



1 **Estimation of soil water holding capacity with Random Forests for**
2 **drought monitoring**

3

4

5 Yves Tramblay ^{1*}

6 Pere Quintana Seguí ²

7

8

9 ¹ HydroSciences Montpellier (University Montpellier, CNRS, IRD), France

10 ² Observatori de l'Ebre (OE), Ramon Llull University – CSIC, 43520 Roquetes, Spain

11

12

13 *Corresponding author : yves.tramblay@ird.fr, 300 avenue du Pr. Emile Jeanbreaux,
14 34090, Montpellier, France. +33 4 67 14 33 59

15

16 **Abstract**

17

18 Soil moisture is a key variable for drought monitoring but soil moisture measurements
19 networks are very scarce. Land-surface models can provide a valuable alternative to
20 simulate soil moisture dynamics, but only a few countries have such modelling schemes
21 implemented for monitoring soil moisture at high spatial resolution. In this study, a soil
22 moisture accounting model (SMA) was regionalized over the Iberian Peninsula, taking as
23 a reference the soil moisture simulated by a high-resolution land surface model. To
24 estimate soil water holding capacity, the parameter required to run the SMA model, two
25 approaches were compared: the direct estimation from European soil maps using
26 pedotransfer functions, or an indirect estimation by a Machine Learning approach,
27 Random Forests, using as predictors altitude, temperature, precipitation,
28 evapotranspiration and land use. Results showed that the Random Forest model
29 estimates are more robust, especially for estimating low soil moisture levels.
30 Consequently, the proposed approach can provide an efficient way to simulate daily soil
31 moisture and therefore monitor soil moisture droughts, in contexts where high-resolution
32 soil maps are not available, as it relies on a set of covariates that can be reliably estimated
33 from global databases.

34

35

36

37

38 **Keywords:** soil moisture, droughts, random forests

39

40



41 **1. Introduction**

42

43 Soil moisture droughts have strong impacts on vegetation and agricultural production
44 (Raymond et al., 2019; Trambly et al., 2020; Vicente-Serrano et al., 2014; Pena-Gallardo
45 et al., 2019). There is a growing interest for simple indicators to monitor drought events
46 at short timescales that could be related to impacts (Li et al., 2020; Noguera et al., 2021).
47 In particular, soil moisture indicators could be more relevant than climatic ones to monitor
48 potential impacts of droughts on agriculture and natural vegetation (Piedallu et al., 2013).
49 Since actual soil moisture measurements remain very scarce, soil moisture simulated
50 from land-surface models are an interesting proxy to develop simplified methodologies
51 that could be applied on data-sparse regions. Land-surface models (LSM) are valuable
52 tools for a fine scale monitoring of drought events; however, their implementation requires
53 accurate forcing data and computational resources (Almendra-Martín et al., 2021;
54 Quintana-Seguí et al., 2019; Barella-Ortiz and Quintana-Seguí, 2019). Global
55 implementation also exists but with a coarser resolution and driven by reanalysis data
56 (Rodell et al., 2004; Muñoz Sabater, 2020) that may not be adequate for local-scale
57 applications. Only very few countries have land-surface schemes implemented at the
58 national level to monitor droughts (Habets et al., 2008).

59

60 Remote Sensing is another option which allows monitoring soil moisture (Dorigo et al.,
61 2017; Brocca et al., 2019). Microwave sensors allow monitoring of surface soil moisture
62 (first 5 cm for L-band based products, skin for C-band based products), without the
63 interference of clouds. However, surface soil moisture is not enough for most applications,
64 which require root zone soil moisture, which is the water resource in the soil available to
65 plants. Furthermore, passive L-band products, such as SMOS (Martínez-Fernández et
66 al., 2016) or SMAP (Mishra et al., 2017), have a low resolution and active C-band
67 products, such as Sentinel 1 (Bauer-Marschallinger et al., 2019), which have higher
68 resolution, suffer from higher noise and are more sensitive to vegetation. Thus, even
69 though remote sensing is very useful, it still has problems to be surmounted. The
70 resolution of passive L-band products can be increased using optical data (NDVI, LST),
71 by means of downscaling algorithms (Merlin et al., 2013; Fang et al., 2021), but then the
72 resulting product is sensitive to cloud cover. Also, some progress has been made in
73 deriving root zone soil moisture from surface soil moisture estimations using an
74 exponential filter (Stefan et al., 2021) calibrated using the SURFEX LSM (Masson et al.,
75 2013), but these products are in early stages and are not operational yet.

76

77 Simplified methodologies to estimate and monitor the status of soil moisture, are needed
78 in contexts where LSM data is not available and where remote sensing products fall short,
79 such as areas and time periods with dense vegetation, or high soil roughness which may
80 affect their accuracy (Escorihuela and Quintana-Seguí, 2016). Different modelling



81 approaches have been proposed, either with conceptual soil moisture accounting models
82 or computational variants of the antecedent precipitation index (Willgoose and Perera,
83 2001; Javelle et al., 2010; Brocca et al., 2014; Zhao et al., 2019; Li et al., 2020). The
84 general availability of spatial estimates of soil moisture content would help introduce soil
85 moisture in drought monitoring systems, improving their scope and usefulness.
86 Furthermore, this would also facilitate the creation of long-term reanalysis, based on
87 meteorological forcing data, and future climate change studies, without the need of
88 running LSM models. However, to apply this type of models at regional or national scale,
89 there is a need to estimate their parameters over the area of interest. For that purpose,
90 regionalization methods have been employed in hydrology for decades to estimate the
91 parameters of hydrological models in ungauged basins (Blöschl and Sivapalan, 1995; He
92 et al., 2011; Hrachowitz et al., 2013). Several methods exist, based either on catchment
93 similarity or the direct estimation of model parameters using regression techniques with
94 physiographic attributes. For soil moisture modelling, up to now only very few studies
95 have considered these approaches to apply soil moisture accounting models at ungauged
96 locations (Grillakis et al., 2021) or estimate root zone soil moisture using machine learning
97 methods (Carranza et al., 2021).

98
99 The goal of the present study is to regionalize (ie. to estimate from surrogate data without
100 calibration) the soil water holding capacity that is the sole parameter required in a simple
101 soil moisture model to monitor soil moisture droughts. Two different approaches are
102 compared: the direct estimation of the soil water holding capacity with soil maps or an
103 estimation with machine learning techniques, namely Random Forests.

104

105 **2. Study area and Data**

106

107 The study area of this work is the Iberian Peninsula, which is located between the
108 Mediterranean Sea and the Atlantic Ocean and thus is influenced by both synoptic scale
109 systems, that often come from the Atlantic side, and mesoscale heavy precipitation
110 events, that often come from the Mediterranean side. The Iberian Peninsula presents a
111 marked relief, with a large and high central plateau and different mountain ranges, which
112 heavily influence the spatial patterns of precipitation, enhancing it windward and
113 decreasing it leeward, generating areas of high precipitation on the west, north-west and
114 north, and very dry areas on the central plains and, specially, on the South-east, as a
115 consequence the Iberian Peninsula has a heterogeneous distribution of average annual
116 rainfall, with values ranging from 2000 mm/y to less than 100 mm/y. All this has a strong
117 influence on the spatial and temporal variability of soil moisture and soil moisture regimes,
118 having wet regimes on the west and north, where the soil is hardly stressed and, and
119 semi-arid areas elsewhere, with a wet (energy limited) and a dry (water limited) season,



120 with a dry down that might be interrupted by convective events. All this makes the
121 modelling of soil moisture in Iberian a rather challenging task.

122

123 Daily precipitation, temperature and evapotranspiration were retrieved from the SAFRAN-
124 Spain database (Quintana-Seguí et al., 2017). SAFRAN (Durand et al., 1993) is a
125 meteorological reanalysis that produces gridded datasets by combining the outputs of a
126 meteorological model and all available observations using an optimal interpolation
127 algorithm. It has been implemented over France (Quintana-Seguí et al., 2008) and
128 recently over the Iberian Peninsula (Quintana-Seguí et al., 2017) with a 5kmx5km spatial
129 resolution. The SAFRAN dataset used in this study not only includes observations from
130 the Spanish part of the Iberian Peninsula, it has also ingested data from Portugal. The
131 SURFEX LSM (Masson et al., 2013) has been run using SAFRAN-Spain as the
132 meteorological forcing dataset and on the same grid, as it was done in Quintana-Seguí
133 et al., (2020). SURFEX uses the ECOCLIMAP2 (Faroux et al., 2013) physiographic
134 database and it uses the ISBA (Interaction Sol-Biosphère-Atmosphère) scheme (Noilhan
135 and Mahfouf, 1996) for natural surfaces. ISBA has different options; we have used ISBA-
136 DIF, the multi-layer diffusion version (Boone 2000; Habets et al. 2003). From this
137 simulation, we have extracted the soil moisture of the first 60 cm of the soil, by performing
138 the weighted average of the soil layers that fall within this range. This simulated soil
139 moisture over the Iberian Peninsula is considered herein as the observed reference, in
140 the absence of dense monitoring networks of soil moisture (Martínez-Fernández et al.,
141 2015). From the ECOCLIMAP2 database, elevation and land cover data have also been
142 retrieved and aggregated in the following nine categories: water, bare, ice/snow, urban,
143 forest, grass, dry crops, irrigated crops, wetlands.

144

145 We use the European Soil database (ESDB) produced by the European Soil Data Centre
146 (Panagos et al., 2012). The ESDB contains information on soil characteristics, including
147 soil depth and texture for topsoil (0-30cm) and subsoil (30-70cm) layers at a grid
148 resolution of 1 km. The total available water content (TAWC) is a volumetric parameter
149 describing the water content between field capacity and permanent wilting point, as a
150 function of available water content, presence of coarse fragments and depth (Reynolds
151 et al., 2000). In ESDB, water content at field capacity and permanent wilting point were
152 determined following the equation from (van Genuchten, 1980) to estimate the soil water
153 retention curve (Hiederer, 2013). The parameters of the equation are provided by a
154 pedotransfer function (Wösten et al., 1999) for volumetric soil water content computed
155 from the soil water retention curve. The pedotransfer function uses soil texture, organic
156 carbon content and bulk density to determine the parameters of the soil water retention
157 curve (Hiederer, 2013). In the present work, the TAWC of subsoil and topsoil layers have
158 been added and averaged at the scale of 5km x 5km, matching the spatial resolution of



159 the SAFRAN grid. Then, these estimates have been used to set the A parameter of the
160 SMA model.

161

162 **3. Methods**

163

164 **3.1 Soil moisture accounting model**

165

166 We use a soil moisture accounting model (SMA) driven by precipitation and PET, with
167 one single parameter A, representing the soil water holding capacity. The soil moisture
168 model considered here has been previously applied in several studies for applications
169 related to soil moisture monitoring (Ancil et al., 2004; Javelle et al., 2010; Trambly et
170 al., 2012, 2014), it consists in the SMA part of the GR4J model (Perrin et al., 2003). The
171 output of the model is daily normalized soil moisture, allowing to detect the days close to
172 saturation (1) or to complete soil moisture depletion (0).

173

174 The SMA model is calibrated using soil moisture simulated with SURFEX covering the
175 full Iberian Peninsula domain. The Nelder-Mead simplex algorithm is used for the
176 calibration with the Nash efficiency criterion. The outputs of SURFEX soil moisture are
177 first normalized with the maximum and minimum values prior to the calibration to compute
178 the SWI consistent with the SMA model output. To regionally estimate the values of A,
179 two different methods are compared: the direct estimation of A with TAWC from ESDB
180 soil maps or its indirect estimation with machine learning methods, namely Random
181 Forests using 5kmx5km grid physiographic properties.

182

183 **3.2 Random forests for regionalization of soil water holding capacity**

184

185 Random Forests (Breiman, 2001) belong to the class of Machine Learning techniques.
186 RF are based on a bootstrap aggregation (Breiman, 1996) of Classification and
187 Regression Trees (Breiman et al., 2017). It generates a bootstrap sample from the original
188 data and trains a tree model using this sample. The procedure is repeated many times
189 and the bagging's prediction is the average of the predictions. Among the many
190 advantages of RF, they are fast, non-parametric, robust to noise in the predictor variables,
191 able to capture nonlinear dependencies between predictors and dependent variables and
192 they can simultaneously incorporate continuous and categorical variables (Tyrallis et al.,
193 2019). The drawbacks are they are complex to interpret and they cannot extrapolate
194 outside the training range. Given their advantages, this algorithm is particularly suited for
195 the estimation of spatial variables such as soil properties (Booker and Woods, 2014;
196 Hengl et al., 2018; Gagkas and Lilly, 2019; Stein et al., 2021). In the present work, a RF
197 model is generated to estimate the values of the A parameter of the SMA model,



198 representing soil water holding capacity, with the properties of the 5x5km grid cells using
199 Random Forests.

200

201 To estimate the reliability of the method, the 5km x 5km grid cells covering the Iberian
202 Peninsula have been split randomly into a training sample containing 70% of the cells
203 and a testing sample with the 30% remaining cells. The random selection of the training
204 and testing sets have been performed using a Latin Hypercube Sampling (McKay et al.,
205 1979) to ensure a homogeneous sampling over the Iberian Peninsula. Given that the RF
206 trees cannot be interpreted directly, as for example the weights in a linear regression, we
207 additionally implemented an out-of-bag predictor importance estimation by permutation
208 (Loh and Shih, 1997), to measure how influential the predictor variables in the model are
209 at predicting the response. The influence of a predictor increases with the value of this
210 measure. If a predictor is influential in prediction, then permuting its values should affect
211 the model error. If a predictor is not influential, then permuting its values should have little
212 to no effect on the model error.

213

214 **3.3 Validation on the ability to detect dry soil moisture conditions**

215

216 To compare the efficiency of the two methods compared to estimate the A parameter of
217 the SMA model, the SMA model was run using the two methods and all daily values of
218 soil moisture below the 10th percentile were extracted, corresponding to dry soil
219 conditions. Only the grid cells in the testing sample were considered for this validation.
220 We computed different verification scores to assess the relative efficiency of the two
221 methods to reproduce daily soil moisture below the 10th percentile using the ISBA
222 simulated soil moisture as a benchmark; the Probability of Detection (POD), the False
223 Alarm Ratio (FAR) and the Heidke Skill Score (HSS) summarizing the global efficiency to
224 detect dry periods (Jolliffe and Stephenson, 2011). These scores are based on the
225 contingency table between forecasts (or simulated values in the case of the present
226 study) and observations (Table 1).

227

228 POD is the probability of detection (equation 1), FAR is the number of false alarms per
229 the total number of warnings or alarms (equation 2) and HSS is a skill score ranging from
230 $-\infty$ to 1 (equation 3), for categorical forecasts where the proportion of correct measure is
231 scaled with the reference value from correct forecasts due to chance.

232

$$233 \quad POD = a / (a + c) \quad \text{eq.1}$$

234

$$235 \quad FAR = b / (a + b) \quad \text{eq.2}$$

236

$$237 \quad HSS = 2(ad - bc) / (a + b)(b + d) + (a + c)(c + d) \quad \text{eq.3}$$



238

239

240

4. Results

241

242

4.1 Calibration of the SMA model

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

4.2 Regional estimation of the A parameter

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

To estimate the robustness of the method, we applied a split-sample validation into a testing and a training sample. 70% of the grid cells (15636 data points) were selected for



278 training the RF model, and the remaining 30% (6701 data points) for testing. The results
279 are presented for the testing set (Figure 4). The performance in terms of Nash for the
280 SMA model with A estimated by Random Forests or soil map is very similar, with mean
281 Nash equal to 0.86 (median = 0.89) with RF and 0.81 (median = 0.85) with soil maps. The
282 Nash values in validation (testing set) are low, or even negative, only for mountainous
283 ranges, as expected. Overall, the spatial patterns of the Nash coefficients obtained with
284 RF or ESDB are very similar too. There are no significant relationships between model
285 efficiency and the aridity index or the presence of irrigated areas, as identified in the
286 ECOCLIMAP2 land cover database.

287

288 **4.3 Estimation of dry soil conditions**

289

290 A further validation is made for daily soil moisture below the 10th percentile corresponding
291 to dry soil conditions. We computed the Probability of Detection (POD), the False Alarm
292 Ratio (FAR) and the Heidke Skill Score (HSS) summarizing the global efficiency to detect
293 dry periods. For both approaches to estimate A, the mean POD is very high, close to
294 97%, while the FAR is close to 3%. But these average results hide some discrepancy in
295 the different regions (Figure 5 and 6): the efficiency is the highest for the North-Western
296 region, the wettest areas of Spain, while in the South and Central parts of Spain the
297 performance is lower on average. For the wettest parts of the Iberian Peninsula, the POD
298 remains higher than 94% and the FAR lower than 6% and it is the region where the main
299 improvements with RF are observed. On average, the RF estimation method outperforms
300 the approach based on ESDB (Figure 7), with more stable results in terms of HSS since
301 all values obtained with RF are above 0.4 while with ESDB for the grid cells the HSS
302 scores drops to values close to zero.

303

304 **5. Summary and conclusions**

305

306 In this study, a simple model allowing the monitoring of the soil saturation level was
307 regionalized over the entire Iberian Peninsula, taking as a reference the soil moisture
308 simulated by a high-resolution land surface model. Two different regionalization methods
309 have been compared, either the direct estimation of soil water holding capacity from
310 european soil maps or by Random Forests, using covariates such as altitude,
311 temperature, precipitation, evapotranspiration and land cover. Results have shown that
312 the estimation by Random Forest is more robust notably to estimate low soil moisture
313 levels. Despite similar average performance between the two methods, the use of soil
314 maps to set the water holding capacity reveals less stable results in some cases, most
315 probably related to the uncertainties in the pedo-transfer functions used. While these
316 pedo-transfer functions are process-based predictive functions of certain soil properties,
317 Random Forest are not based on physical processes and are tailored to provide the best



318 estimates in a statistical sense. Therefore, they provide a valuable alternative in contexts
319 where high-resolution soil maps are not available since they rely on a set of covariates
320 that can be reliably estimated from global databases, such as satellite or reanalysis
321 products (Funk et al., 2015; Hersbach et al., 2020; Muñoz Sabater, 2020).

322

323 It should be noted that the results presented herein are highly dependent on the quality
324 of land surface simulations, in the absence of dense monitoring networks of in situ soil
325 moisture data, thus these results suffer from the same limitations as LSMs, notably, the
326 lack of human processes (irrigation). However, new remote sensing irrigation estimates
327 are being developed (Massari et al., 2021), as a consequence, once the RF model is
328 trained, irrigation estimations could be added to the precipitation forcing data in order to
329 include the human impacts on soil moisture estimations. The results show that this
330 approach allows us to cheaply extend the value of high resolution LSM simulations to
331 areas where no LSM is implemented (ie. north Africa), as long as the climate conditions
332 belong to the range of values used to train the model, mostly in terms of precipitation and
333 evapotranspiration ranges. Thus, the model train over the Iberian Peninsula could be
334 applied to other similar areas such as North Africa, Italy or Greece. As a perspective,
335 other simulations from countries where high resolution LSM simulations are available,
336 such as France or the USA, could be added to the database in order to expand the
337 coverage over different physiographic and climate contexts. Consequently, the benefits
338 of LSM simulations of soil moisture could be expanded to other areas, provided that
339 suitable forcing datasets are available. Furthermore, if public meteorological and
340 hydrological organizations were to create soil moisture observation networks, cleverly
341 designed to cover the most relevant climates of their countries, this approach could be
342 used to train the model using these observations and then regionalize the results to the
343 rest of the territory, thus, converting an *in-situ* observation dataset into a gridded dataset
344 with a much greater spatial coverage.

345

346

347 **Acknowledgements**

348 This work is a contribution to the HyMeX programme through the HUMID project
349 (CGL2017-85687-R, AEI/FEDER, UE).

350

351

352

353

354

355

356

357



358 References

- 359
360 Boone A (2000) Modélisation des processus hydrologiques dans le schéma de surface ISBA:
361 Inclusion d'un réservoir hydrologique, du gel et modélisation de la neige. PhD thesis, Université
362 Paul Sabatier (Toulouse III), http://www.cnrm.meteo.fr/IMG/pdf/boone_thesis_2000.pdf
363
364 Habets F, Boone A, Noilhan J (2003) Simulation of a Scandinavian basin using the diffusion
365 transfer version of ISBA. *Glob Planet Chang* 38(1-2):137–149
366
367 Masson, V., Le Moigne, P., Martin, E., Faroux, S., Alias, A., Alkama, R., Belamari, S., Barbu, A.,
368 Boone, A., Bouyssel, F., Brousseau, P., Brun, E., Calvet, J. C., Carrer, D., Decharme, B., Delire,
369 C., Donier, S., Essaouini, K., Gibelin, A. L., ... Voltaire, A. (2013). The SURFEXv7.2 land and
370 ocean surface platform for coupled or offline simulation of earth surface variables and fluxes.
371 *Geoscientific Model Development*, 6(4), 929–960. <https://doi.org/10.5194/gmd-6-929-2013>
372
373 Vivien-Georgiana Stefan, Gianfranco Indrio, Maria-José Escorihuela, Pere Quintana-Seguí, and
374 Josep Maria Villar: High-Resolution SMAP-Derived Root-Zone Soil Moisture Using an
375 Exponential Filter Model Calibrated per Land Cover Type, *Remote Sensing* 2021, 13(6).
376 <https://doi.org/10.3390/rs13061112>
377
378 Almendra-Martín, L., Martínez-Fernández, J., González-Zamora, Á., Benito-Verdugo, P., and
379 Herrero-Jiménez, C. M.: Agricultural Drought Trends on the Iberian Peninsula: An Analysis
380 Using Modeled and Reanalysis Soil Moisture Products, *Atmosphere*, 12, 236,
381 <https://doi.org/10.3390/atmos12020236>, 2021.
382
383 Antcil, F., Michel, C., Perrin, C., and Andréassian, V.: A soil moisture index as an auxiliary ANN
384 input for stream flow forecasting, *Journal of Hydrology*, 286, 155–167,
385 <https://doi.org/10.1016/j.jhydrol.2003.09.006>, 2004.
386
387 Barella-Ortiz, A. and Quintana-Seguí, P.: Evaluation of drought representation and propagation
388 in regional climate model simulations across Spain, *Hydrol. Earth Syst. Sci.*, 23, 5111–5131,
389 <https://doi.org/10.5194/hess-23-5111-2019>, 2019.
390
391 Bauer-Marschallinger, B., Freeman, V., Cao, S., Paulik, C., Schaufler, S., Stachl, T., Modanesi,
392 S., Massari, C., Ciabatta, L., Brocca, L., and Wagner, W.: Toward Global Soil Moisture
393 Monitoring With Sentinel-1: Harnessing Assets and Overcoming Obstacles, *IEEE Trans.*
394 *Geosci. Remote Sensing*, 57, 520–539, <https://doi.org/10.1109/TGRS.2018.2858004>, 2019.
395
396 Blöschl, G. and Sivapalan, M.: Scale issues in hydrological modelling: A review, *Hydrol.*
397 *Process.*, 9, 251–290, <https://doi.org/10.1002/hyp.3360090305>, 1995.
398
399 Booker, D. J. and Woods, R. A.: Comparing and combining physically-based and empirically-
400 based approaches for estimating the hydrology of ungauged catchments, *Journal of Hydrology*,
401 508, 227–239, <https://doi.org/10.1016/j.jhydrol.2013.11.007>, 2014.
402
403 Breiman, L.: Bagging predictors, *Mach Learn*, 24, 123–140,
404 <https://doi.org/10.1007/BF00058655>, 1996.
405
406 Breiman, L.: Random Forests, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
407



- 408 Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J.: Classification And Regression
409 Trees, 1st ed., Routledge, <https://doi.org/10.1201/9781315139470>, 2017.
410
- 411 Brocca, L., Camici, S., Melone, F., Moramarco, T., Martínez-Fernández, J., Didon-Lescot, J.-F.,
412 and Morbidelli, R.: Improving the representation of soil moisture by using a semi-analytical
413 infiltration model, *Hydrol. Process.*, 28, 2103–2115, <https://doi.org/10.1002/hyp.9766>, 2014.
414
- 415 Brocca, L., Filippucci, P., Hahn, S., Ciabatta, L., Massari, C., Camici, S., Schüller, L., Bojkov, B.,
416 and Wagner, W.: SM2RAIN–ASCAT (2007–2018): global daily satellite rainfall data from
417 ASCAT soil moisture observations, *Earth Syst. Sci. Data*, 11, 1583–1601,
418 <https://doi.org/10.5194/essd-11-1583-2019>, 2019.
419
- 420 Carranza, C., Nolet, C., Pezij, M., and van der Ploeg, M.: Root zone soil moisture estimation
421 with Random Forest, *Journal of Hydrology*, 593, 125840,
422 <https://doi.org/10.1016/j.jhydrol.2020.125840>, 2021.
423
- 424 Dorigo, W., Wagner, W., Albergel, C., Albrecht, F., Balsamo, G., Brocca, L., Chung, D., Ertl, M.,
425 Forkel, M., Gruber, A., Haas, E., Hamer, P. D., Hirschi, M., Ikonen, J., de Jeu, R., Kidd, R.,
426 Lahoz, W., Liu, Y. Y., Miralles, D., Mistelbauer, T., Nicolai-Shaw, N., Parinussa, R., Pratola, C.,
427 Reimer, C., van der Schalie, R., Seneviratne, S. I., Smolander, T., and Lecomte, P.: ESA CCI
428 Soil Moisture for improved Earth system understanding: State-of-the art and future directions,
429 *Remote Sensing of Environment*, 203, 185–215, <https://doi.org/10.1016/j.rse.2017.07.001>,
430 2017.
431
- 432 Durand, Y., Brun, E., Merindol, L., Guyomarc'h, G., Lesaffre, B., and Martin, E.: A
433 meteorological estimation of relevant parameters for snow models, *A. Glaciology.*, 18, 65–71,
434 <https://doi.org/10.1017/S0260305500011277>, 1993.
435
- 436 Escorihuela, M. J. and Quintana-Seguí, P.: Comparison of remote sensing and simulated soil
437 moisture datasets in Mediterranean landscapes, *Remote Sensing of Environment*, 180, 99–114,
438 <https://doi.org/10.1016/j.rse.2016.02.046>, 2016.
439
- 440 Fang, B., Kansara, P., Dandridge, C., and Lakshmi, V.: Drought monitoring using high spatial
441 resolution soil moisture data over Australia in 2015–2019, *Journal of Hydrology*, 594, 125960,
442 <https://doi.org/10.1016/j.jhydrol.2021.125960>, 2021.
443
- 444 Faroux, S., Kaptué Tchuenté, A. T., Roujean, J.-L., Masson, V., Martin, E., and Le Moigne, P.:
445 ECOCLIMAP-II/Europe: a twofold database of ecosystems and surface parameters at 1 km
446 resolution based on satellite information for use in land surface, meteorological and climate
447 models, *Geosci. Model Dev.*, 6, 563–582, <https://doi.org/10.5194/gmd-6-563-2013>, 2013.
448
- 449 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., Husak, G., Rowland,
450 J., Harrison, L., Hoell, A., and Michaelsen, J.: The climate hazards infrared precipitation with
451 stations—a new environmental record for monitoring extremes, *Sci Data*, 2, 150066,
452 <https://doi.org/10.1038/sdata.2015.66>, 2015.
453
- 454 Gagkas, Z. and Lilly, A.: Downscaling soil hydrological mapping used to predict catchment
455 hydrological response with random forests, *Geoderma*, 341, 216–235,
456 <https://doi.org/10.1016/j.geoderma.2019.01.048>, 2019.
457



- 458 van Genuchten, M. Th.: A Closed-form Equation for Predicting the Hydraulic Conductivity of
459 Unsaturated Soils, *Soil Science Society of America Journal*, 44, 892–898,
460 <https://doi.org/10.2136/sssaj1980.03615995004400050002x>, 1980.
461
- 462 Grillakis, M. G., Koutroulis, A. G., Alexakis, D. D., Polykretis, C., and Daliakopoulos, I. N.:
463 Regionalizing Root-Zone Soil Moisture Estimates From ESA CCI Soil Water Index Using
464 Machine Learning and Information on Soil, Vegetation, and Climate, *Water Res*, 57,
465 <https://doi.org/10.1029/2020WR029249>, 2021.
466
- 467 Habets, F., Boone, A., Champeaux, J. L., Etchevers, P., Franchistéguy, L., Leblois, E., Ledoux,
468 E., Le Moigne, P., Martin, E., Morel, S., Noilhan, J., Quintana Seguí, P., Rousset-Regimbeau,
469 F., and Viennot, P.: The SAFRAN-ISBA-MODCOU hydrometeorological model applied over
470 France, *J. Geophys. Res.*, 113, D06113, <https://doi.org/10.1029/2007JD008548>, 2008.
471
- 472 He, Y., Bárdossy, A., and Zehe, E.: A review of regionalisation for continuous streamflow
473 simulation, *Hydrol. Earth Syst. Sci.*, 15, 3539–3553, <https://doi.org/10.5194/hess-15-3539-2011>,
474 2011.
475
- 476 Hengl, T., Nussbaum, M., Wright, M. N., Heuvelink, G. B. M., and Gräler, B.: Random forest as
477 a generic framework for predictive modeling of spatial and spatio-temporal variables, 6, e5518,
478 <https://doi.org/10.7717/peerj.5518>, 2018.
479
- 480 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J.,
481 Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo,
482 G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., Chiara, G., Dahlgren, P., Dee, D.,
483 Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L.,
484 Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C.,
485 Radnoti, G., Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., and Thépaut, J.: The ERA5
486 global reanalysis, *Q.J.R. Meteorol. Soc.*, 146, 1999–2049, <https://doi.org/10.1002/qj.3803>, 2020.
487 Hiederer, R.: Mapping soil properties for Europe: spatial representation of soil database
488 attributes., Publications Office, LU, 2013.
489
- 490 Hrachowitz, M., Savenije, H. H. G., Blöschl, G., McDonnell, J. J., Sivapalan, M., Pomeroy, J. W.,
491 Arheimer, B., Blume, T., Clark, M. P., Ehret, U., Fenicia, F., Freer, J. E., Gelfan, A., Gupta, H.
492 V., Hughes, D. A., Hut, R. W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P. A., Uhlenbrook,
493 S., Wagener, T., Winsemius, H. C., Woods, R. A., Zehe, E., and Cudennec, C.: A decade of
494 Predictions in Ungauged Basins (PUB)—a review, *Hydrological Sciences Journal*, 58, 1198–
495 1255, <https://doi.org/10.1080/02626667.2013.803183>, 2013.
496
- 497 Javelle, P., Fouchier, C., Arnaud, P., and Lavabre, J.: Flash flood warning at ungauged
498 locations using radar rainfall and antecedent soil moisture estimations, *Journal of Hydrology*,
499 394, 267–274, <https://doi.org/10.1016/j.jhydrol.2010.03.032>, 2010.
500
- 501 Jolliffe, I. T. and Stephenson, D. B. (Eds.): *Forecast Verification: A Practitioner’s Guide in*
502 *Atmospheric Science*, John Wiley & Sons, Ltd, Chichester, UK,
503 <https://doi.org/10.1002/9781119960003>, 2011.
504
- 505 Li, J., Wang, Z., Wu, X., Xu, C.-Y., Guo, S., and Chen, X.: Toward Monitoring Short-Term
506 Droughts Using a Novel Daily Scale, Standardized Antecedent Precipitation Evapotranspiration
507 Index, 21, 891–908, <https://doi.org/10.1175/JHM-D-19-0298.1>, 2020.



- 508
509 Loh, W. Y. and Shih, Y. S.: Split Selection Methods for Classification Trees, 7, 815–840, 1997.
510 Martínez-Fernández, J., González-Zamora, A., Sánchez, N., and Gumuzzio, A.: A soil water
511 based index as a suitable agricultural drought indicator, *Journal of Hydrology*, 522, 265–273,
512 <https://doi.org/10.1016/j.jhydrol.2014.12.051>, 2015.
- 513
514 Martínez-Fernández, J., González-Zamora, A., Sánchez, N., Gumuzzio, A., and Herrero-
515 Jiménez, C. M.: Satellite soil moisture for agricultural drought monitoring: Assessment of the
516 SMOS derived Soil Water Deficit Index, *Remote Sensing of Environment*, 177, 277–286,
517 <https://doi.org/10.1016/j.rse.2016.02.064>, 2016.
- 518
519 Massari, C., Modanesi, S., Dari, J., Gruber, A., De Lannoy, G. J. M., Girotto, M., Quintana-
520 Seguí, P., Le Page, M., Jarlan, L., Zribi, M., Ouadi, N., Vreugdenhil, M., Zappa, L., Dorigo, W.,
521 Wagner, W., Brombacher, J., Pelgrum, H., Jaquot, P., Freeman, V., Volden, E., Fernandez
522 Prieto, D., Tarpanelli, A., Barbetta, S., and Brocca, L.: A Review of Irrigation Information
523 Retrievals from Space and Their Utility for Users, *Remote Sensing*, 13, 4112,
524 <https://doi.org/10.3390/rs13204112>, 2021.
- 525
526 McKay, M. D., Beckman, R. J., and Conover, W. J.: Comparison of Three Methods for Selecting
527 Values of Input Variables in the Analysis of Output from a Computer Code, *Technometrics*, 21,
528 239–245, <https://doi.org/10.1080/00401706.1979.10489755>, 1979.
- 529
530 Merlin, O., Escorihuela, M. J., Mayoral, M. A., Hagolle, O., Al Bitar, A., and Kerr, Y.: Self-
531 calibrated evaporation-based disaggregation of SMOS soil moisture: An evaluation study at 3
532 km and 100 m resolution in Catalunya, Spain, *Remote Sensing of Environment*, 130, 25–38,
533 <https://doi.org/10.1016/j.rse.2012.11.008>, 2013.
- 534
535 Mishra, A., Vu, T., Veetil, A. V., and Entekhabi, D.: Drought monitoring with soil moisture active
536 passive (SMAP) measurements, *Journal of Hydrology*, 552, 620–632,
537 <https://doi.org/10.1016/j.jhydrol.2017.07.033>, 2017.
- 538
539 Muñoz Sabater, J.: ERA5-Land hourly data from 1981 to present. Copernicus Climate Change
540 Service (C3S) Climate Data Store (CDS), 10.24381/cds.e2161bac, 2020.
- 541
542 Noguera, I., Domínguez-Castro, F., and Vicente-Serrano, S. M.: Flash Drought Response to
543 Precipitation and Atmospheric Evaporative Demand in Spain, *Atmosphere*, 12, 165,
544 <https://doi.org/10.3390/atmos12020165>, 2021.
- 545
546 Noilhan, J. and Mahfouf, J.-F.: The ISBA land surface parameterisation scheme, *Global and
547 Planetary Change*, 13, 145–159, [https://doi.org/10.1016/0921-8181\(95\)00043-7](https://doi.org/10.1016/0921-8181(95)00043-7), 1996.
- 548
549 Panagos, P., Van Liedekerke, M., Jones, A., and Montanarella, L.: European Soil Data Centre:
550 Response to European policy support and public data requirements, *Land Use Policy*, 29, 329–
551 338, <https://doi.org/10.1016/j.landusepol.2011.07.003>, 2012.
- 552
553 Pena-Gallardo, M., Vicente-Serrano, S. M., Domínguez-Castro, F., and Beguería, S.: The
554 impact of drought on the productivity of two rainfed crops in Spain, *Nat. Hazards Earth Syst.
555 Sci.*, 19, 1215–1234, <https://doi.org/10.5194/nhess-19-1215-2019>, 2019.
- 556
557 Perrin, C., Michel, C., and Andréassian, V.: Improvement of a parsimonious model for



- 558 streamflow simulation, *Journal of Hydrology*, 279, 275–289, <https://doi.org/10.1016/S0022->
559 1694(03)00225-7, 2003.
- 560
- 561 Piedallu, C., Gégout, J.-C., Perez, V., and Lebourgeois, F.: Soil water balance performs better
562 than climatic water variables in tree species distribution modelling: Soil water balance improves
563 tree species distribution models, *Global Ecology and Biogeography*, 22, 470–482,
564 <https://doi.org/10.1111/geb.12012>, 2013.
- 565
- 566 Quintana-Seguí, P., Le Moigne, P., Durand, Y., Martin, E., Habets, F., Baillon, M., Canellas, C.,
567 Franchisteguy, L., and Morel, S.: Analysis of Near-Surface Atmospheric Variables: Validation of
568 the SAFRAN Analysis over France, 47, 92–107, <https://doi.org/10.1175/2007JAMC1636.1>,
569 2008.
- 570
- 571 Quintana-Seguí, P., Turco, M., Herrera, S., and Miguez-Macho, G.: Validation of a new
572 SAFRAN-based gridded precipitation product for Spain and comparisons to Spain02 and ERA-
573 Interim, *Hydrol. Earth Syst. Sci.*, 21, 2187–2201, <https://doi.org/10.5194/hess-21-2187-2017>,
574 2017.
- 575
- 576 Quintana-Seguí, P., Barella-Ortiz, A., Regueiro-Sanfiz, S., and Miguez-Macho, G.: The Utility of
577 Land-Surface Model Simulations to Provide Drought Information in a Water Management
578 Context Using Global and Local Forcing Datasets, *Water Resour Manage.*,
579 <https://doi.org/10.1007/s11269-018-2160-9>, 2019.
- 580
- 581 Raymond, F., Ullmann, A., Tramblay, Y., Drobinski, P., and Camberlin, P.: Evolution of
582 Mediterranean extreme dry spells during the wet season under climate change, *Reg Environ*
583 *Change*, 19, 2339–2351, <https://doi.org/10.1007/s10113-019-01526-3>, 2019.
- 584
- 585 Reynolds, C. A., Jackson, T. J., and Rawls, W. J.: Estimating soil water-holding capacities by
586 linking the Food and Agriculture Organization Soil map of the world with global pedon
587 databases and continuous pedotransfer functions, *Water Resour. Res.*, 36, 3653–3662,
588 <https://doi.org/10.1029/2000WR900130>, 2000.
- 589
- 590 Rodell, M., Houser, P. R., Jambor, U., Gottschalck, J., Mitchell, K., Meng, C.-J., Arsenault, K.,
591 Cosgrove, B., Radakovich, J., Bosilovich, M., Entin, J. K., Walker, J. P., Lohmann, D., and Toll,
592 D.: The Global Land Data Assimilation System, *Bull. Amer. Meteor. Soc.*, 85, 381–394,
593 <https://doi.org/10.1175/BAMS-85-3-381>, 2004.
- 594
- 595 Stein, L., Clark, M. P., Knoben, W. J. M., Pianosi, F., and Woods, R. A.: How Do Climate and
596 Catchment Attributes Influence Flood Generating Processes? A Large-Sample Study for 671
597 Catchments Across the Contiguous USA, *Water Res*, 57,
598 <https://doi.org/10.1029/2020WR028300>, 2021.
- 599
- 600 Tramblay, Y., Bouaicha, R., Brocca, L., Dorigo, W., Bouvier, C., Camici, S., and Servat, E.:
601 Estimation of antecedent wetness conditions for flood modelling in northern Morocco, *Hydrol.*
602 *Earth Syst. Sci.*, 16, 4375–4386, <https://doi.org/10.5194/hess-16-4375-2012>, 2012.
- 603
- 604 Tramblay, Y., Amoussou, E., Dorigo, W., and Mahé, G.: Flood risk under future climate in data
605 sparse regions: Linking extreme value models and flood generating processes, *Journal of*
606 *Hydrology*, 519, 549–558, <https://doi.org/10.1016/j.jhydrol.2014.07.052>, 2014.
- 607



- 608 Tramblay, Y., Koutroulis, A., Samaniego, L., Vicente-Serrano, S. M., Volaire, F., Boone, A., Le
609 Page, M., Llasat, M. C., Albergel, C., Burak, S., Cailleret, M., Kalin, K. C., Davi, H., Dupuy, J.-L.,
610 Greve, P., Grillakis, M., Hanich, L., Jarlan, L., Martin-StPaul, N., Martínez-Vilalta, J., Mouillot, F.,
611 Pulido-Velazquez, D., Quintana-Seguí, P., Renard, D., Turco, M., Türkeş, M., Trigo, R., Vidal,
612 J.-P., Vilagrosa, A., Zribi, M., and Polcher, J.: Challenges for drought assessment in the
613 Mediterranean region under future climate scenarios, *Earth-Science Reviews*, 210, 103348,
614 <https://doi.org/10.1016/j.earscirev.2020.103348>, 2020.
- 615
- 616 Tyralis, H., Papacharalampous, G., and Langousis, A.: A Brief Review of Random Forests for
617 Water Scientists and Practitioners and Their Recent History in Water Resources, *Water*, 11,
618 910, <https://doi.org/10.3390/w11050910>, 2019.
- 619
- 620 Vicente-Serrano, S. M., Lopez-Moreno, J.-I., Beguería, S., Lorenzo-Lacruz, J., Sanchez-
621 Lorenzo, A., García-Ruiz, J. M., Azorin-Molina, C., Morán-Tejeda, E., Revuelto, J., Trigo, R.,
622 Coelho, F., and Espejo, F.: Evidence of increasing drought severity caused by temperature rise
623 in southern Europe, *Environ. Res. Lett.*, 9, 044001, <https://doi.org/10.1088/1748-9326/9/4/044001>, 2014.
- 624
- 625
- 626 Willgoose, G. and Perera, H.: A simple model of saturation excess runoff generation based on
627 geomorphology, steady state soil moisture, *Water Resour. Res.*, 37, 147–155,
628 <https://doi.org/10.1029/2000WR900265>, 2001.
- 629
- 630 Wösten, J. H. M., Lilly, A., Nemes, A., and Le Bas, C.: Development and use of a database of
631 hydraulic properties of European soils, *Geoderma*, 90, 169–185, [https://doi.org/10.1016/S0016-7061\(98\)00132-3](https://doi.org/10.1016/S0016-7061(98)00132-3), 1999.
- 632
- 633
- 634 Zhao, B., Dai, Q., Han, D., Dai, H., Mao, J., Zhuo, L., and Rong, G.: Estimation of soil moisture
635 using modified antecedent precipitation index with application in landslide predictions,
636 *Landslides*, 16, 2381–2393, <https://doi.org/10.1007/s10346-019-01255-y>, 2019.
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645



646 **TABLE**

647

648 Table 1: Contingency table of the comparison between forecasts and observations or
649 any two analyses. The symbols a–d are the different numbers of cases observed to
650 occur in each category.

651

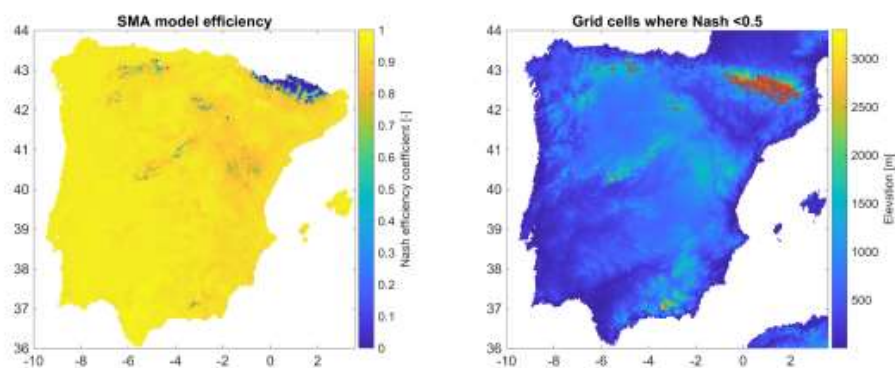
	Observations	
Forecast	1	0
1	a (hit)	b (false alarm)
0	c (miss)	d (correct rejection)

652

653 **FIGURES**

654

655



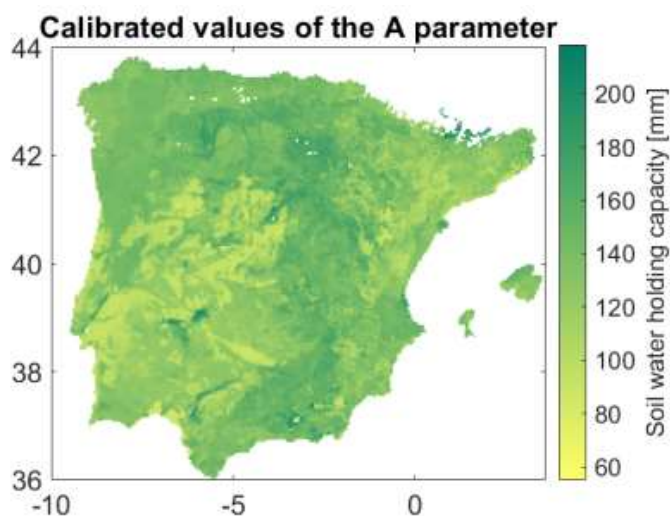
656

657 Figure 1: Efficiency of the SMA model to reproduce soil moisture from SURFEX

658

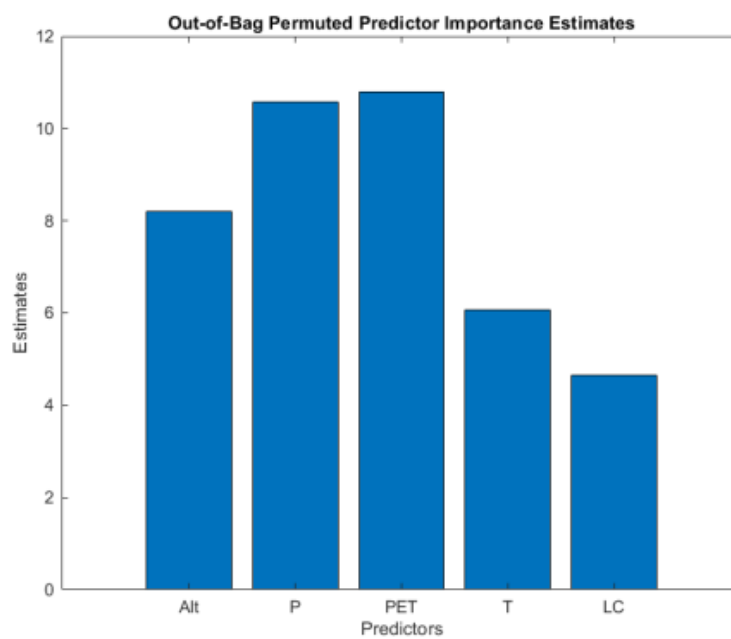
659

660



661
662
663
664

Figure 2: Map of the calibrated values of the A parameter of the SMA model

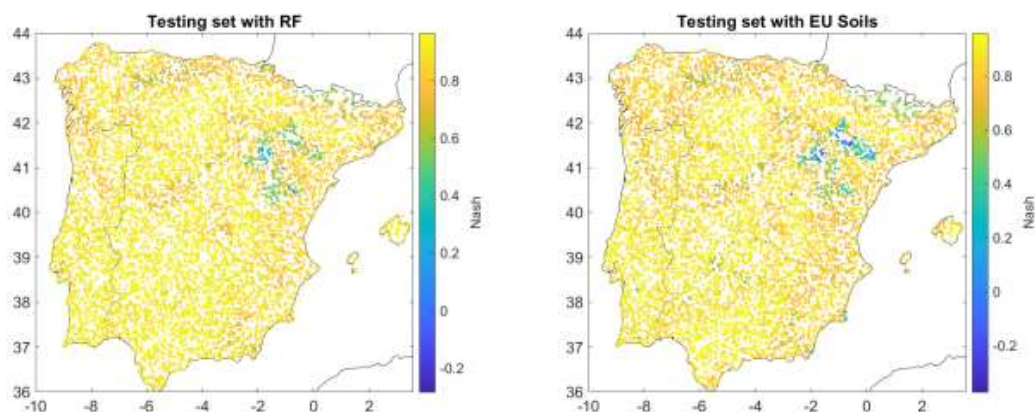


665
666
667
668
669

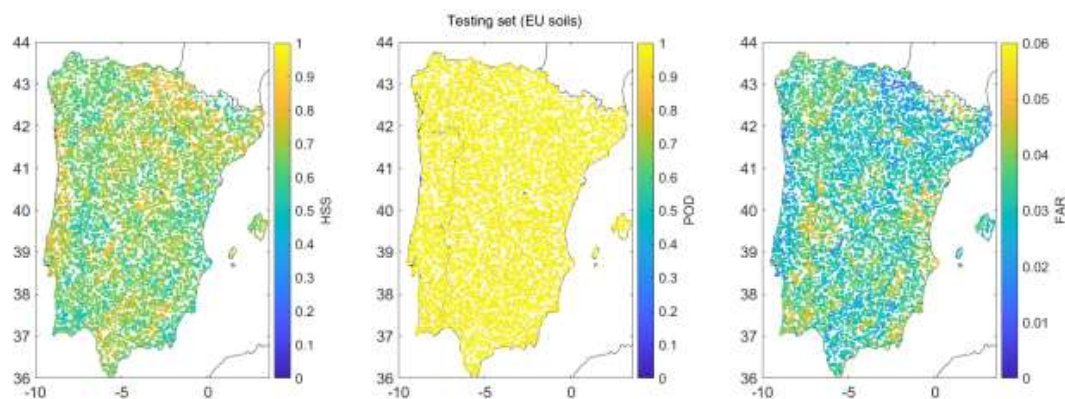
Figure 3: Relative importance of each predictor (Alt= altitude, P= precipitation, PET= potential evapotranspiration, T=temperature, LC=land cover classes) in the Random Forest method



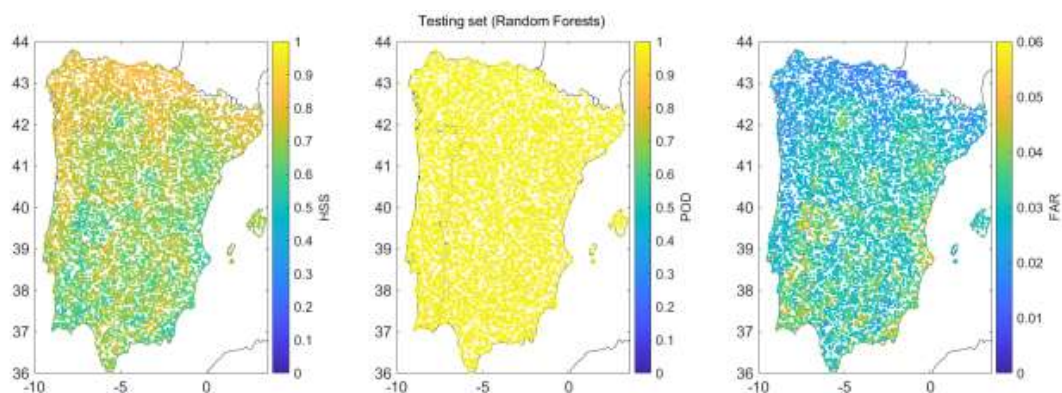
670
671



672
673 Figure 4: Nash efficiency coefficient obtained for the testing set, with the A parameter of
674 the SMA model estimated by RF (left) or ESDB (right)
675
676

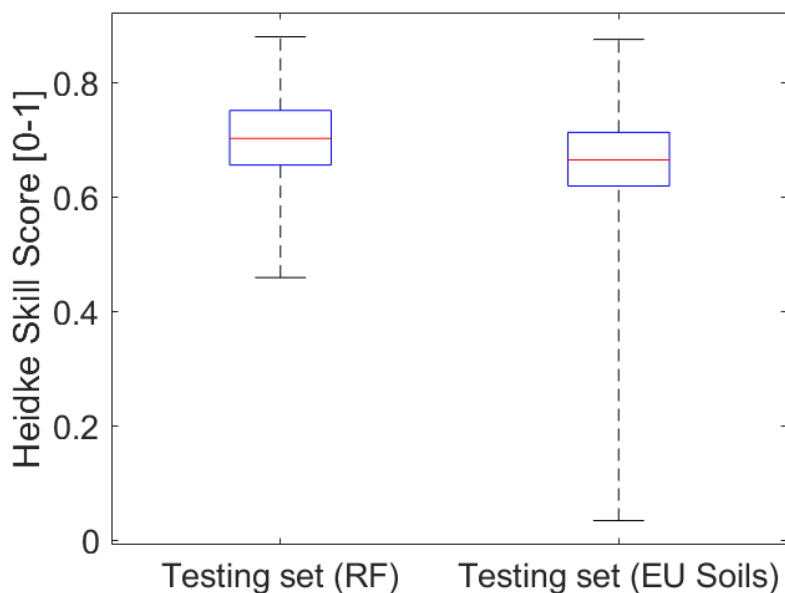


677
678 Figure 5: Validation results in terms of HSS, POD and FAR with A estimated with ESDB
679



680
681
682

Figure 6: same as figure 5 but with A estimated with RF



683
684
685
686
687
688
689
690
691
692
693

Figure 7: Boxplot of the HSS obtained with RF or EU soil maps. The limits of the box represent the 25th and 75 percentiles, the line in the middle refers to the median, and the limits of the whiskers extend to the minimum and maximum values.