fold). Models of soil properties with  $R^2 > 0.6$  are shown (see Table 3 for more detailed model summaries).



**Figure 4.** *Top:* Root mean square error (RMSE) of mid-infrared estimates of total C that CUBIST produced at the respective subsets of spectral variables. The performance profile was obtained with a recursive feature elimination (RFE) procedure using predefined subset sizes, which was embedded into a nested and by site grouped cross-validation scheme that excluded selection bias and ensured generalization of the estimated performance of variable filtering. The 50 repeated sets of outer holdouts gave uncertainties of the procedure (error bars represent standard deviations of the test RMSE). *Middle:* Average importance ranks across the spectrum. Lower rank values indicate higher importance for the estimation of total C. Ranks were determined with RFE. *Bottom:* Mid-infrared absorbance spectra of the Swiss soil spectral library ( $n = 4295$ ; with corresponding total carbon (C) measurements determined by dry combustion). The spectra are annotated with the 17 most influential spectral variables (wavenumbers) in the CUBIST model (average importance rank  $< 15$ ); these had the highest mean importance ranking determined by the recursive feature elimination procedure.



**Table 3.** Cross-validated mid-IR estimates of 18 soil properties derived from rule-based CUBIST models that were developed using all available data. The 10-fold cross-validation procedure was grouped by the site and repeated 5 times to achieve many test-train data combinations and provide realistic model assessment with generalization. #C and #N denote the number of committees and neighbors for CUBIST. C = total carbon, OC = organic carbon, N = total nitrogen, and CEC<sub>pot</sub> = potential cation exchange capacity. Elemental nutrients were extracted with 1:10 ammonium acetate-EDTA (AAE) solution.

	<i>n</i>	Range	#C	#N	#Rules	RMSE	ME	SDE	<i>R</i> <sup>2</sup>
Total C [%]	4295	0.1–58.3	100	9	6–26	0.84 ± 0.01	0.03 ± 0.00	0.84 ± 0.01	0.97 ± 0.00
OC [%]	1090	0.0–27.3	100	9	3–14	1.2 ± 0.0	0.1 ± 0.0	1.2 ± 0.0	0.87 ± 0.01
Humus [%]	3664	0–115	100	9	3–12	3 ± 0	0 ± 0	3 ± 0	0.91 ± 0.00
N [%]	4296	−0.00–2.64	100	9	9–21	0.06 ± 0.00	0.00 ± 0.00	0.05 ± 0.00	0.95 ± 0.00
pH	4339	1.0–8.0	100	9	4–18	0.3 ± 0.0	0.0 ± 0.0	0.3 ± 0.0	0.93 ± 0.00
CaCO <sub>3</sub> [%]	1565	−0.2–88.5	50	9	1–5	6.6 ± 0.0	0.1 ± 0.0	6.6 ± 0.0	0.79 ± 0.00
Clay [%]	867	0–60	100	9	2–6	5 ± 0	0 ± 0	5 ± 0	0.81 ± 0.00
Silt [%]	868	1–80	100	9	1–14	7 ± 0	−0 ± 0	7 ± 0	0.51 ± 0.01
Sand [%]	867	1–82	100	9	1–6	9 ± 0	−0 ± 0	9 ± 0	0.72 ± 0.00
CEC <sub>pot</sub> [cmol(+) kg <sup>−1</sup> ]	795	0–136	50	9	1–10	7 ± 0	0 ± 0	7 ± 0	0.72 ± 0.01
P (AAE) [mg kg <sup>−1</sup> ]	417	1–1047	50	9	1–8	77 ± 1	−3 ± 1	77 ± 1	0.05 ± 0.00
K (AAE) [mg kg <sup>−1</sup> ]	417	1–1255	100	9	1–12	111 ± 2	−3 ± 1	111 ± 2	0.10 ± 0.01
Ca (AAE) [mg kg <sup>−1</sup> ]	417	1–96250	100	9	6–14	3327 ± 120	−59 ± 56	3326 ± 120	0.97 ± 0.00
Mg (AAE) [mg kg <sup>−1</sup> ]	417	1–3196	100	9	1–13	123 ± 5	4 ± 1	123 ± 5	0.79 ± 0.02
Cu (AAE) [mg kg <sup>−1</sup> ]	417	1–73	50	9	1–7	5 ± 0	−0 ± 0	5 ± 0	0.10 ± 0.01
Zn (AAE) [mg kg <sup>−1</sup> ]	417	1–131	100	9	1–13	9 ± 0	−0 ± 0	9 ± 0	0.06 ± 0.01
Fe (AAE) [mg kg <sup>−1</sup> ]	417	1–1640	50	5	1–7	167 ± 3	1 ± 1	167 ± 3	0.28 ± 0.02

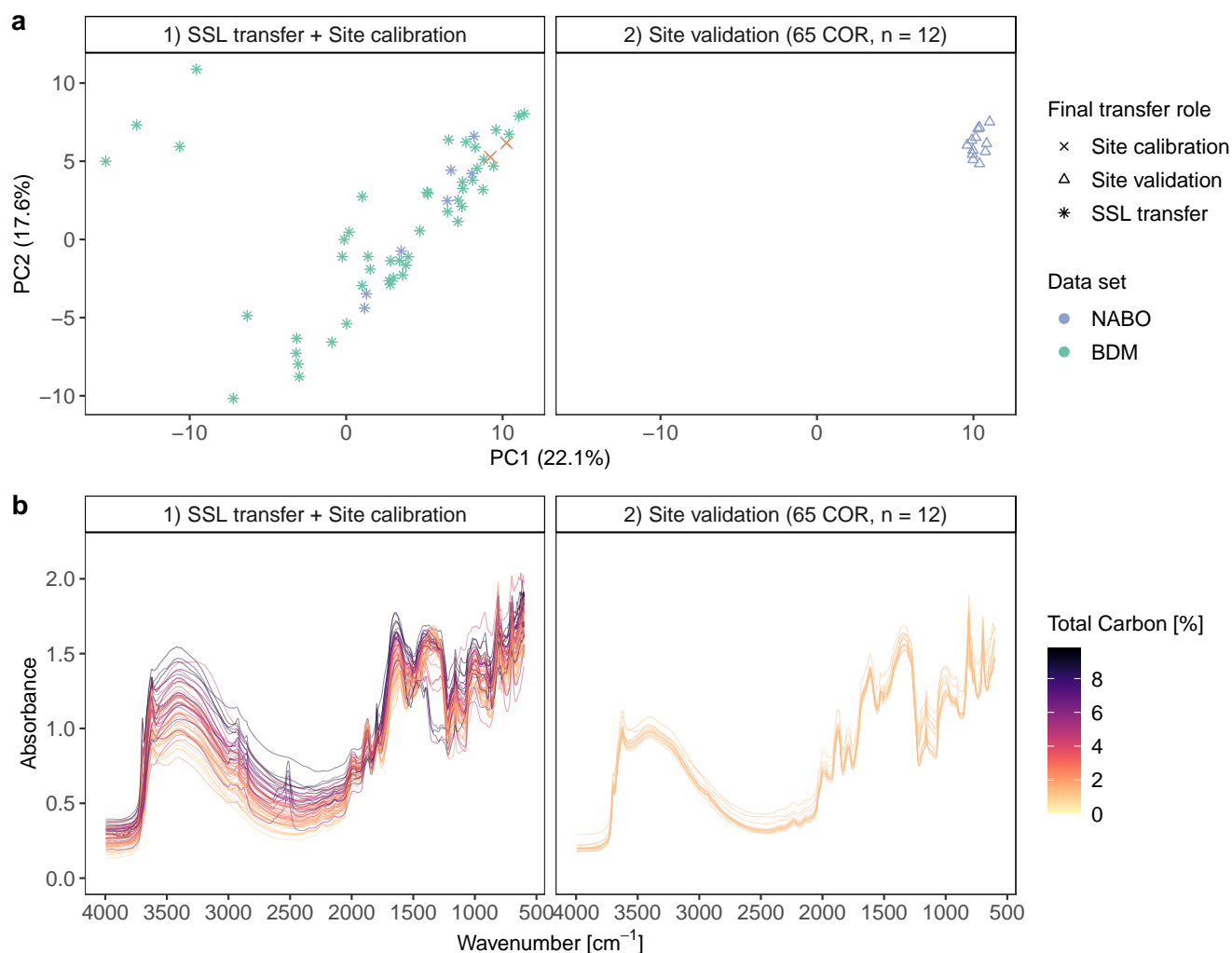
rived from the subset sizes were (starting with best): 1754 cm<sup>−1</sup> (mean(rank) = 1.04), 1786 cm<sup>−1</sup>, 1770 cm<sup>−1</sup>, 2010 cm<sup>−1</sup>, 2506 cm<sup>−1</sup>, 1850 cm<sup>−1</sup>, 1370 cm<sup>−1</sup>, 2522 cm<sup>−1</sup>, 1818 cm<sup>−1</sup>, 1866 cm<sup>−1</sup>, 2058 cm<sup>−1</sup>, 1386 cm<sup>−1</sup> (mean(rank) = 12.7; Figure 4).

### 3.3 Accuracy of the local transfer models compared to the general model

- For the example site 65 COR, the best performance of RS-LOCAL was achieved with 55 samples from the SSL (*K*), 10 sampling events (*B*) of size *K* at each iteration, and 10 % reduction (*r*) at each iteration (Figure 5). This effectively yielded 52 transfer samples from the SSL that were combined with two site calibration samples previously used to supervise the selection from the data source, to form a PLS regression calibration model for the estimation of the site validation samples (Figure 5 panel *a* right). Compared to the target observations from the site (right part of panels *a* and *b*; measured range = 1.19–1.60 % C), the selected



instances were heterogeneous with regard to their spectral appearance, their input feature space, and their measurements (range = 0.87–9.77 % C). The selected instances covered a significant proportion of the first two components in the feature space of the entire SSL.



**Figure 5.** Illustration of the site-specific transfer modeling of total carbon (C), using RS-LOCAL for the example site 65 COR of the Swiss soil monitoring network (NABO). Panel *a* contains the principal components subspace (PC1 and PC2) of the Savitzky-Golay first derivative mid-IR spectra, and panel *b* outlines the corresponding absorbance spectra, which are coloured by the total C content. The left subplots show the SSL transfer samples ( $n = 55$ ) that were selected from the soil spectral library ( $n = 4281$ ; excluding all NABO calibration samples). This subset was most accurate when predicting the two calibration samples under the mechanisms RS-LOCAL and their optimal tuning configuration for the site ( $\{K = 50; B = 10; r = 0.1\}$ ). The right panels shows the time series data for the validation samples of the NABO site "65 COR".





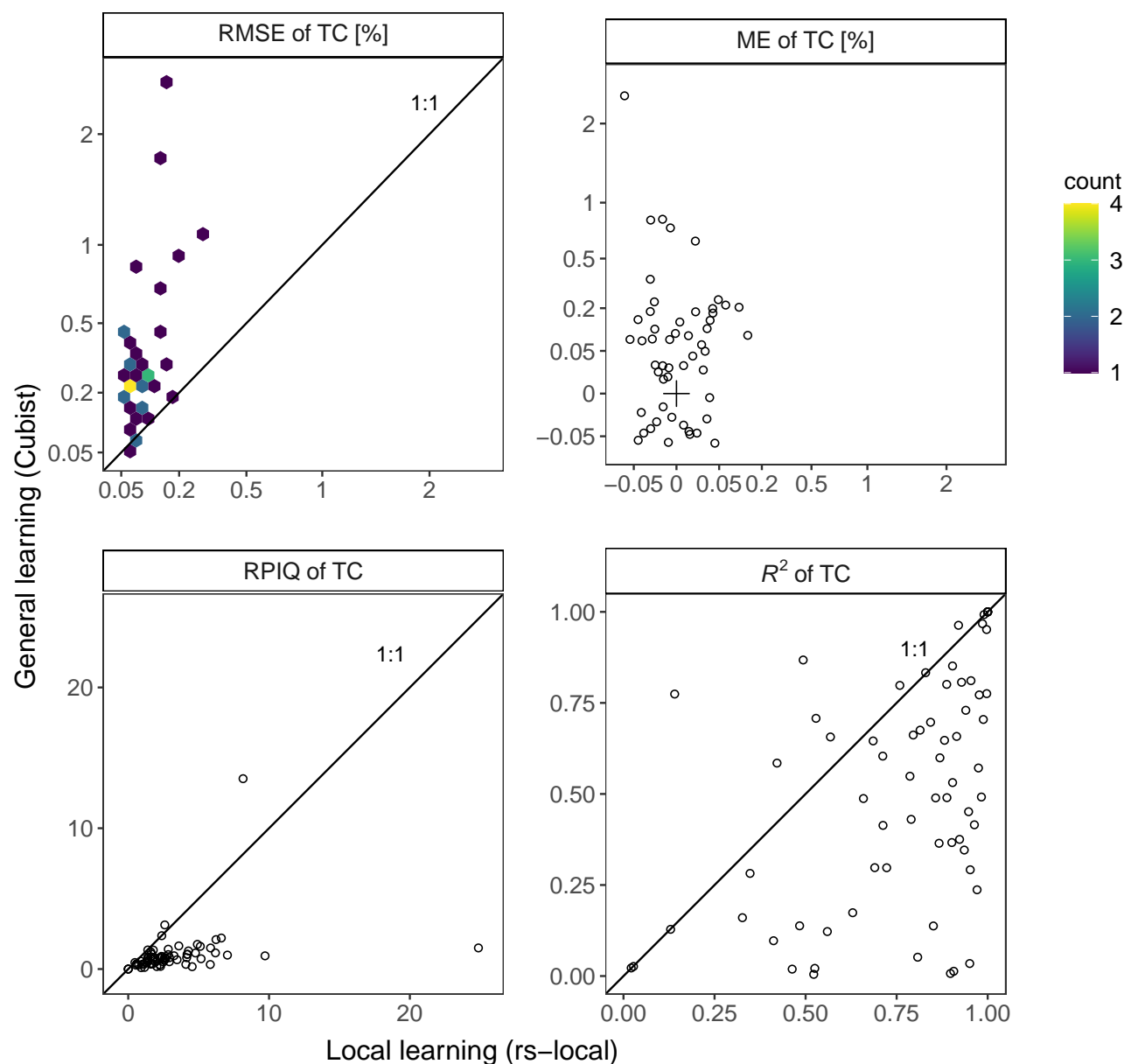
The RMSE on the site validation samples ( $RMSE_{N_i}$ ) at the final subsets varied between 0.001 and 1.073 % C and for all tuning parameter combinations and sites, and between 0.001 and 0.302 % C for the best subsets per site (Figure A1).

The local approach reduced the error of the rule-based approach on average by factor 4.4 (Figure 6;  $\text{mean}(RMSE_{rs-local}) = 0.07\% \text{ C}$ ;  $\text{mean}(RMSE_{cubist}) = 0.31\% \text{ C}$ ). The local transfer was more accurate for the majority of NABO sites (69 out of 71 sites). The linear dependency between modeled and measured values was higher for the local transfer compared to the general model (53 out of 71 sites). Moreover, RS-LOCAL produced on average 1.3 times less biased estimates of total C per site for 52 out of 69 sites in terms of absolute values ( $|ME| = 0.01\% \text{ C}$  vs.  $0.05\% \text{ C}$ ). The ratio of performance to inter-quartile distance (RPIQ) confirmed that local learning in the mid-infrared was able to better discriminate developments of total C over time, relative to its measured distribution. Overall, local learning with two local calibration samples and targeted SSL selections allowed for better estimations than the generic CUBIST (RPIQ = 3.08 vs. 1.00; RPIQ larger for 66 out of 71 sites).

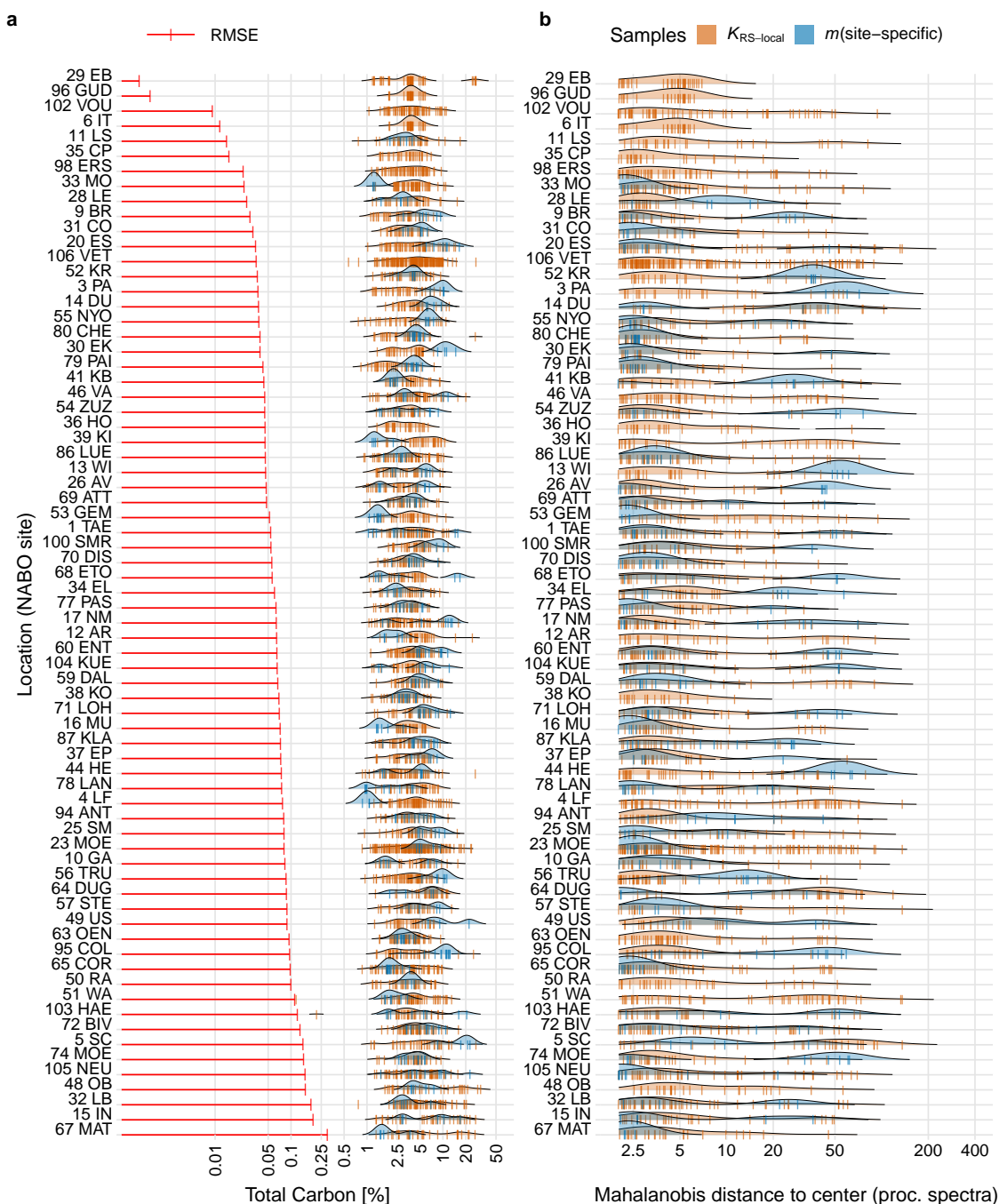
### 3.3.1 Predictive mechanisms behind the local transfer

The samples used for the transfer process (RS-LOCAL data) of the example site *COR 65* showed high spectral dissimilarity along the first 2 PCs, explaining 39.8 % of preprocessed spectral variance (Figure 5). Compared to the entire SSL with total carbon (C) measurements available (the source domain prior selection; range PC1: -41.4 to 13.0; range PC2: -19.0 to 30.0), the selected transfer samples of this site occupied a region of major variation in the PC space (range PC1: -15.4 to 11.4; range PC2: -10.2 to 10.9). The two local calibration samples and the 12 validation samples on the upper right corner were close to each other in the PC1-PC2 subspace (Figure 5, panel a, left and right; range PC1: 9.2 to 11.0; range PC2: 4.9 to 7.5). Not only the absorbance spectra but also the corresponding C reference values were highly variable compared to the exemplary NABO site (Figure 5, panel b; 0.73–11.78 % C for  $K_{rs-local}$ , and 1.19–1.60 % C for the plot of this site). This particular target monitoring site indicated that RS-LOCAL selected soils from the SSL with a relatively large spectral diversity and a wide range of total C.

The instances selected by RS-LOCAL filled a substantial proportion of the SSL's feature space (Figure 7), confirming the trend of site *65 COR*. We found that RS-LOCAL yielded a quite wide selection of relevant samples from the SSL with reference to both the total C range and a wide coverage of spectral features expressed with robust multivariate locations. The spectral estimations of the site validation sets that resulted from RS-LOCAL-based transfers did neither show trends in the mode or spread for distributions of C measurements nor in the ones from their spectral distances. The measured distributions of  $K_{i,site}$  SSL subsets and  $N_{i,site}$  local validation samples for further key soil properties related to the chemical composition (OC, pH,  $CEC_{pot}$ , clay and  $CaCO_3$ ) were also markedly different, confirming the local transfer of quite heterogeneous soils (Table 5). For example, standard deviations of the 0 %, 25 %, 50 %, and 75 % percentile differences between the transfer sets selected the SSL and the samples from the respective NABO site were on average between 1.8 % and 6.6 % for measured C and OC. Further, the measured clay and  $CaCO_3$  contents was markedly different between the RS-LOCAL selection and the local validation sets ( $> 6.8\%$ ). This findings correspond with the dissimilar selection compared to the local target samples found in the PCA space of preprocessed spectra.



**Figure 6.** Model assessment of the estimated total carbon (TC) of 71 NABO sites for the general learning with CUBIST (y-axis) vs. local learning transfer with RS-LOCAL (x-axis). The four panels depict the root-mean-square-error (RMSE), the mean error (ME), the ratio of performance to interquartile distance (RPIQ) and  $R^2$ . The 1:1-line emphasizes the difference between the two approaches.



**Figure 7.** Analyzing the mechanisms behind the individual adaptive transfer realized with RS-LOCAL. Panel *a*: The left horizontal bars show the root mean square error (RMSE) of mid-infrared predictions of the temporal validation set of time series of total Carbon (C) for each of the 71 NABO sites, which was calculated without the two respective calibration samples. The blue density plots depict the distribution of the site-specific validation samples, and the brown vertical bars show the measured values of C for the final subsets of SSL used for the transfer ( $K$ ). Panel *b*: The distribution of the robust distances from the PCA center of Savitzky-Golay preprocessed spectra of the entire soil spectral library compared to the subset of instances involved in the individual transfer modeling (similarity in site-specific vs. final RS-LOCAL selection), computed with the Minimum Covariance Determinant (MCD) estimator.



**Table 5.** Standard deviations (SD) of the absolute differences of percentiles ( $P_0, P_{25}, P_{50}, P_{75}, P_{100}$ ) of final RS-LOCAL subsets ( $K_{i,\text{site}}$ ) and corresponding site validation samples ( $N_{i,\text{site}}$ ) the across 71 long-term monitoring sites. The aggregated values for six measured soil properties are shown. C = total Carbon, OC = organic Carbon,  $\text{CEC}_{\text{pot}}$  = potential cation exchange capacity.

	$\text{SD}\left( P_X(K_{i,\text{site}}) - P_X(N_{i,\text{site}}) \right)$				
	$\text{SD}\left( \Delta P_0 \right)$	$\text{SD}\left( \Delta P_{25} \right)$	$\text{SD}\left( \Delta P_{50} \right)$	$\text{SD}\left( \Delta P_{75} \right)$	$\text{SD}\left( \Delta P_{100} \right)$
C [%]	1.8	2.2	2.5	3.1	6.1
OC [%]	2.5	2.3	2.2	2.2	6.6
pH	0.8	0.7	0.5	0.6	0.9
$\text{CEC}_{\text{pot}}$ [cmol(+) kg <sup>-1</sup> ]	9.8	10.3	10.6	10.6	23.1
Clay [%]	11.0	9.7	8.5	6.8	9.0
CaCO <sub>3</sub> [%]	7.4	7.9	8.9	9.8	17.8

## 4 Discussion

### 4.1 General soil estimation with the Swiss SSL

Many of the chemical properties with distinct links to soil organic matter and the key minerals (e.g., clays and quartz) were discriminated well with mid-IR CUBIST models (Table 3; Figure 3). Specifically, the models estimated total C, OC, N, pH, texture, AAE10-Ca, and AAE10-Mg with  $R^2 > 0.7$ . This suggests that the majority of developed models are useful for applications that require soil proxies in order to manage land resources. For example,  $\text{CEC}_{\text{pot}}$  (RMSE = 7.0 cmol(+) kg<sup>-1</sup>), as well as pH, have high ecological importance for nutrient availability in ecosystems. In agriculture, both measures are key factors for fertilization recommendations.

The accuracy of our estimates for the properties that have direct chemical links, through compound-associated absorptions, were mostly comparable to established continental or country-specific mid-IR SSLs. For example, Clairotte et al. (2016) achieved RMSE = 0.2 % for OC using mid-IR and the spectrum-based learner for local predictions, while Sila et al. (2016) reported RMSE = 0.4 %. The accuracy of our general OC estimates was lower (RMSE = 1.2 %), which we explain by the relatively large range of measured values and variable mineralogy (Stenberg and Rossel (2010)). We found that total C had more CUBIST rules per committee than OC (Table 3), indicating total C leverages more chemical constituents and latent absorptions for its estimation. Further, for total C, which also included inorganic C (mostly CaCO<sub>3</sub>), we had about four times more training data than OC and a higher soil diversity (1079 vs. 472 BDM sites). These differences might in addition explain why we yielded higher accuracy for total C (RMSE = 0.84 %).

The variable importance assessment of the spectroscopic models revealed five major regions of features with particularly high predictive influence for total C: 2890 cm<sup>-1</sup>, 2522 cm<sup>-1</sup>, 2010 cm<sup>-1</sup>, 1754 cm<sup>-1</sup>, and 1370 cm<sup>-1</sup> (Figure 4). We attribute the



two absorption peaks near  $2890\text{ cm}^{-1}$  to C–H stretching vibrations of organic matter (Skjemstad and Dalal, 1987), which were also relatively important for estimating C in other studies (e.g., Janik and Skjemstad (1995); Viscarra Rossel and McBratney (1998)). The important variable at  $2522\text{ cm}^{-1}$  is indicative of C=O absorption due to the carbonyl group present in carbonates (e.g., calcite) (Nguyen et al., 1991; Soriano-Disla et al., 2014). The three important absorptions between  $2010\text{ cm}^{-1}$  and  $1786\text{ cm}^{-1}$  result from three consecutive Si–O–Si (overtone and combination) absorptions, which are indicative of quartz. However, the most important absorptions near  $1754\text{ cm}^{-1}$  showed no distinct peak but an edge feature. This is in accordance with Sila et al. (2016), which identified this region as being most relevant for estimating total C with a general model (random forest) developed from the SSL of the African Soil Information project. This region is close to the C=O stretching vibration of the carboxyl group that occurs around  $1725\text{--}1720\text{ cm}^{-1}$  (Madari et al., 2006), which is further confirmed by the high importance of these vibrations found by Janik and Skjemstad (1995). The last relatively important region around  $1370\text{ cm}^{-1}$  was also an edge feature with no distinctly visible peak of chemical group assigned, which, however, might be influenced by the adjacent carboxylate ( $\text{COO}^-$ ) or  $\text{C}\text{--}\text{H}$  absorptions at  $1400\text{--}1350\text{ cm}^{-1}$  of aliphatic compounds such as humic acids (Madari et al., 2006; Parikh et al., 2014). In summary, the CUBIST-RFE variable importance analysis enabled us to link characteristic absorptions of typically prominent functional groups of soil organic and inorganic C compounds, and as well quartz absorptions as indirect correlative features of predictive relevance, with our general-model based estimates of total C.

Since the rule-based models we developed can estimate 11 soil properties (Figure 3), the Swiss SSL will be useful for new soils when new reference measurements for model adaptation are relatively scarce or not available. Thereby, the Swiss SSL will be cost and time efficient for characterizing soils of similar composition in the near future. The new predictions can further be augmented with straightforward model interpretation, which allows chemical inference of pedological aspects to provide means of model applicability. Although the combined BDM and NABO set comprises a large soil variability in Switzerland, the diversity of subsoils at depths greater than  $20\text{ cm}$  — mostly in terms of the mineral composition — as well as peat and forest soils are probably not yet represented sufficiently in the SSL. We therefore must continuously update the present SSL with more and deeper soil horizons in the near future.

#### 4.2 Local transfer from the SSL for soil monitoring at plot-scale

The local estimates of total C that were derived with RS-LOCAL selection were substantially better on average ( $\text{RMSE} = 0.07\text{ \% C}$ ) as those derived using all of the data and general CUBIST models ( $\text{RMSE} = 0.31\text{ \% C}$ ; Figure 6). The data-driven estimation at plot-scale further considerably reduced bias and increased  $R^2$  compared to the general CUBIST.

Our third goal was to analyze the characteristics of soils that were selected from the SSL and used for establishing locally-adaptive models tailored to the respective long-term monitoring sites. Surprisingly, the RS-LOCAL subsets selected from the SSL had rather dissimilar spectra in the robust PCA space (Figure 5; Figure 7); their distances to the center had a wide distribution compared to the local samples. The  $K_{\text{site},i}$  subsets accordingly covered a large proportion of the spectral input space. The likely dissimilar chemical composition of soils was also reflected in the reference measurements of total C. We conducted a broader analysis to interpret the soil context of the selected samples with further soil compositional covariates (OC, pH,  $\text{CEC}_{\text{pot}}$ , clay,  $\text{CaCO}_3$ ), which also did not resemble the soil characteristics of the local monitoring sites (see Figure



5). These findings together with the accurate validation results clearly indicate that dissimilarity and diversity in soils can also provide the means for fitting locally-adaptive models.

Nevertheless, we can yet only speculate about how and why such diverse calibration sets are able to leverage accurate local calibrations. One hypothesis is that by increasing the range and variability in spectral variables and measurements a model can become quite stable in the central range of local reference measurements because a larger range of input variables is considered; thereby, the RS-LOCAL subsets that are selected from the SSL and used for PLS regression would stabilize and reduce the errors of the local samples. We imagine that we leverage a similar mechanism as in simple linear regression, where narrowing the range of the independent variable ( $x$ ) in the training samples would decrease the accuracy of intermediate values of the independent variable. We therefore need to further look into the details of spectral dissimilar learning, for example, also investigating the relevance of specific spectral features for local spectral transfers. The inherent working principles of RS-LOCAL are in contrast to the spectrum-based learner (SBL) or other forms of memory-based learning that utilize similar samples to infer sample-specific predictions based on existing training data (Lin and Vitter, 1994; Ramirez-Lopez et al., 2013). Our approach could describe a data-driven phenomenon, which implies that spectra can help to estimate a set of unrelated new soils. Another possibility is that there is in fact a pedological explanation that could be elucidated with more soil covariates such as mineralogy.

Local soil characterization is simpler, quicker and cheaper when a large proportion of properties of new soils are estimated by spectroscopy. Our results suggest the importance of optimizing the transfer of relevant information present in large SSLs to minimize the required amount of conventional laboratory analyses of new soils. Soil chemical and physical heterogeneity can be substantial in large SSLs. Therefore, such data variation can be beneficial for future predictions of properties of soils. However, learning a single general model over a heterogeneous training set, and obtaining parameter estimates optimized with a global measure of goodness of fit can introduce bias and inaccuracy to local (soil) estimation (Hand and Vinciotti, 2003; Ramirez-Lopez et al., 2013; Lobsey et al., 2017). Although the highest estimation accuracy could be achieved with only soils of the target study area (Stenberg and Rossel, 2010; Guerrero et al., 2016), it is impractical and inefficient to derive a single spectral prediction model with those. It requires 1) a large volume of reference measurements for a reasonably accurate multivariate calibration, and 2) it does not utilize already existing soil information.

Currently, the Swiss long-term soil monitoring uses a spatially representative sampling and then bulking the soils into four replicates for reference measurements (Desaules et al., 2010; Gubler et al., 2019). When the long-term monitoring would be augmented with mid-IR spectroscopy, one could make spectral measurements on all subsamples, rather than only on bulked samples, which would deliver spatially-explicit information and reduce nuisance factors from different sampling conditions. If not constrained economically (separate drying and sieving of sub-samples), a spectral workflow could thus allow to account for small-scale soil variability and reduce bias in measurements to robustly estimate temporal soil changes. For example, there is currently a relatively large variability in C measurements between the bulked replicate samples at one time point (Gubler et al., 2019). Our results suggest that unbiased spectral measurements eventually mediate such inconsistencies.

Relatively precise and unbiased geographically-local estimates of soil properties from diverse and large SSLs can be achieved by a handful of data-driven statistical approaches that are currently popular in the soil science community (Viscarra Rossel





and Webster, 2012; Ramirez-Lopez et al., 2013; Guerrero et al., 2014; Lobsey et al., 2017; Tsakiridis et al., 2020). Among the methods, we tested RS-LOCAL Lobsey et al. (2017) in our local soil monitoring scenario. Compared to memory-based learning, such as SBL (Ramirez-Lopez et al. (2013)), RS-LOCAL does not precondition the choice of useful subsets based on similarity in the input dimensions, here spectra, when performing the selection of SSL samples. The RS-LOCAL method is applied to exhaustively sample instances from the SSL without replacement, while it preferably selects those that perform well on the local target set, using PLS regression. An advantage of the method is that it can further deal with erroneous spectra as well as inaccurate and imprecise analytical reference measurements in the SSL, because it filters them as irrelevant instances. Besides chemometric and classical machine learning approaches, convolutional neural networks are being popularized for modeling SSLs with large soil variability (e.g., Liu et al. (2018); Padarian et al. (2019a, b); Tsakiridis et al. (2020)). There seems to be a small performance gain of a multi-output CNN with a similarity-based error correction using neighbors compared to the SBL (Tsakiridis et al. (2020); RMSE = 10.96 g kg<sup>-1</sup> vs. 11.74 g kg<sup>-1</sup> for OC). Despite the current development of interpretation methods in deep learning, CUBIST and PLSR modeling employed in both in the SBL and RS-LOCAL offer easier interpretation with comparable accuracy to CNNs.

Transfer learning or local learning introduces a new paradigm to supervised learning: model building is governed by the intended model application and thus coupled to it (Hand and Vinciotti, 2003). This contrasts general-model application, where the inference process is separated from the prediction of new data. Including local samples and their local data characteristics is necessary in order that a combined search and learning algorithm has a chance to capture predictive mechanisms. At the same time, the selection process and the partial data dependence within the predictive unit, the site, requires a careful assessment scheme to prevent a potential selection bias in the assessment of the approach. To account for this, we kept the respective site-specific local tuning and calibration set — whose hold-out performance directed the iterative search process and the reduction of the SSL — at minimum size of two observations at  $t_0$  or in addition  $t_1$  when only one measurement was available from the first sampling (see Figure 2).

### 4.3 Future applications and updates of the SSL

We found that data-driven modeling with selection of spectral dissimilar soils (see Figure 7) is accurate for inducing local predictions of total C (Figure 6). Hence, there is the need to further improve data-driven selection using RS-LOCAL, i.e., by further optimizing the current version of the algorithm. To address this need, we could use combined memory-based or lazy learning strategies (Stanfill and Waltz, 1986; Lin and Vitter, 1994; Ramirez-Lopez et al., 2013) to optimize with more data-driven transfer methods (Pan and Yang, 2010) in terms of reducing the time needed to evaluate suitable subsets of the SSL for a new application. To give an example, some similarity criteria or clustering before doing calibration sampling could be used as prior information for reducing the SSL size to obtain the final subsets. In principle, the sample reduction could also be done with algorithms that can deal with non-linear relationships between spectra and soil properties, such as random forest or Cubist. Another extension is to further filter spectral features and to do data compression to make the local modeling faster and even more adaptive to local conditions.



Our results showed that a transfer of the SSL to individual monitoring sites yielded very low bias and was accurate. This indicates that mid-IR spectroscopy and SSLs have the potential to give quick and relatively precise soil property estimates for soil monitoring. Nevertheless, the sites of the NABO long-term-monitoring program has not undergone substantial changes in OC (Gubler et al., 2019). Up to now, although major changes carbon content and organic composition should yield a spectral response, spectral changes in OC have mostly been reported along chronosequences (i.e., Awiti et al. (2008)), and only rarely for changes within individual plots over time (Deng et al., 2013). Hence, to address this, we propose to further investigate to what extent mid-IR spectroscopy can detect changes of OC considering small-scale variability and different agronomic management practices. This could for example be achieved with a study using soils from a long-term field trial, that shows sufficient temporal changes to be detected with spectroscopy.

The current SSL includes soils that contain between 0.1 and 58.3 % C, and 0.0 and 27.3 % OC (Table 1). Since organic soils can have up to 50 % OC, organic soils might be underrepresented in the current Swiss SSL. For this reason, we also tested the present Swiss SSL with a case study and augmented it to better represent organic soils, using new soils from two peat land regions in Switzerland (Helfenstein et al., 2020, submitted).

Our results suggest that the present mid-IR SSL has great potential for applications that require soil data in high temporal and spatial coverage (i.e., for deriving quantitative indicators of soil quality for spatial planning or for soil-related environmental research). Mid-infrared spectral modeling was able to estimate many soil properties accurately with rather large variation in measurements explained (Figure 3), making them suitable for agronomic diagnosis and the assessment of soil functions in various landscapes. Currently, fine grained soil information of properties and function across agricultural lands in Switzerland is still scarce and often challenging to harmonize (i.e., measurement methods) because legacy maps are at varying levels of detail and quality (Keller et al., 2018; Grêt-Regamey et al., 2018). For example, only 13 % (127000 ha) of soil in agricultural land has been mapped with soil attributes of sufficient quality to evaluate its potential for crop production (Rehbein et al., 2020). Soil properties are also insufficiently mapped nationwide from point into space, depth and over time to regionally model soil processes, or to evaluate site-specific effects of agricultural practices on soils (i.e., soil C dynamics). Therefore, we suggest to couple infrared spectral estimation with traditional soil surveys and digital soil mapping to speed up the collection of soil information in Switzerland and elsewhere. This will offer means to test and further extend this SSL, so that only minimal amounts of costly and time consuming traditional laboratory analyses will be needed for characterizing and mapping soils' properties and functions in the next decades.

## 5 Conclusions

We developed the Swiss mid-IR SSL ( $n = 4374$ ), using legacy soils and reference measurements of 17 properties, from 71 long-term monitoring sites (national soil monitoring; NABO) and 1094 locations sampled from a regular grid over Switzerland (biodiversity monitoring program; BDM). The trained CUBIST models — a general modeling approach using all data — were able to explain a relatively large proportion ( $R^2 > 0.6$ ) of measured variance for 11 out of 17 properties. Total C, OC, total N, pH,  $\text{CEC}_{\text{pot}}$ , and clay content were estimated with high discrimination capacity ( $R^2 > 0.8$ ). Total C was estimated with a



cross-validated RMSE = 0.84 % at a measured range of 0.1 – 58.3 %, and OC with RMSE = 1.20 % at a measured range of 0.0 – 27.3 %. Compared to the general CUBIST approach, the local transfer yielded on average 4.4 times more accurate estimates of total C with the mean RMSE = 0.07 %, which is a substantial improvement of local estimates at plot-scale. Our similarity analysis revealed that local learning with subset selection based on RS-LOCAL produced a chemically diverse calibration set rather than narrowing down soil diversity for local modeling, as it is for example the case in memory-based learning. The developed national mid-IR SSL offers rapid soil estimates which are key inputs for many applications requiring soil information. These include for example, digital soil mapping (i.e., spatial planning, assessing soil functions), agronomic diagnostics and precision farming, soil C accounting and monitoring, as well as ecological and soil research (i.e., biogeochemical modeling, soil physics). The created mid-IR SSL and both local and general models can be updated with new soil records, which will allow to cover more soils conditions and will require less and less soil laboratory reference measurements in relation to spectral measurements for monitoring, mapping and modeling new soils.

## Appendix A: Figures and tables in appendices

### A1 Recursive feature elimination for interpreting general soil estimation with CUBIST

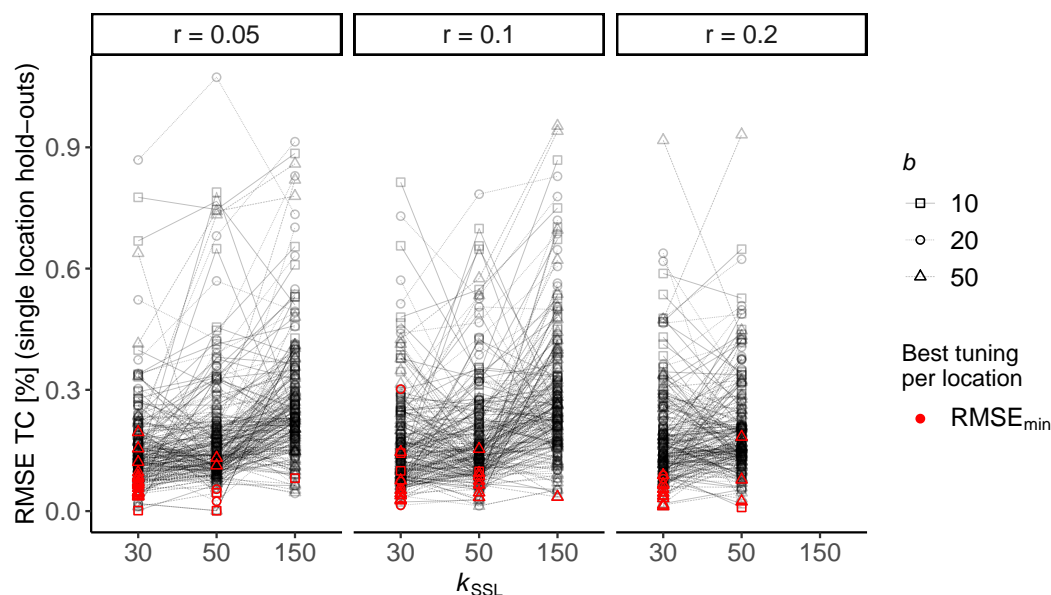
The recursive feature elimination (RFE) procedure started with the initial set of  $S_1 = 1665$  predictive variables that resulted after processing the spectra (see section 2.3). The following subset sizes  $S_i$  representing the number of spectral variables that are retained after each  $i^{\text{th}}$  variable elimination step were defined and evaluated within the RFE procedure:

$$S_i = \{209, 150, 120, 105, 90, 75, 60, 50, 40, 35, 30, 25, 20, 17, 14, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1\} \quad (\text{A1})$$

The first variable elimination step ( $i = 1$ ) started with tuning a full CUBIST model derived from  $S_1 = 209$  possible predictors using 10-fold cross-validation, then calculating the CUBIST model usage statistics for all predictors, next sorting all predictors from highest to lowest importance, and lastly dropping  $S_1 - S_2 = 59$  of the least important predictors. For the next iteration ( $i = 2$ ) and the following ones, we repeated this model fitting and variable reduction procedure with  $S_2 = 150$  predictors and the preceding subsets, until the most important predictive variable ( $S_{30} = 1$ ) was left at the last iteration ( $i = 30$ ).

Variable selection is in addition prone to overoptimistic model assessment when resampling subsets (i.e., cross-validation) are used for two purposes, here model building and selection. This selection bias due to data leakage is well-documented for so-called *wrapper methods* of variable selection like RFE (Ambroise and McLachlan, 2002; Kuhn and Johnson, 2013), and occurs if these two tasks are not sufficiently separated by using independent data sets for each of them; this becomes especially more important when many predictive variables in relation to relatively few observations are used, as it the case for our spectra.

To provide realistic predictive generalization of the RFE method, the aforementioned iterative selection procedure was done within an internal cross-validation scheme so that independent data were used to test the performance of the variable selection on the outer data segments. These outer cross-validation segments served external validation. To quantify the uncertainty of



**Figure A1.** Performance profile of the 27 empirical parameter combinations of RS-LOCAL tested on each of the 71 NABO sites. The errors (root mean squared error (RMSE)) of the plot-level transfer was assessed with the first two calibration samples for each time series of total C (see Figure 1 for an illustration of the setup of the local predictive transfer)

the models using the reduced variable sets and specifically variable selection, the outer cross-validation layer that served cross-validation was repeated five times, leading to five independent estimations per sample.

## A2 Tuning profile of the RS-LOCAL parameters for local predictive transfers

The most relevant samples from the SSL at each respective NABO long-term monitoring plot were empirically selected at the RS-LOCAL configuration that yielded the lowest RMSE on two calibration samples per plot (Figure A1; performance profile). Time-series validation on the remaining samples of each site was separated from the optimization in the transfer workflow (see Figure 2).

*Code and data availability.* The data and the code to reproduce the results of this manuscript are available upon reasonable request.

*Author contributions.* Philipp Baumann drafted the manuscript, conceptualized spectral modeling scenario, carried out the data analysis, and wrote the code. Juhwan Lee and Johan Six developed the concept of having a Swiss Spectral Library. All co-authors contributed to the manuscript. NABO provided the milled soil archive of NABO and BDM, together with the chemical reference data and metadata.



Andreas Gubler created the map of locations. Raphael Viscarra Rossel helped guide the spectral modeling experiments and provided the implementation of RS-LOCAL in R.

*Competing interests.* The authors declare no competing interests.

*Acknowledgements.* We express thanks to Leonardo Ramirez-Lopez, Laura Summerauer, and Jonas Anderegg for their valuable inputs that helped to improve the manuscript. This research was funded by ETH grants. We also would like to express gratitude to the Federal Office of Environment for commissioning and funding the soil analyses conducted within the BDM.



## References

- Agroscope: Referenzmethoden von Agroscope, 1996.
- Ambrose, C. and McLachlan, G. J.: Selection Bias in Gene Extraction on the Basis of Microarray Gene-Expression Data, *Proceedings of the National Academy of Sciences*, 99, 6562–6566, <https://doi.org/10.1073/pnas.102102699>, 2002.
- 5 Angelopoulou, T., Balafoutis, A., Zalidis, G., and Bochtis, D.: From Laboratory to Proximal Sensing Spectroscopy for Soil Organic Carbon Estimation—A Review, *Sustainability*, 12, 443, <https://doi.org/10.3390/su12020443>, 2020.
- Awiti, A. O., Walsh, M. G., Shepherd, K. D., and Kinyamario, J.: Soil condition classification using infrared spectroscopy: A proposition for assessment of soil condition along a tropical forest-cropland chronosequence, *Geoderma*, 143, 73–84, <https://doi.org/10.1016/j.geoderma.2007.08.021>, <http://www.sciencedirect.com/science/article/pii/S0016706107002625>, 2008.
- 10 Baumann, P.: philipp-baumann/simplerspec: Beta release simplerspec 0.1.0 for zenodo, <https://doi.org/10.5281/zenodo.3303637>, <https://doi.org/10.5281/zenodo.3303637>, 2019.
- Bellman, R.: *Adaptive Control Processes: A Guided Tour*, Princeton University Press, <http://www.jstor.org/stable/j.ctt183ph6v>, 1961.
- Bellon-Maurel, V., Fernandez-Ahumada, E., Palagos, B., Roger, J.-M., and McBratney, A.: Critical review of chemometric indicators commonly used for assessing the quality of the prediction of soil attributes by NIR spectroscopy, *TrAC Trends in Analytical Chemistry*, 29, 1073–1081, <https://doi.org/10.1016/j.trac.2010.05.006>, <http://www.sciencedirect.com/science/article/pii/S0165993610001585>, 2010.
- 15 Briedis, C., Baldock, J., de Moraes Sá, J. C., dos Santos, J. B., and Milori, D. M. B. P.: Strategies to Improve the Prediction of Bulk Soil and Fraction Organic Carbon in Brazilian Samples by Using an Australian National Mid-Infrared Spectral Library, *Geoderma*, 373, 114–140, <https://doi.org/10.1016/j.geoderma.2020.114401>, 2020.
- Bundesamt für Umwelt (BAFU): *Biodiversitätsmonitoring Schweiz BDM*, 2014.
- 20 Clairotte, M., Grinand, C., Kouakoua, E., Thébaud, A., Saby, N. P. A., Bernoux, M., and Barthès, B. G.: National Calibration of Soil Organic Carbon Concentration Using Diffuse Infrared Reflectance Spectroscopy, *Geoderma*, 276, 41–52, <https://doi.org/10.1016/j.geoderma.2016.04.021>, 2016.
- Dangal, S. R. S., Sanderman, J., Wills, S., and Ramirez-Lopez, L.: Accurate and Precise Prediction of Soil Properties from a Large Mid-Infrared Spectral Library, *Soil Systems*, 3, 11, <https://doi.org/10.3390/soilsystems3010011>, <https://www.mdpi.com/2571-8789/3/1/11>, 25 2019.
- Deng, F., Minasny, B., Knadel, M., McBratney, A., Heckrath, G., and Greve, M. H.: Using Vis-NIR Spectroscopy for Monitoring Temporal Changes in Soil Organic Carbon, *Soil Science*, 178, 389–399, <https://doi.org/10.1097/SS.0000000000000002>, [https://journals.lww.com/soilsci/Fulltext/2013/08000/Using\\_Vis\\_NIR\\_Spectroscopy\\_for\\_Monitoring\\_Temporal.2.aspx](https://journals.lww.com/soilsci/Fulltext/2013/08000/Using_Vis_NIR_Spectroscopy_for_Monitoring_Temporal.2.aspx), 2013.
- Desaules, A., Ammann, S., and Schwab, P.: Advances in long-term soil-pollution monitoring of Switzerland, *Journal of Plant Nutrition and Soil Science*, 173, 525–535, <https://doi.org/10.1002/jpln.200900269>, <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jpln.200900269>, 30 2010.
- Dietterich, T. G., Wettschereck, D., Atkeson, C. G., and Moore, A. W.: Memory-Based Methods for Regression and Classification, in: *Advances in Neural Information Processing Systems 6*, [7th NIPS Conference, Denver, Colorado, USA, 1993], edited by Cowan, J. D., Tesauero, G., and Alspector, J., pp. 1165–1166, Morgan Kaufmann, 1993.
- 35 Dokuchaev, V.: Report to the Transcaucasian Statistical Committee on Land Evaluation in General and Especially for the Transcaucasia. Horizontal and Vertical Soil Zones. (In Russian.) Off. Press Civ. Affairs Commander-in-Chief Cacasus, Tiflis, Russia, 1899.





- Dowle, M. and Srinivasan, A.: data.table: Extension of 'data.frame', <https://CRAN.R-project.org/package=data.table>, r package version 1.12.8, 2019.
- England, J. R. and Viscarra Rossel, R. A.: Proximal Sensing for Soil Carbon Accounting, *SOIL*, 4, 101–122, <https://doi.org/10.5194/soil-4-101-2018>, <https://www.soil-journal.net/4/101/2018/>, 2018.
- 5 Friedman, J., Hastie, T., and Tibshirani, R.: The elements of statistical learning, Springer series in statistics Springer, Berlin, second edition edn., <http://statweb.stanford.edu/~tibs/book/preface.ps>, 2008.
- Grêt-Regamey, A., Kool, S., Bühlmann, L., and Kissling, S.: Eine Bodenagenda für die Raumplanung, Thematische Synthese TS3 des Nationalen Forschungsprogramms «Nachhaltige Nutzung der Ressource Boden» (NFP 68), Swiss National Science Foundation (SNF), Bern, 2018.
- 10 Gubler, A., Wächter, D., Schwab, P., Müller, M., and Keller, A.: Twenty-five years of observations of soil organic carbon in Swiss croplands showing stability overall but with some divergent trends, *Environmental Monitoring and Assessment*, 191, 277, <https://doi.org/10.1007/s10661-019-7435-y>, 2019.
- Guerrero, C., Stenberg, B., Wetterlind, J., Viscarra Rossel, R. A., Maestre, F. T., Mouazen, A. M., Zornoza, R., Ruiz-Sinoga, J. D., and Kuang, B.: Assessment of Soil Organic Carbon at Local Scale with Spiked NIR Calibrations: Effects of Selection and Extra-Weighting on the Spiking Subset: Spiking and Extra-Weighting to Improve Soil Organic Carbon Predictions with NIR, *European Journal of Soil Science*, 65, 248–263, <https://doi.org/10.1111/ejss.12129>, 2014.
- 15 Guerrero, C., Wetterlind, J., Stenberg, B., Mouazen, A. M., Gabarrón-Galeote, M. A., Ruiz-Sinoga, J. D., Zornoza, R., and Viscarra Rossel, R. A.: Do We Really Need Large Spectral Libraries for Local Scale SOC Assessment with NIR Spectroscopy?, *Soil and Tillage Research*, 155, 501–509, <https://doi.org/10.1016/j.still.2015.07.008>, <http://linkinghub.elsevier.com/retrieve/pii/S0167198715001567>, 2016.
- 20 Guyon, I., Weston, J., Barnhill, S., and Vapnik, V.: Gene Selection for Cancer Classification Using Support Vector Machines, *Machine Learning*, 46, 389–422, <https://doi.org/10.1023/A:1012487302797>, 2002.
- Hand, D. J. and Vinciotti, V.: Local Versus Global Models for Classification Problems: Fitting Models Where it Matters, *The American Statistician*, 57, 124–131, <https://doi.org/10.1198/0003130031423>, 2003.
- Helfenstein, A., Baumann, P., Viscarra Rossel, R., Gubler, A., Oechlin, S., and Six, J.: Predicting Soil Carbon by Efficiently Using Variation in a Mid-IR Soilspectral Library, submitted, 2020.
- 25 Hong, Y., Chen, S., Liu, Y., Zhang, Y., Yu, L., Chen, Y., Liu, Y., Cheng, H., and Liu, Y.: Combination of Fractional Order Derivative and Memory-Based Learning Algorithm to Improve the Estimation Accuracy of Soil Organic Matter by Visible and near-Infrared Spectroscopy, *CATENA*, 174, 104–116, <https://doi.org/10.1016/j.catena.2018.10.051>, <http://www.sciencedirect.com/science/article/pii/S0341816218304867>, 2019.
- 30 Hubert, M. and Debruyne, M.: Minimum Covariance Determinant, *WIREs Computational Statistics*, 2, 36–43, <https://doi.org/10.1002/wics.61>, 2010.
- Janik, L. J. and Skjemstad, J. O.: Characterization and analysis of soils using mid-infrared partial least-squares .2. Correlations with some laboratory data, *Australian Journal of Soil Research*, 33, 637–650, <https://doi.org/10.1071/sr9950637>, <https://www.publish.csiro.au/sr/sr9950637>, publisher: CSIRO PUBLISHING, 1995.
- 35 Janik, L. J., Skjemstad, J. O., and Merry, R. H.: Can mid infrared diffuse reflectance analysis replace soil extractions?, *Australian Journal of Experimental Agriculture*, 38, 681, <https://doi.org/10.1071/EA97144>, <http://www.publish.csiro.au/?paper=EA97144>, 1998.
- Jenny, H.: Factors of Soil Formation, McGraw-Hill Book Co., New York, 1941.



- Keller, A., Franzen, J., Knüsel, P., Papritz, A., and Zürrer, M.: Bodeninformations-Plattform Schweiz (BIP-CH), Thematische Synthese TS4 des Nationalen Forschungsprogramms «Nachhaltige Nutzung der Ressource Boden» (NFP 68), Swiss National Science Foundation (SNF), Bern, 2018.
- Kuhn, M.: caret: Classification and Regression Training, <https://CRAN.R-project.org/package=caret>, r package version 6.0-85, 2020.
- 5 Kuhn, M. and Johnson, K.: Applied Predictive Modeling, Springer New York, New York, NY, 2013.
- Lin, J.-H. and Vitter, J. S.: A Theory for Memory-Based Learning, Machine Learning, 17, 143–167, <https://doi.org/10.1023/A:1022667616941>, 1994.
- Liu, L., Ji, M., and Buchroithner, M.: Transfer Learning for Soil Spectroscopy Based on Convolutional Neural Networks and Its Application in Soil Clay Content Mapping Using Hyperspectral Imagery, Sensors (Basel, Switzerland), 18, <https://doi.org/10.3390/s18093169>, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6165490/>, 2018.
- 10 Lobsey, C. R., Viscarra Rossel, R. A., Roudier, P., and Hedley, C. B.: RS-LOCAL Data-Mines Information from Spectral Libraries to Improve Local Calibrations: RS-LOCAL Improves Local Spectroscopic Calibrations, European Journal of Soil Science, <https://doi.org/10.1111/ejss.12490>, 2017.
- Madari, B. E., Reeves, J. B., Machado, P. L., Guimarães, C. M., Torres, E., and McCarty, G. W.: Mid- and near-Infrared Spectroscopic Assessment of Soil Compositional Parameters and Structural Indices in Two Ferralsols, Geoderma, 136, 245–259, <https://doi.org/10.1016/j.geoderma.2006.03.026>, 2006.
- 15 Meuli, R. G., Wächter, D., Schwab, P., Kohli, L., and Zimmermann, R.: Connecting Biodiversity Monitoring with Soil Inventory Data – A Swiss Case Study, BGS Bulletin, 38, 3, 2017.
- Nguyen, T., Janik, L., and Raupach, M.: Diffuse Reflectance Infrared Fourier Transform (DRIFT) Spectroscopy in Soil Studies, Soil Research, 29, 49, <https://doi.org/10.1071/SR9910049>, 1991.
- 20 Nocita, M., Stevens, A., van Wesemael, B., Brown, D. J., Shepherd, K. D., Towett, E., Vargas, R., and Montanarella, L.: Soil Spectroscopy: An Opportunity to Be Seized, Global Change Biology, 21, 10–11, <https://doi.org/10.1111/gcb.12632>, 2015.
- Ogen, Y., Zaluda, J., Francos, N., Goldshleger, N., and Ben-Dor, E.: Cluster-Based Spectral Models for a Robust Assessment of Soil Properties, Geoderma, 340, 175–184, <https://doi.org/10.1016/j.geoderma.2019.01.022>, <http://www.sciencedirect.com/science/article/pii/S0016706118316045>, 2019.
- 25 Padarian, J., Minasny, B., and McBratney, A. B.: Transfer Learning to Localise a Continental Soil Vis-NIR Calibration Model, Geoderma, 340, 279–288, <https://doi.org/10.1016/j.geoderma.2019.01.009>, <http://www.sciencedirect.com/science/article/pii/S0016706118305639>, 2019a.
- Padarian, J., Minasny, B., and McBratney, A. B.: Using Deep Learning to Predict Soil Properties from Regional Spectral Data, Geoderma Regional, 16, e00 198, <https://doi.org/10.1016/j.geodrs.2018.e00198>, <http://www.sciencedirect.com/science/article/pii/S2352009418302785>, 2019b.
- 30 Pan, S. J. and Yang, Q.: A Survey on Transfer Learning, IEEE Transactions on Knowledge and Data Engineering, 22, 1345–1359, <https://doi.org/10.1109/TKDE.2009.191>, 2010.
- Parikh, S. J., Goynes, K. W., Margenot, A. J., Mukome, F. N. D., and Calderón, F. J.: Chapter One - Soil Chemical Insights Provided through Vibrational Spectroscopy, in: Advances in Agronomy, edited by Sparks, D. L., vol. 126, pp. 1–148, Academic Press, <https://doi.org/10.1016/B978-0-12-800132-5.00001-8>, 2014.
- Pratt, L. and Thrun, S.: Guest Editors' Introduction, Machine Learning, 28, 5–5, <https://doi.org/10.1023/A:1007322005825>, 1997.



- Pratt, L. Y., Pratt, L. Y., Hanson, S. J., Giles, C. L., and Cowan, J. D.: Discriminability-Based Transfer between Neural Networks, in: Advances in Neural Information Processing Systems 5, pp. 204–211, Morgan Kaufmann, 1993.
- Quinlan, J.: Combining Instance-Based and Model-Based Learning, in: Machine Learning Proceedings 1993, pp. 236–243, Elsevier, <https://doi.org/10.1016/B978-1-55860-307-3.50037-X>, <https://linkinghub.elsevier.com/retrieve/pii/B978155860307350037X>, 1993.
- 5 Quinlan, J. R.: Learning with Continuous Classes, in: 5th Australian Joint Conference on Artificial Intelligence, vol. 92, pp. 343–348, World Scientific, 1992.
- R Core Team: R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>, 2019.
- Ramirez-Lopez, L., Behrens, T., Schmidt, K., Stevens, A., Demattê, J. A. M., and Scholten, T.: The Spectrum-Based  
10 Learner: A New Local Approach for Modeling Soil Vis–NIR Spectra of Complex Datasets, *Geoderma*, 195–196, 268–279, <https://doi.org/10.1016/j.geoderma.2012.12.014>, <http://www.sciencedirect.com/science/article/pii/S0016706112004314>, 2013.
- Rehbein, K., Sprecher, C., and Keller, A.: Übersicht Stand Bodenkartierung in Der Schweiz. Ergänzung Des Bodenkartierungskataloges Schweiz Um Bodeninformationen Aus Meliorationsprojekten, Servicestelle NABODAT, Agroscope, Zürich, 2020.
- Rousseeuw, P. J.: Least Median of Squares Regression, *Journal of the American Statistical Association*, 79, 871–880,  
15 <https://doi.org/10.1080/01621459.1984.10477105>, 1984.
- Savitzky, A. and Golay, M. J. E.: Smoothing and Differentiation of Data by Simplified Least Squares Procedures., *Analytical Chemistry*, 36, 1627–1639, <https://doi.org/10.1021/ac60214a047>, publisher: American Chemical Society, 1964.
- Seidel, M., Hutengs, C., Ludwig, B., Thiele-Bruhn, S., and Vohland, M.: Strategies for the Efficient Estimation of Soil Organic Carbon at the Field Scale with Vis–NIR Spectroscopy: Spectral Libraries and Spiking vs. Local Calibrations, *Geoderma*, 354, 113–125,  
20 <https://doi.org/10.1016/j.geoderma.2019.07.014>, <http://www.sciencedirect.com/science/article/pii/S0016706119304537>, 2019.
- Sila, A. M., Shepherd, K. D., and Pokhariyal, G. P.: Evaluating the Utility of Mid-Infrared Spectral Subspaces for Predicting Soil Properties, *Chemometrics and Intelligent Laboratory Systems*, 153, 92–105, <https://doi.org/10.1016/j.chemolab.2016.02.013>, <http://linkinghub.elsevier.com/retrieve/pii/S0169743916300351>, 2016.
- Skjemstad, J. and Dalal, R.: Spectroscopic and Chemical Differences in Organic Matter of Two Vertisols Subjected to Long Periods of  
25 Cultivation, *Soil Research*, 25, 323, <https://doi.org/10.1071/SR9870323>, 1987.
- Solomatine, D.: Combining Machine Learning and Domain Knowledge in Modular Modelling, in: Practical Hydroinformatics: Computational Intelligence and Technological Developments in Water Applications, edited by Abrahart, R. J., See, L. M., and Solomatine, D. P., Water Science and Technology Library, pp. 333–345, Springer, Berlin, Heidelberg, [https://doi.org/10.1007/978-3-540-79881-1\\_24](https://doi.org/10.1007/978-3-540-79881-1_24), 2008.
- Soriano-Disla, J. M., Janik, L. J., Viscarra Rossel, R. A., Macdonald, L. M., and McLaughlin, M. J.: The Performance of Visible, Near-,  
30 and Mid-Infrared Reflectance Spectroscopy for Prediction of Soil Physical, Chemical, and Biological Properties, *Applied Spectroscopy Reviews*, 49, 139–186, <https://doi.org/10.1080/05704928.2013.811081>, 2014.
- Stanfill, C. and Waltz, D.: Toward Memory-Based Reasoning, *Communications of the ACM*, 29, 1213–1228, <https://doi.org/10.1145/7902.7906>, <http://portal.acm.org/citation.cfm?doid=7902.7906>, 1986.
- Stenberg, B. and Rossel, R. V.: Diffuse Reflectance Spectroscopy for High-Resolution Soil Sensing, in: Proximal Soil Sensing, edited  
35 by Viscarra Rossel, R. A., McBratney, A. B., and Minasny, B., Progress in Soil Science, pp. 29–47, Springer Netherlands, Dordrecht, [https://doi.org/10.1007/978-90-481-8859-8\\_3](https://doi.org/10.1007/978-90-481-8859-8_3), 2010.
- Stevens, A. and Ramirez-Lopez, L.: An introduction to the prospectr package, r package version 0.1.3, 2013.
- Thrun, S. and Pratt, L., eds.: Learning to Learn, Springer US, Boston, MA, <https://doi.org/10.1007/978-1-4615-5529-2>, 1998.



- Tsakiridis, N. L., Keramaris, K. D., Theocharis, J. B., and Zalidis, G. C.: Simultaneous prediction of soil properties from VNIR-SWIR spectra using a localized multi-channel 1-D convolutional neural network, *Geoderma*, 367, 114–208, <https://doi.org/10.1016/j.geoderma.2020.114208>, <http://www.sciencedirect.com/science/article/pii/S0016706119308870>, 2020.
- Tziolas, N., Tsakiridis, N., Ben-Dor, E., Theocharis, J., and Zalidis, G.: A Memory-Based Learning Approach Utilizing Combined Spectral Sources and Geographical Proximity for Improved VIS-NIR-SWIR Soil Properties Estimation, *Geoderma*, 340, 11–24, <https://doi.org/10.1016/j.geoderma.2018.12.044>, <http://www.sciencedirect.com/science/article/pii/S0016706118307006>, 2019.
- Varmuza, K. and Filzmoser, P.: *Introduction to Multivariate Statistical Analysis in Chemometrics*, CRC Press, 0 edn., <https://doi.org/10.1201/9781420059496>, <https://www.taylorfrancis.com/books/9781420059496>, 2016.
- Viscarra Rossel, R., Behrens, T., Ben-Dor, E., Brown, D., Demattê, J., Shepherd, K., Shi, Z., Stenberg, B., Stevens, A., Adamchuk, V., Aichi, H., Barthès, B., Bartholomeus, H., Bayer, A., Bernoux, M., Böttcher, K., Brodský, L., Du, C., Chappell, A., Fouad, Y., Genot, V., Gomez, C., Grunwald, S., Gubler, A., Guerrero, C., Hedley, C., Knadel, M., Morras, H., Nocita, M., Ramirez-Lopez, L., Roudier, P., Campos, E. R., Sanborn, P., Sellitto, V., Sudduth, K., Rawlins, B., Walter, C., Winowiecki, L., Hong, S., and Ji, W.: A Global Spectral Library to Characterize the World's Soil, *Earth-Science Reviews*, 155, 198–230, <https://doi.org/10.1016/j.earscirev.2016.01.012>, 2016.
- Viscarra Rossel, R. A. and McBratney, A. B.: Soil chemical analytical accuracy and costs: implications from precision agriculture, *Australian Journal of Experimental Agriculture*, 38, 765, <https://doi.org/10.1071/EA97158>, <http://www.publish.csiro.au/?paper=EA97158>, 1998.
- Viscarra Rossel, R. A. and Webster, R.: Predicting Soil Properties from the Australian Soil Visible-near Infrared Spectroscopic Database, *European Journal of Soil Science*, 63, 848–860, <https://doi.org/10.1111/j.1365-2389.2012.01495.x>, 2012.
- Viscarra Rossel, R. A., Lobsey, C. R., Sharman, C., Flick, P., and McLachlan, G.: Novel Proximal Sensing for Monitoring Soil Organic C Stocks and Condition, *Environmental Science & Technology*, 51, 5630–5641, <https://doi.org/10.1021/acs.est.7b00889>, 2017.
- Wang, Y. and Witten, I. H.: Induction of model trees for predicting continuous classes, Working Paper 96/23, University of Waikato, Department of Computer Science, Hamilton, New Zealand, <https://researchcommons.waikato.ac.nz/handle/10289/1183>, accepted: 2008-10-29T02:09:15Z ISSN: 1170-487X, 1996.
- Wickham, H.: tidyverse: Easily Install and Load the 'Tidyverse', <https://CRAN.R-project.org/package=tidyverse>, r package version 1.3.0, 2019.
- Wolpert, D. and Macready, W.: No Free Lunch Theorems for Optimization, *IEEE Transactions on Evolutionary Computation*, 1, 67–82, <https://doi.org/10.1109/4235.585893>, 1997.
- Wolpert, D. H.: The Lack of A Priori Distinctions Between Learning Algorithms, *Neural Computation*, 8, 1341–1390, <https://doi.org/10.1162/neco.1996.8.7.1341>, 1996.