

1 **A high resolution map of soil types and physical properties for**
2 **Cyprus: a digital soil mapping optimization**

3 Corrado Camera¹, Zomenia Zomeni², Jay S. Noller³, Andreas M. Zissimos², Irene C.
4 Christoforou², Adriana Bruggeman¹

5 ¹ *Energy, Environment and Water Research Centre, The Cyprus Institute, Aglantzia, Lefkosia, 2121, Cyprus*

6 ² *Cyprus Geological Survey, Ministry of Agriculture, Rural Development and the Environment, Strovolos,*
7 *Lefkosia, 2064, Cyprus*

8 ³ *Department of Crop and Soil Science, Oregon State University, Corvallis, Oregon, 97331, United States*

9
10 *Corresponding author: Corrado Camera*

11 *Email: c.camera@cyi.ac.cy*

12 *Telephone: +357 22208691*

13 *Fax: +357 22208625*

14
15 **Highlights**

- 16 • We created a digital soil map of Cyprus (1:25,000) using Random Forest
17 • We developed a reliability index to qualitatively show prediction uncertainties
18 • The model showed greater sensitivity to the number of input points than trees grown
19 • Soil geochemistry predictors had a prevalent role in identifying soil properties

20
21 Accepted manuscript version of the paper published in Geoderma

22 (<https://doi.org/10.1016/j.geoderma.2016.09.019>)

A high resolution map of soil types and physical properties for Cyprus: a digital soil mapping optimization

1 Introduction

Environmental factors such as climate, organisms, relief, parent material, and time (*clorpt*) drive soil genesis (Jenny, 1941). Following this hypothesis, traditional soil survey maps are developed based on an empirical model (the soil-landscape model) derived from inductive reasoning from field and laboratory data, which represent the interchangeable relationships between soils and environmental factors. Digital soil mapping (DSM) techniques (McBratney et al., 2003) are based on the same hypotheses and aim at predicting soil types and properties linking field soil observations to environmental predictors. In DSM, inductive reasoning for developing the relationships among the soil-landscape model factors is replaced by different machine learning techniques (i.e., decision tree, fuzzy logic, neural network etc.) (Lagacherie et al., 1995; Scull et al., 2003; Lagacherie et al., 2007; Grinand et al., 2008; Heung et al., 2016). However, pedological expert knowledge remains a key factor in model building to ensure both statistically and pedologically sound outputs (Kempen et al., 2009).

Digital soil mapping as a discipline has experienced a continuous expansion in the last two decades, mainly due to its increased efficiency in comparison to conventional field soil mapping techniques (Kempen et al., 2012). Reasons are the ever growing computational capacities coupled with the development of data-mining algorithm and GIS tools, and the increased availability of spatial remote-sensing data (Minasny and McBratney, 2016). Due to their numerical nature, digital soil maps also allow handling continuous spatial variations of soils, for example through class membership values, as presented by Burrough et al. (1997). This overcomes the problem of soil spatial patterns being traditionally captured and displayed as choropleth maps with discrete lines representing the boundaries between soil map units,

1 which implies homogeneity within map units (Burrough, 1986; Bolstad et al., 1990). Digital
2 soil mapping expands the notion of the soil-forming equation to that of a soil-mapping
3 equation, the *scorpan* equation, which adds preexisting soil information and spatial location
4 into Jenny's *clorpt* equation.

5 Random Forest (RF) is a fairly recent data mining algorithm (Breiman, 2001) that has been
6 increasingly used for digital soil mapping applications in recent years. Its success is related to
7 several advantages over other statistical (e.g. linear regression or generalized linear models),
8 geo-statistical (e.g. regression or co-kriging), and machine learning (e.g. neural networks,
9 logistic regression, support vector machines, classification trees) techniques. These
10 advantages have been summarized by Grimm et al. (2008): ability of modeling high
11 dimensional non-linear relationships; simultaneous handling of categorical and continuous
12 predictors; robustness against over-fitting; measures of error rate and variable importance;
13 requirement of only three user-defined input parameters; and relatively low sensitivity to
14 parameter values. In particular, the measure of variable importance has proved to be, in many
15 circumstances, a useful tool for enlightening soil-environment relationships to allow authors
16 to infer the effects of possible future environmental changes on soil characteristics (e.g.
17 Barthold et al., 2013).

18 Among its many applications, RF has been used for predicting the spatial distribution of
19 various soil properties, such as soil organic and/or inorganic carbon (e.g. Grimm et al., 2008;
20 Wiesmeier et al., 2011; Poggio et al., 2013; Akpa et al., 2016; Sreenivas et al., 2016); soil
21 texture and cation--exchange capacity (e.g. Lagacherie et al., 2013; Chagas et al., 2016); and
22 soil taxonomic units in unmapped areas (Stum et al., 2010; Barthold et al., 2013; Pahlavan
23 Rad et al., 2014; Brungard et al., 2015; Taghizadeh-Mehrjardi et al., 2015; Heung et al.,
24 2016; Láng et al., 2016). Heung et al. (2014) pointed out that few studies have applied RF for

1 mapping categorical soil properties, specifically referring to soil taxonomic units. In the last
2 3-4 years this gap has started to get filled.

3 Recent soil classification studies mainly deal with the comparison of the performance of
4 many different algorithms (including RF) and sampling techniques (e.g. Brungard et al.,
5 2015; Taghizadeh-Mehrjardi et al., 2015; Heung et al., 2016), partially disregarding model
6 building and model optimization. Models often include continuous variables representing
7 topography, climate, vegetation or land use from remote sensing products (e.g. Pahlavan Rad
8 et al., 2014; Brungard et al., 2015; Heung et al., 2016), sometimes categorical variables
9 representing parent material (e.g. Berthold et al., 2013; Taghizadeh-Mehrjardi et al., 2015),
10 but only Taghizadeh-Mehrjardi et al. (2015) and Láng et al. (2016) also consider soil
11 information and properties. However, none of these authors clearly quantifies the role and the
12 importance of these different predictors in the model. An optimization of the number of trees
13 in the forest and of the number of variables to be used to split branches is quite typical for RF
14 (e.g. Grimm et al., 2008; Barthold et al., 2013; Heung et al., 2016). Conversely, the
15 investigation of the effect of tree pruning, which is common for classification tree and
16 boosted classification tree modelling approaches (Scull et al., 2005; Schmidt et al., 2008;
17 Lemercier et al., 2012), has not been extensively reported in RF soil mapping applications.
18 Finally, Barthold et al. (2013) and Taghizadeh-Mehrjardi et al. (2015) are the only authors
19 explicitly relating soil groups to major soil properties such as soil depth and soil texture,
20 although they do not derive each of them independently. These properties are of major
21 importance for application studies such as agricultural crop modelling and soil erosion (Bird
22 et al., 2016; Djuma et al., under review). There is also a paucity of studies dealing with soil
23 prediction in complex topographical and pedological environments and in the Eastern
24 Mediterranean region.

25

1 The main aim of this study is to develop digital soil maps of the soil groups, depth and
2 texture classes of a topographical and pedological complex area of the eastern Mediterranean,
3 namely the island of Cyprus, based on extensive soil legacy data. Specific objectives are: (i)
4 to analyze the role and importance of a large data set of environmental predictors, covering
5 all the soil formation factors considered in the *scorpan* formula, both for single and groups of
6 predictors; (ii) to investigate the effect of number of training points, forest size (number of
7 trees), number of predictors sampled at each node, and tree size (terminal node) in RF; (iii) to
8 compare RF-derived maps with maps derived with a Multinomial Logistic Regression model,
9 in terms of validation error and map uncertainty, using the confusion index and a newly
10 developed reliability index.

11 **2 Materials and Methods**

12 **2.1 Study area**

13 Cyprus is the third largest island in the Mediterranean and is located between 34-36°N and
14 32-35°E. The main physical characteristics of the island are represented by the two mountain
15 chains, the Troodos, located in the central-west part with the highest peak at Mount Olympus
16 (1951 m a.s.l.), and the Pentadaktylos Range along the north coast with its highest peak at
17 Mount Kyparissovouno (1,024 m a.s.l.). The main agricultural area of the country is the
18 Mesaoria Plain, which lies in between the two mountain ranges and the coastal lowlands.

19 Soils on Cyprus are exceptional due to the geological complexity of the island, the
20 Mediterranean climate and the long presence of man on the landscape. The Troodos
21 Ophiolite, a fragment of fully developed oceanic crust, consisting of Turonian plutonic,
22 intrusive and volcanic rocks and chemical sediments dominates the central topographic high
23 of the island. Older allochthonous rocks are juxtaposed in the southwest (Mamonia Terrane,
24 Middle Triassic – Middle Cretaceous) and the long east-west Pentadaktylos range in the north

1 coast (Keryneia Terrane, Carboniferous – Middle Miocene). Autochthonous carbonate
2 sediments cover the slopes and plains. Quaternary deposits are predominately of gravity and
3 fluvial origins inland and of marine and aeolian origins on the coast. The soils on Cyprus
4 vary between leptosols, regosols, solonchaks, solonetz, vertisols, luvisols, fluvisols, and
5 cambisols based on the World Reference Base of the FAO (Food and Agriculture
6 Organization of the United Nations) soil classification system (IUSS, 2015). They are
7 generally poor in organic matter (Koudounas and Makin, 1978; Grivas, 1988) and closely
8 associated to parent material and landscape position (Zomeni, 2012; Zomeni and Bruggeman,
9 2013). Thin (leptic) and stony (lithic) soils dominate the mountainous areas developing
10 mostly as residuum. Other soils form on transported materials such as alluvial deposits
11 (alluvial fans, fluvial terraces and deltas), colluvial deposits, aeolian deposits, marine deposits
12 (sands and gravels) and lake and estuarine deposits (hydromorphic silts and clays).

13 The geochemistry of the island also reflects the geological complexity and the impact of
14 humans. A recent high sampling density (5,350 sites on a 1 km² grid), multi-element (60
15 elements) and multi-method analysis soil geochemical survey has resulted in the compilation
16 of the Geochemical Atlas of Cyprus (Cohen et al., 2011, 2012a). The survey was carried out
17 at two depths. Surface soil samples were collected at a depth of 0-20 cm and bottom samples
18 at a depth of 50-70 cm. The survey has demonstrated that chemical processes and element
19 concentrations are dominated by parent lithology. Other processes such as the physical
20 concentration of heavy minerals (Ren et al., 2015), ocean influences along the coastal plains,
21 and human activities also affect the spatial geochemical patterns of the soils on the island
22 (Cohen et al., 2012b; Zissimos et al., 2014). For the purpose of this study we have calculated
23 geochemical parameters using data from surface soil samples.

1 The present study covers the areas under the effective control of the government of the
2 Republic of Cyprus, with the addition of the UN buffer zone (tot. 5,979 km²), where data are
3 available (Fig. 1).

4

5 **Fig. 1. The island of Cyprus with its main physical characteristics and the location of the study area.**

6

7 **2.2 Soil data**

8 The most detailed soil references on the island are ten 1:25,000 scale soil sheets prepared
9 between 1967 and 1985 by the Soil Section of the Department of Agriculture, using
10 traditional survey methods (Soteriades and Georgiades, 1967, Soteriades and Grivas, 1968,
11 Soteriades et al. 1968, Soteriades and Markides, 1969, Grivas and Georgiades, 1972,
12 Markides, 1975, Koumis, 1980a, Koumis, 1980b, Koumis, 1980c, Markides, 1985a). The
13 soils were mapped and classified based on their development stage, origin and parent
14 material. These sheets are always accompanied by an agricultural land suitability map and
15 two of them - the Pafos sheet (Soteriades and Koudounas, 1968) and the Polemi sheet
16 (Markides, 1973) - have an extensive soil memoir.

17 The 10 sheets cover around 1,600 km², classified in 369 soil sub-series (52 soil series with
18 local soil names). Based on soil-profile descriptions, provided in the legend of the map sheets
19 or in the available soil memoirs, we independently harmonized soil (sub)series, soil depth,
20 and soil texture. The harmonization of the soil series was based on the World Reference Base
21 (WRB) for soil resources (IUSS, 2015) and led to recognize 31 soil classes (soil groups
22 accompanied by one or two qualifiers) and two miscellaneous classes, namely Water Bodies
23 and No Data (e.g. quarries and excavated areas). Soil depth was harmonized in four classes,
24 based on the depth ranges available in the original maps (Table 1). Soil texture (Table 2) was

1 harmonized in four classes (None, Coarse, Medium, Fine) following the guidelines provided
2 by FAO (2008). As only differences, we grouped all sandy loams as coarse textured and all
3 clay loams as medium textured. In FAO (2008) these two classes were split between coarse-
4 medium and medium-fine, based on the relative abundance of sand and clay, which
5 information is rarely available in our legacy data. The soil profiles usually included: depth
6 class; texture class; color; lithic properties and their quantity and composition; bedrock
7 geology; horizon descriptions; and chemical properties (the latter two available only in the
8 two soil memoirs). The soil classes, soil depth and soil texture maps were digitised and
9 converted to raster format with cell size of 25x25 m² (2,561,849 cells in total) and used as
10 training data for DSM purposes. The training areas, shown in Fig. 2, cover 17% of the island
11 and 27% of the study area.

12 Additional mapping activities were carried for different sheets (Morphou, Polis-Tylliria,
13 Krasochoria, Pissouri-Paramali, Limassol, and Larnaka) of the Land Suitability Map of
14 Cyprus (Grivas, 1969; Markides, 1969; 1985b; Koumis, 1970a; 1970b; 1970c). These studies
15 provided 126 profiles, described in terms of both soil depth and soil texture, outside the
16 training areas that we used as independent validation observations (see Section 2.4.4 - Map
17 extrapolation, prediction uncertainty, map reliability, and independent profile evaluation).
18 Additional 199 validation profiles, reporting data only on soil depth, were derived from
19 Robins (2004). The author assessed the soil depth variability of the Troodos Mountains at
20 road cuts in three valleys on the north slope of the mountain chain. The location of the
21 validation profiles are presented in Fig. 2.

22

23 **Fig. 2. Location of the study area, the existing 1:25,000 scale soil sheets used for training, the 126 soil**
24 **profiles derived from Land Suitability Maps (LSM), and the 199 soil profiles from Robins (2004) used to**
25 **evaluate the digital soil map of Cyprus.**

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Since Water Bodies and No Data are not proper soils and represent mainly human modification of the environment, we superimposed these data layers (of known location) on the map resulting from the applied computational methods (see section 2.4). For both Water Bodies and No Data, soil depth and soil texture were considered not applicable (N/A).

Table 1. Soil depth classes identified from the 1:25,000 scale soil sheets and reclassified in four consistent classes.

Table 2. Soil texture classes identified from the 1:25,000 scale soil sheets and reclassified in four consistent classes. The harmonization has been performed using the guidelines provided in FAO (2008).

2.3 Environmental covariates

A general summary of the methodology applied in this study is presented in this and the following sections (Fig. 3). For the present study, predictors are environmental covariates that have been selected based on the *scorpan* formula (McBratney et al., 2003). The selected 20 variables (17 continuous, 3 categorical) are described in the following paragraphs. A correlation analysis was performed on the 17 continuous predictors to reduce the dimensionality of the input dataset by removing redundant features. Building a parsimonious model brings two different advantages, as explained by Behrens et al. (2010). First, it helps recognize soil formation processes more clearly, and second it usually yields more reliable predictions.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Fig. 3. Summary of the methods used in the study. OOB is the Out Of Bag error, MESS is the Multivariate Environmental Similarity Surface, and CI is the Confusion Index.

2.3.1 Soil chemistry and organism covariates

To describe soil chemistry we selected Electrical Conductivity (EC), pH, Total Carbon (TotC), Organic Carbon (OrgC), Loss on Ignition (LOI), Mafic Index of Alteration (MIA), Chemical Index of Alteration (CIA), and Vogt Ratio (VR). All these covariates were obtained from the Geochemical Atlas of Cyprus (Cohen et al., 2011; 2012a) by using the raw data, which were collected on a regular 1 km² grid, and re-interpolating them by Inverse Distance Weighting (IDW) over the 25 x 25 m² grid. Both TotC and OrgC were measured using an Eltra CS-800 automatic Carbon–Sulfur analyser, while LOI was measured on soil sub samples at 1000° C. The Mafic Index of Alteration (MIA) was selected to account for the chemical changes during weathering of mafic lithologies (Babechuk et al., 2014), while the CIA and the VR are ratios that quantify the loss of mobile major elements relative to immobile elements during weathering. The three indices were calculated as follows:

$$MIA = 100 \cdot [(Al_2O_3) / (Al_2O_3 + Fe_2O_3 + MgO + CaO + Na_2O + K_2O)] \tag{1}$$

$$CIA = Al_2O_3 / (Al_2O_3 + Na_2O + K_2O + CaO) \tag{2}$$

$$VR = (Al_2O_3 + K_2O) / (MgO + CaO + Na_2O) \tag{3}$$

For further information on the analytical methods used to obtain the geochemical data as well as the spatial variability of the used soil chemistry covariates see Cohen et al. (2011, 2012a) and Zissimos et al. (2014).

1 OrgC data can be considered a proxy for the presence of soil biota like micro-organisms and
2 plants (Yeomans 1988; Soon 1991; Panagos et al 2008), In addition, we considered land use
3 as an organism covariate. We derived this parameter from the CORINE land cover data set of
4 2006, which has a 1:250,000 scale (Büttner and Kosztra, 2007). We use categories from the
5 Level 1 of the database except for “Forest and semi-natural areas”, which have been
6 differentiated to Level 2.

7 2.3.2 Climate covariates

8 As climate covariates, we included minimum, maximum and average temperature (T_{\min} , T_{\max} ,
9 T_{ave}). We calculated these parameters over each 25 x 25 m² cell from the daily gridded
10 temperature data set developed by Camera et al. (2013). The data set has a horizontal
11 resolution of 1 km and covers the period 1980-2010. Therefore, the smaller cells of the
12 present study were assigned the value of the coarser cell by nearest neighbour resampling.
13 T_{\min} and T_{\max} were calculated as the mean of the daily minimum and maximum temperatures
14 recorded in January and July, respectively. They represent the influence of the hottest and
15 coldest temperatures on soil genesis (Scull et al., 2005). T_{ave} is the mean annual average
16 temperature.

17 2.3.3 Relief covariates

18 Relief covariates included: elevation from a Digital Elevation Model (DEM); aspect;
19 curvature; planar curvature; profile curvature; slope; and landscape units. The DEM has a
20 horizontal resolution of 25 x 25 m² and was created from digitised contour and point data of
21 published 1:50,000 scale topographical maps using the Topogrid command in ArcGIS®,
22 which is based on the ANUDEM gridding application by Hutchinson (1993). Landscape units
23 were derived on the basis of the DEM, by generating the Shary classes (Shary, 2008) of
24 topographic curvature and form, wetness index, topographic complexity index, in order to

1 classify the landscape in 21,794 polygons and 13 categorical landscape units based only on
2 surface morphology. We calculated all the other covariates from this digital elevation model
3 by also using ArcGIS® functions.

4 *2.3.4 Parent material and age covariates*

5 In this group of covariates we selected to use a map of geological formation grouped in 10
6 classes based on lithology, depositional environment and relative age. This map was prepared
7 specifically for this study as a summary of two different layers: bedrock geology, as derived
8 from a number of geological maps of the Cyprus Geological Survey, of various scales,
9 published and unpublished, digitised, merged and harmonised; and quaternary geology, as
10 derived from recent field, aerial images and satellite remote mapping (Noller, 2009).

11 *2.4 Computational methods*

12 The soil groups and soil property maps were derived for the study area using RF, as
13 implemented in R (www.r-project.org) within the randomForest package (Liaw and Wiener,
14 2014) and multinomial logistic regression, as implemented in package nnet (Venables and
15 Ripley, 2002). All the random subsets used as model input were created using the caret
16 package (Khun, 2008), preserving the original class distribution of the soil classes. All
17 computations were performed on a single node of the CyTera HPC facility
18 (<http://cytera.cyi.ac.cy/>) consisting of 12 cores with 4 GB RAM each.

19 *2.4.1 The Random Forest algorithm*

20 Random Forest is a multiple tree classification and regression algorithm (Breiman, 2001) that
21 can be used to predict a target variable at locations where it is unknown, on the basis of
22 previously defined relationships between the target variable itself and identified predictors. In
23 this study, we applied RF as a classification tool. A clear overview of the method's
24 functioning is presented by Diaz-Uriate and de Anders (2006) and Boulesteix et al. (2012).

1 To construct the relationships between the target variable and the predictors, many
2 classification trees are grown. Besides the number of trees (*ntree*), only two other parameters
3 can be defined by the user before running a RF classification: the minimum size of terminal
4 nodes (*nodesize*), and the number of variables randomly sampled as candidates at each split
5 (*mtry*). Each tree is a standard classification tree. At each node the code randomly samples a
6 number (*mtry*) of predictors and among these it picks the predictor that ensures the best split.
7 Each target point is then classified by aggregating the trees and picking the class that received
8 the maximum number of votes.

9 Bootstrap samples from the original data set are used to build the trees; this means that some
10 observational values are not used to construct the trees and can be grouped to form the out-of-
11 bag (OOB) sample. This sample can be used for validation purposes by comparing it to the
12 model outputs and calculating the corresponding relative error (OOB error). An additional
13 feature of RF is the capacity to rank the relative importance of the variables in the prediction.
14 In particular, we used the Mean Decrease of Accuracy (MDA) value. It is calculated as
15 follows: first, the original OOB sample is run down a tree and the number of votes for the
16 correct target class is kept. Second, the values of a variable are randomly permuted in the
17 OOB sample, which is then run down a tree, saving the number of votes of the correct class.
18 Finally, the difference between the two vote numbers is calculated, and the results from the
19 different trees are averaged (Breiman, 2001).

20 2.4.2 Random Forest model parameterization and optimization

21 The RF classification model optimization was carried out in four steps. In the **first step**, we
22 identified the best combination of number of trees to be grown and number of training points
23 (raster cells) to be used, considering the evolution of the average OOB error and the
24 computational resources available. The number of training points was expressed in terms of
25 percentage of the total number of observation points from the 10 available maps.

1 As a **second step**, we quantified the role of each predictor and group of predictors. Firstly, we
2 evaluated the importance matrix (MDA) calculated by the RF algorithm run using all the
3 covariates not excluded in the correlation analysis. For this purpose we did not scale the
4 variable importance (as default in RF), since this can affect the interpretability of the MDA
5 investigation, as suggested in Strobl et al. (2007). Secondly, we calculated the modification in
6 the OOB error when removing predictors one by one, and group by group (e.g. relief, climate,
7 geochemistry predictors etc.), from the complete set. In addition, similar to Xiong et al.
8 (2014) and Brungard et al. (2015), we used recursive co-variable elimination (based on
9 importance and OOB error modification) to derive the subset of covariates leading to the
10 lowest OOB error. The *nodesize* and *mtry* parameters were kept to their default values for the
11 first two steps.

12 In the **third step**, we calibrated the *mtry* parameter leaving again *nodesize* to its default value.
13 We attributed to *mtry* all the odd values between 1 and the maximum number of predictors
14 coming from the optimized subset of covariates and selected the value giving the minimum
15 OOB error.

16 In the **fourth step**, we investigated the effect of tree pruning by modifying the value of the
17 *nodesize* parameter setting it to 1 (default), 4, 12, and 20. We selected the best value based on
18 the minimum OOB error, minimum average uncertainty value, and minimum independent
19 validation error. The computation of uncertainty value and validation error is explained in
20 section 2.4.4 (Map extrapolation, prediction uncertainty, map reliability, and independent
21 profile evaluation).

22 2.4.3 Multinomial Logistic Regression model

23 To determine the quality of the maps derived with RF, we compared its outputs with those of
24 an alternative, more standard technique from the generalized linear model family. We

1 selected a multinomial logistic regression (MLR) model, since it has been proven a reliable
2 method by many authors (Debella-Gilo and Etzelmüller, 2009; Kempen et al., 2009; Jafari et
3 al., 2012; Collard et al., 2014).

4 MLR is the generalization of the binomial logistic regression, where multiple possible
5 outcomes are derived from a set of predictors to which coefficients are applied. In the
6 binomial case, two possible outcomes (A_1, A_2) are given and the probability of occurrence (π)
7 of the second outcome is $\pi_2 = 1 - \pi_1$. Logistic regression relates probability π_1 to a set of
8 predictors using the logit link function:

$$9 \quad \text{logit}(\pi_1) = \ln\left(\frac{\pi_1}{\pi_2}\right) = \ln\left(\frac{\pi_1}{1 - \pi_1}\right) = \mathbf{x}'\boldsymbol{\beta} \quad (4)$$

10 where \mathbf{x} is the vector of predictors and $\boldsymbol{\beta}$ is the vector of model coefficients. In case of
11 categorical variables, every class level is fitted with a different coefficient. Model coefficients
12 are usually estimated through maximum likelihood. From eq. 4 it is possible to obtain:

$$13 \quad \pi_1 = \frac{\exp(\mathbf{x}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})} \quad (5)$$

14 which can be generalized, for a multinomial case with K target classes, to:

$$15 \quad \pi_k = \frac{\exp(\mathbf{x}'\boldsymbol{\beta}_k)}{1 + \sum_{k=2}^K \exp(\mathbf{x}'\boldsymbol{\beta}_k)}. \quad (6)$$

16 One class is selected as reference and logits are then calculated that compare other classes to
17 it. The constraint is that the summation of the probabilities of the K classes must add up to 1.

18 In this study, for MLR model building we used the same training points and set of predictors
19 optimized for RF. To compare the two methods we analyzed prediction error in comparison
20 to original soil map data, prediction error based on independent soil profiles, and prediction
21 uncertainty. Map errors were calculated based on a test set of points not used for model
22 generation for MLR and based on OOB errors for RF, since OOB errors can be considered a

1 valid alternative to independent validating data sets (Grimm et al., 2008). The computation of
2 the predictive uncertainty and independent validation of both maps are described in the next
3 section.

4 2.4.4 Map extrapolation, prediction uncertainty, map reliability, and independent profile 5 evaluation

6 We analysed different indices to evaluate and discuss the quality of the obtained soil maps.
7 Firstly, we calculated the Multivariate Environmental Similarity Surface (MESS), as
8 presented by Elith et al. (2010), based on all continuous predictors and using the dismo R
9 package (Hijmans, 2013). MESS identifies and quantifies the areas where model predictions
10 are extrapolations by measuring the similarity of any prediction point to the used training
11 points, with respect to the chosen predictor variables. The higher the score, the more common
12 the environment of the point is. Negative values indicate a novel environment, i.e. the
13 presence of at least one variable outside its range in the training points. Categorical predictors
14 are not handled by the MESS algorithm. However, they would not have affected the resulting
15 map, because we had reclassified all categorical variables to have all classes represented in
16 the training areas. Secondly, we evaluated the uncertainty of the prediction with the
17 confusion index (CI) as presented in Burrough et al. (1997):

$$18 \quad CI = [1 - (\mu_{\max i} - \mu_{(\max-1)i})] \quad (7)$$

19 where μ_{\max} represent the probability value of the class with the maximum odds at cell i and
20 $\mu_{\max-1}$ is the second largest probability at the same cell. The higher the CI is, the higher is the
21 confusion, i.e. increasing similar probabilities for the two classes. For RF the probabilities
22 were calculated from the number of tree votes, while for MLR they came from a probability
23 prediction.

1 We developed a reliability index for the created soil maps (groups, depth, texture) with the
2 aim of qualitatively evaluating a combination of model error and extent of extrapolation
3 limited to original data density. The reliability index we derived is based on the ratio of the
4 number of training points and the number of points in the final map per class (unknown areas
5 potential extrapolation error), and the error in the training map areas not used for model
6 construction (OOB error for RF and test set error for MLR). We implemented a two-step
7 approach to derive three reliability classes (High, Medium, Low). We calculated a reliability
8 index for each class of the three maps produced, according to the following equation:

$$9 \quad R_{ij} = (1 - MAP_{ij}) \cdot 2 \frac{P_{Tij}}{P_{Mij}} \quad (8)$$

10 where R is the reliability index, MAP is the error in the training map areas not used for model
11 construction, P_T is the number of training points, P_M is the number of predicted cells in the
12 derived maps, i is one of the three maps (soil groups, soil depth, soil texture), and j is a class
13 of one of the three maps. The coefficient 2 is introduced because only half the points in the
14 training map areas were used to train the models. A perfect reliability of 1 would therefore be
15 achieved with a MAP error equal to zero and a P_{Mij} value double P_{Tij} , meaning that all the
16 points inside the training area were well predicted, and the specific class is not extrapolated to
17 unknown areas. We decided to limit the value of R_{ij} between 0 (minimum reliability) and 1
18 (maximum reliability), meaning that also classes under-represented by the prediction in
19 comparison with the training areas have a maximum reliability of 1. For mapping purposes,
20 we then classified the R values into three classes: Low ($R < 0.25$), Medium ($0.25 \leq R \leq 0.50$),
21 High ($R > 0.50$).

22 The created soil depth and soil texture maps were also evaluated through 325 and 126
23 independent soil profiles (Fig. 2), respectively. The profiles are those described in Section 2.2
24 (Soil data). For each profile from the Land Suitability Map, we identified the maximum depth

1 and then we attributed it a soil depth class according to Table 1. In these profiles, clay, silt,
2 fine sand, and coarse sand percentages were specified for three or four depth intervals. We
3 classified the texture of each depth according to the soil textural triangle (FAO, 2006; IUSS,
4 2015), and then we attributed texture classes, according to Table 2. Robins (2004) provides
5 depth ranges for each road cut analysed. In case the range fell in two or more depth classes as
6 presented in Table 1, we assigned a range of classes to the validation point. Depth and texture
7 of the profile were then compared with those of the predicted maps, at the exact location. For
8 soil depth from Land Suitability map profiles, we evaluate the point as a perfect match (1) or
9 a complete miss (0). For texture we evaluated the point as a perfect match (1) if the prevailing
10 texture was correctly mapped, a partial match (0.5) if the predicted texture matched the
11 texture of at least one of the depth intervals, and a complete miss (0) in case of no match at
12 all. Similarly, for soil depth over the profiles of Robins (2004) we evaluated a perfect match
13 (1) if the average depth of the range was falling in the predicted depth class, a partial match
14 (0.5) if the mapped depth class covered, at least partially, the depth range, and a complete
15 miss (0), if the predicted class was out of range. An average hit ratio (HR) was then derived
16 summing the contribution of each profile and dividing by their total number.

17

18 **3 Results and discussion**

19 **3.1 Covariates correlation analysis**

20 An assessment of a 17 x 16 matrix of correlation scatter plots and Pearson correlation
21 coefficients led to the discard of seven out of the 17 continuous covariates presented in
22 Section 2.3. We removed single covariates out of couples presenting a correlation coefficient
23 higher than 0.6. Loss on Ignition, Total Carbon, the Chemical Index of Alteration, the Vogt
24 Ratio and the Mafic Index of Alteration (MIA) are all highly correlated with each other. We

1 preferred keeping only the MIA because it allowed a broader differentiation of the main
2 lithologies of the study area. We removed planar curvature and profile curvature because they
3 were both highly correlated with curvature. Finally, we removed average temperature
4 because it was highly correlated with minimum temperature and showing a larger correlation
5 coefficient with T_{\max} (0.56) than T_{\min} (0.24). Fig. 4 shows six out of the seven pairs of
6 correlations just described. The relation between TotC and MIA is not shown because very
7 similar to that of LOI and MIA, also in terms of correlation coefficient (-0.88).

8

9 **Fig. 4. Scatter plots, drawn with a 0.05% subset of the original input dataset, showing highly correlated**
10 **covariates. The R value shown in each graph is the Pearson correlation coefficient calculated using the**
11 **full input dataset. The covariates on the x-axis were removed from the data set used as input for the**
12 **digital mapping techniques. CIA is the chemical index of alteration; MIA is the Mafic Index of**
13 **Alteration; VR is the Vogt Ratio; and LOI is the Loss on Ignition.**

14

15 **3.2 Model parameterization and optimization**

16 The results of the **first optimization step** (number of target points and number of trees) are
17 presented in Fig. 5. The results concern the soil groups model. On the one hand, the graph
18 shows how the error depends mainly on the number of target points rather than on the number
19 of trees, especially for forests larger than 200-250 trees. On the other hand, using too many
20 training points (>50%) reduces the maximum number of trees that we can calculate without a
21 memory failure of the system to less than 300. For this reason, we performed our analysis
22 using 50% of the training points and building forests of 350 trees. This guaranteed a
23 satisfactory low average OOB error (around 10%) and a certain robustness of the
24 classification model.

25

1 **Fig. 5. Average OOB error for increasing number of trees in the Random Forest classification model and**
2 **different number (% of the total) of training points as derived from the rasterized 1:25,000 scale soil**
3 **maps.**

4
5 The results of the **second step**, for the systematic removal of single covariates are reported in
6 Table 3. Table 3 shows that there is generally a good agreement between the MDA value
7 calculated by the RF algorithm for each variable and the decrease/increase of the average
8 OOB error when removing a certain predictor. The most illustrative case is that of the
9 elevation (DEM) covariate. In fact, it appears as the most important covariate in terms of
10 MDA, and the predictor causing the largest increase of OOB error when removed for all the
11 three maps. However, a high decrease/increase in the OOB error does not always correspond
12 to a low/high value of MDA. The most enlightening case (soil groups prediction) is the one
13 of the pH, which ranked seventh based on the MDA results, while being the variable causing
14 the third highest increase of OOB average error when removed. The opposite can be seen for
15 Maximum Temperature, which was for all the three maps among the four most important
16 variables for MDA. However, if removed it caused a decrease in the OOB error.

17

18 **Table 3. Average OOB error for forests (350 trees, 50% data) built removing single variables, and**
19 **importance value (Mean Decrease of Accuracy - MDA) calculated by removing the variable from the**
20 **complete model. OOB values are presented as differences from the model derived with the complete set of**
21 **variables. Between parenthesis the rank of each variable.**

22

23 Mean Decrease of Accuracy has a great advantage over the analysis of the OOB error
24 variations due to variable removal; it can cover the impact of each predictor also in terms of
25 multivariate interactions with the other variables (Strobl et al., 2008). However, the same

1 authors point out how often the MDA of correlated variables is overestimated due to a
2 preference for correlated predictors in the early splits. This can explain why Maximum
3 Temperature, which from the previously performed correlation analysis showed a Pearson's
4 correlation coefficient with elevation of -0.46, presents such differences in behaviour in Table
5 3. When using the original RF algorithm as implemented in the randomForest package in R,
6 it is therefore very useful combining an analysis of the MDA values with a correlation
7 analysis of the predictors involved and an analysis of OOB errors following single variable
8 removal, to identify potential anomalies.

9 Table 4 shows the importance of the group of soil chemical covariates. The removal of the
10 five selected soil chemical properties from the complete model caused an increase in the error
11 of more than 9% for soil groups, 12% for soil depth, and 9% for soil texture. The removal of
12 the other four groups of covariates leads to maximum error increases of 3.5%, 3.6%, and
13 2.9% for soil groups, soil depth and soil texture, respectively. However, it is worth noting
14 how the second most important group is always relief and the rank of the groups of variables
15 is consistent for the three maps. In addition, the removal of the climate covariates (T_{\min} and
16 T_{\max}) leads to a decrease in the OOB error, for all the maps, of around 0.5%. This is probably
17 due to an overfitting of the model, when temperatures are included, which can be related to
18 correlation with other covariates (e.g. elevation).

19 The relatively high importance of relief (terrain) attributes was not a surprise. Many authors
20 conducting DSM studies using different algorithms defined them, or demonstrated them to be
21 useful predictors (e.g. Lemercier et al., 2012; Barthold et al., 2013; Taghizadeh-Mehrjardi et
22 al., 2015). Additionally, the impact of their scale and resolution has been extensively studied
23 (Behrens et al., 2010). The high impact of the soil geochemical properties is more surprising.
24 Lawley and Smith (2008) discussed how geochemical surveys can help to describe some soil
25 and weathered-zone characteristics and so be used to improve geological maps, which was

1 their final aim. The results of our study confirm that differences in soil characteristics can be
2 spotted and highlighted by these variables. Láng et al. (2016) used similar variables (organic
3 carbon content, pH, electric conductivity, sand, silt, clay and gravel content, bulk density, and
4 cation exchange capacity) to derive major soil types (WRG soil groups) in Africa. These
5 variables allowed a RF prediction with 68% classification success, thus confirming again the
6 reliability of this type of variables for soil type prediction. However, their study does not use
7 any other environmental predictor and therefore it is not possible to define a relative
8 importance for soil properties predictors. Taghizadeh-Mehrjardi et al. (2015) used variables
9 such as clay index, carbonate index, gypsum index and salinity ratio to predict USDA-family
10 soil groups in Iran. These variables were derived from Landsat 8 products along with other
11 indices. They conclude that Landsat spectral data and terrain parameters are the most useful
12 auxiliary variables, although no specific reference is made to the soil property indices.
13 Considering our results and those of these recent studies, we suggest further investigation of
14 the role of soil properties (geochemical in particular) in the prediction of soil groups, in
15 comparison with other categories of variables, and across different environments.

16

17 **Table 4. Average OOB error for forests (350 trees, 50% data) built removing groups of variables**
18 **according to the *scorpan* formula. Values are presented as differences from the model derived with the**
19 **complete set of variables. Between parenthesis the rank of each group.**

20

21 Based on the recursive elimination analysis, we removed Curvature, Aspect, and Slope for
22 the soil groups map. The resulting 10-variable classification model had an average OOB error
23 of 8.7%. Thus, the three covariates with the lowest MDA - and three out of the five covariates
24 that gave a decrease of the OOB average error when singularly removed - were excluded
25 (Table 3). For soil depth and soil texture recursive elimination led to a classification model

1 with 10 covariates each. The removed variables were the same for both maps: slope,
2 curvature, and minimum temperature. Using only the remaining predictors the OOB error
3 was reduced to 10.7% and 9.0% for soil depth and soil texture, respectively. These results
4 demonstrate that a quick analysis of the model errors derived from single variable removal
5 can provide more insight for the construction of the most parsimonious model than a simple
6 analysis of MDA. In fact, a removal of variables based only on the MDA values of covariates
7 could lead to keeping uninformative predictors in the model (due to preference for correlated
8 variables, as discussed before) or removing layers with useful information.

9 In the **third step** of the model optimization, we defined the best *mtry* parameter value. In
10 Table 5, we present the OOB errors calculated for the different classification models built
11 with varying *mtry* values, 50% of the target points, and 350 trees. Besides an *mtry* value
12 equal to 1, which gives fairly poor results for all the maps, all other *mtry* values give very
13 similar OOB error. Poor performances with low *mtry* values are likely to be expected due to
14 reduced probabilities of a relevant variable to be selected at each split (Hastie et al., 2008).
15 As discussed in Breiman (2001), the forest error rate decreases with decreasing correlation
16 between trees and increasing strength of each individual tree. Both correlation and strength
17 depend on the *mtry* parameter and increase with it. The optimum range of the *mtry* parameter
18 can be quite wide (Breiman 2001) as is also demonstrated by the results in Table 6 and
19 previous literature (e.g., Grimm et al., 2008; Barthold et al., 2013; Heung et al., 2014). For
20 the prediction of all maps, an *mtry* value of 5 was selected, being the one giving the lowest
21 OOB error and ensuring a good balance between correlation among trees and strength of the
22 trees.

23

24 **Table 5. Average OOB errors calculated using different *mtry* values. The forest is made up of 350 trees.**

25

1 In the **fourth step** the effect of tree pruning, through the setting of the *nodesize* parameter
2 was investigated. The results are presented in Table 6. An increase in *nodesize* led to an
3 increase in the OOB error for all the three maps, while it almost did not affect map
4 uncertainty, expressed in terms of CI, and validation HR calculated using independent
5 profiles. Hastie et al. (2009), in a general introduction of the RF algorithm, discussed that
6 fully grown trees (*nodesize* equal to 1) can lead to model unnecessary variance (i.e. to
7 overfitting), although in their experience no relevant performance increase is obtained from
8 *nodesize* tuning. They showed an example in which the tuning of the *nodesize* parameter
9 slightly reduced the model error, having used RF in its regression form. Our results confirm
10 those of Hastie et al. (2009) and add an example for RF used as a classifier. For our final RF
11 prediction, we kept *nodesize* to the default value (1).

12

13 **Table 6. Average OOB (for Random Forest, RF) or MAP (for Multinomial Logistic Regression, MLR,**
14 **from 50% unused points in 1:25,000 soil maps) errors, validation hit rate (HR) from profiles (not**
15 **applicable to soil groups for lack of profile data) and confusion index (CI) calculated using different**
16 **nodesize values (n). The forest is made up of 350 trees and the mtry parameter is fixed to 5.**

17

18 The OOB errors of the optimized models are 8.6%, 10.5% and 8.8% for soil groups, soil
19 depth, and soil texture respectively. Compared with other studies, these can be considered
20 very low values: Stum et al. (2010) obtained an OOB error of 55.2% in predicting soil series;
21 Barthold et al. (2013) reported an OOB error of 51.6%; Pahlavan Rad et al. (2014) calculated
22 OOB errors of 48.5%, 51.5% and 56.5% for great group, subgroup and series levels,
23 respectively; Taghizadeh-Mehrjardi et al. (2015) obtained a OOB error of 78% for family soil
24 groups. OOB error values closer to the ones obtained by our study were obtained in other
25 research fields. As an example Peters et al. (2008) achieved a minimum OOB error of 19%
26 while modelling wetland vegetation distribution.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

3.3 Comparison with Multinomial Logistic Regression

Consistent with many other studies (e.g., Brungard et al., 2015; Taghizadeh-Mehrjardi et al., 2015; Heung et al., 2016), RF provided a better predictive model than MLR for all three maps, based on average scores. The validation error calculated for MLR using the 50% of the 1:25,000 map points are much larger than the comparable OOB errors obtained from the final RF models, as presented in Table 6. Conversely, the HR from independent validation of soil depth and soil texture from profiles is comparable with RF, which performs a little better for soil depth and a little worse for soil texture (Table 6).

A spatial comparison of CI results is presented in Fig. 6. The average CI values are usually lower for RF than MLR, with the exception of soil groups where MLR performs comparably. However, average values in this case do not explain all, especially for soil depth and soil texture. The CIs calculated from the RF output for these two maps show generally low CIs in the training areas (< 0.4) and higher values (> 0.4) in the rest of the study area. Zooming in, we can see large areas with low CIs bordered by narrow bands with higher CI values, perfectly representing diffuse boundaries among classes. Since this level of detail is not kept in all the prediction areas, it can be interpreted as an overfitting of the RF model, although no benefit appeared from tree pruning. CI maps derived from the RF built with *nodesize* equal to 20 (not presented here) show exactly the same pattern. Conversely, The CIs calculated from MLR are more homogeneous throughout the study area. Additional visual comparison between RF and MLR shows, in any case, similar trends of CI over the areas not included in training maps, with rather low values over the Troodos Mountains, in the centre of the island, and higher value along the borders of the study area. Considering these results, and above all the large difference in the validation from independent map values, we decided to present as final maps those predicted using RF.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Fig. 6. Confusion Index (CI) maps calculated for soil groups, soil depth and soil texture prediction from both Random Forest (left) and Multinomial Logistic Regression (right).

3.4 Final maps, map extrapolation, map reliability, and evaluation

Fig. 7, shows the final map of soil groups, whereas Fig. 8 shows the soil depth and the soil texture maps. Three main areas can be recognized: the very shallow (< 10 cm), lithic, and poorly developed soils of the Troodos Mountains; the young, (weakly) developed, shallow to moderately deep (between 10 cm and 100 cm), calcaric soils of the southern flank of the Troodos Mountains; and the well-developed, deep to very deep (> 100 cm), fertile soils of the western Mesaoria Plain and coastal plains in the southeast. Among the latter, it is worth mentioning the characteristic red soils here classified as Chromic Luvisols (with additional qualifiers according to the substrate). The good detail of the predicted maps is demonstrated along the southern coast, where differences between valley bottoms (Calcaric Cambisols) and slopes (Calcaric Regosols) can be recognized, and the most fertile soils of the country, located in the area of Pafos, can be spotted as Calcaric and Eutric Fluvisols. However, small errors, especially over the Troodos Mountains, can be found. Here, small lenses of Gypsic and Calcaric Regosols have been predicted although, considering the igneous, mafic and ultramafic nature of the bedrock, these can be considered impossible soil types for the region.

Fig. 7. Digital soil map of Cyprus. The map is presented with WRB soil group names accompanied by one or two qualifiers as predicted using Random Forest.

Fig. 8. Digital soil depth and soil texture maps of Cyprus as derived from Random Forest.

1

2 The MESS and the reliability of the produced maps are presented in Fig. 9 and Fig. 10,
3 respectively. Although slightly different set of predictors were used for soil groups, soil depth
4 and soil texture, only the MESS related to soil groups is presented, since the differences
5 between these maps were minimal. It is worth noting how both in terms of MESS and
6 reliability the study area can be roughly split into two main regions. On one hand, the MESS
7 and the reliability of the predicted soils in the main agricultural areas of the country
8 (Mesaoria plain, southeastern coast, and mountain foothills) show medium to high values,
9 confirming that the RF prediction is legitimate (not an extrapolation) and its quality is good.
10 On the other hand, the predominant soils of the Troodos Mountains show a low reliability and
11 are located over areas with a highly negative MESS. This low reliability is driven by the large
12 size of the extrapolated area rather than a high OOB error. An important limiting factor for
13 the prediction of soils over the Troodos region is the lack of a training area covering the
14 highest elevations of the mountains. Limited training areas are only available along the
15 foothills of the mountains, where the typical mountain soils start to form, but none of these
16 areas cover regions characterized by steep valleys, as this terrain is found at higher elevations
17 only. In addition, some of these valleys are rather peculiar, because they have been modified
18 by human activities since the Bronze Age, especially by the construction of dry stone terraces
19 to favour agricultural practices (Fall et al., 2012). Thus, while based on the topographic,
20 geological, geomorphological, and climatic conditions, the predicted prevalent soil depth (0-
21 10 cm) and soil texture (none) of the Troodos region (Fig. 8) can be considered reasonable,
22 some of the natural variability and the human modifications of the terrain are not captured.
23 The produced soil depth and soil texture maps generally fail to represent the peculiarities of
24 the mountainous agricultural land, which could be also of interest for applied environmental
25 studies (e.g. water needs for agriculture, soil erosion).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Fig. 9. Multivariate Environmental Similarity Surface (MESS) showing areas where model predictions are extrapolations in comparison to the training data set (values lower than zero).

Fig. 10. Reliability maps derived for the soil groups, soil depth, and soil texture maps of Cyprus predicted with Random Forest.

Table 7 summarizes the validation HRs based on independent profiles for each soil depth and soil texture class and for areas where the MESS values is larger or lower than zero (i.e. non extrapolation or extrapolation). The soil profiles cover all soil depth and soil texture classes, although with a non-uniform distribution. For soil depth, the moderate to deep (50-100 cm) class is highly under-represented in comparison with the other three, while the class shallow (10-50 cm) has the largest number of profiles. For soil texture, the two classes Medium and Fine are similarly represented and comprise 83% of all the profiles.

At the profile locations, the calculated overall validation HR for soil depth and soil texture are 55% and 49%, respectively. Heung et al. (2016) obtained a 54% HR predicting soil great groups with RF using the same method for training data sampling. Barthold et al. (2013) got much better results mapping WRB soil groups in the Mongolian grasslands, having a misclassification error of 29%. Comparably with the latter authors, Lang et al. (2016) obtained a 32% classification error for WRB soil groups in Africa. Considering that the error rate is a little higher than typical literature value, it is worth some discussion.

Table 7 shows that there is not much difference in the results of points predicted in the known or the extrapolated areas, but that there are large differences in the error between different classes, which have different reliabilities. In terms of soil depth, we can state that the model

1 is masking short distance heterogeneity, which is represented by the profiles. For example,
2 the whole area of the Troodos Mountains is represented as a zero depth soil. This is due to the
3 fact that 90% of the limited training areas presenting the main environmental characteristic of
4 the Troodos is mapped as zero depth soil. Variability in this region exists, as testified by the
5 soil-profile analyses of Robins (2004), who found a soil depth range between 0 and 0.96 m,
6 with a median of 0.15 m and an average value of 0.18 m. Our model has not been trained to
7 match Robins' soil depth data, as this would require larger training areas in the Troodos
8 Region and a significantly finer resolved terrain model. Similarly, very-deep soils are more
9 common than moderate-to-deep soils in the training areas and are better represented in the
10 prediction, with less error. For texture, the smallest HR is associated with the class Coarse,
11 which is the least represented in the training data set (6%). The texture class Medium, which
12 is the most represented with 45% of the training points, is also the one with the best HR. The
13 model therefore shows a tendency to predict the classes that occur more frequently in the
14 training areas.

15

16 **Table 7. Number of independent soil profiles (N. prof.) and validation hit rate (HR), calculated based on**
17 **positive or negative Multivariate Environmental Surface Similarity (MESS), for soil depth classes and soil**
18 **texture classes.**

19

20 **4 Conclusions**

21 This study produced a digital soil map of Cyprus (including soil groups, soil depth classes
22 and soil texture classes) with RF using training areas from ten 1:25,000 published soil maps
23 (27% of the study area). The study also focused on the optimization of the RF model in terms
24 of environmental variables used as predictors and of the model's adjustable parameters

1 (number of training points, number of trees in the forest, and number of variables sampled at
2 each branch split). The study proved how the average OOB error, calculated by the RF
3 model, increases at a higher rate with the reduction of the number of training points than with
4 the reduction of the number of trees grown in the forest. This can be especially noticed for
5 forests larger than 200-250 trees (when an asymptotic behaviour of the OOB error is
6 reached). Considering this asymptotic behaviour, in this study forests of 350 trees were
7 calculated.

8 For the selection of the optimal set of environmental variables to run the model (lowest error
9 and lowest number of predictors), the recursive feature elimination based on OOB error was
10 found to be generally superior to elimination based on the lowest Mean Decrease of Accuracy
11 (MDA) value per predictor, as calculated by the RF algorithm. This is mainly due to MDA
12 being influenced by variable correlation. In the optimum models derived for the prediction of
13 soil groups, soil depth and soil texture, a prevalent role in keeping a low model error is
14 played by the variables linked with the geochemistry of the soil. They appeared as the most
15 important group of variables in our models. The importance of these variables should be
16 further investigated in different geographic and pedological areas to confirm their quality as
17 soil group predictors.

18 The calibration of the *mtry* parameter (number of variables sampled at each branch split)
19 revealed a wide range of optimum values, distributed around half the number of
20 environmental variables used to drive the prediction. The calibration of the *nodesize*
21 parameter showed no relevant performance increase and was kept at its default value (1).

22 The produced soil groups, soil depth and soil texture maps showed very low OOB error:
23 8.6%, 10.5%, and 8.8%, respectively. However, when evaluated for soil depth and soil
24 texture outside the training areas by means of independent soil profiles, the average error
25 shows values equals to 45% and 51%, respectively. In particular, the maps show a medium to

1 high reliability over the major agricultural areas of the country (plain and mountain foothills)
2 but a low reliability over the mountainous region. This low reliability is mainly due to a lack
3 of training data over the mountains, whose soils are largely extrapolated from few training
4 areas on the mountain foothills. Despite low OOB errors, peculiar characteristics of steep
5 mountain environments, like medium to deep soils typical of terraced agricultural land, are
6 missed by the model. Also, the model showed a tendency to predict more commonly the
7 classes that occur more frequently in the training areas, therefore masking local variance in
8 soil properties.

9 Maps derived with MLR had comparable prediction uncertainty and higher validation error
10 than RF, indicating the better performance of the latter for DSM in topographical and
11 pedological complex regions. If properly applied, RF was found to be a reliable instrument
12 for DSM activities. Future research will focus on the development of a specific digital soil
13 map for the Troodos Mountains of Cyprus, including both field and modelling studies.

14 **5 Acknowledgements**

15 This study is part of the AGWATER project (ΑΕΙΦΟΡΙΑ/ΓΕΩΡΓΟ/0311(BIE)/06), co-
16 financed by the European Regional Development Fund and the Republic of Cyprus through
17 the Research Promotion Foundation. All the simulations have been run at the Cy-Tera High
18 Performance Computing (HPC) facility of the Cyprus Institute, whose staff, Thekla Loizou in
19 particular, we would like to thank for the support to our study. Special thanks also to Hakan
20 Djuma and Marinos Eliades for the digitization of the data of the soil profiles analysed in the
21 1970s and 1980s.

1 **6 References**

- 2 Akpa, S.I.C., Odeh, I.O.A., Bishop, T.F.A., Hartemink, A.E., Amapu, I.Y., 2016. Total soil organic
3 carbon and carbon sequestration potential in Nigeria. *Geoderma* 271, 202-215.
4 doi:10.1016/j.geoderma.2016.02.021.
- 5 Allen, R.G., Pereira, L.S., Raes, D., Smith, M., 1998. *Crop Evapotranspiration: Guidelines for*
6 *computing crop water requirements*, FAO Irrigation and Drainage Paper No 56. Food and Agriculture
7 Organisation, Land and Water. Rome, Italy.
- 8 Babechuk, M.G., Widdowson, M., Kamber, B.S., 2014. Quantifying chemical weathering intensity
9 and trace element release from two contrasting basalt profiles, Deccan Traps, India. *Chemical*
10 *Geology* 363, 56–75. doi:10.1016/j.chemgeo.2013.10.027.
- 11 Barthold, F.K., Wiesmeier, M., Breuer, L., Frede, H.-G., Wu, J., Blank, F.B., 2013. Landuse and
12 climate control the spatial distribution of soil types in the grasslands of Inner Mongolia. *J. Arid*
13 *Environ.* 88, 194–205. doi:10.1016/j.jaridenv.2012.08.004.
- 14 Behrens, T., Zhu, A., Schmidt, K., Scholten, T., 2010. Multi-scale digital terrain analysis and feature
15 selection for digital soil mapping. *Geoderma* 155, 175-185. doi:10.1016/j.geoderma.2009.07.010.
- 16 Bird, D.N., Benabdallah, S., Gouda, N., Hummel, F., Koeberl, J., La Jeunesse, I., Meyer, S.,
17 Pretenthaler, F., Soddu, A., Woess-Gallasch, S., 2016. Modelling climate change impacts on and
18 adaptation strategies for agriculture in Sardinia and Tunisia using AquaCrop and value-at-risk. *Sci.*
19 *Tot. Environ.* 543, 1019-1027. doi:10.1016/j.scitotenv.2015.07.035.
- 20
- 21 Bolstad, P.V., Gessler, P., Lillesand, T.M., 1990. Positional uncertainty in manually digitized map
22 data. *Intl. Jr. Geog. Info.* 4 399-412. doi: 10.1080/02693799008941555.
- 23 Boulesteix, A.L., Janitza, S., Kruppa, J., and König, I.R., 2012. Overview of random forest
24 methodology and practical guidance with emphasis on computational biology and bioinformatics.
25 *WIREs Data Mining Knowl. Discov.* 2, 493–507. doi:10.1002/widm.1072.
- 26 Breiman, L., 2001. Random forests. *Machine Learning* 45, 5–32. doi:10.1023/A:1010933404324.

1 Brungard, C.W., Boettinger, J.L., Duniway, M.C., Wills, S.A., Edwards Jr., T.C., 2015. Machine
2 learning for predicting soil series in three semi-arid landscapes. *Geoderma* 239–240, 68–63.doi:
3 10.1016/j.geoderma.2014.09.019.

4 Burrough, P.A., 1986. Principles of geographical information systems for land resources assessment.
5 Clarendon Press, Oxford, UK.

6 Burrough, P.A., van Gaans, P.F.M., Hootsmans, R., 1997. Continuous classification in soil survey:
7 spatial correlation, confusion and boundaries. *Geoderma* 77, 115-135. doi:10.1016/S0016-
8 7061(97)00018-9.

9 Büttner, G., Kosztra, B., 2007. CORINE Land Cover 2006 Technical guidelines. Technical Report
10 No. 17/2007.EEA, 2007. Available at:
11 http://www.eea.europa.eu/publications/technical_report_2007_17.

12 Camera, C., Bruggeman, A., Hadjinicolaou, P., Pashiardis, S., Lange, M.A., 2013. High resolution
13 gridded datasets for meteorological variables: Cyprus, 1980-2010 and 2020-2050. AGWATER
14 Scientific Report 5, The Cyprus Institute, Nicosia, 70 pp.

15 Chagas, C.D., de Carvalho, W., Bhering, S.B., Calderano, B., 2016. Spatial prediction of soil surface
16 texture in a semiarid region using random forest and multiple linear regressions. *Catena* 139, 232-240.
17 doi:10.1016/j.catena.2016.01.001.

18 Cohen, D.R., Rutherford, N.F., Morisseau, E., and Zissimos, A.M., 2011. *Geochemical Atlas of*
19 *Cyprus*. Sydney: UNSW Press.

20 Cohen, D.R., Rutherford, N.F., Morisseau, E., Zissimos, A.M., 2012a. Geochemical patterns in the
21 soils of Cyprus. *Sci. Total Environ.* 420, 250–262. doi:10.1016/j.scitotenv.2012.01.036.

22 Cohen, D.R., Rutherford, N.F., Morisseau, E., Christoforou, E., Zissimos, A.M., 2012b.
23 Anthropogenic versus lithological influences on soil geochemical patterns in Cyprus. *Geochem.*
24 *Explor. Environ. Anal.* 12, 349–360. doi:10.1144/geochem2011-111.

1 Collard, F., Kempen, B., Heuvelink, G.B.M., Saby, N.P.A., Richer de Forges, A.C., Lehmann, S.,
2 Nehlig, P., Arrouays, D., 2014. Refining a reconnaissance soil map by calibrating regression models
3 with data from the same map (Normandy, France). *Geoderma Reg.* 1, 21–30.

4 Debella-Gilo, M., Eitzel Müller, B., 2009. Spatial prediction of soil classes using digital terrain analysis
5 and multinomial logistic regression modeling integrated in GIS: examples from Vestfold County,
6 Norway. *Catena* 77, 8–18.

7 Diaz-Uriate, R., de Andres, S.A., 2006. Gene selection and classification of microarray data using
8 random forest. *Bmc Bioinformatics* 7. doi:10.1186/1471-2105-7-3.

9 Djuma, H., Bruggeman, A., Camera, C., Zoumides, C., under review. Combining qualitative and
10 quantitative methods for soil erosion assessments: an application in a sloping Mediterranean
11 watershed, Cyprus. *Land Degrad. Dev.*

12 Elith, J., Kearney, M., Phillips, S., 2010. The art of modelling range-shifting species. *Methods in*
13 *Ecology and Evolution* 1: 330-342. doi:10.1111/j.2041-210X.2010.00036.x.

14 Fall, P.L., Falconer, S.E., Galletti, C.S., Shirmang, T., Ridder, E., Klinge J., 2012. Long-term agrarian
15 landscapes in the Troodos foothills, Cyprus. *J. Archaeol. Sci.* 39, 2335–2347. doi:
16 10.1016/j.jas.2012.02.010.

17 FAO, 2006. Guidelines for soil description, Fourth edition. Food and Agriculture Organization of the
18 United Nations, Rome.

19 FAO/IIASA/ISRIC/ISS-CAS/JRC, 2008. Harmonized World Soil Database (version 1.0). FAO,
20 Rome, Italy and IIASA, Laxenburg, Austria.

21 Grimm, R., Behrens, T., Märker, M., Elsenbeer, H., 2008. Soil organic carbon concentrations and
22 stocks on Barro Colorado Island — Digital soil mapping using Random Forests analysis. *Geoderma*
23 146, 102-113. doi:10.1016/j.geoderma.2008.05.008.

24 Grinand, C., Arrouays, D., Laroche, B., Martin, M.P., 2008. Extrapolating regional soil landscapes
25 from an existing soil map: sampling intensity, validation procedures, and integration of spatial
26 context. *Geoderma* 143 180–190. doi:10.1016/j.geoderma.2007.11.004.

1 Grivas, G.C., 1969. Report on the soils of the Morphou Watershed. Soil Section, Department of
2 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 53 pp.

3 Grivas, G.C., and Georgiades, M., 1972. 1:25,000 Sheet 30 Lakatamia. Soil Section, Department of
4 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus.

5 Grivas, G., 1988. Development of land resources in Cyprus. In: Zomenis, S.L. et al. (ed.) Proceedings
6 - Workshop on conservation and development of natural resources in Cyprus - case studies - soils -
7 groundwater - mineral resources. Ministry of Agriculture and Natural Resources, Cyprus and Federal
8 Institute for Geosciences and Natural Resources, W. Germany, 13-18 Oct 1986, Nicosia. p.7-16.

9 Hastie T., Tibshirani R., Friedman, J., 2009. The elements of statistical learning: Data mining,
10 inference, and prediction, Second Edition. Springer Series in Statistics, Springer, 763 pp.

11 Hijmans, R.J., 2013. Dismo, Species Distribution Modeling. R package version 1.0-15 [Available at
12 <https://cran.r-project.org/web/packages/dismo/dismo.pdf>]

13 Heung, B., Bulmer, C.E., Schmidt, M.G., 2014. Predictive soil parent material mapping at a regional-
14 scale: A Random Forest approach. *Geoderma* 214-215, 141-154. doi:
15 10.1016/j.geoderma.2013.09.016.

16 Hueng, B., Ho, H.C., Zhang, J., Knudby, A., Bulmer, C.E., Schmidt, M.G., 2016. An overview and
17 comparison of machine-learning techniques for classification purposes in digital soil mapping.
18 *Geoderma* 265, 62-77. doi:10.1016/j.geoderma.2015.11.014.

19 Hutchinson, M. F., 1993. Development of a continent-wide DEM with applications to terrain and
20 climate analysis. In *Environmental Modeling with GIS*, ed. M. F. Goodchild et al., 392–399. New
21 York: Oxford University Press

22 IUSS Working Group WRB, 2015. World Reference Base for Soil Resources 2014, update 2015.
23 International soil classification system for naming soils and creating legends for soil maps. World Soil
24 Resources Reports No. 106. FAO, Rome.

1 Jafari, A., Finke, P.A., Van deWauw, J., Ayoubi, S., Khademi,H., 2012. Spatial prediction of USDA-
2 great soil groups in the Arid Zarand region, Iran: comparing logistic regression approaches to predict
3 diagnostic horizons and soil types. *Eur. J. Soil Sci.* 63, 284–309.

4 Jenny, H., 1941. *Factors of Soil Formation*. McGraw-Hill, New York.

5 Kempen B., Brus, D.J., Heuvelink, G.B.M., Stoorvogel, J.J., 2009. Updating the 1:50,000 Dutch soil
6 map using legacy soil data: A multinomial logistic regression approach. *Geoderma* 151, 311-326.
7 doi:10.1016/j.geoderma.2009.04.023.

8 Kempen, B., Brus, D.J., Stoorvogel, J.J., Heuvelink, G.B.M., De Vries, F., 2012. Efficiency
9 comparison of conventional and digital soil mapping for updating soil maps. *Soil Sci. Soc. Am. J.* 76
10 (6), 2097–2115. doi:10.2136/sssaj2011.0424.

11 Kuhn, M., 2008. Building predictive models in R using the caret package. *J. Stat. Softw.* 28, 1–26.

12 Koudounas, C., Makin, J., 1978. A study of representative soil profiles from the Limassol Paphos
13 districts, published report, Department of Agriculture, Nicosia, Cyprus, 61 pp.

14 Koumis, C.I., 1970a. Report of the soils of the Limassol Watershed. Soil Section, Department of
15 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 229pp.

16 Koumis, C.I., 1970b. Report of the soils of the Larnaka Watershed. Soil Section, Department of
17 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 119pp.

18 Koumis, C.I., 1970c. Report of the soils of the Pissouri-Paramali Watershed. Soil Section, Department
19 of Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 59pp.

20 Koumis, C.I., 1980a. 1:25,000 Sheet 53 Ypsonas, Soil Section, Department of Agriculture, Ministry
21 of Agriculture and Natural Resources, Cyprus.

22 Koumis, C.I., 1980b. 1:25,000 Sheet 54 Limassol, Soil Section, Department of Agriculture, Ministry
23 of Agriculture and Natural Resources, Cyprus.

24 Koumis, C.I., 1980c. 1:25,000 Sheet 58&59 Akrotiri, Soil Section, Department of Agriculture,
25 Ministry of Agriculture and Natural Resources, Cyprus.

1 Lagacherie, P., Legros, J.P., Burfough, P.A., 1995. A soil survey procedure using the knowledge of
2 soil pattern established on a previously mapped reference area. *Geoderma* 65 283–301.
3 doi:10.1016/0016-7061(94)00040-H

4 Lagacherie, P., McBratney, A.B., Voltz, M., 2007. *Digital soil mapping: An introductory perspective*.
5 Elsevier, Amsterdam, 658 pp.

6 Lagacherie, P., Sneep, A.R., Gomez, C., Bacha, S., Coulouma, G., Hamrouni, M.H., Mekki, I., 2013.
7 Combining vis–NIR hyperspectral imagery and legacy measured soil profiles to map subsurface soil
8 properties in a Mediterranean area (Cap-Bon, Tunisia). *Geoderma* 209–210, 168–176.
9 doi:10.1016/j.geoderma.2013.06.005

10 Lang, V., Fuchs, M., Szegi, T., Csorba, A., Micheli, E., 2016. Deriving World Reference Base
11 Reference Soil Groups from the prospective Global Soil Map product - A case study on major soil
12 types of Africa. *Geoderma* 263, 226-233. doi:10.1016/j.geoderma.2015.07.005.

13 Lawley, R., Smith, B., 2008. Digital soil mapping at a national scale: A knowledge and GIS based
14 approach to improving parent material and property information. In: Hartemink, A.E., McBratney, A.,
15 Mendoça-Santos, M., Eds, *Digital Soil Mapping with Limited Data*, Springer, 173-182.

16 Lemercier, B., Lacoste, M., Loum, M., Walter, C., 2012. Extrapolation at regional scale of local soil
17 knowledge using boosted classification trees: a two-step approach. *Geoderma* 171-172 75–84.
18 doi:10.1016/j.geoderma.2011.03.010.

19 Liaw, A., Wiener, M., 2014. randomForest: Breiman and Cutler’s random forests for classification
20 and regression. R package version 4.6-10. [Available at [http://cran.r-](http://cran.r-project.org/web/packages/randomForest/randomForest.pdf)
21 [project.org/web/packages/randomForest/randomForest.pdf](http://cran.r-project.org/web/packages/randomForest/randomForest.pdf)]

22 Markides, L., 1973. *Soils Memoirs of Polemi*, Sheet no. 44 & 45, includes 1:25,000 map, Ministry of
23 Agriculture and Natural Resources, Department of Agriculture, 138 p.

24 Markides, L., 1969. *Report on the soils of the Polis-Tylliria Watershed*. Soil Section, Department of
25 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 55 pp.

1 Markides, L., 1975. 1:25,000 Sheet 41 Ormidhia. Soil Section, Department of Agriculture, Ministry of
2 Agriculture and Natural Resources, Cyprus.

3 Markides, L., 1985a. 1:25,000 Sheet 50&56 Kiti. Soil Section, Department of Agriculture, Ministry of
4 Agriculture and Natural Resources, Cyprus.

5 Markides, L., 1985b. Report on the soils of the Krasochoria integrated rural development project area.
6 Soil Section, Department of Agriculture, Ministry of Agriculture and Natural Resources, Cyprus, 35
7 pp.

8 McBratney, A.B., Santos, M.L.M., Minasny, B., 2003. On digital soil mapping. *Geoderma* 117, 3–52.
9 doi:10.1016/S0016-7061(03)00223-4.

10 Minasny, B., McBratney, A.B., 2016. Digital soil mapping: A brief history and some lessons.
11 *Geoderma* 264, 301–311. doi:10.1016/j.geoderma.2015.07.017.

12 Noller, J., 2009. The Geomorphology of Cyprus. Cyprus Geological Survey, Open File Report, 269
13 pp.

14 Pahlavan Rad, M.R., Toomanian, N., Khormali, F., Brungard, C.W., Komaki, C.B., Bogaert, P., 2014.
15 Updating soil survey maps using random forest and conditioned Latin hypercube sampling in the
16 loess derived soils of northern Iran. *Geoderma* 232–234, 97–106.
17 doi:10.1016/j.geoderma.2014.04.036.

18 Panagos, P., Van Liedekerke, M., Montanarella, L., Jones, R.A., 2008. Soil organic carbon content
19 indicators and web mapping applications. *Environmental Modelling and Software* 23, 1207–1209.
20 doi:10.1016/j.envsoft.2008.02.010.

21 Peters, J., Verhoest, N.E.C., Samson, R., Boeckx, P., De Baets, B., 2008. Wetland vegetation
22 distribution modelling for the identification of constraining environmental variables. *Landsc. Ecol.* 23,
23 1049-1065.

24 Poggio, L., Gimona, A., Brewer, M.J., 2013. Regional scale mapping of soil properties and their
25 uncertainty with a large number of satellite-derived covariates. *Geoderma* 209–210, 1–14.
26 doi:10.1016/j.geoderma.2013.05.029.

1 Ren, L., Cohen, D.R., Rutherford, N.F., Zissimos, A.M., Morisseau, E.G., 2015. Reflections of the
2 geological characteristics of Cyprus in soil rare earth element patterns. *Applied Geochemistry* 56, 80–
3 93. doi:10.1016/j.apgeochem.2015.02.011.

4 Robins, C.R., 2004. Spatial Analysis of soil depth variability and pedogenesis along toposequences in
5 the Troodos Mountains, Cyprus. Master of Science Thesis, Oregon State University, Corvallis, 119
6 pp.

7 Saxton, K.E., Rawls, W.J., 2006. Soil water characteristic estimates by texture and organic matter for
8 hydrologic solutions. *Soil Sci. Soc. Am. J.* 70, 1569–1578. doi:10.2136/sssaj2005.0117.

9 Schmidt, K., Behrens, T., Scholten, T., 2008. Instance selection and classification tree analysis for
10 large spatial datasets in digital soil mapping. *Geoderma* 146, 138-146.
11 doi:10.1016/j.geoderma.2008.05.010.

12 Scull, P., Franklin, J., Chadwick, O.A., McArthur, D., 2003. Predictive soil mapping: a review.
13 *Progress Phys. Geogr.* 27, 171–197. doi:10.1191/0309133303pp366ra.

14 Scull, P., Franklin, J., Chadwick, O.A., 2005. The application of classification tree analysis to soil
15 type prediction in a desert landscape. *Ecol. Model.* 181, 1–15. doi:10.1016/j.ecolmodel.2004.06.036.

16 Shary, P.A., 2008. Models of topography. In: Zhou, Q., Lees, B., Tang G.A. (eds.), *Advances in*
17 *digital terrain analysis*, Berlin, Springer, 29-57.

18 Soon, Y. K., Abboud, S., 1991. A comparison of some methods for soil organic carbon determination.
19 *Communications in Soil Science and Plant Analysis* 22, 943–954. doi:10.1080/00103629109368465.

20 Soteriades, C., Georgiades, M., 1967. 1:25,000 Sheet 22 Kythrea, Soil Section, Department of
21 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus.

22 Soteriades, C., Grivas, G., 1968. 1:25,000 Sheet 20 Kokkinotrimithia, Soil Section, Department of
23 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus.

24 Soteriades, C., Koudounas, C, Markides, L., 1968. 1:25,000 Sheet 51 Paphos, Soil Section,
25 Department of Agriculture, Ministry of Agriculture and Natural Resources, Cyprus.

1 Soteriades, C., Koudounas C., 1968. Soils Memoirs of Pafos, Sheet no. 51, includes 1:25,000 map,
2 Ministry of Agriculture and Natural Resources, Department of Agriculture, 96 p.

3 Soteriades, C., Markides, L., 1969. 1:25,000 Sheet 44&45 Polemi, Soil Section, Department of
4 Agriculture, Ministry of Agriculture and Natural Resources, Cyprus.

5 Sreenivas, K., Dadhwal, V.K., Kumar, S., Harsha, G.S., Mitran, T., Sujatha, G., Suresh, G.J.R.,
6 Fyzee, M.A., Ravisankar, T., 2016. Digital mapping of soil organic and inorganic carbon status in
7 India. *Geoderma* 269, 160-173. doi:10.1016/j.geoderma.2016.02.002.

8 Strobl, C., Boulesteix, A. L., Zeileis, A., Hothorn, T., 2008. Bias in random forest variable importance
9 measures: Illustrations, sources and a solution. *BMC bioinformatics* 8, 25. doi:10.1186/1471-2105-8-
10 25.

11 Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., Zeileis, A., 2008. Conditional variable
12 importance for random forests. *BMC bioinformatics* 9, 307. doi:10.1186/1471-2105-9-307.

13 Stum, A.K., Boettinger, J.L., White, M.A., Ramsey, R.D., 2010. Random forests applied as a soil
14 spatial predictive model in Arid Utah. In: Boettinger, J.L., Howell, D.W., Moore, A.C., Hartemink,
15 A.E., Kienast-Brown, S. (Eds.), *Digital Soil Mapping: Bridging Research, Environmental*
16 *Application, and Operation*. Springer, Dordrecht, pp. 179–190.

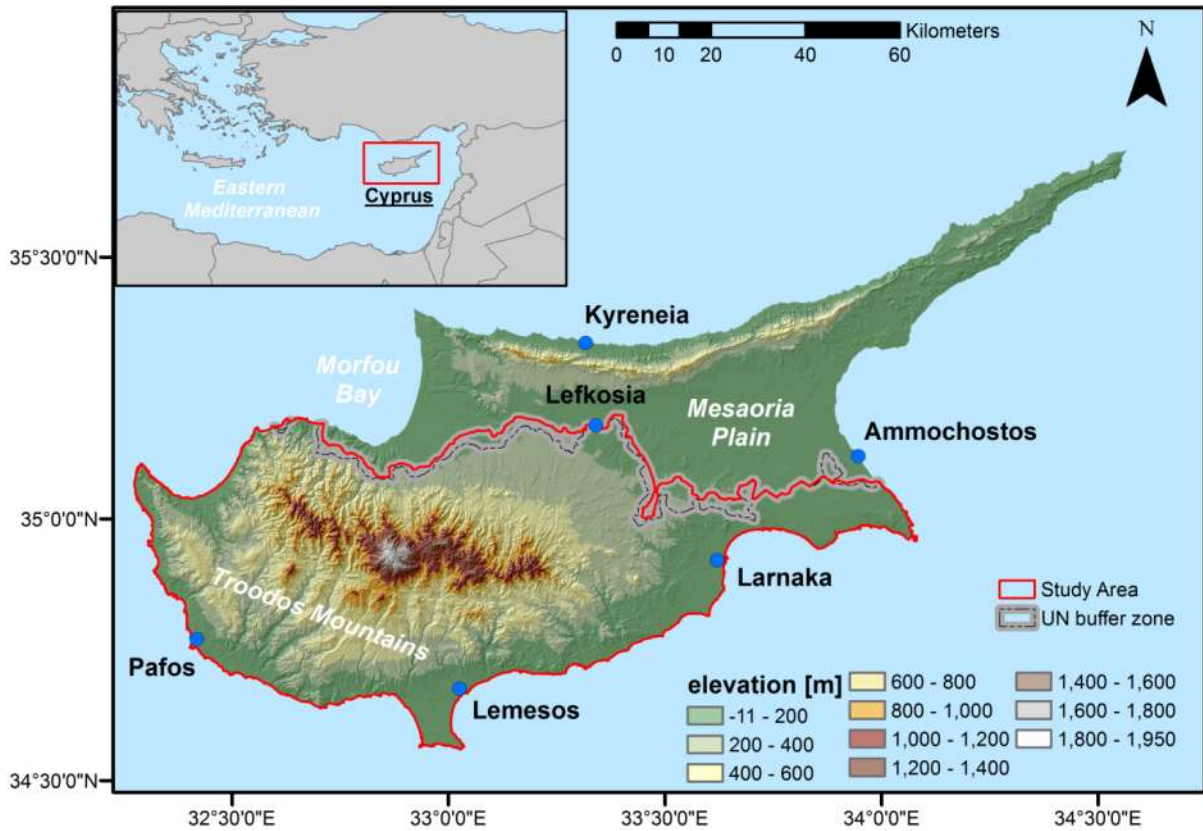
17 Taghizadeh-Mehrjardi, R., Nabiollahi, K., Minasny, B., Triantafilis, J., 2015. Comparing data mining
18 classifiers to predict spatial distribution of USDA-family soil groups in Baneh region, Iran. *Geoderma*
19 253, 67-77. doi: 10.1016/j.geoderma.2015.04.008.

20 Venables, W.N., Ripley, B.D., 2002. *Modern Applied Statistics with S*, Fourth ed. Springer, New
21 York.

22 Wiesmeier, M., Barthold, F., Blank, B., Kögel-Knabner, I., 2011. Digital mapping of soil organic
23 matter stocks using random forest modeling in a semi-arid steppe ecosystem. *Plant Soil* 340, 7–24.

24 Xiong, X., Grunwald, S., Myers, D.B., Kim, J., Harris, W.G., Comerford, N.B., 2014. Holistic
25 environmental soil–landscape modeling of soil organic carbon. *Environ. Model. Softw.* 57, 202–215.
26 doi:10.1016/j.envsoft.2014.03.004.

- 1 Yeomans, J.C., Bremner, J.M., 1988. A rapid and precise method for routine determination of organic
2 carbon in soil. *Communications in Soil Science and Plant Analysis* 19, 1467–1476.
3 doi:10.1080/00103628809368027.
- 4 Zissimos, A.M., Christoforou, I.C., Morisseau, E., Cohen, D.R., Rutherford, N.F., 2014. Distribution
5 of water-soluble inorganic ions in the soils of Cyprus. *Journal of Geochemical Exploration* 146, 1–8.
6 doi:10.1016/j.gexplo.2014.07.004
- 7 Zomeni, Z., 2012. Quaternary marine terraces on Cyprus: constraints on uplift and pedogenesis, and
8 the geoarchaeology of Palaipafos, PhD Dissertation, Oregon State University.
- 9 Zomeni, Z., Bruggeman, A., 2013. Chapter III, Soil resources of Cyprus, in *Soil Resources of*
10 *Mediterranean and Caucasus Countries*, edited by Yigini, Y., Panagos, P., Montanarella, L.,
11 Luxembourg: Publications Office of the European Union, EUR25988EN Scientific and Technical
12 Research series, pp. 37-59.
- 13



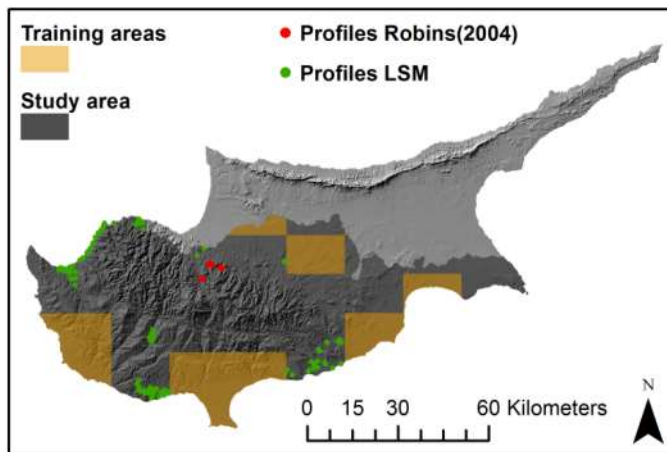
1

2

Fig. 1. The island of Cyprus with its main physical characteristics and the location of the study area.

3

4



5

6

Fig. 2. Location of the study area, the existing 1:25,000 scale soil sheets used for training, the 126 soil

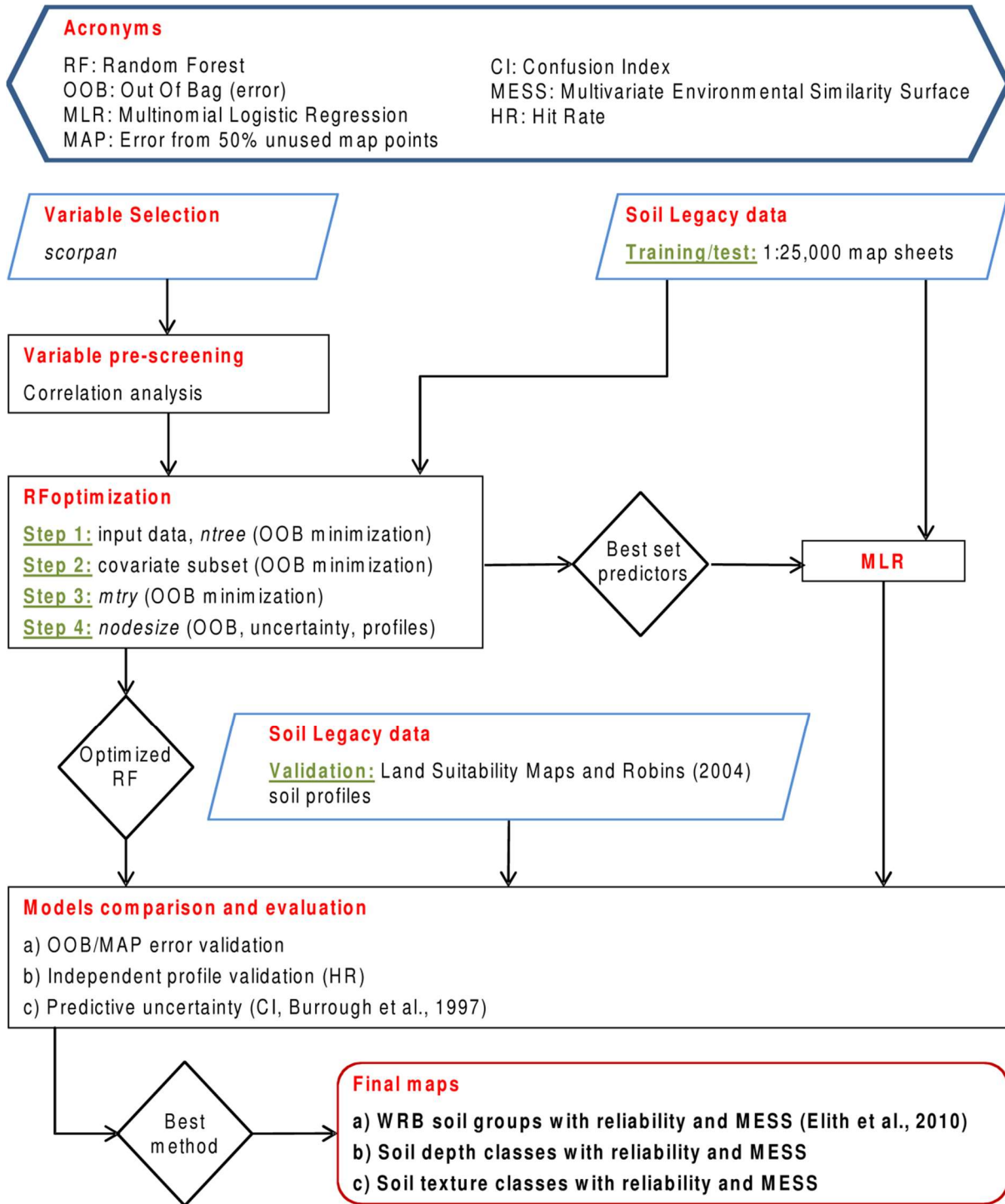
7

profiles derived from Land Suitability Maps (LSM), and the 199 soil profiles from Robins (2004) used to

8

evaluate the digital soil map of Cyprus.

1

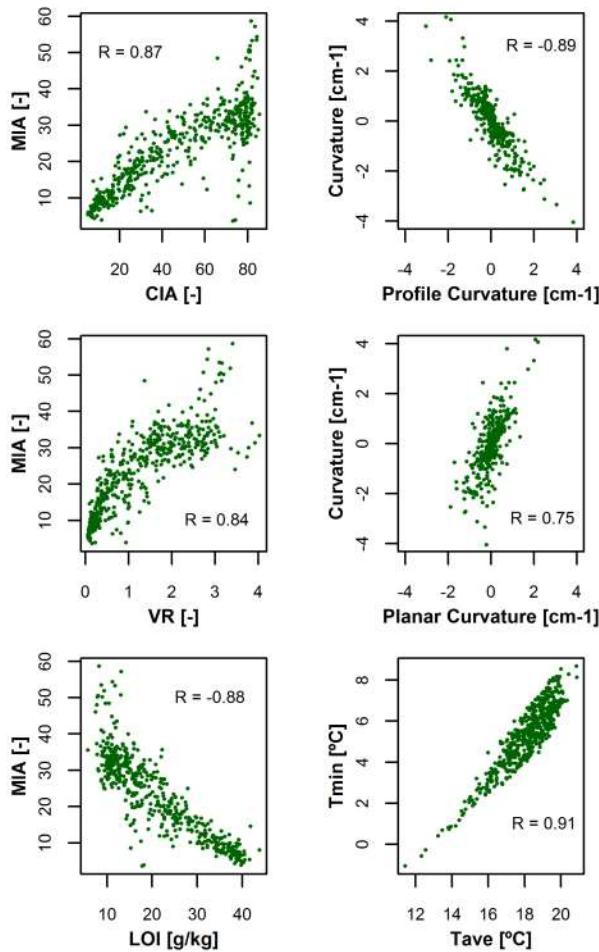


2

3 Fig. 3. Summary of the methods used in the study. OOB is the Out Of Bag error, MESS is the
4 Multivariate Environmental Similarity Surface, and CI is the Confusion Index.

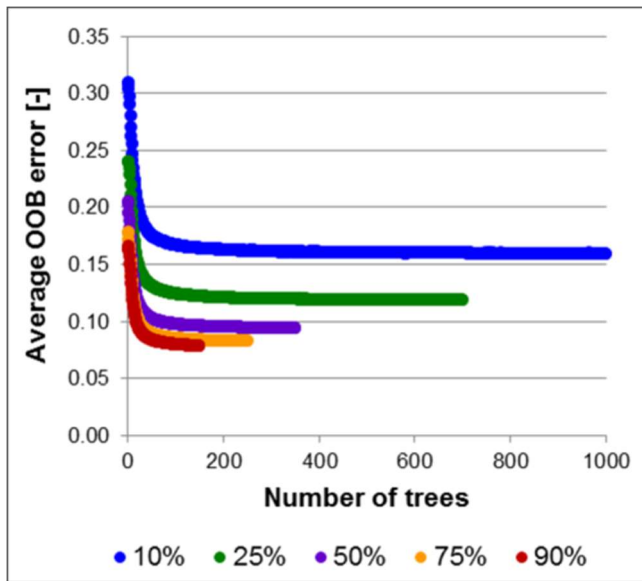
5

1



2

3 Fig. 4. Scatter plots, drawn with a 0.05% subset of the original input dataset, showing highly correlated
4 covariates. The R value shown in each graph is the Pearson correlation coefficient calculated using the
5 full input dataset. The covariates on the x-axis were removed from the data set used as input for the
6 digital mapping techniques. CIA is the chemical index of alteration; MIA is the Mafic Index of
7 Alteration; VR is the Vogt Ratio; and LOI is the Loss on Ignition.



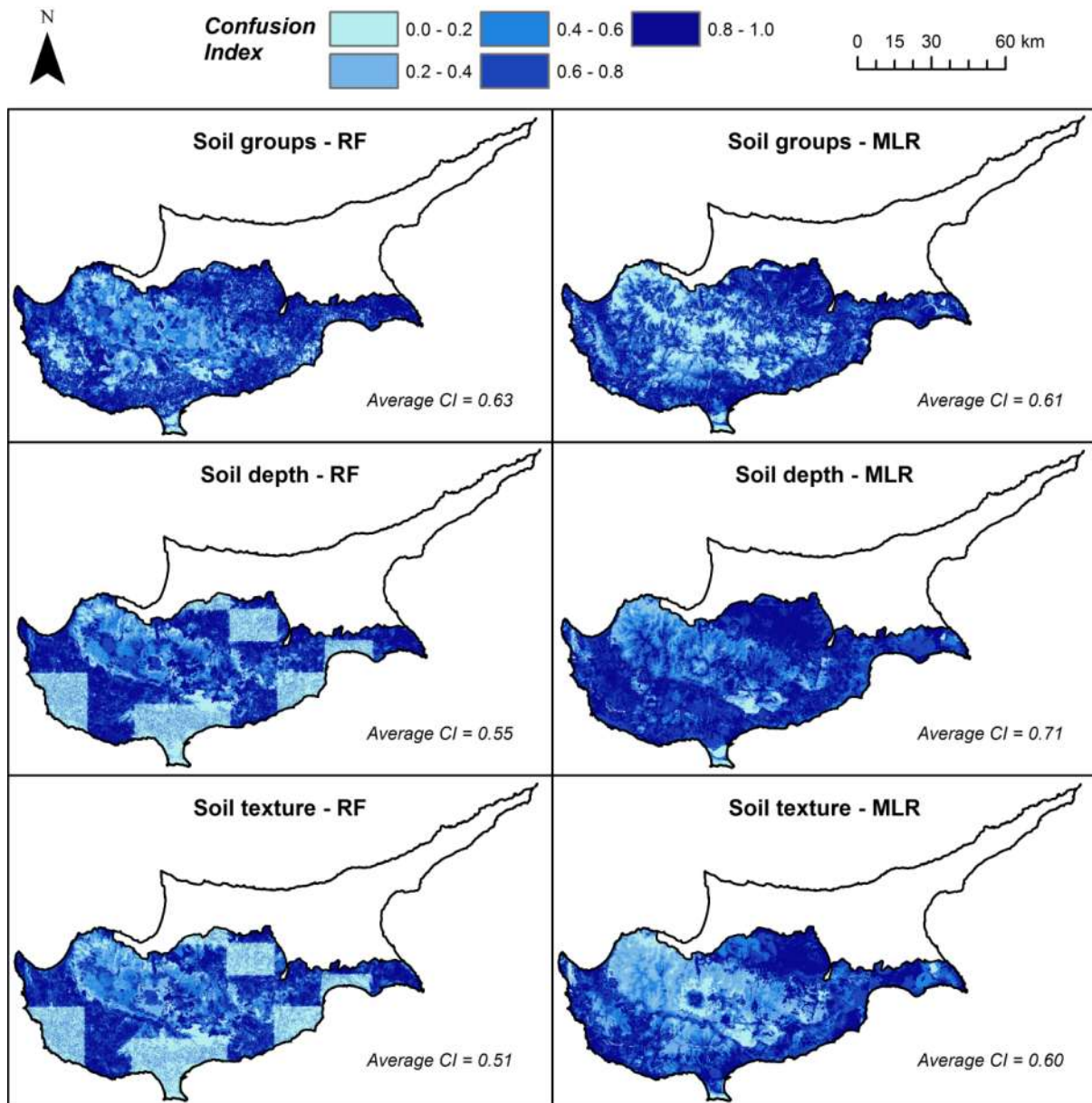
1

2

Fig. 5. Average OOB error for increasing number of trees in the Random Forest classification model and different number (% of the total) of training points as derived from the rasterized 1:25,000 scale soil maps.

4

5



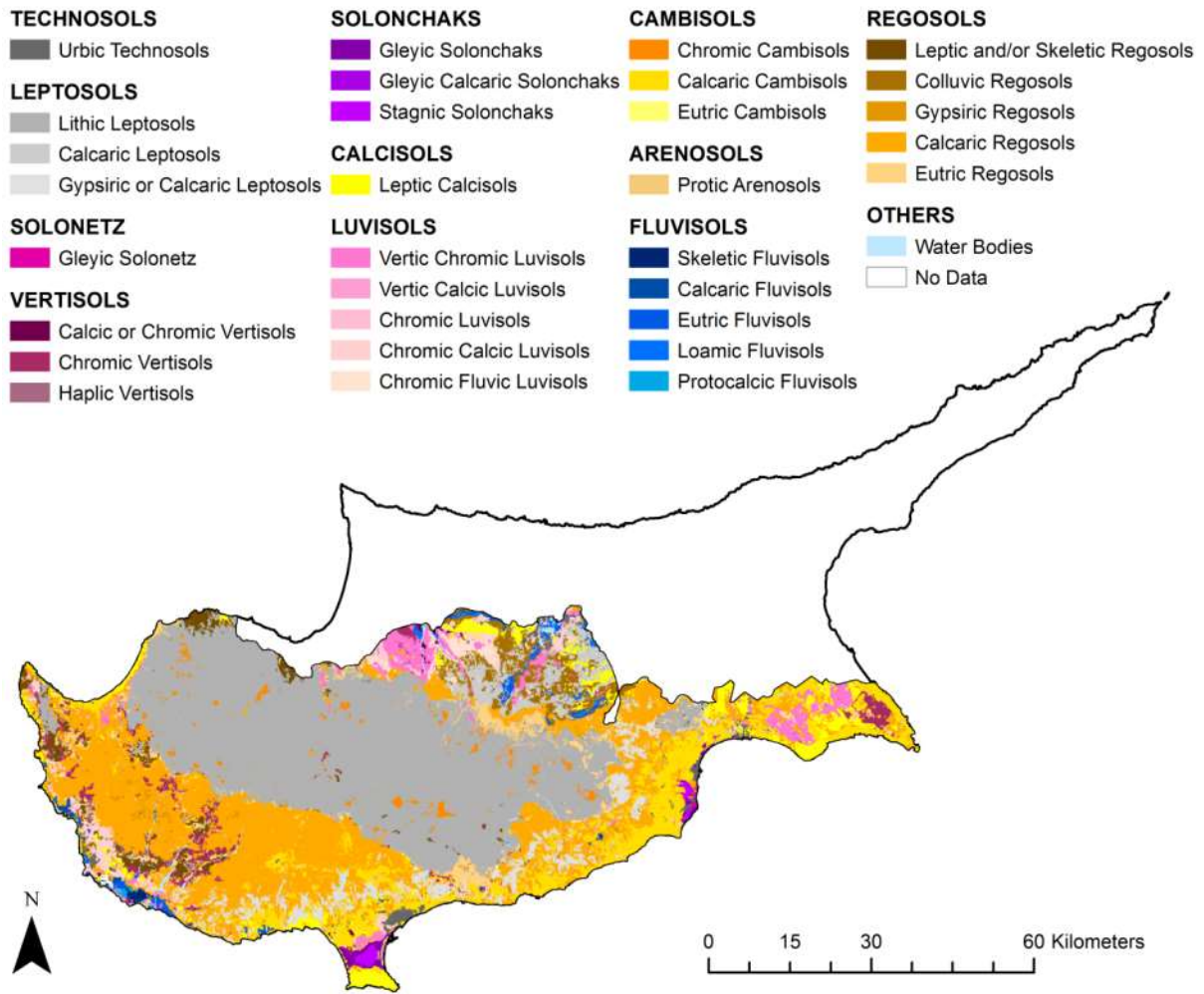
1

2 **Fig. 6. Confusion Index (CI) maps calculated for soil groups, soil depth and soil texture prediction from**
 3 **both Random Forest (left) and Multinomial Logistic Regression (right).**

4

5

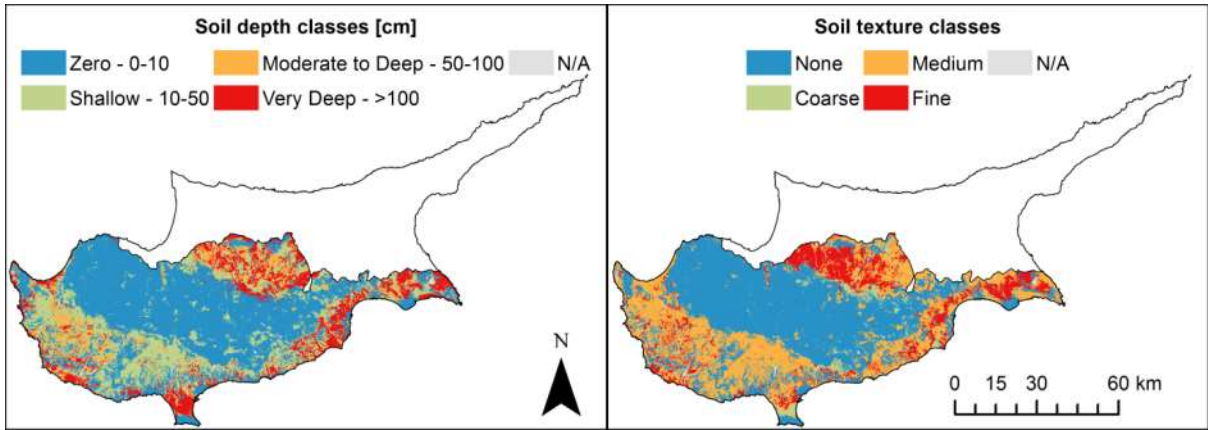
6



1

2 **Fig. 7. Digital soil map of Cyprus. The map is presented with WRB soil group names accompanied by one**
 3 **or two qualifiers, as predicted using Random Forest.**

4



1

2

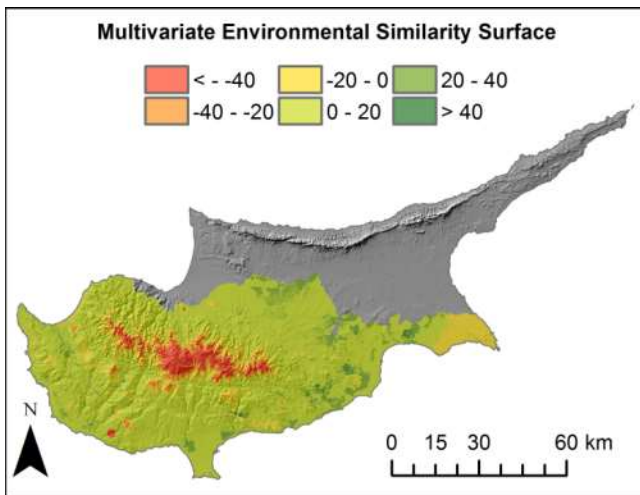
Fig. 8. Digital soil depth and soil texture maps of Cyprus, as derived using Random Forest.

3

4

5

6



7

8

Fig. 9. Multivariate Environmental Similarity Surface (MESS) showing areas where model predictions

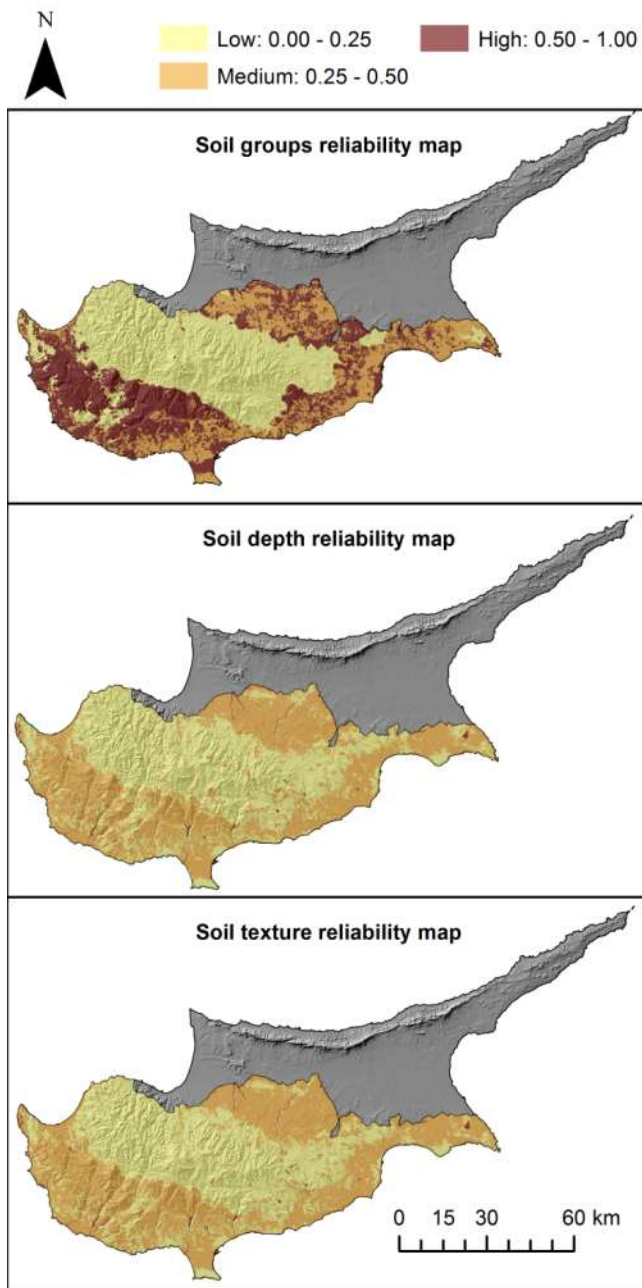
9

are extrapolations in comparison to the training data set (values lower than zero).

10

11

1



2

3 **Fig. 10. Reliability maps derived for the soil groups, soil depth, and soil texture maps of Cyprus predicted**
4 **with Random Forest.**

5

1 **Table 1. Soil depth classes identified from the 1:25,000 scale soil sheets and reclassified in four consistent**
 2 **classes..**

Depth class (from soil maps)	Depth [cm] (from soil maps)	Depth [cm] (new class)
Zero	0-10	0-10 Zero
Very shallow	10-25	10-50
Very shallow to shallow	10-50	
Shallow	25-50	Shallow
Shallow to moderate	25-75	50-100
Moderately deep	50-75	
Deep	75-100	Moderate to Deep
Very deep	> 100	> 100 Very Deep

3

4

1 **Table 2. Soil texture classes identified from the 1:25,000 scale soil sheets and reclassified in four consistent**
 2 **classes. The harmonization has been performed using the guidelines provided in FAO (2008).**

Texture classes (from soil maps)	Texture 4 classes (reclassified FAO, 2008)	AWC¹ [%]
Rock	None	
Gravel	Stony and gravelly	
Sand		
Light to medium		
Coarse to medium		
Coarse	Coarse	5-15
Light	Sand and loamy sand, and sandy loam	
Sandy loam gravelly		
Sandy loam		
Clay loam	Medium	
Medium	Sandy clay loam, loam, clay loam,	10-20
Medium to fine	silty clay loam, silt loam, silt	
Moderately heavy		
Medium heavy		
Heavy	Fine	11-20
Fine	Clay, sandy clay, silty clay	
Clayee		

3 ¹AWC: Available water capacity, range of values taken from Saxton and Rawls (2006) and Allen et al. (1998).
 4
 5

1 **Table 3. Average OOB error for forests (350 trees, 50% data) built removing single variables, and**
2 **importance value (Mean Decrease of Accuracy - MDA) calculated by removing the variable from the**
3 **complete model. OOB values are presented as differences from the model derived with the complete set of**
4 **variables. Between parenthesis the rank of each variable.**

Removed variable	Soil groups		Soil depth		Soil texture	
	OOB	MDA	OOB	MDA	OOB	MDA
	error [%]		error [%]		error [%]	
None	9.4	N/A	12.0	N/A	9.9	N/A
Electrical conductivity	+0.5 (5)	0.22 (6)	+0.6 (5)	0.19 (5)	+0.5 (5)	0.18 (5)
pH	+0.7 (3)	0.20 (7)	+0.8 (3)	0.17 (7)	+0.7 (3)	0.17 (7)
Organic Carbon	+0.9 (2)	0.25 (5)	+1.1 (2)	0.22 (2)	+1.0 (2)	0.21 (3)
Mafic Index of Alteration	+0.7 (3)	0.31 (3)	+0.8 (3)	0.22 (2)	+0.7 (3)	0.24 (2)
Maximum Temperature	-0.3 (12)	0.33 (2)	-0.2 (10)	0.20 (4)	-0.1 (10)	0.21 (3)
Minimum Temperature	-0.2 (10)	0.18 (8)	-0.2 (10)	0.16 (10)	-0.1 (10)	0.15 (9)
Land Use	+0.3 (7)	0.18 (8)	+1.2 (7)	0.17 (7)	+0.2 (7)	0.17 (7)
Elevation (DEM)	+1.9 (1)	0.38 (1)	+2.0 (1)	0.29 (1)	+1.6 (1)	0.26 (1)
Aspect	+0.0 (8)	0.07 (12)	+0.1 (8)	0.07 (12)	+0.1 (8)	0.06 (12)
Curvature	-0.2 (10)	0.01 (13)	-0.2 (10)	0.02 (13)	-0.1 (10)	0.01 (13)
Slope	-0.6 (13)	0.17 (11)	-0.7 (13)	0.17 (7)	-0.4 (13)	0.12 (11)
Landscape Units	+0.4 (6)	0.18 (8)	+0.4 (6)	0.16 (10)	+0.3 (6)	0.14 (10)
Geology	+0.0 (8)	0.27 (4)	+0.1 (8)	0.19 (5)	+0.1 (8)	0.18 (5)

5
6

1 **Table 4. Average OOB error for forests (350 trees, 50% data) built removing groups of variables**
 2 **according to the scorpan formula. Values are presented as differences from the model derived with the**
 3 **complete set of variables. Between parenthesis the rank of each group.**

Removed group	Removed variables	Soil groups	Soil depth	Soil texture
		OOB error [%]	OOB error [%]	OOB error [%]
None	None	9.4	12.0	9.9
Soil chemistry	EC, pH, OrgC, MIA	+13.6 (1)	+13.0 (1)	+11.6 (1)
Climate	Tmin, Tmax	-0.5 (5)	-0.4 (5)	-0.4 (5)
Organisms	Land use, OrgC	+1.3 (3)	+1.5 (3)	+1.3 (3)
Relief	Elevation, aspect, curvature, slope, landscape units	+3.5 (2)	+3.6 (2)	+2.9 (2)
Parent material (Age)	Geology	+0.0 (4)	+0.1 (4)	+0.1 (4)

4

5 **Table 5. Average OOB errors calculated using different *mtry* values. The forest is made up of 350 trees.**

<i>mtry</i>	Soil groups	Soil depth	Soil texture
	OOB error [%]	OOB error [%]	OOB error [%]
1	22.5	26.2	23.0
3	8.7	10.7	9.0
5	8.6	10.5	8.8
7	8.7	10.6	8.9
9	8.8	10.7	9.1

6

7

1 **Table 6. Average OOB (for Random Forest, RF) or MAP (for Multinomial Logistic Regression, MLR,**
 2 **from 50% unused points in 1:25,000 soil maps) errors, validation hit rate (HR) from profiles (not**
 3 **applicable to soil groups for lack of profile data) and confusion index (CI) calculated using different**
 4 **nodesize values (n). The forest is made up of 350 trees and the *mtry* parameter is fixed to 5.**

<i>Method</i>	Soil groups			Soil depth			Soil texture		
	OOB/MAP [%]	CI [-]	HR [%]	OOB/MAP [%]	CI [-]	HR [%]	OOB/MAP [%]	CI [-]	HR [%]
RF n1	8.6	0.63	N/A	10.5	0.55	55	8.8	0.51	49
RF n4	8.7	0.64	N/A	10.7	0.56	55	9.0	0.51	49
RF n12	9.7	0.63	N/A	11.8	0.57	53	9.9	0.52	46
RF n20	10.6	0.63	N/A	12.8	0.57	54	10.6	0.53	53
MLR	48	0.61	N/A	52	0.71	58	44	0.60	46

5

6

7

8

1 **Table 7. Number of independent soil profiles (N. prof.) and validation hit rate (HR), calculated based on**
 2 **positive or negative Multivariate Environmental Surface Similarity (MESS), for soil depth classes and soil**
 3 **texture classes.**

Soil depth					Soil texture				
Class	N. prof.	N. prof.	HR	HR	Class	N. prof.	N. prof.	HR	HR
[cm]	MESS+	MESS-	MESS+	MESS-		MESS+	MESS-	MESS+	MESS-
0-10	80	2	0.97	1.00	None	5	0	1.00	N/A
10-50	129	1	0.43	0.50	Coarse	11	5	0.09	0.20
50-100	17	1	0.03	0.00	Medium	47	8	0.74	0.50
> 100	76	19	0.43	0.42	Fine	42	8	0.31	0.37

4

5

6