# Performance of linear mixed models and random forests for spatial prediction of Soil pH.

Mirriam Makungwe[1], Lydia Mumbi Chabala[1],  Benson H Chishala[1], R Murray Lark[2]

[1]Department of Soil Science, University of Zambia, School of Agricultural Sciences, P.O. Box 32379, Lusaka, Zambia.

[2]School of Biosciences, University of Nottingham, Sutton Bonington, Nottinghamshire, LE12 5RD, UK.

Email address:

mirriammakungwe.tolopu@gmail.com (M. Makungwe), lchabala@unza.zm (L. M. Chabala), bchishala@unza.zm (B. H. Chishala), Murray.Lark@nottingham.ac.uk ( R Murray Lark).

**Author notes**

Correspondence concerning this article should be addressed to Mirriam Makungwe, Department of Soil Science, School of Agricultural Sciences, University of Zambia P.O. Box 32379, Lusaka, Zambia. E-mail: mirriammakungwe.tolopu@gmail.com

**ABSTRACT**

Digital soil maps describe the spatial variation of soil and provide important information on spatial variation of soil properties which provides policy makers with a synoptic view of the state of the soil. This paper presents a study to tackle the task of how to map the spatial variation of soil pH across Zambia. This was part of a project to assess suitability for rice production across the country. Legacy data on the target variable were available along with additional exhaustive environmental covariates as potential predictor variables. We had the option of undertaking spatial prediction by geostatistical or machine learning methods. We set out to compare the approaches from the selection of predictor variables through to model validation, and to test the predictors on a set of validation observations. We also addressed the problem of how to robustly validate models from legacy data when these have, as is often the case, a strongly clustered spatial distribution. The validation statistics results showed that the empirical best linear unbiased predictor (EBLUP) with the only fixed effect a constant mean (ordinary kriging) performed better than the other methods. Random forests had the largest model-based estimates of the expected squared errors. We also noticed that the random forest algorithm was prone to select as "important" spatially correlated random variables which we had simulated.

**Keyword**: Linear Mixed Models; REML-EBLUP; Random forests, Spatial prediction of soil pH.

## 1. INTRODUCTION

Soil maps describe the spatial variation of soil types and provide important information on spatial variation of soil properties (Kempen et al., 2010). Mapping of soil properties is important as it provides policy makers with a synoptic view of the state of the soil, and agricultural stakeholders with information about where soil problems might occur (Lark et al., 2019). Soil maps are

2

54 generated using various soil mapping methods which can be divided into conventional and

55 pedometric approaches (Kienast-Brown et al., 2010; Hengl, 2003).

56 Conventional soil survey represents soil variation in terms of profile classes and corresponding

57 map legend units. It can provide a basis for spatial prediction of soil properties and may also serve

58 as a structure for recording substantial information on soil management and for systematizing

59 knowledge of the distribution of soils in the landscape. Conventional approaches were based

60 largely on manual processes which are costly and time consuming (Kienast-Brown et al., 2010)

61 mainly because of long fieldwork periods (Moonjun et al., 2010). Pedometric approaches are based

62 on the application of mathematical and statistical methods for the primary purpose of predicting

63 the values of soil properties where these have not been observed directly (McBratney et al., 2000).

64 A well-established statistical approach to doing this is the application of model-based geostatistics

65 (Stein, 1999; Diggle and Ribeiro, 2007). In this approach the variation of the soil is represented in

66 a linear mixed model (LMM) as a combination of fixed effects (which may be a constant unknown

67 mean, or a function of predictive covariates such as remote sensor data), and random effects,

68 including Gaussian random fields which exhibit spatial correlation. The parameters of the LMM

69 model can be estimated by Residual Maximum Likelihood (REML) method developed by

70 Patterson and Thompson (1971), which allows parameters of the random effects to be estimated

71 with small bias arising from uncertainty in the fixed effects (Kitanidis, 1987; Swallow and

72 Monahan 1984; Zimmerman and Zimmerman, 1991; Lark and Cullis, 2004) . When the model is

73 fixed, values of the soil property at unsampled sites can be obtained by the empirical best linear

74 unbiased predictor (EBLUP) (Stein, 1999; Lark et al., 2006; Lark and Webster, 2006; Minsay and

75 McBratney, 2007).

76 There has been a growing interest in the potential of machine learning methods (e.g. Breiman,

77 2001) as an alternative to statistical modelling for spatial prediction of soil properties (Hengl et al.,

78 2015; Behrens and Scholten, 2007). The main difference between geostatistical approaches and

79 random forest is that geostatistics is based on a statistical model. This provides a basis for formal

80 inference about the validity of the model (including the task of selecting which covariates to use in

81 prediction), and for producing a prediction distribution at unsampled sites of interest. One may

82 then derive point predictions from this distribution (typically the mean), and measures of

83 uncertainty. On the other hand, machine learning methods such as random forests, are predictive

84 tools applied to identify empirical relationships between the target variable in a training data set

85 and associated predictive covariates and to extrapolate these to unsampled sites. With no model

86 there can be no formal inference, but empirical approaches, based on internal cross-validation are

87 used, for example, to evaluate the evidence that a particular variable is predictive. One particular

88 strength of the geostatistical approach is that the estimation of coefficients for predictor variables,

89 and inferences about them, are based on a model of the spatial dependence of the random variation.

90 This accounts for the fact that data which are strongly spatially clustered are likely to be correlated,

91 and so do not provide independent evidence to support the fitted model.

92 In the study reported here our objective was to assess approaches for digital mapping of soil pH at

93 national scale across Zambia to support evaluation of land potential for rice production. Legacy

94 data on soil pH were available from a previous national survey. As with many such surveys, this

95 followed a two-stage design, and so the observations were spatially clustered. In addition we had

96 access to various exhaustive environmental covariates which could be regarded as potential

97 predictor variables for soil pH.. We compared different forms of linear mixed model, and

98 prediction with the random forest using a validation subset of the data. Prediction errors were

99 evaluated at the validation locations by comparing predictions with observed values. The selection

100 of the validation subset, and the quantification of the uncertainty from the observed prediction

101 errors had to take account of the spatial clustering of the observations in the legacy data. Because

102 of this clustering, no subset could be regarded as independent random observations.

103 **2. THEORY**

## 2.1 Linear Mixed Model

The theory of residual maximum likelihood (REML) in combination with the empirical best linear unbiased predictor (EBLUP) for spatial interpolation has been illustrated and described in detail by Lark et al., 2006. The LMM takes the form

$$\mathbf{z} = \mathbf{M}\boldsymbol{\beta} + \mathbf{S}\boldsymbol{\eta} + \boldsymbol{\varepsilon}, \tag{1}$$

where $\mathbf{z}$ is a set of observations of the random variable at sampled locations, $\mathbf{M}$ is the design matrix of fixed effects, which could include covariates such as topographic attributes, $\boldsymbol{\beta}$ is the vector of the fixed effects parameters or regression coefficients, $\mathbf{S}$ is the design matrix of random effects (which is an identity matrix unless analytical duplicate observations are included), $\boldsymbol{\eta}$ is a random effect, a Gaussian random variable which has a mean of zero and, in the spatial setting, a covariance matrix which expresses spatial dependence, $\boldsymbol{\varepsilon}$ is an independently and identically distributed Gaussian residual of mean zero and variance $\sigma^2$. These two random components have a joint distribution

$$\begin{bmatrix} \boldsymbol{\eta} \\ \boldsymbol{\varepsilon} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 \xi \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \sigma^2 \mathbf{I} \end{bmatrix} \right), \tag{2}$$

where $\mathbf{I}$ is the identity matrix and $\mathbf{G}$ is the correlation matrix of the random effect $\boldsymbol{\eta}$. Element $[i,j]$ of $\mathbf{G}$, at locations $\mathbf{x}_i$ and $\mathbf{x}_j$ depends only on the interval in space between them under an assumption of second-order stationarity. The lag vector $\mathbf{x}_i - \mathbf{x}_j$, under the assumption of isotropy, depends only on the scalar part of this vector, the lag distance and so

$$\mathbf{G}[i,j] = \rho\big(|\mathbf{x}_i - \mathbf{x}_j|; \alpha\big), \tag{3}$$

5

126      where $\rho(h; \alpha)$ is a correlation function of lag distance $h$ with spatial parameters $\alpha$ which control

127      how the correlation decreases with increasing distance. The term $\xi$ is the ratio of the variance

128      of the random effect $\boldsymbol{\eta}$ to $\sigma^2$, the variance of the residual term.

129      The residuals depend on the fixed effects parameters $\boldsymbol{\beta}$ in the model, and in ordinary maximum

130      likelihood estimation the uncertainty in the estimates of the fixed effects parameters biases the

131      estimates of the random effects parameters. To avoid this, we use residual maximum likelihood

132      (REML) which is based on the principle that a new random variable, independent of the fixed

133      effects, is computed by projecting the original data $\mathbf{z}$ into a residual space where the fixed

134      effects can be filtered out (Chai et al., 2008). The log likelihood of the new random variable

135      which we now call the residual log-likelihood because its independent of fixed effects can be

136      expressed as;

137
$$\ell_R(\sigma^2, \xi, \alpha | \mathbf{z}) = -\frac{1}{2}\{\log|\mathbf{H}| + \log|\mathbf{M}^T\mathbf{H}\mathbf{M}| + (n-p)\sigma^2 + \frac{1}{\sigma^2}\mathbf{z}^T(\mathbf{I} - \mathbf{W}\mathbf{C}^{-1}\mathbf{W}^T)\mathbf{z}, \quad (4)$$

138      Where $\mathbf{H} = \xi\mathbf{M}\mathbf{G}\mathbf{Z}^T + \mathbf{I}$, $\mathbf{W} = [\mathbf{M}, \mathbf{S}]$ and $\mathbf{C} = \begin{bmatrix} \mathbf{M}^T\mathbf{M} & \mathbf{M}^T\mathbf{S} \\ \mathbf{S}^T\mathbf{M} & \mathbf{S}^T\mathbf{S} + \xi^{-1}\mathbf{G}^{-1} \end{bmatrix}$.

139      Once the covariance parameters $\sigma^2, \xi, \alpha$ have been estimated by REML, they are used to

140      compute the estimated covariance matrix at sampled points. With the estimated covariance

141      matrix computed, the estimated fixed effects parameter, $\widehat{\boldsymbol{\beta}}$, and predicted random effects, $\widetilde{\boldsymbol{\eta}}$, are

142      then computed by solution of mixed model equation:

143
$$\mathbf{C}\begin{bmatrix} \widehat{\boldsymbol{\beta}} \\ \widetilde{\boldsymbol{\eta}} \end{bmatrix} = \begin{bmatrix} \mathbf{M}^T\mathbf{z} \\ \mathbf{S}^T\mathbf{z} \end{bmatrix} \quad (5)$$

144      With the covariance matrix for the error of the estimates being:

145
$$\text{Cov}\begin{bmatrix} \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \widetilde{\boldsymbol{\eta}} - \boldsymbol{\eta} \end{bmatrix} = \sigma^2\mathbf{C}^{-1}. \quad (6)$$

146 Now that the covariance matrix and the fixed effects parameters have been estimated, they are

147 used in EBLUP to predict the soil property , $\tilde{z}_p$ , at unsampled locations:

$$\tilde{z}_p = \mathbf{M}_\mathrm{p}^\mathrm{T}\widehat{\boldsymbol{\beta}} + \widetilde{\boldsymbol{\eta}}_\mathbf{p} = \mathbf{M}_\mathrm{p}^\mathrm{T}\widehat{\boldsymbol{\beta}} + \mathbf{g}_{\mathrm{o,p}}^\mathrm{T}\mathbf{G}^{-1}\widetilde{\boldsymbol{\eta}}, \tag{7}$$

149 where $\mathbf{M}_\mathrm{p}$ is the design matrix for the prediction sites, $\mathbf{g}_{\mathrm{o,p}}$ is a vector computed from the

150 covariance matrix of $\boldsymbol{\eta}$ with the $\boldsymbol{\eta}_\mathrm{p}$ values at the unsampled locations $(\mathrm{Cov}[\boldsymbol{\eta},\boldsymbol{\eta}_\mathrm{p}] = \xi\sigma^2\mathbf{g}_{\mathrm{o,p}})$.

151 The variance of the prediction errors, $Var[\tilde{z}_p - z_p]$, which accounts for the uncertainty in

152 predicting the fixed effects and uncertainty in predicting the random effects is expressed as:

153 $$\mathrm{Var}[\tilde{z}_p - z_p] \;\; = \sigma^2\left\{[\mathbf{M}_\mathrm{p}, \mathbf{g}_{\mathrm{o,p}}^\mathrm{T}\mathbf{G}^{-1}]^\mathrm{T}\mathbf{C}^{-1}[\mathbf{M}_\mathrm{p}, \mathbf{g}_{\mathrm{o,p}}^\mathrm{T}\mathbf{G}^{-1}] + \xi(g_{\mathrm{p,p}} - \mathbf{g}_{\mathrm{o,p}}^\mathrm{T}\mathbf{G}^{-1}\mathbf{g}_{\mathrm{o,p}}) + 1\right\}. \tag{8}$$

154 There are many variables that researchers can use as fixed effects in linear mixed models for

155 spatial prediction of soil properties. However, it is unwise to include variables without regard

156 for evidence that they are of predictive value, the inclusion of predictors unrelated to the target

157 variable may inflate the prediction error variance. To avoid this, variable selection is an

158 important step in model development. One approach to the problem is to base the inclusion or

159 rejection of a predictor based on a hypothesis test in the LMM framework (e.g. by a log-

160 likelihood ratio test) (Verbeke and Molenberghs, 2000). To reduce the risk of including excess

161 predictors because of multiple hypothesis testing, one may use false discovery rate control

162 (Lark, 2017). The false discovery rate (FDR) is the probability that a null hypothesis is true,

163 given that it has been rejected. False discovery rate control can reduce the power to detect real

164 predictors, and Lark, (2017) demonstrated how this problem can be reduced, while maintaining

165 FDR control, by the method of alpha investment (Foster and Stine, 2008). This entails an initial

166 ordering of the predictors starting with the one which, a priori (and not based on inspection of

167 the data) is thought most likely to be of predictive value and adding in predictors in declining

168 order of expected predictive power. In this approach the power to detect a predictor is increased
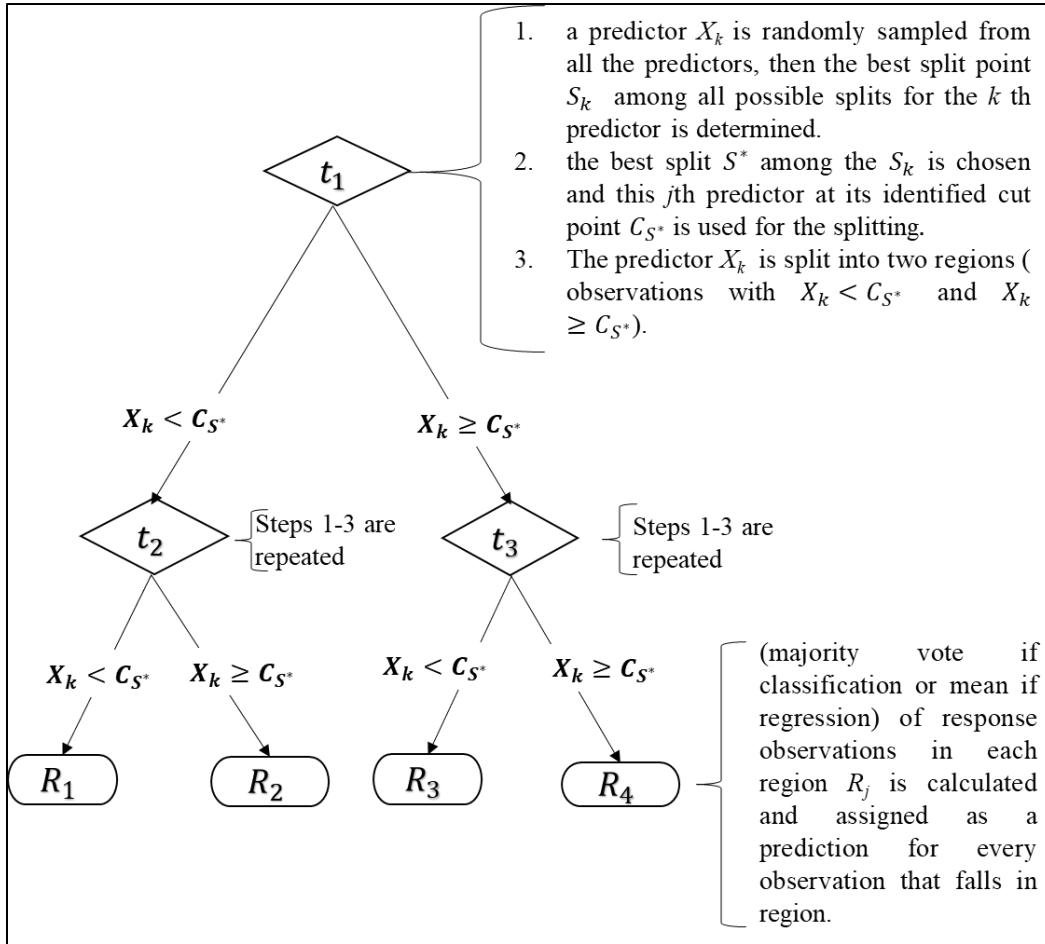
169 by the rejection of null hypotheses early in the sequence while maintaining control of FDR.

170 This approach has been used elsewhere for spatial prediction (Gashu et al., 2020).

171 A disadvantage of the LMM approach is that it assumes that the fixed effects are linear in the

172 parameters. Such a model can represent complex and non-linear relations between soil

173 properties and predictors, for example through the use of polynomial terms in the predictor

174 variables, or spline basis functions, but there has been increasing interest in more flexible

175 methods to predict soil properties from covariates, in particular the machine learning method

176 known as the random forest.

## 177 2.2 Random Forests

178 The random forest is an ensemble tree-based method that combines multiple decision trees

179 (classification or regression) to give a prediction (Breiman, 2001). A decision tree is an

180 algorithm that involves recursive partitioning of data into several simple regions using a series

181 of splitting rules. It is called a decision tree because these series of splitting rules can be

182 summarised into an upside-down tree structure as illustrated in Figure 1. Figure 1 shows a

183 structure made up of predictors ($X_1, X_2, ... ..., X_k$) which are split into $J$ distinct and non-

184 overlapping regions ($R_1, R_2, ..., R_j$) at test node $t$, and the mean of the response values for the

185 training observations in each region $R_j$ is calculated and assigned as a prediction for every

186 observation that falls in region $R_j$ (James et al., 2013). When growing a decision tree, the

187 following steps are taken; (1) at each test node $t$, a predictor $X_k$ is randomly sampled from all

188 the predictors, then the best split point $S_k$ among all possible splits for the $k$ th predictor is

189 determined; (2) the best split $S^*$ among the $S_k$ is chosen and this $j$th predictor at its identified

190 cut point $C_{S^*}$ is used for the splitting at test node $t$. (3) The predictor $X_k$ is split into two regions

191 ( observations with $X_k < C_{S^*}$ and $X_k \geq C_{S^*}$) at test node $t$. Steps 1-3 are repeated on all

192 descendant nodes to grow a tree $\hat{f}(x)$ (Archer and Kimes, 2008; James et al., 2013).

The figure contains the following annotations:

1. a predictor $X_k$ is randomly sampled from all the predictors, then the best split point $S_k$ among all possible splits for the $k$ th predictor is determined.
2. the best split $S^*$ among the $S_k$ is chosen and this $j$th predictor at its identified cut point $C_{S^*}$ is used for the splitting.
3. The predictor $X_k$ is split into two regions ( observations with $X_k < C_{S^*}$ and $X_k \geq C_{S^*}$).

$t_1$

$X_k < C_{S^*}$   $X_k \geq C_{S^*}$

$t_2$ — Steps 1-3 are repeated

$t_3$ — Steps 1-3 are repeated

$X_k < C_{S^*}$   $X_k \geq C_{S^*}$   $X_k < C_{S^*}$   $X_k \geq C_{S^*}$

$R_1$   $R_2$   $R_3$   $R_4$

(majority vote if classification or mean if regression) of response observations in each region $R_j$ is calculated and assigned as a prediction for every observation that falls in region.

193

194    *Figure 1: illustration of a decision tree*

195    One major limitation with decision trees is that using only one tree for prediction, results in

196    highly unstable predictions, a small change in the data can result into a large change in the final

197    estimated tree. To improve the performance of decision trees, Breiman (1996) introduced an

198    algorithm called Bagging, also known as bootstrap aggregation which takes repeated

199    (bootstrap) samples (where $B$ is the number of bootstrap samples) from training set with

200    replacement and builds a total of $B$ trees ($\hat{f}^1(x), \hat{f}^2(x), \ldots \ldots \hat{f}^b(x)$) which the average of all the

201    prediction trees $\hat{f}_{bag}(x)$ is calculated:

202
$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^b(x)$$
(9)

203  One disadvantage of bagging is that a single predictor may dominate all trees in the bag,

204  meaning that their outputs are strongly correlated.  As a result of this the reduction in variance

205  from the use of multiple trees is very limited (James et al., 2013).  To address this, Breiman

206  (2001) developed a random forest algorithm which is an improvement of bagging. Like

207  bagging, random forest also takes repeated (bootstrap) samples from the training data and

208  builds $B$ decision trees. But in the case of random forest, when building these trees, to avoid

209  using one strong predictor for all bagged trees, at every test node $t$, when splitting, every bagged

210  tree is forced to consider only a random subset of predictors by randomly sampling a fresh $m$

211  predictors from a set of $k$ predictors, and the split is only allowed to use  one of these $m$

212  predictors. For regression trees, the number of $m$ predictors considered at each split is

213  approximately the total number of predictors divided by three ($m \approx {k}/{3}$ ) and for classification

214  trees, $m \approx \sqrt{k}$. Because of this, random forest results in many uncorrelated trees which give a

215  large reduction in variance when averaged.

216  The random forest algorithm has three important outputs. These are the out-of-bag Mean

217  Squared Error (OOB Mean Squared Error), the out-of-bag R-squared and the variable

218  importance. The RF model does not use all the data for building the tree. In each bootstrap

219  training set, about one-third of the data are left out. The data that are left out when building the

220  trees is called out-of-bag (OOB) data and after the trees are grown, the OOB data are used as

221  test set to measure the strength (OOB Mean Squared Error)  and correlation (OOB R-squared)

222  of the model. In short, random forest has an inbuilt cross-validation. Variable importance is

223  defined as the increase in prediction error when OOB data for that variable is randomly

224  permuted while all others are left unchanged (Liaw, 2002). It analyses the contributions of each

225  predictor to the overall results (Breiman, 2001). The algorithm randomly permutes the predictor

226  $X_m$ several times, breaking its original association with the response variable and asses the

227  relevance of the predictor by using the permuted predictor together with the other unpermuted

10

228 predictors to predict the response variable for the out-of-bag observations giving the difference

229 in prediction accuracy before and after permuting. The result is a vector of importance measures

230 for each predictor equivalent to the number of permutations. The algorithm then computes a p-

231 value as a measure of the evidence that a variable is predictive (Strobl et al., 2007; Altmann et

232 al., 2010). This permutation p-value is the probability of observing a permuted model (from

233 the several number of permutations) that is equal to or better than the unpermuted model

234 (Cummings et al., 2004).

235 Equation (7) presents the E-BLUP from a linear mixed model. The second term on the right-

236 hand side corresponds to the spatial interpolation of the correlated random effect in the model.

237 In this way the E-BLUP combines a regression-type prediction based on the predictor variables

238 with a spatial interpolation component. As described above the random forest predicts a soil

239 property from the predictor variables only, making no use of spatial dependence through

240 interpolation. An attempts has been made to include spatially weighted local observations in

241 prediction by random forests by including coordinates as predictors and using weighted buffer

242 distances (Hengl et al., 2018), neighbouring observations and their distances to the prediction

243 location were used by Sekulić et al., (2020). Li et al., (2011) and Viscarra Rossel et al., (2014)

244 combined random forest with kriging, just like in regression-kriging, by calculating the random

245 forest residuals and then kriged them to all prediction positions and then added to the results of

246 the prediction positions. .

247 As described above, inferences in the random forest approach are based on an internal cross-

248 validation procedure. This might lead to overoptimistic conclusions about a random forest

249 model, or about the value of a particular predictor if observations from the same clusters appear

250 in the OOB sample and in the data used to develop the trees. That is because the observations

251 within a cluster can be expected to be strongly correlated, and so the validation of a model fitted

252    to data on strongly correlated observations will give an unduly optimistic impression of the

253    model's capacity to predict at an independent location.
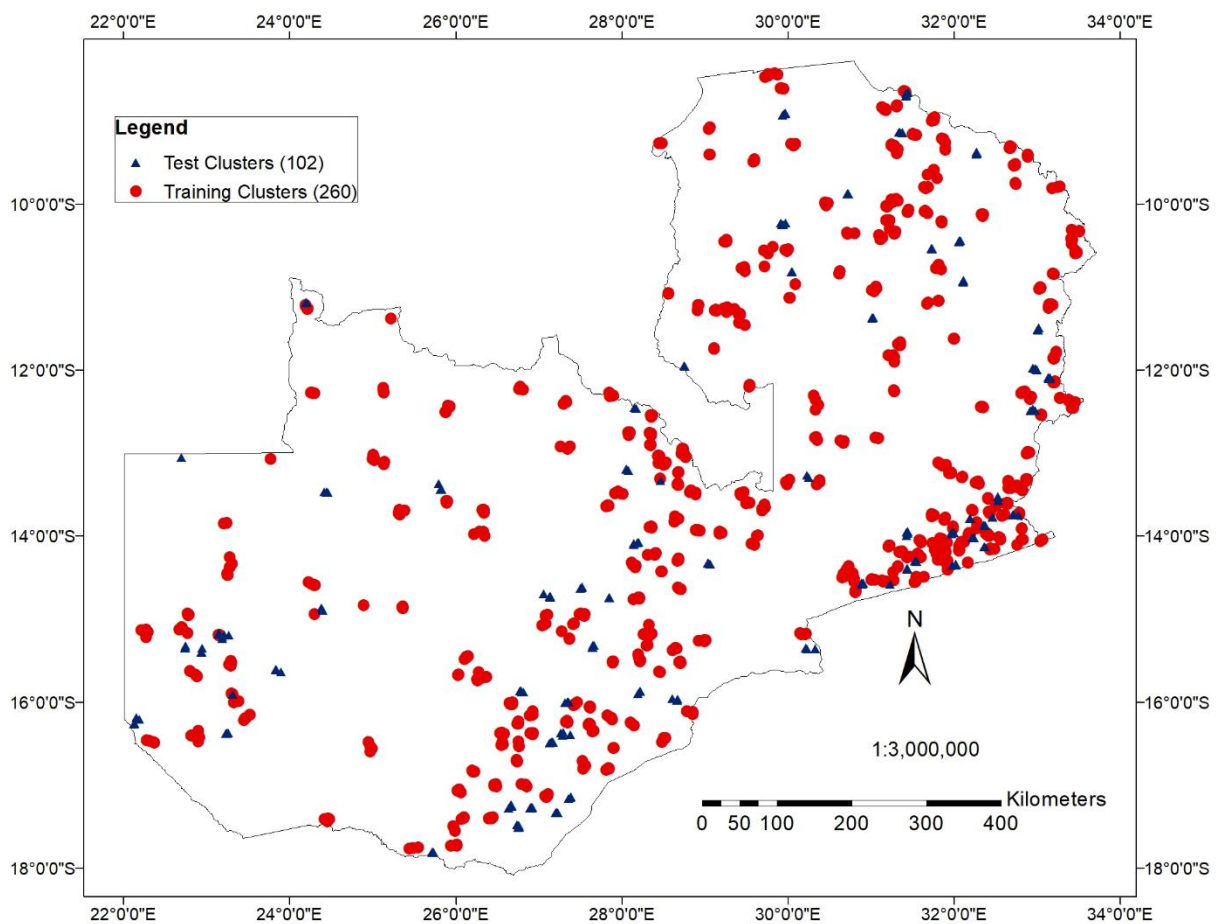
254    **3.  CASE STUDY**

255    **3.1.    Data**

256    **Soil data**

257    This case study uses Rural Agricultural Livelihoods Survey (RALS) of 1713 soil pH data

258    collected by Indaba Agricultural Policy Research Institute (IAPRI) in collaboration with

259    Central Statistical Office (CSO) and Ministry of Agriculture. The sampling frame for the RALS

260    2012 survey was based on the 2010 Census of Housing and Population (CSO/MAL/IAPRI,

261    2015). Full detail of the stratified two-stage sampling design is provided by (CSO 2012). Four

262    households were randomly selected in each Standard Enumeration Area (SEA) and soil samples

263    were collected from the largest maize field. A composite of 10–20 sub-samples of soil collected

264    throughout each field and each sub-sample was a composite of equal parts soil in the 0–10cm

265    and10–20cm depth horizons.  Full details on the soil collection and laboratory analysis for soil

266    pH (determined for a soil suspension in $CaCl_2$ with a standard pH meter) are provided by Burk

267    et al., (2019) and Chapoto et al., (2016). The spatial prediction of soil pH for Zambia using this

268    data has been studied by Chapoto et al., (2016) who only used ordinary kriging for the

269    prediction.

270    Data cleaning involved removal of spurious values in the x and y coordinates.  The need for

271    this was indicated when the raw data were first plotted, showing points lying outside the borders

272    of Zambia. The mean coordinates of all households were computed in each Standard

273    Enumeration Area (SEA) centroid, and then the households were removed from the data set if

274    the notional distance to the SEA centroid exceeded 10km.  This threshold value was decided

275    after discussion with IAPRI staff about plausible values for the distance between a village in

276    the EA and the centroid.  After data cleaning, a total of 1202 soil samples were used for analysis.

12

277    The sampling pattern for the RALS survey was not designed for spatial interpolation of soil.

278    Due to the sampling pattern, the data is strongly clustered at the level of SEAs (the SEAs are

279    the clusters) with a total of 362 clusters. For this reason, splitting of the dataset into training set

280    (80%) and test set (20%) was done at cluster level (the 362 clusters were split into 260 (80%)

281    for training and 102 (20%) for validation).  Figure 2 Shows the training and test clusters with

282    the red solid dots being the training clusters  and the blue solid triangles being the test clusters.

283



285    *Figure 2: Cluster locations for the RALS 2012 soil data. red solid dots being the training*

286    *clusters used for spatial prediction of soil pH and the blue solid triangles being the test clusters*

287    *left out for validation.*

288    **Environmental Covariates**

289 The environmental predictors available for use in this study were Soilclass, Landcover, mean

290 annual rainfall, elevation, slope, aspect, valley depth, LS-Factor (a combination of slope and

291 slope length, relative slope position (RSP), channel network base Level (CNBL) and

292 Normalized difference vegetation index (NDVI).

293 Soilclass information was obtained from the 1:1,000,000 scale exploratory soil map of Zambia

294 compiled by the Zambian Ministry of Agriculture, Zambia Agricultural Research Institute

295 (ZARI) - Soil Survey Section in 1991(GRZ, 1991) and then digitized to raster format. Map

296 units are allocated to suborders of the FAO-UNSECO classification as used in the Third Draft

297 of the legend to the Soil Map of the World (Jahn et al, 2006)).  A total of 96 soil classes were

298 represented in the data available for model development, but these do not comprise all the

299 classes on the map of Zambia, and so some generalization is required to develop models for

300 spatial prediction. We, therefore, reduced the number of classes, by aggregating the classes

301 from suborder to order level, this reduced the number of classes to 18 and all the classes in the

302 prediction grid where represented in the training set. Land cover data for the years between

303 2000 and 2015 with spatial resolution of 300m were downloaded from the European Space

304 Agency (ESA), (2017). The data presented a similar situation as that of soilclass with landcover

305 classes in the prediction sites not being represented in the training set. We also reduced the

306 number of landcover classes, by aggregating them as shown in Table 1.

307 *Table 1: Aggregated landcover classes based on ESA, (2017)*

| New Class | ESA Class | Description |
|---|---|---|
| | 10 | rainfed cropland |
| | 20 | irrigated or post-flooding cropland |
| 1 | 30 | Mosaic cropland (>50%) / natural vegetation (tree |
| | 11 | Herbaceous cover |
| | 40 | herbaceous cover) (>50%) / cropland (<50%) |
| 2 | 110 | Mosaic herbaceous cover (>50%) / tree and shrub (<50%) |
| 3 | 12 | Tree or shrub cover |

| | | 100 | Mosaic tree and shrub (>50%) / herbaceous cover (<50%) |
|---|---|---|---|
| 4 | | 50 | closed to open (>15%), evergreen, broadleaved, tree cover |
| | | 60 | closed to open (>15%), deciduous, broadleaved, tree cover |
| | | 61 | closed (>40%), deciduous, broadleaved, tree cover |
| | | 62 | open (15-40%), deciduous, broadleaved, tree cover |
| | | 120 | Shrubland |
| 5 | | 122 | Shrubland deciduous |
| 6 | | 130 | Grassland |
| 7 | | 160 | fresh or brackish water, flooded, tree cover |
| | | 170 | saline water, flooded, tree cover |
| | | 180 | fresh/saline/brackish water, flooded Shrub or herbaceous cover |
| 8 | | 190 | Urban areas |
| 9 | | 200 | Consolidated bare areas |
| | | 202 | Unconsolidated bare areas |

308

309 Mean annual rainfall data (averages from 1970 to 2000) with a spatial resolution of 1km was

310 downloaded from WorldClim website (Fick and Hijmans, 2017). A 90-m resolution NASA

311 Shuttle Radar Topography Mission (SRTM3) Digital elevation model (DEM) was downloaded

312 from USGS (2019) and projected to WGS 84 UTM 35 S. The DEM was pre-processed by filling

313 sinks using the fill sinks (Planchon/Darboux, 2001) tool in Saga GIS, and then elevation, slope,

314 aspect, valley depth, LS-Factor (a combination of slope and slope length), relative slope

315 position (RSP) and channel network base Level (CNBL) data was extracted from the DEM

316 using basic terrain analysis tool in Saga GIS. MODIS land surface reflectance (MOD009GA

317 V6) was downloaded from USGS (2019). After downloading the respective data sets, Quantum

318 GIS was used to project the data sets to WS 84 UTM 35s and then converted to the Integrated

319 Land and Water Information System (ILWIS) format. Then ILWIS was used to harmonise all

320 the raster files to the same extent and cell size of 1km. Normalized difference vegetation index

321 (NDVI) was extracted from the remote sensing images using the imageIndices of the

322 soilassessment package for the R platform (Omuto, 2020).

15

## 3.2.    Spatial Prediction of soil pH

Soil pH is an important chemical property of the soil that affects its fertility status. This is because the availability of most essential plant nutrients is influenced by the levels of pH in the soil (Jones, 2012). There are two principal processes that affect the levels of soil pH in the soil (1) the production of $H^+$ ions and (2) the loss of basic cations from the soil. Eleven variables were available to be considered as possible predictors for soil pH.

In section 2.1 we explained how variable selection for the LMM included false discovery rate control, to avoid over-fitting, with alpha-investment to improve the probability of retaining covariates which are predictive as predictor variables. The alpha-investment approach is most effective if the predictors can be ordered, a priori, with the one thought most likely to be predictive ranked first and so on as shown in Table 3. It must be emphasized that this ranking is based on prior considerations about the underlying process, and not on exploration of the data. We did this ordering of exhaustive environmental covariates based on how they influence the production of $H^+$ ions and the loss of basic cations from the soil. Rainfall was proposed as the most influential factor at national scale. Soils in environments with large annual rainfall tend to have relatively low pH due to reduced based saturation resulting from loss of basic cations by leaching (McCauley, et al., 2009; Brady and Weil, 2014). For this reason more acid soils are expected in the northern parts of Zambia (Agroecological Region Three) and in the south (Agroecological Region One) where annual rainfall is much smaller (Veldkamp et al., 1984; GRZ and UNDP, 2009). Soil class was ranked second because the soil classes represent variations in soil parent material, weathering and rejuvenation of land surfaces and development of the soils. The old, highly weathered plateau soils in the northern part of the country have lost most of the basic cations. The sandy soils in the western part are easily leached with little accumulation of basic cations. On the on the hand, the Karoo group materials in the valleys are rich in basic cations resulting in high pH values. After soil class we included topographic

16

variables Slope, Elevation and Valley Depth. These should reflect processes such as the movement of water which carries with it dissolved basic cations from steep slope to flat areas, and the rejuvenation of weathered land surfaces which entails the removal of old highly-weathered material to reveal material with a larger content of weatherable minerals. We then included Landcover and the Normalized Difference Vegetation Index (NDVI). These will reflect effects of land management practices, including agricultural inputs, and decisions on land use which may depend on how local pH limits crop performance. The NDVI will also reflect local vigor of vegetation growth, which may be pH limited. Finally we included some further topographic variables which may reflect differences in the soil-forming environment (length-slope factor, channel network base level, relative slope position and aspect.

The data points were first projected from WGS 1984 to WSG UTM 35s. A total of 19 observations had duplicate coordinates, which were jittered by adding 100m to each of the coordinates for one site. An exploratory model was fitted to the data with all predictors included, using the likfit function of geoR package (Ribeiro and Diggle, 2001) with residual maximum likelihood (REML) as the likelihood method. The only output from this model which was examined were the residuals, for which summary statistics were calculated, and exploratory plots to check the plausibility of the assumption of normally distributed errors. In addition, the correlation model type (exponential or spherical) was identified for which the residual likelihood was largest, and this model was then used in all further analyses. During the sequence testing of hypothesis, first the null model, $m_0$, (with the only fixed effect a constant mean) was fitted with the likfit function and ML as the likelihood method. Then the next model, $m_1$, with the first predictor in the sequence was fitted in the same way. The likelihood ratio was then calculated:

$$L = 2(L_{m1} - L_{mo}), \tag{10}$$

where $L_{m1}$ is the likelihood for model $m_1$ and $L_{m0}$ is the likelihood for the null model. If the null model is correct, then the asymptotic distribution of $L$ is $\chi^2$ with degrees of freedom equal to the number of additional parameters in model $m_1$ by comparison with $m_0$. If $L$ provided evidence to reject the null model with P<0.05, then the additional predictor in model $m_1$ was retained. The second predictor in the list was then considered. When all predictors had been examined the *P*-values at each step were reassessed in sequence using alpha-investment as described by Lark (2017) and controlling the FDR at 0.05. Details of this approach are provided by Lark (2017), but in summary successive tests are made against a threshold *P*-value which depends on a quantity, the alpha-wealth, which is either augmented when null hypotheses are rejected or augmented when they are rejected. If the hypotheses are ordered so that the variables which, a priori, are expected to be good predictor variables are considered early, then this alpha investment method increases the probability of selecting predictive covariates while controlling FDR.

After variable selection, the likfit function of geoR package in R with REML as the likelihood method was used to fit two linear mixed models. One with the selected predictors as fixed effects (Kriging with an external drift) and the other with a constant mean as fixed effect (ordinary kriging). The E-BLUP prediction for both models was then calculated at the validation points.

The ranger function of ranger package (Wright and Ziegler, 2017) was used to fit the random forest model. Because random forest has an inbuilt variable selection that occurs within the model by randomly selecting variables to be used at splitting nodes, two models were fitted, one with all variables and the other with the two variables that were selected during the alpha-investment variable selection procedure. In this study, we use the ranger package in R to compute the permutation variable importance according to Altmann et al., (2010).

When predicting soil properties in space, the random forest algorithm can find apparently predictive relationships between the target variable and arbitrary spatial variables (such as digital images of human faces) when these are presented as candidate predictor variables alongside covariates which pedometricians might reasonably expect to be predictive of soil properties, (Wadoux et al., 2020). This shows that pattern recognition should not be equated to knowledge discovery. It may also suggest that the random forest algorithm is prone to overfitting, as a result of which its predictions at independent locations may be unreliable. To investigate this effect, we generated entirely random spatially autocorrelated candidate predictor variables, independent of our data, which we call null predictors. We used six spatially correlated but mutually independent null predictors, specifying a spherical variogram with a distance parameter of 100 km, nugget variance of 0 and correlated variance of 1 for each. We used the function `RFsimulate` from the `RandomFields` package for R (Schlather et al., 2015) to simulate valuess of these null predictors at the calibration locations. We then used the ranger package (Wright and Ziegler, 2017) to fit a random forest model with all predictors including the null variables as predictors and then computed the permutation variable importance according predictor p-values from the model result.

To examine the possibility of improving RF predictions by an additional kriging step (following Li et al., 2011 and Viscarra Rossel et al., 2014), residuals of the models at training points were derived (subtracting the predicted values from the observed values) and then a variogram was fitted to the residuals using likfit with a constant mean as the only fixed effect. The evidence for spatial dependence in the residuals was assessed by comparing the Akaike Information Criterion (AIC) for the fitted model and for a non-spatial alternative which are reported in likfit output.

**3.3.    Validation**

420  Validation of each selected model or random forest was done using the validation data set.  At

421  each validation location the predicted soil pH was computed, and the prediction error was

422  calculated as the difference between the predicted and observed soil pH (so a positive error is

423  when the predicted pH exceeds the observed value).  As exploratory summary statistics the

424  mean error, median and mean square error were computed.

425  The validation data belong to a subset of SEA from the original survey, and as such are strongly

426  clustered.  Because of this the sample average of the squared errors may not be a good estimate

427  of the mean square error, because the observations are not independent.  A model-based

428  approach was therefore taken to compute the expected squared error of prediction.  A LMM

429  was fitted to the prediction errors at the validation site (with a constant mean the only fixed

430  effect).  The expected square error (ESE) for each set of predictions was then computed as the

431  sum of the squared mean error and the variances (nugget and spatially correlated) from the

432  LMM.  This is the *a priori* mean square error, i.e. the expected square error at a random location,

433  and as such is likely to exceed the MSE computed directly from the errors of clustered data.
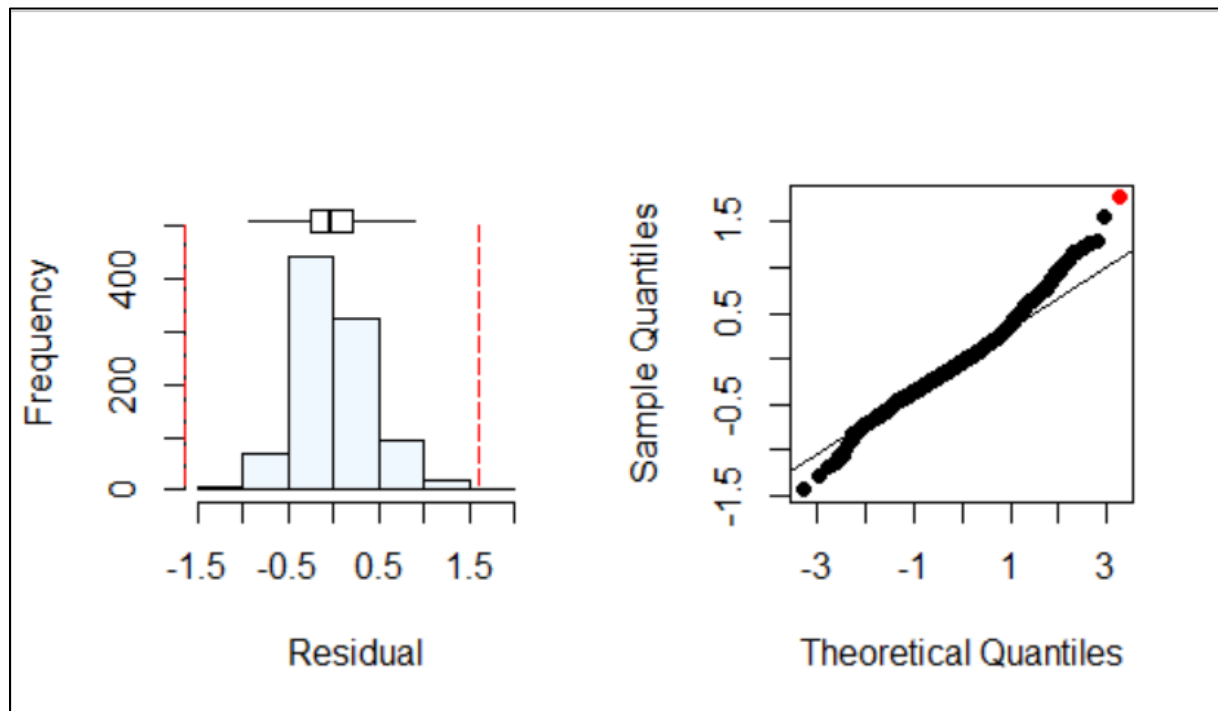
434  **4.  RESULTS**

435  **4.1 Variable selection**

436  Table 2 and Figure 3 show the distribution of the residuals from the exploratory model. The

437  histogram appears symmetrical and normal and the points on the QQ plot are close to a straight

438  line.  The residuals have octile skewness inside the range $[-0.2, 0.2]$ and skewness inside $[-1,1]$,

439  which would mean that a transformation is not normally considered necessary (Rawlins et al.,

440  2005, Webster and Oliver, 2007).

441 *Table 2: Statistical summary of residuals from the exploratory model*

| mean | Median | Variance | SD | Skewness | Octile skewness | Kurtosis |
|---|---|---|---|---|---|---|
| 2.366e-16 | -0.040 | 0.167 | 0.408 | 0.487 | 0.159 | 1.044 |

442



443

444 *Figure 3: Histogram and quantile plot of residuals from exploratory model*

445

446 *Table 3: REML estimates of parameters and AIC for the exploratory model, null model and the*

447 *hypothesis tests.*

| Test | Predictors | Partial Sill | Range | Nugget | AIC | |
|---|---|---|---|---|---|---|
| | | | | | Max.likelihood | Non-spatial |
| | Exploratory (all predictors) | 0.1367 | 21.08 | 0.218 | 1570 | 1664 |
| 0 | mean | 0.259 | 68.28 | 0.224 | 1618 | 1946 |
| 1 | Rainfall | 0.217 | 51.87 | 0.222 | 1611 | 1865 |
| 2 | Rainfall + Soilclass | 0.197 | 52.71 | 0.221 | 1622 | 1817 |
| 3 | Rainfall +Slope | 0.214 | 51.98 | 0.222 | 1610 | 1856 |
| 4 | Rainfall + Elevation | 0.180 | 39.70 | 0.220 | 1598 | 1774 |
| 5 | Rainfall + Elevation + Valley depth | 0.173 | 33.48 | 0.218 | 1598 | 1770 |
| 6 | Rainfall+ Elevation +  Landcover | 0.174 | 40.09 | 0.220 | 1603 | 1769 |
| 7 | Rainfall + Elevation + NDVI | 0.181 | 39.83 | 0.220 | 1600 | 1776 |
| 8 | Rainfall + Elevation + LS | 0.174 | 38.84 | 0.221 | 1596 | 1760 |

| 9 | Rainfall + Elevation + LS + CNBL | 0.171 | 37.40 | 0.221 | 1597 | 1783 |
|---|---|---|---|---|---|---|
| 10 | Rainfall + Elevation + LS + RSP | 0.171 | 38.13 | 0.221 | 1598 | 1751 |
| 11 | Rainfall + Elevation + LS + Aspect | 0.173 | 39.18 | 0.221 | 1598 | 1760 |

448

449 Table 4 shows the log likelihood ratio and p-values of each test at respective degree of freedom

450 df and chi-square distribution values. The likelihood ratios of tests 1,4 and 8 were greater than

451 the chi-square distribution value values. Therefore, the null hypothesis for these cases were

452 rejected, and the predictors retained during the sequential testing.  The rest of the tests had log

453 likelihood ratio less than their respective chi-square distribution values. Hence the null

454 hypothesis was accepted for these predictors and they were dropped.

455 *Table 4:  likelihood ratio and p-values of each hypothesis test at respective degree of*

456 *freedom(df) and chi-square distribution values.*

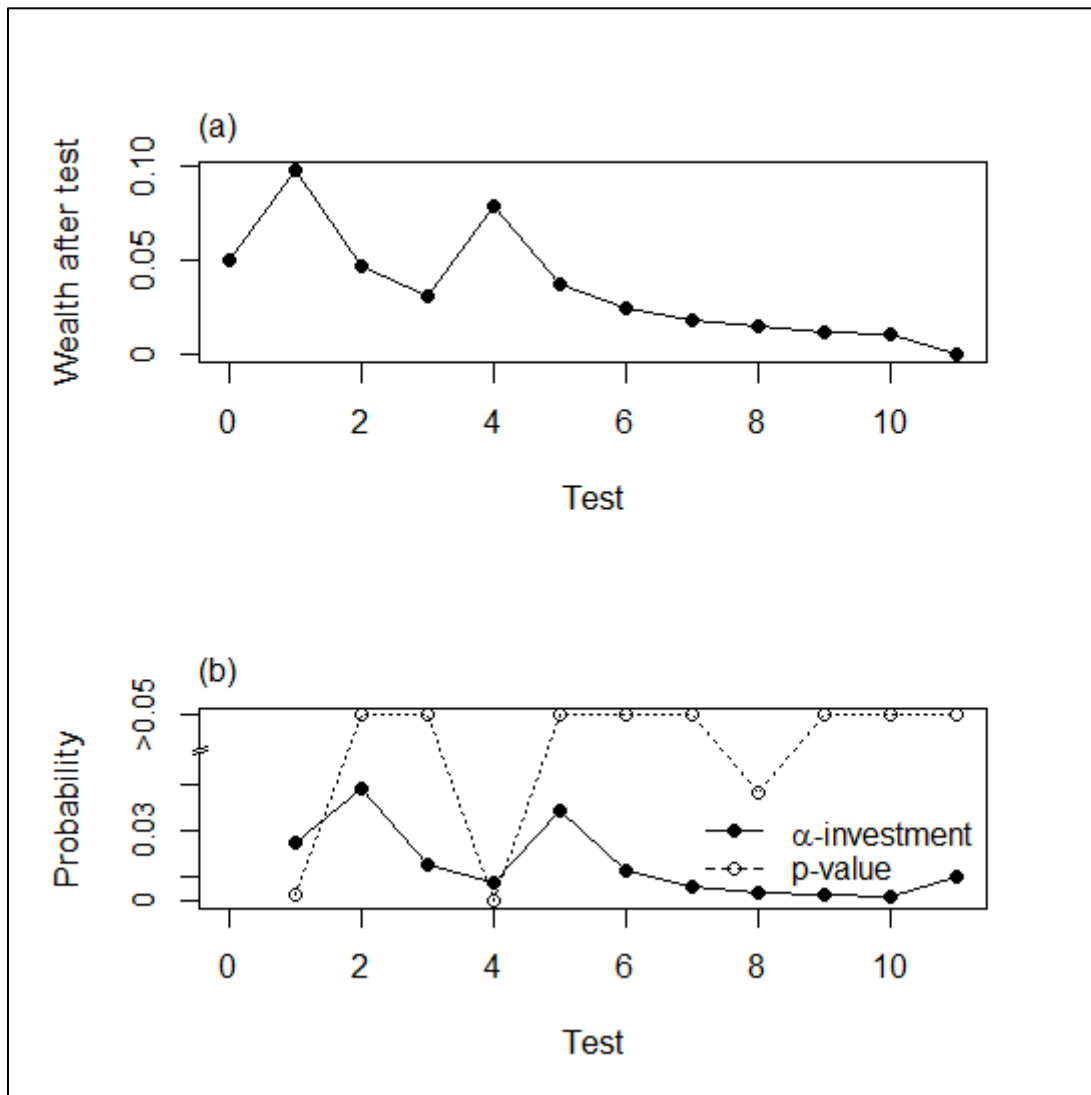| Test | | df | Chi-square | Likelihood ratio | p-value |
|---|---|---|---|---|---|
| 1 | Rainfall | 1 | 3.841 | 9.533 | 0.002 |
| 2 | Rainfall + Soilclass | 17 | 27.587 | 22.506 | 0.166 |
| 3 | Rainfall +Slope | 1 | 3.841 | 2.429 | 0.119 |
| 4 | Rainfall + Elevation | 1 | 3.841 | 14.416 | 0.000 |
| 5 | Rainfall + Elevation + Valley depth | 1 | 3.841 | 1.817 | 0.177 |
| 6 | Rainfall+ Elevation +  Landcover | 8 | 15.507 | 10.836 | 0.211 |
| 7 | Rainfall + Elevation + NDVI | 1 | 3.841 | 0.598 | 0.439 |
| 8 | Rainfall + Elevation + LS | 1 | 3.841 | 3.946 | 0.047 |
| 9 | Rainfall + Elevation + LS + CNBL | 1 | 3.841 | 1.047 | 0.306 |
| 10 | Rainfall + Elevation + LS + RSP | 1 | 3.841 | 0.403 | 0.525 |
| 11 | Rainfall + Elevation + LS + Aspect | 1 | 3.841 | 0.213 | 0.644 |

457 *Df= Degree of freedom which is the difference between the total degree of freedom for the*

458 *target model and that of the model being compared to, LS= LS-Factor, CNBL= Channel*

459 *Network Base Level, RSP= relative slope position.*

460

461 Figure 4.a shows the alpha wealth after each test and it can be observed that the quantity of the

462 wealth is increased when the null hypothesis is rejected and depleted when the null hypothesis

463    retained, and it goes to zero at the end of the sequence. Figure 4.b shows the p-values (open

464    symbols) for the successive tests of additional predictors, as in Table 4, and the threshold (solid

465    symbols) against which each successive p-value is tested to achieve FDR.  On this basis rainfall

466    and elevation were selected as predictors.

467



469    *Figure 4: a. alpha wealth after each test. b. probability of alpha investment and p-values*

470

471    **4.2  Spatial prediction of soil pH**

23

472  The estimated covariance parameters for the linear mixed models to be used for spatial

473  prediction of soil pH by the E-BLUP with elevation and rainfall as fixed effects for prediction

474  (Method A) and a constant mean as the only fixed effect (Method B, equivalent to ordinary

475  kriging) are shown in Table 5. The nugget, partial sill and range for the model with rainfall and

476  elevation as fixed effects are 0.220, 0.195 and 33.95 respectively.  These values are smaller

477  than the corresponding parameters for the model with a constant mean as the only fixed effect

478  (0.224, 0.269 and 72.83). AIC values for both models are less than those of respective non-

479  spatial AIC.  On this basis we may conclude that there is evidence for spatial dependence in the

480  random component of the LMM, and so potentially benefits in computing the E-BLUP for

481  spatial prediction at unsampled sites.

482  *Table 5: Covariance parameters for A= REML-EBLUP with elevation and rainfall as fixed*

483  *effects selected through alpha-investment (kriging with external drift), B=REML-EBLUP with*

484  *the only fixed effect a constant mean (ordinary kriging).*

| Method | Partial Sill | Range | Nugget | AIC | |
|---|---|---|---|---|---|
| | | | | Max.likelihood | Non-spatial |
| A | 0.195 | 33.95 | 0.220 | 1184 | 1310 |
| B | 0.269 | 72.83 | 0.224 | 1614 | 1945 |

485

486  Table 6 shows the number of trees, number of predictors, number of variables considered at

487  each split, target node size and out-of-bag cross validation of the two random forest methods.

488  The out-of-bag MSE and R-squared show that there is a slight reduction in performance of the

489  random forest model with rainfall and elevation

490  *Table 6: ntree = number of trees in the forest; mtry = number of variables considered at each*

491  *split*

| Method | ntree | predictors | mtry | Target node size | Out-of-Bag MSE | Out-of-Bag R-squared |
|---|---|---|---|---|---|---|

| | | | | | | |
|---|---|---|---|---|---|---|
| C | 200 | 11 | 3 | **5** | 0.31 | 0.30 |
| D | 200 | 2 | 1 | 5 | 0.32 | 0.27 |

492

493 Table 7 shows the permutation variable importance for each predictor when a random forest

494 model is fit with all predictors alone and when we include null predictors (sim1 to sim6) which

495 were generated by simulation to examine how random forest variable importance performs with

496 predictors that have no relation to the data. For the random forest model, the most important

497 variable is elevation with importance value of 0.166, followed by Channel Network Base Level

498 with value of 0.155. some variable importance values for some predictors are almost equal or

499 even less, but their p-values are much smaller. Null variables sim1 and sim6 despite have low

500 variable importance values, but very small p-values (less than 0.01). The inclusion of these

501 null variables has a substantial effect on the p-values of some predictors such as soilclass, slope,

502 landcover.

503 *Table 7: Permutation variable importance and p-values when a random forest model is fit with*

504 *all predictors alone and when we include null predictors (sim1 to sim6).*

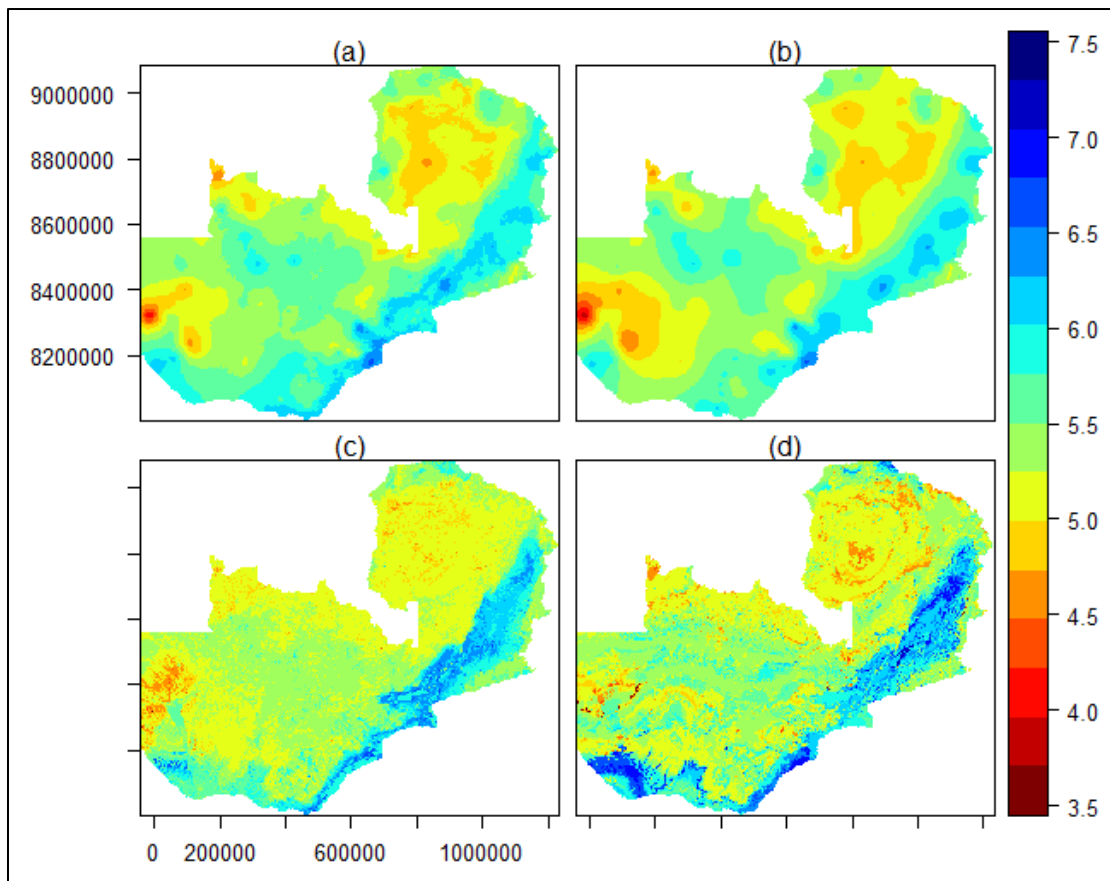| Predictor | No null predictors | | Null predictors included | |
|---|---|---|---|---|
| | Importance | p-value | Importance | p-value |
| rain | 0.0889 | 0.0099 | 0.0912 | 0.0099 |
| soilclass | 0.0231 | 0.0198 | 0.0153 | 0.0099 |
| slope | 0.0318 | 0.8218 | 0.0268 | 0.2079 |
| elevation | 0.1657 | 0.0099 | 0.1718 | 0.0099 |
| valley | 0.0761 | 0.0099 | 0.0544 | 0.0099 |
| landcover | 0.0063 | 0.4752 | 0.0074 | 0.0792 |
| NDVI | 0.0499 | 0.0099 | 0.0470 | 0.0099 |
| ls | 0.0467 | 0.3267 | 0.0286 | 0.1287 |
| cnbl | 0.1554 | 0.0099 | 0.1477 | 0.0099 |
| rsp | 0.0554 | 0.0198 | 0.0330 | 0.0495 |
| aspect | 0.0148 | 0.1188 | 0.0066 | 0.3663 |
| Sim1 | | | 0.0279 | 0.0099 |
| Sim2 | | | 0.0232 | 0.0198 |
| Sim3 | | | 0.0187 | 0.0594 |
| Sim4 | | | 0.0204 | 0.0198 |
| Sim5 | | | 0.0160 | 0.1782 |
| Sim6 | | | 0.0274 | 0.0099 |

505

506    Table 8 shows the estimated parameters of the random forest residuals for the exponential,

507    spherical and pure nugget correlation models. the non-spatial model was preferred because the

508    AIC values for the spatial component was higher than that of the non-spatial component in both

509    random forest predictions. Indeed, the fitted correlated variance for the spatial covariance

510    function was zero.  On this basis there is no scope to improve the RF predictions by a kriging

511    step.

512    *Table 8: Estimated parameters of the exponential, spherical and pure nugget correlation*

513    *functions for the residuals of the two random forest predictions.*

| Method | Parameter | Exponential | Spherical | Pure.nugget |
|---|---|---|---|---|
| | Partial Sil | 0 | 0 | 0.188 |
| | range | 0 | 0 | 50 |
| **RF (dem + rain)** | Nugget | 0.188 | 0.188 | 0 |
| | $AIC_{max.lelihood}$ | 1123 | 1123 | 1123 |
| | $AIC_{non\text{-}spatial}$ | 1119 | 1119 | 1119 |
| | Partial Sil | 0 | 0 | 0.157 |
| | range | 0 | 0 | 50 |
| **RF (all predictors)** | Nugget | 0.157 | 0.157 | 0 |
| | $AIC_{max.likelihood}$ | 951.4 | 951.4 | 951.5 |
| | $AIC_{non\text{-}spatial}$ | 947.4 | 947.4 | 947.4 |

514

515    Figure 5 shows the predicted spatial variability of soil pH. The spatial pattern is similar for all

516    the models with low pH values (less than 5.5) in the Western and Northern parts and higher

517    values (above 6) in the Southern and Eastern parts.

*Figure 5: Prediction maps of soil pH for (a) REML-EBLUP with elevation and rainfall as fixed effects selected through alpha-investment (kriging with external drift), (b)REML-EBLUP with the only fixed effect a constant mean (ordinary kriging (c) random forest with all predictors (d). random forest with elevation and rainfall as predictors selected through alpha-investment.*

**4.3 Map Validation**

Table 9 shows the summary validation statistics for (A) REML-EBLUP with elevation and rainfall as fixed effects selected through alpha-investment (kriging with external drift), (B)REML-EBLUP with the only fixed effect a constant mean (ordinary kriging), (C) random forest with all predictors (D). random forest with elevation and rainfall as predictors selected through alpha-investment. The mean and median error values were smallest for the REML-EBLUP (ordinary kriging) method while that of the REML-EBLUP (kriging with external drift) was larger than that of the two random forest methods. The MSE and RMSE for the two REML-

531  EBLUP methods were smaller than those of the two random forest methods with REML-

532  EBLUP (ordinary kriging) having the smallest values. There was spatial dependency in the

533  prediction error in all the cases with the two REML-EBLUP cases having the smaller partial

534  sill values of 0.15 compared to 0.30 and 0.22 values for Random forest (elevation and rainfall

535  predictors) and random forest (all predictors) respectively. The ESE values for all the cases

536  were larger than the MSE values because the bias (ME) for the models is greater than zero.

537  REML-EBLUP (ordinary kriging) had the smallest ESE value.

538  *Table 9: Summary Validation statistics for(A) REML-EBLUP with elevation and rainfall as*

539  *fixed effects selected through alpha-investment (kriging with external drift), (B)REML-EBLUP*

540  *with the only fixed effect a constant mean (ordinary kriging), (C) random forest with all*

541  *predictors (D). random forest with elevation and rainfall as predictors selected through alpha-*

542  *investment.*

| Variable | | A | B | C | D |
|---|---|---|---|---|---|
| Prediction error | Mean | 0.168 | 0.094 | 0.116 | 0.128 |
| | Median | 0.212 | 0.148 | 0.200 | 0.200 |
| MSE | | 0.417 | 0.388 | 0.463 | 0.551 |
| Corr.Model | | Exponential | Exponential | Exponential | Exponential |
| Partial Sil | | 0.154 | 0.145 | 0.218 | 0.299 |
| Range | | 48.380 | 44.670 | 40.200 | 36.250 |
| Nugget | | 0.240 | 0.240 | 0.237 | 0.238 |
| ESE | | 0.422 | 0.393 | 0.468 | 0.553 |

543

## 5.  DISCUSSION

545  The mapped soil pH by all approaches is shown in Figure 5. The optimum pH ($CaCl_2$) for plant

546  growth is between 5.2 – 7.5. bellow the pH of 5.2, the levels of Aluminum, Manganese and

547  Copper are toxic for plant growth, Phosphorous and Magnesium are not available to plant.

548  Above pH of 7.5, the interactions between Calcium, Magnesium and Potassium have a negative

549  impact on root absorption. Copper, Iron, Manganese, Zinc, Boron and Phosphorous are

550   deficient (Lake, 2000). The maps in Figure 5 show pH values less than 5.2 in the western and

551   northern parts of the country meaning we expect these areas to have challenges of Aluminum,

552   manganese and copper toxicity as well as Phosphorous and Magnesium deficiencies. In the

553   Southern parts of the country the pH values, according to all the maps in Figure 5, range from

554   5.2 to 7.5 which are optimal for plant growth. There are few areas in the southern part of Zambia

555   with pH above 7.5. Similar spatial variations were observed by Chapoto et al., 2016. The

556   southern parts where the pH is high is a valley area, the northern parts receives high rainfall and

557   the western part despite receiving the same amount of rainfall as the eastern part, the areas is

558   characterized by Kalahari sand. Our results show a similar spatial pattern for soil pH as that

559   presented in the SoilGrids map (www.soilgrids.org) of  Hengl et al., (2017). The main

560   difference is that our map shows low pH values in the west of the country, whereas the SoilGrids

561   map shows larger values there.  Our results are more plausible pedologically given the parent

562   material, and it has been long-established that the soils formed over the Kalahari sands of

563   western Zambia are weakly to extremely acidic (Brammer, 1976).  A more thorough assessment

564   of the SoilGrids predictions using the RALS data would be of interest.

565   Predictions by the E-BLUP from the LMM with the only fixed effect a constant mean

566   (equivalent to ordinary kriging) were better than other predictions in the sense that the mean

567   and median errors were closest to zero and the mean square error and expected square error

568   were the smallest. This is unexpected, given the evidence provided in the model-fitting stage

569   for a significant relationship between soil pH and the selected covariates.  This might be

570   expected to result in better predictions from the LMM which includes these covariates as fixed

571   effects.  However, one may note (Table 5) that the correlated random variance in the LMM with

572   rainfall and elevation as fixed effects is only about 25% smaller than the corresponding variance

573   in the LMM with a constant mean the only fixed effect.  The fact that a covariate is significantly

574   related to a soil property does not necessarily mean that it will allow improved prediction of

575  that property relative to a model without that covariate. That is because the corresponding fixed

576  effect coefficient must be estimated, and this estimation is a source of error in the prediction.

577  Furthermore, Zimmerman et al., (1999) found that ordinary kriging performed better than

578  universal kriging (UK, equivalent to the E-BLUP with some covariates) with a spatially

579  clustered data set, while UK performed better when the data were not clustered. This may be

580  because, in a strongly clustered data set, the effective degrees of freedom with which the fixed

581  effects coefficients are estimated may be relatively small.

582  The use of random forests to include the environmental covariates in spatial prediction was less

583  successful than the LMM and E-BLUP, with larger values of ESE. This could be due to over-

584  fitting. It is notable that the residuals from the fitted RF at the calibration data points showed

585  no spatial dependence, while the RF prediction errors at the validation points (Table 9) do show

586  spatial dependence. This could arise because the RF algorithm, given its flexibility and ability

587  to fit non-linear relationships, generates a model which closely fits the variations within the

588  training data set, but which is not representative of the relationship between the predictor

589  variables and target variable in the underlying population. This would lead both to marked bias

590  in models of the random variation based on the residuals, as can also occur with ordinary least

591  squares (Lark et al., 2006) and also in poor performance of the RF on a separate validation data

592  set. These data may also provide a problem for the RF methodology because of the strong

593  spatial clustering. If some data from a cluster are used in the development of trees while others

594  are in the OOB subset then the assessment of the model and the value of the predictors may be

595  over-optimistic. A predictor variable overfitted to a clustered data set might well fail to predict

596  effectively at independent validation points. This emphasizes the importance of a genuinely

597  independent validation of spatial predictions (Brus et al., 2011).

598  Spatial clustering of the observations may also be a contributing factor to the small p-values

599  attributed to the entirely random, although spatially autocorrelated, null predictors which we

30

evaluated. This gives reason for caution when interpreting RF output. It is consistent with the findings of Wadoux et al., (2020) that the RF algorithm may select as predictor variables covariates which are not related to the target properties of interest by any direct or indirect causal relations. A spatially dependent predictor variable of this nature may indeed support spatial prediction of a variable to which it has no underlying relationship, but if this is the case then one might prefer to use a properly-designed set of orthogonal polynomial basis functions for the model rather than arbitrary variables. Furthermore, with strongly spatially clustered data, it is even more likely that a uncausally-related predictor will result in poor predictions at independent validation sites.

Many digital soil mapping studies use legacy data sets, rather than new samples collected for the purpose. As legacy data sets may originate in local surveys, or from networks of experimental stations, they may show marked spatial clustering, as do the RALS data because of their two-stage cluster sampling design. We note that such clustering may cause difficulties for the RF algorithm but that it is also important to account for it when dividing data into prediction and validation subsets. There is a risk of bias in the validation of a map if validation and training data are drawn from common clusters. The estimation of validation statistics from a validation set which is strongly clustered may also result in bias, which is why we have used a model-based approach to compute these statistics in this study. The expected squared error, computed from the model, in this case is not very different from the mean squared error computed directly although in each case it exceed the mean square error as expected (Table 9). This could be because the clusters are reasonably balanced (similar numbers of observations in each), and are themselves selected independently and at random. The model-based method to quantify prediction uncertainty is, nonetheless, a more general approach for use with validation data from locations not selected by probability sampling.

31

## 6. CONCLUSION

Spatial variability of soil pH was mapped using REML-EBLUP with elevation and rainfall as fixed effects selected through alpha-investment (kriging with external drift), REML-EBLUP with the only fixed effect a constant mean (ordinary kriging), random forest with all predictors and random forest with elevation and rainfall as predictors selected through alpha-investment models. The soil pH maps from these models showed similar patterns with pH values less than 5.2 in the Western and Northern parts of the country. In the Southern parts of the country the pH values range from 5.2 to 7.5 which are optimal levels of soil pH for plant growth.

The ME, MSE and ESE (computed as the sum of the squared mean error and the nugget and spatially correlated variances from the LMM) were used to validate the performance of the models for spatial prediction of soil pH. The values of the ME, MSE and ESE from the validation statistics showed that REML-EBLUP with the only fixed effect a constant mean (ordinary kriging) performed better than the other methods.

Random forests had the largest values of MSE and ESE. This may result from over-fitting, and from the strongly spatially clustered distribution of the observations in the legacy data set which could affect the internal cross-validation in the RF algorithm.

We also noticed that the algorithm appeared susceptible to wholly random "null predictors" which we had simulated. Other studies, notably by Wadoux et al. (2020) have shown this, but we believe this to be the first example where mutually independent random spatially autocorrelated candidate predictor variables have been selected alongside pedologically plausible ones. The selection of such null predictors should give pause as it suggests that the random forest algorithm may be prone to overfitting. We suggest that this problem warrants further study as pedometricians should always aim to generate insight from their analyses, as well as predictions.

648  We note that legacy data, often used in digital soil mapping, may be strongly clustered like

649  ours. We emphasize again the importance of splitting data into prediction and validation

650  subsets at cluster level (i.e. allocating all data in any one cluster to the prediction or to the

651  validation set). In this case there was not a large difference between the ESE (model-based

652  estimate of the expected squared error) and the MSE (average of the squared errors), but this

653  would not be true in general, and the use of a model-based approach to the analysis of validation

654  errors at locations not selected independently and at random is most appropriate.

655  Finally, we note that we found no evidence for spatial correlation in the residuals from the fit of the

656  random forest to the prediction data set, although there was correlation in the prediction errors by this

657  method at the validation sites. This is an important reminder that such residuals given us little if any

658  insight into the actual behavior of the prediction error, and are good reasons to avoid using kriging

659  methods in combination with modelling methods which, unlike REML and EBLUP, do not have a built-

660  in methodology to estimate parameters of the error without bias.

**8. REFERENCES**

668
669  Altmann, A., Tolosi, L., Sander, O., Lengauer, T., 2010. Permutation importance: a corrected

670      feature importance measure, Bioinform. 26,1340-1347,

671      https://doi.org/:10.1093/bioinformatics/btq134

672 Archer, K.J., Kimes, R.V., 2008. Empirical characterization of random forest variable

673   importance measures. Comput. Stat. Data Anal., 52(4), .2249-2260,

674   https://doi.org/10.1016/j.csda.2007.08.015

675 Behrens, T., Scholten, T., 2007. A Comparison of Data-Mining Techniques in Predictive Soil

676   Mapping, in: Lagacherie, P., McBratney, A, B., Voltz, M. (Eds.), Developments in

677   Soil Sci. Elsevier B.V., 353.

678 Brady, N.C. and Weil, R.R., 2014. The nature and properties of soil. (14th New Int. Ed.).

679 Brammer, H. 1976. Soils of Zambia. Department of Agriculture, Land Use Branch. Lusaka.

680 Breiman, L., 1996. Bagging Predictors. Mach. learn., 26(2), 123-140.

681 Breiman, L., 2001. Random forests. Mach. learn., 45(1), 5-32.

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

683   maps. Eur. J. Soil Sci. 62(3), pp.394-407.

684 Burke, W.J., Frossard, E., Kabwe, S., Jayne, T.S., 2019. Understanding fertilizer adoption and

685   effectiveness on maize in Zambia. Food policy, 86, 101721,

686   https://doi.org/10.1016/j.foodpol.2019.05.004

687 Central Statistical Office, Ministry of Agriculture and Livestock, and Indaba Agricultural

688   Policy Research Institute (CSO/MAL/IAPRI)., 2015. Rural Agricultural Livelihoods

689   Survey. Lusaka, Zambia: CSO/MAL/IAPRI.

690 Central Statistical Office (CSO)., 2012. Rural Agricultural Livelihoods Survey; Instruction

691   Manual for Listing, Sample Selection, and Largest Maize Field Data Collection.

692 Chai, X., Shen, C., Yuan, X., Huang, Y., 2008. Spatial prediction of soil organic matter in the

693   presence of different external trends with REML-EBLUP. Geoderma, 148(2), 159-

694   166,  https://doi.org/10.1016/j.geoderma.2008.09.018

695     Chapoto, A., Chabala, L. M., Lungu, O. N., 2016. "A Long History of Low Productivity in

696          Zambia: Is it Time to Do Away with Blanket Recommendations?" Zambia Social Sci.

697          J. 6: (2), Article 6.

698     Cummings, M.P., Myers, D.S., Mangelson, M., 2004. Applying permutation tests to tree-

699          based statistical models: extending the R package rpart. In Tech Rep CS-TR-4581,

700          UMIACS-TR-2004-24, Center for Bioinformatics and Computational Biology,

701          Institute for Advanced Computer Studies, University of Maryland.

702     Diggle, P.J., Ribeiro Jr., P.J., 2007. Model-based Geostatistics. Springer, New York.

703     ESA. Land Cover CCI Product User Guide Version 2. Tech. Rep. 2017.

704          http://maps.elie.ucl.ac.be/CCI/viewer/download.php (Accessed on 20[th] July 2019).

705     Fick, S.E., Hijmans R.J., 2017. WorldClim 2: new 1km spatial resolution climate surfaces for

706          global land area. Intl. J. Climatol. 37 (12): 4302-4315. https://www.worldclim.org

707     Foster, D.P., Stine, R.A., 2008. α-investing: a procedure for sequential control of expected false

708          discoveries. J. R. Stat. Soc.: Series B (Stat. Methodol.), 70(2), 429-444.

709     Gashu, D., Lark, R.M., Milne, A.E., Amede, T., Bailey, E.H., Chagumaira, C., Dunham, S.J.,

710          Gameda, S., Kumssa, D.B., Mossa, A.W., Walsh, M.G., 2020. Spatial prediction of the

711          concentration of selenium (Se) in grain across part of Amhara Region, Ethiopia. Sci.

712          Total Environ. p.139231.

713     Government of the Republic of Zambia (GRZ). 1991. Exploratory Soil Map of Zambia (1:

714          1.000, 000).

715     Government of the Republic of Zambia (GRZ) and UNDP., 2009. Adaptation to the effects of

716          drought and climate change in Agro-ecological Regions I and II in Zambia. Project

717          document. Ministry of Agriculture and Cooperatives. Accepted by: Ministry of finance

718          and National Planning, and UNDP.

719      Hengl, T., 2003. Pedometric mapping: bridging the gaps between the conventional and

720          pedometric approaches. Wageningen University.

721      Hengl, T., Heuvelink, G.B.M., Kempen, B., Leenaars, J.G.B., Tamene, L., Tondoh, J.E.,

722          2015. Mapping Soil Properties of Africa at 250 m Resolution: Random Forests

723          Significantly Improve Current Predictions. PLoS One 10, 1–26,

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

725      Hengl, T., Nussbaum, M., Wright, M.N., Heuvelink, G.B., Gräler, B., 2018. Random forest as

726          a generic framework for predictive modeling of spatial and spatio-temporal variables.

727          PeerJ, *6*, p.e5518, https://doi.org/10.7717/peerj.5518

728      Jahn, R., Blume, H.P., Asio, V.B., Spaargaren, O., Schad, P., 2006. Guidelines for soil

729          description. FAO.

730      James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. An introduction to statistical learning

731          112, 3-7. New York: springer, https://doi.org/10.1007/978-1-4614-7138-7.

732      Jones Jr, J. Benton., 2012. Plant nutrition and soil fertility manual. CRC press.

733      Kempen, B., Heuvelink, B, G., Brus, J, D., Stoorvogel, J, J., 2010. Pedometric mapping of

734          soil organic matter using a soil. Eur. J. Soil Sci. 61, 333–347,

735          https://doi.org/10.1111/j.1365-2389.2010.01232.x

736      Kitanidis, P.K. 1987. Parametric estimation of covariances of regionalized variables. Water

737          Resour. Bull., 23, 557–567.

738      Kienast-Brown, Suzann Libohova, Z., Janis, B., 2010. Digital Soil Mapping, in: Soil Survey

739          Manual. USDA-NRCS, 429–436, https://doi.org/10.1007/978-90-481-8863-5

740      Lake, B., 2000. Understanding soil pH. Acid Soil Action. Leaflet, (2).

741      Lark, R.M., 2017. Controlling the marginal false discovery rate in inferences from a soil

742          dataset with $\alpha$-investment. Eur. J. Soil Sci. 68(2),221-234.

743    Lark, R. M., Ander, E. L., Broadley, M. R. 2019. Combining two national-scale datasets to

744         map soil properties, the case of available magnesium in England and Wales. Eur. J.

745         Soil Sci. 70(2), 361-377.

746    Lark, R.M., Cullis, B.R., 2004. Model-based analysis using REML for inference from

747         systematically sampled data on soil. Eur. J. Soil Sci. *55*(4),799-813,

748         https://doi.org/10.1111/j.1365-2389.2004.00637.x

749    Lark, R.M., Cullis, B.R., Welham, S.J. 2006. On spatial prediction of soil properties in the

750         presence of a spatial trend: the empirical best linear unbiased predictor (E-BLUP) with

751         REML. Eur. J. Soil Sci. 57(6): 787-799, https://doi.org/10.1111/j.1365-

752         2389.2005.00768.x

753    Lark, R.M., Webster, R., 2006. Geostatistical mapping of geomorphic variables in the

754         presence of trend. Earth Surface Processes and Landforms: The J. Br. Geomorphol.

755         Res. Group, 31(7), 862-874, https://doi.org/10.1002/esp.1296

756    Liaw, A., Wiener, M., 2002. Classification and regression by randomForest. R news, 2(3), 18-

757         22.

758    Li, J., Heap, A.D., Potter, A., Huang, Z., Daniell, J.J., 2011. Can we improve the spatial

759         predictions of seabed sediments? A case study of spatial interpolation of mud content

760         across the southwest Australian margin. Cont. Shelf Res. 31(13), 1365-1376.

761    McBratney, A.B., Odeh, I.O.A., Bishop, T.F.A., Dunbar, M.S., Shatar, T.M., 2000. An

762         overview of pedometric techniques for use in soil survey q. Geoderma 97, 293–327.

763    McCauley, A., Jones, C., Jacobsen, J., 2009. Soil pH and organic matter. Nutr. Manag.

764         Modul. 8(2), 1-12.

765    Minasny, B., McBratney, A.B., 2007. Spatial prediction of soil properties using EBLUP with

766         the Matérn covariance function. Geoderma, 140(4), 324-336,

767         https://doi.org/10.1016/j.geoderma.2007.04.028

768  Moonjun, R., Farshad, A., Shrestha, D, P., Vaiphasa, C., 2010. Artificial Neural Network and

769    Decision Tree in Predictive Soil Mapping of Hoi Num Rin Sub-Watershed, Thailand,

770    in: Boettinger, J, L., Kienast-Brown, S., Howell, D, W., Moore, A, C., Hartemink, A,

771    E. (Eds.), Digital Soil Mapping (Bridging Research, Environmental Application, and

772    Operation). Springer Science+Business Media B.V.

773  Omutu, T. C., 2020. soilassessment: Assessment Models for Agriculture Soil Conditions and

774    Crop Suitability, https://CRAN.R-project.org/package=soilassessment

775  Patterson, H. D., Thompson, R., 1971. Recovery of inter block information when block sizes

776    are unequal. Biom. 58, 545-554.

777  Rawlins, B.G., Lark, R.M., O'donnell, K.E., Tye, A.M., Lister, T.R., 2005. The assessment of

778    point and diffuse metal pollution of soils from an urban geochemical survey of

779    Sheffield, England. Soil Use Manag. 21(4), 353-362,

780    https://doi.org/10.1079/SUM2005335

781  Ribeiro, P. J., Diggle, P.J., 2001. geoR: a package for geostatistical analysis. R-News 1, 15–

782    18.

783  Schlather, M., Malinowski, A., Menck, P.J., Oesting, M., Strokorb, K., 2015. Analysis,

784    simulation and prediction of multivariate random fields with package random fields. J.

785    Stat. Softw. 63(8), 1-25, https://doi.org/10.18637/jss.v063.i08

786  Stein, M.L. 1999. Interpolation of Spatial Data: Some Theory for Kriging. Springer, New

787    York.

788  Strobl, C., Boulesteix, A.L., Zeileis, A., Hothorn, T., 2007. Bias in random forest variable

789    importance measures: Illustrations, sources and a solution. BMC bioinform. 8(1), 25,

790    https://doi.org/10.1186/1471-2105-8-25

791  Sekulić, A., Kilibarda, M., Heuvelink, G., Nikolić, M., Bajat, B., 2020. Random forest spatial

792    interpolation. Remote Sens. 12(10), 1687,  https://doi.org/10.3390/rs12101687

793   Swallow, W. H., Monahan, J. F., 1984. Monte-Carlo Comparison of ANOVA, MIVQUE,

794       REML, and ML Estimators of Variance Components. Technometr. 26, 47-57

795   United States Geological Survey (USGS)., 2019 NASA Shuttle Radar Topography Mission

796       (SRTM3) data available on the World Wide Web https://earthexplorer.usgs.gov

797       (accessed 10 July, 2019)

798   Veldkamp, W.J., Muchinda, M. and Delmotte, A.P., 1984. Agro-climatic zones in

799       Zambia. Soil Survey Bulletin (Zambia).

800   Verbeke, G., Molenberghs, G., 2000. Linear mixed models for longitudinal data. Springer-

801       Verlag, New York.

802   Viscarra Rossel, R.A., Webster, R., Kidd, D., 2014. Mapping gamma radiation and its

803       uncertainty from weathering products in a Tasmanian landscape with a proximal

804       sensor and random forest kriging. Earth Surf. Process. Landf., 39(6), 735-748.

805   Wadoux, A.M.C., Samuel-Rosa, A., Poggio, L., Mulder, V.L., 2020. A note on knowledge

806       discovery and machine learning in digital soil mapping. Eur. J. Soil Sci. 71(2), 133-

807       136.

808   Webster, R. and Oliver, M.A., 2007. Geostatistics for environmental scientists. John Wiley &

809       Sons.

810   Wright, M. N., Ziegler, A., 2017. ranger: A fast implementation of random forests for high

811       dimensional data in C++ and R. J. Stat. Softw. 77, 1-

812       17, https://doi.org/10.18637/jss.v077.i01

813   Zimmerman, D.L., Zimmerman, M.B., 1991. A comparison of spatial semivariogram

814       estimators and corresponding ordinary kriging predictors. Technometr., 33(1), 77-91,

815       https://doi.org/10.1080/00401706.1991.10484771

816    Zimmerman, D., Pavlik, C., Ruggles, A., Armstrong, P, M., 1999. An Experimental

817          Comparison of Ordinary and Universal Kriging and Inverse Distance Weighting.

818          Math. Geol. 31, 375–390, https://doi.org/10.1023/A:100758650743