

# The relationship between diffuse spectral reflectance of the soil and its cation exchange capacity is scale-dependent

A. Savvides<sup>a</sup>, R. Corstanje<sup>b,d</sup>, S.J. Baxter<sup>a</sup>, B.G. Rawlins<sup>c</sup>, R.M. Lark<sup>b</sup>

<sup>a</sup> *Department of Soil Science, The University of Reading, Whiteknights, Reading, RG6 6DW, UK*

<sup>b</sup> *Rothamsted Research, Harpenden, Hertfordshire AL5 2JQ, UK*

<sup>c</sup> *British Geological Survey, Keyworth, Nottingham NG12 5GG, UK*

<sup>d</sup> *Now at National Soil Resources Institute, Cranfield University, Cranfield, Bedford, MK43 0AL, UK*

---

## 1 Abstract

2 Diffuse reflectance spectroscopy (DRS) is increasingly being used to predict nu-  
3 merous soil physical, chemical and biochemical properties. However, soil properties and  
4 processes vary at different scales and, as a result, relationships between soil properties  
5 often depend on scale. In this paper we report on how the relationship between one  
6 such property (CEC) and the DRS of the soil depends on spatial scale. We show this  
7 by means of a nested analysis of covariance of soils sampled on a balanced nested design  
8 in a 16 km × 16 km area in eastern England. We used principal components analysis  
9 on the DRS to obtain a reduced number of variables while retaining key variation. The  
10 first principal component accounted for 99.8 % of the total variance, the second for  
11 0.14 %. Nested analysis of the variation in the CEC and the two principal components  
12 showed that the substantial variance components are at the > 2000-m scale. This is  
13 probably the result of differences in soil composition due to parent material.

14 We then developed a model to predict CEC from the DRS and used partial least  
15 squares (PLS) regression to do so. Leave-one-out cross-validation results suggested a  
16 reasonable predictive capability ( $R^2 = 0.71$  and  $RMSE = 0.048 \text{ mol}_c\text{kg}^{-1}$ ). However,

17 the results from the independent validation were not as good, with  $R^2= 0.27$ , RMSE  
18  $= 0.056 \text{ mol}_c\text{kg}^{-1}$  and an overall correlation of 0.52. This would indicate that DRS  
19 may not be useful for predictions of CEC. When we applied the analysis of covariance  
20 between predicted and observed we found significant scale-dependent correlations at  
21 scales of 50 and 500 m (0.82 and 0.73 respectively). DRS measurements can therefore  
22 be useful to predict CEC if predictions are required, for example, at the field scale (50  
23 m). This study illustrates that the relationship between DRS and soil properties is scale  
24 dependent and that this scale dependency has important consequences for prediction  
25 of soil properties from DRS data.

26 *Keywords:* Pedometrics; Nested Sampling; Diffuse Reflectance Spectra.

27

---

28 \*Corresponding author: *E-mail address:* roncorstanje@cranfield.ac.uk (R. Corstanje).

## 29 **1. Introduction**

30 There is a large demand for data on soil for quality assessment, environmen-  
31 tal monitoring and precision agriculture. Diffuse reflectance spectroscopy, in the near  
32 infrared (NIR, 750 – 2500 nm) or visible and near infrared (VNIR, 350 – 2500 nm) fre-  
33 quency bands, has been proposed as a rapid, cost-effective and non-destructive method  
34 to obtain predictions of soil properties which would be too expensive to measure di-  
35 rectly on many samples (e.g. Viscarra Rossel et al., 2006; Zornoza et al., 2008). The  
36 method is based on the premise that the variation of the diffuse spectra of soil at these  
37 wavelengths is due to variation in the composition of the soil (Cohen et al., 2005).  
38 Predictive relationships between the spectra and target soil properties are obtained by  
39 chemometric methods such as partial least-squares regression (PLSR; Viscarra Rossel  
40 et al., 2006). Diffuse reflectance spectroscopy has been used to predict various proper-  
41 ties of the soil including pH, cation exchange capacity, organic matter or organic carbon

42 content, composition of the microbial community and extracellular enzyme activities  
43 (see Viscarra Rossel et al., 2006 and Zornoza et al., 2008).

44 Variation of soil properties depends on factors such as parent material, climate,  
45 land use and topography. These factors all operate at different scales and will therefore  
46 influence soil processes and soil variation at different scales. As a consequence the  
47 relationship between soil variables might also be scale dependent. This has been shown  
48 in previous studies on heavy metals in soil (Goovaerts and Webster, 1994), organic  
49 carbon and urease activity (Corstanje et al., 2007), and the effects of various soil  
50 properties on the rates of emission of trace gases (Lark et al., 2004). One result of  
51 such scale-dependence is that the overall correlation between two soil properties might  
52 mask underlying relationships at different scales.

53 We are not aware of any previous studies on how the relationship between diffuse  
54 reflectance spectra and soil properties depends on spatial scale. Published studies  
55 report that DRS is effective for predicting some soil properties such as soil organic C,  
56 but is less effective at predicting others such as CEC or nutrient content such as N  
57 or P (Viscarra Rossel et al., 2006). However, in these studies the measurements are  
58 made on specimens collected on some support (e.g. a core) distributed across a field or  
59 landscape according to some sampling scheme. The covariation of DRS measurements  
60 and soil properties therefore contains unresolved contributions from a range of spatial  
61 scales. Diffuse reflectance spectra are surrogates for soil properties that determine,  
62 more or less directly, the nature of the interaction between soil and electromagnetic  
63 radiation. It is therefore possible that the variations of the spectra at some spatial  
64 scales are dominated by variation of one soil property, and other soil properties cause  
65 most of the spectral variation at other scales. A weak overall correlation of the spectra  
66 with measurements of the soil property might therefore mask a strong relationship at  
67 some particular scale, and if this scale coincides with the scale at which information on  
68 the soil property is actually needed (e.g. field averages) then the DRS measurements

69 may be of considerably more practical value than is indicated by the simple statistics  
70 on the basic sample support.

71 We therefore need to study the scale-dependence of the relationship between  
72 DRS measurements and soil properties. Spatially nested sampling is an efficient way  
73 to study scale-dependent variation over large areas and very disparate scales (Youden  
74 & Melich, 1937; Webster and Oliver, 1990). The variance of soil properties sampled  
75 this way can be partitioned into scale-specific components by a nested analysis of  
76 variance, and the covariances of properties can similarly be analysed by scale (Lark,  
77 2005). In a recent study (Corstanje et al., 2008) we investigated the covariance of  
78 soil properties such as pH, CEC and bulk density with rates of ammonia volatilization  
79 from soils collected by a nested scheme across a variable region in the eastern midlands  
80 of England. The resulting collection of soil offered the opportunity to investigate the  
81 relationship between the DRS of the soils and a basic soil property at disparate scales.  
82 That analysis is the subject of this paper.

83 We chose cation exchange capacity (CEC) as the target soil property for this  
84 study. The CEC of the soil is a basic physico-chemical property. It is laborious to  
85 measure, since different cations must be extracted and determined. Nonetheless, it is  
86 important if we are to predict the behaviour of the soil since it affects, among other  
87 things, the behaviour of various pollutants in soil (Wang and Keller, 2008), the ability  
88 of the soil to retain and supply important plant nutrients (Bailey et al., 2008) and the  
89 rates of important steps in biogeochemical cycles including the emission of trace gases  
90 from the soil (Jarecki et al., 2008). It would therefore be useful if DRS could be used  
91 to predict CEC, and this has been attempted previously. Viscarra Rossel et al. (2006)  
92 report several studies in which CEC was predicted with reasonable success (coefficients  
93 of determination between 0.64 and 0.88).

94 In this paper we report analyses to investigate the scale-dependent relationship  
95 between DRS and CEC across our study area, and an evaluation by scale of the efficacy

96 of DRS-based predictions of CEC.

## 97 **2. Materials and Methods**

### 98 *2.1 The study region and the sampling scheme*

99 A detailed account of the study region and the sampling is given by Corstanje et  
100 al. (2008). The region is approximately 16 km  $\times$  16 km and lies between the towns  
101 of Luton (south) and Bedford (north) in the eastern midlands of England. Most of  
102 this region is over Cretaceous formations: Chalk, Gault Clay and Lower Greensand,  
103 but there is also Oxford Clay (Jurassic) in the north. Superficial material, including  
104 chalky boulder clay, and other glacial drift of variable texture cover the country rock  
105 over much of the region.

106 We used a balanced nested sampling design, in which  $n_1$  sampling main stations  
107 are chosen on a grid or transect of interval  $d_1$ . Two substations (level 2) were then  
108 chosen about each main station, separated from each other by fixed distance,  $d_2$ , on  
109 a line on a random bearing. We repeated this procedure until, about each substation  
110 at level  $m - 1$ , two sample points (level  $m$ ) separated from each other by distance  
111  $d_m$  were selected. A nested analysis of the variances and covariances of variables  
112 measured on the sample points is possible, and components associated with the spatial  
113 scales determined by the distances,  $d_1, d_2, \dots, d_m$  can be estimated. As described by  
114 Corstanje et al. (2008), our main stations were on nodes of a 2-km grid, chosen so that  
115 the associated (co)variance components would be dominated by differences between  
116 the major parent materials. The substations were separated by 500 m, 50 m and 2  
117 m. We selected 36 main stations, each with eight sample points on the nested scheme  
118 giving 288 sample points in total.

### 119 *2.2 Soil preparation and analysis*

120 This sampling exercise was done as part of a study on ammonia volatilization  
121 from soil, and this is reflected in the sample treatment. The soils were air-dried, large

122 plant fragments were removed, sieved to pass 0.5 mm, and then 1-kg portions were  
123 washed in 1.5 dm<sup>3</sup> of 10 mM CaCl<sub>2</sub> to remove nitrate and replace exchangeable cations  
124 with Ca<sup>2+</sup>. The soils were then air-dried again and re-sieved to 0.5 mm.

### 125 *2.3 Measurement of Cation Exchange Capacity*

126 The cation exchange capacity of each sample was determined as described by  
127 Rowell (1994). The exchangeable calcium, magnesium and potassium ions were ex-  
128 tracted from a weighed subsample of the air-dried soil into 1 M ammonium ethanoate  
129 buffer (pH 7). This was then displaced with ethanol and then flame photometry was  
130 used to measure the concentrations of the three ions. Ammonium was then extracted  
131 from the soil with acidified 1 M KCl and measured by steam distillation and titration.  
132 The CEC was then expressed as mol<sub>c</sub> kg<sup>-1</sup> air-dry soil.

### 133 *2.3 Measurement of diffuse reflectance spectra*

134 Soil samples were scanned in the visible–near infrared region (350–2500 nm) using  
135 an ASD (Analytical Spectral Devices, Boulder, CO) Agri-Spec NIR Spectrometer.  
136 A 20-g subsample of each soil sample was placed in a holder with a quartz window  
137 for scanning. Soils were illuminated and scanned from below using the spectrometer  
138 connected to an ASD muglight with an internal tungsten-quartz-halogen light source  
139 and a 12 mm spot size. Data were collected every 1 nm and every spectrum was an  
140 average of 25 readings. Each sample was scanned twice; the second scan was made  
141 after rotating the sample in its holder through 90° whilst placed on the muglight.  
142 During scanning, a Spectralon 99% reflectance panel was used to optimize and white-  
143 reference the spectrometer after scanning every set of ten samples. Before further  
144 statistical analysis, we obtained an average of two spectra for each sample, truncated  
145 by removing the values below 450 nm and above 2450 nm.

### 146 *2.3 Statistical analysis*

147 We used nested analysis of covariance to study the correlation of the DRS measurements

148 and soil CEC at different scales, and to assess predictions of CEC by partial least-  
 149 squares regression on the DRS at a set of validation sites. In the following section we  
 150 describe the general nested analysis. We then describe how this was used to compare  
 151 the DRS measurements and CEC data and then to assess predictions of CEC.

152 *2.3.1 Nested analysis* Nested analysis of covariance is described by Lark (2005) and we  
 153 give here only a summary for the balanced case. The randomization of directions in  
 154 the nested sampling scheme allows us to treat the values of two soil properties,  $u$  and  
 155  $v$ , as random variables  $Z^u$  and  $Z^v$ , which comprise the following components;

$$\begin{aligned} Z_{ij\dots m}^u &= \mu_u + A_i^u + B_{ij}^u + \dots + \varepsilon_{ij\dots m}^u \\ Z_{ij\dots m}^v &= \mu_v + A_i^v + B_{ij}^v + \dots + \varepsilon_{ij\dots m}^v. \end{aligned} \quad (1)$$

156 The values  $\mu_u, \mu_v$  are the overall means of  $u$  and  $v$ , respectively. The random variables  
 157  $A_i^u, A_i^v$  are, respectively, the differences between the corresponding overall means,  $\mu_u$   
 158 and  $\mu_v$ , and the corresponding means of the  $i$ th main station. Similarly  $B_{ij}^u, B_{ij}^v$  are  
 159 the differences, within the  $i$ th main station, between the mean values of the  $i$ th main  
 160 station and  $j$ th substation. The variables  $A_i^u, B_{ij}^u, \dots$  and  $A_i^v, B_{ij}^v, \dots$  have zero mean,  
 161 and the variables associated with each scale in the nested scheme (e.g.  $A_i^u$  and  $B_{ij}^u$ )  
 162 have covariance matrices  $\mathbf{C}_i, \mathbf{C}_j, \dots$ . The objective of multivariate nested analysis is to  
 163 estimate these covariance matrix components, which are additive components of  $\mathbf{C}$ , the  
 164 overall covariance matrix of the two random variables, since they are associated with the  
 165 scales of interest in the sampling scheme. Because estimates of the covariance matrix  
 166 components by method-of-moments are not guaranteed to be non-negative definite, and  
 167 therefore admissible as covariance matrices for real random variables, Lark (2005) used  
 168 a residual maximum likelihood (REML) algorithm due to Calvin & Dykstra (1992).  
 169 We used this method in the present study.

170 The estimated covariance matrix components were then converted to correlation  
 171 matrices by dividing each element,  $\mathbf{C}_{k,l}$  by the square-root of the product of the cor-  
 172 responding elements on the main diagonal,  $\mathbf{C}_{k,k}$  and  $\mathbf{C}_{l,l}$ . We obtained confidence

173 intervals for the scale-dependent correlations with Fisher's  $z$ -transform, following Lark  
174 (2005).

### 175 *2.3.2 Spectral reduction and correlation with CEC*

176 We used principal components analysis on the spectra to obtain a smaller number of  
177 variables that represent the key variations in spectral variation among our soils. This  
178 was done using GenStat (Payne et al., 2008) to analyse the correlation matrix of the  
179 spectral reflectance in the 2001 channels. Principal component analysis finds  $p$  linear  
180 combinations of a set of  $k$  variables that are uncorrelated (see, for example, Webster and  
181 Oliver, 1990). The first component has the largest variance of any possible such linear  
182 combination. The second component has the largest variance of any linear combination  
183 that is orthogonal to the first and so on. The sum of the variances of the principal  
184 components is equal to the sum of the variances of the original variables, but if there are  
185 correlations between the latter then a large proportion of the total variance of the full  
186 data set is represented by substantially fewer than  $p$  of the principal components. In  
187 fact in our case the first principal component accounted for 99.8% of the total variance,  
188 and the second for 0.14%. This shows that the spectra are very redundant. We used  
189 the nested analysis of covariance to investigate the scale-dependent correlation between  
190 these two principal components of the spectra and soil CEC.

### 191 *2.3.3 Prediction and validation*

192 The fact that more than 99% of the variance of the observations in 2001 channels  
193 can be accounted for by the first principal component indicates that there is a good deal  
194 of redundancy in the spectra, that is to say different channels are so strongly correlated  
195 that they present little independent information. This is a common situation in the  
196 analysis of spectra, and partial least squares (PLS) methods are widely used to obtain  
197 predictive regressions of variables of interest (such as soil properties) on such very  
198 redundant predictor variates.

199 In this study we used PLS to obtain predictive relationships between the diffuse



200 reflectance spectra and soil CEC, using a subset of the data for estimation of the  
 201 regression model, and the remainder to test the predictions. We used the PLS regression  
 202 (PLSR) algorithm in the ParLes package (Viscarra Rossel, 2008). In PLSR we have  $n$   
 203 observations of  $k$  soil variables (predictands) in the  $n \times k$  matrix  $\mathbf{Y}$  and  $n$  values of  
 204 a  $p$ -variate predictor (e.g. the DRS) in the  $n \times p$  matrix  $\mathbf{X}$ . In PLSR these variables  
 205 are decomposed into common orthogonal factors (similar to the principal components  
 206 discussed above) from which the original variates can be reconstituted by means of  
 207 loading matrices for the predictands and predictors. The algorithm finds an orthogonal  
 208 decomposition such that the first few factors account for as much variation in the  
 209 predictands and predictors as possible. The decomposition can be expressed by the  
 210 following equation

$$\begin{aligned}\mathbf{X} &= \mathbf{TP}' + \mathbf{E} \\ \mathbf{Y} &= \mathbf{TQ}' + \mathbf{F},\end{aligned}\tag{2}$$

211 where  $\mathbf{T}$  is an  $n \times l$  matrix of factor scores, and  $\mathbf{P}$  and  $\mathbf{Q}$  are, respectively,  $p \times$   
 212  $l$  and  $k \times l$  matrices of loadings. The number  $l$  is the number of factors that are  
 213 assumed to be informative, and is selected according to a criterion such as the Akaike  
 214 Information Criterion, see Viscarra Rossel (2008) for details. The matrices  $\mathbf{E}$  and  $\mathbf{F}$   
 215 contain residuals, i.e. the contributions of the excluded factors. The number of  $l$  used  
 216 in the PLSR model was determined through leave-one-out cross-validation. We selected  
 217  $l=4$  on the basis of the RMSE and AIC criterion.

218 In this study we randomly divided the main-stations of our nested sampling  
 219 scheme into a prediction set of 29 (232 observations) for estimation of the regression  
 220 model to predict CEC and a validation set of 7 (56 observations) to test the predic-  
 221 tions. We used the PLSR algorithm in the ParLes software to fit the predictive model.  
 222 In Table 1 we report the results from the cross-validation on the prediction set and  
 223 prediction at the separate validation sites.

224 Since the validation observations have a nested structure we were able to analyse

225 the covariation of the observed and predicted values of CEC by nested analysis of  
226 covariance. Correlations at each scale and their confidence intervals were computed.

### 227 **3. Results**

228 Summary statistics of the soil cation-exchange capacities are shown in Table 1.  
229 In Figure 1 we present the computed scale dependent variance components for CEC  
230 and the first two principal components of the DRS; PC1 and PC2. These show that  
231 variation in both the DRS and CEC is scale-dependent. The largest component of  
232 variance in CEC, which comprises 66% of the total, was at  $> 2000$  m. Variation at this  
233 scale will be predominantly due to differences between parent materials. The observed  
234 values of CEC were smaller over Lower Greensand, where the average CEC was  $0.104$   
235  $\text{mol}_c \text{ kg}^{-1}$ . Most soils over the Lower Greensand are sandy loams or loamy sands, so  
236 the CEC is relatively small because the soil contains relatively little clay. In contrast  
237 the average CEC over the Gault clay, Chalk and chalky till were larger:  $0.219$ ,  $0.161$   
238 and  $0.188 \text{ mol}_c \text{ kg}^{-1}$  respectively. The second-largest component of variance was at  
239 the 500-m scale (15 % of the total), followed by the 50-m scale (12 %) and  $<2$ -m scale  
240 (6 %). In general, then, the variation of CEC appears to be dominated by parent  
241 material differences. The smaller variance components at the finer scales will be due  
242 to variation in factors such as organic matter content of the soil.

243 The largest components of variance of both the first two principal components of  
244 the DRS were at the coarsest scale of  $> 2000$  m — 60 and 50 % of the total for PC1  
245 and PC2 respectively. As with CEC, this suggests that the variation of the DRS in this  
246 data set is dominated by differences between the parent materials. The second-largest  
247 component of variance of both principal components was at the 500-m scale, with 22  
248 % for both PC1 and PC2. Variation at this scale will be due to differences in land use  
249 and management practices as well as some short-range variation in parent materials  
250 such as superficial deposits. The components of variance for PC1 at the finer scales of  
251 50 m and 2 m were small (6 and 8 %, respectively) for PC1, but slightly larger (12 and

252 16 %, respectively) for PC2.

253 The results from the nested analysis of covariance of CEC and each of the first  
254 two principal components of DRS are presented in Figure 2. The correlations of CEC  
255 and PC1 are weak, and the correlation at 500 m only (0.39) is significantly different  
256 from zero. The overall correlation between these two variables was weak and positive  
257 (0.19). In the case of PC2, we found significant, strong and positive correlations at  
258 all scales except the finest scale of 2-m. The observed scale dependent correlations  
259 were 0.60, 0.65 and 0.86 for scales 50-, 500- and >2000-m respectively and the overall  
260 correlation was 0.47.

261 The results of the PLS model fitting and validation are summarized in Table 2.  
262 Cross-validation of the fitted PLS model suggested that the predictions are reasonable  
263 ( $R^2 = 0.71$  and  $RMSE = 0.048 \text{ mol}_c\text{kg}^{-1}$ ). However, the tests on the separate vali-  
264 dation data are less encouraging, with  $R^2 = 0.27$  and  $RMSE = 0.056 \text{ mol}_c\text{kg}^{-1}$ . If we  
265 consider the scale-dependent correlations (Figure 3), then these show that the overall  
266 correlation of the predicted and measured values of CEC in our validation set (0.52)  
267 masks stronger correlations (0.82 and 0.73 respectively) at the scales 50 and 500 m,  
268 while the correlation at 2m is zero. A correlation at the coarsest scale is not reported  
269 because the estimated covariance matrix was positive semi-definite, but not positive  
270 definite, so the estimated correlation is 1.0, see Lark (2005).

#### 271 4. Discussion and Conclusions

272 We have seen that both CEC and the principal components of the DRS show  
273 scale-dependent variation with the variance components increasing with distance. This  
274 indicates that the variation of the DRS is dominated by aspects of the composition of  
275 the soil associated with parent material. By contrast, for example, Corstanje et al.  
276 (2008) found that about 20% of the variance in urease activity in this soil occurred at  
277 the 2-m scale. In this environment, broad-scale variations in parent material have a  
278 larger impact on the DRS than factors operating at finer scales, such as geomorpho-

279 logical or biological controls on soil composition. This was also observed for topsoil  
280 geochemistry in eastern England (Rawlins et al., 2003). This is likely to explain the  
281 scale-dependent relationship between CEC and the spectra. The soil spectrum responds  
282 to components of the soil which are themselves correlated with its CEC.

283 In this landscape, for example, the large iron oxide content of the Lower Green-  
284 sand is likely to be spectrally distinctive, and we have noted that soils on the Lower  
285 Greensand are generally lighter-textured with smaller CEC than those located over  
286 different parent materials. It is therefore likely that the better predictive relationships  
287 between CEC and the DRS will be seen at scales where these surrogate relationships  
288 are expressed, while the correlations at other scales are very poor. This is the behaviour  
289 that our nested sampling and analysis reveals in this case study.

290 The important point that these results illustrate is that a poor overall correla-  
291 tion between DRS and a target soil property, or poor overall validation statistics for  
292 predictions, do not necessarily indicate that the spectra are not suitable for predictive  
293 purposes. For example, although the overall correlation of the predicted and measured  
294 spectra in our validation set was only 0.52, the correlations of the components at the  
295 50- and 500-m scales were much stronger. The very weak correlation at the 2-m scale  
296 masks the relationship at coarser scales.

297 In practice we might often be interested in predicting a soil property only at  
298 coarser scales. For example, if we want to estimate the mean CEC for each of a set of  
299 fields, then the relatively good correlation of DRS and CEC at the 50-m scale suggests  
300 that the spectral measurements might be useful, and the average predicted CEC for a  
301 set of soil specimens collected within a field should give a reasonable prediction of the  
302 true field mean. Similarly, cokriging estimates for blocks with sides 50 m or longer, from  
303 a set of measurements of CEC and a denser set of DRS spectra should be reasonable,  
304 because the variation at fine scales, where the DRS and CEC are weakly correlated, is  
305 averaged out.

306 Some additional issues are raised by this study. First, it provides evidence for the  
307 importance of assessing predictions from spectra on separate validation data sets, and  
308 not giving undue weight to cross-validation assessments. Second, the scale-dependence  
309 in the errors from our PLSR predictions suggests that there is a need to develop  
310 the algorithm to allow for models where the errors are not assumed, as in standard  
311 implementations of PLSR as we used here, to be independent random variables. While  
312 the regression coefficients are still unbiased, they are not necessarily the ones that give  
313 us minimum variance predictions. There may therefore be advantages in extending the  
314 PLSR algorithm to deal with such circumstances.

315 To conclude, the relationship between DRS and soil properties has been shown  
316 to be scale-dependent for one case study. An important consequence of this is that  
317 assessments of the predictive value of statistical models that use DRS to predict soil  
318 properties should account for scale-dependence. If this is not done then weak rela-  
319 tionships between spectral properties and the target soil property at one or more scale  
320 might obscure strong relationships at other scales, which might well be scales at which  
321 the soil information is needed. The nested sampling and analysis scheme used in this  
322 study is one way to identify such behaviour.

### 323 **Acknowledgements.**

324 The field work described here, along with RC's contribution, was undertaken as  
325 part of a project funded by the United Kingdom's Biotechnology and Biological Sci-  
326 ences Research Council (BBSRC), grant no. BB/C506813/1. RML's contribution was  
327 part of the programme of the Centre for Mathematical and Computational Biology at  
328 Rothamsted Research, funded by the BBSRC. AS acknowledges the technical assis-  
329 tance of Karen Gutteridge and Chris Speed at the University of Reading. This paper  
330 is published with the permission of the Executive Director of the British Geological  
331 Survey (Natural Environment Research Council).

## References

- Bailey, J.S., Ramakrishna, A. and Kirchof, G., 2008. Relationships between important soil variables in moderately acidic soils ( $\text{pH} \leq 5.5$ ) in the highlands of Papua New Guinea and management implications for subsistence farmers. *Soil Use Manage.*, 24: 281–291.
- Calvin, J.A. and Dykstra, R.L., 1992. An algorithm for restricted maximum likelihood estimation in balanced multivariate variance components models. *J. Stat. Comput. Simul.*, 40: 233–246.
- Cohen, M.J., Prenger, J.P. and DeBusk, W.F., 2005. Visible-near infrared reflectance spectroscopy for rapid, non-destructive assessment of wetland soil quality. *J. Environ. Qual.*, 34: 1422–1434.
- Corstanje, R., Schulin R. and Lark, R.M., 2007. Scale-dependent relationships between urease activity and soil organic carbon. *Eur. J. Soil Sci.*, 58: 1087–1095.
- Corstanje, R., Kirk, G.J.D., Pawlett, R. M., Read, R. and Lark, R.M., 2008. Spatial variation of ammonia volatilization from soil and its scale dependent correlation with soil properties. *Eur. J. Soil Sci.*, 59: 1260–1270.
- Goovaerts, P. and Webster, R., 1994. Scale-dependent correlation between topsoil copper and cobalt concentrations in Scotland. *Eur. J. Soil Sci.*, 45: 79–95.
- Jarecki, M.K., Parkin, T.B., Chan, A.S.K., Hatfield, J.L. and Jones, R., 2008. Greenhouse gas emissions from two soils receiving nitrogen fertilizer and swine manure slurry. *J. Environ. Qual.*, 37: 1432–1438.
- Lark, R.M., 2005. Exploring scale-dependent correlation of soil properties by nested sampling. *Eur. J. Soil Sci.*, 56: 307–317.

- Lark, R. M., Milne, A. E., Addiscott, T. M., Goulding, K. W. T., Webster, C. P. and OFlaherty, S., 2004. Scale- and location-dependent correlation of nitrous oxide emissions with soil properties: an analysis using wavelets. *Eur. J. Soil Sci.*, 55: 611–627.
- Payne, R.W. (Editor). GenStat<sup>®</sup> for Windows<sup>™</sup> 11<sup>th</sup> Edition Introduction.
- Rawlins, B. G., Webster, R. and Lister, T.R., 2003. The influence of parent material on topsoil geochemistry in eastern England. *Earth Surf. Processes Landforms*, 28: 1389–1409.
- Rowell, D.L., 1994. *Soil Science: Methods and Applications*. Longman, London.
- Viscarra Rossel, R.A., 2008. ParLes: Software for chemometric analysis of spectroscopic data. *Chemometrics Intell. Lab. Syst.*, 90: 72–83.
- Viscarra Rossel, R.A., Walvoort, D.J.J., McBratney, A.B., Janik, L.J. and Skjemstad, J.O., 2006. Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil properties. *Geoderma*, 131: 59–75.
- Wang, P. and Keller, A.A., 2008. Soil particle-size dependent partitioning behavior of pesticides within water-soil-cationic surfactant systems. *Water Res.*, 42: 3781–3788.
- Webster, R. and Oliver, M.A., 1990. *Statistical Methods in Soil and Land Resource Survey*. Oxford University Press, Oxford.
- Youden, W.J. and Mehlich, A., 1937. Selection of efficient methods for soil sampling. *Contrib. Boyce Thompson Inst. Plant Res.*, 9: 59-70.
- Zornoza, R., Guerrero, C., Mataix-Solera, J., Scow, K.M., Arcenegui, V. and Mataix-Beneyto, J., 2008. Near infrared spectroscopy for determination of various phys-

ical, chemical and biochemical properties of Mediterranean soils. *Soil Biol. Biochem.*, 40: 1923–1930.



Table 1. Summary statistics - soil cation exchange capacity (n=288).

Statistic	value /mol <sub>c</sub> kg <sup>-1</sup>
10 <sup>th</sup> percentile	0.074
Median	0.18
90 <sup>th</sup> percentile	0.28
Mean	0.18
Variance	0.0067

Table 2. Statistics for prediction performance by the partial least squares model. These are determined by cross-validation on the prediction set, or on the independent validation set.

---

Model	$R^2$	$R^2_{\text{adj}}$	RMSE
Cross-validation	0.71	0.71	0.048
Independent	0.27	0.24	0.056

---

## Figures

Figure 1. Accumulated components of variation for CEC (a) and the first two principal components of the diffuse reflectance spectra (PC1 and PC2; represented by b and c respectively). Components for the largest scale are plotted here against 2000 m but apply to distances  $> 2000$  m.

Figure 2. Scale-dependent correlations between CEC and the first two principal components of the diffuse reflectance spectra (PC1 and PC2; represented by a and b respectively). The correlations are plotted with 95 % confidence intervals (bars). Components for the largest scale are plotted here against 2000 m but apply to distances  $> 2000$  m.

Figure 3. Scale-dependent correlations between predicted and measured CEC from the validation set. The predicted CEC were obtained using a PLS model on the DRS. See text for details. The correlations are plotted with 95 % confidence intervals (bars). Components for the largest scale are plotted here against 2000 m but apply to distances  $> 2000$  m.

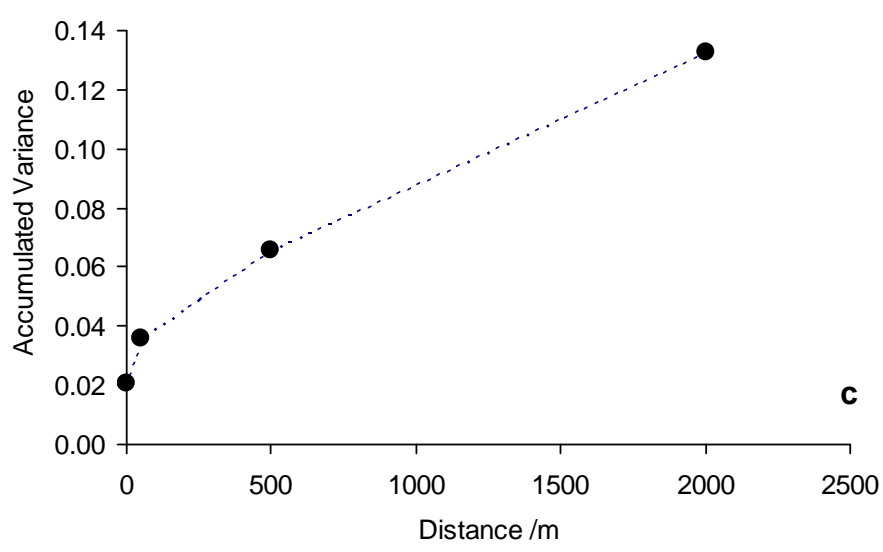
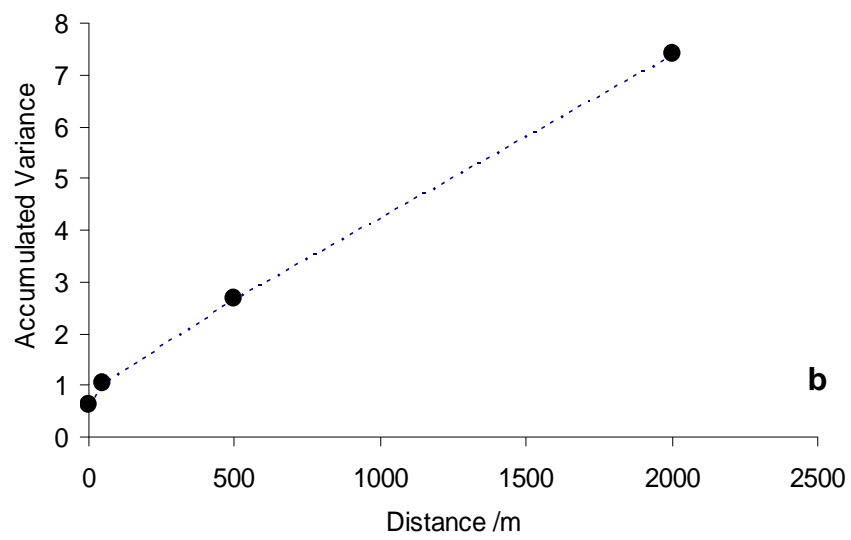
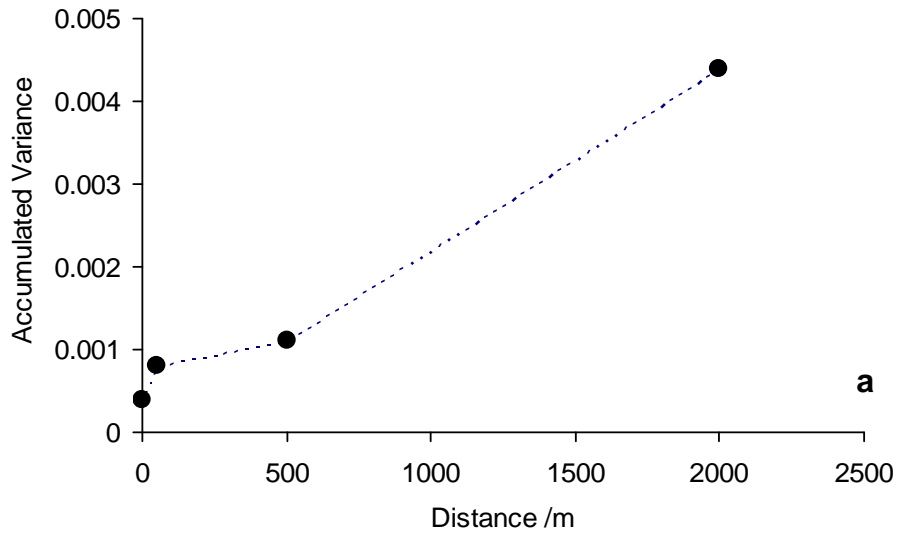
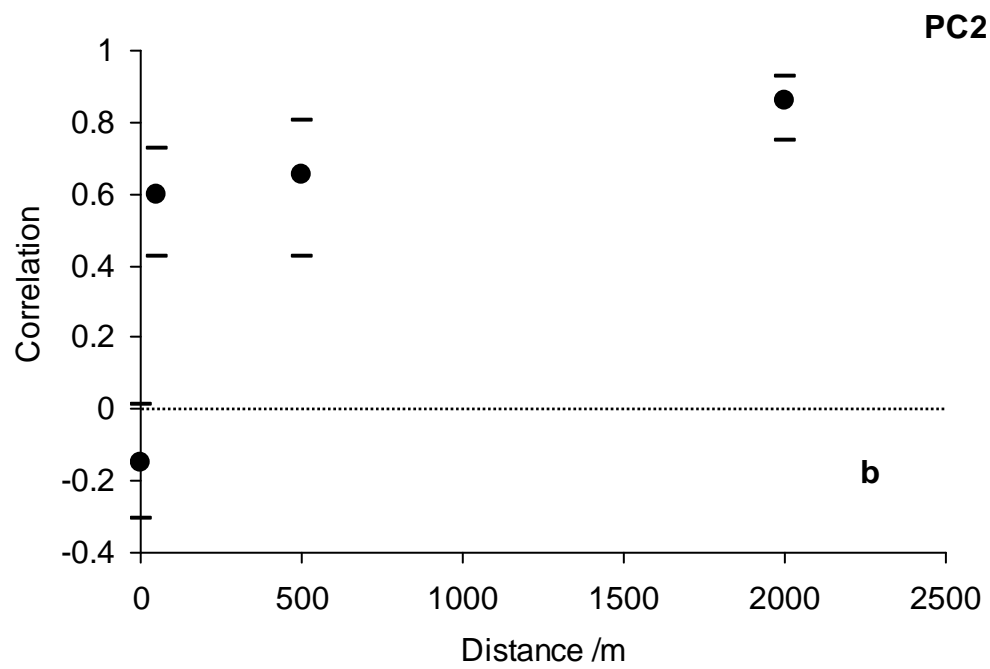
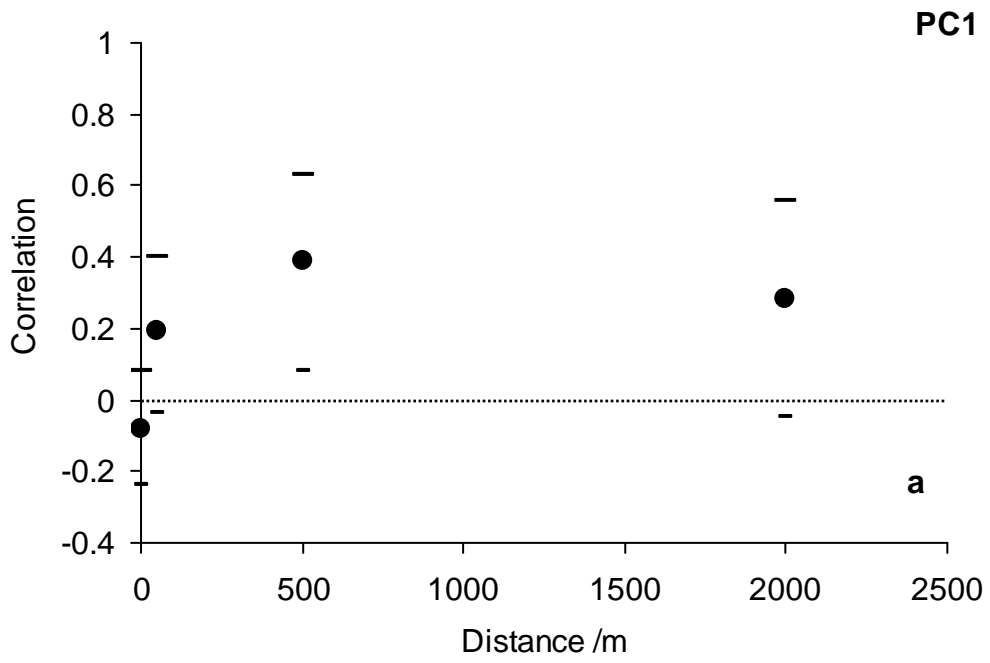


Fig1



**Fig2**

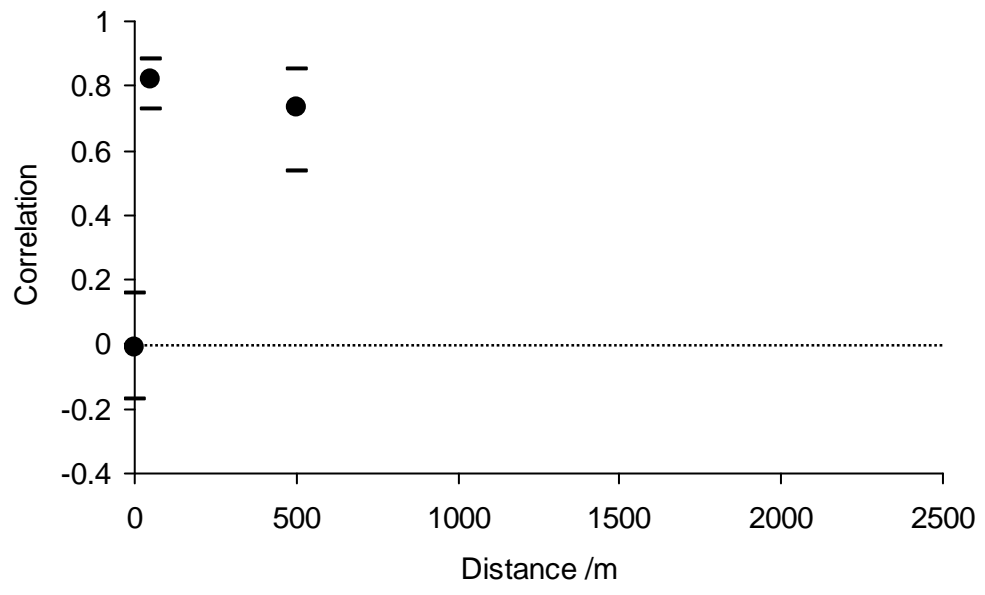


Fig3