

# 1 **Validation of uncertainty predictions in digital soil mapping**

2 Jonas Schmidinger<sup>1,\*</sup>, Gerard B.M. Heuvelink<sup>1,2</sup>

3 <sup>1</sup> Wageningen University and Research, Soil Geography and Landscape Group, Wageningen, the  
4 Netherlands

5 <sup>2</sup> ISRIC-World Soil Information, Wageningen, the Netherlands

6  
7 \* Corresponding author. Email: Jonas.Schmidinger@outlook.com; Address: Soil Geography and  
8 Landscape group, Wageningen University, Droevendaalsesteeg 3, 6708BP Wageningen, the  
9 Netherlands

10

## 11 **Abstract**

12 It is quite common in digital soil mapping (DSM) to quantify the uncertainty of issued predictions, that  
13 is to make probabilistic predictions. Yet, little attention has been paid to its validation. Probabilistic  
14 predictions are only of value for end users if they are reliable and ideally also sharp. The prediction  
15 interval coverage probability (PICP) is currently used in DSM to validate the reliability of prediction  
16 intervals but it is ignorant of a potential one-sided bias of its boundaries. Therefore, we propose to extend  
17 the current validation procedure with metrics used in the broader probabilistic literature. These metrics  
18 not only evaluate probabilistic predictions in prediction interval format but also quantiles or full  
19 conditional probability distributions. We suggest the quantile coverage probability (QCP) and  
20 probability integral transform (PIT) histogram as alternative to PICP and proper scoring rules for relative  
21 comparisons of competing probabilistic models. As scoring rules, we present the interval score (IS) and  
22 the continuous ranked probability score (CRPS), which can be decomposed into a reliability part (RELI).  
23 We illustrated the use of these metrics in a case study using soil pH and soil organic carbon from the  
24 LUCAS-soil database. Thereby, probabilistic predictions of five different models were compared: a  
25 reference null model (NM), quantile regression forest (QRF), quantile regression post-processing of a  
26 random forest (QRPP RF), kriging with external drift (KED) and quantile regression neural network  
27 (QRNN). For KED and QRNN, one-sided bias was found. This was not apparent from PICP but was

28 shown by use of the PIT histogram and QCP. RELI summarized the trends found in QCP, PICP and PIT  
29 histograms to one numerical value. CRPS and IS were especially harsh to outliers and low sharpness.  
30 According to CRPS and IS, the best probabilistic predictions were obtained by QRF and QRPP RF and  
31 the worst by NM.

32 **Keywords** · Validation · Digital soil mapping · Uncertainty · Machine learning · Proper scoring rules  
33 · Quantile regression

34

## 35 **1 Introduction**

36 Soils are of great importance to humankind since they provide various ecosystem services that contribute  
37 to food production, climate mitigation and air- and water quality (Keesstra et al., 2016). In order to  
38 maintain these services, soil as a resource has to be adequately managed and protected. This requires  
39 quantitative soil information at high spatial resolution, prompting the increasing popularity of digital  
40 soil mapping (DSM) (Chen et al., 2022). DSM creates soil maps through statistical inferences from a  
41 prediction model, using exhaustively accessible environmental covariates as predictors and soil sample  
42 data for model training (McBratney et al., 2003). Unavoidably, these predictions and thus the generated  
43 soil maps are not error-free. Map error originates from a variety of sources but most importantly it comes  
44 from the inability of the covariates to explain all soil spatial variation (Nelson et al., 2011). Other sources  
45 of error include the limited ability of a model to exploit all information provided by the covariates, a  
46 too-small training sample size, and measurement errors in the training data.

47 Estimation of the overall error can be done using a design-based approach, with independent test data  
48 obtained from probability sampling (Brus et al., 2011). Using this approach, map predictions are  
49 compared to the independent observations. The model performance, i.e. map accuracy, can then be  
50 quantified by well-established validation metrics such as the mean error (ME), root-mean-squared error  
51 (RMSE), and Nash-Sutcliffe model efficiency coefficient (MEC) (Piikki et al., 2021).

52 End users might not only be interested in the overall map accuracy but might require information  
53 about the accuracy at each and every location in the mapped study area. In such case, a design-based  
54 statistical inference is not suitable because it only provides summary measures of the map accuracy.  
55 However, with a model-based approach (Heuvelink, 2018) location-specific information about the  
56 prediction accuracy can be derived through the use of a probabilistic prediction model. A probabilistic  
57 prediction model goes beyond point prediction and estimates the entire conditional probability  
58 distribution of a soil property of interest, either directly or from a large set of conditional quantiles  
59 (Lauret et al., 2019). We refer to them as predictive distributions. They are generated for every location  
60 in the area of interest, in which the mean of the predictive distribution is typically used as a point  
61 prediction. The predictive distribution defines the probability of obtaining a large or small prediction  
62 error. A narrow, also called sharp, predictive distribution indicates that the point prediction is likely

63 close to the true value. In such case we are confident about the obtained point prediction and do not  
64 expect to have a large prediction error. With a wider predictive distribution, it cannot be ruled out that  
65 the prediction error is large, meaning that we are more uncertain if the true value is close to the predicted  
66 value. In DSM, we usually refer to this general concept as uncertainty (Heuvelink, 2018). In the  
67 following, we will use the more general term ‘probabilistic prediction’ as a synonym to uncertainty  
68 prediction.

69 While uncertainty is completely characterized by a predictive distribution, often it is summarized  
70 through a prediction interval (PI) for a more intuitive and practical interpretation. PI indicates a range  
71 in which the true value is expected to be found, given an assigned probability. Usually, a 90% prediction  
72 interval (PI) is used in DSM (Chen et al., 2022). For instance, if the 90% PI of the pH of the soil at some  
73 location is given by [5.3, 7.1], we claim that there is a chance of 90% that the true soil pH is between  
74 the lower limit of 5.3 and the upper limit of 7.1.

75 Indicating the uncertainty, usually in the form of a PI, has multiple advantages, such as: (i) preventing  
76 end users from getting a wrong sense about the accuracy of a soil map, thus allowing them to decide if  
77 the quality of the map is sufficient for the intended purpose (Heuvelink, 2018); (ii) allowing uncertainty  
78 propagation if the soil map is further used as input in other simulations or models (Heuvelink, 1998);  
79 and (iii) enabling to consider uncertainty in decision making (Breure et al., 2022; Lark et al., 2022).  
80 Because of these reasons, it is strongly encouraged to deliver the underlying uncertainty next to the  
81 actual predicted soil attributes.

82 Traditionally, much attention has been paid in DSM to quantify uncertainty of kriging models  
83 (Goovaerts, 2001), such as ordinary kriging or kriging with external drift (Webster and Oliver, 2007).  
84 However, machine learning algorithms that are able to predict conditional quantiles are getting  
85 increasingly popular in DSM (Kasraei et al., 2021; Lagacherie et al., 2019; Vaysse and Lagacherie,  
86 2017). Two examples of such techniques are quantile regression forest (Meinshausen, 2006) and  
87 quantile regression neural network (Cannon, 2011), which are probabilistic adaptations of a random forest  
88 and an artificial neural network, respectively. Recently, Kasraei et al. (2021) introduced quantile  
89 regression post-processing, which makes use of a quantile regression (Koenker and Hallock, 2001)  
90 implemented in a two-step algorithm.

91 Even though often disregarded, probabilistic predictions must be validated too since a poor  
92 uncertainty map could encourage harmful decisions if used in practice. The validation of probabilistic  
93 predictions is more complicated than the validation of point predictions as the former are characterized  
94 by probabilities and can occur in different forms such as PIs, quantiles, predictive distribution functions,  
95 or a mixture of them. Two general attributes that are usually evaluated for probabilistic predictions are  
96 reliability and sharpness (Gneiting and Raftery, 2007). Reliability, also known as probabilistic  
97 calibration, is a measure of the consistency between the predicted probability and the empirical  
98 frequency of independent test data. A probabilistic prediction should also be informative, which can be  
99 expressed by its sharpness. Sharpness is measured by the concentration of probabilistic information.  
100 Hence, high sharpness is characterized by a narrow PI or predictive distribution. Gneiting and Raftery  
101 (2007) define the goal for probabilistic prediction to “maximize the sharpness of the predictive  
102 distributions subject to reliability”. In DSM we usually measure sharpness by the prediction interval  
103 width and validate reliability with the prediction interval coverage probability (PICP) (Goovaerts, 2001).  
104 PICP evaluates if the probability assigned to the PIs is equal to the frequency of empirical test data  
105 within the PIs. Various studies compared the reliability of probabilistic prediction models frequently  
106 used in DSM based on PICP (e.g. Kasraei et al., 2021; Szatmári and Pásztor, 2019; Vaysse and  
107 Lagacherie, 2017). Vaysse and Lagacherie (2017) and Szatmári and Pásztor (2019) reported suboptimal  
108 PICPs for kriging as they were outperformed by quantile regression forest. Contrarily, Kasraei et al.  
109 (2021) had mixed results for quantile regression forest and obtained more stable results with quantile  
110 regression post processing combined to various machine learning models.

111 One reason for these ambiguous results may lie in the use of PICP as a validation metric. Pinson and  
112 Tastu (2014) pointed out that PICP is not an optimal metric to measure reliability, since it cannot account  
113 for one-sided bias. Therefore, it is of interest to expand the current validation procedure in DSM with  
114 more validation metrics for probabilistic predictions. In other academic fields, in which probabilistic  
115 predictions have a longer tradition, its validation is more comprehensive (Bracher et al., 2021; Lauret et  
116 al., 2019; Pinson et al., 2007; Zhang et al., 2014). These fields naturally include a broader set of  
117 validation metrics such as proper scoring rules (Gneiting and Raftery, 2007), probability integral  
118 transform histograms and the quantile coverage probability (Lauret et al., 2019).

119 The overall objective of this study is to introduce well-established concepts for the validation of  
120 probabilistic predictions from other academic fields to DSM. In a case study, using Land Use and  
121 Coverage Area Frame Survey (LUCAS) data, their added value will be illustrated for several  
122 probabilistic prediction models relevant to DSM. The performance of these models will then be  
123 compared on the basis of the old- and the newly introduced metrics.

124

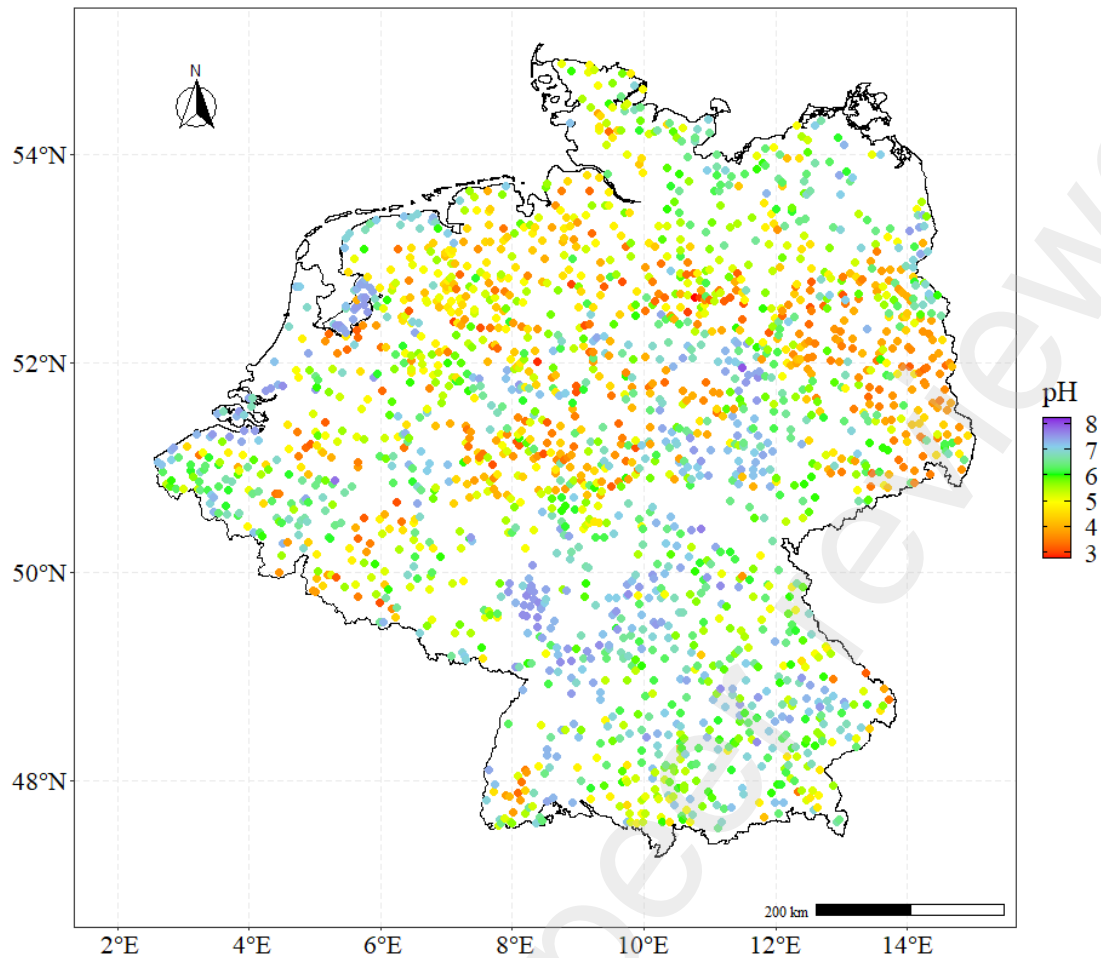
## 125 **2 Materials and Methods**

### 126 **2.1 Study area, soil data and covariates**

127 We used the soil pH and soil organic carbon (SOC) data from LUCAS-soil 2015 in the area of Germany  
128 and Benelux, which consisted of 2,018 data points. LUCAS-soil contains various soil attributes for all  
129 countries of the European Union. The large size, open license and consistent methodology makes  
130 LUCAS-soil attractive for testing new methods in DSM. Each soil sample of LUCAS-soil is a composite  
131 sample of five topsoil subsamples (0 – 20 cm) (Orgiazzi et al., 2018). pH was measured according to  
132 ISO 10390:1994 (ISO, 1994) with a glass electrode both in water and in a calcium-chloride solution. In  
133 this study we only used pH measured in calcium-chloride. SOC was obtained according to ISO  
134 10694:1995 (ISO, 1995) through the determination of the sample weight loss after dry combustion and  
135 removal of carbonates. Further, we log-transformed SOC to  $\log(\text{SOC})$  to remove skewness. The study  
136 area and sampling locations and values of pH can be found in Fig. 1. Those of  $\log(\text{SOC})$  are given in  
137 Fig. A1 in the Supplementary Information (SI).

138 We used the pre-processed covariates from Poggio et al. (2021), which consisted of various  
139 environmental factors important in the context of soil formation. These were, among others, vegetation  
140 indices, climate variables, land cover and lithology. Additionally, we addressed multicollinearity by  
141 randomly dropping one covariate of each covariate pair that had a Pearson correlation bigger than 0.9  
142 and eliminated covariates with near-zero variance. This led to a list of 89 covariates. In a few places, the  
143 covariate space was incomplete. We excluded two soil samples that fell in such areas, leading to a total  
144 of  $N = 2,016$  soil samples used in this study.

145



146

147 Fig. 1. LUCAS-soil sampling sites in Germany and Benelux with color-coded values of soil pH.

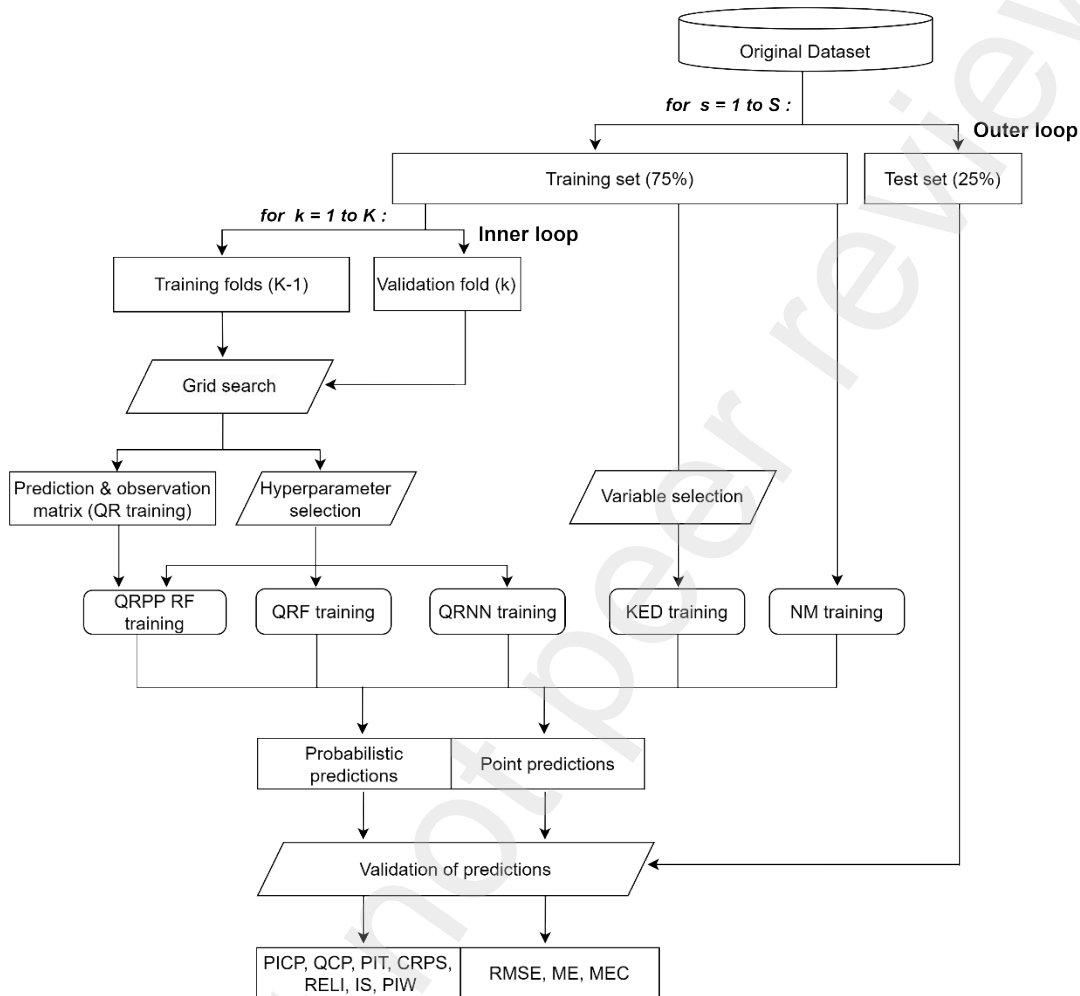
148

## 149 2.2 Study design

150 In an outer loop, the original dataset of size  $N$  was fully randomly split for  $S$  times ( $S = 25$ ) into a training  
 151 set (75%) for model fitting and a test set (25%) for the validation of model performance. Five prediction  
 152 models with probabilistic capabilities were trained (Section 2.3) to issue both point predictions and  
 153 probabilistic predictions (Section 2.4) for pH and log(SOC) at every sample site of the test set. Prior to  
 154 model training, for some of these models, a hyperparameter selection or a step-wise variable selection  
 155 grounded on the Akaike Information Criterion was implemented. The hyperparameter selection was  
 156 based on a grid search within an inner loop with  $K$ -fold cross-validation ( $K = 5$ ) of the training set. Note,  
 157 that the selection was executed on the basis of optimizing point prediction performances, not  
 158 probabilistic prediction performances. After the training, the test set was used to validate the point  
 159 predictions with standard validation metrics (Section 2.5) and the probabilistic predictions with the PICP

160 and newly proposed validation metrics (Section 2.6). Finally, the performances obtained from the outer  
 161 loop were stored and aggregated over the  $S$  repetitions. The whole study design is conceptualized in  
 162 Fig. 2.

163  
 164



165  
 166 Fig. 2. Conceptualization of the study design.

167

## 168 2.3 Probabilistic Prediction models

### 169 2.3.1 Null Model

170 The null model (NM) uses the mean of the training set as a point prediction and the empirical cumulative  
 171 distribution function (CDF) of the training set as a predictive CDF for all prediction points. Note that  
 172 this implies that both the point prediction and the probabilistic prediction are spatially invariant. NM  
 173 acts as a reference compared to the other models.



### 174 2.3.2 Quantile Regression Forest

175 Quantile regression forest (QRF) (Meinshausen, 2006) is a probabilistic adaption of the random forest  
176 (RF) algorithm (Breiman, 2001). QRF and RF make use of an ensemble of decision trees. Each single  
177 decision tree of the ensemble is grown with a recursive partitioning of the feature space on an individual  
178 bootstrapped training dataset, in which different nodes are created. However, only a random subset of  
179 covariates is used in the partitioning of each node. The information given at a prediction site is then run  
180 through each decision tree to obtain the corresponding terminal nodes, also known as leaves. The RF  
181 point prediction is then the weighted mean of the training observations stored in the corresponding  
182 leaves of every tree. QRF makes use of the fact that the RF prediction is a linear combination of the  
183 training data. It uses the RF weights and indicator transforms of the training data to estimate the CDF at  
184 multiple thresholds, from which the quantile prediction is inferred.

185 QRF was implemented in the statistical language R (R Core Team, 2023) via the *quantregForest*  
186 R-package (Meinshausen, 2017). The maximum node size (*nodesize*) and the number of randomly  
187 selected covariates in the partitioning of decision trees (*mtry*) were selected on the basis of a grid-search  
188 parameter tuning. For the number of trees fitted in the ensemble (*ntree*), we used the default of 500.

189

### 190 2.3.3 Quantile regression post-processing with a random forest

191 Kasraei et al. (2021) introduced quantile regression post-processing (QRPP) to DSM, which originates  
192 from the field of hydrology. It makes use of linear quantile regression (QR) (Koenker and Hallock,  
193 2001), which is comparable to standard linear regression but with the difference that it predicts  
194 conditional quantiles instead of a conditional mean. For that, the quantile loss function, also known as  
195 pinball loss function, is minimized in the training process. In QRPP, a QR is fitted on the relationship  
196 between point prediction values obtained by a model and observed values. Therefore, the fundamental  
197 difference in comparison to other probabilistic prediction models is that the actual probabilistic  
198 prediction is not embedded within the model algorithm, making it a two-step procedure. Thus, it has a  
199 flexible usage and can be combined with any point prediction model. In this study, we combined QRPP  
200 with a RF model (QRPP RF).

201 RF was modeled with the *randomForest* R-package (Liaw and Wiener, 2022). The parameters for RF  
202 were selected with the same parameter tuning procedure as in QRF. QR was implemented through the  
203 *quantreg* R-package (Koenker, 2022).

204

#### 205 **2.3.4 Kriging with external drift**

206 Kriging with external drift (KED) (Webster and Oliver, 2007) is a hybrid interpolation technique, based  
207 on a geostatistical model that represents the dependent variable as the sum of a non-constant trend, i.e.  
208 external drift, and a zero-mean stochastic residual. The external drift is usually modeled with multiple  
209 linear regression, while the stochastic residual is interpolated with kriging. The trend parameters are  
210 estimated with generalized least squares, in order to account for autocorrelation of the residuals. KED  
211 prediction error variances, which are used to generate the predictive distribution, can then be calculated  
212 from both the kriging variance and the estimation variance of the trend parameters (Brus and Heuvelink,  
213 2007). Finally, it is assumed that the predictive distribution follows a normal distribution (Goovaerts,  
214 2001; Heuvelink, 2018).

215 For each training set, a stepwise-variable selection based on the Akaike Information Criterion was  
216 implemented. The variograms needed for kriging were fitted with the *automap* R-package (Hiemstra,  
217 2022), and KED was executed with the *gstat* R-package (Pebesma, 2022).

218

#### 219 **2.3.5 Quantile Regression Neural Network**

220 Quantile regression neural networks (QRNN) (Cannon, 2011) is a probabilistic adaption of an artificial  
221 neural network (ANN). QRNN has the classic multilayer perceptron architecture, which consists of  
222 multiple layers, including an input layer, at least one hidden layer and an output layer. Layers are  
223 composed of neurons that are connected to the neurons of the previous and following layer. These  
224 connections have an associated weight term and each neuron possesses a bias term, apart from neurons  
225 in the input layer. The neurons of the input layer supply covariate input data to the first hidden layer, in  
226 which in every neuron an output is computed from the weights, bias and a defined hidden layer transfer  
227 function which introduces non-linearity. The output then serves as input to the next layer. This continues  
228 until the output layer is reached, in which the conditional quantiles are computed with an output layer

229 transfer function. The numeric values of the biases and the weights are determined with a  
230 backpropagation algorithm. The loss function used in QRNN is a differentiable approximation of the  
231 quantile loss function originating from QR.

232 QRNN was modeled with the *qrnn* R-package (Cannon, 2019). By default, we used only one hidden  
233 layer and the identity function for the output layer transfer function. The number of neurons in the hidden  
234 layer (*n.hidden*), the hidden layer transfer function (*Th* and *Th.prime*) and the weight decay (*penalty*)  
235 were determined using a hyperparameter selection.

236

## 237 **2.4 Estimation of the predictive distribution**

238 A predictive CDF can be generated with either a parametric or a nonparametric approach (Lauret et al.,  
239 2019). In a parametric approach, assumptions about the shape of the distribution are made beforehand.  
240 Hence, in order to generate a predictive distribution, one first has to determine the desired distribution  
241 from a parametric family (e.g. Gaussian, exponential, Weibull) and next estimate the parameters of the  
242 predictive CDF. Kriging and hence also KED uses a parametric approach, in which a normal distribution  
243 is typically assumed (Section 2.3.4).

244 With QRNN, QRPP RF and QRF, a non-parametric approach is used based on predicted quantiles.  
245 Nonparametric distributions do not have to adhere to restrictive assumptions imposed by the preselected  
246 parametric family (Lauret et al., 2019). However, these methods only predict a finite number of  
247 quantiles. Hence, if it is desired to generate a predictive CDF from quantile regression methods, one has  
248 to approximate the CDF from these quantiles (Zamo and Naveau, 2018). We did this using a quantile  
249 set consisting of 199 quantiles at the 0.5% to 99.5% percentile. Additionally, QRF and QRNN did not  
250 provide point prediction values directly, so they were obtained by taking the mean of the quantile set.

251

## 252 **2.5 Point prediction validation metrics**

253 Next to probabilistic predictions, we also issued and validated point predictions, to show the  
254 performance of the models outside of the probabilistic context. The root mean square error (RMSE) is  
255 the most commonly used validation measure in DSM (Piikki et al., 2021) and indicates how much the  
256 predictions deviate from the observations:

257 
$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad (1)$$

258 where  $n$  is the size of the test set,  $y_i$  ( $i = 1, \dots, n$ ) are the predicted values and  $y_i$  the observed values of  
 259 the test data.

260 The mean error (ME) is a bias indicator. Other than RMSE, it can have both positive and negative  
 261 values (Piikki et al., 2021). A value close to zero indicates that the point predictions are free from bias.

262 
$$ME = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i). \quad (2)$$

263 The Nash–Sutcliffe model efficiency coefficient (MEC) (Nash and Sutcliffe, 1970) is a relative  
 264 error measure:

265 
$$MEC = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (3)$$

266 in which  $\bar{y}$  is the arithmetic mean of the test data. In case of perfect agreement between test observations  
 267 and predictions, the MEC is equal to one. The NM is expected to have a MEC close to zero. Note that  
 268 MEC will be negative for models that perform worse than the NM.

269

## 270 **2.6 Uncertainty validation metrics**

### 271 **2.6.1 Prediction interval width**

272 The prediction interval width (PIW) is a measure for the sharpness of a probabilistic prediction. PIW  
 273 indicates the width of a certain  $\tau \cdot 100$  per cent PI, for any value of  $\tau$  between 0 and 1:

274 
$$PIW(\tau) = \frac{1}{n} \sum_{i=1}^n (u_i - l_i), \quad (4)$$

275 where  $l_i$  is the lower bounding quantile and  $u_i$  the upper bounding quantile that together define a  $\tau \cdot 100$   
 276 per cent PI. Usually, central PIs are of interest, meaning that the probability mass below  $l_i$  and above  $u_i$   
 277 are equal. Therefore,  $l_i$  and  $u_i$  are determined by the chosen  $\tau$  value:

278 
$$l_i = q_{(1-\tau)/2}^i, \quad u_i = q_{(1+\tau)/2}^i, \quad (5)$$

279 where  $q_{(1-\tau)/2}^i$  and  $q_{(1+\tau)/2}^i$  are the  $(1 - \tau)/2$  and  $(1 + \tau)/2$  quantiles of the predictive distribution of  $y_i$ .

280 Lower PIW values imply higher sharpness, i.e. lower uncertainty (Kasraei et al., 2021; Pinson et al.,  
 281 2007). Therefore, lower PIW values are preferred, given the constraints of reliability. The degree of

282 sharpness is also related to the point prediction performance. For example, when we have a small RMSE,  
283 then our probabilistic predictions will have high sharpness, given they are reliable. Although PIW  
284 formally is not a validation metric, because its value is independent of the test data, it should be included  
285 in the evaluation of probabilistic predictions.

286

## 287 **2.6.2 Prediction interval coverage probability**

288 To assess the reliability of PIs, PICP is commonly adopted in DSM (Piikki et al., 2021). Most analyses  
289 rely on PICP as a single reliability metric (Kasraei et al., 2021; Lagacherie et al., 2019; Szatmári and  
290 Pásztor, 2019; Vaysse and Lagacherie, 2017). The underlying idea is to evaluate what percentage of soil  
291 samples from the test set lies in the  $\tau \cdot 100$  per cent PIs:

$$292 \quad PICP(\tau) = \frac{1}{n} \sum_i^n \delta(l_i \leq y_i \leq u_i) \cdot 100, \quad (6)$$

293 where  $\delta$  is an indicator function, with a Boolean argument:

$$294 \quad \delta(t) = \begin{cases} 1 & \text{if } t \text{ is TRUE} \\ 0 & \text{else} \end{cases}. \quad (7)$$

295 In an ideal case,  $PICP(\tau)$  is equal to  $\tau \cdot 100$  per cent, e.g., for a 90% PI we desire a PICP of 90%.  
296 Multiple PICPs are usually calculated for different PI levels ( $\tau$  values). This then represents the  
297 reliability over the whole predictive distribution. A reliability plot allows a visual evaluation of the  
298 reliability by plotting the PICP against the associated PI level. It is then desired that the points are on or  
299 close to the 1:1 line. Values below and above the 1:1 line indicate over-pessimistic or over-optimistic  
300 PIs, respectively. Note, that in DSM the reliability plot was introduced in Goovaerts (2001) and referred  
301 to as ‘accuracy plot’. However, ‘reliability plot’ is a more generally accepted term within other academic  
302 fields (Lauret et al., 2019). Schematic examples of PICP reliability plots are shown in Fig. 3.

303 One clear advantage of PICP is that its value has an intuitive interpretation. Nonetheless, PICP has  
304 also a disadvantage, which has not yet been addressed in DSM. As demonstrated in Pinson and Tastu  
305 (2014), PICP does not account for a systematic one-sided bias. This occurs when the quantile predictions  
306 of the lower and upper boundary of a PI are both either positively or negatively shifted. For example,  
307 for a central 90% PI we expect that 5% of the test data are below the lower boundary and 5% above the  
308 upper boundary. However, if we have a one-sided bias, in which both boundaries are shifted by +4%,

309 we would observe that 9% of the test data are below the lower boundary and 1% above the upper  
310 boundary. In this case, we would still obtain a PICP of 90%, yet it ignores the asymmetrical coverage.  
311 The effect of one-sided bias on PICP is conceptualized in Fig. 3.

312

### 313 2.6.3 Quantile coverage probability

314 A simple solution that overcomes the shortcoming of PICP but otherwise has similar properties is the  
315 use of the quantile coverage probability (QCP). It has the same underlying logic as PICP but it evaluates  
316 single quantile predictions. It computes which fraction of the test set is below a quantile:

$$317 \quad QCP(\tau) = \frac{1}{n} \sum_i^n \delta(y_i \leq q_i^{\tau}) \cdot 100. \quad (8)$$

318 This has the advantage that a potential bias will not be hidden. In some studies only the coverage based  
319 on quantiles is computed and the PICP is left out entirely (Lauret et al., 2019; Vasseur and Aznarte,  
320 2021). Examples of QCP reliability plots are also shown in Fig. 3.

321

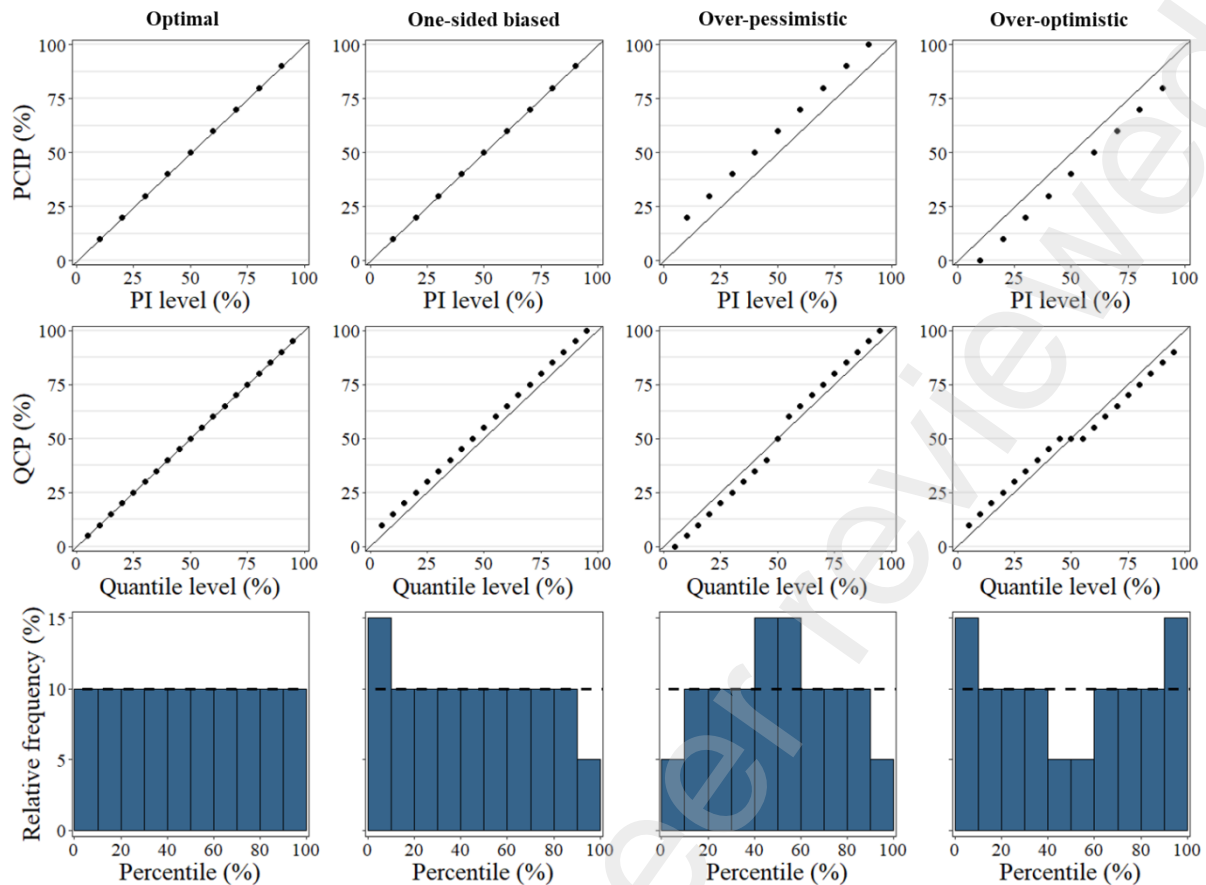
### 322 2.6.4 Probability integral transform

323 The probability integral transform (PIT) histogram (Gneiting et al., 2007) is another visual tool for the  
324 assessment of reliability. We found one example of PIT usage within DSM in Nussbaum et al. (2014).  
325 It evaluates if the test observations  $y_i$  cover the whole range of the predictive CDFs. It starts by  
326 computing the percentiles  $P_i$  associated to the test observations  $y_i$  in the predictive distributions  $F_i$ :

$$327 \quad P_i = F_i(y_i). \quad (9)$$

328 When plotting the obtained  $P_i$  as a histogram, this should ideally be a uniform distribution. An uneven  
329 distribution indicates that some parts of  $F_i$  are disproportionately often or sparsely covered. A sloped,  
330 convex or concave shape of the histogram indicates one-sided biased, over-optimistic or over-  
331 pessimistic probabilistic predictions, respectively. These cases are schematically exemplified in Fig. 3.

332



333

334 Fig. 3. Reliability plots of PICP and QCP and PIT histograms for four hypothetical scenarios: an optimal, one-  
 335 sided biased, over-pessimistic and over-optimistic scenario.

336

### 337 2.6.5 Scoring rules

338 The so far presented validation metrics (PICP, QCP and PIT) indicate the reliability of a single  
 339 probabilistic prediction model without a required reference or comparison. However, one may also be  
 340 interested in a relative comparison of competing probabilistic prediction models (Gneiting and Raftery,  
 341 2007). For this purpose, scoring rules can be used. Scoring rules are measures that evaluate the quality  
 342 of a probabilistic prediction model and return a numeric score value. Based on the obtained score values  
 343 the performance of competing probabilistic prediction models can be compared and ranked. Further, it  
 344 is desired that scoring rules are proper. The term proper refers to the concept that there is no incentive  
 345 to report any predictive distribution other than the one of true belief of the model (Gneiting and Raftery,  
 346 2007). In the next subsections we suggest two proper scoring rules. Both are negatively oriented, so that  
 347 smaller values indicate a better score. Further, they consider both sharpness and reliability.

### 348 2.6.6 Interval Score

349 The Interval Score (IS) is a scoring rule that evaluates PIs (Bracher et al., 2021; Gneiting and Raftery,  
350 2007). Therefore, IS depends on the chosen  $\tau$  value. It is calculated for each test observation  $y_i$  and  
351 subsequently averaged over the whole test set, to get one final score value:

$$352 \quad IS(\tau) = \frac{1}{n} \sum_{i=1}^n (u_i - l_i) + \frac{2}{1-\tau} \cdot \frac{1}{n} \sum_{i=1}^n (l_i - y_i) \cdot \delta(y_i < l_i) + \frac{2}{1-\tau} \cdot \frac{1}{n} \sum_{i=1}^n (y_i - u_i) \cdot \delta(y_i > u_i). \quad (10)$$

353 The first term in Eq. 10 is the average width of the PIs, meaning that sharper PIs receive a lower penalty.  
354 The other two terms only consider test observations that are below  $l_i$  and above  $u_i$ . These observations  
355 get a penalty with the distance to the boundaries of the PI. Hence, unlike PICP and QCP, IS also  
356 penalizes how far outside a PI the observations are. This may seem contradictory at first because unless  
357 we issue a 100% PI, it is desired that a fraction  $1 - \tau$  of the observations are outside the PI. Yet, as  
358 already mentioned, IS additionally considers the width of a PI in its scoring. Therefore, a wider PI may  
359 lead to fewer observations outside its boundaries but it simultaneously gets punished for its low  
360 sharpness.

361 To our knowledge, IS was not yet used in DSM. The fact that IS is a scoring rule that evaluates  
362 probabilistic predictions in PI format can be an advantage (Bracher et al., 2021) since probabilistic  
363 predictions in DSM are usually issued as PIs. Fig. 4 shows a schematic visualization of how single IS(  
364  $\tau$ ) values are calculated.

365

### 366 2.6.7 Continuous Ranked Probability Score

367 The Continuous Ranked Probability Score (CRPS) is a widely used scoring rule for continuous variables  
368 (Lauret et al., 2019). We know of two instances where CRPS was used in DSM (Caubet et al., 2019;  
369 Nussbaum et al., 2014). Other than IS, CRPS directly evaluates the whole predictive CDF. The  
370 calculation of CRPS is comparable to that of point prediction metrics like the mean squared error. It is  
371 defined as the integral of the squared difference between the predictive CDFs  $F_i$  and empirical CDFs  
372 from observed test data. The latter can be also interpreted as a Heaviside step function  $H$ , since its CDFs  
373 are generated from single test samples  $y_i$ :

$$374 \quad CRPS = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} (F_i(y) - H(y - y_i))^2 dy, \quad (11)$$



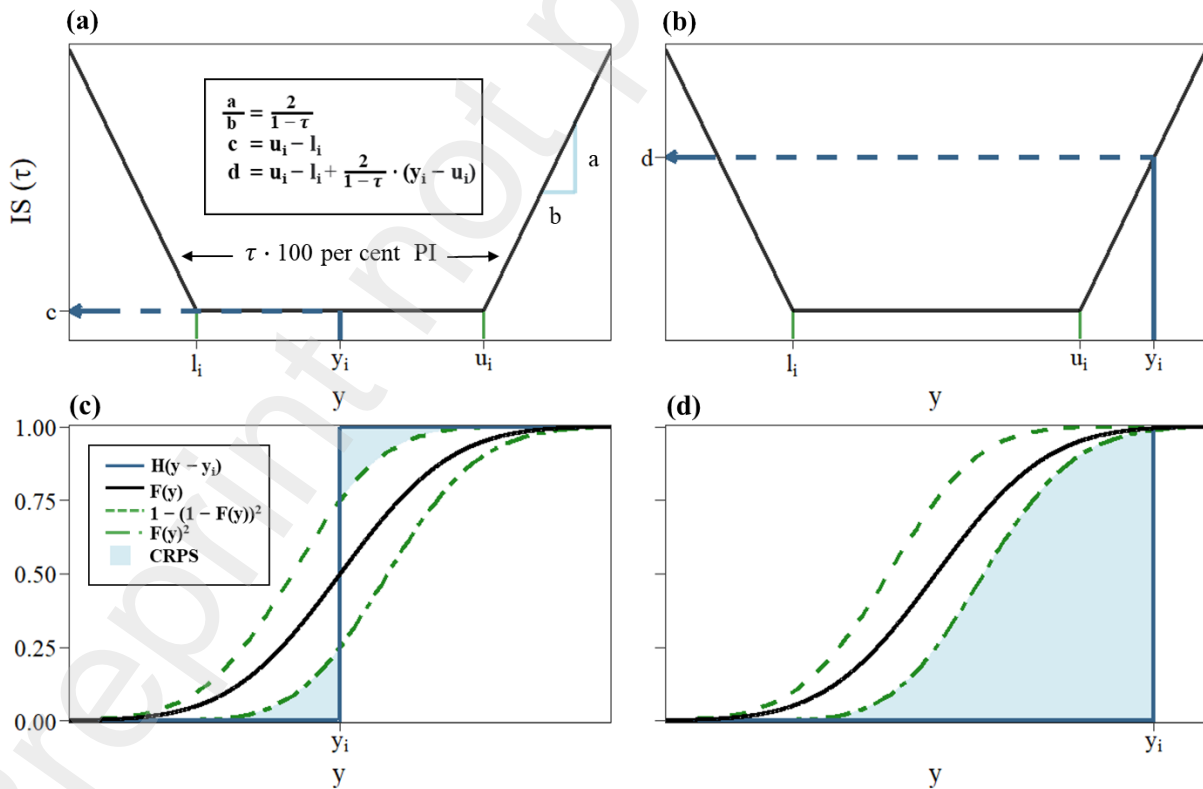
375 where

376 
$$H(t) = \begin{cases} 1 & \text{if } t \geq 0 \\ 0 & \text{else} \end{cases} \quad (12)$$

377 A schematic visualization of how a single CRPS value is calculated is given in Fig. 4. Additionally, the  
 378 median CRPS can be evaluated to reduce the influence of outliers in the scoring.

379 For probabilistic models that return quantiles, there is no continuous predictive CDF as required for  
 380 Eq. 11. Instead, a step function is generated from a quantile set which approximates the continuous CDF  
 381 (Section 2.4). In this context, the equations from Hersbach (2000) can be applied, see SI. To calculate  
 382 CRPS, *crpsDecomposition* from the *verification* R-package (NCAR - Research Applications  
 383 Laboratory, 2015) was used. This package applies an equation from Hersbach (2000), i.e. Eq. A1 in SI.  
 384 We also used it in the case of KED, even though KED provides a continuous predictive CDF. However,  
 385 in case of KED, Eq. A1 functions as numerical integration and the approximation error in comparison  
 386 to Eq. 11 will be small due to the large number of quantiles used in the approximation.

387



388  
 389 Fig. 4. Schematic representation of two examples of how individual values of  $IS(\tau)$  (a-b) and CRPS (c-d) are  
 390 calculated. The examples on the left show a case in which the mean of the predictive distribution is the same as

391 the observed value; the examples on the right show a predictive distribution where the observed value is at the  
392 extreme of the distribution. Note, that the area of CRPS is squared, hence the added  $1 - (1 - F(y))^2$  and  $F(y)^2$ . This  
393 figure was inspired by illustrations in Bracher et al. (2021).

394

### 395 **2.6.8 Reliability decomposition**

396 CRPS can be decomposed into different parts (Hersbach, 2000). In this study, we only introduce and  
397 use the reliability part (RELI). In RELI, the mean coverage of the quantiles used to approximate the  
398 predictive CDFs are evaluated. Therefore, it is closely related to QCP and PIT but it returns a single  
399 numerical value. It further considers the distance between the quantiles in its weighting. More technical  
400 information about the computation of RELI is given in SI. As for CRPS, RELI was computed by the  
401 function *crpsDecomposition* from the *verification* R-package.

402

## 403 **3 Results**

404 In the following sections, only figures for pH are shown. Figures for log(SOC) can be found in SI. We  
405 refer to both soil properties in the text but prioritize pH.

406

### 407 **3.1 Point prediction performance**

408 According to MEC and RMSE presented in Table 1, the point prediction performances of pH from QRF,  
409 QRPP RF, KED and QRNN were very similar. KED had the single best RMSE of 0.81 but QRPP RF  
410 and QRF were close with an RMSE of 0.82 and QRNN with 0.83. As expected, by far the worst point  
411 predictions were produced by NM, whose RMSE was with 1.26 around 35% bigger than that of the  
412 other models. QRF, QRPP RF and QRNN obtained negative ME values that deviated most from 0.  
413 Nevertheless, with a ME of  $-0.013$ ,  $-0.015$  and  $-0.017$  respectively, the differences to zero were  
414 very small compared to RMSE which indicates that systematic prediction errors were negligible.

415 For log(SOC), point prediction performances were considerably different (Table A1 in SI). Here,  
416 QRF and QRPP RF generated the best results and slightly outperformed QRNN. KED was apart from

417 NM the worst model. Overall, with a maximum MEC of 0.43, log(SOC) point predictions were poorer  
418 compared to pH, where the highest MEC was 0.58.

419

420 Table 1. Point prediction performances for pH.

|      | NM      | QRF    | QRPP RF | KED    | QRNN   |
|------|---------|--------|---------|--------|--------|
| MEC  | -0.0016 | 0.57   | 0.58    | 0.58   | 0.56   |
| RMSE | 1.26    | 0.82   | 0.82    | 0.81   | 0.83   |
| ME   | -0.000  | -0.013 | -0.015  | -0.007 | -0.017 |

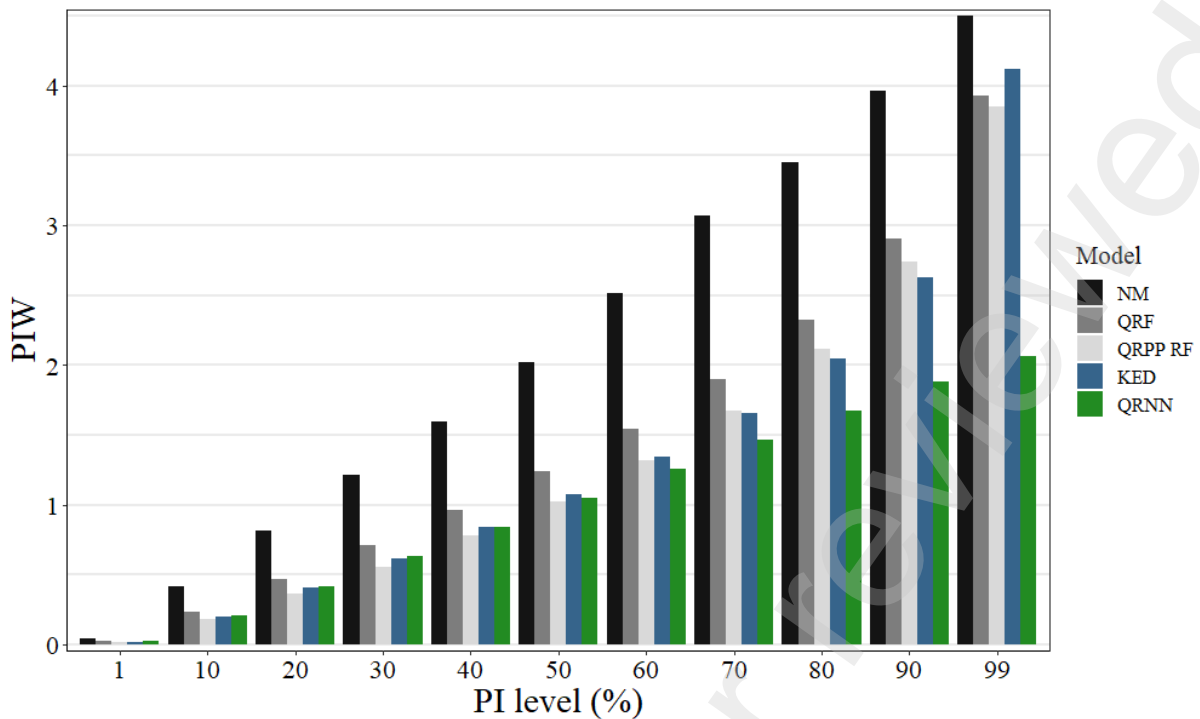
421

## 422 3.2 Probabilistic prediction performance

### 423 3.2.1 Prediction interval width

424 Fig. 5 illustrates PIW for pH at various PI levels to indicate the sharpness of the five prediction models:  
425 NM, QRF, QRPP RF, KED and QRNN. Throughout the predictive distribution, NM received the highest  
426 PIW values, which were on average 36% larger compared to the other models. KED and QRPP RF  
427 shared similar PIW values. For PI between 10% and 90%, PIW of QRF was slightly larger than those  
428 of KED and QRPP RF but at the extremes, their values had a similar level. Also, QRNN had similar  
429 PIW values as KED and QRPP RF up to the 60% PI. Thereafter, QRNN obtained PIW values that were  
430 much smaller compared to the other models. For example, for the 99% PI, the PIW of QRNN was about  
431 46% smaller than that of QRPP RF, which had the second lowest value. Almost the same patterns were  
432 found for log(SOC) (Fig. A2), except for KED. Here, KED was less sharp than QRPP RF and had PIW  
433 values that were similar to that of QRF.

434



435

436 Fig. 5. Sharpness diagram indicating the PIW for multiple PI levels of the predictive distribution for pH.

437

### 438 3.2.2 Prediction interval coverage probability & quantile coverage probability

439 Fig. 6 and Fig. 7 show PICP and QCP reliability plots of pH, respectively. The evaluated quantiles in

440 Fig. 7 correspond to the PIs in Fig. 6. The PICP and QCP values of NM, QRPP RF and KED were close

441 to the 1:1 line, which is an indicator of good reliability. PICP values of QRF were fairly over-pessimistic,

442 so that in some instances its PICP was around 5% above the 1:1 line. Yet, good concordance was found

443 at the extremes, more specifically above the 90% PI and below the 10% PI. This corresponds to good

444 agreement in terms of QCP at the extremes, i.e. below the 5% quantile and above the 95% quantile, and

445 in the center, i.e., between the 45% and 55% quantile. Between the 5% to 45% quantile, QCP of QRF

446 tended to be below and between the 55 to 95% above the 1:1 line. Since PICP combines the deviation

447 from its lower and upper boundaries the deviation was more visible for QRF in Fig. 6 compared to Fig.

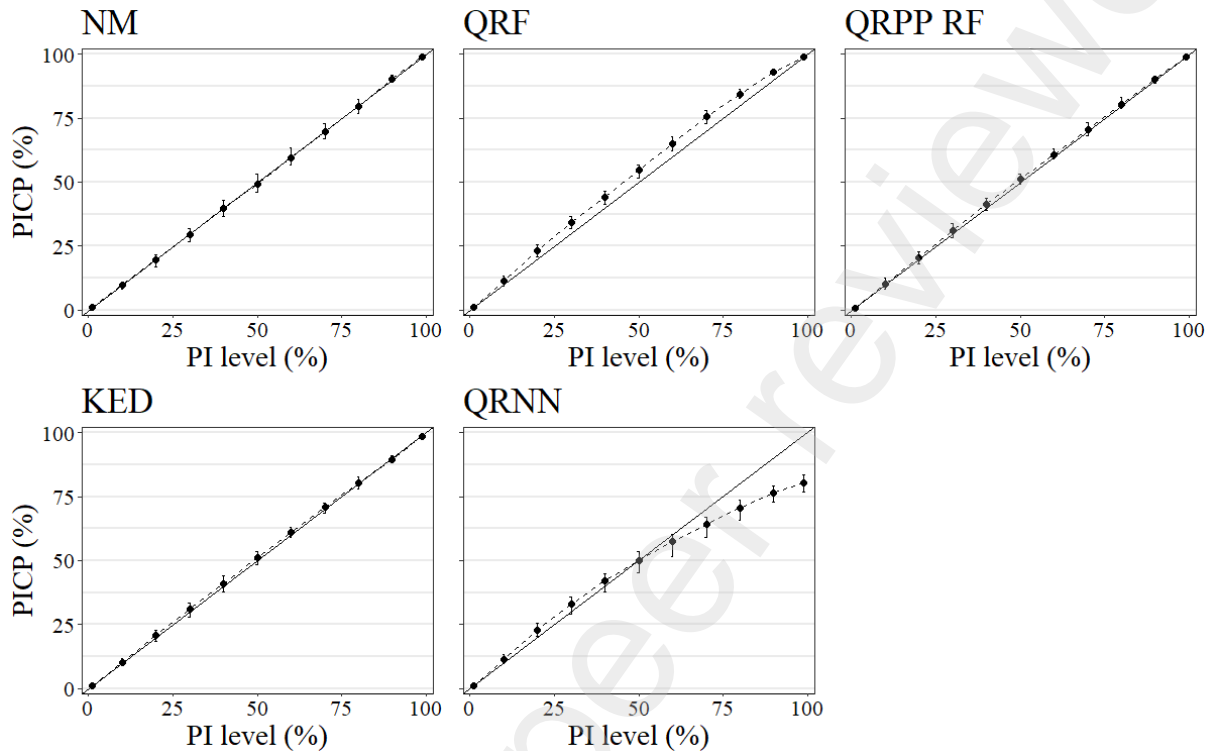
448 7. QRNN had a very different outcome compared to the other models. In regards to PICP it performed

449 well until the 60% PI. This corresponds to good agreement in terms of QCP between the 20% and 80%

450 quantile. However, QCP at the edges and PICP for the large PIs showed strong deviations from the 1:1

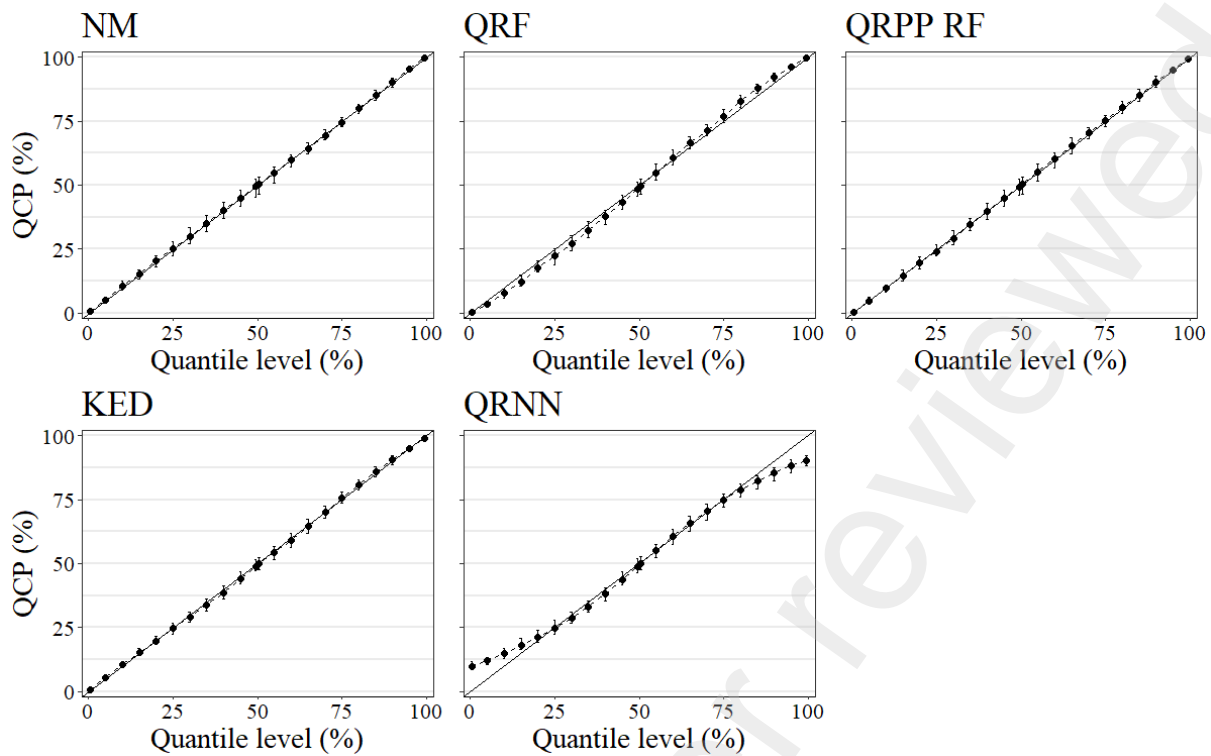
451 line. For instance, at the 99% PI, PICP was about 20% below the 1:1 line. There was no meaningful

452 one-sided bias in any model for probabilistic predictions of pH according to the reliability plot of QCP  
 453 (Fig. 7).  
 454



455  
 456 Fig. 6. Mean PICP reliability plots for pH. Error bars are retrieved from the 80% confidence interval of the 25  
 457 repetitions in the outer loop. The 1:1 black line indicates the desired outcome.

458  
 459 For log(SOC), the trends found for PICP (Fig. A3) and QCP (Fig. A4) of NM, QRPP RF and QRF were  
 460 similar to those of pH. In contrast, KED obtained over-pessimistic results that were comparable to QRF,  
 461 judging based on PICP. However, when looking at QCP, additionally one-sided bias was found for  
 462 KED, so that the deviation from the 1:1 line was more pronounced for KED than for QRF (Fig. A4). In  
 463 the range between the 30% to 85% quantile, QCP of KED was systematically above the 1:1 line. For  
 464 example, the PICP corresponding to the 10% PI was around 11%. Yet, for the corresponding 45% and  
 465 55% quantile, a QCP of around 49% and 60% was achieved, respectively. QRNN also showed one-  
 466 sided bias between the 30% and 70% quantile but it was less than for KED.



467

468 Fig. 7. Mean QCP reliability plots. Error bars are retrieved from the 80% confidence interval of the 25 repetitions  
 469 in the outer loop. The 1:1 black line indicates the desired outcome.

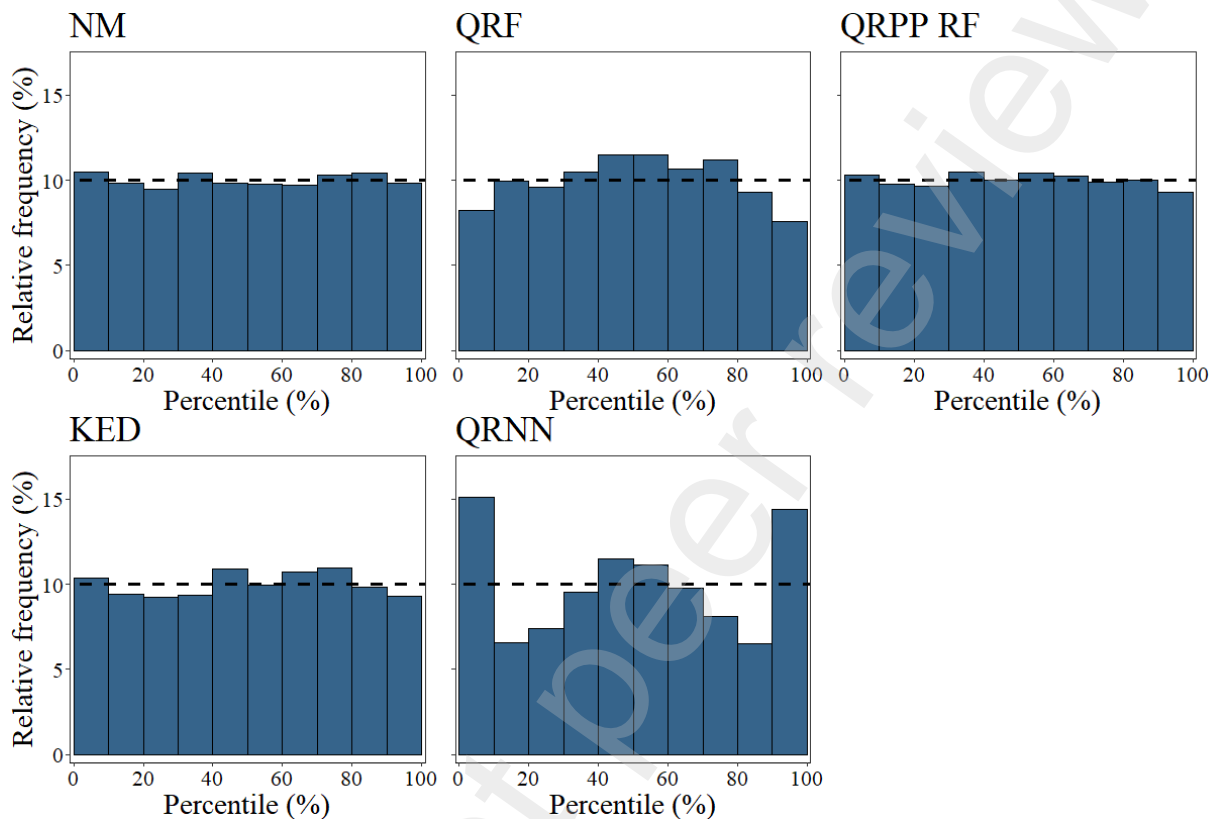
470

### 471 3.2.3 Probability integral transform

472 PIT histograms for pH are provided in Fig. 8. Since the bin width was 10%, for good reliability it is  
 473 desirable that the frequency in each bin is close to 10%, as indicated by the horizontal dashed lines in  
 474 Fig. 8. This was more or less achieved for NM, QRPP RF and KED. In these cases, the relative frequency  
 475 neither exceeded 11% nor fell below 9%. A concave histogram was obtained for QRF, meaning that  
 476 lower relative frequency values were obtained at the edges, i.e. the 0% to 10% and 90% to 100%  
 477 percentile range. Such a shape is characteristic for an over-pessimistic performance. A somewhat convex  
 478 distribution was achieved for QRNN. Here, the relative frequencies at the edges were around 15%. On  
 479 the other hand, the bins before and after the edges (10% to 30% and 70% to 90% percentiles), had very  
 480 diminished relative frequencies.

481 The PIT histograms obtained for log(SOC) (Fig. A5), were similar to those of pH, except for KED.  
 482 What stands out is that the PIT histogram of KED was the only histogram without any symmetrical  
 483 structure. Between the 10% and 90% percentile, frequencies started at a high level above the 10% line

484 but decreased steadily so that after the 60% percentile values were below 10%. Contrary to that trend,  
 485 the edges did not reflect the same behavior. For the 0% to 10% bin, a decreased frequency was obtained.  
 486 For the 90% to 100% percentile, the frequency was slightly above 10%.  
 487



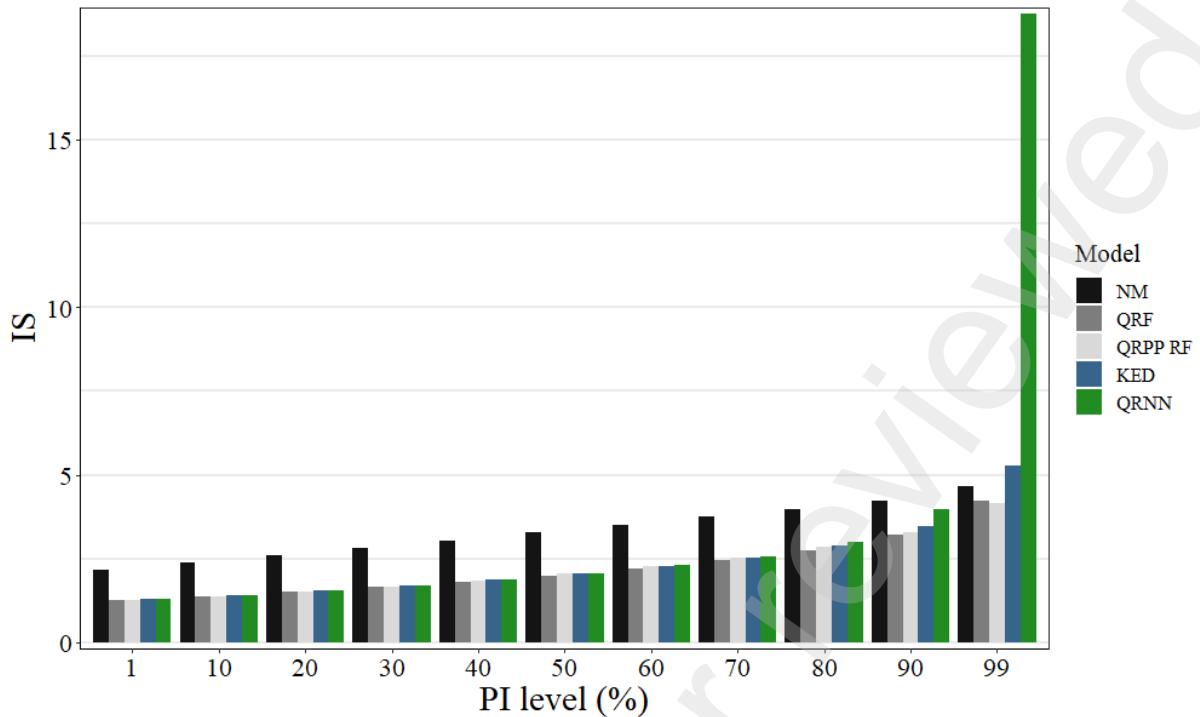
488 Fig. 8. PIT histograms of pH. The dashed line indicates the desired frequency for a flat and uniform PIT.

489

### 490 3.2.4 Interval score

491 Fig. 9 shows the IS over the whole predictive distribution of pH. Between the 1% and 90% PI, the same  
 492 ranking order can be found: NM obtained the largest and thus worst scores. The IS of the other models  
 493 were overall similar but the best scores were obtained for QRF, followed by QRPP RF, KED and lastly  
 494 QRNN. With increasing PI levels, all acquired IS values of the models rose. Yet, IS of QRNN and KED  
 495 rose disproportionately more. As a consequence, QRNN and KED surpassed NM at the 99% PI, at which  
 496 the score of QRNN overshadowed all other models. For log(SOC) (Fig. A6), more or less the same  
 497 trends were found as for pH, with the only difference that KED was slightly worse than QRNN up to  
 498 the 80% PI. Just thereafter, QRNN scored worse.

499



500

501 Fig. 9. IS diagram of pH indicating the IS values for multiple PI levels of the predictive distribution.

502

503 **3.2.5 Continuous ranked probability score & reliability decomposition**

504 RELI, CRPS and median CRPS for pH are given in Table 2. All single CRPS values (25 x 2,016)  
 505 obtained in the outer loop are shown as boxplots in Fig. 10. CRPS was the lowest and thus the best for  
 506 QRF and QRPP RF. The performance was followed by KED, QRNN and lastly NM. This indicates that  
 507 based on CRPS, QRF and QRPP RF were able to obtain better uncertainty predictions for pH than KED,  
 508 QRNN and NM. Yet, the median CRPS of QRF was worse than that of KED and QRNN (Table 2) but  
 509 the spread of KED and QRNN seemed to be larger (Fig. 10). The ranking order of RELI was very  
 510 different compared to the ranking of CRPS and median CRPS, as RELI only evaluates reliability.  
 511 According to RELI, the best reliability was achieved with KED, QRPP RF and NM. RELI of QRF was  
 512 considerably worse and QRNN scored by far the worst.

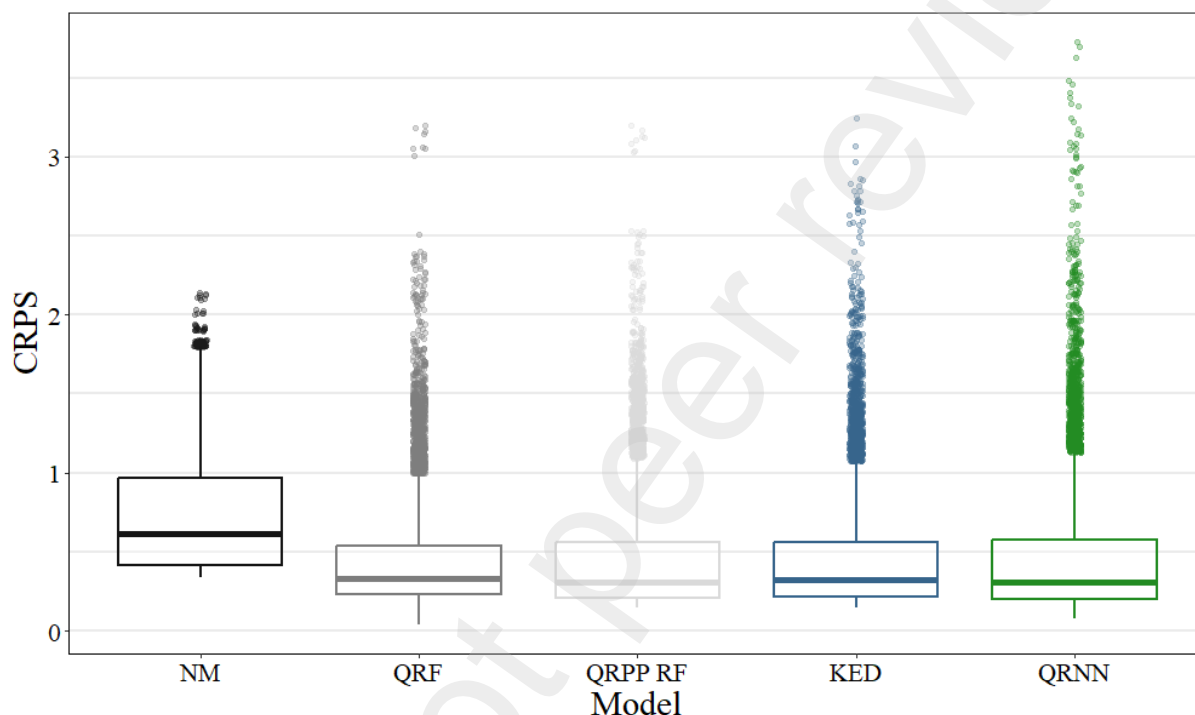
513 For log(SOC), the ranking order with respect to CRPS was slightly different compared to pH (Table  
 514 A2). Here, QRF and QRPP RF again scored the best but this time, QRNN scored better than KED. KED  
 515 was also worst in terms of median CRPS. The spread of CRPS was very similar between the models  
 516 (Fig. A7). Furthermore, KED achieved no longer the best reliability according to RELI. Better results  
 517 were obtained by NM, QRPP RF and QRF but KED still achieved a lower RELI value than QRNN.



518 Table 2. Scoring outcomes of CRPS, median CRPS and RELI for pH.

|             | NM     | QRF    | QRPP RF | KED    | QRNN   |
|-------------|--------|--------|---------|--------|--------|
| CRPS        | 0.72   | 0.44   | 0.44    | 0.46   | 0.47   |
| Median CRPS | 0.61   | 0.33   | 0.31    | 0.32   | 0.32   |
| RELI        | 0.0016 | 0.0071 | 0.0012  | 0.0010 | 0.0117 |

519



520

521 Fig. 10. Single CRPS values portrayed as boxplots for pH. Each boxplot was based on 25 x 2,016 values.

522

## 523 4 Discussion

### 524 4.1 Comparison of model performance

#### 525 4.1.1 Point prediction performance

526 NM used the mean value of the training set as point predictions. Therefore, it was not surprising that it  
 527 delivered the worst point prediction performance, both for pH and log(SOC). QRF and QRPP RF  
 528 provided stable point prediction results. This is in line with findings of other DSM studies, where good  
 529 performances of the associated RF were reported (Khaledian and Miller, 2020). In case of pH, QRF and  
 530 QRPP RF were on par with - or slightly worse - than KED. It thus appears that strong linear relationships

531 between pH and the available covariates were present next to some degree of spatial autocorrelation.  
532 Contrarily, non-linear relationships or interactions seemed to be present between log(SOC) and the  
533 covariates because QRPP RF, QRF and also QRNN outperformed KED. There was a small mismatch  
534 between QRPP RF and QRF point prediction results, even though both models in theory should achieve  
535 a similar outcome. This might be the result of different implementations in the R-packages or small  
536 errors with QRF at the conversion of quantile predictions to point predictions. However, the mismatch  
537 was very small and insignificant.

538

#### 539 **4.1.2 Probabilistic prediction performance**

540 It is no surprise that NM provided very reliable results since the uncertainty was modeled by the  
541 empirical distribution of the training set. Nonetheless, it suffered from very low sharpness as can be seen  
542 from PIW. Due to its inability to issue predictions with high sharpness, it was ranked last according to  
543 CRPS and IS.

544 QRF is a very commonly used method to estimate the uncertainty of an RF in DSM. Nonetheless,  
545 slightly over-optimistic probabilistic predictions were found in the center of the distribution. This seems  
546 to be a common outcome for QRF predictions, as similar problems were reported in Kasraei et al. (2021),  
547 Szatmári and Pásztor (2019) and Vaysse and Lagacherie (2017), or outside the DSM literature in for  
548 example Vasseur and Aznarte (2021). Yet, according to CRPS and IS, QRF achieved along with QRPP  
549 RF the best probabilistic predictions. The reasons for QRF performing well compared to other models  
550 despite having slightly over-pessimistic probabilistic predictions are discussed later in this section.

551 QRPP RF performed most consistently on all validation metrics for log(SOC) and pH. The good  
552 performance is in agreement with the reported PICP values in Kasraei et al. (2021). This outcome of  
553 QRPP RF may be surprising because it uses a simple structure, centered around a linear QR fitted on  
554 predicted and observed values. The method of QRPP strongly depends on the residual structure, so in  
555 future studies it could be further investigated in what way the residual structure influences the  
556 probabilistic predictions.

557 Probabilistic predictions of KED were inconsistent in terms of reliability. Good reliability and sharp  
558 distributions were found for pH, whereas for log(SOC) considerably worse reliability and unsharp

559 distributions were obtained that additionally had a one-sided bias. Suboptimal probabilistic predictions  
560 with different forms of kriging were earlier reported in Vaysse and Lagacherie (2017) and Szatmári and  
561 Pásztor (2019). KED had worse CRPS and IS values compared to QRF and QRPP RF, and also for  
562 QRNN in case of log(SOC). For pH this may seem counterintuitive at first. Here, KED was sharper and  
563 more reliable than QRF. Therefore, one might expect KED to obtain a better score. However, CRPS and  
564 IS are more sensitive to outliers. This is specifically reflected by IS at the 99% PI, i.e. IS(0.99) (Fig. 9).  
565 For IS(0.99), KED scored even worse than NM, despite KED having better sharpness and a PICP close  
566 to 99%. When evaluating how IS is calculated (Eq. 10), it is clear that IS imposes a linear penalty  
567 depending on how far outside the PI boundaries the test observations fall. Consequently, in the case of  
568 KED, the test observations that were not within the 99% PI had a large distance to the boundaries of the  
569 PI. It is also noteworthy that KED performed better than QRF in terms of the median CRPS. This is due  
570 to the fact that the median CRPS ignores issues with spread and outliers.

571 QRNN obtained strongly over-optimistic probabilistic predictions. Therefore, QRNN had the worst  
572 outcomes among all reliability measures and was also apart from NM the worst with respect to the  
573 scoring rules. We do not know what caused these issues. A possible explanation might be that the  
574 hyperparameter selection of QRNN was based on its point prediction performance and that the selected  
575 hyperparameters influenced the probabilistic prediction performance. The relationship between  
576 hyperparameters and probabilistic performance might need more caution considering that potential  
577 issues with overfitting of QRNN were discussed in Zhang et al. (2019). Nonetheless, it has to be noted  
578 that in one example outside DSM, decent reliability was achieved with QRNN using the same R-package  
579 (David et al., 2018). Therefore, we encourage more research with QRNN, in which its application is  
580 tested and if needed optimized. Furthermore, there are other noteworthy probabilistic adaptations of neural  
581 networks that are popular and eventually useful for DSM. Two examples are the lower upper bound  
582 estimation (LUBE) neural network, which predicts PIs (Khosravi et al., 2011), and the 'improved'  
583 version of QRNN (iQRNN), which aims to prevent overfitting and is reportedly faster (Zhang et al.,  
584 2019).

585

## 586 **4.2 Value of the proposed uncertainty metrics**

587 Much attention within DSM has been devoted to the importance of strict and rigorous validation  
588 procedures for point predictions, in order to provide an honest estimate of the quality of a soil map  
589 (Piikki et al., 2021). However, the validation of probabilistic predictions has not yet received the same  
590 amount of attention. So far, most analyses relied entirely on PICP when validating the reliability of  
591 probabilistic predictions. However, as introduced in the methodological framework but also now shown  
592 in a real-world case with KED and QRNN for log(SOC), PICP can hide one-sided bias of the estimated  
593 quantiles that set the boundaries of a PI. In this instance, PICP did not capture the actual probabilistic  
594 performance well and indicated better results than actually present. Such kind of bias may not occur  
595 frequently in practice as the other models did not show the same issues. Nonetheless, in order to truly  
596 validate the reliability of probabilistic predictions, PICP must be accompanied by a metric that can  
597 compensate for that weakness, such as the PIT histogram and QCP. Additionally, the PIT histogram and  
598 QCP have an intuitive interpretation similar to that of PICP. That means that they are not only useful  
599 for comparing models but can inherently show the reliability of a probabilistic predictions. Therefore,  
600 we strongly encourage to adopt these two metrics whenever the reliability of an uncertainty map is  
601 validated in DSM, so that they either supplement or replace PICP.

602 Scoring rules allow to rank the probabilistic performance of the models based on a returned numeric  
603 score value. RELI solely evaluated the reliability of probabilistic predictions and thus mostly reflected  
604 the trends found with PICP, QCP and PIT histograms. Hence, RELI is a useful metric to summarize  
605 probabilistic prediction performances with regards to reliability. CRPS and IS were more sensitive to  
606 outliers and sharpness. Since sharpness and point prediction performances are related, CRPS and IS  
607 reflect trends observed with point prediction validation metrics. This was especially apparent for NM,  
608 which scored last due to its low sharpness explained by the poor point prediction performance.  
609 Furthermore, the strong effect of outliers on IS and CRPS may be seen as a nuisance because it means  
610 that the outcome can be influenced by a few bad predictions. On the other hand, if test observations are  
611 too often found in areas of the predictive CDF that had a low probability density, it can confidently be  
612 interpreted as a major flaw of the probabilistic prediction model. However, at least for CRPS, the effect  
613 of outliers could be removed by evaluation of the median CRPS. For IS, the returned score additionally

614 depended on the PI level. The penalty for observations outside the PI was more severe with increasing  
615  $\tau$ , i.e. it was bigger at larger PI levels. Hence, the final judgment should not be grounded on one PI level  
616 evaluated by IS.

617 Usually, when we validate a single PIW-uncertainty map in DSM, we are mainly interested in the  
618 reliability of the uncertainty map. In such case, QCP and PIT histograms next to PICP are preferred  
619 metrics, as they have an intuitive interpretation. They can inherently show if an uncertainty map is  
620 reliable and thus safe to use for an end-user. One can also provide numerical scores from scoring rules  
621 but there might not be a great benefit in doing so, because scoring rules are mainly useful for the purpose  
622 of ranking model performance. Users are not necessarily interested in comparing performances of  
623 multiple models. Researchers on the other hand are more often interested in ranking competing  
624 probabilistic prediction models (e.g. Kasraei et al., 2021; Szatmári and Pásztor, 2019; Vaysse and  
625 Lagacherie, 2017). For such purposes, scoring rules are a very useful tool by which new information  
626 about the weaknesses and strengths of models may become apparent.

627

### 628 **4.3 Beyond the proposed metrics**

629 As demonstrated and shown in this study, the validation of probabilistic predictions is more complicated  
630 than for point predictions. In this paper, we proposed metrics from the broader probabilistic literature  
631 (e.g. Bracher et al., 2021; Brown et al., 2010; Gneiting et al., 2007; Gneiting and Raftery, 2007; Lauret  
632 et al., 2019; Pinson et al., 2007) that we deemed useful, intuitive and easy to implement for the context  
633 of DSM. Nonetheless, a short outlook is given about concepts and metrics that we did not address but  
634 may be worth exploring.

635 There is a large pool of other scoring rules which can be used to validate probabilistic predictions  
636 (Gneiting and Raftery, 2007). For example, the logarithmic score, also known as ignorance score, is  
637 another very popular and commonly used scoring rule. It takes the logarithm of the predictive probability  
638 density. Therefore, it is even more sensitive to outliers than CRPS because if the predictive probability  
639 density is close to zero, it converges to infinity or minus infinity, depending on whether a negative or  
640 positive orientation is chosen. However, Bracher et al. (2021) argued that the edges of a probability

641 density may not be reliably approximated from a set of predicted quantiles, which can be detected with  
642 the logarithmic score.

643 In Section 4.2 we stated that it may not be necessary to include scoring rules when validating an  
644 uncertainty map. However, in other academic fields, scoring rules are sometimes expressed in form of  
645 skill scores (Gneiting and Raftery, 2007; Lauret et al., 2019). A skill score measures the relative  
646 performance, where an obtained score is compared to a reference, for which NM is a natural choice. As  
647 such it has a similar logic as the MEC for point predictions, for which the NM produces a value close to  
648 zero. Therefore, if it is desired to further deliver the probabilistic performance in terms of a scoring rule,  
649 a skill score allows for a more general interpretability than the pure absolute numeric scores presented  
650 in this study.

651 The PIT-histograms presented in this study were interpreted upon subjective visual judgment.  
652 Goodness-of-fit test could be used to test the ‘flatness’ of a histogram (Elmore, 2005). It tests the null  
653 hypothesis of whether the obtained percentiles used for the PIT histogram follows a uniform distribution.  
654 It can be used to reduce the subjectivity of a purely visual assessment.

655 We restricted our analysis to sharpness and reliability. Yet, there are also other attributes that can be  
656 evaluated. For example, in the context of ensemble forecasts, multiple studies advocated to also evaluate  
657 probabilistic predictions with the concept of resolution (Brown et al., 2010; Lauret et al., 2019; Pinson  
658 et al., 2007). Resolution measures how case-dependent the resulting probabilistic predictions are,  
659 meaning that different predictive distributions are generated depending on the covariate condition.

660 In this study, we only focused on how probabilistic predictions and thus uncertainty maps could be  
661 evaluated within DSM with respect to validation metrics. We did not address the importance of a  
662 validation strategy (e.g. independent sampling, cross-validation or data-splitting) and sampling design  
663 (e.g. probability or non-probability sampling). These are important aspects to obtain unbiased validation  
664 results. Yet, they were already studied and discussed in great depth for point predictions (Brus et al.,  
665 2011; Piikki et al., 2021) and the same rules apply to the validation of probabilistic predictions. Note  
666 that the LUCAS-soil dataset used in this study is based on a multistage stratified random sampling design  
667 (Orgiazzi et al., 2018).

668

## 669 **5 Conclusion**

670 New metrics and concepts for the validation of probabilistic predictions were introduced and their  
671 relevance for DSM studies was illustrated in a case study with five different prediction models for pH  
672 and log(SOC). The methodical framework can be used to improve currently used validation procedures  
673 in DSM. Our conclusions are:

- 674 - PICP cannot truly validate the reliability of a probabilistic predictions because it is incapable to  
675 test for one-sided bias of the lower- and upper boundaries of a PI. Considerable one-sided bias  
676 may occur in practice as shown for KED and QRNN. Therefore, other validation tools like QCP  
677 or PIT histogram should complement or replace PICP.
- 678 - Scoring rules such as CRPS, IS and RELI allow for a ranking of probabilistic prediction model  
679 performances based on a numeric value. CRPS and IS were sensitive to outliers and sharpness.  
680 RELI summarized the trend found with reliability metrics such as PICP, QCP and PIT  
681 histograms and was less sensitive to outliers. Yet, it has to be acknowledged that scoring rules  
682 are mostly useful for comparing probabilistic performances of competing models.
- 683 - Depending on the metrics evaluated, different outcomes can be perceived in terms of  
684 probabilistic performance. Therefore, including a set of different validation metrics allows for  
685 a more critical evaluation.
- 686 - Considering all metrics for evaluating the probabilistic prediction performance of the five  
687 prediction models: NM showed high reliability but suffered from low sharpness; QRF had over-  
688 pessimistic uncertainty in the center of the distribution but performed well on the edges, QRPP  
689 RF was the most consistent; KED obtained inconsistent results and QRNN suffered from low  
690 reliability due to over-optimistic probabilistic predictions.

691 We strongly encourage to use the recommended tools in future studies for a more comprehensive and  
692 honest validation of probabilistic predictions in the field of DSM.

693

## 694 **Acknowledgements**

695 We thank Dr. James Brown for useful comments and ISRIC - World Soil Information for providing  
696 the covariate data.

697

## 698 **Declaration of competing interest**

699 The authors declare that they have no known competing financial interests or personal relationships  
700 that could have appeared to influence the work reported in this paper.

701

## 702 **Funding**

703 This research did not receive any specific grant from funding agencies in the public, commercial, or  
704 not-for-profit sectors.

705

## 706 **6 References**

707 Bracher, J., Ray, E.L., Gneiting, T., Reich, N.G., 2021. Evaluating epidemic forecasts in an interval  
708 format. *PLOS Computational Biology* 17 (2), e1008618.

709 Breiman, L., 2001. Random Forests. *Machine Learning* 45 (1), 5–32.

710 Breure, T.S., Haefele, S.M., Hannam, J.A., Corstanje, R., Webster, R., Moreno-Rojas, S., Milne, A.E.,

711 2022. A loss function to evaluate agricultural decision-making under uncertainty: a case study of

712 soil spectroscopy. *Precision agriculture* 23 (4), 1333–1353.

713 Brown, J.D., Demargne, J., Seo, D.-J., Liu, Y., 2010. The Ensemble Verification System (EVS): A

714 software tool for verifying ensemble forecasts of hydrometeorological and hydrologic variables at

715 discrete locations. *Environmental Modelling & Software* 25 (7), 854–872.

716 Brus, D.J., Heuvelink, G.B.M., 2007. Optimization of sample patterns for universal kriging of

717 environmental variables. *Geoderma* 138 (1-2), 86–95.

718 Brus, D.J., Kempen, B., Heuvelink, G.B.M., 2011. Sampling for validation of digital soil maps.

719 *European Journal of Soil Science* 62 (3), 394–407.



720 Cannon, A.J., 2011. Quantile regression neural networks: Implementation in R and application to  
721 precipitation downscaling. *Computers & Geosciences* 37 (9), 1277–1284.

722 Cannon, A.J., 2019. qrmn: Quantile Regression Neural Network. R-package Version 2.0.5.

723 Caubet, M., Dobarco, M.R., Arrouays, D., Minasny, B., Saby, N.P., 2019. Merging country,  
724 continental and global predictions of soil texture: Lessons from ensemble modelling in France.  
725 *Geoderma* 337, 99–110.

726 Chen, S., Arrouays, D., Mulder, V.L., Poggio, L., Minasny, B., Roudier, P., Libohova, Z., Lagacherie,  
727 P., Shi, Z., Hannam, J., Meersmans, J., Richer-de-Forges, A.C., Walter, C., 2022. Digital mapping  
728 of GlobalSoilMap soil properties at a broad scale: A review. *Geoderma* 409, 115567.

729 David, M., Luis, M.A., Lauret, P., 2018. Comparison of intraday probabilistic forecasting of solar  
730 irradiance using only endogenous data. *International Journal of Forecasting* 34 (3), 529–547.

731 Elmore, K.L., 2005. Alternatives to the Chi-Square Test for Evaluating Rank Histograms from  
732 Ensemble Forecasts. *Wea. Forecasting* 20 (5), 789–795.

733 Gneiting, T., Balabdaoui, F., Raftery, A.E., 2007. Probabilistic forecasts, calibration and sharpness. *J*  
734 *Royal Statistical Soc B* 69 (2), 243–268.

735 Gneiting, T., Raftery, A.E., 2007. Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of*  
736 *the American Statistical Association* 102 (477), 359–378.

737 Goovaerts, P., 2001. Geostatistical modelling of uncertainty in soil science. *Geoderma* 103 (1-2), 3–  
738 26.

739 Hersbach, H., 2000. Decomposition of the Continuous Ranked Probability Score for Ensemble  
740 Prediction Systems. *Wea. Forecasting* 15 (5), 559–570.

741 Heuvelink, G.B.M., 1998. *Error Propagation in Environmental Modelling with GIS*. CRC Press.

742 Heuvelink, G.B.M., 2018. *Uncertainty and Uncertainty Propagation in Soil Mapping and Modelling*,  
743 in: *Pedometrics*. Springer, Cham, pp. 439–461.

744 Hiemstra, P.H., 2022. automap: Automatic Interpolation Package. R-package version 1.0-16.

745 ISO, 1994. *Soil quality — Determination of pH*. International Organization for Standardization,  
746 Geneva.

747 ISO, 1995. Soil quality — Determination of organic and total carbon after dry combustion (elementary  
748 analysis). International Organization for Standardization, Geneva.

749 Kasraei, B., Heung, B., Saurette, D.D., Schmidt, M.G., Bulmer, C.E., Bethel, W., 2021. Quantile  
750 regression as a generic approach for estimating uncertainty of digital soil maps produced from  
751 machine-learning. *Environmental Modelling & Software* 144, 105139.

752 Keesstra, S.D., Bouma, J., Wallinga, J., Tiftonell, P., Smith, P., Cerdà, A., Montanarella, L., Quinton,  
753 J.N., Pachepsky, Y., van der Putten, W.H., Bardgett, R.D., Moolenaar, S., Mol, G., Jansen, B.,  
754 Fresco, L.O., 2016. The significance of soils and soil science towards realization of the United  
755 Nations Sustainable Development Goals. *SOIL* 2 (2), 111–128.

756 Khaledian, Y., Miller, B.A., 2020. Selecting appropriate machine learning methods for digital soil  
757 mapping. *Applied Mathematical Modelling* 81, 401–418.

758 Khosravi, A., Nahavandi, S., Creighton, D., Atiya, A.F., 2011. Lower upper bound estimation method  
759 for construction of neural network-based prediction intervals. *IEEE Trans. Neural Netw.* 22 (3),  
760 337–346.

761 Koenker, R., 2022. *quantreg: Quantile Regression*. R-package version 5.94.

762 Koenker, R., Hallock, K.F., 2001. Quantile Regression. *Journal of Economic Perspectives* 15 (4), 143–  
763 156.

764 Lagacherie, P., Arrouays, D., Bourennane, H., Gomez, C., Martin, M., Saby, N.P., 2019. How far can  
765 the uncertainty on a Digital Soil Map be known?: A numerical experiment using pseudo values of  
766 clay content obtained from Vis-SWIR hyperspectral imagery. *Geoderma* 337, 1320–1328.

767 Lark, R.M., Chagumaira, C., Milne, A.E., 2022. Decisions, uncertainty and spatial information. *Spatial*  
768 *Statistics* 50, 100619.

769 Lauret, P., David, M., Pinson, P., 2019. Verification of solar irradiance probabilistic forecasts. *Solar*  
770 *Energy* 194, 254–271.

771 Liaw, A., Wiener, M., 2022. *randomForest: Classification and Regression by randomForest*. R-  
772 package version 4.7-1.1.

773 McBratney, A., Mendonça Santos, M., Minasny, B., 2003. On digital soil mapping. *Geoderma* 117 (1-  
774 2), 3–52.

775 Meinshausen, N., 2006. Quantile regression forests. *Journal of machine learning research* 7 (6).  
776 Meinshausen, N., 2017. *quantregForest: Quantile Regression Forests*. R-package version 1.3-7.  
777 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I — A  
778 discussion of principles. *Journal of Hydrology* 10 (3), 282–290.  
779 NCAR - Research Applications Laboratory, 2015. *verification: Weather Forecast Verification*  
780 *Utilities*. R-package version 1.42.  
781 Nelson, M.A., Bishop, T.F.A., Triantafilis, J., Odeh, I.O.A., 2011. An error budget for different  
782 sources of error in digital soil mapping. *European Journal of Soil Science* 62 (3), 417–430.  
783 Nussbaum, M., Papritz, A., Baltensweiler, A., Walthert, L., 2014. Estimating soil organic carbon  
784 stocks of Swiss forest soils by robust external-drift kriging. *Geosci. Model Dev.* 7 (3), 1197–1210.  
785 Orgiazzi, A., Ballabio, C., Panagos, P., Jones, A., Fernández-Ugalde, O., 2018. LUCAS Soil, the  
786 largest expandable soil dataset for Europe: a review. *European Journal of Soil Science* 69 (1), 140–  
787 153.  
788 Pebesma, E., 2022. *gstat: Spatial and Spatio-Temporal Geostatistical Modelling, Prediction and*  
789 *Simulation*. R-package version 2.1-0.  
790 Piikki, K., Wetterlind, J., Söderström, M., Stenberg, B., 2021. Perspectives on validation in digital soil  
791 mapping of continuous attributes—A review. *Soil Use Manage* 37 (1), 7–21.  
792 Pinson, P., Nielsen, H.A., Møller, J.K., Madsen, H., Kariniotakis, G.N., 2007. Non-parametric  
793 probabilistic forecasts of wind power: required properties and evaluation. *Wind Energ.* 10 (6),  
794 497–516.  
795 Pinson, P., Tastu, J., 2014. Discussion of “Prediction Intervals for Short-Term Wind Farm Generation  
796 Forecasts” and “Combined Nonparametric Prediction Intervals for Wind Power Generation”. *IEEE*  
797 *Trans. Sustain. Energy* 5 (3), 1019–1020.  
798 Poggio, L., Sousa, L.M. de, Batjes, N.H., Heuvelink, G.B.M., Kempen, B., Ribeiro, E., Rossiter, D.,  
799 2021. SoilGrids 2.0: producing soil information for the globe with quantified spatial uncertainty.  
800 *SOIL* 7 (1), 217–240.  
801 R Core Team, 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for  
802 Statistical Computing, Vienna.

803 Szatmári, G., Pásztor, L., 2019. Comparison of various uncertainty modelling approaches based on  
804 geostatistics and machine learning algorithms. *Geoderma* 337, 1329–1340.

805 Vasseur, S.P., Aznarte, J.L., 2021. Comparing quantile regression methods for probabilistic  
806 forecasting of NO<sub>2</sub> pollution levels. *Sci Rep* 11 (1), 11592.

807 Vaysse, K., Lagacherie, P., 2017. Using quantile regression forest to estimate uncertainty of digital  
808 soil mapping products. *Geoderma* 291, 55–64.

809 Webster, R., Oliver, M.A., 2007. Kriging in the Presence of Trend and Factorial Kriging, in: Webster,  
810 R., Oliver, M.A. (Eds.), *Geostatistics for environmental scientists*, Second Edition ed. *Statistics in*  
811 *practice*. Wiley, Chichester, pp. 195–218.

812 Zamo, M., Naveau, P., 2018. Estimation of the Continuous Ranked Probability Score with Limited  
813 Information and Applications to Ensemble Weather Forecasts. *Math Geosci* 50 (2), 209–234.

814 Zhang, W., Quan, H., Srinivasan, D., 2019. An Improved Quantile Regression Neural Network for  
815 Probabilistic Load Forecasting. *IEEE Trans. Smart Grid* 10 (4), 4425–4434.

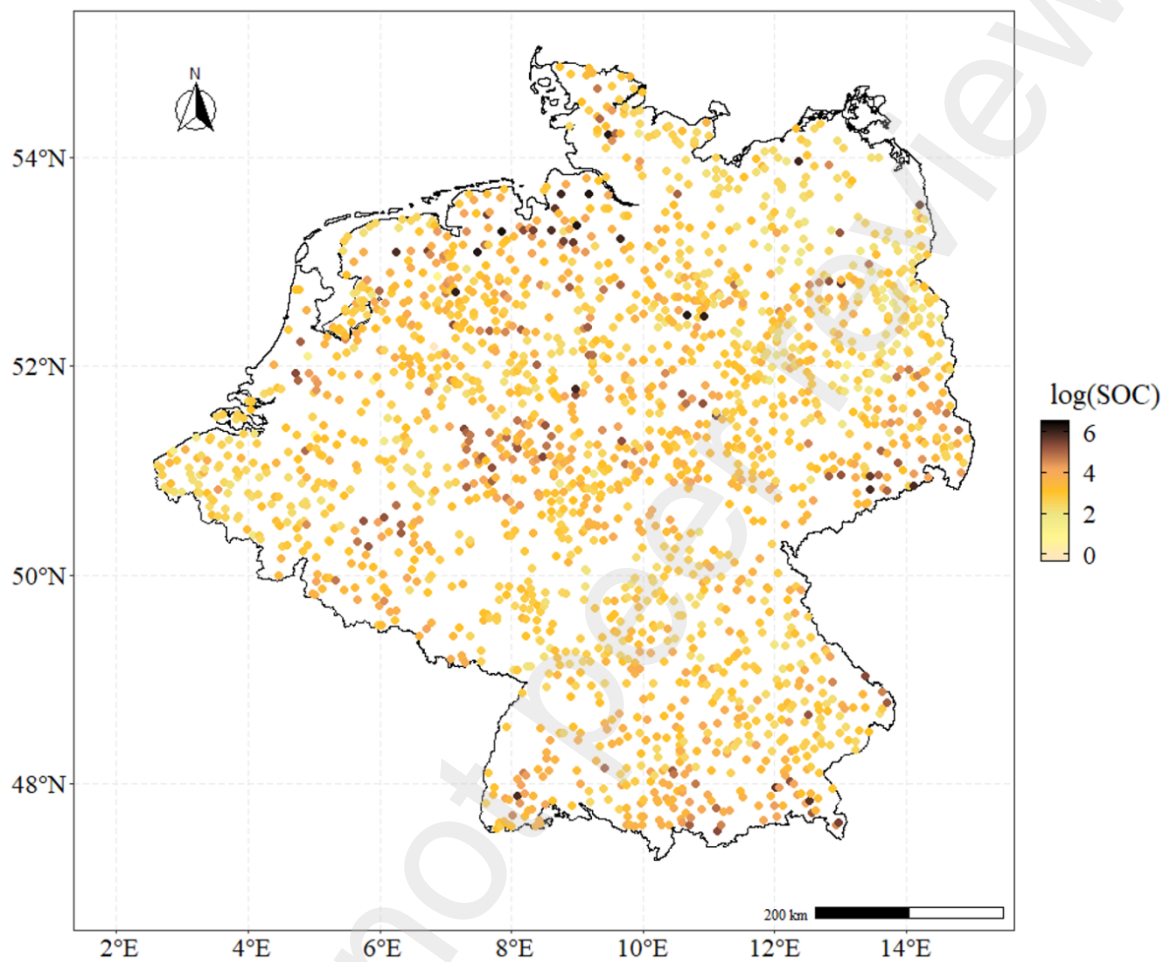
816 Zhang, Y., Wang, J., Wang, X., 2014. Review on probabilistic forecasting of wind power generation.  
817 *Renewable and Sustainable Energy Reviews* 32, 255–270.

818

819 **Supplementary Information**

820 **log(SOC) results**

821



822

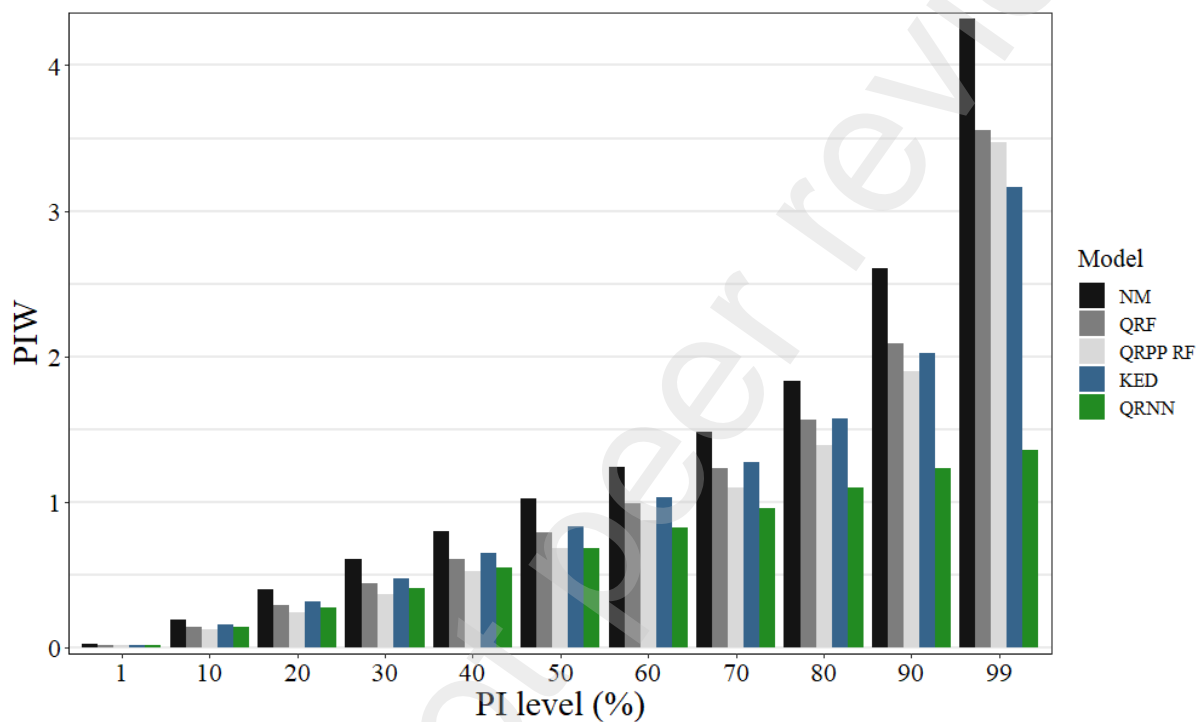
823 Fig. A1. Sampling sites from the LUCAS-dataset of log(SOC) for the study area of interest.

824

825 Table A1. Point prediction performances for log(SOC).

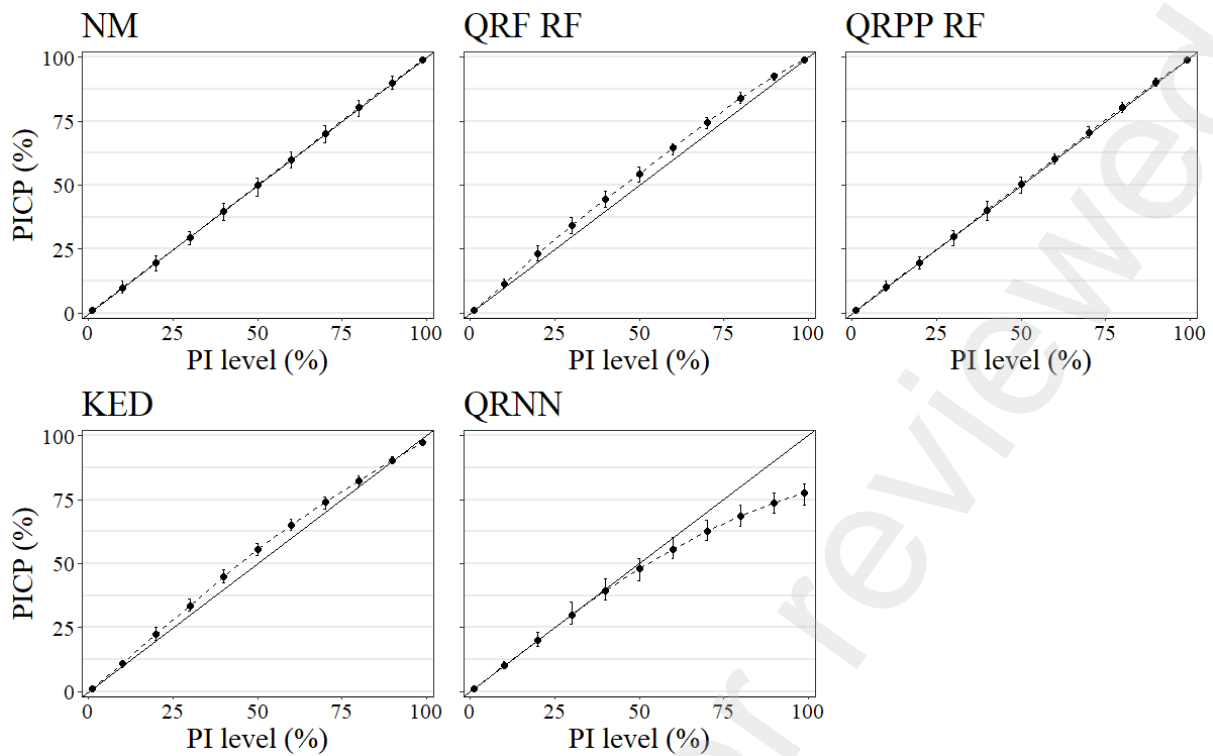
|      | NM      | QRF   | QRPP RF | KED    | QRNN   |
|------|---------|-------|---------|--------|--------|
| MEC  | -0.0034 | 0.42  | 0.43    | 0.36   | 0.41   |
| RMSE | 0.80    | 0.61  | 0.61    | 0.64   | 0.62   |
| ME   | -0.004  | 0.002 | 0.013   | -0.001 | -0.017 |

826



827

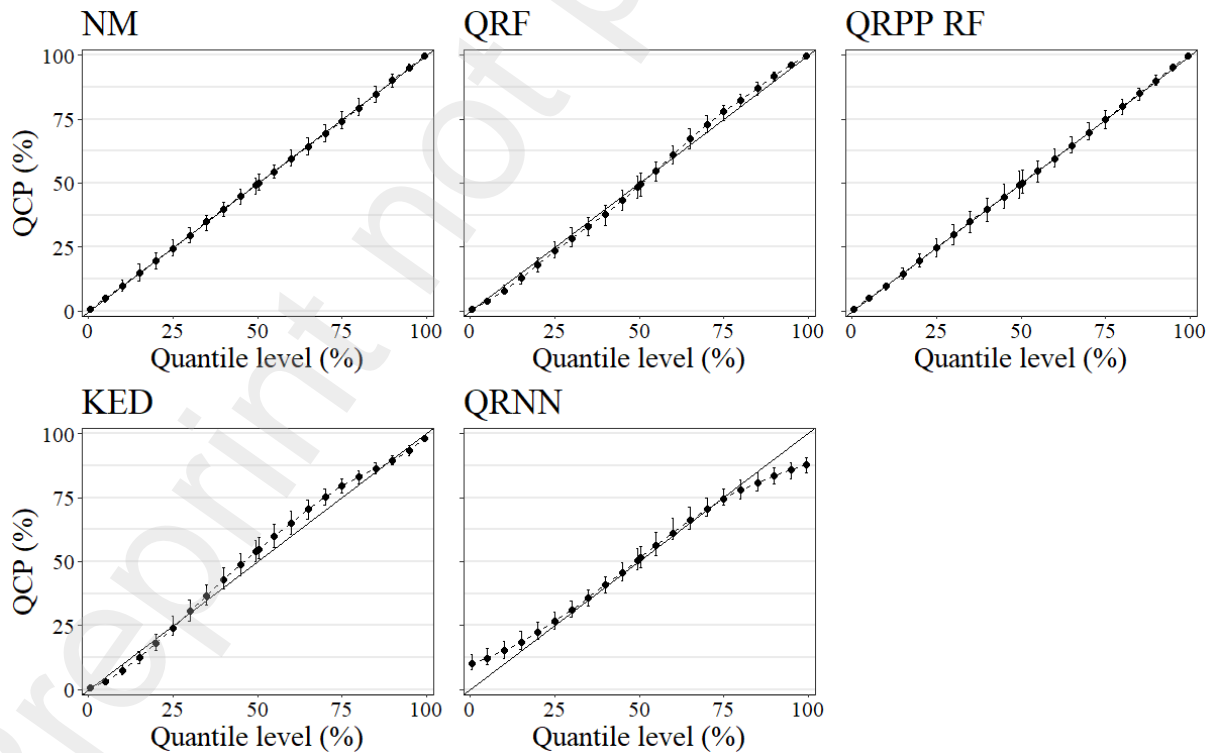
828 Fig. A2. Sharpness diagram indicating the PIW throughout the predictive distribution for log(SOC).



829

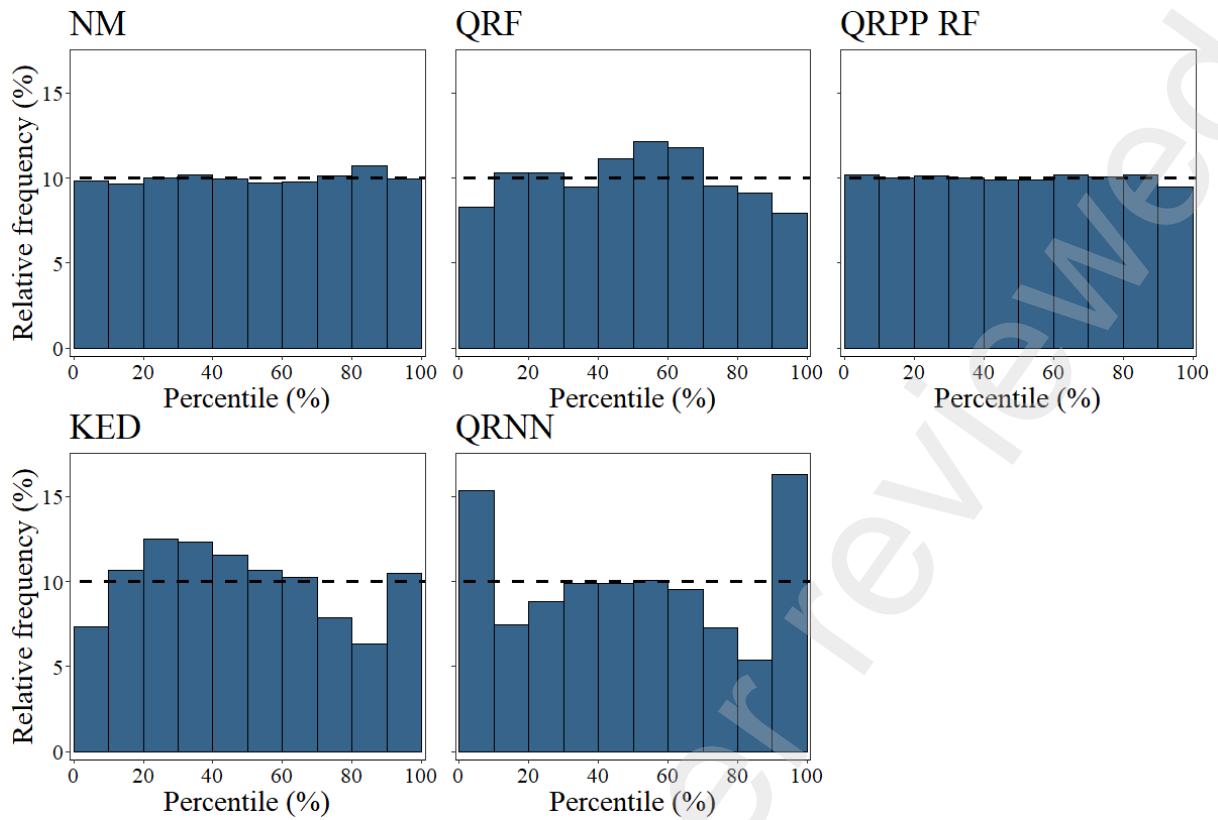
830 Fig. A3. Mean PICP reliability plots for log(SOC). Error bars are retrieved from the 80% confidence interval of  
 831 the 25 repetitions in the outer loop. The 1:1 black line indicates the desired outcome.

832



833

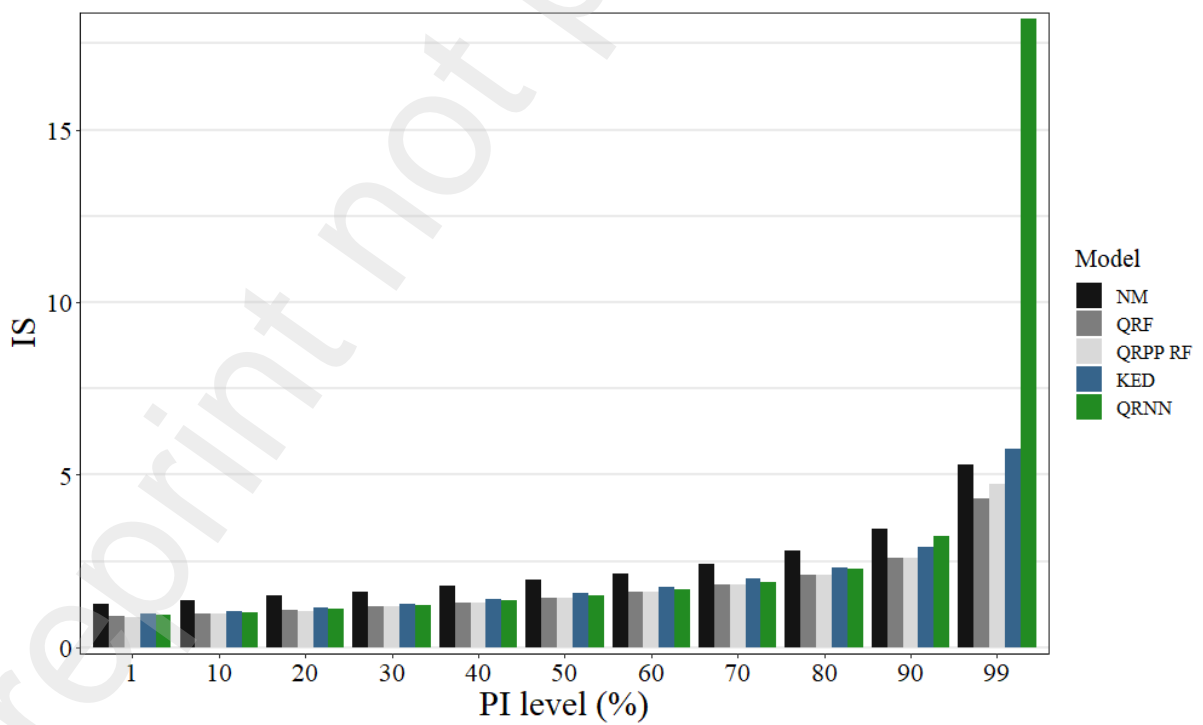
834 Fig. A4. Mean QCP reliability plots for log(SOC). Error bars are retrieved from the 80% confidence interval of  
 835 the 25 repetitions in the outer loop. The 1:1 black line indicates the desired outcome.



836

837 Fig. A5. PIT histograms of log(SOC). The dashed line indicates the desired frequency for a flat and uniform PIT.

838



839

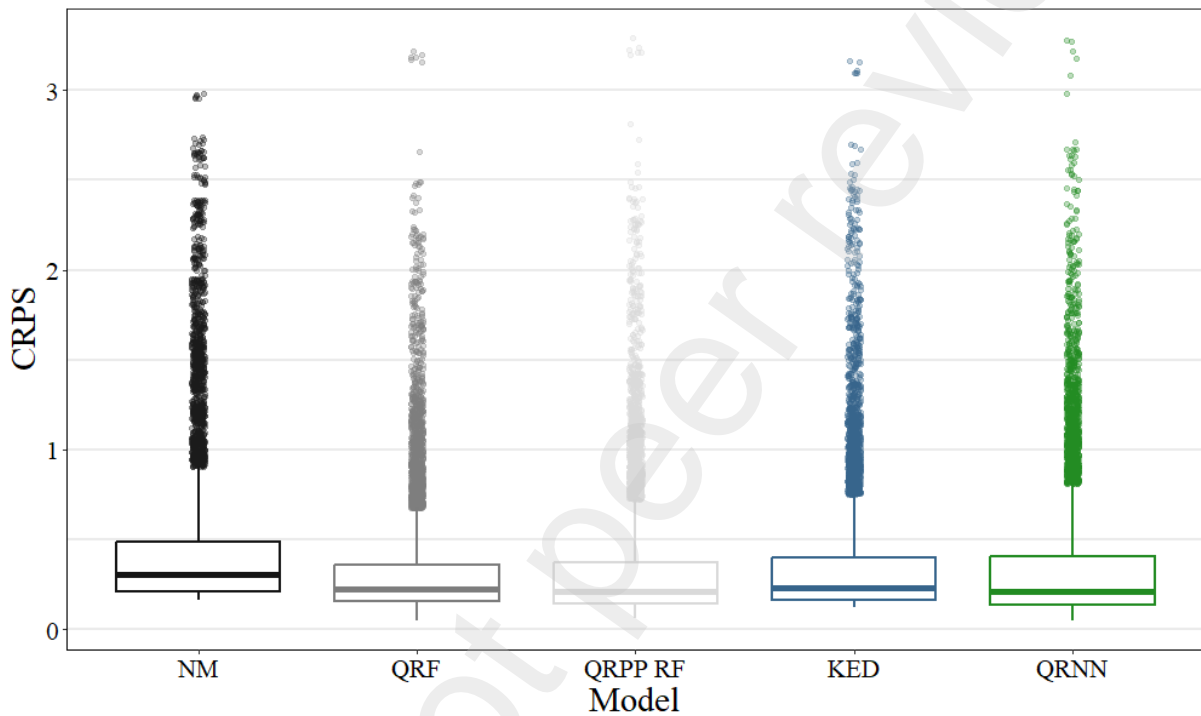
840 Fig. A6. IS diagram of log(SOC) indicating the IS values throughout the predictive distribution.



841 Table A2. Scoring outcomes of CRPS, median CRPS and RELI for log(SOC).

|             | NM     | QRF    | QRPP RF | KED    | QRNN   |
|-------------|--------|--------|---------|--------|--------|
| CRPS        | 0.44   | 0.32   | 0.32    | 0.35   | 0.34   |
| Median CRPS | 0.30   | 0.22   | 0.21    | 0.23   | 0.21   |
| RELI        | 0.0014 | 0.0029 | 0.0013  | 0.0036 | 0.0135 |

842



843

844 Fig. A7. Single CRPS values portrayed as boxplots for log(SOC). Each boxplot was based on 25 x 2,016 values.

845

### 846 **Continuous ranked probability score**

847 Hersbach (2000) introduced a method to calculate CRPS for a finite number of members in the context

848 of ensemble forecasting, but the concepts can also be applied to probabilistic models in quantile format.

849 Here CRPS is computed from the sum of  $M$  squared rectangles, where  $M$  is the number of quantiles, i.e.

850 steps that were used to approximate a predictive CDF. For that, we used the 199 quantile set as described

851 in Section 2.4 and additionally a 99.99% quantile to approximate the 100% quantile (Zamo and Naveau,

852 2018), leading to  $M = 200$ . The last step was done to evaluate outliers, to which we refer later in this

853 section. CRPS is estimated as follows:

$$CRPS = \frac{1}{n} \sum_{i=1}^n \sum_{j=0}^M [\alpha_{ij} \tau_{ij}^2 + \beta_{ij} (1 - \tau_{ij})^2], \quad (A1)$$

$$\tau_{ij} = \frac{j}{M}, \quad (A2)$$

854 where  $\tau_{ij}$  is the associated probability to the quantile and  $\alpha_{ij}$  and  $\beta_{ij}$  are the bin widths of each step in the  
 855 predictive CDF. Both  $\alpha_{ij}$  and  $\beta_{ij}$  are determined via Table A3. Test samples  $y_i$  that are lying either in the  
 856 0% range ( $y_i < q_{i(j=1)}$ ) or 100% range ( $q_{i(j=M)} < y_i$ ) of the CDF are defined as outliers, where  $q_{ij}$  is the  
 857 predicted quantile with the probability level  $\tau_{ij}$ .

858

859 Table A3. Determination of  $\alpha_{ij}$  and  $\beta_{ij}$ .

| $0 < j < M$                    | $\alpha_{ij}$         | $\beta_{ij}$          |
|--------------------------------|-----------------------|-----------------------|
| $y_i > q_{i(j+1)}$             | $q_{i(j+1)} - q_{ij}$ | 0                     |
| $q_{i(j+1)} > y_i > q_{ij}$    | $y_i - q_{ij}$        | $q_{i(j+1)} - y_i$    |
| $y_i < q_{ij}$                 | 0                     | $q_{i(j+1)} - q_{ij}$ |
| Outlier for $j = 0$ or $j = M$ | $\alpha_{ij}$         | $\beta_{ij}$          |
| $y_i < q_{i(j=1)}$             | 0                     | $q_{i(j=1)} - y_i$    |
| $q_{i(j=M)} < y_i$             | $y_i - q_{i(j=M)}$    | 0                     |

860

## 861 Reliability decomposition

862 RELI can be calculated by defining  $\bar{\alpha}_j$  and  $\bar{\beta}_j$ , where  $\bar{\alpha}_j = \frac{1}{n} \sum_{i=1}^n \alpha_{ij}$  and  $\bar{\beta}_j = \frac{1}{n} \sum_{i=1}^n \beta_{ij}$ :

$$RELI = \sum_{j=0}^M \bar{g}_j (\bar{\sigma}_j - \tau_j)^2, \quad (A3)$$

$$\bar{g}_j = \sum_{j=0}^M \bar{\alpha}_j + \bar{\beta}_j, \quad (A4)$$

$$\bar{\sigma}_j = \sum_{j=0}^M \frac{\bar{\beta}_i}{\bar{\alpha}_j + \bar{\beta}_j}, \quad (A5)$$

863 To include outliers at  $j = 0$  and  $j = M$  we define:

$$\bar{\sigma}_0 = \frac{1}{n} \sum_{i=1}^n H(q_{i(j=1)} - y_i) \quad \bar{g}_0 = \frac{\bar{\beta}_0}{\bar{\sigma}_0}, \quad (A6)$$

$$\bar{o}_M = \frac{1}{n} \sum_{i=1}^n H(q_{i(j=M)} - y_i) \quad \bar{g}_M = \frac{\bar{\alpha}_M}{(1 - \bar{o}_M)}. \quad (A7)$$

864 Thereby,  $\bar{o}_j$  can be interpreted very similarly to the average empirical frequency of data being below a  
 865 specific quantile with probability  $\tau_j$ . Hence, RELI is related to PIT and QCP. But RELI additionally  
 866 considers the average width between a quantile  $q_j$  and its neighboring quantile  $q_{j+1}$ , which is  $\bar{g}_j$ . For more  
 867 detailed explanations and interpretations see Hersbach (2000).