



HAL
open science

Prediction of soil organic and inorganic carbon concentrations in Tunisian samples by mid-infrared reflectance spectroscopy using a French national library

Cécile Gomez, Tiphaine Chevallier, Patricia Moulin, Imane Bouferra, Kaouther Hmaidi, Dominique D. Arrouays, Claudy C. Jolivet, Bernard Barthès

► To cite this version:

Cécile Gomez, Tiphaine Chevallier, Patricia Moulin, Imane Bouferra, Kaouther Hmaidi, et al.. Prediction of soil organic and inorganic carbon concentrations in Tunisian samples by mid-infrared reflectance spectroscopy using a French national library. *Geoderma*, 2020, 375, pp.114469. 10.1016/j.geoderma.2020.114469 . hal-02890493

HAL Id: hal-02890493

<https://hal.inrae.fr/hal-02890493v1>

Submitted on 23 Aug 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

1 **Prediction of soil organic and inorganic carbon concentrations in**
2 **Tunisian samples by mid-infrared reflectance spectroscopy using a**
3 **French national library**

4
5 Cécile Gomez ^{1,2}, Tiphaine Chevallier ³, Patricia Moulin ⁴, Imane Bouferra ³, Kaouther
6 Hmaldi ^{3,5}, Dominique Arrouays ⁶, Claudy Jolivet ⁶, Bernard G. Barthès ³

7
8 ¹ LISAH, University of Montpellier, INRAE, IRD, Montpellier SupAgro, 34060 Montpellier,
9 France. Corresponding author. cecile.gomez@ird.fr

10 ² Indo-French Cell for Water Sciences, IRD, Indian Institute of Science, Bangalore
11 560012, India

12 ³ Eco&Sols, University of Montpellier, CIRAD, INRAE, IRD, Montpellier SupAgro, 34060
13 Montpellier, France

14 ⁴ US IMAGO, IRD, BP1386, Dakar, Senegal

15 ⁵ UR Pédologie, Faculté des Sciences de Tunis, El Manar Tunis, Tunisia

16 ⁶ INRAE, US 1106 InfoSol, F45000, Orléans, France

17

18

19 **Keywords:** Soil Organic and Inorganic Carbon; Mid-Infrared reflectance spectroscopy;
20 natural logarithm transformation; Local and Global calibration; Partial least squares
21 regression (PLSR); National dataset.

22

23 **Abstract**

24 Mid-infrared reflectance spectroscopy (MIRS, 4000–400 cm⁻¹) is being considered to
25 provide accurate estimations of soil properties, including soil organic carbon (SOC) and

26 soil inorganic carbon (SIC) contents. This approach has mainly been demonstrated by
27 using datasets originating from the same area *A*, with similar geopedological conditions, to
28 build, validate and test prediction models. The objective of this study was to analyse how
29 MIRS performs when applied to predict SOC and SIC contents, from a calibration
30 database collected over a region *A*, to predict over a region *B*, where *A* and *B* have no
31 common area and different soil and climate conditions. This study used a French MIRS
32 soil dataset including 2178 topsoil samples to calibrate SIC and SOC prediction models
33 with partial least squares regression (PLSR), and a Tunisian MIRS topsoil dataset
34 including 96 soil samples to test them. Our results showed that when using the French
35 MIRS soil database, *i*) the SOC and SIC of French validation samples were successfully
36 predicted using global models ($R^2_{val} = 0.88$ and 0.98 , respectively), *ii*) the SIC of Tunisian
37 samples was also predicted successfully both using a global model and using a selection
38 of spectral neighbours from the French calibration database (R^2_{test} of 0.96 for both), *iii*) the
39 SOC of Tunisian samples was predicted moderately by global model (R^2_{test} of 0.64) and a
40 transformation by natural logarithm of the calibration SOC values significantly improved
41 the SOC prediction of Tunisian samples (R^2_{test} of 0.97), and *iv*) a transformation by natural
42 logarithm of SOC values provided more benefit than a selection of spectral neighbours
43 from the French calibration database for predicting Tunisian SOC values. Therefore, in the
44 future, MIRS might replace conventional physico-chemical analysis techniques, or at least
45 be considered as an alternative technique, especially when optimally exhaustive
46 calibration databases will become available.

47

48 **1. Introduction**

49 Soil is the largest reservoir of continental carbon and participates in the global carbon
50 cycle (Jacobson et al., 2000; Scharlemann et al., 2014). Soil emits carbon dioxide (CO₂)
51 through autotrophic and heterotrophic respiration and acts as a sink of atmospheric CO₂

52 through photosynthesis; then, organic decomposition products are integrated into the soil
53 as organic matter, which is composed of approximately 50% to 58% of carbon (Gregorich
54 et al., 1994; Pribyl, 2010). Moreover, soil organic carbon (SOC) has a long-acknowledged
55 and key role in soil physical, chemical and biological fertility (Reeves, 1997). Thus,
56 understanding the dynamics of soil carbon is a major issue for soil fertility and for climate
57 change mitigation.

58 Soil carbon is not only found in organic form. At the global scale, approximately one-
59 third of the total soil carbon is inorganic (soil inorganic carbon, SIC) (Batjes, 1996), and
60 calcareous soils cover more than 30% of the earth's land surface (Chen and Barak, 1982;
61 Romanyà and Rovira, 2011). SIC is mainly in the form of calcium carbonate (CaCO_3), and
62 high SIC contents are often localised in dry areas where SOC stocks are low. In these
63 soils, SIC levels can be 2 to 10 times higher than SOC levels (Bernoux and Chevallier,
64 2014). SIC is made of primary minerals derived from the fragmentation of carbonate
65 bedrock (lithogenic carbonates) or secondary minerals from inorganic carbon precipitation
66 in soil pores or around roots (pedogenic carbonates). Pedogenic carbonates have various
67 forms—nodules, lamellae or crystals—and have varying solubility.

68 Several analytical methods have been developed to quantify SIC and SOC contents in
69 soils, but they are tedious and/or costly. In particular, SOC determination in carbonated
70 soils often requires hazardous reagents. SIC content has usually been measured by
71 calcimetry (ISO, 1995b) but can also be measured by dry combustion with a CNH
72 elemental analyser equipped with a specific module ($\text{CO}_3\text{-C}$ module) after phosphoric acid
73 dissolution of the SIC (Mc Crea, 1950; Hannam et al. 2016). Quantifying SOC in
74 calcareous soils has either been carried out directly, by wet oxidation (Walkley and Black,
75 1934) or dry combustion after removing SIC by acid pretreatment (Harris et al., 2001), or
76 indirectly, by subtracting the SIC, measured by calcimetry, from the total carbon content
77 determined by dry combustion. Due to indirect determination, incomplete oxidation and the

78 use of hazardous reactants with both of these methods, alternative methods based on the
79 thermal lability of SIC and SOC have also been tested (Wang et al., 2012; Apesteguia et
80 al., 2018).

81 For one to two decades, visible near-infrared (Vis-NIR, 400–2500 nm) and mid-
82 infrared reflectance spectroscopy (MIRS, 4000–400 cm^{-1}), which measures diffuse
83 reflectance, have been proposed as alternative methods to these physico-chemical
84 analytical methods (e.g., Viscarra Rossel et al., 2006; Cécillon et al., 2008; Bellon-Maurel
85 and McBratney, 2011). MIRS is based on the study of absorption bands corresponding to
86 fundamental molecular vibrations, and NIRS is based on the study of absorption bands
87 corresponding to overtones and combinations of fundamental vibrations (Williams and
88 Norris, 1987). Several studies highlighting the potential of Vis-NIR and/or MIRS for
89 predicting various soil attributes, including SIC and SOC, have been listed by Viscarra
90 Rossel et al. (2006) and then by Soriano-Disla et al (2014).

91 In the MIR range, carbonates may be identified by strong and numerous absorption
92 bands, for instance, bands at approximately 820-750 cm^{-1} , 1800 cm^{-1} , 2520 cm^{-1} and
93 2900-2990 cm^{-1} (e.g., Tatzber et al., 2007; Du and Zhou, 2009, Comstock et al., 2019). In
94 the NIR range, carbonates may be identified by peaks at 2341 and 2480 nm (Lagacherie
95 et al., 2008; Barthès et al., 2016). In the MIR range, organic carbon may also be identified
96 by numerous absorption peaks, for instance, peaks at approximately 2920 and 1230 cm^{-1}
97 (Grinand et al., 2012). In the NIR range, organic carbon may be identified by absorption
98 peaks at approximately 1910 nm (Viscarra Rossel et al., 2006; Viscarra Rossel and
99 Webster, 2012) and 2050–2150 nm (Workman and Weyer, 2008). Finally, MIRS has
100 generally been reported to provide more accurate performance in terms of SOC and SIC
101 predictions compared to NIRS (e.g., McCarty et al., 2002; Reeves, 2010; Bellon-Maurel &
102 McBratney, 2011; Clairotte et al., 2016). Nevertheless, better NIRS than MIRS predictions
103 of SOC have been reported in tropical and Mediterranean regions due to the overlap of

104 absorption regions related to metal oxides and organic compounds ([Rabenarivo et al.,](#)
105 [2013](#); [Barthès et al., 2016](#)).

106 Following the emergence of infrared spectroscopy technologies for soil
107 characterisation, soil spectral libraries covering extensive areas have recently been
108 developed for estimating soil properties, especially SOC; most libraries are in the Vis-NIR
109 range at the national scale (e.g., in Australia, [Viscarra Rossel and Webster, 2012](#); in
110 Denmark, [Knadel et al., 2012](#); in France, [Gogé et al., 2014](#); in China, [Shi et al., 2015](#)),
111 continental scale (in Europe, [Stevens et al., 2013](#)) and even global scale ([Brown et al.,](#)
112 [2006](#); [Viscarra Rossel et al., 2016](#)), but libraries also exist in the MIR range (e.g., in
113 France, [Grinand et al., 2012](#), and [Clairotte et al., 2016](#); in the US, [Wijewardane et al.,](#)
114 [2018](#), [Comstosck et al., 2019](#) and [Dangal et al., 2019](#)). All these libraries use spectra
115 collected using dried and ground samples in laboratory conditions.

116 Most of the studies dealing with large soil spectral databases (national or
117 continental) have aimed to calibrate prediction models with samples from a region *A* to
118 predict the properties of soil samples from the same region *A*. Therefore, the soil and
119 climate conditions are similar between the calibration and validation databases. [McCarty et](#)
120 [al. \(2002\)](#) calibrated SOC and SIC prediction models by using two-thirds of 257 soil
121 samples collected from 14 geographically diverse locations over eight states in the west
122 central US and validated the models by using the remaining one-third of soil samples,
123 obtaining good accuracy and low bias. [Stevens et al. \(2013\)](#) used the LUCAS database,
124 which includes samples from 23 member states of the European Union, to predict SOC. In
125 their study, each calibration dataset and associated validation dataset was composed of
126 samples from a similar soil type (organic or mineral) and land use (cropland, grassland or
127 woodland), and the most spectrally representative samples were selected as calibration
128 samples. [Shi et al. \(2015\)](#) used a Chinese database, which included samples from 20
129 provinces, to predict SOC. In their study, the calibration dataset was selected to minimise

130 both the spectral distance and geographical distance to each validation sample, which also
131 belonged to the database (region *A*). [Clairotte et al. \(2016\)](#) successfully calibrated
132 prediction models by using subsets of the French national database and then tested these
133 models on a test subset (10% of this French database, either spectrally representative or
134 not; samples were not selected based on their pedo-climatic context).

135 Some studies have focused on potential and limitation analysis of prediction models
136 calibrated with a large database composed of samples collected over a region *A* to predict
137 properties for soil samples collected from a small region *b* within *A*. [McCarty et al.](#)
138 [\(2002\)](#) calibrated SOC and SIC prediction models by using 257 soil samples collected
139 from 14 geographically diverse locations over 8 states in the west central US and obtained
140 biased predictions when their models were applied to 16 independent soil samples
141 collected in another state (Nebraska) that was also in the west central US. [Gogé et al.](#)
142 [\(2014\)](#) calibrated local prediction models from the French national database, collected
143 over 550 000 km² (region *A*), to predict soil properties in soil samples collected from a
144 small French area of 24 km² (Occitanie region, south of France; this region *b* was included
145 in *A* but under-represented). Additionally, [Comstosck et al. \(2019\)](#) calibrated prediction
146 models by using a US national database (region *A*) to predict carbonate in soil samples
147 collected from two states (New York and Iowa, regions *b*) that were poorly represented in
148 the US database.

149 Finally, to the best of our knowledge, few studies have focused on potential and
150 limitation analysis of prediction models calibrated with a large database composed of
151 samples collected over a region *A* to predict properties for soil samples collected over a
152 region *B*, where *A* and *B* have no common area and hence have potential differences in
153 soil and climatic conditions. Only [Jauss et al. \(2017\)](#) and [Ahmed et al. \(2017\)](#) used MIR
154 spectroscopy as a routine method for predicting pyrogenic carbon and for predicting Total

155 Carbon and SOC, respectively, on United States soils (region *B*), from a soil spectral
156 database of Australia (region *A*).

157 The objective of this study was to analyse how MIRS may be used to predict SOC
158 and SIC contents when using a national database collected over a region *A* to predict
159 values for soil samples collected over a region *B*, where *A* and *B* have no common area
160 and have different soil and climate conditions. This study used a French MIRS soil dataset
161 (region *A*) including 2178 topsoil samples collected from an area of 550 000 km² (French
162 metropolitan territory, composed of temperate and Mediterranean soils) to calibrate SIC
163 and SOC prediction models. These models were tested on a Tunisian MIRS soil dataset
164 (region *B*), including 96 soil samples collected from an area of approximately 80 000 km²
165 (northern half of Tunisia, mainly Mediterranean and arid soils).

166

167 **2. Materials and methods**

168

169 **2.1. Soil datasets**

170 **2.1.1. The French national soil collection**

171 The national soil collection provided by the French national soil quality monitoring network
172 (RMQS; [Arrouays et al., 2002](#)) and called *DB_RMQS* was used in this study to calibrate
173 the SIC and SOC models. This RMQS collection is composed of 2178 soil samples
174 representing all main soil types encountered over the sampled 552 000 km² of the French
175 metropolitan territory (Corsica included): Cambisols, Calcisols, Luvisols, Leptosols,
176 Andosols, Albeluvisols, Podzolsols, *etc.* ([IUSS Working Group WRB, 2014](#)). The latitude of
177 sample sites ranges from 41 to 51°N, and their longitude ranges from 5.0°W to 9.5°E. The
178 sampling design was based on a square grid with 16-km spacing. At the centre of each
179 square, 25 individual core samples were taken from 0 to 30 cm depth using an unaligned

180 sampling design within a 20 × 20 m area. Core samples were bulked to obtain a composite
181 sample for each site (Arrouays et al., 2002).

182

183 **2.1.2. The Tunisian soil samples**

184 Ninety-six soil samples were used as the test set, called *DB_Tunisia*. These samples were
185 collected from 45 localities, covering approximately 80 000 km² (from 35 to 37°N and 08 to
186 11°E), with the aim of representing the main soil types and land uses of the northern half of
187 Tunisia. This was done based on previous studies carried out at the Tunis El Manar
188 University, without particular design. Field samples within the same locality were
189 kilometres apart and under different land uses. Soil samples were collected at 0-10 cm
190 using a spade, and the sampling campaign was carried out within a few months in late
191 2010. This Tunisian set was previously studied in Barthès et al. (2016).

192

193 **2.2. Laboratory Analysis**

194 **2.2.1. Physico-chemical analyses**

195 The 2178 RMQS samples were air dried, 2-mm sieved and then finely ground (< 0.25 mm)
196 using mortar and pestle. The SIC content of the RMQS samples was calculated as 0.12
197 times the soil calcium carbonate content, which was determined using these finely ground
198 (< 0.25 mm) air-dried samples using a Bernard calcimeter according to the standard
199 procedure ISO 10693 (ISO, 1995a). The carbonate content was calculated after calibration
200 with a pure calcium carbonate standard and was expressed as the equivalent calcium
201 carbonate content. The SIC content of the RMQS samples (*DB_RMQS_SIC*) ranged from
202 0 to 103.9 g kg⁻¹, averaged 6.4 g kg⁻¹, and had a median of 0 g kg⁻¹ and a skewness value
203 close to 3.1 (Table 1).

204 Then, the SOC content of the RMQS samples was calculated as the difference
205 between the total carbon (TC) and inorganic carbon contents. The TC content was

206 determined by dry combustion with an elemental analyser (Thermo Fisher Scientific CHN
207 NA2000, Waltham, MA, US) using approximately 25–30-mg aliquots of finely ground
208 (< 0.25 mm) air-dried soil samples that were sealed into tin capsules, according to the
209 standard procedure ISO 10694 (ISO, 1995b). The SOC content of the RMQS samples
210 (*DB_RMQS_SOC*) ranged from 0.6 to 411.3 g kg⁻¹, averaged 25.8 g kg⁻¹, and had a
211 median of 19.6 g kg⁻¹ and a skewness value close to 4.9 (Table 1).

212

213

[Table 1]

214

215

216 The 96 Tunisian soil samples were also air-dried, sieved to 2 mm and then finely ground
217 (< 0.2 mm) using mortar and pestle. The SOC content of the Tunisian samples was
218 analysed by dry combustion after decarbonisation using chlorhydric acid, following the
219 standard procedure ISO 10694 (1995b), with the same elemental analyser as that used for
220 RMQS samples but using silver capsules. Soils were decarbonated prior to SOC
221 determination: 10 mL of water were added to 1 g of soil and 0.5 M HCl solution was then
222 dripped onto the sample until there was no more effervescence; then the samples were
223 washed in water until pH reached 7. The SOC concentration was then determined on
224 finely ground 25–30 mg aliquots by dry combustion using an elemental analyser (Thermo
225 Fisher Scientific CHN NA2000, Waltham, MA, USA). The SOC content of the Tunisian
226 samples (*DB_Tunisia_SOC*) ranged from 2.0 to 121.0 g kg⁻¹, averaged 20.1 g kg⁻¹, and
227 had a median of 14.6 g kg⁻¹ and a skewness value close to 3 (Table 1).

228

229

230

The soil inorganic carbon content of the Tunisian samples was calculated as the
difference between the TC (determined by dry combustion using the same CHN analyser
as that used for the RMQS samples) and SOC contents. The SIC content of the Tunisian

231 samples (*DB_Tunisia_SIC*) ranged from 0.0 to 92.9 g kg⁻¹, averaged 43.3 g kg⁻¹, and had
232 a skewness value close to -0.2 (Table 1).

233 As the SIC contents of the RMQS and Tunisian samples were analysed by two
234 different methods, 29 test samples from the Tunisian set (i.e., 30% of the set) were re-
235 analysed following the same method as that used for the RMQS samples (i.e., directly, by
236 calcimetry; ISO, 1995a). The Pearson correlation coefficient (R), root mean square error
237 (RMSE) and bias between the SIC values determined by both approaches (calcimetry vs.
238 difference between TC and SOC) were 0.997, 2.4 g kg⁻¹ and 1.1 g kg⁻¹, respectively
239 (Figure 1). The values of SIC calculated by both approaches could thus be considered
240 equivalent.

241

242 [Figure 1]

243

244

245 **2.2.2. Mid-infrared spectroscopy**

246 Mid-infrared spectroscopic analysis was performed following the same procedure for both
247 spectral libraries. First, air-dried, 2-mm sieved, and 0.2-mm ground samples were oven-
248 dried at 40°C for twelve hours. Reflectance spectra were acquired using a Fourier
249 transform Nicolet 6700 spectrophotometer (Thermo Fischer Scientific, Madison, WI, US) in
250 the MIR region. Reflectance was acquired at 934 wavenumbers between 4000 and
251 400 cm⁻¹ with a 3.86 cm⁻¹ spectral resolution. This spectrophotometer is equipped with a
252 silicon carbide source, a Michelson interferometer as a dispersive element, and a
253 deuterated triglycine sulfate detector. Soil samples were placed in a 17-well plate. The soil
254 surface was flattened with the flat section of a glass cylinder, and samples were then
255 scanned using an auto-sampler (soil surface area scanned: ca. 10 mm²). Each spectrum
256 resulted from 32 co-added scans, and the body of the plate (next to the wells) was used as

257 a reference standard and scanned once per plate (i.e., every 17 samples). Twenty
258 wavenumbers were removed due to frequent noise in the spectrum, so MIR spectra in the
259 range from 4000 to 478 cm^{-1} , with 914 wavenumbers, were used.

260

261 **2.3 PLSR model calibration**

262 All procedures were performed using R software ([R Core Team, 2012](#)), and both the `ade4`
263 ([Dray and Dufour, 2007](#)) and `pls` packages ([Mevik and Wehrens, 2007](#)) were used.

264

265 **2.3.1. Dataset preparation**

266 Both *DB_RMQS_SIC* and *DB_RMQS_SOC* were divided into a calibration set (3/4 of the
267 dataset) and a validation set (1/4 of the dataset). The samples of each dataset were
268 ranked according to ascending reference value (observed SIC or SOC). The sample with
269 the lowest reference value was put in the calibration set, the next sample was put in the
270 validation set, and then the next three samples were put in the calibration set. The
271 procedure was continued by alternately placing the next sample in the validation set and
272 the following three samples in the calibration set. Following this process, the distributions
273 of the *DB_Calib_RMQS_SOC* calibration and *DB_Valid_RMQS_SOC* validation datasets
274 were similar. As well, the distributions of the *DB_Calib_RMQS_SIC* calibration and
275 *DB_Valid_RMQS_SIC* validation datasets were similar.

276 As the SOC values of the RMQS dataset followed a non-normal distribution (Table
277 1), the SOC values of the *DB_Calib_RMQS_SOC* calibration dataset were transformed
278 with natural logarithm ($\ln(\text{SOC})$) to reach a normal distribution, giving rise to new dataset
279 for calibration *DB_Calib_RMQS_InSOC*. The SIC distribution in the French dataset was
280 also non normal, but due to the very large number of null values, \ln -transformation was
281 hardly possible ([Bellon-Maurel et al., 2010](#); [Terra et al., 2015](#)).

282 The reflectance was converted into “absorbance” ($\log_{10} [1/\text{reflectance}]$), and a
283 standard normal variate correction was applied to remove additive and multiplicative
284 effects (Barnes et al., 1989). Spectral outliers, which are defined as samples spectrally
285 different than the rest of the samples (e.g., Pearson, 2002), were removed from the
286 calibration datasets. The spectral outliers were identified by applying the Mahalanobis
287 distance (Mark and Tunnell, 1985) to data condensed by principal component analysis
288 (PCA). In the present study, a Mahalanobis distance of 3.5 was selected as the threshold
289 for the identification of spectral outliers.

290

291 **2.3.2. Partial least squares regression**

292 Partial least squares regression (PLSR) is a multivariate approach that specifies a linear
293 relationship between a dependent (response) variable (Y-variable, i.e., SIC or SOC
294 content in the present case), and a set of predictor variables (X-variables, i.e., MIR spectra
295 in the present case; Tenenhaus, 1998). The general concept of PLSR is to extract a small
296 number of orthogonal variables (called the latent variables) that account for the maximum
297 variation in the X-variables. A detailed description of the PLSR procedure can be found in
298 Wold et al. (2001). This method is commonly used for NIRS or MIRS prediction of soil
299 properties (e.g., Viscarra-Rossel et al., 2006; Bellon-Maurel et al., 2010).

300 The maximum number of latent variables of PLSR was defined as 30. A leave-one-
301 out cross-validation (LOOCV) procedure was adopted to verify the prediction capability of
302 the PLSR model for the calibration set. Each time, $n - 1$ samples were used to build a
303 regression model, which was applied to the sample not used in developing the model. This
304 procedure was repeated for all n samples, resulting in predictions for all n samples.

305

306 **2.3.3. Global models**

307 Global calibration is a common calibration procedure where all calibration samples are
308 used to build a unique prediction model that is applied identically to all validation or test
309 samples. One global prediction model (denoted GM_{SIC}) was built for SIC prediction based
310 on $DB_Calib_RMQS_SIC$, validated on $DB_Valid_RMQS_SIC$ and tested on
311 $DB_Tunisia_SIC$. As well, a global prediction model (denoted GM_{SOC}) was built for SOC
312 prediction based on $DB_Calib_RMQS_SOC$, validated on $DB_Valid_RMQS_SOC$ and
313 tested on $DB_Tunisia_SOC$. Finally, a global prediction model (denoted GM_{InSOC}) was built
314 for SOC prediction based on $DB_Calib_RMQS_InSOC$, applied to spectra of
315 $DB_Valid_RMQS_SOC$ and $DB_Tunisia_SOC$, and the output predictions $\ln(SOC)$ were
316 back-transformed into SOC values using $\exp(\ln(SOC))$.

317 The optimal number of latent variables of GM_{SIC} , GM_{SOC} and GM_{InSOC} was
318 determined using prediction residual error sum of squares (PRESS) analysis of LOOCV
319 results to avoid under- and over-fitting. Then, all calibration samples were used to build the
320 prediction model with the appropriate number of latent variables, and this model was
321 applied to validation and test sets.

322

323 **2.3.4. Local models**

324 A local regression approach was implemented based on PLSR to predict the SOC and SIC
325 content of Tunisian samples. Given a sample p_i from $DB_Tunisia$ to predict:

326 1- The Pearson coefficient of correlation between the spectrum of the Tunisian sample
327 p_i and each RMQS calibration spectrum (from DB_Calib_RMQS) was calculated;

328 2- The N samples from DB_Calib_RMQS with spectra that correlated to the spectrum
329 of p_i beyond a cut-off value of 0.95 were considered spectral neighbours of the Tunisian
330 sample p_i , without maximum limit for N .

331 3- A PLSR model was built using the N spectral neighbours of the Tunisian sample p_i .

332 If a Tunisian soil sample had less than 30 spectral neighbours among the
333 *DB_Calib_RMQS* set, this soil sample was not predicted.

334 One local prediction model (denoted LM_{SIC}) was built for SIC prediction based on
335 *DB_Calib_RMQS_SIC*, validated on *DB_Valid_RMQS_SIC* and tested on
336 *DB_Tunisia_SIC*. One local prediction model (denoted LM_{SOC}) was built for SOC prediction
337 based on *DB_Calib_RMQS_SOC*, validated on *DB_Valid_RMQS_SOC* and tested on
338 *DB_Tunisia_SOC*. Finally, one local prediction model (denoted LM_{lnSOC}) was built for SOC
339 prediction based on *DB_Calib_RMQS_InSOC*, applied to spectra of
340 *DB_Valid_RMQS_SOC* and *DB_Tunisia_SOC*, and the output predictions $\ln(SOC)$ were
341 back-transformed into SOC values using $\exp(\ln(SOC))$.

342 As the calibration sets for SOC and SIC predictions did not include the same
343 samples (*DB_Calib_RMQS_SOC* and *DB_Calib_RMQS_SIC*, respectively), the nearest
344 calibration neighbours of a given French validation and Tunisian sample were not the same
345 for SOC and SIC predictions. The optimal number of latent variables was finally
346 determined using PRESS analysis of LOOCV on the selected spectral neighbours to avoid
347 under- and over-fitting. Regardless of the type of local model, the spectral outliers were not
348 investigated because the selection of nearest neighbours was considered an implicit
349 rejection of outliers.

350

351 **2.4 PLSR model evaluation**

352 The performance of global models was evaluated according to figures of merit described in
353 [Bellon Maurel et al. \(2010\)](#), from cross-validation, validation and test databases.

354 The coefficient of determination of cross-validation (R^2_{cv}) and root mean square
355 error of cross-validation ($RMSECV$) for *DB_Calib_RMQS* were used. R^2_{cv} was computed
356 as $1-ESS/TSS$, where ESS is the error sum of squares and TSS the total sum of squares.

357 The coefficient of determination and root mean square error of prediction for
358 *DB_Valid_RMQS*, R^2_{val} and $RMSE_{val}$ respectively, were used. R^2_{val} was also computed as
359 1-ESS/TSS. The ratio of performance to deviation in *DB_Valid_RMQS* (RPD_{val}), which is
360 the ratio between the standard deviation in *DB_Valid_RMQS* and $RMSE_{val}$, was
361 calculated. The ratio of performance to interquartile range of *DB_Valid_RMQS* ($RPIQ_{val}$),
362 which is the ratio between interquartile range (difference between the third and first
363 quartiles) of *DB_Valid_RMQS* and $RMSE_{val}$, was also calculated. This parameter has been
364 proposed for variables with non-normal distributions (Bellon-Maurel et al., 2010). And the
365 bias, which is the mean difference between observations and predictions, was calculated
366 for *DB_Valid_RMQS* ($bias_{val}$).

367 The coefficient of determination and root mean square error of prediction for
368 *DB_Tunisia*, R^2_{test} and $RMSE_{test}$ respectively, were used. R^2_{test} was computed as 1-
369 ESS/TSS. The ratio of performance to deviation in *DB_Tunisia* (RPD_{test}), which is the ratio
370 between the standard deviation in *DB_Tunisia* and $RMSE_{test}$, was calculated. The ratio of
371 performance to interquartile range of *DB_Tunisia* ($RPIQ_{test}$), which is the ratio between the
372 interquartile range of *DB_Tunisia* and $RMSE_{test}$, was also calculated. And the bias, which is
373 the mean difference between observations and predictions, was calculated for *DB_Tunisia*
374 ($bias_{test}$).

375 Finally, a wavelength was considered a significant contributor in a global model
376 when the values of both the regression coefficient and variable importance in the
377 projection (VIP) were sufficiently large: the threshold for the VIP was set to 1 (Chong and
378 Jun, 2005; Wold et al., 1993, 2001), and the thresholds for the regression coefficients were
379 their standard deviations (Viscarra-Rossel et al., 2008).

380 The performances of local models were based on the same figures of merit as
381 those used in global calibration, calculated on *DB_Valid_RMQS* and *DB_Tunisia*.

382 Concerning the models built from the *DB_Calib_RMQS_InSOC*, independent validation
383 and test statistics were calculated from back-transformed data.

384

385

386 **4. Results**

387 **4.1. Preliminary analysis of soil properties and spectra**

388 The RMQS soil sampling covered the French territory, and ranges of SOC and SIC
389 contents in the *DB_RMQS_SOC* and *DB_RMQS_SIC*, respectively, are large (Table 1).
390 According to the French soil classification, 33 soil reference groups were sampled, with a
391 dominance of Cambisols ([IUSS Working Group WRB, 2014](#); 27% of the sample set),
392 calcareous soils (Calcosols, 22%) and Luvisols (16%). High SIC values are mainly located
393 in three French areas: 1) the southeast (Prealps), mainly with Leptosols and Calcosols, 2)
394 the northeast (chalk Champagne) also mainly with Leptosols, and 3) a transect from west
395 (the Aquitanian Basin) to south (Mediterranean Sea), mainly with Calcosols (Figure 2A1).
396 High SOC values are mainly located in 1) mountain areas (Alps in the southeast, Pyrenees
397 in the extreme southwest, Massif Central in the south-centre, Jura in the centre-east), 2)
398 cool regions covered by forests and pastures (centre-east), and 3) intensive livestock
399 production areas (northwest; Figure 2B1).

400 The Tunisian soil samples covered the northern half of Tunisian territory, and the
401 SIC contents range in the *DB_Tunisia_SIC* is as large as the one in the *DB_RMQS_SIC*,
402 whereas the SOC contents range in the *DB_Tunisia_SOC* is lower than the one of the
403 *DB_RMQS_SOC* (Table 1). The sampled Tunisian soils were mainly Calcaric Cambisols
404 and Regosols, Kastanozems, and Chromic and Vertic Cambisols.

405

406

[Figure 2]

407

408 Principal component analyses were performed on pre-treated spectra of
409 *DB_Calib_RMQS_SIC* and *DB_Calib_RMQS_SOC*, respectively, and pre-treated Tunisian
410 spectra were projected onto the plans made by the first and second components. Most
411 Tunisian spectra overlapped a subset of RMQS spectra for SIC (Figure 3a) and for SOC
412 (Figure 3b). So most Tunisian soil samples had similar spectral signatures than a subset of
413 RMQS spectra used for calibrating prediction models.

414

415 [Figure 3]

416

417

4.2. Global models

418

a. Soil inorganic carbon content

419 For SIC prediction, 52 spectral outliers were identified within the initial calibration dataset,
420 so 1582 RMQS samples were ultimately kept and constituted the *DB_Calib_RMQS_SIC*
421 dataset. The SIC content of these 1582 RMQS samples contained in
422 *DB_Calib_RMQS_SIC* ranged from 0.0 to 103.9 g kg⁻¹, averaged 6.4 g kg⁻¹, and had a
423 skewness value close to 3.1 (Table 1). The SIC content of the 544 RMQS samples
424 contained in *DB_Valid_RMQS_SIC* ranged from 0.0 to 95.2 g kg⁻¹, averaged 6.3 g kg⁻¹,
425 and also had a skewness value close to 3.1 (Table 1).

426 The GM_{SIC} was built from the *DB_Calib_RMQS_SIC* dataset using an optimal
427 number of 15 latent variables, validated on the *DB_Valid_RMQS_SIC* dataset and then
428 tested on the *DB_Tunisia_SIC* dataset. The performance of the GM_{SIC} prediction model
429 was accurate, with an R^2_{cv} of 0.97 and $RMSECV$ of 2.8 g kg⁻¹ in the calibration step (Table
430 2) and an R^2_{val} of 0.98 and $RMSE_{val}$ of 2.1 g kg⁻¹ in the validation step (Table 2, Figure 4b).
431 When applied to the *DB_Tunisia_SIC* dataset, this GM_{SIC} prediction model provided
432 accurate and unbiased predictions ($R^2_{test} = 0.96$, $RMSE_{test} = 5.2$ g kg⁻¹ and $bias_{test} = 0.2$
433 g kg⁻¹; Table 3, Figure 4c).

434 A total of 96 spectral bands might be considered significant based on the analysis of
435 the VIP and regression coefficients of GM_{SOC} (Figure 4d). Among the 96 significant spectral
436 bands, those at approximately 2500, 1800 and 860 cm^{-1} had regression coefficients higher
437 than 3 times the standard deviation and therefore might be considered the most significant
438 ones.

439

440 [Figure 4]

441

442 [Table 2]

443

444 [Table 3]

445

446 **b. Soil organic carbon content**

447 For SOC prediction, 48 spectral outliers were identified within the initial calibration dataset,
448 so 1586 RMQS samples were ultimately kept and constituted the *DB_Calib_RMQS_SOC*
449 and *DB_Calib_RMQS_InSOC* datasets. The SOC contents of the 1586 RMQS samples of
450 *DB_Calib_RMQS_SOC* ranged from 0.6 to 411.3 g kg^{-1} , averaged 25.6 g kg^{-1} , and had a
451 skewness value close to 5.5 (Table 1). The $\ln(\text{SOC})$ values of the 1586 RMQS samples of
452 *DB_Calib_RMQS_InSOC* ranged from -0.5 to 6 $\ln(\text{g kg}^{-1})$, averaged 3 $\ln(\text{g kg}^{-1})$, had a
453 median of 3 $\ln(\text{g kg}^{-1})$ and a skewness value close to 0.3. The SOC content of the 544
454 RMQS samples contained in *DB_Valid_RMQS_SOC* ranged from 1.5 to 159 g kg^{-1} ,
455 averaged 25.5 g kg^{-1} , and had a skewness value close to 2.5 (Table 1).

456 The GM_{SOC} was built from the *DB_Calib_RMQS_SOC* dataset using 23 latent
457 variables, validated on the *DB_Valid_RMQS_SOC* dataset and then tested on the
458 *DB_Tunisia_SOC* dataset. The performance of the GM_{SOC} prediction model was modest,
459 with an R^2_{cv} of 0.80 and $RMSECV$ of 9.9 g kg^{-1} in the calibration step (Table 2) and an R^2_{val}

460 of 0.88 and $RMSE_{val}$ of 7.2 g kg⁻¹ in the validation step (Table 2, Figure 5b). When applied
461 to the Tunisian test set, this GM_{SOC} prediction model provided low accuracy ($R^2_{test} = 0.64$,
462 $RMSE_{test} = 16.0$ g kg⁻¹) and biased predictions ($bias_{test} = -5.2$ g kg⁻¹) (Table 3, Figure 5c).

463 A total of 92 spectral bands might be considered significant based on the analysis of
464 the VIP and regression coefficients of GM_{SOC} (Figure 5d). Nevertheless, among these 92
465 spectral bands, none was associated to very high regression coefficients.

466

467

[Figure 5]

468

469 A GM_{lnSOC} prediction model was built from the *DB_Calib_RMQS_InSOC* dataset using 10
470 latent variables and this model was applied to spectra of both *DB_Valid_RMQS_SOC* and
471 *DB_Tunisia_SOC* datasets. Finally, the ln(SOC) predictions were back-transformed into
472 SOC values for calculating the figures of merit. The performance of the GM_{lnSOC} prediction
473 model was accurate, with an R^2_{cv} of 0.89 and $RMSECV$ of 0.2 g kg⁻¹ in the calibration step
474 (Table 2) and an R^2_{val} of 0.90 and $RMSE_{val}$ of 6.6 g kg⁻¹ in the validation step (Table 2,
475 Figure 6b). When applied to the Tunisian test set, the GM_{lnSOC} prediction model provided
476 high accuracy ($R^2_{test} = 0.97$, $RMSE_{test} = 4.2$ g kg⁻¹) and very slightly biased predictions
477 ($bias_{test} = 0.7$ g kg⁻¹) (Table 3, Figure 6c).

478 A total of 95 spectral bands might be considered significant based on the analysis of
479 the VIP and regression coefficients in the GM_{lnSOC} (Figure 6d). Among these 95 significant
480 spectral bands, those at approximately 2915, and 1800 cm⁻¹ had regression coefficients
481 higher than 3 times the standard deviation and therefore might be considered the most
482 significant ones.

483

484

[New Figure 6]

485

486

4.3. Local models

487

a. Soil inorganic carbon content

488 All validation soil samples had more than 30 spectral neighbours within the
489 *DB_Calib_RMQS_SIC* dataset, so SIC could be predicted from LM_{SIC} for all samples of
490 *DB_Valid_RMQS_SIC*. The LM_{SIC} provided very accurate and unbiased SIC predictions in
491 validation ($R^2_{val} = 0.99$ and $RMSE_{val} = 1.8 \text{ g kg}^{-1}$; Table 2). Therefore, this LM_{SIC} provided
492 validation performance slightly better than that of GM_{SIC} (Table 2).

493 All Tunisian soil samples had more than 30 spectral neighbours within the
494 *DB_Calib_RMQS_SIC* dataset, so SIC could be predicted from LM_{SIC} for all Tunisian
495 samples. The LM_{SIC} provided accurate and slightly biased SIC predictions on Tunisian
496 samples ($R^2_{test} = 0.96$, $RMSE_{test} = 5.6 \text{ g kg}^{-1}$ and $bias_{test} = 1.7 \text{ g kg}^{-1}$; Table 3, Figure 7a).
497 Therefore, LM_{SIC} provided test performance slightly lower than that of GM_{SIC} , mainly due to
498 bias (Table 3). The number of latent variables selected for LM_{SIC} on Tunisian samples
499 varied depending on the sample predicted and followed a relatively normal distribution
500 centred at approximately 13, which was close to the optimal number of latent variables
501 selected by the GM_{SIC} (Figure 7c).

502 The number of spectral neighbours of Tunisian samples varied from 65 to 1293
503 (Figure 7b). Only a slight trend was observed between the number of neighbours and the
504 prediction error, with a lower error when the number of neighbours increased (Figure 7b).
505 This trend could be expected, as the use of a higher number of neighbour samples for
506 calibration should result in more accurate predictions. The spectra from
507 *DB_Calib_RMQS_SIC* used for building the 96 local individual Tunisian models LM_{SIC}
508 were selected from 0 to 87 times. So no spectrum from *DB_Calib_RMQS_SIC* was
509 systematically selected, whereas only 1.6% of spectra from *DB_Calib_RMQS_SIC* were
510 never selected. The frequently selected samples were mainly located in SIC-richest areas
511 (Calcosols and Leptosols) such as the southeast (Prealps), the northeast (chalk

512 Champagne) and the transect from west (the Aquitanian Basin) to south (Mediterranean
513 Sea) (Figure 2A1 and 2A2).

514

515

[Figure 7]

516

517

b. Soil organic carbon content

518 As *DB_Calib_RMQS_SOC* and *DB_Calib_RMQS_InSOC* datasets contained same
519 predictors X-variables (MIR spectra), spectral neighbours of validation samples were the
520 same for LM_{SOC} and LM_{InSOC} models. All validation soil samples had more than 30 spectral
521 neighbours within the *DB_Calib_RMQS_SOC* dataset and *DB_Calib_RMQS_InSOC*
522 datasets, so SOC and ln(SOC) could be predicted from LM_{SOC} and LM_{InSOC} , respectively,
523 for all samples of *DB_Valid_RMQS_SOC*. The LM_{SOC} provided accurate and very slightly
524 biased SOC predictions in validation ($R^2_{val} = 0.93$, $RMSE_{val} = 5.4 \text{ g kg}^{-1}$ and $bias_{val} = -0.7$
525 g kg^{-1} ; Table 2). Therefore, LM_{SOC} provided validation performance higher than that of
526 GM_{SOC} (Table 2). The LM_{InSOC} provided accurate and unbiased SOC predictions in
527 validation ($R^2_{val} = 0.92$, $RMSE_{val} = 5.7 \text{ g kg}^{-1}$ and $bias_{val} = -0.1 \text{ g kg}^{-1}$; Table 2). Therefore,
528 LM_{InSOC} provided validation performance higher than that of GM_{SOC} and almost similar to
529 that of LM_{SOC} (Table 2).

530 All Tunisian soil samples also had more than 30 spectral neighbours within the
531 *DB_Calib_RMQS_SOC* and *DB_Calib_RMQS_InSOC* datasets, so SOC and ln(SOC)
532 could be predicted from LM_{SOC} and LM_{InSOC} , respectively, for all Tunisian samples. The
533 LM_{SOC} provided accurate and very slightly biased SOC predictions on Tunisian samples
534 ($R^2_{test} = 0.89$, $RMSE_{test} = 6.9 \text{ g kg}^{-1}$ and $bias_{test} = 0.5 \text{ g kg}^{-1}$; Table 3, Figure 8a). Therefore,
535 this LM_{SOC} provided test performance markedly higher than that of GM_{SOC} (Table 3). The
536 number of latent variables selected by this LM_{SOC} depending on the sample followed a
537 bimodal distribution centred at approximately 16 and 22 (Figure 8c). The second peak of

538 number of latent variables, approximately 22, was close to the number of latent variables
539 selected by the GM_{SOC} (Figure 8c).

540 The LM_{InSOC} model provided accurate and very slightly biased SOC predictions on
541 Tunisian samples ($R^2_{test} = 0.93$, $RMSE_{test} = 5.8 \text{ g kg}^{-1}$ and $bias_{test} = -0.6 \text{ g kg}^{-1}$; Table 3,
542 Figure 9a). The LM_{InSOC} provided test performance markedly higher than that of GM_{SOC} but
543 lower than that of GM_{InSOC} (Table 3). The number of latent variables selected by this
544 LM_{InSOC} depending on the sample followed a relatively normal distribution centred at
545 approximately 11 (Figure 9c), which was close to the optimal number of latent variables
546 selected by the GM_{InSOC} .

547 As $DB_Calib_RMQS_SOC$ and $DB_Calib_RMQS_InSOC$ datasets contained same
548 predictors X-variables (MIR spectra), spectral neighbours of Tunisian samples were the
549 same for LM_{SOC} and LM_{InSOC} models. The number of spectral neighbours of Tunisian
550 samples varied from 65 to 1292 (Figure 8b and 9b), and no clear trend was observed
551 between the number of neighbours and the prediction error obtained by LM_{SOC} and
552 LM_{InSOC} models (Figure 8b and 9b, respectively). As for SIC prediction with LM_{SIC} , the
553 spectra from $DB_Calib_RMQS_SOC$ and $DB_Calib_RMQS_InSOC$ datasets used for
554 building the 96 individual local Tunisian models were selected from 0 to 87 times. So no
555 spectrum from $DB_Calib_RMQS_SOC$ and $DB_Calib_RMQS_InSOC$ datasets was
556 systematically selected, whereas 1.7% of spectra from $DB_Calib_RMQS_SOC$ and
557 $DB_Calib_RMQS_InSOC$ datasets were never selected. The frequently selected samples
558 were mainly located in SOC-poor areas (Calcisols and Leptosols) in the northeast (chalk
559 Champagne) and on the transect from west (Aquitainian Basin) to south (Mediterranean
560 Sea) and in soils richer in SOC in the southeast (Prealps) (Figure 2B1 and 2B2).

561

562

[Figure 8]

563

564

565

566

567 **5. Discussion**

568 **5.1. Global models built on region A for application to region A.**

569 Before being applied to the Tunisian database, the global models were calibrated with an
570 RMQS subset (*DB_Calib_RMQS_SIC*, *DB_Calib_RMQS_SOC* and
571 *DB_Calib_RMQS_InSOC*) and validated on an RMQS subset (*DB_Valid_RMQS_SIC* and
572 *DB_Valid_RMQS_SOC*), which means that the global models were calibrated using soil
573 samples collected over a region A to predict values for soil samples collected over this
574 same region A.

575 The validation performance of GM_{SIC} was in accordance with results reported in the
576 literature. [Grinand et al. \(2012\)](#) obtained similar performances using the same RMQS
577 MIRS database, with R^2_{val} and RPD_{val} values of 0.97 and 7.6, respectively, when 3/4 of the
578 set, selected at random, was used for calibration and 1/4 was used for validation. [Barthès
579 et al. \(2016\)](#) obtained similar performances for SIC using the Tunisian MIRS database for
580 both calibration and prediction, with R^2_{cv} and RPD_{cv} values of 0.98 and 7.8, respectively.
581 [Mc Carty et al. \(2002\)](#) obtained similar performances using another MIRS database
582 collected in the US, with $R^2_{val} = 0.98$ (RPD_{val} was not mentioned and could not be
583 calculated). The most significant spectral bands for GM_{SIC} , located at 2500, 1800 and 860
584 cm^{-1} (Figure 4d), might be attributed to stretching or bending vibrations in carbonate
585 molecules, as suggested by [Du and Zhou \(2009\)](#) and then by [Grinand et al. \(2012\)](#).

586 The validation performance of the GM_{SOC} was also in accordance with some
587 literature results. [Clairotte et al. \(2016\)](#) obtained very similar performances using the full
588 RMQS MIRS database (including two depth layers, 0-30 and 30-50 cm, instead of one in
589 the present study), with R^2_{val} and RPD_{val} values of 0.88 and 2.7, respectively. [Barthès et al.](#)

590 (2016) obtained better performances using the Tunisian MIRS database with cross-
591 validation, with R^2_{cv} and RPD_{cv} values of 0.95 and 4.3, respectively. Moreover, [Mc Carty et](#)
592 [al. \(2002\)](#) obtained slightly higher performances using a MIRS database collected in the
593 US for both SOC calibration and validation, with $R^2_{val} = 0.94$ (RPD_{val} was not mentioned
594 and could not be calculated).

595 Following previous researches dealing with non-normal distribution of soil properties
596 (e.g., [Waruru et al., 2014](#); [Dangal et al., 2019](#)), natural logarithm transformation was
597 applied to the highly skewed SOC values of the RMQS database to reach a normal
598 distribution in the calibration dataset (*DB_Calib_RMQS_InSOC*). Thanks to this normal
599 distribution, the performance of the GM_{InSOC} on *DB_Valid_RMQS_SOC* was slightly better
600 than the one of the GM_{SOC} .

601 Finally, validation performance was higher for SIC with GM_{SIC} (R^2_{val} and RPD_{val}
602 values of 0.98 and 7.6, respectively; Table 2) than for SOC prediction with GM_{SOC} (R^2_{val}
603 and RPD_{val} values of 0.88 and 2.7, respectively; Table 2) and GM_{InSOC} (R^2_{val} and RPD_{val}
604 values of 0.90 and 2.9, respectively; Table 2), confirming that MIRS allows markedly more
605 accurate predictions of SIC than SOC as also shown by [McCarty et al. \(2002\)](#), [Grinand et](#)
606 [al. \(2012\)](#) and [Barthès et al. \(2016\)](#).

607

608 **5.2. Models built on region A for application to region B.**

609 The models were calibrated by using a RMQS subset (*DB_Calib_RMQS_SIC*,
610 *DB_Calib_RMQS_SOC* and *DB_Calib_RMQS_InSOC*) and tested on the Tunisian
611 dataset, which means that the models were calibrated by using soil samples collected over
612 a region A to predict values for soil samples collected over a region B, where A and B had
613 no common area, so the soil and climate conditions were different between the calibration
614 and test datasets. Our results showed that GM_{SIC} provided accurate test performance
615 (R^2_{test} and RPD_{test} values of 0.96 and 4.9, respectively; Table 3), which was however lower

616 when applied to region *B* than to region *A* (Figure 4c, Tables 2 and 3). The GM_{SOC} also
617 provided markedly lower performance (R^2_{test} and RPD_{test} values of 0.64 and 1.3,
618 respectively; Table 3) when applied to our region *B* than to our region *A* (Figures 5c, Table
619 2 and 3).

620 The RMQS spectra used for SIC and SOC predictions by local models were
621 selected using a similarity measure (Pearson's coefficients of correlation), following the
622 same approach than [Shenk et al. \(1997\)](#) and [Nocita et al. \(2014\)](#). So the driver of
623 neighbours selection was the spectral similarity between Tunisian and French spectra. The
624 LM_{SIC} did not improve the SIC prediction accuracy compared to the GM_{SIC} (Table 3).
625 Therefore, GM_{SIC} seemed robust and did not need to be adjusted to spectral particularities
626 of region *B*. So rather than spectral similarity between French and Tunisian samples, the
627 main reason for accurate SIC predictions in region *B* seemed to be the strong spectral
628 features of SIC in the MIR region, as suggested by [Gogé et al. \(2014\)](#). The LM_{SOC}
629 improved SOC prediction accuracy compared to the GM_{SOC} (Table 3). Therefore, GM_{SOC}
630 seemed poorly robust and the calibration over region *A* needed to be adjusted to spectral
631 particularities of region *B* using spectral neighbours (e.g., [Shenk et al., 1997](#); [Nocita et al.,](#)
632 [2014](#)). The increase in the performance of MIRS-based SOC prediction when shifting from
633 global to local PLSR is in accordance with literature (e.g. [Ramirez-Lopez et al., 2013](#); [Shi](#)
634 [et al., 2015](#); [Clairotte et al., 2016](#) and [Dangal et al., 2019](#)). Finally, the frequently selected
635 spectral neighbours of Tunisian samples by the LM_{SOC} were mainly located in SOC-poor
636 areas in the northeast (chalk Champagne) and on the transect from west (Aquitainian
637 Basin) to south (Mediterranean Sea) and in soils richer in SOC in the southeast (Prealps)
638 (Figure 2B1 and 2B2). So these frequently selected spectral neighbours of Tunisian
639 samples by SOC local model were not located only over the more similar climatic and
640 pedological contexts such as the Mediterranean context (southeast of France).

641

642 **5.3. Impact of the SOC In-transformation on models built on region A for application**
643 **to region B.**

644 The global model calibrated on region A with log-transformed SOC values (GM_{InSOC})
645 provided accurate performance on region B (R^2_{test} and $RMSE_{test}$ values of 0.97 and 4.9 g
646 kg^{-1} , respectively; Figure 6c, Table 3). The local model provided better performance when
647 calibrated on region A with log-transformed SOC values (LM_{InSOC}) than with SOC values
648 (LM_{SOC}), but less accurate than GM_{InSOC} on region B (R^2_{test} and $RMSE_{test}$ values of 0.93
649 and 3.6 g kg^{-1} , respectively; Figure 9c, Table 3). So whatever the model using log-
650 transformed SOC data in calibration database (GM_{InSOC} or LM_{InSOC}), the SOC predictions
651 on region B were improved compared to models using highly skewed SOC values in
652 calibration database (GM_{SOC} or LM_{SOC} ; Table 3), as showed by [Jaconi et al. \(2019\)](#).

653 Finally, and unexpectedly, GM_{InSOC} provided better performance than LM_{SOC} when
654 applied on Tunisian samples. So a transformation of calibration SOC values improved
655 SOC model performance more clearly than spectral selection of calibration samples
656 (neighbours). Therefore SOC model performance was more sensitive to the distribution of
657 the explained variable of the calibration samples than to spectral similarity between
658 calibration and test spectra.

659

660 **5.4. Perspectives**

661 This study, which used MIRS to predict SOC and SIC contents by using a database
662 collected over a region A to predict values over a region B, where A and B have no
663 common area, could be continued with a study to develop predictions based on selected
664 spectral bands. Indeed, spectral band selection remains to be explored, as several studies
665 testing such an approach have obtained different results. [Viscarra Rossel and Lark \(2009\)](#)
666 successfully used wavelets and a variable selection technique to improve SOC calibration
667 using Vis–NIR and MIRS data. Additionally, [Volhand et al. \(2016\)](#) outperformed SOC

668 predictions based on Vis–NIR spectra by using band selection. However, [Stevens et al.](#)
669 [\(2013\)](#) tested recursive feature elimination based on the random forest approach and
670 obtained no overall increase in the accuracy of soil property prediction using the LUCAS
671 (European) Vis-NIR soil database compared to models using all spectral bands. In
672 addition, [Yang et al. \(2019\)](#) tested a generic algorithm for spectral band selection but
673 obtained no overall increase in SOC prediction accuracy. As previously tested by, e.g.,
674 [Guerrero et al. \(2014\)](#) and [Guy et al. \(2016\)](#), spiking could be another useful approach to
675 improve prediction accuracy when applying large-scale calibrations to small regions.
676 Spiking consists of adding a small subset of samples from region *B* (spiking subset) to the
677 dataset from region *A* to recalibrate a model.

678 As well, this study could be continued with an impact analysis of the selection of
679 spectral neighbours. Both the number of spectral neighbours and the procedure to select
680 them could be analysed. Several approaches are available for selecting representative
681 calibration samples ([Shetty et al., 2012](#)) and could also be tested. For example, to analyse
682 the spectral similarity between calibration and test spectra, the Pearson correlation
683 coefficient between spectra could be replaced by the Mahalanobis distance between
684 spectra ([Nocita et al., 2014](#)) or the Pearson correlation coefficient distance based on Fast
685 Fourier Transform of spectra ([Gogé et al., 2012](#)). Finally, some covariates could be added
686 in local regression to improve prediction accuracy, as previously tested by [Nocita et al.](#)
687 [\(2014\)](#) who used clay contents of samples as covariates to predict SOC content.

688

689

690 **6. Conclusion**

691 This work highlighted that, as expected, the SOC and SIC contents of French samples
692 were successfully predicted from the French MIRS soil database using a global model
693 based on PLS regression. Predictions of SIC and SOC are accurate when the calibration

694 and validation samples come from same pedologic and climatic contexts. This work also
695 highlighted that when the calibration and validation samples come from different pedologic
696 and climatic contexts, the SOC prediction performance over validation samples decreases,
697 whereas the SIC prediction performance remains accurate. Finally, this work showed that
698 prediction models were more sensitive to the distribution of the explained variables of
699 calibration samples than to the spectral similarity between calibration and test spectra.
700 This study confirmed the very high applicability of MIRS for SIC determination and the
701 robustness of SIC prediction models, even when the calibration and validation samples
702 come from different contexts.

703

704

705 **Acknowledgment**

706 RMQS soil sampling and physico-chemical analyses were supported by the GIS Sol,
707 which is a scientific group of interest on soils involving the French Ministry for ecology and
708 sustainable development, the French Ministry of agriculture, the French National institute
709 for geographical and forest information (IGN), the French Environment and Energy
710 Management Agency (ADEME), the French agency for biodiversity, French Institute for
711 Research and Development (IRD) and the French Institute for Agronomic Research (INRA,
712 which is a French public research organisation dedicated to agriculture, food and
713 environment). We thank all the people involved in sampling RMQS sites and in samples
714 preparation. D.A. is coordinator of the Research Consortium GLADSOILMAP, supported
715 by LE STUDIUM Loire Valley Institute for Advanced Studies. And we acknowledge the
716 contribution of three anonymous reviewers for their help in improving the manuscript.

717

718

719 **References**

- 720 Apestegui, M., Plante, A.F., Virto, I., 2018. Methods assessment for organic and
721 inorganic carbon quantification in calcareous soils of the Mediterranean region.
722 *Geoderma Regional* 12, 39–48.
- 723 Ahmed, Z.U., Woodbury, P.B., Sanderman, J., Hawke, B., Jauss, V., Solomon, D.,
724 Lehmann, J., 2017. Assessing soil carbon vulnerability in the Western USA by
725 geospatial modeling of pyrogenic and particulate carbon stocks. *Journal of*
726 *Geophysical Research: Biogeosciences* 122, 354–369.
- 727 Arrouays, D., Jolivet, C., Boulonne, L., Bodineau, G., Saby, N., Grolleau, E., 2002. A new
728 initiative in France: a multi-institutional soil quality monitoring network. *Comptes*
729 *Rendus de l'Académie d'Agriculture de France* 88, 93–105.
- 730 Barnes, R.J., Dhanoa, M.S., Lister, S.J., 1989. Standard normal variate transformation and
731 de-trending of near-infrared diffuse reflectance spectra. *Applied Spectroscopy*. 43,
732 772–777.
- 733 Barthès, B.G., Kouakoua, E., Moulin, P., Hmaid, K., Gallali, T., Clairotte, M., Bernoux, M.,
734 Bourdon, E., Toucet, J., Chevallier, T., 2016. Studying the physical protection of soil
735 carbon with quantitative infrared spectroscopy. *Journal of Near Infrared Spectroscopy*
736 24, 199–214.
- 737 Batjes, N. H., 1996. Total carbon and nitrogen in the soils of the world. *European Journal*
738 *of Soil Science* 47, 151–163.
- 739 Bellon-Maurel, V., Fernandez-Ahumada, E., Palagos, B., Roger, J.M., McBratney, A., 2010.
740 Prediction of soil attributes by NIR spectroscopy. A critical review of chemometric
741 indicators commonly used for assessing the quality of the prediction. *Trends in*
742 *Analytical Chemistry (TRAC)* 29(9), 1073–1081.
- 743 Bellon-Maurel, V., McBratney, A.B., 2011. Near-infrared (NIR) and mid-infrared (MIR)
744 spectroscopic techniques for assessing the amount of carbon stock in soils — Critical
745 review and research perspectives. *Soil Biology and Biochemistry* 43, 1398–1410.

746 Bernoux, M., Chevallier, T., 2014. Carbon in Drylands. Multiple Essential Functions. Les
747 Dossier Thématiques du CSFD. N°10. CSFD/Agropolis International, Montpellier,
748 France.

749 Brown, D.J., Shepherd, K.D., Walsh, M.G., Dewayne Mays, M., Reinsch, T.G., 2006.
750 Global soil characterization with VNIR diffuse reflectance spectroscopy. *Geoderma*
751 132, 273–290.

752 Cécillon, L., Barthès, B.G., Gomez, C., Ertlen, D., Genot, V., Hedde, M., Stevens, A., Brun,
753 J.J., 2009. Assessment and monitoring of soil quality using near-infrared reflectance
754 spectroscopy (NIRS). *European Journal of Soil Science* 60(5), 770–784.

755 Chen, Y., Barak, P., 1982. Iron nutrition of plants in calcareous soils. *Advances in*
756 *Agronomy* 35, 217–240.

757 Chong, I.G., Jun, C.H., 2005. Performance of some variable selection methods when
758 multicollinearity is present. *Chemometrics and Intelligent Laboratory Systems* 78,
759 103–112.

760 Clairotte, M., Grinand, C., Kouakoua, E., Thébault, A., Saby, N.P.A., Bernoux, M., Barthès,
761 B.G., 2016. National calibration of soil organic carbon concentration using diffuse
762 infrared reflectance spectroscopy. *Geoderma* 276, 41–52.

763 Comstock, J.P., Sherpa, S.R., Ferguson, R., Bailey, S., Beem-Miller, J.P., Lin, F., Lehmann,
764 J., Wolfe, D.W., 2019. Carbonate determination in soils by mid-IR spectroscopy with
765 regional and continental scale models. *PLoS ONE* 14(2): e0210235.

766 Dangal, S.R.S., Sanderman, J., Wills, S., Ramirez-Lopez, L., 2019. Accurate and precise
767 prediction of soil properties from a large mid-infrared spectral library. *Soil Systems* 3,
768 11.

769 Dray, S., Dufour, A.B., 2007. The ade4 package: implementing the duality diagram for
770 ecologists. *Journal of Statistical Software* 22, 1–20.

771 Du, C., Zhou, J., 2009. Evaluation of soil fertility using infrared spectroscopy: a review.
772 Environmental Chemistry Letters 7, 97–113.

773 Gogé, F., Joffre, R., Jolivet, C., Ross, I., Ranjard, L., 2012. Optimization criteria in sample
774 selection step of local regression for quantitative analysis of large soil NIRS
775 database. Chemometrics and Intelligent Laboratory Systems 110, 168–176.

776 Gogé, F., Gomez, C., Jolivet, C., Joffre, R., 2014. Which strategy is best to predict soil
777 properties of a local site from a national Vis–NIR dataset? Geoderma 213, 1–9.

778 Gregorich E.G., Carter M.R., Angers D.A., Monreal C.M. et Ellert B.H., 1994. Towards a
779 minimum data set to assess soil organic matter quality in agricultural soils. Canadian
780 Journal of Soil Science 74, 367–385.

781 Grinand, C., Barthès, B.G., Brunet, D., Kouakoua, E., Arrouays, D., Jolivet, C., Caria, G.,
782 Bernoux, M., 2012. Prediction of soil organic and inorganic carbon contents at a
783 national scale (France) using mid-infrared reflectance spectroscopy (MIRS).
784 European Journal of Soil Science 63(2), 141–151.

785 Guerrero, C., Stenberg, B., Wetterlind, J., Viscarra Rossel, R.A., Maestre, F.T., Mouazen,
786 A.M., Zornoza, R., Ruiz-Sinoga, J.D., Kuang, B., 2014. Assessment of soil organic
787 carbon at local scale with spiked NIR calibrations: effects of selection and extra-
788 weighting on the spiking subset. European Journal of Soil Science 65, 248–263.

789 Guy, A.L., Siciliano, S.D., Lamb, E.G., 2015. Spiking regional vis-NIR calibration models
790 with local samples to predict soil organic carbon in two High Arctic polar deserts using
791 a vis-NIR probe. Canadian Journal of Soil Science 95, 237–249.

792 Hannam, K.D., Kehila, D., Millard, P., Midwood, A.J., Neilson, D., Neilson, G.H., Forge, T.
793 A., Nichol, C., Jones, M.D., 2016. Bicarbonates in irrigation water contribute to
794 carbonate formation and CO₂ production in orchard soils under drip irrigation,
795 Geoderma 266, 120–126.

796 Harris, D., Horwath, W.R., van Kessel, C., 2001. Acid fumigation of soils to remove
797 carbonates prior to total organic carbon or carbon-13 isotopic analysis. *Soil Science*
798 *Society of America Journal* 65, 1853–1856.

799 ISO (International Organization for Standardisation), 1995a. ISO 10693:1995 –
800 Determination of Carbonate Content – Volumetric Method. ISO, Geneva.

801 ISO (International Organization for Standardisation), 1995b. ISO 10694:1995 - Soil Quality
802 - Determination of Organic and Total Carbon after Dry Combustion (Elementary
803 Analysis). ISO, Geneva.

804 IUSS Working Group WRB, 2014. International Union of Soil Sciences, Working Group
805 World Reference Base for Soil Resources. World Reference Base for Soil Resources
806 2014. International Soil Classification System for Naming Soils and Creating Legends
807 of Soil Maps. FAO, Rome.

808 Jacobson, M., Charlson, R.J., Rodhe, H., Orians, G.H., 2000. *Earth System Science: From*
809 *Biogeochemical Cycles to Global Changes*. Academic Press, London.

810 Jaconi, A., Poeplau, C., Ramirez-Lopez, L., Van Wesemael, B., Don, A., 2019. Log-ratio
811 transformation is the key to determining soil organic carbon fractions with near-
812 infrared spectroscopy. *European Journal of Soil Science*, 70, 127–139.

813 Janik, L.J., Skjemstad, J.O., Merry, R.H., 1998. Can mid infrared diffuse reflectance
814 analysis replace soil extractions? *Australian Journal of Experimental Agriculture* 38,
815 681–696.

816 Jauss, V., Sullivan, P.J., Sanderman, J., Smith, D.B., and Lehmann, J., 2017. Pyrogenic
817 carbon distribution in mineral topsoils of the northeastern United States. *Geoderma*
818 296, 69–78.

819 Knadel, M., Deng, F., Thomsen, A., Greve, M.H., 2012. Development of a Danish national
820 vis-NIR soil spectral library for soil organic carbon determination. In: Minasny, B.,

821 Malone, B.P., McBratney, A.B. (Eds.), Digital Soil Assessments and Beyond. CRC
822 Press, Boca Raton, FL.

823 Lagacherie, P., Baret, F., Feret, J.-B., Madeira Netto, J., Robbez-Masson, J.-M., 2008.
824 Estimation of soil clay and calcium carbonate using laboratory, field and airborne
825 hyperspectral measurements. *Remote Sensing of Environment* 112(3), 825–835.

826 McCrea, J.M. , 1950. On the isotopic chemistry of carbonates and a paleotemperature
827 scale. *Journal of Chemical Physics* 18, 849–857.

828 Mark, H.L., Tunnell, D., 1985. Qualitative near-infrared reflectance analysis using
829 Mahalanobis distances. *Analytical Chemistry* 57, 1449–1456.

830 McCarty, G.W., Reeves III, J.B., Reeves, V.B., Follett, R.F., Kimble, J.M., 2002. Mid-
831 infrared and near-infrared diffuse reflectance spectroscopy for soil carbon
832 measurement. *Soil Science Society of America Journal* 66: 640–646.

833 Mevik, B.-H., Wehrens, R., 2007. The pls Package: Principal Component and Partial Least
834 Squares Regression in R. *Journal of Statistical Software* 18, 1–24.

835 Nocita, M., Stevens, A., Toth, G., Panagos, P., van Wesemael, B., Montanarella, L., 2014.
836 Prediction of soil organic carbon content by diffuse reflectance spectroscopy using a
837 local partial least square regression approach. *Soil Biology & Biochemistry* 68, 337–
838 347

839 Pearson, R.K., 2002. Outliers in process modeling and identification. *IEEE Transactions on*
840 *Control Systems Technology* 10(1), 55–63.

841 Pribyl, D.W., 2010. A critical review of the conventional SOC to SOM conversion factor.
842 *Geoderma* 156, 75–83.

843 R Development Core Team, 2012. R: A Language and Environment for Statistical
844 Computing. R foundation for Statistical Computing, Vienna. <http://www.R-project.org/>

845 Rabenarivo, M., Chapuis-Lardy, L., Brunet, D., Chotte, J.-L., Rabeharisoa, L.,
846 Barthès, B.G., 2013. Comparing near and mid-infrared reflectance spectroscopy for

847 determining properties of Malagasy soils, using global or LOCAL calibration. Journal
848 of Near Infrared Spectroscopy 21(6): 495–509.

849 Ramirez-Lopez, L., Behrens, T., Schmidt, K., Stevens, A., Demattê, J.A.M., Scholten, T.
850 2013. The spectrum-based learner: A new local approach for modeling soil vis–NIR
851 spectra of complex datasets. *Geoderma* 195–196, 268–279.

852 Reeves, D.W., 1997. The role of soil organic matter in maintaining soil quality in
853 continuous cropping systems. *Soil & Tillage Research* 43, 131–167.

854 Reeves III, J.B., 2010. Near- versus mid-infrared diffuse reflectance spectroscopy for soil
855 analysis emphasizing carbon and laboratory versus on-site analysis: where are we
856 and what needs to be done? *Geoderma* 158 (1-2), 3–14.

857 Romanyà, J., Rovira, P., 2011. An appraisal of soil organic C content in Mediterranean
858 agricultural soils. *Soil Use and Management* 27, 321–332.

859 Scharlemann, J.P., Tanner, E.V., Hiederer, R., Kapos, V., 2014. Global soil carbon:
860 understanding and managing the largest terrestrial carbon pool. *Carbon*
861 *Management* 5, 81–91.

862 Shenk, J., Westerhaus, M., Berzaghi, P., 1997. Investigation of a LOCAL calibration
863 procedure for near infrared instruments. *J. Near Infrared Spectrosc.* 5 (1), 223-232.

864 Shetty, N., Rinnan, A., Gislum, R., 2012. Selection of representative calibration sample
865 sets for near-infrared reflectance spectroscopy to predict nitrogen concentration in
866 grasses. *Chemometrics and Intelligent Laboratory Systems* 111, 59–65.

867 Shi, Z., Ji, W., Viscarra-Rossel, R.A., Chen, S., Zhou, Y., 2015. Prediction of soil organic
868 matter using a spatially constrained local partial least squares regression and the
869 Chinese vis–NIR spectral library. *European Journal of Soil Science* 66, 679–687.

870 Soriano-Disla, J.M., Janik, L.J., Viscarra Rossel, R.A., MacDonald, L.M., McLaughlin, M.J.,
871 2014. The performance of visible, near-, and mid-infrared reflectance spectroscopy

872 for prediction of soil physical, chemical, and biological properties. *Applied*
873 *Spectroscopy Reviews* 49(2), 139–186.

874 Stevens, A, Nocita, M, Tóth, G, Montanarella, L, van Wesemael, B., 2013. Prediction of
875 soil organic carbon at the European scale by visible and near infrared reflectance
876 spectroscopy. *PLoS ONE* 8(6): e66409.

877 Tatzber, M., Mutsch, F., Mentler, A., Leitger, E., Englich, M., Gerzabek, M., 2010.
878 Determination of organic and inorganic carbon in forest soil samples by mid-infrared
879 spectroscopy and partial least squares regression. *Applied Spectroscopy*
880 64(10):1167–75.

881 Tenenhaus, M., 1998. *La Régression PLS*. Editions Technip, Paris.

882 Terra, F.S., Demattê, J.A.M., Viscarra Rossel R.A. (2015). Spectral libraries for quantitative
883 analyses of tropical Brazilian soils: Comparing vis–NIR and mid-IR reflectance data.
884 *Geoderma* 255–256 (2015) 81–93

885 Viscarra Rossel, R.A., Jeon, Y.S., Odeh, I.O.A., McBratney, A.B., 2008. Using a legacy soil
886 sample to develop a mid-IR spectral library. *Soil Research* 46, 1–16.

887 Viscarra Rossel, R.A., Walvoort, D.J.J., Mc Bratney, A.B., Janik, L.J., Skjemstad, J.O.,
888 2006. Visible, near-infrared, mid-infrared or combined diffuse reflectance
889 spectroscopy for simultaneous assessment of various soil properties. *Geoderma* 131,
890 59–75.

891 Viscarra Rossel, R.A., Webster, R., 2012. Predicting soil properties from the Australian soil
892 visible-near infrared spectroscopic database. *European Journal of Soil Science* 63,
893 848–860.

894 Viscarra Rossel, R.A., Behrens, T., Ben-Dor, E., Brown, D.J., Demattê, J.A.M., Shepherd,
895 K.D., Shi, Z., Stenberg, B., Stevens, A., Adamchuk, V., Aïchi, H., Barthès, B.G.,
896 Bartholomeus, H.M., Bayer, A.D., Bernoux, M., Böttcher, K., Brodsky, L., Du, C.W.,
897 Chappell, A., Fouad, Y., Genot, V., Gomez, C., Grunwald, S., Gubler, A., Guerrero,

898 C., Hedley, C.B., Knadel, M., Morras, H.J.M., Nocita, M., Ramirez-Lopez, L., Roudier,
899 P., Campos, E.M. Rufasto, Sanborn, P., Sellitto, V.M., Sudduth, K.A., Rawlins, B.G.,
900 Walter, C., Winowiecki, L.A., Hong, S.Y., Ji, W., 2016. A global spectral library to
901 characterize the world's soil. *Earth-Science Reviews* 155, 198–230.

902 Viscarra Rossel, R.A., Lark, R.M., 2009, Improved analysis and modelling of soil diffuse
903 reflectance spectra using wavelets. *European Journal of Soil Science* 60, 453–464.

904 Vohland, M., Ludwig, M., Harbich, M., Emmerling, C., Thiele-Bruhn, S., 2016. Using
905 variable selection and wavelets to exploit the full potential of visible–near infrared
906 spectra for predicting soil properties. *Journal of Near Infrared Spectroscopy* 24, 255–
907 269.

908 Walkley, A., Black, I.A., 1934. An examination of the Degtjareff method for determining soil
909 organic matter, and a proposed modification of the chromic acid titration method. *Soil*
910 *Science* 37(1), 29–38.

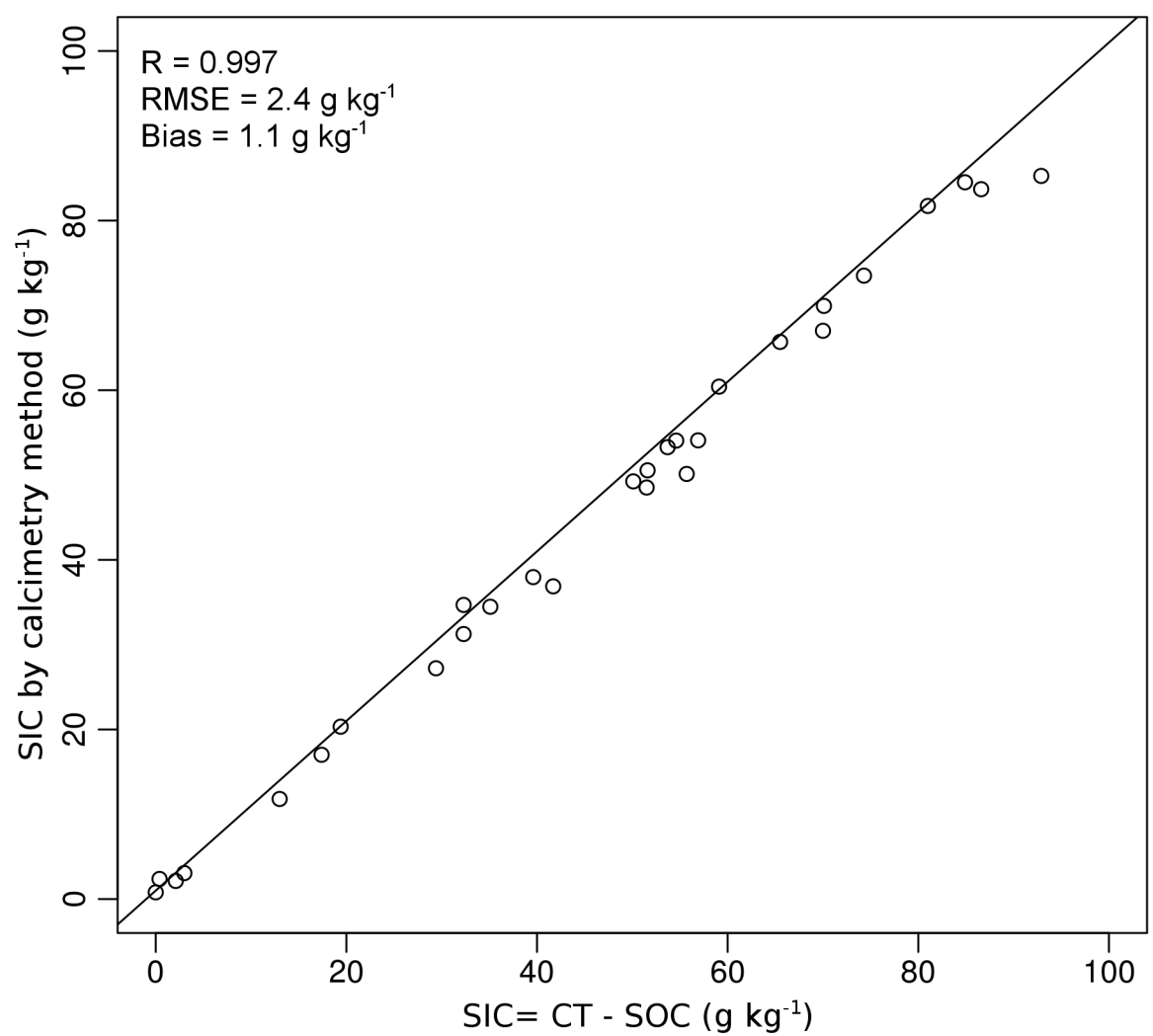
911 Wang, X., Wang, J., Zhang, J. (2012) Comparisons of three methods for organic and
912 inorganic carbon in calcareous soils of northwestern China. *PLoS ONE* 7(8): e44334.

913 Wijewardane, N.K., Ge, Y., Wills, S., Libohova, Z., 2018. Predicting physical and chemical
914 properties of US soils with a mid-infrared reflectance spectral library. *Soil Science*
915 *Society of America Journal* 82, 722–731

916 Williams, P.C., Norris, K.H., 1987. Qualitative applications of near-infrared reflectance
917 spectroscopy. In: Williams, P., Norris, K. (Eds.), *Near-Infrared Technology in the*
918 *Agricultural and Food Industries*. American Association of Cereal Chemists, St. Paul,
919 MN, pp. 241–246.

920 Wold, S., Johansson, E., Cocchi, M., 1993. PLS - partial least squares projections to latent
921 structures. In: Kubinyi, H. (Ed.), *3D-QSAR in Drug Design, Theory, Methods, and*
922 *Applications*. ESCOM Science Publishers, Leiden, pp. 523–550.

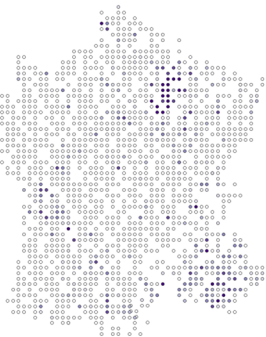
- 923 Wold, S., Sjöström, M., Eriksson, L., 2001. PLS-regression: a basic tool of chemometrics.
924 Chemometrics and Intelligent Laboratory Systems 58, 109–130.
- 925 Workman, J. Jr., Weyer, L., 2008. Practical Guide to Interpretive Near-Infrared
926 Spectroscopy. CRC Press, Boca Raton, FL.
- 927 Waruru, B.K., Shepherd K.D., Ndegwa G.M., Kamoni P.T., Sila A.M. 2014. Rapid
928 estimation of soil engineering properties using diffuse reflectance near infrared
929 spectroscopy. Biosystems Engineering, 121, pp. 177-185.
- 930 Yang, M., Xu, D., Chen, S., Li, H., Shi, Z., 2019. Evaluation of machine learning
931 approaches to predict soil organic matter and pH using vis-NIR spectra. Sensors 19,
932 263.



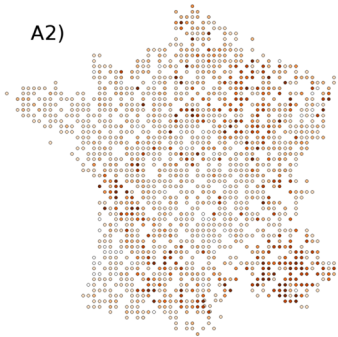
A1)

Observed
SIC (g kg⁻¹):

- 0 - 5.4
- 5.4 - 17.8
- 17.8 - 32.6
- 32.6 - 50
- 50 - 71.5
- 71.5 - 103.9

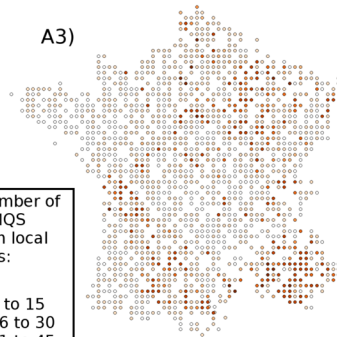


A2)

Selection number of
each RMQS
spectrum in local
models:

- 0
- From 1 to 15
- From 16 to 30
- From 31 to 45
- From 46 to 60
- From 61 to 75
- From 76 to 90

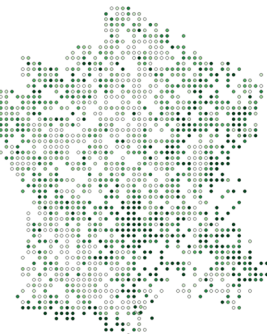
A3)



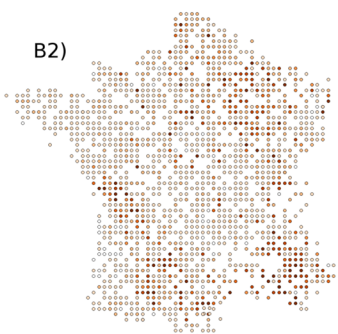
B1)

Observed
SOC (g kg⁻¹):

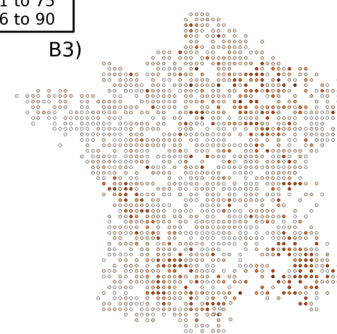
- 0 - 11.2
- 11.2 - 15.1
- 15.1 - 19.6
- 19.6 - 1
- 26.1 - 38.6
- 38.6 - 411.3

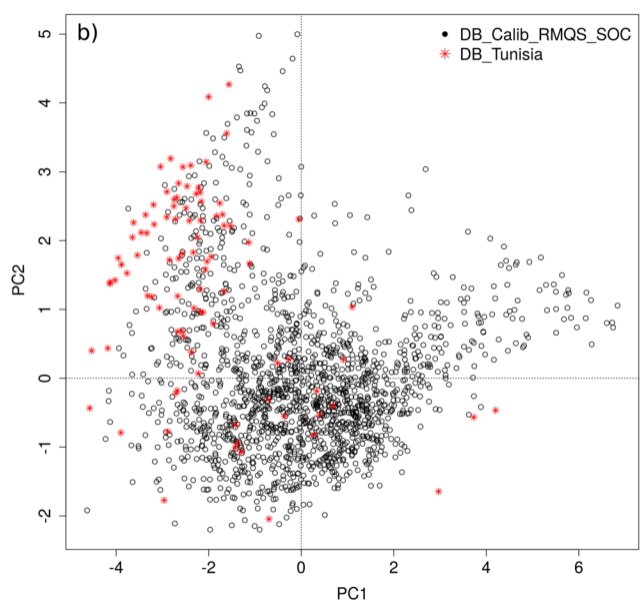
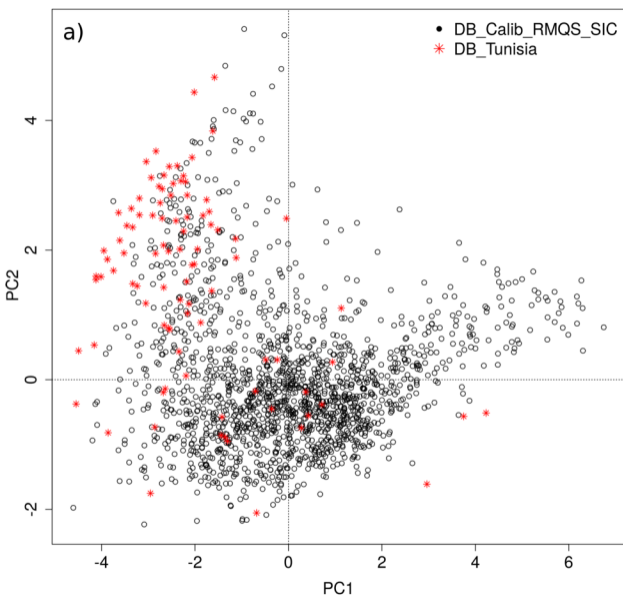


B2)

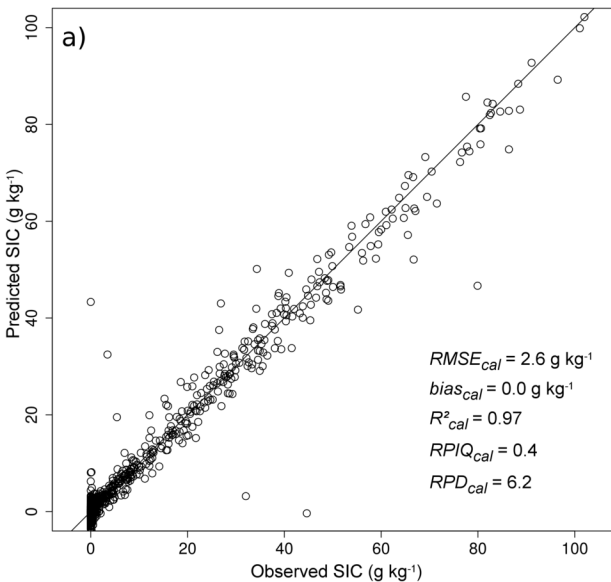


B3)

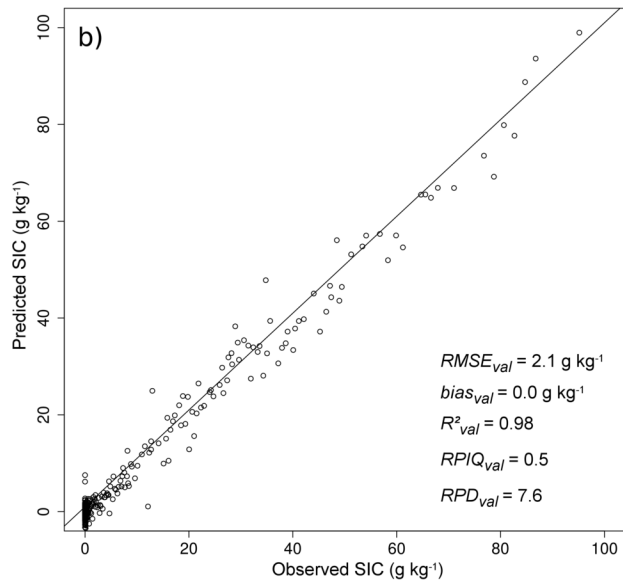




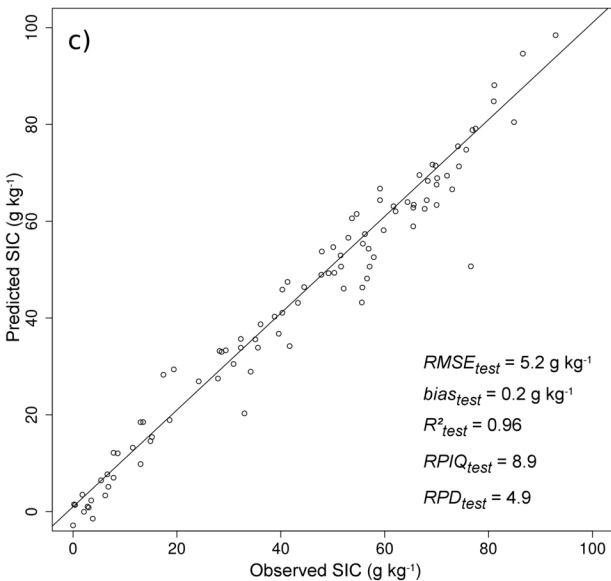
Calibration data



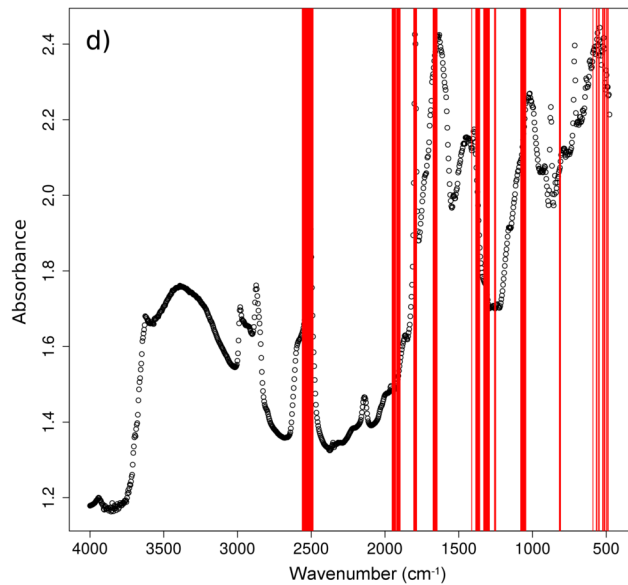
Validation data



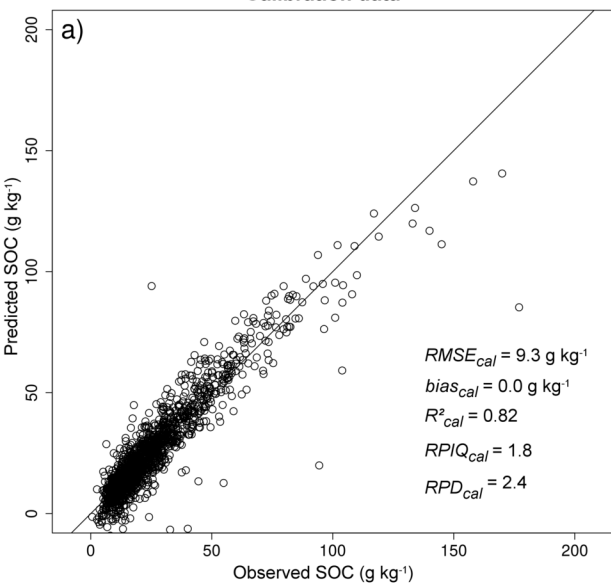
Test data



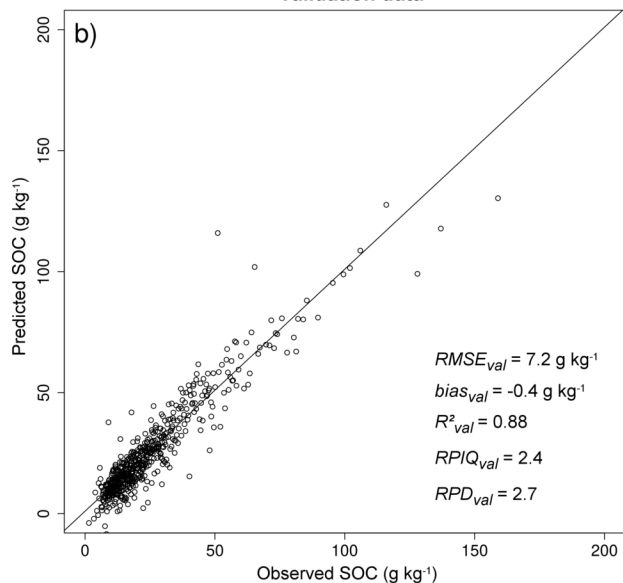
Number of significant wavenumbers = 96



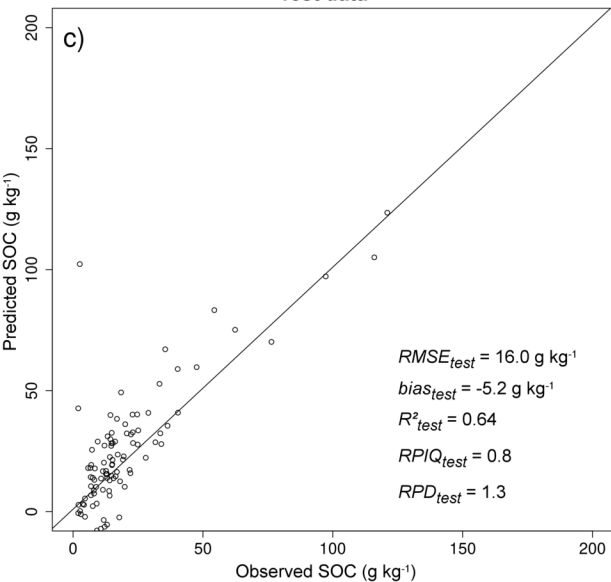
Calibration data



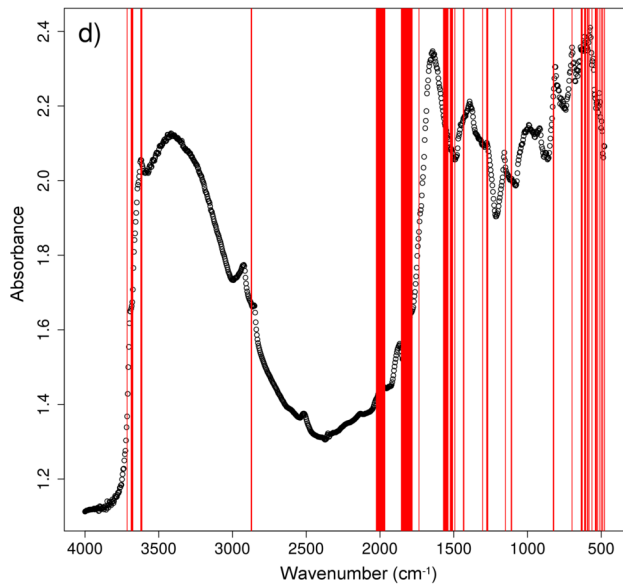
Validation data



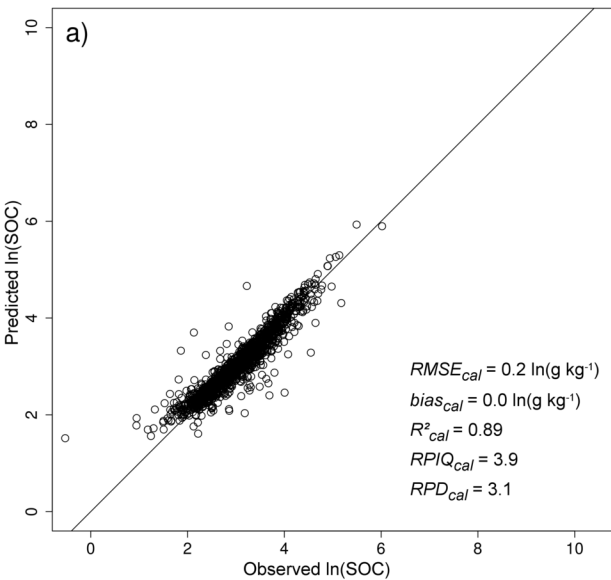
Test data



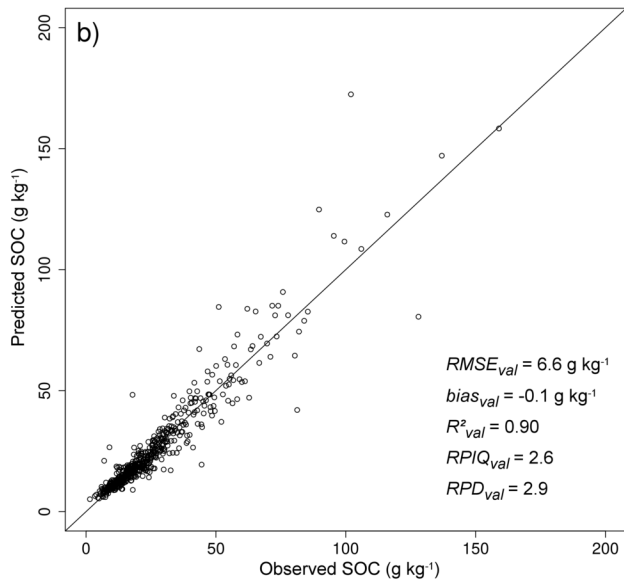
Number of significant wavenumbers = 92



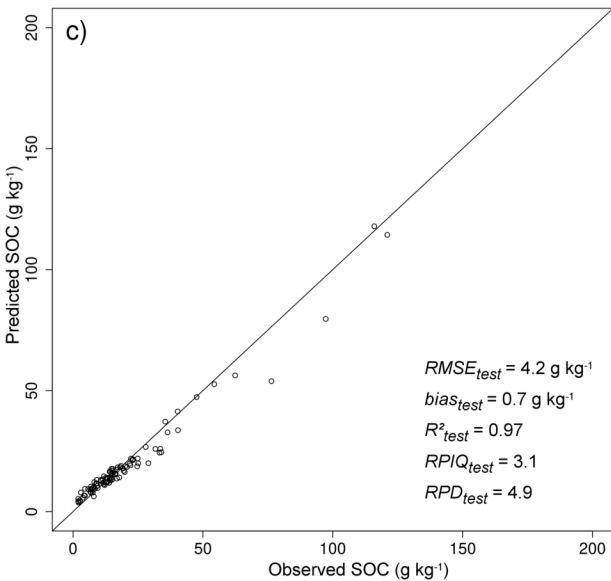
Calibration data



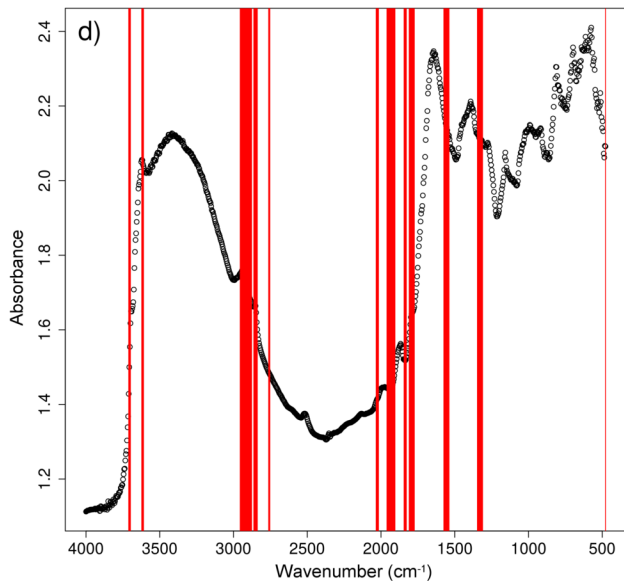
Validation data

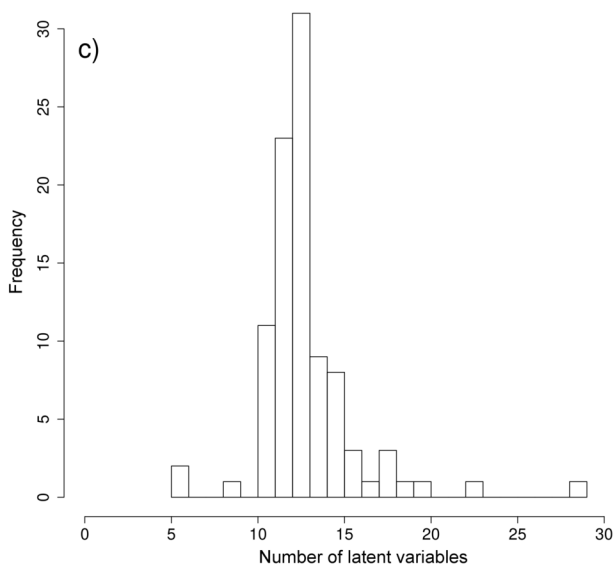
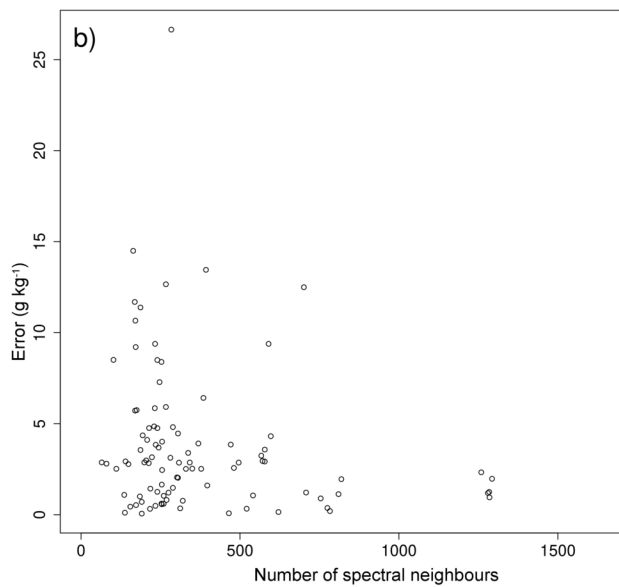
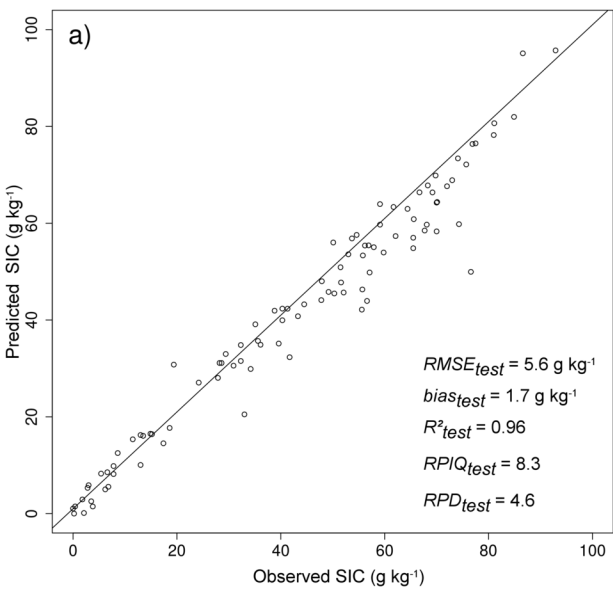


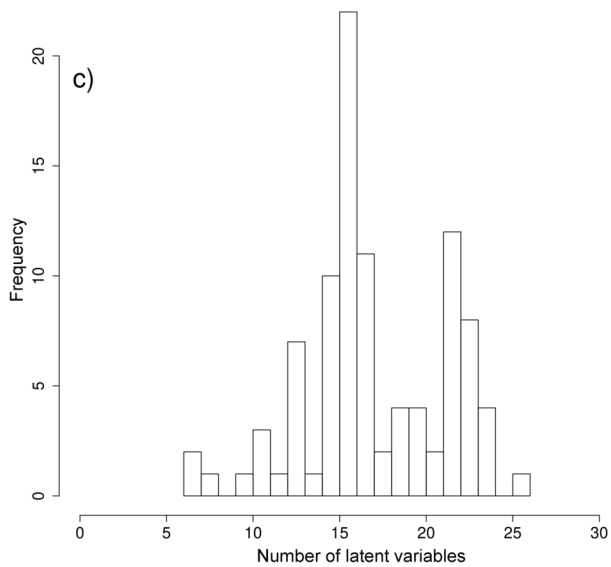
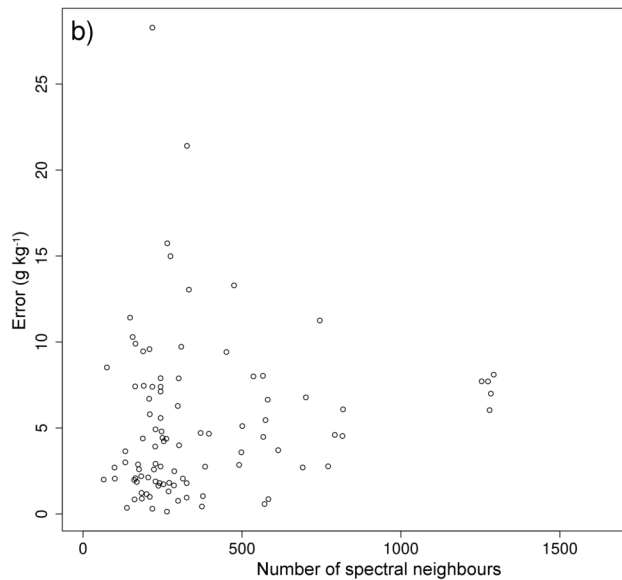
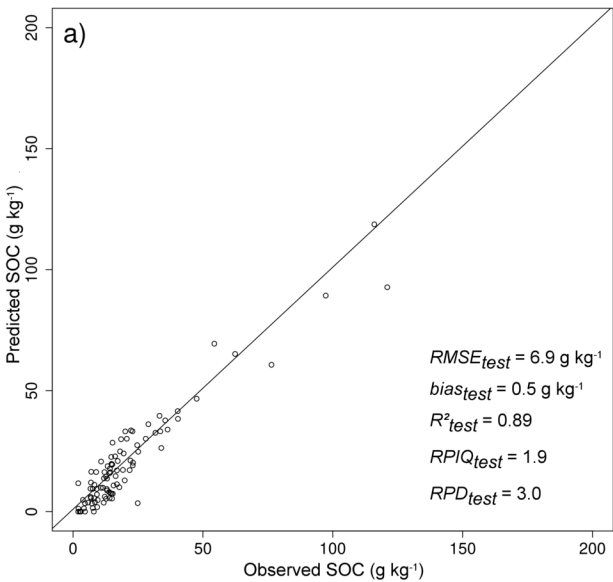
Test data



Number of significant wavenumbers = 95







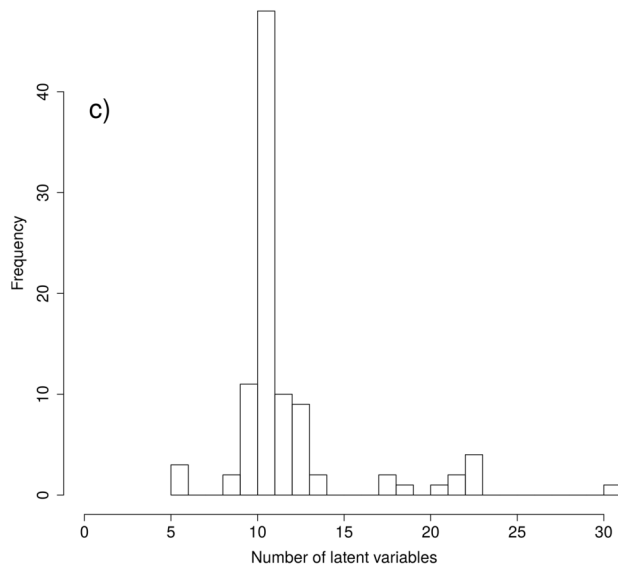
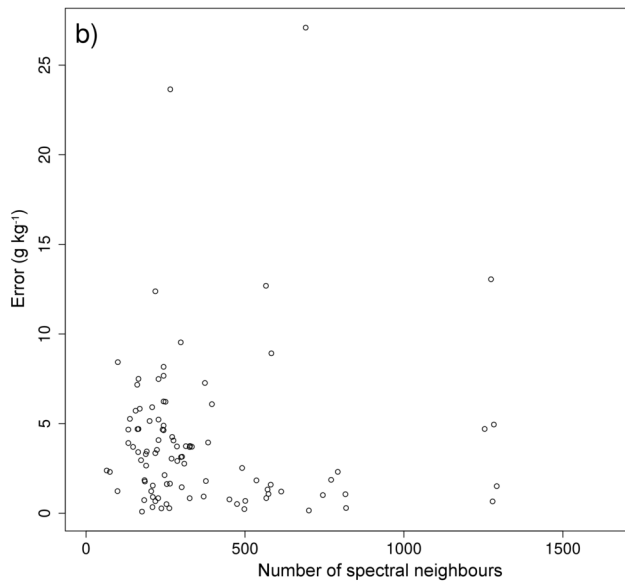
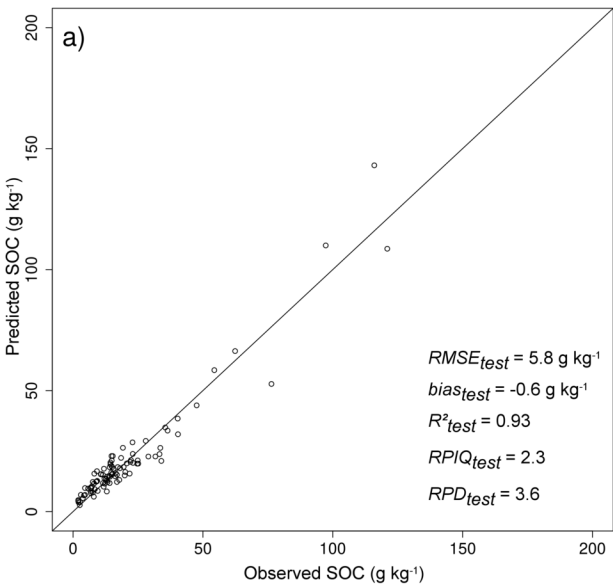


Table 1: Soil datasets statistics. The SIC values set to zero correspond to values under the laboratory quantification limit ($< 0.1 \text{ g kg}^{-1}$).

Dataset	Number of soil samples	Min g kg ⁻¹	Max g kg ⁻¹	Mean g kg ⁻¹	Median g kg ⁻¹	SD^a g kg ⁻¹	Skewness
<i>DB_RMQS_SOC</i>	2178	0.6	411.3	25.8	19.6	21.8	4.9
<i>DB_RMQS_SIC</i>	2178	0.0	103.9	6.4	0.0	16.1	3.1
<i>DB_Tunisia_SOC</i>	96	2.0	121.0	20.1	14.6	21.0	3.0
<i>DB_Tunisia_SIC</i>	96	0.0	92.9	43.3	48.5	25.6	-0.2
<i>DB_Calib_RMQS_SOC</i>	1586 ^b	0.6	411.3	25.6	19.4	22.3	5.5
<i>DB_Calib_RMQS_SIC</i>	1582 ^c	0.0	103.9	6.4	0	16.1	3.1
<i>DB_Valid_RMQS_SOC</i>	544	1.5	159.0	25.5	19.6	19.5	2.5
<i>DB_Valid_RMQS_SIC</i>	544	0.0	95.2	6.3	0.0	15.7	3.1

^a SD: standard deviation

^b after removing 48 spectral outliers

^c after removing 52 spectral outliers

Table 2: Figures of merit obtained with global and local models over the French calibration and validation databases.

Models	R^2_{cv}	$RMSE_{CV}$ (g kg ⁻¹)	R^2_{val}	$RMSE_{val}$ (g kg ⁻¹)	$bias_{val}$ (g kg ⁻¹)	RPD_{val}	$RPIQ_{val}$
GM_{sic}	0.97	2.8	0.98	2.1	0.0	7.6	0.5
LM_{sic}	<i>nd</i>	<i>nd</i>	0.99	1.8	0.0	8.8	0.6
GM_{soc}	0.80	9.9	0.88	7.2	-0.4	2.7	2.4
LM_{soc}	<i>nd</i>	<i>nd</i>	0.93	5.4	-0.7	3.6	3.2
GM_{insoc}	0.89	0.2*	0.90	6.6	-0.1	2.9	2.6
LM_{insoc}	<i>nd</i>	<i>nd</i>	0.92	5.7	-0.1	3.4	3

* $RMSE_{cv}$ calculated on $\ln(SOC)$

nd: Not determined.

Table 3: Figures of merit obtained with global and local models over the Tunisian soil samples.

Prediction model	R^2_{test}	$RMSE_{test}$ (g kg ⁻¹)	$bias_{test}$ (g kg ⁻¹)	RPD_{test}	$RPIQ_{test}$
<i>GM_{sic}</i>	0.96	5.2	0.2	4.9	8.9
<i>LM_{sic}</i>	0.96	5.6	1.7	4.6	8.3
<i>GM_{soc}</i>	0.64	16.0	-5.2	1.3	0.8
<i>LM_{soc}</i>	0.89	6.9	0.5	3.0	1.9
<i>GM_{insoc}</i>	0.97	4.2	0.7	4.9	3.1
<i>LM_{insoc}</i>	0.93	5.8	-0.6	3.6	2.3