**RESPONSE TO REVIEWER COMMENTS. REVIEWER #1**

We thank the reviewer for their comments, which we have used to clarify the research strategy and generally improve the manuscript.

Comment-1: Samples are taken in clusters covering only 12 fields in the study area (containing presumably a few 102 fields), with rather a poor spatial distribution. This will most probably affect the distribution of the data in the multi-dimensional space, and hence, do not cover enough the associated landscape complexity within the study area. These 2 concerns (poor spatial distribution and poor distribution in multi-dimensional space) are a major problem when the model is used for extrapolation / predicting pixels elsewhere in study area, i.e. outside the variable range covered by the calibration dataset.

Response-1: Although the samples are grouped spatially in fields for logistical reasons, they were carefully located to capture the distribution of conditions across the feature space of the agricultural fields. Specifically, the samples include the variety of parent materials in the study area, form transects across topographic positions, as well as cover regional highs and lows. If the modelling technique relied on spatial autocorrelation, such as with kriging, the spatial position of these samples would indeed have been a problem. However, because the modelling method uses spatial regression, which relies on the principle of spatial association, the important space to cover was the feature space. Therefore, the samples were taken to encompass the variable range of the agriculture fields as best could be determined prior to sampling. Notably, only agricultural fields were sampled and thus non-agricultural areas are masked out of prediction maps because they were outside the variable range covered by the calibration dataset. Text has been added to emphasize these points.

Comment-2: As a consequence, it's quite possible that the differences in SOC stock maps between the two methods are more the consequence of the fact that the 2 modelling approaches (i.e. direct versus indirect) are reacting differently on this shortcoming (inappropriate multidimensional data cover) then it is actually reflecting a real difference in model output just/purely caused by the fact that 2 different approaches were used.

Response-2: We understand this concern, but consider the variables used to calculate the SOC stock (indirect) and the SOC stock itself (direct) to be intertwined due to their definitional relationship. Because of this relationship, covering the feature space of one approach increases the probability that the feature space of the other approach is also covered.

Comment-3: Finally, it's clear how the authors calculated errors on SOC stocks by using classical error propagation techniques for individual pixels (i.e. for both the direct and the indirect method (including error predictions on components)), but it's not clear if/how spatial autocorrelation was taken into account when mapping these errors. It's important to integrate this effect of spatial autocorrelation in order to make a fair comparison between the error maps obtained by the two methods.

Response-3: Spatial autocorrelation was used in the grouping of rule condition zones, on which error estimations are applied. Spatial autocorrelation minimizes how different unsampled areas within a condition zone could be from the sampled locations. This closely resembles the approaches of Shrestha and Solomatine (2006) and Malone et al. (2006) for taking autocorrelation into account when mapping estimated errors from spatial regression models.

## RESPONSE TO REVIEWER COMMENTS. REVIEWER #2

We thank the reviewer for their comments and suggestions for improving the paper. The following is our response to the specific comments.

Comment-1: I wonder about the effect of the 10 closely-paired samples…in addition to ensuring the estimated errors capture random variation, could it induce a bias towards these sample points?

Response-1: The closely-paired samples were spread across the feature space, which should help them be representative of the random variation across the spectrum of conditions and minimize emphasis on a particular set of conditions. Some text was added to explain this part of the sampling strategy.

The Cubist software does have an option for dealing with data sets that are biased towards a grouping of target variable values. Our experimentation with this feature did not produce regression equations or estimated errors that were very different from those produced without that feature activated. This gives us confidence that the closely-paired samples were distributed in a way to not bias the models built.

Comment-2: I am a little bit worried about the large number of potential predictors in the pool compared to the number of data on the target variable (117). I appreciate that the authors use the discussion to suggest explanations of why particular predictors were selected, and thus partly validate their selection, but am still not totally convinced that the same could not be done even with junk data and this many potential predictors. I wonder if some acknowledgement of the potential of data-mining software to overfit should be included and commented on. I don't think Cubist does anything to deal with the size of predictor pool (in a multiple hypothesis testing kind of way): : :a comment on this issue could be useful.

Response-2: Text has been added to provide more detail about how we used Cubist to reduce the predictor pool. Specifically, multiple passes were used so that the final models were built from predictor pools equal to the size of the sets used in the models (3-19 predictors divided between multiple rules). This is explained in greater detail in the referenced paper Miller et al. (2015).

Comment-3: As you state in the methods section, for the propagation of error in the indirect approach, the variances and covariances should be those of the residuals from the fitted models, not of the data themselves. It seems that this is what was done, but lines 5 and 6 of page 770 made me wonder if the variances and covariances of the raw data had been used. Could you clarify this, as this could be an alternative explanation of the larger uncertainties resulting from the indirect approach?

Response-3: The text has been corrected to specify the use of the residual's covariance.

Comment-4: Define f in Equations 2 and 3, and explain exactly what |f| is.

Response-4: |f| is now defined with an explanation for why it is used. Equation 3 was also corrected.

Comment-5: In the cross validation (Table 5), I am not sure why the results for predictions of SOC stock by the indirect approach are omitted. I think the table should include these.

Response-5: Cross-validation was only conducted on the products of the Cubist models themselves. To cross-validate the SOC stock predictions from the indirect approach would be the cumulative cross-validations of the component models. We argue that cross-validations only offer a measure of

the models' robustness (how much it would change if different points were used), but does not measure the quality of the best model's performance. For these reasons, calculating a cross-validation of the SOC stock from cross-validations of the component models would be intensive without a large gain in information.

Comment-6: I think it would also be good to provide some validation of the uncertainties...I appreciate the difficulties of validating with a small, clustered dataset such as this, but I think it would be worth including some measure of the adequacy of uncertainty assessment in the cross validation. One possibility is the mean of the theta-statistic, which should be close to 1 (see e.g. Lark, RM. 2000. A comparison of some robust estimators of the variogram for use in soil survey. European Journal of Soil Science, 51, 137-157).

Response-6: We contend that the cross-validation statistics are not really a validation of the model based on all of the sample points, especially for MLR models of soil systems. For this reason, our focus in this paper is to compare the results of the two modelling approaches with each other. Nonetheless, the change in the MAE between the models used (based on all points) and the cross-validation models are compared. Additional text has been added to point out that the stability of the MAEs suggests that the estimated uncertainties are also robust.

Comment-7: It is quite interesting that although the ME for all subsoil component variables was <1, the resulting predictions of SOC stock gave a ME of 1.67. This is worth commenting on in Section 3.1.1.

Response-7: We agree this is an intriguing result and text has been added to section 3.1.1. to highlight it.

Comment-8: I think that the residuals for all variables are assumed normal...however, depending on the dataset, it may be more appropriate to model log SOC % as normally distributed. Some comment about this, and about the effect that this could have on predictions and uncertainties in the indirect approach could be useful.

Response-8: This is indeed an important point. Text, along with a table of residual skewness coefficients, has been added to describe the distribution of the residuals and their potential effect on the results.

Comment-9: Is a conservative estimate of the spatial distribution the best thing? The most conservative would be to use the mean across the entire study area, but this would not be very useful. I am not sure whether the paper is recommending that the more conservative approach should be used, or just saying that the direct approach is more conservative than the indirect approach.

Response-9: Recognizing that different situations will have different needs, we were careful to chose the term 'conservative' to avoid judging which approach may be best for a given set of goals or purpose. Specifically, sometimes representing variation can be more useful than minimizing the amount of error, and vice versa. A sentence has been added to the conclusions to emphasize this point.

Comment-10: What exactly is meant by the 'spatial association approach'?

Response-10: Spatial association is a term parallel to spatial autocorrelation. Like the specific method of kriging is often described when applying a spatial autocorrelation approach, the term

154    spatial regression is a popular method for applying the concept of spatial association. We use the
155    term spatial association frequently to emphasize the difference between our study and a similar
156    study that compared different spatial autocorrelation methods. Additional text has been added to
157    clarify spatial association for those who may not be familiar with the term.
158
159    Comment-11: Were all soil profiles deeper than 2 m?
160
161    Response-11: Thank you for catching this omission. Indeed field logistics prevented the full 2 m from
162    being sampled for some of the profiles and some assumptions needed to be made. Text explaining
163    this has been added.
164
165    Comment-12: Page 767, sentence starting on line 27: 'correlation…of R2 = 0.59'. Correlation should
166    be measured by r, not R2…reword this sentence.
167
168    Response-12: Wording has been corrected.
169
170    Comment-13: Page 768, line 6: direct R2 = 0.14, but in Table 4 is 0.19…is this correct?
171
172    Response-13: The typographical error has been fixed.
173
174    Comment-14: Figures 2, 3 and 4: I am not sure that the hillshade effect helps. I found it difficult to
175    distinguish between the effect of the hillshade and the SOC stock differences. I would suggest
176    removing this effect.
177
178    Response-14: Agreed. The hillshade effect has been removed from the respective figures.
179

180

181   **Comparison of spatial association approaches for landscape mapping of soil organic carbon stocks**

182   Bradley A. Miller*, Sylvia Koszinski, Marc Wehrhan, and Michael Sommer

183   Leibniz Centre for Agricultural Landscape Research (ZALF) e.V., Institute of Soil Landscape Research,

184   Eberswalder Straße 84, 15374 Müncheberg, Germany

185   *Corresponding author

186   Email addresses*:* miller@zalf.de (B.A. Miller), skoszinski@zalf.de (S. Koszinski), wehrhan@zalf.de (M.

187   Wehrhan), sommer@zalf.de (M. Sommer)

188

189   **Abstract**

190   The distribution of soil organic carbon (SOC) can be variable at small analysis scales, but

191   consideration of its role in regional and global issues demands the mapping of large extents. There

192   are many different strategies for mapping SOC, among which are to model the variables needed to

193   calculate the SOC stock indirectly or to model the SOC stock directly. The purpose of this research is

194   to compare direct and indirect approaches to mapping SOC stocks from rule-based, multiple linear

195   regression models applied at the landscape scale via spatial association. The final products for both

196   strategies are high-resolution maps of SOC stocks (kg m$^{-2}$), covering an area of 122 km$^2$, with

197   accompanying maps of estimated error. For the direct modelling approach, the estimated error map

198   was based on the internal error estimations from the model rules. For the indirect approach, the

199   estimated error map was produced by spatially combining the error estimates of component models

200   via standard error propagation equations. We compared these two strategies for mapping SOC

201   stocks on the basis of the qualities of the resulting maps as well as the magnitude and distribution of

202   the estimated error. The direct approach produced a map with less spatial variation than the map

203     produced by the indirect approach. The increased spatial variation represented by the indirect

204     approach improved $R^2$ values for the topsoil and subsoil stocks. Although the indirect approach had a

205     lower mean estimated error for the topsoil stock, the mean estimated error for the total SOC stock

206     (topsoil + subsoil) was lower for the direct approach. For these reasons, we recommend the direct

207     approach to modelling SOC stocks be considered a more conservative estimate of the SOC stocks'

208     spatial distribution.

209     *Keywords:* digital soil mapping, organic carbon, spatial association, estimated error, uncertainty

210

211     Highlights

212         1. Spatial association methods for mapping SOC stock directly and indirectly were compared.

213         2. Data mining produced models that could be interpreted by expert knowledge.

214         3. The indirect approach map had greater spatial variation and higher $R^2$ values.

215         4. The direct approach map had less spatial variation and a lower total estimated error.

## 1. Introduction

The storage of carbon in soil is a critical point of information for several environmental issues. Globally, soil carbon, which is about 60% organic carbon, accounts for 3.3 times more carbon than that found in the atmosphere (Lal, 2004). The high amount of carbon stored in the soil, makes soil carbon an important factor for understanding the carbon cycle and dynamics influencing global climate change (Grace, 2004; Johnston et al., 2004; Powlson et al., 2011). In addition, higher concentrations of soil organic carbon (SOC) are associated with better water storage capacity, regulation of nutrients, and stabilization of soil aggregates resulting in improved soil structure and resistance to erosion (Neemann, 1991; Angers and Carter, 1996; Rawls et al., 2003; Snyder and Vazquez, 2005; Johnston et al., 2009; Kay, 1998; Wilhelm et al., 2004). Each of these factors has important roles in issues of water management and crop productivity.

Although SOC management has far reaching implications, the distribution of SOC is highly variable and dynamic at the field-scale (Cambardella et al., 1994; McBratney and Pringle, 1999; Walter et al., 2003; Kravchenko et al., 2006b; Simbahan et al., 2006). Differing conditions, such as hydrology or management practices, greatly impact the SOC content (Kravchenko et al., 2006a). The combination of global implications and high spatial variability make high-resolution maps of SOC for large extents desirable for both policy decisions and land-owner response. This situation creates the need to accurately and efficiently assess the spatial distribution of SOC stocks at a high-resolution. High-resolution mapping captures information essential for assessing field-specific conditions, which can later be aggregated as need to provide summary information.

Many studies have tested a variety of strategies for predicting the spatial distribution of SOC (Minasny et al., 2013 and references therein). The various studies on SOC mapping have analyzed different soil depths, which has large implications for the consideration of the complete SOC stock (Richter and Markewitz, 1995; Batjes, 1996, Jobbágy and Jackson, 2000; Sombroek et al, 2000; Schwartz and Namri, 2002; Meersmans et al., 2009). For example, some have focused on spatially modelling the topsoil to depths of 20-30 cm (e.g. Ungaro et al., 2010; Zhang et al., 2010; Martin et

242 al., 2011). Other variations of strategies for digital SOC mapping differ in which variables are

243 modelled in order to predict SOC. For instance, some studies have modelled the SOC stock (e.g. kg

244 $m^{-2}$, T $ha^{-1}$, kg $m^{-3}$) directly (Simbahan et al., 2006; Lufafa et al., 2008; Nyssen et al., 2008; Mishra et

245 al., 2010; Phachomphon et al., 2010; Kempen et al., 2011), while others have separately modelled

246 the variables needed to calculate the SOC stock and then combined them (Grimm et al., 2008; Khalil

247 et al., 2013; Lacoste et al., 2014). The usual component variables are total bulk density (BD), particles

248 > 2 mm (SK), SOC concentration ($SOC_\%$), and stock thickness (H), which are then combined by:

249 $$SOC_{stock} = \left(\frac{SOC_\%}{100}\right) * (BD * 1000) * \left(\frac{100-SK}{100}\right) * H \qquad (1)$$

250 where, $SOC_{stock}$ is in kg $m^{-2}$, $SOC_\%$ is in percent, BD in g $cm^{-3}$, SK in percent, and H in m.

251     Irrespective of the approach used, an important output of digital soil mapping is a measure of

252 uncertainty. Orton et al. (2014) compared uncertainties resulting from directly modelling the SOC

253 stock (direct = calculate-then-model) with modelling component variables for calculating the SOC

254 stock (indirect = model-then-calculate), based on geostatistical approaches that rely on spatial

255 autocorrelation. In the present study, we made a similar assessment for rule-based, multiple linear

256 regression (MLR) models, which rely on spatial association.

257     With the spatial association (i.e. spatial regression) approach to soil mapping, the empirical

258 model error can be transferred along with the model itself (Lemercier et al., 2012). For digital soil

259 mapping, Malone et al. (2011) adapted the Shrestha and Solomatine (2006) approach for empirically

260 summarizing model error and extending that information to prediction areas. In those previous

261 studies, areas expected to have similar errors were grouped by cluster analysis. Because similar sites

262 are already grouped together in rule-based, MLR models, the estimated errors can be applied to the

263 areas meeting the same rule conditions and thus mapped. The ability to map predictions of soil

264 properties and the confidence in those predictions via spatial association is important for landscape

265 to national extents because of the common limitation of sampling density (Martin et al., 2014).

266    The purpose of this study was to compare the maps of SOC stocks produced from direct and

267    indirect modelling approaches, using rule-based MLR. The resulting maps were compared in terms of

268    their predicted spatial patterns, coefficient of determination ($R^2$), as well as the magnitude and

269    spatial distribution of the estimated errors. The predictors selected for the models via the data

270    mining procedure were evaluated in the context of known landscape processes. In addition, the

271    separate assessment of topsoil and subsoil stocks tested the models' ability to predict SOC storage at

272    depths to two meters.

273    **2.   Methods**

274    *2.1. Study Area and Sampling*

275    A dominantly agricultural area located near Wulfen, Saxony-Anhalt, Germany, which has been

276    examined by several previous studies (Selige et al., 2006; Brenning et al., 2008; Kühn et al., 2009;

277    Migdall et al., 2009), was selected for this research. The mapping area extends from 11.86°N,

278    51.74°E to 11.96°N, 51.90°E (Figure 1), covering a total area of 122 km². The landscape includes

279    hummocky till plain, outwash plain, loess, and a broad floodplain (Königlich Preußische Geologische

280    Landesanstalt, 1913a, b). The study area is dominated by Calcaric Cambisols and Luvic Phaeozems,

281    while the depressional area in the floodplain is primarily Dystric Gleysols (European Commission,

282    2014). Between 2005 and 2006, 117 locations were sampled from a variety of landscape positions in

283    12 different agricultural fields, covering the known feature space for agricultural land in this area.

284    Because all models were calibrated and validated on these samples, evaluation of the resulting maps

285    focused on areas with similar land-use (i.e. water bodies and urban areas excluded). Ten of the

286    sample points, also spread across the feature space, were of repeated locations (within 2 m of

287    original), which helped to insure that random error was reflected in the assessment of estimated

288    error.

289    Soil horizons identified in the field were sampled at each sampling location. To avoid biases from

290    horizon classifications and to focus on the two major process zones for SOC, the soil profile of two

291    meters was divided into topsoil and subsoil stocks. The division was defined by the largest decrease

292    in SOC$_\%$, as determined by lab analysis, between field identified horizons. Not all profiles were able

293    to be sampled to the full depth of two meters. In those cases, the properties of the sampled subsoil

294    were assumed to be representative of the remaining depth. Data for the horizons within each stock

295    were combined using a thickness-weighted mean, as appropriate. Descriptive statistics for these

296    observation points are provided in Table 1.

297    *2.2. Modelling*

298       Models for each of the target variables were generated using the Cubist 2.08 software (Quinlan

299    1992, 1993, 1994). Previous studies have demonstrated the utility of this tool for digital soil mapping

300    (Bui et al. 2006; Minasny and McBratney, 2008; Adhikari et al., 2013; Lacoste et al., 2014). Cubist

301    uses a data mining algorithm to build two-tiered models. The top level consists of a series of

302    conditional rules that can utilize both continuous and categorical predictors. For each rule, a MLR

303    equation is produced for predicting the target variable. Cubist's process for selecting predictors and

304    building the models is described in Quinlan (1993) and Holmes et al. (1999) and will not be repeated

305    here. One advantage of this approach is the interpretability of the produced model, which allows the

306    modeler to assess relationships between the model and physical processes (Bui et al., 2006).

307       The results of the data mining process are dependent upon the predictors made available to the

308    data mining software. For this reason, we used the large predictor pool method described by Miller

309    et al. (2015) to identify the optimal models for each of the respective target variables. That method

310    includes a multiple pass test, which reapplies the Cubist algorithms to the limited pool selected by

311    the previous run. This helps to insure that the selected predictors have been optimally reduced by

312    the Cubist software, decreasing the concern of overfitting. The predictor pool for this study included

313    410 base maps covering the full extent of the study area (Table 2). These base maps consisted of a

314    legacy geologic map, a variety of remote sensing/spectral products, and digital terrain analysis

315    (DTA). The spectral products ranged from four bands of Ikonos data to a variety of Landsat data

316    collected at different times in 2006. DTA was conducted on a 2 m resolution, digital elevation model

317    (DEM), created from LiDAR data that was also collected in 2006. The DTA base maps included land-

318    surface derivatives based on a wide range of analysis scales (a-scales) and a suite of hydrologic

319    indicators. Land-surface derivatives were calculated in GRASS 6.4.3 (Geographic Resources Analysis

320    Support System, grass.osgeo.org) and ArcGIS 10.1 (www.esri.com/software/arcgis). Hydrologic

321    indicators were calculated using SAGA 2.1.0 (System for Automated Geoscientific Analysis,

322    http://www.saga-gis.org/en/index.html).

323        The predictors selected by the Cubist software were then used as base maps to generate maps

324    of $SOC_{stock}$. Using the raster calculator in ArcGIS 10.1, the base maps were combined according to the

325    MLR equations produced by Cubist. When base maps of different resolutions were combined, the

326    finest resolution was maintained. The respective MLR equations were only applied in the areas that

327    met the conditions of the Cubist model's first tier. The first experimental approach used this method

328    to directly map $SOC_{stock}$ from the $SOC_{stock}$ calculated at each sample point. The second experimental

329    approach used this method to map each of the component variables. These modelled variables were

330    then used as base maps to create a $SOC_{stock}$ map. The raster calculator was then again used to

331    combine the component variables, but this time according to equation 1. For both experimental

332    approaches, the topsoil and subsoil were mapped separately. After the respective $SOC_{stock}$ maps

333    were produced, they were added together to create total $SOC_{stock}$ maps.

334        Within the extent of the study area, there were a few areas with conditions outside the range

335    observed in the point samples. In these limited cases, extreme predictor values produced model

336    predictions of target variables either far below or above the ranges observed for the respective

337    target variables. To address this issue, spatial predictions were limited to be within 10% of the

338    observed target variable minimum and maximums.

339    *2.3. Propagation of Error*

340    For each of the model rules, estimated error was calculated based on the internal fit of the MLR

341    to the data classified within that rule. This estimation provided a measure for the respective

342    uncertainty under each rule. The conditions for the respective rules were used to spatially classify

343    the base maps, thus allowing the estimated errors to be mapped. Measurement error, positional

344    error, and limitations of the model to predict the target variable were all empirically encapsulated by

345    the estimated error.

346    When the target variable was the end product, the uncertainty was simply represented by the

347    estimated error. However, when multiple variables were modelled and subsequently used to

348    calculate the final product, the estimated errors of the component variables propagated through the

349    combination of those variables in the function. In order to map estimated error for the indirect

350    approach of modelling $SOC_{stock}$, estimated error maps were produced for each of the component

351    variables. These error estimation maps were then combined using standard equations for

352    propagation of error (Mardia et al., 1979; Taylor, 1997; Weisstein, 2014). Although potentially biased

353    by the approximation to a first-order Taylor series expansion, simplified equations for error

354    propagation are more practical and are regularly used in engineering and physical science

355    applications (Goodman, 1960; Ku, 1966). Because covariance between variables has the potential to

356    impact the estimation of $SOC_{stock}$ (Panda et al., 2008; Goidts et al., 2009), we did not assume the

357    variables were independent. The observed residual covariance was thus used to modify the

358    estimated error within the standard equations for propagation of error by multiplication,

359    $$\sigma_f \approx |f|\sqrt{\left(\frac{\sigma_A}{A}\right)^2 + \left(\frac{\sigma_B}{B}\right)^2 + 2\frac{cov_{AB}}{AB}} \qquad (2)$$

360    and by addition,

361    $$\sigma_f \approx \sqrt{\sigma_A{}^2 + \sigma_B{}^2 + 2cov_{AB}} \qquad (3)$$

362    where *f* is the result of the original function (to convert from relative to estimated error), A and B are

363    the real variables, with estimated errors $\sigma_A$ and $\sigma_B$, and their residuals' covariance $cov_{AB}$. In order to

364 calculate a predicted relative error (e.g. $\frac{\sigma_A}{A}$) at unsampled locations, the predicted variable was

365 assumed to accurately represent the variable's magnitude.

366      Locations with small ratios between estimated error and predicted values together with large,

367 negative covariances had the potential to produce a calculation taking the square root of a negative.

368 This issue was addressed by not considering the covariance in those limited circumstances. While

369 this solution may have led to an overestimation of error, it provided a means to mathematically

370 calculate estimated error without declaring it to be zero.

371 **3. Results**

372 *3.1. Models*

373 *3.1.1. Model building and fitting performance*

374      Explicit models were obtained for each of the component variables needed to calculate $SOC_{stock}$

375 indirectly and for predicting $SOC_{stock}$ directly. Models for predicting component variables used a

376 higher quantity of predictors for each of the respective models than the direct modelling approach

377 (Table 3). With the exception of $SOC_{\%}$, the models for component variables included a combination

378 of DTA and spectral variables. The $SOC_{\%}$ models relied solely on DTA predictors for both stocks, but

379 with additional spatial partitioning by geologic map units for the topsoil model. The models for

380 directly predicting the $SOC_{stock}$ used only three DTA predictors for the topsoil and only four Landsat

381 predictors for the subsoil.

382      Fitting performances for the component variable models were better than the fitting

383 performances for the direct modelling of $SOC_{stock}$ (Table 4). For the component variables, $R^2$ values of

384 subsoil models were only slightly less than the topsoil models. $SOC_{\%}$ was the exception by having the

385 lowest fitting performance for the subsoil stock ($R^2$ = 0.55), while the model for the $SOC_{\%}$ topsoil was

386 able to fit observations with an $R^2$ of 0.86. However, it was the aim of this research to examine if the

387 performance of the models was maintained through the calculation of $SOC_{stock}$.

388    Comparison of the $SOC_{stock}$ predictions by the indirect approach to observed values showed

389    better performance for the topsoil stock ($R^2$ = 0.73) than for the subsoil stock ($R^2$ = 0.34). Fitting

390    performance for directly modelling $SOC_{stock}$ showed the same pattern, but was lower than the

391    indirect approach for both stocks. Analysis of the direct approach's ability to fit observed values

392    yielded an $R^2$ of 0.58 for the topsoil and 0.~~14~~ 19 for the subsoil.


393    In general, calculated model efficiencies (ME) showed that the respective models reduced the

394    mean absolute error (MAE) to about half the MAE that would result from simply using the mean of

395    all points as the prediction. The $SOC_{\%}$ model for the topsoil improved upon the mean model more

396    than the other MLR models with a ME of 0.34. However, an intriguing result is the lack of model

397    efficiency for the indirect modelling of the subsoil's $SOC_{stock}$. Despite the component models all

398    having MEs well below one, the indirect approach did not improve upon the mean model for

399    predicting the subsoil $SOC_{stock}$. Although the ME of the direct model for subsoil $SOC_{stock}$ was also not

400    as good as the other models, it was still an improvement over the mean model.


401    *3.1.2. Model Robustness*

402    It is common for digital soil mapping models to be evaluated by cross-validation procedures.

403    However, in the context of this study, the meaning of such an analysis has less utility. Higher sample

404    density increases the robustness of the model (Minasny et al, 2013); thus the popularity of cross-

405    validation procedures over independent validation procedures in order to maintain more points in

406    the calibration set. However, the model generated for each cross-validation run is different because

407    of differences in calibration sets. The performance of each run is dependent on the randomly

408    selected calibration points' ability to represent the variation in the remaining validation points. For a

409    simple data trend, a single outlier would have minimal effect because only the runs in which it is

410    included in the validation set – and not used in calibrating the model – would have lower

411    performance values. However, in a complex landscape where similar soil properties can result from

412    different combinations of factors, the concept of an outlier has many more dimensions (Johnson et

413    al., 1990; Phillips, 1998). A point with a similar value can be an outlier by being a product of a

414    different set of factors. In other words, the problem of induction continues to apply in predictive soil

415    mapping. Further, in the context of error propagation, the error estimation from the actual model

416    used seems more appropriate than the mean of error estimations from a series of less robust

417    models.

418        Nonetheless, the models in this study were cross-validated using the k-fold method with 10

419    iterations. The $R^2$ was naturally reduced in the cross-validation analysis, but the ~~mean absolute error~~

420    ~~(~~MAE~~)~~ was not as severely affected (Table 5). The $R^2$ values for the respective models all decreased

421    greatly in the cross-validation, except for the topsoil $SOC_\%$ and the subsoil $SOC_{stock}$ models. The

422    subsoil $SOC_{stock}$ model already had a low $R^2$ value for the internal fit. In contrast, the MAEs for the

423    cross-validation of the models were not increased enough to present a practical problem. <u>The</u>

424    <u>relative stability of the MAEs also suggests that the estimated uncertainties are also robust.</u> For

425    example, the MAE for both stocks of BD only increased 0.03 g cm$^{-3}$. Also, the MAE for $SOC_\%$ only

426    increased 0.13% and 0.03% for the topsoil and subsoil, respectively. Similarly, the MAE for the direct

427    $SOC_{stock}$ model increased 0.67 kg m$^{-2}$ and 0.05 kg m$^{-2}$ for the topsoil and subsoil, respectively. The

428    MAE for the models of stock H and SK did increase more in cross-validation. However, they had a

429    minor impact on the indirect modelling of $SOC_{stock}$. The increase of 5.9 cm for the topsoil H MAE was

430    only a shift of the depth estimated by topsoil or subsoil models. The larger MAE for SK was more of

431    an issue for the subsoil. However, the majority of the samples had SK below 5%, leaving most of the

432    error due to the difficulty in predicting the limited areas of high SK. While it was possible that a

433    different sampling design could have improved the $R^2$ values for cross-validation, they are not always

434    practical for landscape-scale mapping.

435    *3.1.3. Comparison with previous studies*

436        It is difficult to compare results between SOC mapping studies due to differences in study areas

437    and strategies for defining $SOC_{stock}$ (i.e. map extent and resolution, sampling density, and

438    consideration of depth). Further, the differences between and variability within methods for

439    estimating component variables for calculating $SOC_{stock}$ can have a large impact on results, especially

440    bulk density (Liebens and VanMolle, 2003; Schrumpf et al., 2011) and $SOC_\%$ (Lowther et al., 1990;

441    Soon and Abboud, 1991; Sutherland, 1998; Bowman et al., 2002). Also, because model performance

442    is dependent upon the provided predictors, results of different studies can vary based on the

443    predictors available to and derived by the modeller (Miller et al., 2015). However, because the area

444    in this study has been used for several previous studies, some comparisons between methods can be

445    made.

446         Kühn et al. (2009) examined many of the same samples used in this study and found a

447    correlation coefficient of determination between soil electrical conductivity and soil organic matter

448    to a 1 m depth (kg m$^{-2}$) of $R^2 = 0.59$. Although a slightly different calculation, that correlation

449    coefficient of determination is similar to this study's direct model of topsoil $SOC_{stock}$ ($R^2 = 0.58$),

450    which used three DTA predictors. However, for the topsoil, the indirect approach in this study

451    produced a $SOC_{stock}$ model with less estimated error and an $R^2$ of 0.73. The Kühn et al. (2009) study

452    usually included depths that this study defined as subsoil, where the models in this study did not

453    perform as well (direct $R^2 = 0.1419$, indirect $R^2 = 0.34$).

454         For the same area as this study, Selige et al. (2006) compared MLR and partial least-square

455    regression for predicting $SOC_\%$ from hyperspectral data with a 6 m spatial resolution. Although the

456    study by Selige et al. (2006) utilized a higher spectral resolution, the MLR models produced by both

457    that study and the present study had $R^2$ of 0.86 for the topsoil $SOC_\%$. In the present study, Cubist was

458    able to compensate for the limited spectral information by utilizing several DTA predictors that were

459    available at a high spatial resolution.

460    *3.2. $SOC_{stock}$ maps*

461      Application of the obtained models and aggregation of the component variable maps by

462      equation 1 produced maps of predicted $SOC_{stock}$ for the topsoil and subsoil (Figures 2 and 3). The

463      respective topsoil and subsoil maps were added together to produce a total $SOC_{stock}$ map to a depth

464      of 2 m (Figure 4). Although some field boundaries were observed, the dominant pattern appeared to

465      be associated with terrain features. This interpretation was supported by the number of DTA

466      predictors selected by Cubist for many of the models. However, it would not have been safe to

467      assume this pattern from the list of selected predictors alone. Certain predictors (i.e. spectral data

468      reflecting land use patterns) could have dominated calculations without being the most frequently

469      selected category of predictors.

470      The map derived from the direct approach for modelling the topsoil $SOC_{stock}$ emphasizes

471      drainageways. Whereas the map derived by the same approach for the subsoil $SOC_{stock}$ reflects more

472      patterns of land use, especially in the uplands in the southern part of the study area. The topsoil

473      $SOC_{stock}$ map based on the indirect approach has similar overall patterns to the direct approach's

474      map. However, both the topsoil and subsoil maps produced by the indirect approach display greater

475      spatial variation.

476      Patterns in the topsoil $SOC_{stock}$ map, based on the indirect approach, mostly coincide with terrain

477      features, but do contain some transitions that align with field boundaries. The corresponding map

478      for the subsoil reflects patterns of microtopography and slope gradient. Larger values for the subsoil

479      $SOC_{stock}$ are predicted by the indirect approach for local lows in elevation (smaller a-scales).

480      Predictions of larger subsoil $SOC_{stock}$ on steeper slopes result from the modelling of thinner topsoil

481      stocks in these areas and the consistent calculation of a 2 m profile. Consequently, the subsoil is

482      calculated to be thicker in these areas, substantially increasing the subsoil $SOC_{stock}$ prediction

483      compared to other areas of the subsoil.

484      Maps derived by both approaches for the total $SOC_{stock}$ primarily reflected patterns from the

485      topsoil maps because of the higher concentration of SOC that defined the topsoil stock.

486    Nonetheless, modelled storage for the subsoil stock contributed about one-third of the prediction of

487    total $SOC_{stock}$ and recognized additional complexity in the SOC landscape. Despite the greater

488    variation in the indirect approach's prediction of $SOC_{stock}$, the difference between estimates of total

489    $SOC_{stock}$ by the two approaches were within 5 kg m$^{-2}$ for the majority of the map area (Figure 5). Also,

490    the summed $SOC_{stock}$ for the study area was only 6% more for the indirect (1.9 Mt) versus the direct

491    (1.8 Mt) approach. The mean $SOC_{stock}$ estimate for the study area by the direct approach was 14.7 kg

492    m$^{-2}$, whereas the indirect approach estimated 15.7 kg m$^{-2}$.

493    These aggregated landscape estimates agreed with those made by the Harmonized World Soil

494    Database (HWSD; FAO/IIASA/ISRIC/ISSCAS/JRC, 2012) for this area. The HWSD estimated several soil

495    properties from taxonomic pedotransfer functions for static topsoil (0-30 cm) and subsoil (30-100

496    cm) depth zones. Within the area of the present study, the HWSD has a cell resolution of

497    approximately 765 m. Calculating $SOC_{stock}$ from that data yielded a mean of 8.8 kg m$^{-2}$. Assuming the

498    characteristics of the subsoil to 100 cm extended to 200 cm, the mean $SOC_{stock}$ would be 15.3 kg m$^{-2}$.

499    *3.3 Error estimations*

500    The mapping of estimated errors based on the conditions of rules generated by Cubist resulted

501    in a spatial representation of uncertainty (Figure 6). In order to calculate the final estimated errors

502    for the indirect approach, estimated errors for models of component variables were combined

503    spatially by equations ~~3~~ 2 and ~~4~~3. Due to the known covariance of component variables, the

504    observed covariance of the residuals was included in the calculation of error propagation through

505    the calculation of the total $SOC_{stock}$. Inclusion of covariance reduced relative error estimates in the

506    topsoil because increases in residuals for BD coincided with decreases in the residuals for percent

507    fine-earth, increases in fine-earth BD residuals coincided with decreases in $SOC_{\%}$ residuals, and

508    increases in SOC content (kg m$^{-3}$) residuals coincided with decreases in stock thickness residuals. The

509    influence of covariance was ~~not~~mostly the same in the subsoil calculations. The exception was a

510    positive covariance between the residuals for modelling BD and the percent fine-earth. ~~With the~~

18

511   exception of the covariance between fine-earth BD and SOC~~%~~, which was very small, subsoil

512   ~~covariances were positive. However, overall~~Nonetheless, the covariances were relatively small with

513   respect to the estimated errors and therefore had a minimal impact on the final calculation of

514   estimated error.

515   The application of error estimates based on the full range of predicted values in a rule zone to

516   small values in that zone yielded extremely high relative error values. Although the areal extent for

517   this type of situation was very limited, the issue needed to be addressed in order to maintain the

518   readability of the attribute scale. Therefore relative error was capped at one for the original relative

519   error grids, but not thereafter for the calculation of error propagation.

520   Despite not having as strong of a fitting performance as the indirect approach, the direct

521   approach had lower estimated errors for greater extents of the study area. The mean estimated

522   error for the total $SOC_{stock}$ map derived by the direct approach was 2.81 kg m$^{-2}$, compared to 8.17 kg

523   m$^{-2}$ for the indirect approach. This behavior in the models may be explained by the negative

524   covariance between the residuals for many of the variables influencing the $SOC_{stock}$. The observed

525   covariances did reduce the calculation of error through propagation. However, they did not reduce

526   the estimated error for the indirect approach to as low as the estimated error based on the direct

527   modelling approach. It is also useful to note that the residuals for modelling SK and $SOC_{\%}$ ~~were~~had a

528   negative and positive skew~~ed~~, respectively, for both stocks (Table 6). ~~However, for the~~Of the

529   residuals ~~of~~for the final prediction of $SOC_{stock}$, regardless of approach or stock, only the indirect

530   model for the subsoil had strongly skewed residuals. This suggests that ~~the~~ error for the indirect

531   model of the subsoil $SOC_{stock}$ may have been overestimated.

532   The spatial distribution of model rules was an important factor in the resulting maps' estimated

533   error. The models for the direct approach used fewer rules than the component variable models,

534   resulting in less spatial variation of the estimated error. However, variation in predicted values did

535   introduce additional spatial variation to the mapping of relative error. Nonetheless, the map of

536    relative error from the indirect approach was more complex than that resulting from the direct

537    approach. In addition to using more rules for each model, the combined relative estimated error for

538    the indirect approach was further tessellated by the unique intersections of the different spatial

539    distributions of the rules for each component variable model.

540    **4. Discussion**

541    *4.1. Predictor selection*

542    *4.1.1. Review of relationships between predictors and environmental conditions*

543    Spectral predictors from satellites such as Ikonos and Landsat have most commonly been used

544    to detect characteristics of land use, vegetation, and soil water content (Bannari et al., 1995; Xie et

545    al., 2008). However, they have also been used to detect mineralogy on sparsely vegetated areas

546    (Mulder et al., 2011). Although Ikonos has a finer spatial resolution, it is limited to three bands (band

547    1 = blue, band 2 = green, and band 3 = red) in the visible spectrum, plus a near infrared band (band 4

548    = NIR). Landsat provides additional bands in the shortwave infrared (band 5 = SWIR-1; band 7 =

549    SWIR-2) and thermal infrared (band 6 = TIR). The relative reflectance of a single band can be used to

550    distinguish landscape conditions. For example, the green band can be used to distinguish different

551    vegetation from bare soil. However, combinations of bands - particularly including the red and NIR

552    bands - have been even more useful for distinguishing the spectral signature of different land uses

553    (Richards, 2006) and the condition of the vegetation (Ashley and Rea, 1975; Myneni et al., 1995;

554    Rasmussen, 1998; Daughtry, 2001; Hatfield et al., 2008). Additional use of TIR emission would

555    resemble methods such as the Surface Temperature/Vegetation Index for estimating soil moisture

556    (Bartholic et al., 1972; Heilman et al., 1976; Carlson et al., 1994; Li et al., 2009; Petropoulus et al.,

557    2009). Similarly, use of SWIR wavelengths in concert with red and infrared ~~red~~ bands would be a way

558    of compensating for the changing effect of soil reflection in dry to wet conditions (Huete, 1988;

559    Lobell and Asner, 2002). Relationships between bands in the visible to SWIR range have also been

560    used to predict SOC$_\%$ and its biochemical composition (Bartholomeus et al., 2008; Gomez et al.,

561    2008; Stevens et al., 2010).

562        Spectral predictors have been used for both classification of discrete phenomenon and

563    quantification of continuous phenomenon on the landscape. Because of the rule-based MLR

564    structure of the Cubist models, spectral predictors used for conditional rules were more likely to be

565    distinguishing discrete features (e.g. vegetation/land use type) than when used within an MLR

566    equation. Continuous features (e.g. vegetation health) were more likely to be represented in MLR

567    equations.

568        DTA predictors in this study were all derived from the LiDAR data for elevation. The land-surface

569    derivatives (e.g. slope gradient, relative elevation) described the surface geometry with which the

570    climate interacts. For example, aspect has been shown to influence the amount of solar insolation a

571    hillslope receives (Hunckler and Schaetzl, 1997; Beaudette and O'Geen, 2009). The surface geometry

572    is also known to direct water flow, which affects erosion processes and groundwater recharge

573    (Huggett, 1975; Zevenbergen and Thorne, 1987). Hydrologic predictors (e.g. flow accumulation,

574    catchment slope) provided additional information about the relative volume and energy that the

575    water flow may have (Moore et al., 1991; Wilson and Gallant, 2000).

576    *4.1.2. Topsoil model predictors*

577        All of the topsoil models generated by Cubist relied on DTA predictors the most. Of those

578    predictors, different a-scales of relative elevation, topographic position index (TPI), and aspect were

579    the most commonly used. With the exception of the direct SOC$_{stock}$ model, every topsoil model also

580    included one or two predictors indicative of flow accumulation (i.e. flow path length, SAGA wetness

581    index, or modified catchment area).

582        Aspect at different a-scales influenced predictions for three of the indirect topsoil models. The

583    Cubist generated model identified decreasing topsoil SOC$_\%$ on more north facing slopes (155 m a-

584    scale), which corresponds with a potential decrease in plant productivity due to less solar insolation.

585    Aspect (215 m a-scale) was also used to predict higher topsoil BD on south to west facing slopes,

586    especially on topographic (2000 m a-scale) and micro-topographic (20 m a-scale) highs. Additionally,

587    aspect at a variety of a-scales was used to predict decreasing topsoil SK for low TPI areas facing

588    southeast to southwest. Together, these models suggested a pattern of increased erosion and

589    deposition along the southern sides of hillslopes. This type of pattern has been observed before in

590    other landscapes and has been attributed to topo-climatic differences such as exposure to storms,

591    differences in temperature regime, rainfall effectiveness, or vegetation density (Kennedy, 1976;

592    Churchill, 1981; Cuff, 1985; Weaver, 1991).

593        Although DTA parameters dominated the topsoil models, their predictions were often modified

594    by spectral variables. For example, the primary distinction for predicting topsoil H was between low

595    and high relative elevations. Low relative elevations had a mean topsoil H that was about 20 cm

596    thicker than high relative elevations (1,100 m a-scale). Within most MLR equations, however,

597    predictions were increased by less blue and more green reflectance in early July. This combined use

598    of blue and green bands indicated increasing topsoil H with more productive vegetation on wetter

599    soils. In summary, the dominant pattern identified by the model was between high-low ground

600    (Bushnell, 1943; Sommer et al., 2008), but the degree of topsoil thinning or thickening was predicted

601    by the vegetation's response to soil conditions.

602        Cubist selected a much simpler combination of only DTA predictors to directly model the topsoil

603    $SOC_{stock}$. In general, the model predicted increasing $SOC_{stock}$ with decreasing vertical distance to

604    channel. Areas low in relative elevation (1,100 m a-scale) and not far above the channel network

605    were predicted to have the largest $SOC_{stock}$.  However, for areas low in relative elevation, but

606    sufficiently above the DEM based channel network, the model predicted the opposite trend of the

607    $SOC_{stock}$ *decreasing* with decreasing vertical distance to channel. This pattern identified by the model

608    may be explained by a corresponding pattern observed in the model for the topsoil H. In that model,

609    areas low in relative elevation (1,100 m a-scale) were predicted to have some of the thickest topsoil

610    stocks. However, within a few of those zones the modelled topsoil H decreased with decreasing

611    relative elevation and TPI. This trend in the observed data, as detected by Cubist, was potentially

612    caused by an eroding out of topsoil sediments closer to the center of drainageways. In which case,

613    the vertical distance to channel – used in the topsoil $SOC_{stock}$ model - may have been more an

614    indicator of proximity to the channel than wetness; the threshold was only 0.5 m above the channel

615    modelled from the DEM. Predictors related to surface flow energy would have been expected to be

616    better predictors of this kind of process. However, the upslope drainage network for much of the

617    map area extended beyond the boundaries of the available data. Thus the use of local elevation data

618    may have been a better proxy in this case, compared to the predictors calculated from truncated

619    watersheds.

620    *4.1.3. Subsoil model predictors*

621    With the exception of $SOC_\%$, the subsoil models all used several predictors from Landsat.

622    Selection of Landsat predictors for subsoil models suggested that vegetation characteristics or

623    surface soil moisture at different times of the year indicated subsoil conditions. In contrast, the

624    subsoil $SOC_\%$ model's complete dependence on DTA predictors suggested that soil property was

625    mostly related to hydrology and that vegetation had little response to or effect on the SOC content

626    in the subsoil.

627    An example of spectral predictors detecting vegetation characteristics that likely reflected

628    subsoil conditions was the subsoil SK model. All of the MLR equations were strongly influenced by

629    the predictors of stream power, catchment slope, or SAGA wetness index. However, the ~~skeleton~~ SK

630    predictions were modified by green reflectance in June and additional Landsat predictors collected

631    at different times of the year that related to the vigor of the vegetation. The weaker or drier the

632    vegetation appeared, the higher the prediction of SK content in the subsoil. Assuming soil moisture

633    conditions did not reach detrimental levels that year, these patterns fit known relationships

634    between particle size, soil drainage, and timing to crop maturity (Day and Intalap, 1970; Rawls et al.,

635    1982).

636    The generated model for subsoil BD most likely utilized a relationship with soil moisture as

637    detected by spectral predictors. In all areas, the MLR equations decreased predictions of subsoil BD

638    with increasing reflectance in the blue and SWIR-1 bands along with increasing emission in the TIR

639    band. Increases in the normalized difference vegetation index (NDVI) were used to slightly increase

640    predictions of subsoil BD. The use of the NDVI to offset the decreasing BD predicted by the other

641    Landsat predictors suggested those variables were indicating soil moisture conditions. Locations that

642    are wetter due to surface runoff would have a greater potential for organic material to be

643    translocated deeper in the soil profile (Schaetzl, 1986; Schaetzl, 1990). Also, the association of

644    wetter environments with cooler temperatures and anaerobic conditions would also inhibit

645    decomposition (Gates, 1942; Krause et al., 1959; Frazier and Lee, 1971).

646    The subsoil $SOC_\%$ model was different than the other subsoil models generated. Instead of

647    selecting spectral predictors, the subsoil $SOC_\%$ model relied solely on DTA predictors. The model

648    predicted the highest subsoil $SOC_\%$ on steeper mid-slopes. The pattern of increasing subsoil $SOC_\%$

649    from the upper to middle slope fit the landscape translocation model proposed by Sommer et al.

650    (2000). In that study, the $SOC_\%$ in the Bh horizon increased from the upper slope to the midslope due

651    to lateral translocation. Different than the pattern identified in the present study, the data in

652    Sommer et al. (2000) showed a continued increase in the $SOC_\%$ of Bh horizons in the downslope

653    position. However, this contradiction may be partially explained by aggradation where the slope

654    gradient declines and the topsoil stock has been overthickened by developmental upbuilding

655    (McDonald and Busacca, 1990; Almond and Tonkin, 1999). Also, lateral flow would be expected to

656    return closer to the surface at downslope positions. In Sommer et al. (2000), while the upslope and

657    midslope profiles had E horizons separating the Bh from A horizons, the downslope Bh horizons

658    were exceptionally thick with little to no division between them and the A horizon. In that situation,

24

659    the definition of topsoil used in the present study would have grouped the downslope Bh horizons

660    into the topsoil stock. Therefore, the Cubist generated model may have been a simplification of the

661    complex interaction between topography and lateral flow depth and direction.

662    The rule groups for subsoil $SOC_\%$ also differentiated for the plan curvature where the slope

663    gradient was not too high and the stream power index (SPI) was not too low. Concave plan

664    curvatures (138 m a-scale) were predicted to have increasingly higher and convex plan curvatures

665    were predicted to have increasingly lower subsoil $SOC_\%$. This relationship with plan curvature

666    matches patterns of water movement identified to be important to soil formation by Huggett (1975),

667    where convergent footslopes have the highest deposition rates (Pennock and De Jong, 1987).

668    Assuming the absence of any restrictive layer below, areas with the highest sediment deposition

669    rates would be expected to also have the highest volume of water infiltration.

670    The Cubist generated model for predicting the subsoil $SOC_{stock}$ was simpler than any of the

671    indirect component models. It used only one MLR equation to relate red and infrared predictors to

672    subsoil $SOC_{stock}$. This model predicted more $SOC_{stock}$ storage with increasing reflectance in the red

673    and SWIR-2 bands along with increasing emission in the TIR band – primarily captured on 6 July. Of

674    these variables, model predictions were dominated by increasing reflectance in the red band

675    increasing the estimated subsoil $SOC_{stock}$. This suggested less productive vegetation corresponding

676    with larger subsoil $SOC_{stock}$. This trend was counter to the patterns observed in the topsoil models,

677    but was sensible in the context of how the subsoil stock was defined for this study. Although the

678    *total* $SOC_{stock}$ was less in areas with lower plant productivity, the subsoil $SOC_{stock}$ was larger relative

679    to other subsoil areas due to the inverse relationship between topsoil and subsoil H used in this

680    study. A thicker topsoil stock would mean a thinner subsoil stock – and vice versa – due to the 2 m

681    depth limit. Regarding the other predictors in this model, increases in SWIR-2 reflectance could have

682    indicated more plant productivity. However, its use with the TIR band suggested that together they

683    were indicators of wetter soil conditions.

684    *4.2. Unconventional predictor selections*

685        The Cubist software made some intriguing selections in regards to predictors that were

686    calculated using alternative approaches. One example of this was the selection of alternative types

687    of aspect predictors. The conversion of aspect to northness and eastness is generally considered to

688    be the preferred method for addressing the circular problem of using aspect as a predictor. In our

689    approach of including many different predictors in the available pool, we also experimented with

690    simply rotating the central angle (position of 0°) to each cardinal direction for creating different

691    aspect predictors. In the models generated for this study, northness and eastness were only selected

692    for the topsoil $SOC_{\%}$ model. In contrast, rotated versions of aspect were selected for the topsoil

693    $SOC_{\%}$, topsoil BD, as well as the topsoil and subsoil SK models.

694        Another example of an intriguing predictor selection by Cubist was the use of bands from the

695    LandsatLook products. These images were limited to four bands (SWIR-1, NIR, red, and TIR) and

696    were smoothed by an algorithm to facilitate image selection and visual interpretation. Although the

697    USGS does not recommend the use of these files for data analysis, the Cubist data mining found

698    them to be more useful than the data without LandsatLook processing. Most of these selections can

699    be explained by the greater variety of LandsatLook dates provided in the predictor pool. However,

700    there were a few instances where Cubist chose LandsatLook data over the unprocessed version of

701    the same Landsat data.

702    *4.3. Error propagation*

703        Although both the direct and indirect modelling approaches had base maps with a 2 m

704    resolution available to them, the direct modelling approach produced a more generalized $SOC_{stock}$

705    map. In terms of predicted error, the cost of trying to account for the variation in all of the variables

706    related to the $SOC_{stock}$ appeared to be larger relative errors. The $SOC_{stock}$ model from the direct

707    approach, on the other hand, did not attempt to predict as many variations occurring at small

708    phenomenon scales. Because these very local variations were difficult to predict, the estimated error

709    for the direct approach was less than for the indirect approach for most of the map area. Therefore,

710    it may be appropriate to consider the direct modelling approach to be a conservative approach for

711    estimating the $SOC_{stock}$ for landscapes.

712        Possible sources of error in the base maps included atmospheric conditions for the satellite data

713    and the estimation of bare earth elevation under dense vegetation for the DEM. Several spectral

714    capture dates were made available in the predictor pool to enable Cubist to not only select the

715    optimal changes in seasonal vegetation characteristics, but to also select the image with minimal

716    noise from atmospheric effects such as clouds. Fewer options were available for DTA predictors,

717    because all DTA predictors needed to be derived from the same high-resolution DEM. The effect of

718    anomalies in the elevation data was more pronounced for larger a-scales. For example, a small forest

719    plot – located roughly between the two larger cities in the center of the map area – had not been

720    fully filtered out by the bare-earth algorithm. Any DTA calculation that included this area in its

721    analysis neighborhood was incorrectly influenced by those elevation values. The impact on this

722    study's models was an increased prediction of $SOC_{stock}$ in the surrounding area.

723        The error propagation method used in this study could not directly account for errors in the base

724    maps. Instead, it could only quantify the combined model, base map, and target variable error

725    observed at sample locations. Although none of the sample points were in proximity to the before

726    mentioned error in the DEM, this phenomenon of elevation error affecting scale-dependent

727    predictors would have applied universally, even where the error was less obvious. The higher

728    relative error for both mapping approaches in the area surrounding the known problem in the DEM

729    suggested this potential source of error was at least partially accounted for.

730    **5.  Conclusions**

731        This study demonstrated the use of spatial association to predict the $SOC_{stock}$ and the estimated

732    error at unsampled locations within a 122 km$^2$ landscape at a high-resolution. The Cubist data

733    mining software detected patterns in the observed soil data, which was used to predict soil

734    properties in the greater map region. The ability of the available base maps to predict the variation

735    of those soil properties was quantified for each conditional rule of the respective models. The spatial

736    characteristics of the model rules allowed the uncertainty to be mapped along with the target

737    variable prediction.

738    There were two main advantages to using data mining software to produce relatively simple

739    model structures. First, patterns between the predictors and target variables were objectively

740    identified. Second, the resulting models were simple enough to be interpreted by the user and

741    related to known processes in the soil system.  A relationship between selected predictors and

742    known processes provided confidence that their use in the model was not coincidental.  The

743    separate modelling of topsoil and subsoil stocks identified a general division between useful

744    predictors for predicting soil properties at different depths. The data mining in this study suggested

745    DTA predictors tend to be most useful for topsoil properties, while spectral characteristics of

746    vegetation and soil moisture tend to be more useful for indicating subsoil properties.

747    Direct and indirect approaches were tested for predicting the $SOC_{stock}$ with the rule-based, MLR

748    spatial modelling method. Although the spatial patterns in the two maps were generally similar, the

749    indirect approach produced a map with more spatial variation. While attempting to account for

750    more sources of variability resulted in less estimated error for the topsoil (indirect MAE = 1.69, direct

751    MAE = 2.27), the indirect approach had a higher potential for error in the subsoil (indirect MAE =

752    2.75, direct MAE = 1.37). Because the direct approach accounts for less variation (topsoil: direct $R^2$ =

753    0.58, indirect $R^2$ = 0.73; subsoil: direct $R^2$ = 0.14, indirect $R^2$ = 0.34), but also results in a lower total

754    MAE (direct MAE = 3.64, indirect MAE = 4.44), it should be considered a more conservative

755    prediction of the $SOC_{stock}$'s spatial distribution. The choice of which approach is best will likely

756    depend on a given situation's need to prioritize the representation of spatial pattern or to minimize

757    estimated error.

28

762 **References**

763 Adhikari, K., Kheir, R.B., Greve, M.B., and Greve, M.H.: Comparing kriging and regression approaches
764       for mapping soil clay content in a diverse Danish landscape, Soil Science, 178(9), 505-517,
765       doi:10.1097/SS.0000000000000013, 2013.

766 Almond, P.C. and Tonkin, P.J.: Pedogenesis by upbuilding in an extreme leaching and weathering
767       environment, and slow loess accretion, south Westland, New Zealand, Geoderma, 92(1-2), 1-
768       36, doi: 10.1016/S0016-7061(99)00016-6, 1999.

769 Angers, D.A. and Carter, M.R.: Aggregation and organic matter storage in cool, humid agricultural
770       soils, in: Structure and Organic Matter Storage in Agricultural Soils, Carter, M.R. and Stewart,
771       B.A. (Eds.), CRC Press, Boca Raton, 193-211, 1996.

772 Ashley, M.D. and Rea, J.: Seasonal vegetation differences from ERTS imagery, Journal of American
773       Society of Photogrammetry, 41(6), 713-719, 1975.

774 Bannari, A., Morin, D., Bonn, F., and Huete, A.R.: A review of vegetation indices, Remote Sensing
775       Reviews, 13(1-2), 95-120, doi:10.1080/02757259509532298, 1995.

776 Bartholic, J.F., Namken, L.N., and Wiegand, C.L.: Aerial thermal scanner to determine temperature of
777       soils and of crop canopies differing in water stress, Agronomy Journal, 64(5), 603-608, 1972.

778 Bartholomeus, H.M., Schaepman, M.E., Kooistra, L., Stevens, A., Hoogmoed, W.B., and Spaargaren,
779       O.S.P.: Spectral reflectance based indices for soil organic carbon quantification, Geoderma,
780       145(1-2), 28-36, doi:10.1016/j.geoderma.2008.01.010, 2008.

781 Batjes, N.H.: Total carbon and nitrogen in the soils of the world, European Journal of Soil Science, 47,
782       151–163, doi:10.1111/j.1365-2389.1996.tb01386.x, 1996.

783 Beaudette, D.E. and O'Geen, A.T.: Quantifying the aspect affect: an application of solar radiation
784       modeling for soil survey, Soil Science Society of America Journal, 73(4), 1345-1352, 2009.

785 Bowman, R.A., Reeder, J.D., and Wienhold, B.J.: Quantifying laboratory and field variability to assess
786       potential for carbon sequestration, Communications in Soil Science and Plant Analysis, 33(9-10),
787       1629-1642, doi:10.1081/CSS-120004304, 2002.

788 Brenning, A., Koszinski, S., and Sommer, M.: Geostatistical homogenization of soil conductivity across
789       field boundaries, Geoderma, 143, 254-260, doi:10.1016/j.geoderma.2007.11.007, 2008.

790 Bui, E.N., Henderson, B.L., and Viergever, K.: Knowledge discovery from models of soil properties
791       developed through data mining, Ecological Modelling, 191, 431-446,
792       doi:10.1016/j.ecolmodel.2005.05.021, 2006.

793 Bushnell, T.M.: Some aspects of the soil catena concept, Soil Science Society Proceedings, 7(C), 466-
794       476, 1943.

795 Cambardella, C.A., Moorman, T.B., Novak, J.M., Parkin, T.B., Karlen, D.L., Turco, R.F., and Konopka,
796       A.E.: Field-scale variability of soil properties in central Iowa soils, Soil Science Society of
797       America, 58, 1501-1511, doi:10.2136/sssaj1994.03615995005800050033x, 1994.

798 Carlson, T.N., Gilles, R.R., and Perry, E.M.: A method to make use of thermal infrared temperature
799       and NDVI measurements to infer surface soil water content and fractional vegetation cover,
800       Remote Sensing Reviews, 9(1-2), 161-173, doi:10.1080/02757259409532220, 1994.

801 Churchill, R.R.: Aspect-related differences in badlands slope morphology, Annals of the Association
802 of American Geographers, 71(3), 374-388, doi:10.1111/j.1467-8306.1981.tb01363.x, 1981.

803 Cuff, J.R.I.: Quantifying erosion-causing parameters in a New Zealand watershed, in: Soil
804 Conservation, El-Swaify, S.A., Moldenhauer, W.C., Lo, A. (Eds.), Soil Conservation Society of
805 America, Ankeny, 99-112, 1985.

806 Daughtry, C.S.T.: Discriminating crop residues from soil by shortwave infrared reflectance, Agronomy
807 Journal, 93, 125-131, doi:10.2134/agronj2001.931125x, 2001.

808 Day, A.D. and Intalap, S.: Some effects of soil moisture stress on the growth of wheat (Triticum
809 aestivum L. em Thell), Agronomy Journal, 62(1), 27-29,
810 doi:10.2134/agronj1970.00021962006200010009x, 1970.

811 European Commission: European Soil Database v2, European Soil Data Centre,
812 http://eusoils.jrc.ec.europa.eu, last access: 7 October 2014.

813 FAO/IIASA/ISRIC/ISSCAS/JRC: Harmonized World Soil Database (version 1.2), FAO, Rome, Italy and
814 IIASA, Laxenburg, Austria, 2012.

815 Frazier, B.E. and Lee, G.B.: Characteristics and classification of three Wisconsin Histosols, Soil Science
816 Society of America Journal, 35(5), 776-780, doi:10.2136/sssaj1971.03615995003500050040x,
817 1971.

818 Gates, F.C.: The bogs of northern lower Michigan, Ecological Monographs, 12, 216-254,
819 doi:10.2307/1943542, 1942.

820 GLCF: Global Land Cover Facility, University of Maryland, http://glcf.umd.edu/data, last access: 19
821 February 2014.

822 Goidts, E., Van Wesemael, B., and Crucifix, M.: Magnitude and sources of uncertainties in soil organic
823 carbon (SOC) stock assessments at various scales, European Journal of Soil Science, 60(5), 723-
824 739, doi:10.1111/j.1365-2389.2009.01157.x, 2009.

825 Gomez, C., Viscarra Rossel, R.A., and McBratney, A.B.: Soil organic carbon prediction by
826 hyperspectral remote sensing and field vis-NIR spectroscopy: An Australian case study,
827 Geoderma, 146(3-4), 403-411, doi:10.1016/j.geoderma.2008.06.011, 2008.

828 Goodman, L.A.: On the exact variance of products, Journal of the American Statistical Association,
829 55, 708-713, 1960.

830 Grace, J.: Understanding and managing the global carbon cycle, Journal of Ecology, 92(2), 189-202,
831 doi:10.1111/j.0022-0477.2004.00874.x, 2004.

832 Grimm, R., Behrens, T., Märker, M., and Elsenbeer, H.: Soil organic carbon concentrations and stocks
833 on Barro Colorado Island – digital soil mapping using random forests analysis, Geoderma,
834 146(1-2), 102-113, doi:10.1016/j.geoderma.2008.05.008, 2008.

835 Hatfield, J.L., Gitelson, A.A., Schepers, J.S., and Walthall, C.L.: Application of spectral remote sensing
836 for agronomic decisions, Agronomy Journal, 100(3), S-117 – S-131,
837 doi:10.2134/agronj2006.0370c, 2008.

838 Heilman, J.L., Kanemasu. E.T., Rosenberg, N.J., and Blad, B.L.: Thermal scanner measurement of
839 canopy temperatures to estimate evapotranspiration, Remote Sensing of Environment, 5(C),
840 137-145, doi:10.1016/0034-4257(76)90044-4, 1976.

841 Holmes, G., Hall, M., and Frank, E.: Generating rule sets from model trees, Advanced Topics in
842      Artificial Intelligence, Lecture Notes in Computer Science, 1747, 1-12, doi:10.1007/3-540-
843      46695-9_1, 1999.

844 Huete, A.R.: A soil-adjusted vegetation index (SAVI), Remote Sensing of Environment, 25, 295-309,
845      doi:10.1016/0034-4257(88)90106-X, 1988.

846 Huggett, R.J.: Soil landscape systems: a model of soil genesis, Geoderma, 13, 1-22,
847      doi:10.1016/0016-7061(75)90035-X, 1975.

848 Hunckler, R.V. and Schaetzl, R.J.: Spodosol development as affected by geomorphic aspect, Baraga
849      County, Michigan, Soil Science Society of America Journal, 61(4), 1105-1115,
850      doi:10.2136/sssaj1997.03615995006100040017x, 1997.

851 Jobbágy, E.G. and Jackson, R.B.: The vertical distribution of soil organic carbon and its relation to
852      climate and vegetation, Ecological Applications, 10(2), 423–436, doi:10.1890/1051-
853      0761(2000)010[0423:TVDOSO]2.0.CO;2, 2000.

854 Johnson, D.L., Keller, E.A., and Rockwell, T.K.: Dynamic pedogenesis: new views on some key
855      concepts, and a model for interpreting quaternary soils, Quaternary Research, 33(3), 306-319,
856      doi:10.1016/0033-5894(90)90058-S, 1990.

857 Johnston, A.E., Poulton, P.R., and Coleman, K.: Soil organic matter: its importance in sustainable
858      agriculture and carbon dioxide fluxes, Advances in Agronomy, 101, 1-57, doi:10.1016/S0065-
859      2113(08)00801-8, 2009.

860 Johnston, C.A., Groffman, P., Breshears, D.D., Cardon, Z.G., Currie, W., Emanuel, W., Gaudinski, J.,
861      Jackson, R.B., Lajtha, K., Nadelhoffer, K., Nelson, D., Jr., Mac Post, W., Retallack, G., and
862      Wielopolski, L.: Carbon cycling in soil, Frontiers in Ecology and the Environment, 2(10), 522-528,
863      doi:10.2307/3868382, 2004.

864 Kay, B.D.: Soil structure and organic carbon: a review, in: Soil Processes and the Carbon Cycle, Lal, R.,
865      Kimble, J.M., Follett, R.F., and Stewart, B.A. (Eds.), CRC Press, Boca Raton, 169-197, 1998.

866 Kempen, B., Brus, D.J., and Stoorvogel, J.J.: Three-dimensional mapping of soil organic matter
867      content using soil type-specific depth functions, Geoderma, 162, 107-123,
868      doi:10.1016/j.geoderma.2011.01.010, 2011.

869 Kennedy, B.A.: Valley-side slopes and climate, in: Geomorphology and Climate, Derbyshire, E. (Ed.),
870      John Wiley, London, 171-201, 1976.

871 Khalil, M.I., Kiely, G., O'Brien, P., and Müller, C.: Organic carbon stocks in agricultural soils in Ireland
872      using combined empirical and GIS approaches, Geoderma, 193-195, 222-235,
873      doi:10.1016/j.geoderma.2012.10.005, 2013.

874 Königlich Preußische Geologische Landesanstalt: Geologische Karte von Preußen und benachbarten
875      Bundesstaaten, 1:25,000 (Geological Map of Prussia and adjacent Federal States, 1:25,000),
876      Landesamt f. Geologie und Bergwesen, Halle, Sachsen-Anhalt, Germany, Sheet Wulfen 4137,
877      1913a.

878 Königlich Preußische Geologische Landesanstalt: Geologische Karte von Preußen und benachbarten
879      Bundesstaaten, 1:25,000 (Geological Map of Prussia and adjacent Federal States, 1:25,000),
880      Landesamt f. Geologie und Bergwesen, Halle, Sachsen-Anhalt, Germany, Sheet Cöthen 4237,
881      1913b.

882 Krause, H.H., Rieger, S., and Wilde, S.A.: Soils and forest growth on different aspects in the Tanana
883     watershed of interior Alaska, Ecology, 40, 492-495, doi:10.2307/1929773, 1959.

884 Kravchenko, A.N., Robertson, G.P., Hao, X., and Bullock, D.G.: Management practice effects on
885     surface total carbon: differences in spatial variability patterns, Agronomy Journal, 98, 1559-
886     1568, doi:10.2134/agronj2006.0066, 2006a.

887 Kravchenko, A.N., Robertson, G.P., Snap, S.S., and Smucker, A.J.M.: Using information about spatial
888     variability to improve estimates of total soil carbon, Agronomy Journal, 98, 823-829,
889     doi:10.2134/agronj2005.0305, 2006b.

890 Ku, H.H.: Notes on the use of propagation of error formulas, Journal of Research of the National
891     Bureau of Standards – C. Engineering and Instrumentation, 70C(4), 263-273, 1966.

892 Kühn, J., Brenning, A., Wehrhan, M., Koszinski, S., and Sommer, M.: Interpretation of electrical
893     conductivity patterns by soil properties and geological maps for precision agriculture, Precision
894     Agriculture, 10, 490-507, doi:10.1007/s11119-008-9103-z, 2009.

895 Lal, R.: Soil carbon sequestration impacts on global climate change and food security, Science,
896     304(5677), 1623-1627, doi:10.1126/science.1097396, 2004.

897 Lacoste, M., Minasny, B., McBratney, A., Michot, D., Viaud, V., and Walter, C.: High resolution 3D
898     mapping of soil organic carbon in a heterogeneous agricultural landscape, Geoderma, 213, 296-
899     311, doi:10.1016/j.geoderma.2013.07.002, 2014.

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

903 Li, Z.L., Tang, R., Wan, Z., Bi, Y., Zhou, C., Tang, B., Yan, G., and Zhang, X.: A review of current
904     methodologies for regional evapotranspiration estimation from remotely sensed data, Sensors
905     9(5), 3801-3853, doi:10.3390/s90503801, 2009.

906 Liebens, J. and VanMolle, M.: Influence of estimation procedure on soil organic carbon stock
907     assessment in Flanders, Belgium, Soil Use and Management, 19(4), 364-371,
908     doi:10.1111/j.1475-2743.2003.tb00327.x, 2003.

909 Lobell, D.B. and Asner, G.P.: Moisture effects on soil reflectance, Soil Science Society of America
910     Journal, 66, 722-727, doi:10.2136/sssaj2002.7220, 2002.

911 Lowther, J.R., Smethurst, P.J., Carlyle, J.C., and Nambiar, E.K.S.: Methods for determining organic
912     carbon in podzolic sands, Communications in Soil Science and Plant Analysis, 21(5-6), 457-470,
913     doi:10.1080/00103629009368245, 1990.

914 Lufafa, A., Diédhiou, I., Samba, S.A.N., Séné, M., Khouma, M., Kizito, F., Dick, R.P., Dossa, E., and
915     Noller, J.S.: Carbon stocks and patterns in native shrub communities of Senegal's Peanut Basin,
916     Geoderma, 146, 75-82, doi:10.1016/j.geoderma.2008.05.024, 2008.

917 Malone, B.P., McBratney, A.B., and Minasny, B.: Empirical estimates of uncertainty for mapping
918     continuous depth functions of soil attributes, Geoderma, 160, 614-626,
919     doi:10.1016/j.geoderma.2010.11.013, 2011.

920 Mardia, K.V., Kent, J.T., and Bibby, J.M.: Multivariate Analysis, Academic Press, London, United
921     Kingdom, 521 pp., 1979.

922 Martin, M.P., Orton, T.G., Lacarce, E., Meersmans, J., Saby, N.P.A., Paroissien, J.B., Jolivet, C.,
923     Boulonne, L., and Arrouays, D.: Evaluation of modelling approaches for predicting the spatial
924     distribution of soil organic carbon stocks at the national scale, Geoderma, 223-225, 97-107,
925     doi:10.1016/j.geoderma.2014.01.005, 2014.

926 Martin, M.P., Wattenbach, M., Smith, P., Meersmans, J., Jolivet, C., Boulonne, L., and Arrouays, D.:
927     Spatial distribution of soil organic carbon stocks in France, Biogeosciences, 8, 1053-1065,
928     doi:10.5194/bg-8-1053-2011, 2011.

929 McBratney, A.B. and Pringle, M.J.: Estimating average and proportional variograms of soil properties
930     and their potential use in precision agriculture, Precision Agriculture, 1, 125-152,
931     doi:10.1023/A:1009995404447, 1999.

932 McDonald, E.V. and Busacca, A.J.: Record of pre-late Wisconsin giant floods in the Channeled
933     Scabland interpreted from loess deposits, Geology, 16, 728-731, doi:10.1130/0091-
934     7613(1988)0162.3.CO;2, 1988.

935 Meersmans, J., Van Wesemael, B., De Ridder, F., Fallas Dotti, M., De Baets, S., and Van Molle, M.:
936     Changes in organic carbon distribution with depth in agricultural soils in northern Belgium,
937     1960-2006, Global Change Biology, 15(11), 2739-2750, doi:10.1111/j.1365-2486.2009.01855.x,
938     2009.

939 Migdall, S., Bach, H., Bobert, J., Wehrhan, M., and Mauser, W.: Inversion of a canopy reflectance
940     model using hyperspectral imagery for monitoring wheat growth and estimating yield, Precision
941     Agriculture, 10, 508-524, doi:10.1007/s11119-009-9104-6, 2009.

942 Miller, B.A., Koszinski, S., Wehrhan, M., and Sommer, M., Impact of multi-scale predictor selection
943     for modeling soil properties, Geoderma, 239-240, 97-106,
944     doi:10.1016/j.geoderma.2014.09.018, 2015.

945 Minasny, B. and McBratney, A.B.: Regression rules as a tool for predicting soil properties from
946     infrared reflectance spectroscopy, Chemometrics and Intelligent Laboratory Systems, 94(1), 72-
947     79, doi:10.1016/j.chemolab.2008.06.003, 2008.

948 Minasny, B., McBratney, A.B., Malone, B.P., and Wheeler, I.: Digital mapping of soil carbon, Advances
949     in Agronomy, 118, 1-47, doi:10.1016/B978-0-12-405942-9.00001-3, 2013.

950 Mishra, U., Lal, R., Liu, D., and Van Meirvenne, M.: Predicting the spatial variation of the soil organic
951     carbon pool at a regional scale, Soil Science Society of America Journal, 74, 906-914,
952     doi:10.2136/sssaj2009.0158, 2010.

953 Moore, I.D., Grayson, R.B., and Ladson, A.R.: Digital terrain modelling: a review of hydrological,
954     geomorphological, and biological applications, Hydrological Processes, 5(1), 3-30,
955     doi:10.1002/hyp.3360050103, 1991.

956 Mulder, V.L., de Bruin, S., Schaepman, M.E., and Mayr, T.R.: The use of remote sensing in soil and
957     terrain mapping – a review, Geoderma, 162, 1-19, doi:10.1016/j.geoderma.2010.12.018, 2011.

958 Myneni, R.B., Hall, F.G., Sellers, P.J., and Marshak, A.L.: The interpretation of spectral vegetation
959     indexes, IEEE Transactions on Geoscience and Remote Sensing, 33(2), 481-486,
960     doi:10.1109/36.377948, 1995.

961 Neemann, W.: Bestimmung des Bodenerodierbarkeitsfaktors für winderosionsgefährdete Böden
962     Norddeutschlands (Determination of soil erodibility factors for wind-erosion endangered soils
963     in Northern Germany), Geologisches Jahrbuch Reihe F, 25, 131 pp., 1991.

964 Nyssen, J., Temesgen, H., Lemenih, M., Zenebe, A., Haregeweyn, N., and Haile, M.: Spatial and
965     temporal variation of soil organic carbon stocks in a lake retreat area of the Ethiopian Rift
966     Valley, Geoderma, 146, 261-268, doi:10.1016/j.geoderma.2008.06.007, 2008.

967 Orton, T.G., Pringle, M.J., Page, K.L., Dalal, R.C., and Bishop, T.F.A.: Spatial prediction of soil organic
968     carbon stock using a linear model of coregionalisation, Geoderma, 230-231, 119-130,
969     doi:10.1016/j.geoderma.2014.04.016, 2014.

970 Panda, D.K., Singh, R., Kundu, D.K., Chakraborty, H., and Kumar, A.: Improved estimation of soil
971     organic carbon storage uncertainty using first-order Taylor series approximation, Soil Science
972     Society of America Journal, 72, 1708-1710, doi:10.2136/sssaj2007.0242N, 2008.

973 Pennock, D.J. and De Jong, E.: The influence of slope curvature on soil erosion and deposition in
974     hummock terrain, Soil Science, 144(3), 209-217, doi:10.1097/00010694-198709000-00007,
975     1987.

976 Petropoulus, G., Carlson, T.N., Wooster, M.J., and Islam, S.: A review of Ts/VI remote sensing based
977     methods for the retrieval of land surface energy fluxes and soil surface moisture, Progress in
978     Physical Geography, 33(2), 224-250, doi:10.1177/0309133309338997, 2009.

979 Phachomphon, K., Dlamini, P., and Chaplot, V.: Estimating carbon stocks at a regional level using soil
980     information and easily accessible auxiliary variables, Geoderma, 155(3-4), 372-380,
981     doi:10.1016/j.geoderma.2009.12.020, 2010.

982 Phillips, J.D.: On the relations between complex systems and the factorial model of soil formation
983     (with discussion), Geoderma, 86, 1-21, doi:10.1016/S0016-7061(98)00054-8, 1998.

984 Powlson, D.S., Whitmore, A.P., and Goulding, K.W.T.: Soil carbon sequestration to mitigate climate
985     change: a critical re-examination to identify the true and the false, European Journal of Soil
986     Science, 62, 42-55, doi:10.1111/j.1365-2389.2010.01342.x, 2011.

987 Quinlan, J.R. Learning with continuous classes, Proceedings of the 5[th] Australian Joint Conference on
988     Artificial Intelligence, 343-348, 1992.

989 Quinlan, J.R.: Combining instance-based and model-based learning, in: Proceedings of the Tenth
990     International Conference on Machine Learning, Kaufmann, M. (Ed.), 236-243, 1993.

991 Quinlan, J.R.: C4.5: Programs for machine learning, Machine Learning, 16, 235-240, 1994.

992 Rasmussen, M.S.: Developing simple, operational, consistent NDVI-vegetation models by applying
993     environmental and climatic information: Part I. Assessment of net primary production,
994     International Journal of Remote Sensing, 19(1), 97-117, doi:10.1080/014311698216459, 1998.

995 Rawls, W.J., Brakensiek, D.L., and Saxton, K.E.: Estimation of soil water properties, Transactions of
996     the American Society of Agricultural Engineers, 25(5), 1316-1320, 1982.

997 Rawls, W.J., Pachepsky, Y.A., Ritchie, J.C., Sobecki, T.M., and Bloodworth, H.: Effect of soil organic
998     carbon on soil water retention, Geoderma, 116(1-2), 61-76, doi:10.1016/S0016-7061(03)00094-
999     6, 2003.

1000 Richards, J.A.: Remote Sensing Digital Image Analysis, Springer, 439 pp., 2006.

1001 Richter, D.D. and Markewitz, D.: How deep is soil?, Bioscience, 45(9), 600-609, doi:10.2307/1312764,
1002     1995.

1003    Schaetzl, R.J.: Complete soil profile inversion by tree uprooting, Physical Geography, 7(2), 181-189,
1004        doi:10.1080/02723646.1986.10642290, 1986.

1005    Schrumpf, M., Schulze, E.D., Kaiser, K., and Schumacher, J.: How accurately can soil organic carbon
1006        stocks and stock changes be quantified by soil inventories?, Biogeosciences, 8(5), 1193-1212,
1007        doi:10.5194/bg-8-1193-2011, 2011.

1008    Schwartz, D. and Namri, M.: Mapping the total organic carbon in the soils of the Congo, Global and
1009        Planetary Change, 33(1–2), 77–93, doi:10.1016/S0921-8181(02)00063-2, 2002.

1010    Selige, S. Böhner, J., and Schmidhalter, U.: High resolution topsoil mapping using hyperspectral
1011        image and field data in multivariate regression modeling procedures, Geoderma, 136(1-2), 235-
1012        244, doi:10.1016/j.geoderma.2006.03.050, 2006.

1013    Simbahan, G.C., Dobermann, A., Goovaerts, P., Ping, J., and Haddix, L.: Fine-resolution mapping of
1014        soil organic carbon based on multivariate secondary data, Geoderma, 132, 471-489,
1015        doi:10.1016/j.geoderma.2005.07.001, 2006.

1016    Snyder, V.A. and Vazquez, M.A.: Structure, in: Encyclopedia of Soils in the Environment, Hillel, D.,
1017        Hatfield, J.L., Powlson, D.S., Rozenweig, C., Scow, K.M., Singer, M.J., and Sparks, D.L. (Eds.),
1018        Elsevier Academic Press, 54-68, 2005.

1019    Shrestha, D.L. and Solomatine, D.P.: Machine learning approaches for estimation of prediction
1020        interval for the model output, Neural Networks, 19(2), 225-235,
1021        doi:10.1016/j.neunet.2006.01.012, 2006.

1022    Sombroek, W.G., Fearnside, P.M., and Cravo, M.: Geographic assessment of carbon stored in
1023        Amazonian terrestrial ecosystems and their soils in particular, in: Global Climate Change and
1024        Tropical Ecosystems, Lal, R., Kimble, J.M., and Stewart, B.A. (Eds.),  CRC Lewis, Boca Raton, 375–
1025        389, 2000.

1026    Sommer, M., Gerke, H.H., and Deumlich, D.: Modelling soil landscape genesis – a "time split"
1027        approach for hummocky agricultural landscapes, Geoderma, 145, 480-493,
1028        doi:10.1016/j.geoderma.2008.01.012, 2008.

1029    Sommer, M., Halm, D., Weller, U., Zarei, M., and Stahr, K.: Lateral podzolization in a granite
1030        landscape, Soil Science Society of America Journal, 64(4), 1434-1442,
1031        doi:10.2136/sssaj2000.6441434x, 2000.

1032    Soon, Y.K. and Abboud, S.: A comparison of some methods for soil organic carbon determination,
1033        Communications in Soil Science and Plant Analysis, 22(9-10), 943-954, doi:
1034        10.1080/00103629109368465, 1991.

1035    Stevens, A., Udelhoven, T., Denis, A., Tychon, B., Lioy, R., Hoffmann, L., and van Wesemael, B.:
1036        Measuring soil organic carbon in croplands at regional scale using airborne imaging
1037        spectroscopy, Geoderma, 158(1-2), 32-45, doi:10.1016/j.geoderma.2009.11.032, 2010.

1038    Sutherland, R.A.: Loss-on-ignition estimates of organic matter and relationships to organic carbon in
1039        fluvial sediments, Hydrobiologia, 389(1-3), 153-167, doi: 10.1023/A:1003570219018, 1998.

1040    Taylor, J.R.: An Introduction to Error Analysis: The Study of Uncertainties in Physical Measurements,
1041        2$^{nd}$ ed. University Science Books, Sausalito, California, USA, 1997.

1042 Ungaro, F., Staffilani, F., and Tarocco, P.: Assessing and mapping topsoil organic carbon stock at
1043     regional scale: a scorpan kriging approach conditional on soil map delineations and land use,
1044     Land Degradation & Development, 21, 565-581, doi:10.1002/ldr.998, 2010.

1045 USGS: Earth Explorer, U.S. Geological Survey, http://earthexplorer.usgs.gov, last access 19 February
1046     2014.

1047 Walter, C., Viscarra Rossel, R.A., and McBratney, A.B.: Spatio-temporal simulation of the field-scale
1048     evolution of organic carbon over the landscape, Soil Science Society of America Journal, 67,
1049     1477-1486, doi:10.2136/sssaj2003.1477, 2003.

1050 Weaver, A. van Breda, The distribution of soil erosion as a function of slope aspect and parent
1051     material in Ciskei, Southern Africa, GeoJournal, 23(1), 29-34, doi:10.1007/BF00204406, 1991.

1052 Weisstein, E.W.: Error Propagation, Wolfram MathWorld,
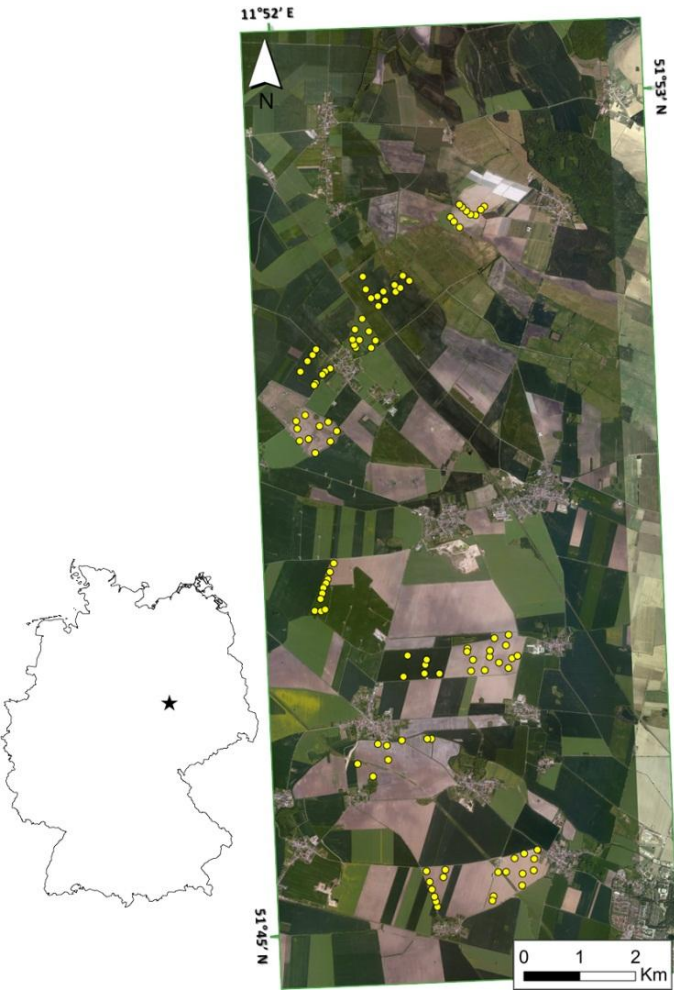1053     http://mathworld.wolfram.com/ErrorPropagation.html, last access: 25 August 2014.

1054 Wilhelm, W.W., Johnson, J.M.F., Hatfield, J.L., Voorhees, W.B., and Linden, D.R.: Crop and soil
1055     productivity response to corn residue removal: a literature review, Agronomy Journal, 96, 1-17,
1056     doi:10.2134/agronj2004.1000, 2004.

1057 Wilson, J.P. and Gallant, J.C. (Eds.): Terrain Analysis: Principles and Applications, John Wiley & Sons,
1058     2000.

1059 Xie, Y., Sha, Z., and Yu, M.: Remote sensing imagery in vegetation mapping: a review. Journal of Plant
1060     Ecology, 1(1), 9-23, doi:10.1093/jpe/rtm005, 2008.

1061 Zevenbergen, L.W. and Thorne, C.R.: Quantitative analysis of land surface topography, Earth Surface
1062     Processes and Landforms, 12(1), 47-56, doi:10.1002/esp.3290120107, 1987.

1063 Zhang, Z., Yu, D., Shi, X., Warner, E., Ren, H., Sun, W., Tan, M., and Wang, H.: Application of
1064     categorical information in the spatial prediction of soil organic carbon in the red soil area of
1065     China, Soil Science and Plant Nutrition, 56, 307-318, doi:10.1111/j.1747-0765.2010.00457.x,
1066     2010.

1067    <u>Figures</u>



1068
1069    Figure 1. Locations of sample points and study area within Germany.

1070
1071　　Figure 2. Topsoil SOC$_{stock}$ modelled by a) the direct approach and b) the indirect approach. Overlaid
1072　　on a hillshade to show relationship with relief and field boundaries.



1073
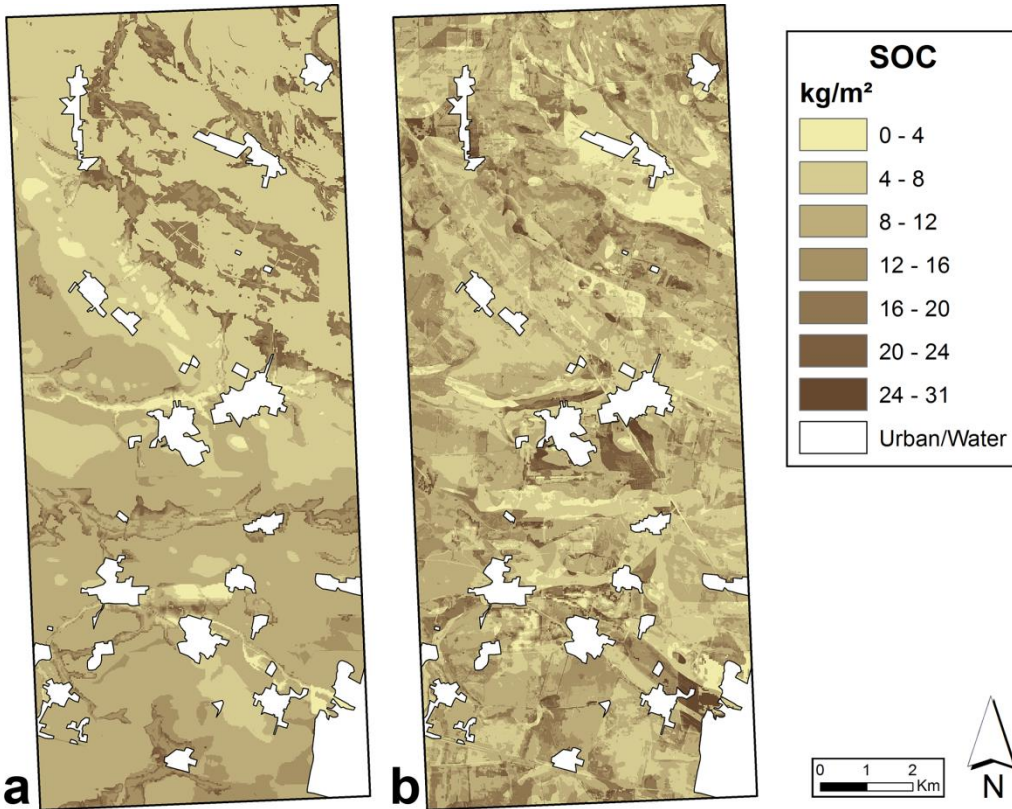1074　　Figure 3. Subsoil SOC$_{stock}$ modelled by a) the direct approach and b) the indirect approach. Overlaid
1075　　on a hillshade to show relationship with relief and field boundaries.

Figure 4. Total SOC$_{stock}$ (topsoil + subsoil) modelled by a) the direct approach and b) the indirect approach. Overlaid on a hillshade to show relationship with relief and field boundaries.

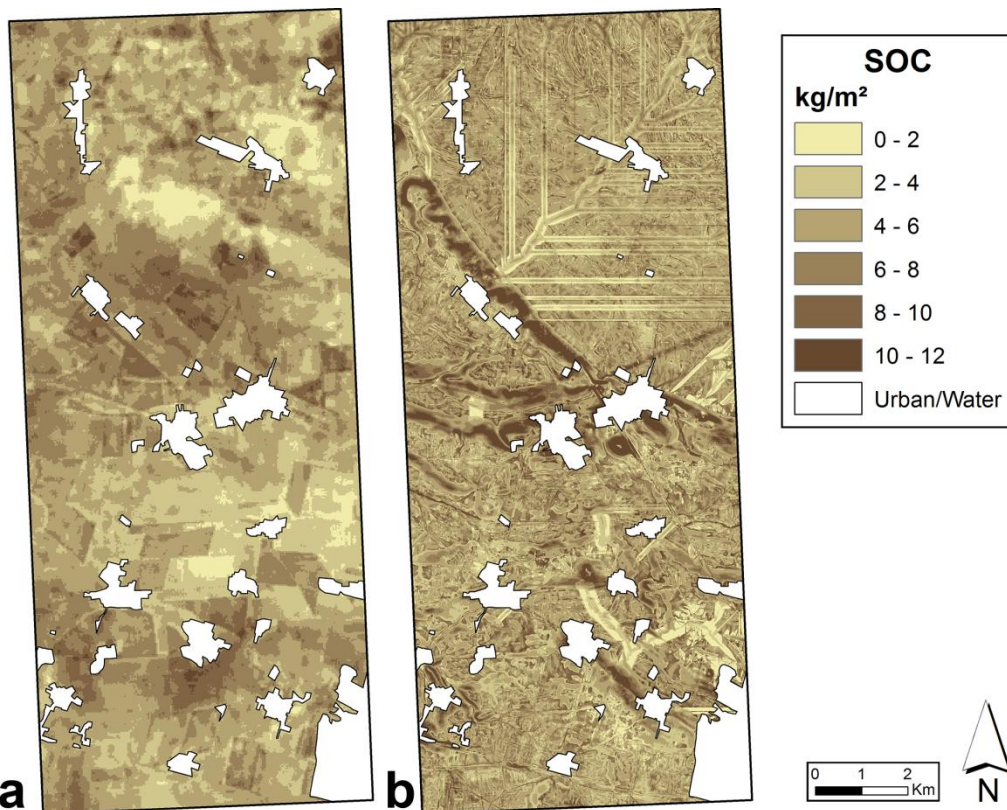

Figure 5. Calculated difference between the direct and indirect approaches of modelling the total SOC$_{stock}$. Negative values are where the indirect approach predicted more SOC$_{stock}$ than the direct approach and positive values are where the indirect approach predicted less.

1083

1084      Figure 6. Estimated relative error for the total $SOC_{stock}$ modelled by a) the direct approach and b) the

1085      indirect approach.

1086 <u>Tables</u>

1087 Table 1. Descriptive statistics for the observed target variables. BD = total bulk density (g cm$^{-3}$), SK =
1088 particles > 2 mm (%), SOC$_{\%}$ = SOC concentration (%), H = stock thickness (cm), and SOC$_{stock}$ = mass of
1089 organic carbon per unit area of soil (kg m$^{-2}$).

| **Topsoil** | **BD** | **SK** | **H** | **SOC$_{\%}$** | **SOC$_{stock}$** |
|---|---|---|---|---|---|
| Min. | 1.18 | 0.00 | 10 | 0.75 | 1.80 |
| Median | 1.50 | 1.30 | 40 | 1.46 | 9.27 |
| Mean | 1.51 | 3.15 | 43.61 | 1.56 | 9.82 |
| Max. | 1.85 | 44.70 | 105 | 4.03 | 28.03 |
| Std. Dev. | 0.11 | 5.50 | 15.35 | 0.53 | 4.49 |
| **Subsoil** | | | | | |
| Min. | 1.33 | 0.00 | 18 | 0.02 | 0.07 |
| Median | 1.63 | 4.07 | 86 | 0.23 | 3.10 |
| Mean | 1.63 | 8.99 | 86.66 | 0.26 | 3.37 |
| Max. | 1.96 | 63.36 | 155 | 0.71 | 9.86 |
| Std. Dev. | 0.13 | 12.28 | 32.60 | 0.13 | 2.04 |

1090

1091    Table 2. Predictor variables considered in this study.

| Predictor | Software | Analysis Scale |
|---|---|---|
| Elevation (LiDAR, bare-earth) | n/a | 2 m |
| Slope gradient | GRASS | 6 - 195 m |
| Profile curvature | GRASS | 6 - 195 m |
| Plan curvature | GRASS | 6 - 195 m |
| Aspect -west {rotated for N, E, and S} | GRASS | 6 - 345 m |
| Aspect (8 classes) | ArcGIS (raster calculator) | 6 - 345 m |
| Northness | transformed from aspect | 6 - 345 m |
| Eastness | transformed from aspect | 6 - 345 m |
| Longitudinal curvature | SAGA | 10 m |
| Cross-section curvature | SAGA | 10 m |
| Convexity | SAGA | 10 m |
| Relative elevation - rect. neighborhood | ArcGIS toolbox | 6 - 4000 m |
| Relative elevation - circ. neighborhood | ArcGIS toolbox | 6 - 4000 m |
| Topographic position index (TPI) | ArcGIS toolbox | 6 - 4000 m |
| TPI - slope position | ArcGIS toolbox | multiple |
| TPI - landform classification | ArcGIS toolbox | multiple |
| Hillslope position | ArcGIS toolbox | multiple |
| Catchment area | SAGA | n/a |
| Catchment slope | SAGA | n/a |
| Channel network base level | SAGA | n/a |
| Convergence index | SAGA | n/a |
| Flow accumulation | SAGA | n/a |
| Flow path length | SAGA | n/a |
| Length-slope factor | SAGA | n/a |
| Modified catchment area | SAGA | n/a |
| Relative slope position | SAGA | n/a |
| SAGA wetness index | SAGA | n/a |
| Stream power | SAGA | n/a |
| Vertical distance to channel | SAGA | n/a |
| Wetness index | SAGA | n/a |
| Geology (1:25,000 legacy map) | n/a | 423 ha (mean) |

1092

1093 Table 2 (cont'd).

| Predictor | Resolution | Date |
|---|---|---|
| AVIS - LAI-green leaf area | 5m | 21 Jun. 2005 |
| AVIS - LAI-brown leaf area | 5m | 21 Jun. 2005 |
| Ikonos | 4 m, 4 bands | 4 Jul. 2006 |
| Ikonos - panchromatic | 1 m | 4 Jul. 2006 |
| Ikonos - LAI | 5m | 4 Jul. 2006 |
| Ikonos - dry matter | 5m | 4 Jul. 2006 |
| Landsat 5 NDVI (USGS, 2014) | 30m | 11 Jun. 2006 |
| Landsat 5 NDVI (USGS, 2014) | 30m | 22 Jul. 2006 |
| Landsat 5 LandsatLook (USGS, 2014) | 30m, 3+1 band | 20 Jun. 2006 |
| Landsat 5 LandsatLook (USGS, 2014) | 30m, 3+1 band | 6 Jul. 2006 |
| Landsat 5 LandsatLook (USGS, 2014) | 30m, 3+1 band | 22 Jul. 2006 |
| Landsat 5 LandsatLook (USGS, 2014) | 30m, 3+1 band | 15 Sep. 2006 |
| Landsat 5 LandsatLook (USGS, 2014) | 30m, 3+1 band | 17 Oct. 2006 |
| Landsat 5 TM (USGS, 2014) | 30m, 6 bands; 60m, 1 band | 11 Jun. 2006 |
| Landsat 5 TM (USGS, 2014) | 30m, 6 bands; 60m, 1 band | 22 Jul. 2006 |
| Landsat 5 SR (GLCF, 2014) | 30m, 7+2 bands | 11 Jun. 2006 |
| Landsat 5 SR (GLCF, 2014) | 30m, 7+2 bands | 22 Jul. 2006 |

1094

1095 Table 3. Relative use (%) of predictors in models derived by Cubist for the topsoil and subsoil stocks.
1096 BD = total bulk density (g cm$^{-3}$), SK = particles > 2 mm (%), SOC$_\%$ = SOC concentration (%), H = stock
1097 thickness (cm), and SOC$_{stock}$ = mass of organic carbon per unit area of soil (kg m$^{-2}$).

| Topsoil | | | Subsoil | | |
|---|---|---|---|---|---|
| *Rules* | *MLR* | *Predictor* | *Rules* | *MLR* | *Predictor* |
| **BD** | | | **BD** | | |
| 100% | 100% | Relative elev. - circ. (2000 m) | 100% | 0% | Geology map units |
| 51% | 100% | Landsat5 SR, band 7 (6 Jun. 2006) | 68% | 100% | LandsatLook, band 5 (6 Jul. 2006) |
| 17% | 100% | Relative elev. - rect. (20 m) | | 100% | Landsat5 NDVI (22 Jul. 2006) |
| | 96% | LandsatLook, band 5 (17 Oct. 2006) | | 100% | LandsatLook, band 6 (6 Jul. 2006) |
| | 87% | Relative elev. - rect. (10 m) | | 100% | Landsat5 TM, band 1 (11 Jun. 2006) |
| | 87% | Aspect, N central angle (215 m) | | 68% | Landsat5 SR, band 7 (22 Jul. 2006) |
| | 83% | Landsat5 SR, band 2 (6 Jun. 2006) | | 32% | Landsat5 SR, band QA (6 Jun. 2006) |
| | 34% | SAGA wetness index | | 32% | Landsat5 SR, band 1 (22 Jul. 2006) |
| | 13% | Relative elev. - circ. (800 m) | | 32% | Landsat5 SR, band 6 (22 Jul. 2006) |
| **SK** | | | **SK** | | |
| 100% | 100% | TPI (70 m) | 100% | 3% | Stream power |
| 94% | 0% | Aspect class (70 m) | 76% | 76% | Landsat5 SR, band 2 (11 Jun. 2006) |
| 39% | 16% | Relative elev. - rect. (550 m) | 21% | 0% | Profile Curvature (118 m) |
| 37% | 14% | LandsatLook, band 6 (17 Oct. 2006) | 15% | 79% | Landsat5 SR, band 4 (6 Jun. 2006) |
| | 94% | Relative elev. - rect. (1800 m) | | 85% | Catchment slope |
| | 84% | Landsat5 NDVI (11 Jun. 2006) | | 76% | LandsatLook, band 3 (20 Jun. 2006) |
| | 80% | Aspect, N central angle (50 m) | | 56% | Landsat5 NDVI (11 Jun. 2006) |
| | 78% | Landsat5 TM, band 4 (20 Jun. 2006) | | 56% | LandsatLook, band 4 (20 Jun. 2006) |
| | 78% | Relative elev. - circ. (3000 m) | | 56% | Aspect, W central angle (70 m) |
| | 64% | Aspect, N central angle (130 m) | | 21% | SAGA wetness index |
| | 64% | Aspect, S central angle (345 m) | | | |
| | 64% | Flow path length | | | |
| | 37% | Aspect, N central angle (295 m) | | | |
| **H** | | | **H** | | |
| 100% | 93% | Relative elev. - rect. (1100 m) | | | |
| 39% | 100% | LandsatLook, band 5 (15 Sept. 2006) | | | *Cubist not used* |
| 34% | 34% | LandsatLook, band 5 (22 Jul. 2006) | | | *(based on 2 m - topsoil thickness)* |
| 25% | 93% | Ikonos, band 2 (4 Jul. 2006) | | | |
| 18% | 7% | LandsatLook, band 4 (17 Oct. 2006) | | | |
| | 100% | Relative elev. - rect. (1200 m) | | | |
| | 93% | Ikonos, band 1 (4 Jul. 2006) | | | |
| | 93% | Relative elev. - rect. (1300 m) | | | |
| | 74% | LandsatLook, band 4 (15 Sept. 2006) | | | |
| | 74% | TPI (1800 m) | | | |
| | 74% | TPI (2600 m) | | | |
| | 74% | Flow path length | | | |
| | 28% | Relative elev. - circ. (650 m) | | | |
| | 7% | Landsat5 TM, band 6 (11 Jun. 2006) | | | |

1098

1099    Table 3 (cont'd).

| Topsoil | | | Subsoil | | |
|---|---|---|---|---|---|
| *Rules* | *MLR* | *Predictor* | *Rules* | *MLR* | *Predictor* |
| **SOC$_{\%}$** | | | **SOC$_{\%}$** | | |
| 100% | 0% | Geology map units | 100% | 100% | Slope gradient (98 m) |
| 49% | 39% | Relative elev. - rect. (3200 m) | 74% | 74% | Stream power |
| 39% | 69% | Relative elev. - rect. (2000 m) | 55% | 55% | Plan curvature (138 m) |
| 33% | 74% | Flow path length | | 74% | Slope gradient (90 m) |
| 21% | 62% | Northness (155 m) | | 74% | Slope gradient (138 m) |
| | 81% | TPI (1200 m) | | 74% | Slope gradient (185 m) |
| | 80% | Relative elev. - rect. (250 m) | | 74% | Relative elev. - rect. (3400 m) |
| | 80% | Northness (345 m) | | 55% | Plan curvature (90 m) |
| | 74% | Aspect, W central angle (90 m) | | 19% | TPI (950 m) |
| | 69% | Relative elev. - circ. (1600 m) | | 19% | Vertical distance to channel |
| | 69% | TPI (1100 m) | | | |
| | 62% | TPI (550 m) | | | |
| | 62% | Northness (215 m) | | | |
| | 62% | Eastness (345 m) | | | |
| | 62% | Modified catchment area | | | |
| | 32% | Aspect, W central angle (110 m) | | | |
| | 21% | TPI (250 m) | | | |
| | 21% | Aspect, W central angle (175 m) | | | |
| | 12% | Northness (6 m) | | | |
| **SOC$_{stock}$** | | | **SOC$_{stock}$** | | |
| 100% | 48% | Relative elev. - rect. (1100 m) | | 100% | LandsatLook, band 5 (6 Jul. 2006) |
| 48% | 100% | Vertical distance to channel | | 100% | LandsatLook, band 3 (6 Jul. 2006) |
| | 80% | Channel network base level | | 100% | LandsatLook, band 6 (6 Jul. 2006) |
| | | | | 100% | Landsat5 TM, band 7 (11 Jun. 2006) |

1100

1101

1102    Table 4. Fitting performance for the respective models. The model's efficiency (ME) is the ratio
1103    between the model's mean absolute error (MAE) and the MAE that would result from only using the
1104    mean value as the model. Cubist reports the ME as relative error, but it is renamed here to avoid
1105    confusion with the more common definition of relative error. An ME of greater than one indicates
1106    that the model is not performing well.

| **Topsoil models** | BD | SK | H | SOC$_{\%}$ | Indirect - SOC$_{stock}$ | Direct - SOC$_{stock}$ |
|---|---|---|---|---|---|---|
| MAE | 0.05 | 1.36 | 5.90 | 0.14 | 1.69 | 2.27 |
| ME | 0.52 | 0.41 | 0.47 | 0.34 | 0.49 | 0.66 |
| $R^2$ | 0.69 | 0.85 | 0.71 | 0.86 | 0.73 | 0.58 |
| **Subsoil models** | | | | | | |
| MAE | 0.06 | 3.77 | 5.90 | 0.06 | 2.75 | 1.37 |
| ME | 0.58 | 0.42 | 0.47 | 0.59 | 1.67 | 0.83 |
| $R^2$ | 0.67 | 0.79 | 0.71 | 0.55 | 0.34 | 0.19 |

1107

1108 Table 5. Cross-validation performance for the respective models. Note that although the $R^2$ was
1109 severely reduced for most models, the MAE was generally only increased a small amount.

| Topsoil models | BD | SK | H | SOC$_\%$ | Direct - SOC$_{stock}$ |
|---|---|---|---|---|---|
| MAE | 0.08 | 2.70 | 11.80 | 0.27 | 2.94 |
| ME | 0.86 | 0.82 | 0.93 | 0.66 | 0.85 |
| $R^2$ | 0.26 | 0.08 | 0.12 | 0.61 | 0.27 |
| **Subsoil models** | | | | | |
| MAE | 0.09 | 7.18 | 11.80 | 0.09 | 1.42 |
| ME | 0.80 | 0.80 | 0.93 | 0.98 | 0.86 |
| $R^2$ | 0.36 | 0.26 | 0.12 | 0.05 | 0.17 |

1110

1111 Table 6. Skewness coefficients for the residuals of each model.

| | BD | SK | H | SOC$_\%$ | Indirect - SOC$_{stock}$ | Direct - SOC$_{stock}$ |
|---|---|---|---|---|---|---|
| Topsoil models | -0.25 | -1.15 | 0.17 | 1.04 | 0.10 | 0.37 |
| Subsoil models | 0.11 | -0.74 | -0.17 | 1.18 | -1.61 | -0.16 |

1112