# Harvesting spatially dense legacy soil datasets for digital soil mapping of available water capacity in Southern France

Quentin Styc, François Gontard, Philippe Lagacherie

1 **Harvesting spatially dense legacy soil datasets for digital soil**
2 **mapping of available water capacity in Southern France**

3 Styc Quentin[1,2]

4 François Gontard[2]

5 Lagacherie Philippe[1]

6 [1] LISAH, Univ Montpellier, INRAE, IRD, Institut Agro, Montpellier, France

7 [2] BRL Exploitation, Nîmes, France

8 Corresponding author: Lagacherie Philippe

9 Abstract

10 Although considerable work has been conducted in recent decades to build soil databases, the

11 legacy data from a lot of former soil survey campaigns still remain unused. The objective of

12 this study was to determine the interest in harvesting such legacy data for mapping the soil

13 available water capacities (SAWCs) at different rooting depths (30 cm, 60 cm, 100 cm) and to

14 the maximal observation depth, over the commune of Bouillargues (16 km$^2$, Occitanie region,

15 southern France)

16 An increasing number of available auger hole observations with SAWC estimations – from 0

17 to 2781 observations – were added to the existing soil profiles to calibrate quantile regression

18 forests (QRFs) using the Euclidean buffer distances from the sites as soil covariates. The

19 SAWC was first mapped separately for different soil layers, and the mapping outputs were

20 pooled to estimate the required SAWC. The uncertainty of the SAWC prediction was

21 estimated from the estimated mapping uncertainties of the individual soil layers by an error

22 propagation model using a first-order Taylor analysis.

23  The performances of the SAWC predictions and their uncertainties were evaluated with a 10-

24  fold cross validation that was iterated 20 times. The results showed that the use of a quantile

25  regression forest that was fed with auger hole observations and that used the Euclidean buffer

26  distances as soil covariates considerably augmented the performances of the SAWC

27  predictions (percentages of explained variance from 0.39 to 0.70) compared to the

28  performance of a classical DSM approach, i.e., a QRF that solely used soil profiles and only

29  environmental covariates (percentages of explained variance from 0.04 to 0.51). The analysis

30  of the results revealed that the performances were also dependent on the spatial patterns of the

31  different examined SAWCs and was limited by the observational uncertainties of the SAWCs

32  determined from auger holes. The best performance tended to also provide the best view of

33  the uncertainty patterns with an overestimation of uncertainty.

34  Despite these gains in performance, the cost-efficiency analysis showed that the augmentation

35  of soil observations was not cost efficient because of the highly time-consuming manual data

36  harvesting protocol. However, this result did not account for the observed gain in map details.

37  Furthermore, the cost efficiency could be further improved by automation.

38

39

40  1.  Introduction

41  Digital soil mapping (DSM) has been recognized as the appropriate solution to provide spatial

42  soil information for land users, scientist communities and policy and decision makers in

43  agriculture and the environment (McBratney et al., 2003; Sanchez et al., 2009). The principle

44  of DSM is to predict a soil property or soil classes and the associated prediction uncertainty

45  by determining the quantitative relationships between the soil information available over a

46  limited set of  locations and the spatial data reflecting the state factors of soil formation

47  (envionmental covariates). DSM has now moved from a largely academic movement toward

48  an operational activity (Minasny & McBratney, 2016, Arrouays et al, 2017).

49  However, the performances of DSM predictions of soil properties often exhibit more

50  uncertainty than initially expected. For example, the percentages of explained variances of

51  less than 0.5 were observed for 95%, 76%, 100% and 86% of the tested soil properties for

52  DSM applications at the catchment scale (Nussbaum et al., 2018), at the regional scale

53  (Vaysse and Lagacherie, 2015), at the national scale (Mulder et al., 2016), and at the global

54  scale (Hengl et al., 2014), respectively.

55  These authors converged toward the conclusion that the density of soil observations used for

56  calibrating the DSM models was the main factor that limited the DSM performances. Most of

57  the soil information used as input in DSM applications has been either soil maps or the spatial

58  sampling of sites with soil property measurements. The average densities used in most

59  operational DSM applications have been low, e.g., 4-12 sites/km² (several study areas in

60  Nussbaum et al., 2018), 0.07 sites/km² (Vaysse and Lagacherie, 2015), 0.03 sites/km² (Mulder

61  et al., 2016), and 0.001 sites/km² (Hengl et al., 2014), which limits the performances of soil

62  prediction, especially when the pattern of variation in the soil property is largely below the

63  spacing of soil profiles (Vaysse and Lagacherie, 2015; Gomez and Coulouma, 2018). In

64  addition, further experiments that consisted of varying the spatial density of soil input

65  confirmed this analysis (Somarathna et al. 2017, Wadoux et al. 2019, Lagacherie et al, 2020).

66  Consequently, it is of paramount importance to increase the density of soil inputs to improve

67  the performance of DSM models in predicting soil properties (Voltz et al., 2020).

68  The most straightforward way to increase the density of DSM model soil inputs involves

69  harvesting the legacy soil data that have not yet been stored in the existing soil databases.

70  Arrouays et al. (2017) showed that during the period 2009-2015, the numbers of legacy soil

71  profiles stored in global and national soil databases increased by 1,046% and 45%,

72    respectively. However, they estimated that a large amount of soil legacy data can still be

73    harvested. This is even more true in some areas across the world where soil surveying has

74    been particularly active in the past.

75    For example, in southern France, the BRL irrigation company conducted detailed soil surveys

76    over its irrigation perimeter between 1957 and 1992, which resulted in detailed soil maps,

77    25,000 soil profiles (5/km²) and 203,000 auger hole observations (31/km²). At this stage, such

78    soil data have not yet been harvested and therefore cannot be used as input for DSM

79    applications. However, this data has great potential for improving DSM performance and

80    should be thoroughly examined.

81    In this paper, a spatially dense set of soil observations harvested from soil survey documents

82    was tested for improving the performances of DSM models in mapping soil available water

83    capacities for different rooting depths (0-30 cm, 0-60 cm, 0-100 cm) and at maximum

84    observation depth, and the associated uncertainties. Our aim was to evaluate the cost-

85    efficiency ratio of using such soil observations and to evaluate the added value of using

86    euclidian buffer distances as additional inputs of DSM models as proposed by Hengl et al

87    (2018). The study is conducted in the commune of Bouillargues, which is one of the

88    communes included in the BRL irrigation perimeter.

89

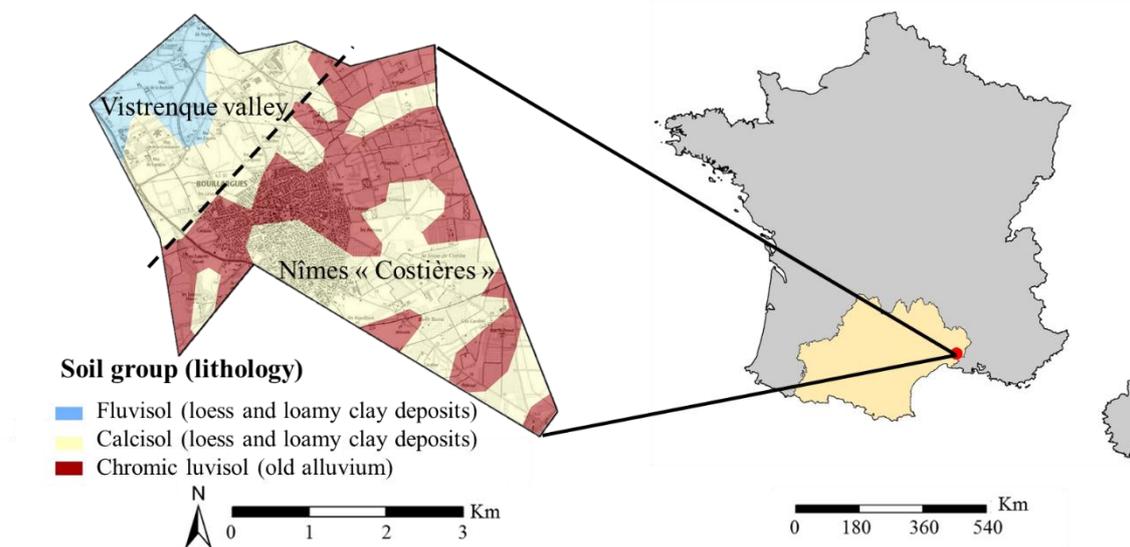90    2.  The case study

91       2.1. The study area

92    This study took place in the administrative commune of Bouillargues in the Occitanie

93    administrative French region (Figure 1). Located in southern France, Bouillargues covers 16

94    km² and is mainly devoted to vineyards, agricultural lands, forests, and scrublands.

95 Bouillargues has a Mediterranean climate characterized by a moderate average annual rainfall

96 (600 mm) and dry and hot summers.

97 The study area is topographically split into two subregions with the large flat valley of the

98 Vistrenque in the northern part and old fluviatile alluvium terraces belonging to the Nîmes

99 "Costière" in the southern part. The two subregions have contrasting parent materials with i)

100 loess and loamy clay deposition in the Vistrenque valley and ii) old alluvium in the Nimes

101 Costière part, covered by some loess deposits. The contrast in parent materials induces

102 variations in soils with i) fluvisols and calcisols developed in loess and loamy clay deposition,

103 characterized by an absence of coarse fragments and a loamy texture, and ii) chromic luvisols

104 developed in old alluvium terraces characterized by important coarse fragment contents and

105 compacted clay accumulations (Figure 1).

106



107

108 *Figure 1. Location of the study area and distribution of soil samples*

109 2.2. Soil data

110 2.2.1. History and content of the BRL soil database

111 The soil data of this study are a part of the soil survey led by the "Compagnie Nationale

112 d'Aménagement de la Région du Bas-Rhône et du Languedoc" (CNARBRL) between 1957

113 and 1992 over the irrigated perimeter of this irrigation company, which covers 6,636 km². The

114 objectives of this survey were to provide suitable soil information for i) improving the

115 development master plan of the irrigation perimeter and estimating the surface area of arable

116 and potentially irrigable lands and ii) supporting the cultural intensification made possible by

117 irrigation, assessing the irrigation supply, and setting technical assistance for landholders to

118 start irrigation and crop conversion.

119 The compilation of those studies resulted in a database of 228,000 soil observations with

120 25,000 soil profile descriptions and laboratory analyses (Figure 2) and 203,000 auger holes

121 (Figure 3), which correspond to average spacings of 515 m and 181 m for the soil profiles and

122 auger holes, respectively.

123



124 *Figure 2. Soil profile a) horizon descriptions with geographical coordinates (black box) and b) laboratory analysis results,*
125 *physical analysis (red box) and chemical analysis (blue box)*

*Figure 3. Auger hole descriptions*

### 2.2.2. Spatial sampling and georeferencing in the study area

Focusing on the commune of Bouillargues, the harvested dataset is composed of 2850 sites with soil observations that include 2781 auger holes and 69 soil profiles, which correspond to average spacings of 76 m and 500 m, respectively (Figure 4). Both the soil profiles and the auger hole observations were fairly evenly distributed over the study area; however, some gaps corresponded to urbanized areas or lands that were not expected to have any agricultural potential.

135

138    The soil profile data records included geographical coordinates (Lambert III, black box in

139    Figure 2a), whereas manual preprocessing was necessary for georeferencing the auger holes.

140    The auger holes were initially located through a non-georeferenced map representing the local

141    sampling scheme (Figure 5a). Each sampling scheme corresponded to an area of water

142    distribution supplied by an irrigation water access point of the BRL irrigation network. This

143    access point was georeferenced and could be positioned onto a georeferenced former cadastre

144    (red box in Figure 5b). To acquire the coordinates of the auger holes, the sampling scheme

145    was first located in the georeferenced cadastre using the coordinates of the irrigation water

146    access point. Its boundaries were then positioned (green dashed perimeter on Figure 5b) using

147    the geometry of the parcels and communication paths. Finally, each auger hole was manually

148    positioned onto the georeferenced cadastre (blue stars on Figure 5b) using the sampling

149    scheme (Figure 5a), and the coordinates of the auger holes were obtained using the

150    coordinates acquisition tool of BRL's web-GIS (Figure 5c).

a) Ungeoreferenced sampling scheme

b) Fitting sampling scheme boundaries on georeferenced former cadastre

c) Auger holes adjustment and coordinates acquistion

★ Auger holes

151

152    *Figure 5. Fitting the non-georeferenced sampling scheme of auger holes in the georeferenced former cadastre*

153

154    2.2.3.  Soil available water capacity determinations at sites with soil observations

155    This study took the mapping of soil available water capacity (SAWC) as an example of

156    applying DSM. SAWC refers to the capacity of the soil to store water for plant growth

157    (Veihmayer and Hendrickson, 1927). This functional property plays a key role in many

158    ecosystem services, such as food production, soil drought or climate and gas regulation.

159    Consequently, it is a crucial parameter used in land evaluations and recently in ecosystem

160    services assessments (Dominati et al., 2014). Information about the SAWC distribution in

161    space is essential for planning and management in agriculture and for ecological modeling. In

162    the present example, SAWC was required for fulfilling the irrigation objectives evoked above

163    (section 2.2.1). Currently, SAWC is computed in the literature as follows (Cousin et al.,

164    2003):

$$SAWC = \sum_{i=1}^{n} dh_i * bd_i * \left(\frac{100 - st_i}{100}\right) * (\theta r_i - \theta w_i)$$    (1)

165

166    where SAWC is the soil available water capacity (cm), $dh_i$ = the thickness of the $i$th horizon

167    (cm), $bd_i$ = the bulk density (g/cm$^3$) of the $i$th horizon, $st_i$ = the coarse fragment content of

168    the $i$th horizon (% volumetric), and $\theta r_i$ and $\theta w_i$ are the gravimetric soil water contents at

169    field capacity (i.e., the soil water content that remains in the soil after water has drained due to

170    gravitational force) and the permanent wilting point (i.e., the soil water retained so strongly

171    that it is no longer available for plant roots, so plants wither and cannot recover their

172    turgidity) of the $i$th horizon (cm$^3$.cm$^{-3}$), respectively.

173    Historically, the CNARBRL had a different approach for expressing the water retention term

174    of the fine earth, i.e., $(\theta r_i - \theta w_i)$, which leads to the following equation:

$$SAWC = \sum_{i=1}^{n} dh_i * bd_i * \left(\frac{100 - st_i}{100}\right) * (b_i * EqW_i)$$

(2)

175

176    The equivalent water content (EqWi) corresponds to $\theta r_i$ of Eq. 1, and the textural coefficient

177    $b_i$ is an expression of the water content at the permanent wilting point that weights $EqW_i$ to

178    account for the water content that is not available for the plant (i.e., beyond the wilting point,

179    defined as $\theta w_i$ in Eq. 1).

180    The values of $bd_i$ and $EqW_i$ were measured at each soil profile; $bd_i$ was determined in the

181    field following the Vergières protocol (Bourrier, 1965) but was estimated as 1.6 times the

182    mass fraction of the fine earth from the ensemble coarse fragment and fine earth, when the

183    coarse fragment phase of the soil sample was too important to perform the Vergières protocol

184    (Legros, 1996).

185    The $EqW_i$ of sieved samples was determined in the laboratory using a centrifuge apparatus set

186    at 100 kPa (pF = 3.0), a reference pressure that was considered, at the time of the CNARBRL

187    soil survey, as yielding the best approximation of the water content at the field capacity (see

188  section 2.2.1) (Baize and Jabiol, 1995). The $EqW_i$ values were estimated on auger hole

189  observations by local pedotransfer functions using the field estimated textural classes.

190  The $b_i$ coefficient was determined both on soil profiles and on auger hole observations by a

191  local pedotransfer function using the textural classes determined from granulometric analyses

192  and field estimation, respectively, for soil profile and auger hole observations.

193  The coarse fragment content and the horizon thicknesses of Eq. 2 were retrieved from the

194  descriptions of the physical analyses and descriptions of the soil profiles and of the auger hole

195  observations, respectively (Figures 1 and 2). Different total soil thicknesses (i.e., $\sum_{i=1}^{n} dh_i$)

196  were considered to determine the different rooting depths related to the different possible

197  crops of the study area (from market gardening to vineyard passing by annual crops). In

198  addition to the maximum soil thicknesses given by the soil observations that were considered

199  for calculating the maximum soil available water capacity (SAWCmax), restricted thicknesses

200  of 30 cm, 60 cm and 100 cm were then considered, leading to different restricted SAWCs,

201  denoted further as SAWC30, SAWC60, and SAWC100.

202  It must be noted that both the profiles and auger holes had limited observation depths of 140

203  and 120 centimeters, respectively, which may cause underestimations of SAWCmax.

204

205    2.3. Environmental covariates

206  The DSM approach, as formalized by the scorpan model (McBratney et al., 2003), considers

207  quantitative relationships between a target soil property and environmental variables, which

208  are also known as "covariates".

209  The selection of environmental covariates depends on two criteria: i) they could be derived

210  from geodatasets freely available at least at the French national level, and ii) they have a

211  logical and process-based relationship with soil properties according to the literature.

212    Following these criteria, we derived covariates related to the scorpan model component, i.e.,

213    topography, organisms, and parent material, that regroups the major landscape types across

214    the study area. Climate data were not considered in this study since we did not find any

215    climate data at a spatial resolution fine enough to represent the climate variations over such a

216    small area. The relief component was described by a set of geomorphometric indicators

217    currently considered in DSM studies: elevation, slope, aspect, multiresolution valley bottom

218    flatness (MRVBF), multiresolution ridge top flatness (MRRTF), topographic wetness index

219    (TWI), topographic position index, plan curvature and profile curvature. These indicators

220    were derived from the French altimetry database (BD ALTI, 25 m resolution) digital elevation

221    model (DEM). They were  computed using the SAGA GIS software (Böhner et al., 2006) and

222    his Terrain Analysis procedures.

223

224    Organisms and parent materials were derived from the Landsat 7 imagery and geological

225    map, respectively, and were both resampled at the native resolution of the DEM (i.e., 25 m).

226    Additionally, parent material covariates were developed by Vaysse and Lagacherie (2015)

227    from the geological map (1:50,000) qualitative descriptions to quantitative indicators

228    describing the hardness, mineralogy and texture of alteration materials.

229    *Table 1. Exhaustive categorical and continuous covariates*

| Variables | Abbreviation | Resolution/Scale | Source | Soil-forming factor[1] | Type[2] |
|---|---|---|---|---|---|
| *Topography* | | | | | |
| Elevation | ELEV | 25 m | BD ALTI | r | Q |
| Multiresolution Valley Bottom Flatness | MRVBF | 25 m | BD ALTI | r | Q |
| Slope | SLOPE | 25 m | BD ALTI | r | Q |
| Topographic Wetness | TWI | 25 m | BD ALTI | r | Q |

| Index | | | | | |
|---|---|---|---|---|---|
| Plan Curvature | PLANCURV | 25 m | BD ALTI | r | Q |
| Profile Curvature | PROCURV | 25 m | BD ALTI | r | Q |
| Multiresolution Ridge Top Flatness | MRRTF | 25 m | BD ALTI | r | Q |
| Topographic Position Index | TPI | 25 m | BD ALTI | r | Q |
| *Geology* | | | | | |
| Hardness | HARDNESS | 25 m | Geological map/soil profile | p | C |
| Texture | TEXTURE | 25 m | Geological map/soil profile | p | C |
| Mineralogy | MINERALOGY | 25 m | Geological map/soil profile | p | C |
| *Organisms* | | | | | |
| Land use | LANDUSE | 25 m | Landsat 7 | o | C |

[1]: SCORPAN factors (o = organisms, r = relief, p=parent material)

[2]: Q = quantitative, C = categorical

230

231    2.4. Acquisition process and cost assessment

232    In section 2.2., we presented the main difference in using soil profiles and auger holes in a

233    DSM application, i.e., the accessibility of the data. While soil profile acquisition is quite

234    straightforward, i.e., recording soil data and locations, auger hole acquisition is more

235    complicated as the locations are not directly available and manual georeferencing is required,

236    thus, the acquisition process is longer. In Table 2, we provide the main information about the

237    acquisition process for soil profiles and auger holes. As the number of auger hole observations

238    is substantially larger than the number of soil profiles and take longer to record, we provided

239    an assessment of the cost of soil data acquisition.

240 <p style="text-align:center">*Table 2. Information to assess the cost of the acquisition process*</p>

|  | Auger holes | Soil profiles |
|---|---|---|
| Recorded time of soil properties* (min/observation) | 0.8 | 0.8 |
| Recorded time of geo-localizations* (min/observation) | 2.2 | 0.2 |
| Number of observations | 2721 | 69 |

241     *Computed from timed sessions of harvesting

242 To compute the cost of the acquisition process, we applied the following formula using the

243 information in Table 2:

$$Cost = \left(\frac{N * rec\_time}{Daytime}\right) * Sal. \tag{5}$$

244 With N the number of harvested soil observations, rec_time the recorded times of harvesting a

245 given soil observation in mn (see table 2), Daytime is 1440 (number of mn in a day) and Sal is

246 the daily salary of the harvester.

247 3. Methods

248     3.1. DSM models for soil profiles

249 In this study, we used several mapping models derived from the random forest algorithm.

250 Hereafter, we provide a general description of random forest and its derivatives used in this

251 study.

252     3.1.1. Random forest

253 Random forest models (RF) (Breiman, 2001) are an ensemble learning method for both

254 classification and regression. A forest, i.e., an ensemble of randomized decision trees, is built

255 and trained based on a bootstrap approach. Individual trees are built using the principle of

256 recursive partitioning. "*The feature space is recursively split into regions containing*

257 *observations with similar response value*" (Strobl et al., 2009). The predictions of the

258 individual trees are finally averaged to obtain a single prediction.

259        3.1.2.  Quantile regression forest

260    The quantile regression forest algorithm (QRF) (Meinshausen, 2006) is an extension of

261    random forests that has become one of the most commonly used algorithms in DSM studies

262    (Hengl et al., 2015; Ugbaje and Reuter, 2013; Vaysse and Lagacherie, 2017). As a RF, QRF

263    provides an ensemble prediction based on $n$ regression trees. However, while RF provides

264    solely the conditional mean, QRF supplies the whole conditional distribution of the target

265    variable by keeping all observations at the terminal nodes. This allows us to infer estimates

266    for the conditional quantiles (Meinshausen, 2006). More details on QRF can be found in

267    Meinshausen (2006).

268    QRF was performed with the ranger package, which is a fast implementation of Breiman's

269    random forest and Meinshausen's quantile regression forest for big data (Wright and Ziegler,

270    2017). QRF was run with the default parameters given by ranger.

271    3.2. Mapping models for dense spatial sampling

272    The usual applications of RF and its derivative to DSM only exploit the relationships between

273    the soil properties to be predicted with landscape elements characterized by a set of covariates

274    derived from the available spatial data. However, they do not consider the spatial relationships

275    between sites or spatial autocorrelation, which allows the spatial interpolations of a given soil

276    property between sites. This can lead to suboptimal predictions and possibly systematic over-

277    and underestimation of predictions, especially if the target variable is spatially autocorrelated

278    and if point patterns show clear sampling bias (Hengl et al., 2018). In the case of dense

279    sampling, such spatial interpolation can be of great interest to overcome the limitations of

280    landscape covariates for predicting soil properties (Lagacherie et al, 2020).

281    To correct the non-spatial approach of RF and its derivative, Hengl et al. (2018) proposed

282    adding new covariates that consider the locations of the sites. These covariates are defined as

283    the Euclidean buffer distances from the observation sites. To limit the number of covariates

284    and the computing time in the case of a large dataset ($> 1,000$ sites), these distances to the

285    nearest points were not calculated for each individual observation site but for $n$ equal classes

286    (from low to high AWC values). As RF is sensitive to the number of classes (Hengl et al.,

287    2018), we performed a trial and error process, which was conducted to choose different

288    classes according to the maximal soil thickness considered and to the density scenario

289    (number of classes varying between 6 and 15). For each targeted SAWC, a map was

290    generated. In this DSM model, we considered soil profile and auger hole observations

291    indifferently as soil inputs, omitting their possible differences of uncertainty on the SAWC

292    determinations. This model will be denoted further $QRF_{dist}$. Euclidean buffer distance

293    mapping was performed using the *GSIF* package (Hengl, 2019).

294

295    3.3. Inference trajectories

296    Since we aimed to map SAWC, which is a soil indicator involving several soil properties and

297    several soil depths, it could be estimated following various possible inferences following the

298    order with which "combining primary soil properties", "aggregating soil layers across depths"

299    and "mapping" were performed to provide the SAWC (Styc and Lagacherie, 2019). Styc and

300    Lagacherie (2019) experienced a total of 18 inference trajectories throughout Languedoc-

301    Roussillon that were performed to obtain the most appropriate SAWC map. From this study,

302    we considered the best-performing inference trajectory, i.e., we mapped the first AWC of four

303    separate layers (0-30, 30-60, 60-100 and 100-200 cm) and then aggregated the maps of the

304    four soil layers to obtain the final SAWC map.

305    3.4. Uncertainty analysis using error propagation

306 In this section, we provide the main details of uncertainty assessment using propagation error.

307 More details of the procedure can be found in (Román Dobarco et al., 2019, Styc and

308 Lagacherie, submitted).

309 The selected inference trajectory, i.e., SAWC estimated as the aggregation of AWC predicted

310 at four depth soil layers, required an error propagation to estimate the variance in SAWC,

311 considered as a proxy of the uncertainty prediction of the target variable (Heuvelink et al.,

312 1989). In this study, we used a first-order Taylor expansion to calculate the error variance of

313 SAWC that results from the error variances of its components (here, the different mapped

314 AWC for the four considered soil layers). This calculation involved i) the error variances of

315 AWC for each soil layer obtained from the conditional distributions provided by QRF for

316 each predicted location (Meinshausen, 2006) and ii) the correlation coefficients between the

317 errors at each soil layer provided by the mapping residuals. Then, the estimate of the SAWC

318 variances was translated into a 90% prediction interval, assuming a normal distribution, by:

$$CIL_i = \hat{y}_i \pm 1.645 \; \sigma_{\hat{y}_i} \qquad (6)$$

319

320 where $CIL_i$ is the interval limits of the prediction, $\hat{y}_i$ is the mean of the distribution, $\sigma_{\hat{y}}$ is the

321 standard deviation and 1.645 is the Student's coefficient for a 90% confidence interval

322 estimation.

323 Error propagation was performed using the *propagate* R package (Spiess, 2018).

324

325     3.5. The experiment

326     The goal of the experiment was two-fold: i) to evaluate the efficiency of the DSM model

327     proposed for dealing with dense spatial sampling of auger holes ($QRF_{dist}$) and ii) to evaluate

328     the cost-efficiency ratio of using auger hole observations with increasing densities.

329     For that, $QRF_{dist}$ was applied to different soil input scenarios with increasing numbers of

330     auger holes. The performances of the $QRF_{dist}$ were compared with those of a baseline QRF

331     application that did not consider any spatial relation between the sites, as practiced in most

332     DSM applications. The four SAWCs presented in section 2.2.3 were considered. In the

333     following, we provide some details about the sampling strategy for selecting auger holes, the

334     evaluation protocol and the cost-benefit analysis.

335        3.5.1.  The sampling procedure of auger holes

336     Different data scenarios were considered, all of which included all the available soil profiles

337     as inputs. An increasing number of auger holes were sampled from the available set and

338     added to the soil profiles in the soil input datasets (from 10% to 100% of the auger hole

339     observations each 10%, e.g., average spacing of 278 m, 556 m, 834 m, 1112 m, 1391 m, 1669

340     m, 1947 m, 2225 m, 2503 m and 2781 m).

341     At each step, the auger holes were selected using a stratified random sampling technique

342     using compact geographical strata (Walvoort et al., 2010), as recommended by (Brus et al.,

343     2011). Thirty-three geographical strata of 0.5 km$^2$ were considered. Spatial stratification

344     sampling was performed using the *spcosa* R package (Walvoort et al., 2018).

345        3.5.2.  Evaluation protocol

346     The performance of the SAWC DSM models was evaluated by k-fold cross validation. This

347     evaluation procedure consisted of randomly dividing the data into k subsets. Then, the holdout

348     method was repeated k times such that one of the k subsets was used as the validation set in

349     each repetition, while the other k-1 subsets were combined to form the calibration set.

350     Following this procedure, every data point was included in a calibration set k-1 times. In this

351     study, we selected k = 10 and to increase the robustness of the evaluation, the 10-fold cross

352     validation was iterated 20 times. The k-fold cross validation was performed using *cvTools*

353     (Alfons, 2012).

354     To avoid uncertain estimations of the model performances due to the inherent uncertainty of

355     SAWC estimations from the auger hole observations, the evaluation protocol presented

356     hereafter was solely applied to the soil profiles.

357     To evaluate the prediction performances, we used classic performance indicators, e.g., the

358     mean square error skill score (Nussbaum et al., 2018), which has the same interpretation as

359     the percentage of variance explained by the model, the root mean square error (RMSE) and

360     the bias.

361     Furthermore, we evaluated the estimation of the prediction uncertainty using the prediction

362     interval coverage probability (PICP; Shrestha and Solomatine, 2006) and error-predicted

363     uncertainty plots. The PICP was computed as follows:

$$PICP = \frac{count(LPL_i \leq y_i \leq UPL_i)}{n} \times 100 \qquad (7)$$

364

365     where $n$ is the total number of observations in the validation set, and the numerator counts if

366     the observation $y_i$ fits within the prediction limits prior to estimation by the error propagation

367     method. For a 90% confidence level, which is usually chosen in DSM studies (Arrouays et al.,

368     2014b), the uncertainty is optimally predicted when the PICP value is close to 90%.

369     The PICP provides an assessment of the overall uncertainty prediction bias (underestimation

370     or overestimation) but does not tell anything about the ability to map differences in

371     uncertainty across the study area. The PICP was therefore completed by error-predicted-

372     uncertainty estimations that materialized the evolution of the cross validation RMSE with the

373    widths of the predicted confidence intervals. To remove noise, the RMSEs were averaged per

374    quartile of prediction interval widths denoted "low/fairly low/fairly high/high predicted

375    uncertainty". It was expected that the RMSE would increase from low to high predicted

376    uncertainty.

377        3.5.3.   The cost efficiency of SAWC Digital Soil Mapping

378    Soil data need to be recorded, but this process can be time consuming and therefore costly. To

379    answer the question, "Is all the data necessary to reach quality predictions?", we set two

380    indicators to assess i) the cost of a unit of gained RMSE and ii) the relative cost efficiency,

381    which were both calculated for each percentage of auger holes added to the soil profiles. The

382    cost of a unit of RMSE was evaluated using the following equation (Eq. 8):

$$Err_{cost} = \frac{cost_i}{RMSE_i} \qquad (8)$$

383

384    where $Err_{cost}$ is the cost of a unit of RMSE (in €/cm) and $RMSE_i$ is the root mean square

385    error of the combination of $i$% of auger hole and soil profile datasets.

386    The relative cost efficiency was assessed following the recommendation of Kish (1965) used

387    by (Viscarra Rossel and Brus, 2018, Eq. 9):

$$CE_r = \frac{cost_{ref} * RMSE_{ref}}{cost_i * RMSE_i} \qquad (9)$$

388    where $CE_r$ is the relative cost-efficiency ratio, $cost_{ref}$ and $RMSE_{ref}$ are the cost and the

389    error of a reference design, respectively, here using solely soil profiles in the SAWC DSM,

390    and $cost_i$ and $RMSE_i$ are the cost and the error, respectively, of the combination of $i$% of

391    auger hole observation and soil profiles. A $CE_r$ larger than one reveals more efficient

392    sampling than the reference (Viscarra Rossel and Brus, 2018).

393

## 4. Results

### 4.1. Preliminary results

In Figure 5, we present the distributions of SAWC30, SAWC60, SAW100 and SAWCmax for the soil profiles (left panel of Figure 5) and auger holes (right panel of Figure 5). We first observed that the distributions of SAWC regardless of the considered soil depth were bimodal for both the soil profiles and auger holes, with i) a higher peak for higher values of SAWC30 and SAWC60 and with ii) a higher peak for lower values of SAWC100 and SAWCmax. Additionally, it is worth noting that both the SAWC ranges and the means of the auger holes were systematically greater than those of the soil profiles. This could be explained by i) possible underestimations of coarse fragments by visual determinations on very small volumes using auger holes compared to real measurements of coarse fragments on larger volumes using soil profiles and ii) possible biases of the field determination of textural class on auger holes compared with laboratory analyses performed on soil profiles.
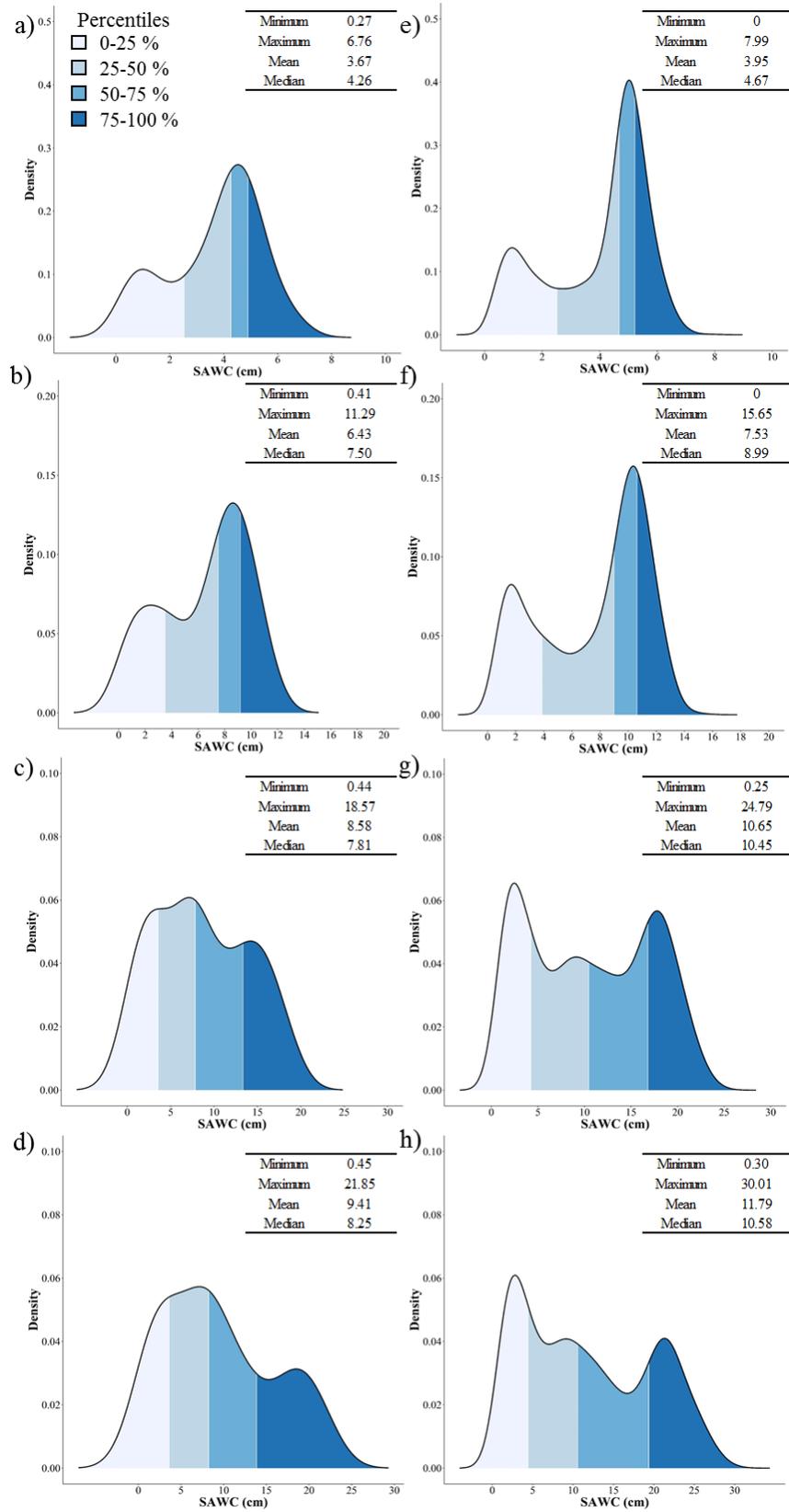
407

408



409

*Figure 5. Distributions of the soil available water capacity of soil profiles at a) 0-30 cm, b) 0-60 cm, c) 0-100 cm and d) 0-depth$_{max}$ and of auger holes at e) 0-30 cm, f) 0-60 cm, g) 0-100 cm and h) 0-depth$_{max}$*
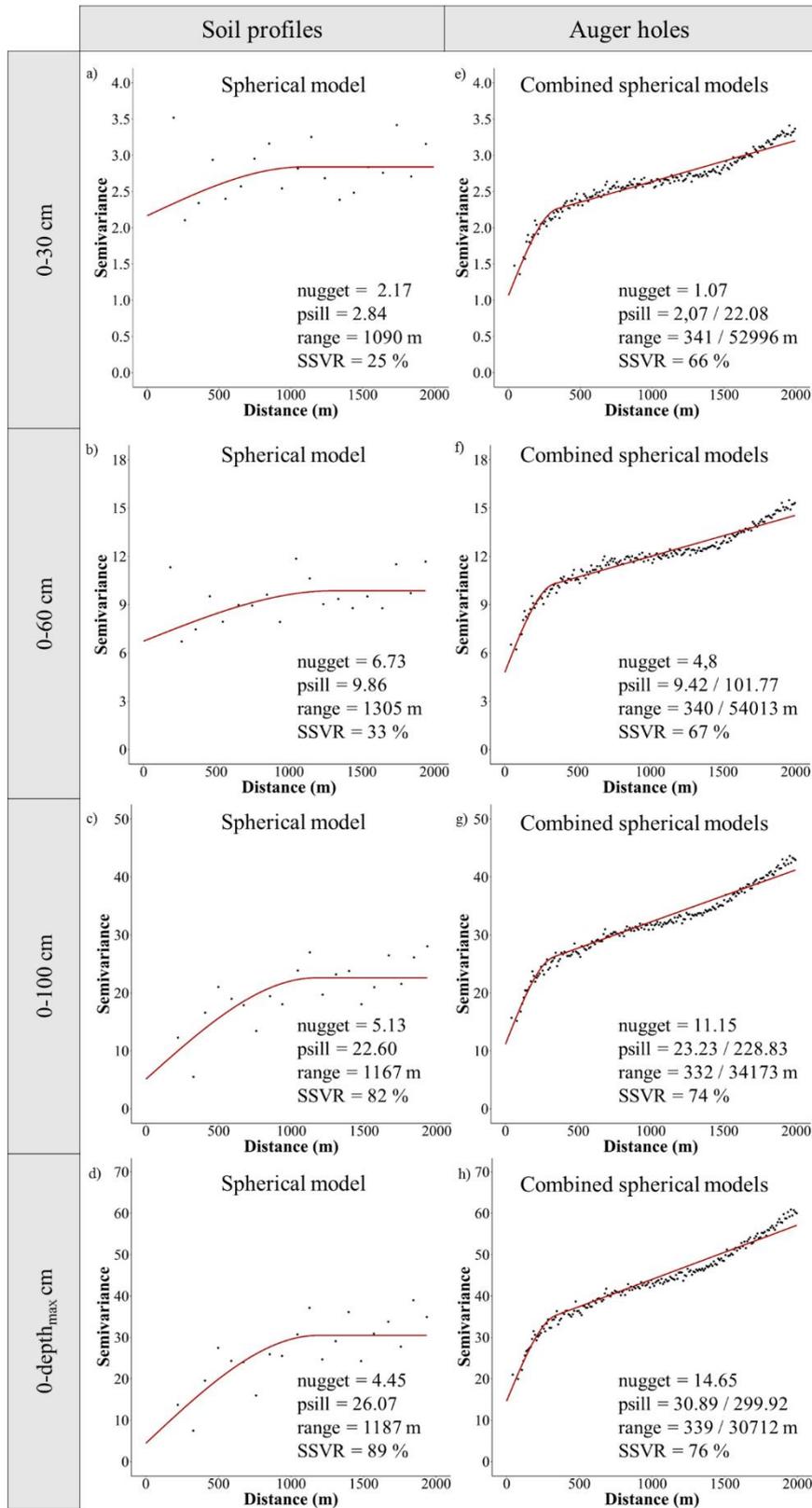
412    In addition, empirical variograms and their fitted models were computed using the *gstat*

413    package (Pebesma, 2004) both from the soil profile data (Figure 6, left panel) and from the

414    auger hole data (Figure 6, right panel), and for the different considered SAWC ( lines of

415    Figure 6). The Spatially structured variance ratio (SSVR, Eq. 10), which estimated the portion

416    of the variance that was spatially structured, was computed from the variograms as follows:

$$SSVR = 1 - \left( nugget \big/ variance \right) x\ 100 \tag{10}$$

417

418    First, we noted that the variogram of the SAWC determined from auger hole observations

419    exhibited clear spatial structures regardless of the maximal depth (SSVR ranging from 66% to

420    76%). The variograms showed a mix of short-range spatial structures (fitted ranges between

421    332 and 341 m) and large-range structures (fitted ranges exceeding 30 km). Conversely, the

422    variograms of SAWC30 and SAWC60 determined from the soil profile empirical variogram

423    exhibited less clear spatial structures (SSVR of 25% and 33%), whereas a clear structure

424    appeared for SAWC100 and SAWCmax (SSVR of 82% and 89%). Because of their larger

425    spacing, the soil profiles did not allow us to see the short-range spatial structures revealed by

426    the auger hole observations. Additionally, significant decreases in nuggets were observed

427    from the variograms of SAWC30 and SAWC60 processed from profiles to those processed

428    from auger holes. This decrease can be interpreted as the result of increasing sampling

429    densities that better captured the short-range spatially structured variance that was otherwise

430    included in the profile variogram nuggets. It is interesting to note that the converse occurred

431    for SAWC100 and SAWCmax. The probable increase in the uncertainty of observations with

432    depth due to the difficulties in observing deep horizons from auger holes yielded a nugget

433    increase that largely counterbalanced the effect of the sampling density evoked previously.

434

435

Figure 6. Empirical variograms computed for SAWC using 69 soil profiles at a) 30 cm, b) 60 cm, c) 100 cm and d) 200 cm
and using 2781 auger hole observations at e) 30 cm, f) 60 cm, g) 100 cm and h) 200 cm, and their theoretical variograms.

436
437

438

439       4.2.    Comparing DSM model prediction and uncertainty prediction performances

440      Table 3 shows the prediction and the uncertainty prediction performances of the two

441      considered DSM models in predicting the SAWCs at four different depths. Only the extreme

442      data scenario, i.e., no auger hole vs. the whole set of auger holes, is shown.

443      First, better performances of SAWC predictions were generally obtained by adding the auger

444      hole observations, with the noticeable exceptions of the predictions of SAWC60, SAWC100

445      and SAWCmax using a classical (nonspatial) QRF. When using QRF$_{dist}$, the performance

446      increases by adding auger hole observations tended to decrease as the maximum considered

447      depth increased.

448      Additionally, using QRF$_{dist}$ that included geographical information led to better prediction

449      performances regardless of the SAWC only when the auger hole observations were added to

450      the soil profiles. Otherwise, (i.e., when only the soil profiles were used for calibrating the

451      model), using QRF yielded equal or slightly better prediction performances.

452      Concerning the ability of the models to provide unbiased estimates of prediction uncertainty,

453      as measured by the PICP, larger PICP values were obtained with QRF$_{dist}$ than with QRF,

454      except for the PICP for SAWC100 with only soil profiles. Furthermore, the effects of

455      including auger holes in QRF calibration were different according to the selected model: the

456      PICP decreased when QRF was selected, whereas the PICP increased when the QRF$_{dist}$ model

457      was selected. As far as the closeness to the nominal value of 90% is concerned, better results

458      were generally obtained when the auger hole observations were not used, with the noticeable

459      exception of the SAWC30 predictions using QRF. Furthermore, QRF$_{dist}$ had more PICP

460      values close to the 90% nominal value ($< 2\%$) than did QRF (4 out of 8 vs. 1 out 8).

461

462

463

Table 3. Prediction and uncertainty prediction performances of SAWC using multiple DSM models

| DSM models | | QRF | | | | $QRF_{dist}$ | | | |
|---|---|---|---|---|---|---|---|---|---|
| SAWC | Auger holes portion (%) | $SS_{MSE}$ | RMSE (cm) | Bias (cm) | PICP (%) | $SS_{MSE}$ | RMSE (cm) | Bias (cm) | PICP (%) |
| SAWC30 | 0 | 0.04 | 1.66 | 0.17 | 86 | -0.02 | 1.71 | 0.32 | 85 |
| | 100 | 0.38 | 1.34 | 0.49 | 86 | 0.49 | 1.22 | 0.37 | 90 |
| | | | | | | | | | |
| SAWC60 | 0 | 0.33 | 2.74 | 1.08 | 87 | 0.3 | 2.79 | 0.35 | 89 |
| | 100 | 0.32 | 2.76 | 1.28 | 83 | 0.54 | 2.26 | 0.82 | 93 |
| | | | | | | | | | |
| SAWC100 | 0 | 0.55 | 3.73 | -0.47 | 92 | 0.46 | 3.97 | 0.22 | 90 |
| | 100 | 0.43 | 4.06 | 1.82 | 85 | 0.63 | 3.27 | 1.09 | 95 |
| | | | | | | | | | |
| SAWCmax | 0 | 0.61 | 4.01 | -0.68 | 90 | 0.53 | 4.41 | -0.56 | 91 |
| | 100 | 0.54 | 4.37 | 1.88 | 85 | 0.7 | 3.54 | 0.18 | 96 |

464

465 As expected, the averaged RMSE tended to increase with the widths of the confidence

466 intervals predicted by $QRF_{dist}$ (Table 4), which demonstrated the overall validity of the

467 uncertainty predictions. However, non-monotonous increases were observed for the SAWC

468 predictions at small depths that also exhibited the weakest performances (Table 3). This non-

469 monotonousness was clearer when the auger hole observations were added. Similar trends

470 were observed for the confidence interval widths predicted by QRF (results not shown).

471

472

| Rooting depth (cm) | Uncertainty | RMSE (cm) | |
|---|---|---|---|
| | | Soil profiles | Soil profiles and auger holes |
| 30 | Low | 1.09 | 1.31 |
| | Fairly low | 1.25 | 0.79 |
| | Fairly high | 2.75 | 1.10 |
| | High | 1.9 | 1.59 |
| | | | |
| 60 | Low | 2.31 | 1.25 |
| | Fairly low | 2.24 | 2.02 |
| | Fairly high | 2.81 | 3.25 |
| | High | 3.46 | 2.08 |
| | | | |
| 100 | Low | 2.81 | 1.52 |
| | Fairly low | 2.82 | 2.81 |
| | Fairly high | 3.49 | 3.69 |
| | High | 5.71 | 4.32 |
| | | | |
| Maximum observation depth | Low | 3.07 | 2.24 |
| | Fairly low | 2.88 | 2.82 |
| | Fairly high | 4.55 | 4.20 |
| | High | 6.09 | 4.37 |

473

474
475
*Table 4. Error-predicted uncertainty results of QRF$_{dist}$ using only soil profiles and using soil profiles and auger hole observations for predicting SAWC at multiple depths*

476    4.3.  Spatial distribution of the SAWC and its associated uncertainty

477 All the predicted maps of SAWC (Figure 7) exhibited spatial patterns of variation that were

478 globally in accordance with the lithological variations shown in Figure 1. The highest values

479 of SAWC were predicted in the northeastern section of the study area with fluvisols

480 developed on loess. The smallest values corresponded to chromic luvisols developed on the
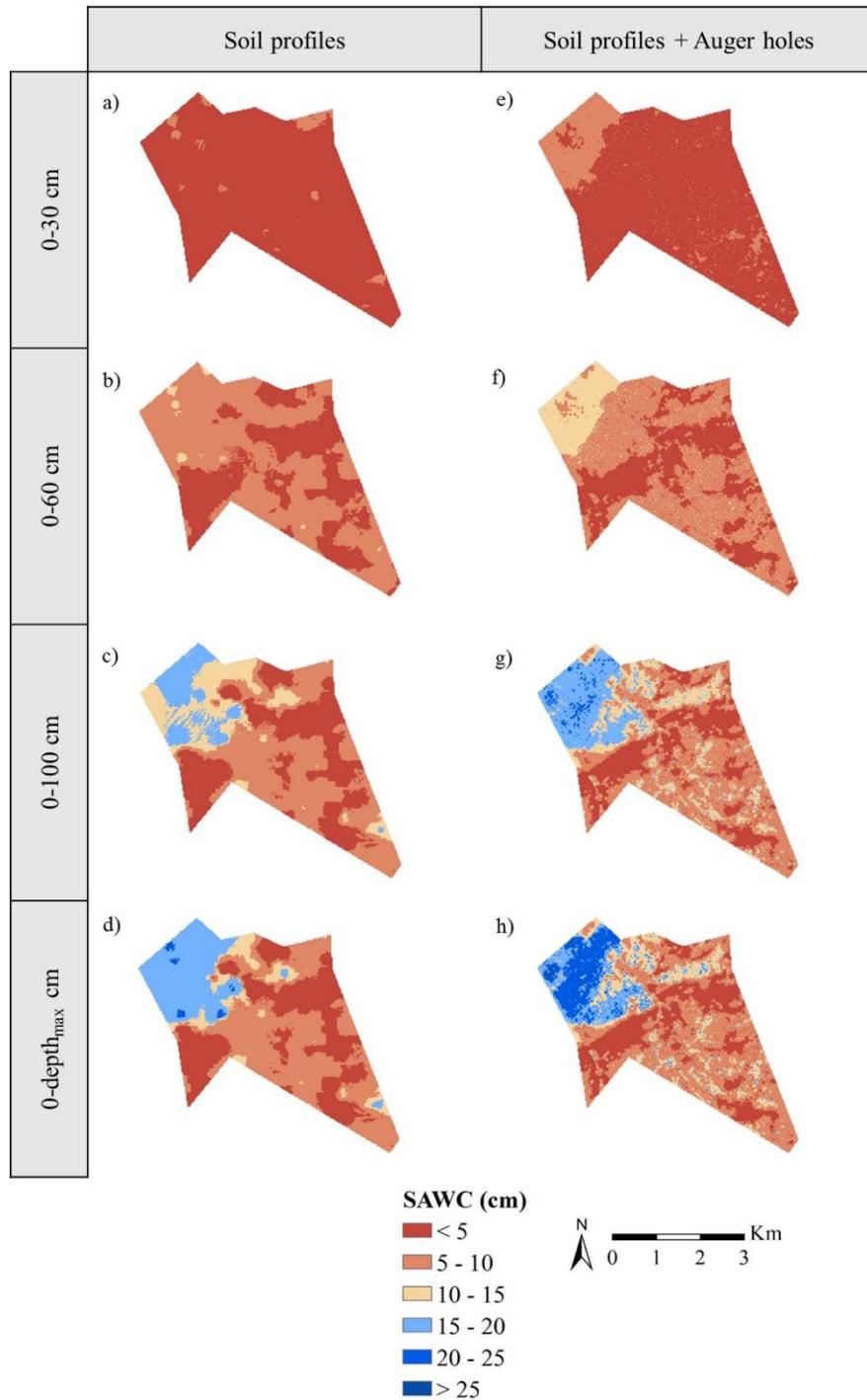
481 old stony alluvial deposits.

482 The spatial pattern became increasingly clear and contrasted as the considered soil depth for

483 calculating the SAWC increased (from the top to the bottom of Figure 7). The incorporation

484    of auger holes (from the left to the right column in Figure 7) led to i) an increase in the

485    predicted variabilities of the SAWC, leading to more contrasted patterns regardless of the

486    predicted SAWC; ii) an increase in the spatial resolution of the SAWC pattern delineations,

487    showing very fine details of variation; iii) the removal of some obvious artifacts of the map of

488    SAWC100 obtained from the soil profiles (Figure 7c); and iv) the addition of some artifacts

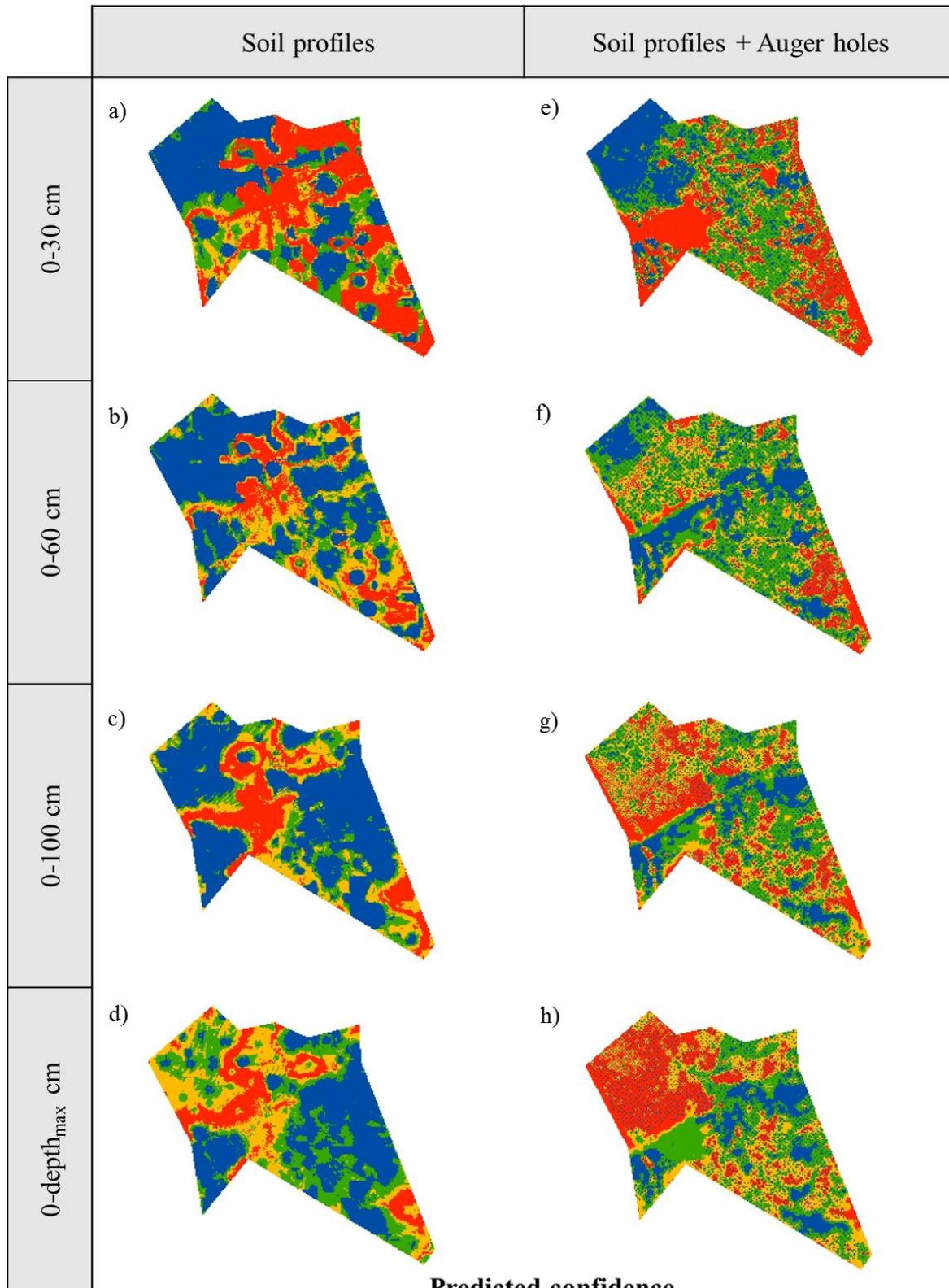489    (isolated pixels) in the SAWC30 and SAWC60 maps (Figure 7e and 7f).

490

*Figure 7. Predicted maps of SAWC over Bouillargues using QRF$_{dist}$ with soil profiles for predicting a) SAWC30, b) SAWC60, c) SAWC100, and d) SAWCmax and using QRF$_{dist}$ with soil profiles and auger hole observations for predicting e) SAWC30, f) SAWC 60, g) SAWC100, and f) SAWCmax*
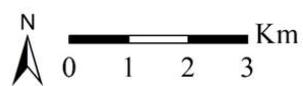
The uncertainty maps of SAWC predictions (Figure 8) obtained from the QRF$_{dist}$ model exhibited spatial patterns that were both complex and very contrasted across the predicted SAWCs and soil inputs. When examining the variations between quartiles of predicted uncertainty that looked significant according to the error-predicted uncertainty results (Table

499    4), some of the maps revealed strong spatial pattern similarities with those of some

500    uncertainty drivers, i.e., the SAWC30 uncertainty map using soil profiles (Figure 8a) with the

501    lithology map (Figure 1), SAWC100 map using soil profiles (Figure 8c) with the spatial

502    density of soil profiles that is observable on the map of soil profiles (Figure 2a), SAWC30

503    uncertainty map using auger hole observations (Figure 8e) with the spatial density of auger

504    hole observations that is observable on the map of auger hole observations (Figure 2b),

505    SAWCmax uncertainty map using auger hole observations (Figure 8h) with the predicted map

506    of SAWCmax. The other uncertainty maps (Figure 8b, 8d, 8f) showed less interpretable

507    patterns, with probably mixed impacts of the above evoked drivers.

|  | Soil profiles | Soil profiles + Auger holes |
|---|---|---|
| 0-30 cm | a) | e) |
| 0-60 cm | b) | f) |
| 0-100 cm | c) | g) |
| 0-depth$_{max}$ cm | d) | h) |

**Predicted confidence interval width**

Low

Fairly low

Fairly high

High

N

0   1   2   3   Km

508

509

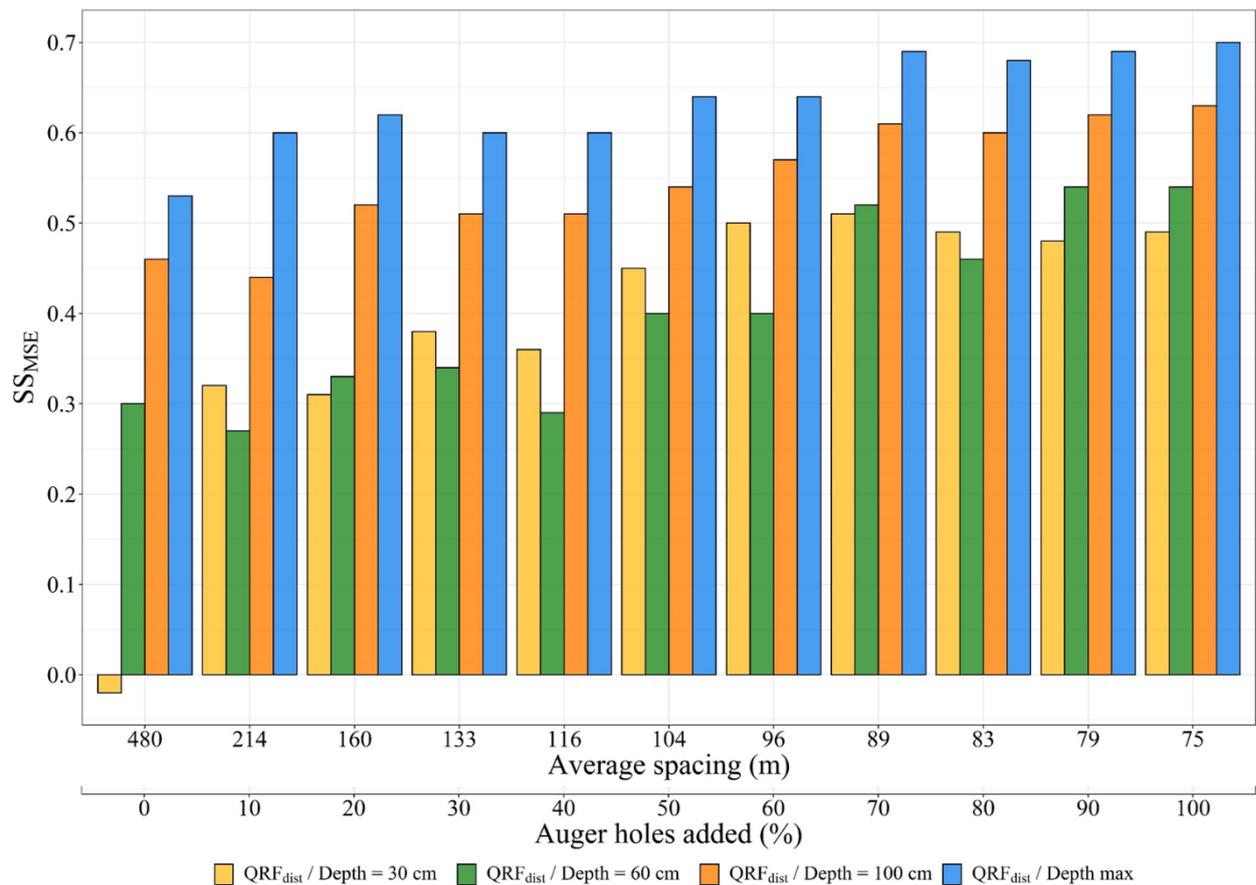### 4.4. Comparing the spatial densities of auger hole observations

In Figure 9, we present the evolution of the $SS_{MSE}$ with the increasing number of auger hole observations in the calibration process. The density in the number of observations/km$^2$ is also expressed as the average spacing between observation sites, which means that the density increases as the average spacing decreases. The average spacing between observation sites was estimated as follows:

$$Average\ spacing = \sqrt{\frac{total\ area}{size}} \tag{11}$$

As already observed from Table 3, the general trend was an increase in performance as the number of auger hole observations increased regardless of the maximal depth at which the SAWCs were calculated. However, some local decreases in performance were observed, e.g., on SAWC60 and 100 predictions when adding 10% auger holes or on SAWC100 and SAWCmax predictions when passing from 20 to 30% auger holes. Conversely, the addition of 10% to 20% auger holes and 60% to 70% auger holes seemed beneficial for all predictions of the SAWC.

527

528

529

530

531

532 When considering the costs of adding new auger hole observations according to the two cost-

533 efficiency indicators described in section 3.5., it appeared that the cost of gaining one unit of

534 RMSE (the error cost, $Err_{cost}$) was important until the first addition of the auger hole and
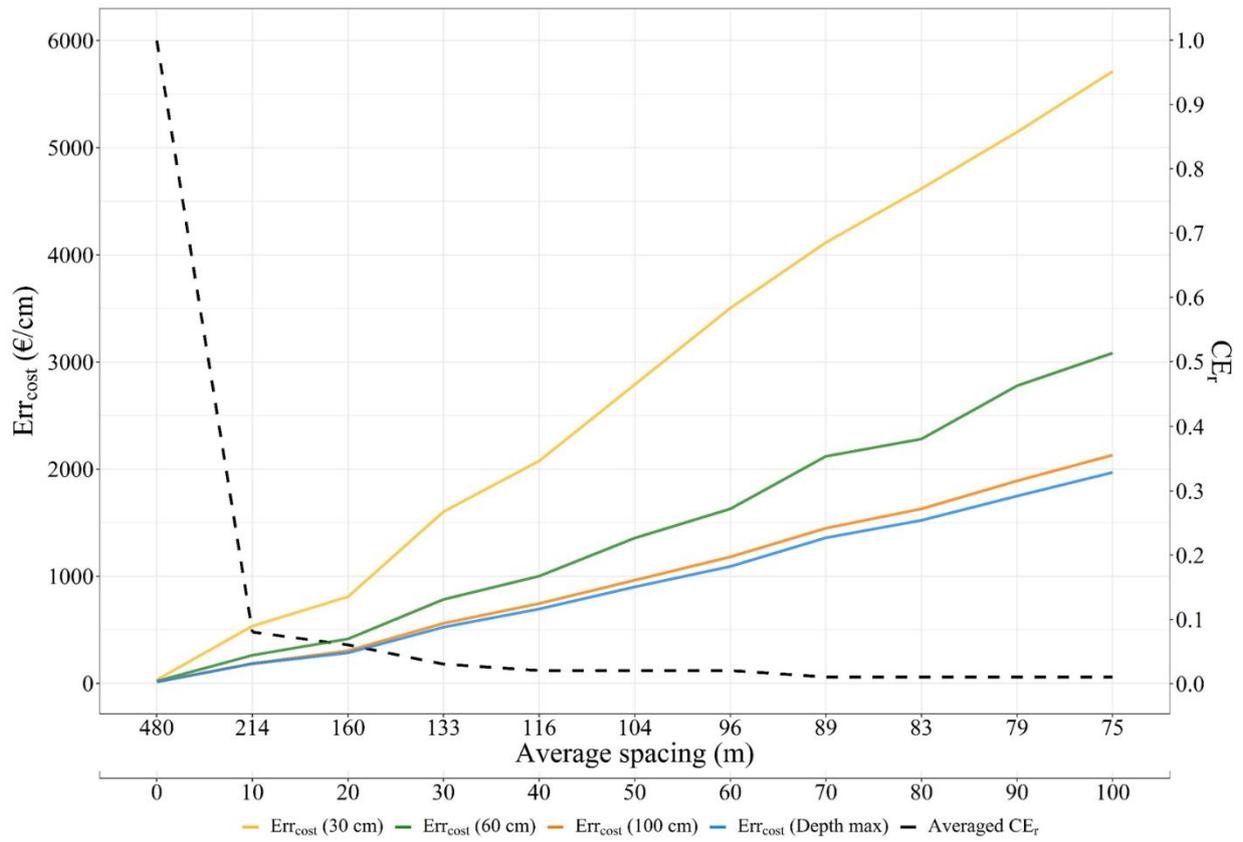
535 further linearly increased as new auger holes were added (Figure 10). This is translated by the

536 relative cost-efficiency ratio ($CE_r$) by a dramatic decrease under the 1:1 ratio when adding the

537 first auger hole observations and then a slow decrease for further additions.

*Figure 10. Cost-efficiency ratios according to the average spacing related to the number of auger holes*

## 5. Discussion

### 5.1. Soil Available Water Capacity

The selected case study considered the soil available water capacity, which is among the most highly demanded properties of end users, as the targeted soil property (Richer de Forges et al, 2019). This paper completes the small set of papers that were devoted to the digital mapping of SAWC (Hong et al., 2013; Malone et al., 2009; Padarian et al., 2014; Poggio et al., 2010; Román Dobarco et al., 2019; Ugbaje and Reuter, 2013, Amirian-Chackan et al., 2019) and the even smaller set of papers that addressed all the SAWC components as defined by the original definition reported by Cousin et al. (2003) (Eq. 1) (Leenaars, 2018; Romàn Dobarco et al., 2019; Styc and Lagacherie, 2019, submitted).

However, as in many DSM applications, the SAWC was determined at local sites without the full measurements of its components. Visual estimations of the coarse fragment content and of the soil depth generated observational uncertainties and, for the latter, right-censored estimations due to the limitation in observation depths. Furthermore, the water retention capacity of each horizon was not fully measured, although it is worth noting that some components of this retention that are usually not measured (bulk density, field capacity) were measured here on the soil profiles. To overcome the measurement limitations, pedotransfer functions were used (see section 2.2.3). It is worth noting that these pedotransfer functions were highly case specific both regarding their input (textural classes + field capacity measurements) and their target (the b coefficient). The addition of all these peculiar uncertainties should result in a significant overall uncertainty of the soil inputs that is well reported by the nuggets of the variograms of the densest datasets (Figure 6, right panel). This uncertainty may greatly explain the limitation of performances that was observed, even for the densest datasets.

### 5.2. The interest of "spatial RFs"

575    Our results showed that the SAWC prediction performances were nearly systematically

576    increased by adding some geographical information, i.e., the n of "scorpan" in McBratney et

577    al.'s (2003) formula, to the set of candidate covariates used in a random forest. This

578    confirmed the results obtained by Hengl et al. (2018) from various case studies. This,

579    however, enriched these results by showing that the gains in performances provided by the

580    addition of geographical covariates depend on the density of the sampling. Indeed, these gains

581    were only effective when the dense sampling of auger hole observations was used (76 m

582    spacing), whereas the low density of soil profiles did not provide clear improvements (Table

583    3). At high density levels, the classical landscape covariates were not sufficient to account for

584    the variability shown in the dataset of soil inputs as represented by the variograms of Figure 6

585    (right panel), whereas the proximity effects brought by the geographical covariates allowed us

586    to overcome this limitation.

587    In digital soil mapping, proximity effects have been traditionally addressed by using

588    regression kriging (Hengl et al., 2004; Malone et al., 2009; Vaysse and Lagacherie, 2015).

589    However, spatial QRF was demonstrated to have similar performances (Hengl et al, 2018)

590    while having some decisive advantages in the context of our case study. Spatial QRF does not

591    require any rigid statistical assumptions about the distribution and the stationarity of the target

592    variable, which allows us to handle the bimodal distributions of SAWCs (Figure 5). It also

593    does not require any geostatistical expertise for the manual fitting of variograms, which opens

594    the possibility to fully automate the procedure so that non pedometrician, such as BRL staff,

595    could use it for the other communes of the irrigation perimeter.

596        5.3.  The interest of adding auger hole observations

597    The addition of dense spatial sets of auger hole observations in the modeling process

598    significantly increased the level of performance when considering the best model ($QRF_{dist}$),

599    which is in accordance with several previous experiments studying the impact of soil

600    sampling densities (Somarathna et al. 2017, Wadoux et al. 2019 and Lagacherie et al., 2020).

601    The performances observed in this case study were better than those in most of the published

602    DSM applications dealing with SAWC (Ugbaje and Reunter, 2013; Styc and Lagacherie,

603    2019, submitted), which was the result of a much greater spatial density of the soil inputs

604    (from $6/km^2$ to $26/km^2$) than in these previous applications (from $0.01/km^2$ to $0.05/km^2$)).

605    However, strong limitations in the SAWC prediction performances were still observed, even

606    when using the most dense set of auger hole observations. These limitations increased as the

607    maximum depth at which the SAWC was calculated decreased (Table 3). This means that

608    significant proportions of the SAWC variabilities were not mapped despite the large densities

609    of the auger hole observations used as input. To explain this fact, it is first interesting to note

610    that for both the soil profiles and the soil profiles plus auger hole inputs, the performances and

611    the spatially structured variance ratios of the input soil datasets were ranked similarly across

612    SAWCs and spatial densities (Figure 6), which was already observed in the same region for

613    different soil properties and study extents by Vaysse and Lagacherie (2015). Concerning the

614    results using solely the profiles, this revealed that a part of the short-range variability shown

615    by the variograms built from auger holes (Figure 6, left panel) was not captured by the soil

616    dataset because of a limitation in spacing. However, this limitation decreased as the

617    considered depth of the SAWC calculation increased, which explained the observed increase

618    in performance from SAWC30 to SAWCmax. Concerning the results using the auger hole

619    observations, a similar trend was observed since the local uncertainty as revealed by the

620    variogram nuggets (Figure 6, right panel) remained important due to observational uncertainty

621    (see section 5.1.), which may induce noise that may perturb the calibration of the QRF model.

622    Finally, it should be recalled that these performances were calculated for predictions of the

623    SAWC at precise locations, whereas SAWC is required for field or in-field management

624      zones for most of the decision making. It could be expected that these performances would

625      increase when the SAWC prediction will be spatially aggregated (Vaysse et al, 2017).

626          5.4.   Uncertainty predictions

627      Since SAWC is a soil functional property composed of several primary soil properties,

628      uncertainty predictions were provided by a specific error model previously proposed by

629      (Román Dobarco et al., 2019) and further refined by Styc and Lagacherie (submitted). The

630      uncertainty predictions were classically evaluated with regard to their unbiasness (PICP,

631      Table 3). They were also evaluated for their ability to identify contrasted uncertainty areas

632      (comparisons between residuals and predicted uncertainty, Table 4), which, to our knowledge,

633      has never been done in the DSM literature before Styc and Lagacherie (submitted). The

634      results were highly variable across models and spatial densities. However, the more accurate

635      models tended to also provide the best pictures of the uncertainty patterns (Figure 6) with an

636      overestimation of uncertainty ($QRF_{dist}$ on Table 3). This overestimation was already observed

637      by Lagacherie et al. (2020) and was assumed to be due to the inclusion of outliers as the

638      average spacing decreased, which probably disturbs the limit estimations of the confidence

639      interval. On the other hand, a part of the inaccuracy of uncertainty predictions may come from

640      the differences (Figure 5)  between the distributions of SAWC values calculated from auger

641      holes (used as calibration data only) and from soil profiles (used as evaluation data). More

642      attention must be paid in the future to uncertainty predictions in view of identifying the

643      possible causes of these uncertainty mispredictions.

644      It is interesting to note that some of the produced uncertainty maps showed strong similarities

645      with possible drivers (see comments of Figure 8), which can be interpreted from our common

646      sense pedological knowledge. The largest uncertainties were estimated i) in chromic Luvisols

647      (Figure 6a) because of the large rates of coarse fragment content that are known to be difficult

648      to quantify in the field, ii) in areas of lower densities of soil observations (Figure 6c and 6e)

649    because of difficulties of model calibration at these locations and iii) for the largest predicted

650    values of SAWCs with the best models (Figure 6h) because the estimates of relative

651    uncertainty reached an unsurmountable floor that is likely related to the observational

652    uncertainty. All these observations reinforce the credibility of the presented uncertainty maps.

653

654          5.5. The level of performance obtained and cost.

655    The use of auger hole observations as complementary soil input to soil profiles led to a

656    substantial increase in performance, but the harvesting process was very time consuming,

657    which resulted in high costs (see section 4.5). Figure 9 curves show that the performance

658    gains were obtained by increasing costs as the density of the auger holes increased. A

659    compromise should then be found, which can be formulated as "the number of auger hole

660    observations that reach an acceptable level of performance while keeping an acceptable cost

661    level". The cost indicator curves of Figure 10 did not reveal a clear compromise. However,

662    such curves could be used with a prior definition of what performance and costs are

663    acceptable. Furthermore, such cost curves could be improved if either more sophisticated

664    sampling is used (e.g., van Groningen et al, 1998) or if the harvesting costs could be reduced

665    by a partial automation of digitizing procedures (Yang and Yang, 2017).

666    Finally, it should be stressed that the quantitative evaluation of prediction performance that

667    served as a basis for building the curve costs should be completed by a qualitative

668    examination of the maps. As revealed by the spatial patterns of the predicted SAWC maps,

669    considerable gains in spatial resolution were obtained by adding auger holes, which may

670    enable field-level decision making. This may constitute a more decisive added value than the

671    moderate gain in precision quantitatively evaluated by the cost indicators.

672         6.    Conclusion

673    In this study, the main lessons were as follows:

674    • A QRF approach using euclidian buffer distances outperformed a classical QRF
675       approach in predicting SAWC with a dense set of profiles and auger holes

676    • The addition of a dense spatial sampling of auger hole observations dramatically
677       increased the performance in predicting SAWCs and increased the spatial resolutions
678       of the SAWC pattern delineations, but there were limitations due to the uncertainty of
679       the auger hole observations.

680    • The performances in predicting SAWC values varied following some drivers that were
681       expected  - average spacing of sites, and type of observations (profiles vs. auger holes)
682       -   and following other drivers that were revealed by the uncertainty maps   –
683       pedological context, local  density of sites, SAWC predicted values – (see section
684       5.4.).

685    • The cost-efficiency analysis did not reveal a clear compromise in terms of limiting the
686       costly harvesting of auger hole data. Rather, the compromise should be user specific
687       and should be updated as soon as partial automation is possible (see section 5.5)

688

689    7.    Acknowledgments

694

695    8.    References

696    Alfons, A. (2012) cvTools: Cross-Validation Tools for Regression Models. R Package.

697    Available online: https://cran.rproject.org/web/packages/cvTools/cvTools.pdf. (accessed on

698    11 May 2012).

699    Amirian-Chakan, Alireza, Budiman Minasny, Ruhollah Taghizadeh-Mehrjardi, Rokhsar

700        Akbarifazli, Zahra Darvishpasand, and Saheb Khordehbin. (2019). Some Practical

701        Aspects of Predicting Texture Data in Digital Soil Mapping. *Soil and Tillage Research*

702        194:104289.

703    Arrouays, D., Lagacherie, P., & Hartemink, A. E. (2017). Digital soil mapping across the

704        globe. *Geoderma Regional*. https://doi.org/10.1016/j.geodrs.2017.03.002

705    Arrouays, D., Leenaars, J. G. B., Richer-de-Forges, A. C., Adhikari, K., Ballabio, C., Greve,

706        M., … Rodriguez, D. (2017). Soil legacy data rescue via GlobalSoilMap and other

707        international and national initiatives. *GeoResJ*. https://doi.org/10.1016/j.grj.2017.06.001

708    Baize, D., & Jabiol, B. (1995). *Guide des sols, Quae*.

709    Böhner, J., McCloy, K.R., Strobl, J., 2006. SAGA — analysis and modelling applications

710    Göttinger Geogr. Abh., 115, p. 130

711    Breiman, L. Random forests. Mach. Learn. 2001, 45, 5–32.

712    Brus, D. J., Kempen, B., & Heuvelink, G. B. M. (2011). Sampling for validation of digital soil

713        maps. *European Journal of Soil Science*, *62*(3), 394–407. https://doi.org/10.1111/j.1365-

714        2389.2011.01364.x

715    Cousin, I.,Nicoullaud, B., Coutadeur, C. (2003). Influence of rock fragments on the water

716        retention and water percolation in a calcareous soil. Catena, 53, pp. 97-114

717    Dominati, E., Mackay, A., Green, S., & Patterson, M. (2014). A soil change-based

718        methodology for the quantification and valuation of ecosystem services from agro-

719    ecosystems: A case study of pastoral agriculture in New Zealand. *Ecological Economics*.

720    https://doi.org/10.1016/j.ecolecon.2014.02.008

721    Gomez, C., & Coulouma, G. (2018). Importance of the spatial extent for using soil properties

722    estimated by laboratory VNIR/SWIR spectroscopy: Examples of the clay and calcium

723    carbonate content. *Geoderma*. https://doi.org/10.1016/j.geoderma.2018.06.006

724    Hengl, T. (2019). GSIF : Global Soil Information Facilities. R package version 0.5-5.

725    https://CRAN.R-project.org/package=GSIF.

726    Hengl, T., De Jesus, J. M., MacMillan, R. A., Batjes, N. H., Heuvelink, G. B. M., Ribeiro, E.,

727    … Gonzalez, M. R. (2014). SoilGrids1km - Global soil information based on automated

728    mapping. *PLoS ONE*, *9*(8). https://doi.org/10.1371/journal.pone.0105992

729    Hengl, T., Heuvelink, G. B., Kempen, B., Leenaars, J. G., Shepherd, K., Sila, A., … Tondoh,

730    J. E. (2015). Mapping soil properties of Africa at 250 m resolution: random forests

731    significantly improve current predictions S1 Regression-kriging in R using the Meuse

732    data         set.         *PlosOne*,         *10*(6),         e0125814.

733    https://doi.org/https://doi.org/10.1371/journal.pone.0125814

734    Hengl, T., Nussbaum, M., Wright, M. N., Heuvelink, G. B. M., & Gräler, B. (2018). Random

735    forest as a generic framework for predictive modeling of spatial and spatio-temporal

736    variables. *PeerJ*, *2018*(8). https://doi.org/10.7717/peerj.5518

737    Heuvelink, G. B. M., Burrough, P. A., & Stein, A. (1989). Propagation of errors in spatial

738    modelling with GIS. *International Journal of Geographical Information Systems*, *3*(4),

739    303–322. https://doi.org/10.1080/02693798908941518

740    Hong, S. Y., Minasny, B., Han, K. H., Kim, Y., & Lee, K. (2013). Predicting and mapping

741    soil available water capacity in Korea. *PeerJ*, *1*, e71. https://doi.org/10.7717/peerj.71

742    Kish, L. (1965). Survey sampling. New York, London: John Wiley

743    Lagacherie, P., Arrouays, D., Bourennane, H., Gomez, C., & Nkuba-Kasanda, L. (n.d.).

744        Analysing the impact of soil spatial sampling on the performances of Digital Soil

745        Mapping models and their evaluation: a numerical experiment using clay contents

746        obtained from Vis-NIR-SWIR hyperspectral imagery. *Geoderma*, 375, 114503

747    Leenaars, J. G. B., Claessens, L., Heuvelink, G. B. M., Hengl, T., Ruiperez González, M., van

748        Bussel, L. G. J., … Cassman, K. G. (2018). Mapping rootable depth and root zone plant-

749        available water holding capacity of the soil of sub-Saharan Africa. *Geoderma*.

750        https://doi.org/10.1016/j.geoderma.2018.02.046

751    Legros, J. P. (1996). *Cartographies des sols: de l'analyse spatiale à la gestion des territoires*.

752        In      *Collection      Gérer      l'environnement*.      Retrieved      from

753        https://books.google.fr/books?id=MiONzDc-jnQC

754    Malone, B. P., McBratney, A. B., Minasny, B., & Laslett, G. M. (2009). Mapping continuous

755        depth functions of soil carbon storage and available water capacity. *Geoderma*, *154*(1–

756        2), 138–152. https://doi.org/10.1016/j.geoderma.2009.10.007

757    McBratney, A. B., Mendonça Santos, M. L., & Minasny, B. (2003). On digital soil mapping.

758        *Geoderma*. https://doi.org/10.1016/S0016-7061(03)00223-4

759    Meinshausen, N. (2006). Quantile regression forests. *Journal of Machine Learning Research*,

760        *7*, 983–999.

761    Minasny, B., & McBratney, A. B. (2016). Digital soil mapping: A brief history and some

762        lessons. *Geoderma*, *264*, 301–311. https://doi.org/10.1016/j.geoderma.2015.07.017

763    Mulder, V. L., Lacoste, M., Richer-de-Forges, A. C., & Arrouays, D. (2016). GlobalSoilMap

764        France: High-resolution spatial modelling the soils of France up to two meter depth.

*Science of the Total Environment*, *573*, 1352–1369. https://doi.org/10.1016/j.scitotenv.2016.07.066

Nussbaum, M., Spiess, K., Baltensweiler, A., Grob, U., Keller, A., Greiner, L., … Papritz, A. (2018). Evaluation of digital soil mapping approaches with large sets of environmental covariates. *Soil*, *4*(1), 1–22. https://doi.org/10.5194/soil-4-1-2018

Padarian, J., Minasny, B., McBratney, A. B., & Dalgliesh, N. (2014). Predicting and mapping the soil available water capacity of Australian wheatbelt. *Geoderma Regional*, *2–3*(C), 110–118. https://doi.org/10.1016/j.geodrs.2014.09.005

Pebesma, E. J. (2004). Multivariable geostatistics in S: The gstat package. *Computers and Geosciences*, *30*(7), 683–691. https://doi.org/10.1016/j.cageo.2004.03.012

Richer-de-forges, A.C., Arrouays, D., Bardy, M., Bispo, A., Lagacherie, P., Laroche, B., Lemercier, B., Sauter, J., Voltz, M., (2019). Mapping of Soils and Land-Related Environmental Attributes in France : Analysis of End-Users ' Needs. Sustainability 11, 1–15.

Román Dobarco, M., Bourennane, H., Arrouays, D., Saby, N. P. A., Cousin, I., & Martin, M. P. (2019). Uncertainty assessment of GlobalSoilMap soil available water capacity products: A French case study. *Geoderma*, *344*(February), 14–30. https://doi.org/10.1016/j.geoderma.2019.02.036

Sanchez, P. A., Ahamed, S., Carré, F., Hartemink, A. E., Hempel, J., Huising, J., … Zhang, G. L. (2009). Digital soil map of the world. *Science*, *325*(5941), 680–681. https://doi.org/10.1126/science.1175084

Somarathna, P. D. S. N., Minasny, B., & Malone, B. P. (2017). More Data or a Better Model? Figuring Out What Matters Most for the Spatial Prediction of Soil Carbon. *Soil Science*

788    *Society       of       America       Journal*,       *81*(6),       1413–1426.

789    https://doi.org/10.2136/sssaj2016.11.0376

790    Shrestha, D.L., Solomatine, D.P., (2006). Machine learning approaches for estimation of

791    prediction interval for the model output. Neural Networks 19 (2), 225–235

792    Spiess, A-N., (2018). Propagate: Propagation of Uncertainty. R package version 1.0-6.

793    Available on : https://CRAN.R-project.org/package-propagate

794    Strobl, C., Malley, J., & Tutz, G. (2009). An Introduction to Recursive Partitioning:

795    Rationale, Application, and Characteristics of Classification and Regression Trees,

796    Bagging,    and    Random    Forests.    *Psychological    Methods*,    *14*(4),    323–348.

797    https://doi.org/10.1037/a0016973

798    Styc, Q., & Lagacherie, P. Uncertainty assessment of soil available water capacity using error

799    propagation : a test in Languedoc Roussillon. Submitted on Geoderma since 24/01/2020.

800    Styc, Q., & Lagacherie, P. (2019). What is the Best Inference Trajectory for Mapping Soil

801    Functions : An Example of Mapping Soil Available Water Capacity over Languedoc

802    Roussillon    (    France    ).    *Soil    Systems*,    *3*(34),    17.

803    https://doi.org/10.3390/soilsystems3020034

804    Ugbaje, S. U., & Reuter, H. I. (2013). Functional Digital Soil Mapping for the Prediction of

805    Available Water Capacity in Nigeria using Legacy Data. *Vadose Zone Journal*, *12*(4), 0.

806    https://doi.org/10.2136/vzj2013.07.0140

807    Van Groenigen, J. W., Stein, A., (1998). Constrained optimization of spatial sampling using

808    continuous simulated annealing. Journal of Environmental Quality 27 (5), 1078−1086.

809    Vaysse, K., & Lagacherie, P. (2015). Evaluating Digital Soil Mapping approaches for

810    mapping GlobalSoilMap soil properties from legacy data in Languedoc-Roussillon

811        (France). *Geoderma Regional*, *4*, 20–30. https://doi.org/10.1016/j.geodrs.2014.11.003

812    Vaysse, Kévin, & Lagacherie, P. (2017). Using quantile regression forest to estimate

813        uncertainty of digital soil mapping products. *Geoderma*, *291*, 55–64.

814        https://doi.org/10.1016/j.geoderma.2016.12.017

815    Veihmeyer, F.J., Hendrickson, A.H. 1927. The relation of soil moisture to cultivation and

816        plant growth. Soil Sci. 3, 498–513.

817    Viscarra Rossel, R. A., & Brus, D. J. (2018). The cost-efficiency and reliability of two

818        methods for soil organic C accounting. *Land Degradation and Development*, *29*(3), 506–

819        520. https://doi.org/10.1002/ldr.2887

820    Voltz, M., Arrouays, D., Bispo, A., Lagacherie, P., Laroche, B., Lemercier, B., Richier-de-

821        Forges, A., Sauter, J., Schnebelen, N. (2020). Disseminating Digital Soil Mapping in

822        national soil mapping programmes: a prospective analysis in France. Submitted in

823        Geoderma Regional.

824    Wadoux, A. M. J. C., Brus, D. J., & Heuvelink, G. B. M. (2019). Sampling design

825        optimization for soil mapping with random forest. *Geoderma*.

826        https://doi.org/10.1016/j.geoderma.2019.113913

827    Walvoort, D. J. J., Brus, D. J., & de Gruijter, J. J. (2010). An R package for spatial coverage

828        sampling and random sampling from compact geographical strata by k-means.

829        *Computers and Geosciences*, *36*(10), 1261–1267.

830        https://doi.org/10.1016/j.cageo.2010.04.005

831    Walvoort, D. J. J., Brus, D. J., & de Gruijter, J. J. (2018). Spatial Coverage Sampling and

832        Random Sampling from Compact Geographical Strata. R package version 0.3-8.

833        https://CRAN.R-project.org/package=spcosa

834    Wright, M. N., & Ziegler, A. (2017). Ranger: A fast implementation of random forests for

835       high dimensional data in C++ and R. *Journal of Statistical Software*, *77*(1), 1–17.

836       https://doi.org/10.18637/jss.v077.i01

837    Yang, C. S., & Yang, Y. H. (2017). Improved local binary pattern for real scene optical

838       character       recognition.       *Pattern*       *Recognition*       *Letters*.

839       https://doi.org/10.1016/j.patrec.2017.08.005

840