

1 **Scope to predict soil properties at within-field scale from small samples using**  
2 **proximally sensed  $\gamma$ -ray spectrometer and EM induction data**

3 J. Huang<sup>a</sup>, J. Triantafilis<sup>a</sup>

4 R.M. Lark<sup>b</sup>,

5 D.A. Robinson,<sup>c</sup> I. Lebron,<sup>c</sup>

6 A.M. Keith<sup>d</sup>

7 B. Rawlins<sup>b</sup>, A. Tye<sup>b</sup>, O. Kuras<sup>b</sup>, M. Raines<sup>b</sup>

8

9 <sup>a</sup>School of Biological, Earth and Environmental Science, The University of New South Wales,  
10 Kensington NSW 2052, Australia

11 <sup>b</sup>British Geological Survey, Keyworth, Nottingham, NG12 5GG, UK

12 <sup>c</sup>NERC- Centre for Ecology and Hydrology, Environment Centre Wales, Deiniol Road, Bangor,  
13 Gwynedd, LL57 2UW, UK

14 <sup>d</sup>NERC- Centre for Ecology and Hydrology, Bailrigg, Lancaster, LA1 4AP.

15

16 Correspondence: E-mail: John Triantafilis [j.triantafilis@unsw.edu.au]

17

18 Key Words: EC<sub>a</sub>, induction, EM, fuzzy k-means, gamma-ray spectrometry,  
19 soil variability, characterization and distribution, digital soil mapping

20

21 Abbreviations: EMI Electromagnetic induction, EC<sub>a</sub> bulk soil electrical conductivity

## 1 **1. Introduction**

2           Spatial predictions of soil properties are needed for various purposes including agriculture and  
3 engineering as well as scientific disciplines such as soil science, ecology and hydrology (Goovaerts,  
4 1997). For example, maps of clay content can be used to ascertain land-use potential, whilst maps of  
5 soil pH can indicate lime requirement to counteract soil acidity, or potential nutrient availability.  
6 However, the costs associated with soil sampling and laboratory analysis are substantial, and spatial  
7 prediction requires considerable sample effort given the observation by Webster and Oliver (1992) that  
8 approximately 100 sample points are required to estimate a spatial statistical model. One way to  
9 improve soil sampling efficiency is to combine direct measurement of soil properties with collection of  
10 cheaper-to-measure ancillary data. Ancillary data can be used to improve precision with which  
11 properties are predicted from relatively few direct observations. Hence the growing interest in proximal  
12 geophysical sensing methods (Robinson et al., 2008) which have been applied to a range of problems  
13 including, soil salinity assessment (Lesch et al., 2005), prediction of depth to clay (Jung et al., 2006),  
14 soil moisture determination (Robinson et al., 2012), determination of soil cation exchange capacity  
15 (Triantafilis et al., 2009a) and deep drainage estimation (Woodforth et al., 2012).

16           In this paper we consider two possible approaches. The first is to use ancillary data to form a set  
17 of land classes by a numerical clustering algorithm. The mean value of the soil property in each class,  
18 estimated from samples within each class, can then be used for prediction. This approach could be  
19 useful because it makes no assumptions about the nature of the relationship between the soil property

1 and the ancillary variables and because precise estimates of class means can be obtained from bulk  
2 samples formed by aggregating individual sample cores within the class thereby reducing analytical  
3 costs. One practical question for the implementation of this approach is how many classes should be  
4 defined. This is usually addressed by considering the distribution of the ancillary variables used to  
5 form the classes, looking for evidence of compact structures in feature space (e.g. Triantafilis et al.,  
6 2009b). The rationale of this approach is that the classes so-identified reflect natural clusters in the  
7 feature space rather than an arbitrary partition, and so should reflect underlying sources of variation in  
8 the soil. Another approach (not used in this context to date) is prediction-based. As we consider more  
9 and smaller classes the within-class variance of the soil properties we wish to predict will, in general  
10 diminish, but the prediction error does not necessarily because the class mean is estimated with less  
11 precision as a fixed sample effort is divided between more classes.

12 A second and more commonly-used approach is linear predictive modeling, essentially a  
13 multiple regression of the target soil property on the ancillary variables. Ideally this is done using data  
14 obtained from a probability sample so residuals can be treated as independent. The model is then used  
15 to form a prediction of the target property at a site where only ancillary data is known. Often data are  
16 not collected according to a probability design, in which case a linear mixed model (LMM) fitted in  
17 which covariates are fixed effects but the residuals are treated as a combination of a spatially correlated  
18 random effect and an independent and identically distributed error (Lark et al., 2006). The prediction of  
19 the soil property at an unsampled site is then a combination of a regression-type prediction from

1 proximally-sensed covariates and a kriging-type prediction of residuals from the fixed effects model at  
2 sampled sites (e.g. Gooley et al., 2013).

3 In this paper we consider both approaches, showing how the question ‘how many classes?’ can  
4 be addressed in terms of the uncertainty of resulting predictions, and compared with the linear mixed  
5 model. We illustrate this with a case study in which  $\gamma$ -ray spectrometry and the apparent electrical  
6 conductivity using an electromagnetic (EM) induction instrument were measured as ancillary data  
7 across two fields located east of the village of Shelford near Nottingham in the UK. We formed classes  
8 from the ancillary data using fuzzy k-means (FKM) analysis. We then analyse data on soil properties  
9 along with the classes formed from the ancillary data and the ancillary data themselves. We show how  
10 the precision of class means as predictors of soil properties (for fixed total sample effort) varies with  
11 the number of classes and compare this criterion for the number of classes with measures based on the  
12 distribution of the ancillary data. We also compare these measures of precision with comparable ones  
13 for direct prediction from the ancillary data by a linear model.

14

15 **Figures 1 and 2 (Near here)**

16

17

1    **2.     Materials and methods**

2    *2.1 Study area*

3           The study fields (Figure 1) are located east of the village of Shelford, which lies approximately  
4    4 km east of Nottingham in the valley of the River Trent. Geologically, the area consists of recent  
5    Holocene alluvial deposits and Pleistocene river terraces formed ~24,000 years BP. These lie  
6    unconformably over a series of Triassic mudstones and sandstones (Mercia Mudstone) to form a  
7    shallow floodplain aquifer. The annual average precipitation is 615 mm and evaporation 517 mm. The  
8    eastern field is used for pasture, whereas the field to the west for no-till arable farming (rape/wheat).

9           Recently a detailed (1:10,000) free soil survey was undertaken (Palmer, 2007). Figure 1 shows  
10   the soil series map. These are all simple map units that correspond to fifteen series of the classification  
11   of the Soil Survey of England and Wales (Avery, 1980). Four are in Trent alluvium (i.e. Wharfe, Trent,  
12   Compton and Stixwould); five on Trent River terraces (i.e. Newport, Arrow, Reaseheath, Quorndon  
13   and Wigton Moor series); and, six on the Triassic Mercia Mudstone (e.g. Worcester, Whimple,  
14   Brockhurst, Melbourne, Salwick and Clifton).

15           In brief, the northern end of the arable field is located on Mercia Mudstone where Clifton (Cu)  
16   and Salwick (So) are found. At the southern end Worcester (Wf) is common; however, slightly lighter-  
17   textured soil occurs in the middle and in a thin strip along the northern part where the Wimple (wM)  
18   map unit is located on the same parent material (i.e. mudstone). The pasture field located a short  
19   distance (~1.3 km) to the west is not in the mapped area, but from the topographic relations we might

1 reasonably expect the southwest side is associated with the deep permeable loams of the Arrow (aO)  
2 series of the Trent River terrace, whilst to the northeast the permeable alluvial medium loams of the  
3 Wharfe (Wv) series, associated with the Trent alluvium, dominate.

4

## 5 *2.2 Ancillary instruments and data collection*

6 Two sources of ancillary data were acquired with proximal sensors. To characterize the topsoil  
7 we collected  $\gamma$ -ray spectrometry data. An Exploranium GR-320 portable  $\gamma$ -ray spectrometer mounted in  
8 a backpack was used in conjunction with a handheld GPX-21 detector containing a  $76 \times 76$ -mm NaI  
9 (TI) scintillation crystal, held at an approximate height of 1 m. At this height the instrument detects  $\gamma$ -  
10 radiation from an area within approximately a 10-metre radius (Atomic Energy Commission, USA,  
11 1972). The detector was energy-stabilised with a small  $^{133}\text{Ba}$  source. The measured environmental  
12 levels of radiation are given in counts per second (cps) for the total counts (TC), percent for potassium  
13 (K %), and parts per million for uranium (U ppm), and thorium (Th ppm). In the study field these were  
14 recorded using integration over a 5 s time window as the operator walked slowly over the ground. Data  
15 were collected using a TDS Ranger palmtop running Pocket GIS software. Positioning data (SBAS  
16 enabled) were obtained from an internal Compact Flash card GPS receiver in the Ranger palmtop.  $\gamma$ -ray  
17 spectrometer transects are shown in Figure 2b.

18 In order to infer subsurface and subsoil variation we also collected EM data with a DUALEM-  
19 1S, because it incorporates single horizontal co-planar (HCP) and perpendicular (PRP) receiver arrays

1 that operate at a low frequency (9 kHz). The transmitter is located at one end with the distance to the  
2 centre of the HCP receiver being 1 m. The depth of  $EC_a$  measurement is approximately 0-1.5 m  
3 (1mHcon). The distance from the transmitter to the PRP receiver is 1.1 m which enables depth of  $EC_a$   
4 of 0-0.5 (1mPcon) (DUALEM-421S Manual, 2008). The DUALEM-1S was connected to an Archer  
5 field computer (Juniper Systems Inc. Logan, UT, USA) and Bluetooth Sirf-III Royaltek BT-GPS  
6 receiver (Royal Tek, Kuei Shan, Tao Yuan, Taiwan). Measurements were integrated using the HGIS  
7 software package (Starpal inc., 2531 wapiti road, Fort Collins, CO 80525, USA).

8         The DUALEM-1S surveys were conducted across the two fields after harvest on Sept 2-3, Oct  
9 14 and Nov 9-10, 2011. The instrument was held 0.2 m above the ground, with the instrument aligned  
10 in parallel with the direction of travel. Measurements were made by traversing the fields across the  
11 prevailing slope and following a predetermined route. In the arable field the transect spacing of ~15m  
12 and ~12 m in the pasture field. DUALEM-1S transects are shown in Figure 2b.  $EC_a$  outliers associated  
13 with the remnants of a buried metallic pipe running north-south in the western field were removed.

14

### 15 *2.3 Soil sampling and laboratory analysis*

16         The soil of both fields was sampled on a square grid with an interval of 25 m (see Figure 2c).  
17 There were 68 sample sites in the pasture and 137 in the arable field. Soil sampling was conducted on  
18 Nov 9 2011 (pasture) and Sept 03 2011 (arable field). Soil samples were collected for the depth interval

1 0–0.15 m with a gouge auger. Each of the 205 samples was stored in an individual water-tight plastic  
2 bag, taken back to the laboratory, weighed and then stored in a cool room.

3 The samples were dried, homogenized and sieved to 2mm. The particle size distribution was  
4 determined from a subsample by first determining the various particle size fractions using a laser  
5 diffractometer. The results were reported in terms of percent sand, silt and clay as defined by the Soil  
6 Survey of England and Wales (Hodgson, 1976). Soil pH was measured using a 1 part soil to 2.5 parts  
7 water dilution. We report results for clay content and pH in terms of differences among the classes  
8 obtained from the FKM analysis, to aid pedological interpretation and consider these soil properties for  
9 evaluating the use of the proximal soil sensor data for prediction of soil mapping units and LMM.

10

#### 11 *2.4 Fuzzy k-means (FKM) analysis*

12 There are an increasing number of papers where ancillary data such as EM measurements of  
13 apparent electrical conductivity are analysed to form classes with a numerical clustering algorithm.  
14 This includes the use of k-means (FKM) to cluster EM with either NDVI (Dang et al., 2011) or  
15 topographic wetness index (Priori et al., 2013) or more commonly FKM analysis of EM with Quickbird  
16 imagery (Guo et al., 2013) or  $\gamma$ -ray spectrometry data (Van Meirvenne et al., 2013). Here we use the  
17 FKM method. It is described in detail by McBratney et al. (1992).



1 In brief, the similarity between an individual  $i$  and a cluster  $c$  is measured to determine how  
 2 much they are alike in multi-variate space (Bezdek, 1981). The best outcome minimizes the objective  
 3 function  $J(\mathbf{M}, \mathbf{C})$ :

$$4 \quad J(\mathbf{M}, \mathbf{C}) = a \sum_{i=1}^n \sum_{c=1}^k m_{ic}^{\phi} d_{ic}^2(x_i, c_c) \quad (1)$$

5 where  $\mathbf{M} = m_{ic}$  is a  $n \times k$  matrix of membership values ( $n$  denoting the number of objects),  $\mathbf{C} = (c_{cv})$  is a  
 6  $k \times p$  matrix of class centers ( $p$  denotes the number of variables),  $c_{cv}$  is the value of the center of class  $c$   
 7 for variable  $v$ ,  $x_i = (x_{i1}, \dots, x_{ip})^T$  is the vector representing individual  $i$ ,  $c_c = (c_{c1}, \dots, c_{cp})^T$  is the vector  
 8 representing the center of class  $c$ , and  $d_{ic}^2(x_i, c_c)$  is the square distance between  $x_i$  and  $c_c$  according to a  
 9 distance measure ( $d_{ic}^2$ ). We chose Euclidean given our local knowledge of the geology of the area  
 10 (Bezdek, 1981).

11 The fuzziness exponent ( $\phi$ ) determines degree of fuzziness. When  $\phi = 1$  this is equivalent to  
 12 the hard partition. As  $\phi$  increases, memberships tend to become uniform. The fuzziness performance  
 13 index (FPI) and the normalized classification entropy (NCE) were then used to identify values for  $\phi$   
 14 and  $k$ . This is because the FPI is a measure of continuity between classes:

$$15 \quad FPI = 1 - \frac{kF - 1}{k - 1} \quad (2)$$

16 where  $F$  is the partition coefficient;

$$17 \quad F = \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^k (m_{ic})^2 \quad (3)$$

1 An FPI value of 1 suggests a very fuzzy classification, whilst a value approaching 0 indicates a  
 2 hard one. The NCE is a measure of disorganization in data partitioning:

$$3 \quad NCE = \frac{H}{\log k} \quad (4)$$

4 where  $H$  is the entropy function;

$$5 \quad H = -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^k m_{ic} \log(m_{ic}) \quad (5)$$

6 Values approaching 0 indicate that the classes are well structured, whilst values approaching 1  
 7 suggest the classes are disorganized. Triantafilis et al. (2013) suggest that values around 0.5 provide a  
 8 balance between continuity and structure. A quantitative can be discerned for  $\phi$  and  $k$  using the  
 9 derivative of  $J(\mathbf{M}, \mathbf{C})$  with respect to  $\phi$  (Bezdek, 1981):

$$10 \quad \frac{dJ(\mathbf{M}, \mathbf{C})}{d\phi} = \sum_{i=1}^n \sum_{c=1}^k m_{ic}^{\phi} \log(m_{ic}) d_{ic}^2 \quad (6)$$

11 To determine the number of  $k$ , the outcome of  $J(\mathbf{M}, \mathbf{C})$  partitioning of the ancillary data into  $k = 2$   
 12 to 10 classes using increments in  $\phi$  of 0.2 and between  $\phi = 1.2$  to 2.4 is considered. A suitable value of  
 13  $\phi$  for a given value of  $k$  is determined when the derivative of  $-J(\mathbf{M}, \mathbf{C})$  with respect to  $\phi$  is largest  
 14 (McBratney and Moore, 1985).

15  
 16  
 17  
 18  
 19

## 1 2.5 Linear mixed model (LMM)

2 When soil data were collected according to a systematic sampling scheme, it would not be  
3 appropriate to fit a linear model by ordinary least squares, since the residuals cannot be treated as  
4 independent random variables. Rather we propose a linear mixed model for the data of the form

$$5 \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta} + \boldsymbol{\varepsilon}, \quad (7)$$

6 where  $\mathbf{y}$  is a  $n \times 1$  vector of values of the target soil variable,  $\mathbf{X}$  is a  $n \times p$  design matrix,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector  
7 of fixed effects coefficients,  $\boldsymbol{\eta}$  is a  $n \times 1$  vector the elements of which are a realization of a spatially  
8 correlated random variable and  $\boldsymbol{\varepsilon}$ , is a  $n \times 1$  vector the elements of which are a realization of an  
9 independent and identically distributed random variable. The elements of the design matrix are the  
10 predictor variables and the fixed effects coefficients correspond to these. For example, if the predictors  
11 are  $p$  classes then element  $\{i,j\}$  of  $\mathbf{X}$  is 1 if the  $I$  th observation corresponds to the  $j$ th class and zero  
12 otherwise. There is therefore exactly one element equal to 1 in each row of the design matrix. In this  
13 case the elements of  $\boldsymbol{\beta}$  are the estimated mean values of the target soil variable in the respective classes.

14 The correlated random variable  $\boldsymbol{\eta}$  is assumed to be normal and has mean zero and variance  
15 parameters which are familiar from the geostatistical literature. These are an overall variance,  $\sigma^2_{\boldsymbol{\eta}}$ , and  
16 a distance parameter for a selected variogram function (e.g. the range of a spherical variogram). The  
17 error variable  $\boldsymbol{\varepsilon}$  also has zero mean and a variance  $\sigma^2_{\boldsymbol{\varepsilon}}$ . In an alternative form of this model the elements  
18 in the design matrix may be, in the first column, a column of ones and in the next  $p-1$  columns the  
19 values of the  $p-1$  proximally sensed auxiliary variable, in which case the vector  $\boldsymbol{\beta}$  contains an intercept

1 and  $p-1$  regression coefficients. We fitted models of the form in Equation (7) for target soil variables  
2 and with the fixed effects either the class of maximum membership for the FKM clustering of the  
3 ancillary variables with  $k = 2-10$  or a subset of the ancillary variables in a regression type model. The  
4 fitting was done using the LME procedure from the NLME library for the R platform (Pinheiro et al.,  
5 2013; R Development Core Team, 2010).

6 Under this procedure, variance parameters for the random effects are first estimated by residual  
7 maximum likelihood (REML) and the fixed effects coefficients are then estimated by weighted least  
8 squares. The null hypothesis that all class means are equal (where the fixed effects are classes) or that  
9 the regression coefficients are zero (where the fixed effects are continuous variables), is tested by the  
10 Wald statistic. These methods are described by Lark et al. (2006). After model fitting summary  
11 statistics and histograms of the residuals were examined to check that these appeared consistent with an  
12 assumption of normality.

13 The subset of ancillary variables for the model with continuous fixed effects was selected by  
14 first fitting a full model with all of the  $\gamma$ -ray (K, U, Th and TC) and DUALEM (1mHcon and 1mPcon)  
15 data as predictors. This was then compared to a series of reduced models by dropping each predictor in  
16 turn, and the full and reduced models were compared by computing their log-likelihood ratio and  
17 testing this against chi-squared with one degree of freedom (Verbeke and Mohlenbergs, 2000). Any  
18 predictor where the reduced model formed by dropping it was not significantly worse than the full  
19 model was rejected. This procedure was repeated until no further predictors were rejected. This

1 procedure requires maximum likelihood rather than REML estimation since residual likelihoods cannot  
2 be compared between models with different fixed effects. Once a predictor set was selected the model  
3 was then re-estimated by REML.

4

## 5 *2.6 Computation of the prediction error variance for class means*

6 The objective of this paper is to compare two approaches to the prediction of soil properties  
7 from the ancilliary variables, given that a relatively small set of direct measurements of the soil  
8 property is available. To do this we compute the expected value of the mean squared prediction error  
9 for the alternative methods:

$$10 \quad \sigma_p^2 = E[\{y - y^*\}^2] \quad (8)$$

11 where  $y$  denotes the value of the target variable at some unsampled location and  $y^*$  denotes the  
12 predicted value.

13 When the predictor is the mean of a class (here obtained by cluster analysis of the ancillary  
14 data) then the mean-squared prediction error in class  $i$  is

$$15 \quad \sigma_{p,i}^2 = \sigma_i^2 (1 + 1/n_i) \quad (9)$$

16 where  $\sigma_i^2$  is the variance of the target property within class  $i$  and the mean of class  $i$  was estimated  
17 from  $n_i$  independently and randomly selected observations within the class (Brus and Lark, 2013). In  
18 this study we use a pooled within-class variance ( $\sigma_w^2$ ). If  $\pi_i$  denotes the relative area of the  $i$ th class out

1 of  $k$  and  $N$  is the total number of observations then the expected value of the mean squared prediction  
2 error for classes is:

$$3 \quad \sigma_{p,C}^2 = \sum_{i=1, \dots, k} \sigma_w^2 \pi_i (1 + 1/N\pi_i) = \sigma_w^2(1+k/N). \quad (10)$$

4 In general, as  $k$  increases we expect the classes to become internally more uniform with respect  
5 to soil properties, so  $\sigma_w^2$  should decrease. However, it is apparent that the term in brackets on the right-  
6 hand side of Equation (10) will increase with increasing  $k$ , and that this increase will be greater the  
7 smaller is  $N$ . In summary,  $\sigma_{p,C}^2$  will only decrease with increasing  $k$  if the reduction in the within-class  
8 variance is large enough to compensate for the fact that the fixed total sample size is spread more thinly  
9 over more classes which contributes to the uncertainty with which the class means are estimated.

10 In this study we computed  $\sigma_{p,C}^2$  for each target soil property for  $k = 2-10$  classes formed by the  
11 FKM algorithm. To do this we require a value of  $\sigma_w^2$ . This was obtained from the LMM, Equation (7),  
12 fitted to the observed soil data for the corresponding classification. The sum of the variances of the  
13 random effects in the model for  $k$  classes as the random effects,  $\sigma_{\eta,k}^2 + \sigma_{\varepsilon,k}^2$ , was treated as the  
14 expected value of the variance for the random variable. This approach is used elsewhere to compute  
15 values for the variances of design-based sample estimates from the results of model-based analyses —  
16 Cochran (1977); Lark (2011) provides an example in soil science. The expected value of the mean  
17 squared prediction error for our classification into  $k$  classes for some sample size  $N$  is computed here as

$$18 \quad * \sigma_{p,C}^2(N | k) = (\sigma_{\eta,k}^2 + \sigma_{\varepsilon,k}^2)(1+k/N). \quad (11)$$

19

1 2.7 Computation of the prediction error variance for regression models

2 For the case of prediction direct from a selected subset of proximally sensed variables by a  
3 multiple regression-type predictor the mean-squared prediction error for a particular prediction at an  
4 unsampled site is:

5 
$$\sigma_{p,R}^2 = \sigma^2 (1 + \mathbf{x}^T \{\mathbf{X}^T \mathbf{X}\}^{-1} \mathbf{x}) \quad (12)$$

6 where  $\mathbf{X}$  is the design matrix for the data set used to predict the model and  $\mathbf{x}$  is a vector in which the  $i$ th  
7 element is the difference between the value of the  $i$ th predictor for the particular prediction and the  
8 overall mean of the  $i$ th predictor (Dudewicz and Mishra, 1988). The term  $\sigma^2$  is the residual variance of  
9 the fitted regression. To compute the expected value of  $\sigma_{p,R}^2$  for some simple random sample of size  $N$   
10 we evaluated Equation (12) for values of the selected predictor variables at each of the  $M=205$   
11 observation sites and computed the average.

12 We obtained the expression  $\mathbf{X}^T \mathbf{X}$  by computing it from the design matrix,  $\mathbf{X}$ , for our original  $M$   
13 observations and rescaling it for a sample of size  $N$  and, as before, we used the sum of the REML  
14 estimates of the variances of the random effects in the corresponding LMM (Equation 7),  $\sigma_{\eta,R}^2 + \sigma_{\epsilon,R}^2$ ,  
15 as the expected value of  $\sigma^2$ . The expected value of the mean squared prediction error from a regression  
16 estimated for some sample size  $N$  is therefore computed here as

17 
$$*\sigma_{p,R}^2(N) = (\sigma_{\eta,R}^2 + \sigma_{\epsilon,R}^2) \{1 + 1/M \sum_{i=1, \dots, M} \mathbf{x}_i^T \{(N/M)\mathbf{X}^T \mathbf{X}\}^{-1} \mathbf{x}_i\} \quad (13)$$

18 where  $\mathbf{x}_i$  is the vector of predictor values for the  $i$ th of our original  $M = 205$  sampled locations and  $\mathbf{X}$  is  
19 the design matrix for that same set.

1 *2.8 Data analysis*

2 The  $\gamma$ -ray spectrometry (i.e. K, U, Th and TC) and DUALEM-1S (i.e. 1mHcon and 1mPcon)  
3 data were first interpolated onto a common 10-m grid. This was done by ordinary kriging (OK) within  
4 a neighborhood of 20-30 and a local variogram. The Vesper program (Minasny et al., 1999) was used.  
5 Numerical clustering of the OK ancillary data was conducted by FKM analysis using FuzME 3.0  
6 (Minasny and McBratney, 2002).

7

8 **Figure 3 (Near here)**

9

10 **3. Results & discussion**

11 *3.1 Spatial distribution of proximally sensed data*

12 Figure 3a) shows the spatial variation of TC. The smallest to intermediate-small counts (<40  
13 cps) are found in the centre of the pasture field. Conversely, intermediate-large (50-60 cps) and large (>  
14 60 cps) TC define the southern parts of the arable field. Similar patterns are evident in K (Figure 3b)  
15 and Th (not shown). It is worth noting the intermediate radioelement values (e.g. K % = 1.5-2 %) at the  
16 southern end of the arable field and in addition the large K readings (>2.5 %) to the south. Figure 3c)  
17 shows the spatial variation of the 1mHcon. For the most part the spatial patterns are similar. This is  
18 particularly the case in the pasture field where  $EC_a$  is smallest (< 6 mS/m). A difference between the  $\gamma$ -



1 ray and  $EC_a$  data is observed in the western half of the southern part of the arable field. Here 1mHcon  
2 is intermediate-large (18-24 mS/m) compared to the eastern half, which is large ( $> 24$  mS/m).

3  
4 **Figures 4 and 5 and Tables 1 (Near here)**

5  
6 *3.2 FKM analysis*

7 Figure 4a) and b) shows the FPI and NCE, respectively for clustering of the sensor data for  
8 different values of  $k$  and  $\phi$ . It is evident that the FPI is at a local minimum when  $k = 2, 3$  and  $4$  when  $\phi$   
9  $= 1.4$  and when  $k = 3$  for values of  $\phi$  from  $1.6$  to  $2.4$ . At several values of  $\phi$  local minima are evident for  
10  $k = 5$  ( $\phi = 1.6$ ),  $k = 7$  ( $\phi = 1.8$ ) and  $k = 8$  ( $\phi = 2.0$ ). As with the FPI, the NCE is at a minimum for each  
11 value of  $\phi$  considered and when  $k = 3$ . Local minima are also evident for  $k$  as indicated above. Given  
12 the equivocal nature of the results, that is the FPI and NCE do not provide a clear indication of what an  
13 appropriate number of  $k$  might be in the ancillary data, we looked to the plot of  $\phi$  versus  $-dJ(\mathbf{M},\mathbf{C})/d\phi$ .  
14 Figure 4c shows that where the derivative of  $-J(\mathbf{M},\mathbf{C})$  is a maximum, McBratney and Moore (1985)  
15 indicate this is where  $\phi$  is optimal.  $\square$  In most cases (i.e.  $k = 5-7$ ) this occurs when  $\phi = 2.0$ . Given these  
16 results we selected an exponent of  $\phi$  equal to  $2.0$  and we compare the results of  $k = 2-10$  classes with  
17 the soil series previously recognized by Palmer (2007).

18  
19 *3.3 Spatial distribution of the FKM classes*

20 Figure 5a) shows the result for  $k = 4$ . Of the four classes, three are identified in the pasture field.  
21 The largest contiguous area is demarcated by 4A which is defined by small radioelement values (e.g. K

1 = 0.86 %) and  $EC_a$  (e.g.  $1mP_{con} = 7.69$  mS/m). Class 4B defines the eastern margin, whilst 4C is  
2 found within 4B as inclusions. In the arable field, 4B defines the northern third, having slightly larger  
3 radioelements (e.g.  $K = 1.90$  %) and  $EC_a$  (e.g.  $1mP_{con} = 14.78$  mS/m) and matches the area ascribed  
4 by the So series. At the southern end, the western half is defined by 4C which was mapped as the wM  
5 series. Of all the classes, 4D has the largest radioelement (e.g.  $K = 2.36$  %) and  $EC_a$  (e.g.  $1mP_{con} =$   
6  $42.11$  mS/m). Moreover the areal extent of 4D coincides with the location of the Wf series.

7 Figure 5b) shows the spatial distribution for  $k = 5$ . The only real difference is that 4A broadly  
8 corresponds to 5A and another class (i.e. 5E), and 5B defines the northern end of the arable field.

9 Figure 6c) shows similar patterns for  $k = 7$ . Here 5D broadly corresponds to 7D and 7F, whereby the  
10 latter has slightly smaller radioelement values (e.g.  $K = 2.31$  %) and  $EC_a$  (e.g.  $1mP_{con} = 40$  mS/m).

11 These results are consistent with Tye et al. (2011), who generated an automated resistivity profile map.

12 In doing so they inferred that the distinct linear variation in the area denoted by 7D is due to outcrops

13 of siltstone beds of the Gunthorpe Member. It is also worth noting that class 5B similarly corresponds

14 to 7B and 7G. The former represents the So series, whilst the latter represents the previously

15 unrecognized Cu series (see Figure 5c). When  $k = 8$  classes are considered 7E corresponds to 8E and

16 8H, otherwise the remaining classes are equivalent (Figure not shown).

17

18

**Figures 6 and 7 (Near here)**

19

### 1 3.4 Mean squared prediction error of a map

2 Figure 6 shows the results of the calculated expected value of mean squared prediction error (i.e.  
3  $\sigma^2_{p,C}$ ) of the estimated class means for different total sample sizes for each of the  $k = 2-10$  maps. Figure  
4 6a) shows the result for clay. It is evident that with increasing  $k$ ,  $\sigma^2_{p,C}$  decreases. This is particularly the  
5 case between  $k = 2$  and 6 classes. A minimum is reached when  $k = 7$  and 8. Beyond this  $\sigma^2_{p,C}$  increases  
6 because, while the classes may become increasingly internally uniform, the total sample effort is now  
7 divided over too many classes to provide adequate estimates of class means.

8 This approach to the selection of the number of classes to use is novel. It differs from  
9 previously-used approaches which focus purely on the internal uniformity of the classes because it  
10 considers the pragmatic question of the precision of predictions based on the classification. In  
11 particular note that the value of the criterion depends on the total sample size that we assume we have  
12 available for calibration. The strength of evidence for a minimum value of the  $\sigma^2_{p,C}$  at some value of  $k$   
13 is also sensitive to the size of the calibration sample.

14 Figure 6b) shows the equivalent plot but for pH. Here  $\sigma^2_{p,C}$  decreases to a local minimum at  $k =$   
15 4, but increases again until a global minimum is attained at  $k = 7$  and 8. This is consistent with the FPI  
16 and MPE metrics of the FKM algorithm; which indicated that  $k = 7$  or 8 might be a suitable number of  
17 classes in the ancillary data.

18

19

### 1 3.5 REML analysis of FKM classes

2 Here we interpret the mean values of soil properties for classes obtained by FKM analysis of the  
3 sensor data with  $k = 7$  because their spatial distribution most closely reflects the soil map (Palmer,  
4 2007). Figure 7 shows the mean values of the soil properties and their standard error for each of the  
5 classes obtained from the LMM estimated by REML. In all cases the Wald statistic allows us to reject  
6 the null hypothesis of no difference among the class means.

7 The particle size fraction appears most revealing about the partitioning of the  $k = 7$  classes.  
8 Figures 7a) shows that 7A has the smallest clay (22.3 %). Conversely, class 7D had the largest clay  
9 (45.7 %) followed closely by 7F. The four remaining classes (i.e. 7B, 7C, 7E and 7G) have similar mean  
10 clay (28.5-34.3 %). Figure 7b) shows that soil pH is largest for 7D (7.1) and the closely related 7F (6.9).  
11 It is worth noting that 7B (5.6) and 7G (6.9) have markedly different pH. This is interesting, given the  
12 FKM analysis of the proximally sensed data did not differentiate these classes until  $k = 7$ . Similarly, and  
13 whilst 7A (6.1) and 7E (6.2) are distinguishable early on in the classification (i.e.  $k = 4$ ) soil pH is not  
14 different between these classes. Worthy of note is that 7B is different from 7D and 7C despite all these  
15 classes being associated with the Mercia mudstone, but it is similar to 7A and 7E which are not.

16

### 17 3.6 Comparison of FKM clustering and regression for predicting soil properties

18 In order to determine whether classifying the  $\gamma$ -ray spectrometer and EM data using FKM  
19 provides better prediction of soil properties than developing regression models we compared the  $\sigma^2_{p,C}$

1 and  $\sigma_{p,R}^2$ , respectively. First we developed the regression models. The reduced models for clay and pH  
2 and the model parameters calculated from R platform are shown in Table 2. In terms of mapping clay,  
3 the best combination of ancillary data is 1mHcon, 1mPcon and K. However, 1mHcon, 1mPcon and U  
4 were selected for mapping pH.

5 The calculated  $\sigma_{p,R}^2$  are shown in Figure 6 and for sample sizes of 20, 40 and 80. The first thing  
6 to note is that as sample size increases,  $\sigma_{p,R}^2$  drops accordingly. This is consistent with the effect of  
7 sampling with regard to the FKM classes (i.e.  $\sigma_{p,C}^2$ ). In order to compare the  $\sigma_{p,C}^2$  using FKM class  
8 means and  $\sigma_{p,R}^2$  using regression models, mean square prediction errors are evaluated for the same  
9 sample size (e.g. 80). In terms of predicting clay content, the regression models generated better  
10 predictions than can be achieved using class means by FKM clustering. However, when predicting pH,  
11 FKM clustering performs better than the linear regression model when  $k = 6 - 9$ .

12

#### 13 **4. Conclusions**

14 A catenary sequence characterised by a broad range of textural variation and associated with  
15 Trent alluvium, Trent River terraces, and Triassic Mercia mudstone, was partitioned into  $k = 2-10$   
16 classes, using fuzzy k-means (FKM) analysis of four radioelement windows (K, U, Th and TC)  
17 acquired from a  $\gamma$ -ray spectrometer and two proximally sensed  $EC_a$  data (1mPcon and 1mHcon). The  
18 use of the FKM algorithm, along with various indices (e.g. FPI and  $-dJ(\mathbf{M},\mathbf{C})/d\phi$ ), suggested that  
19 partitioning the data into  $k = 3, 7$  or  $8$  classes and using a  $\phi$  of 2.0 was most suitable, with the results

1 broadly reflecting a soil series map developed by an experienced soil surveyor using traditional  
2 morphological site descriptions and a pre-existing soil classification scheme (Palmer, 2007).

3 To test this we determined the mean prediction error variance ( $\sigma_{p,C}^2$ ) of the class mean as a  
4 predictor for a soil physical (e.g. clay content) and chemical (i.e. pH) property. Using this independent  
5 approach the results indicated that  $k = 7$  and  $8$  were statistically different and accounted for most of the  
6 soil variation. Prediction of soil properties by FKM analysis and regression models was also compared.  
7 The results of this analysis indicated that for a sample size of  $80$ , the regression models were able to  
8 predict clay content, better than FKM clustering. However, when predicting pH, FKM clustering  
9 performs better than the linear regression model when  $k = 6 - 9$ . It is concluded that both the FKM and  
10 LMM methods have merit. In the case of the clustering, the approach is able to account for soil  
11 properties which have non-linearity with the ancillary data (i.e. pH), whereas the LMM approach is  
12 best when there is a strong linear relationship (e.g. clay).

13 In order to test this further, for example if we wanted to make predictions of soil properties  
14 across the Trent valley where conditions are homologous with our study site, we would be unable to  
15 sample with enough intensity to use kriging. However, we could collect proximally sensed  $\gamma$ -ray  
16 spectrometer and DUALEM-1 data, plus a set of calibration data. Then using the theory for computing  
17 the  $\sigma_{p,R}^2$  for the regression and the  $\sigma_{p,C}^2$  for classes of different sizes we could predict variations of the  
18 properties using regression or by using class means as predictors.

19

## 1 **5. References**

- 2 Avery, B.W., 1980. Soil Classification for England and Wales (Higher Categories). Soil Survey  
3 Technical Monograph No. 14. Harpenden.
- 4 Bezdek, J.C., 1981. Pattern recognition with fuzzy objective function algorithms. Plenum Press, New  
5 York.
- 6 Brus, D., Lark, R.M., 2013. Soil survey. In (A.H. El-Shaarawi, ed.) Encyclopaedia of Environmetrics,  
7 2nd Edition, John Wiley & Sons.
- 8 Cochran, W.G., 1977. Sampling Techniques. Wiley, New York.
- 9 Dang, Y. P., Dalal, R. C., Pringle, M. J., Biggs, A. J. W., Darr, S., Sauer, B., . . . Orange, D. 2011.  
10 Electromagnetic induction sensing of soil identifies constraints to the crop yields of north-eastern  
11 Australia. Soil Research, 49, 559-571.
- 12 DUALEM-421S., 2008. User's Manual. Dualem Inc., Canada.
- 13 Dudewicz, E. J., Mishra, S. N., 1988. Modern mathematical statistics. Wiley, New York.
- 14 Gooley, L., Huang, J., Page, D., Triantafilis, J., 2013. Digital soil mapping available water content  
15 using proximal and remotely sensed data, Soil Use and Management, in press.
- 16 Goovaerts, P., 1997. Geostatistics for Natural Resources Evaluation. Oxford Univ. Press, New York.
- 17 Guo, Y., Shi, Z., Li, H.Y., Triantafilis, J., 2013. Application of digital soil mapping methods to identify  
18 salinity management classes in coastal lands of central China. Soil Use and Management, 29, 445-  
19 456.

- 1 Hodgson, J.M. 1976. Soil survey field handbook, 2nd ed. Technical monograph (Soil Survey of  
2 England and Wales), No. 5. Soil Survey of England and Wales, Harpenden,UK.
- 3 Jung, W. K., Kitchen, N. R., Sudduth, K. A., Anderson, S. H., 2006. Spatial characteristics of claypan  
4 soil properties in an agricultural field. *Soil Science Society of America Journal* 70, 1387-1397.
- 5 Lark, R. M., Cullis, B. R., Welham, S. J., 2006. On spatial prediction of soil properties in the presence  
6 of a spatial trend: the empirical best linear unbiased predictor (E-BLUP) with REML. *European*  
7 *Journal of Soil Science* 57, 787-799.
- 8 Lark, R. M., 2011. Spatially nested sampling schemes for spatial variance components: Scope for their  
9 optimization. *Computers and Geosciences* 37, 1633-1641.
- 10 Lesch, S. M., Corwin, D. L., Robinson, D. A., 2005. Apparent soil electrical conductivity mapping as  
11 an agricultural management tool in arid zone soils. *Computers and Electronics in Agriculture* 46(1-3  
12 SPEC. ISS.), 351-378.
- 13 McBratney, A.B., Moore, A.W., 1985. Application of fuzzy-sets to climatic classification. *Agricultural*  
14 *Forest Meteorology* 35, 165-185.
- 15 McBratney, A.B., De Gruijter, J.J., Brus, D.J. 1992. Spatial prediction and mapping of continuous soil  
16 classes. *Geoderma* 54, 39-64.
- 17 Minasny, B., McBratney, A.B., Whelan, B.M., 1999. VESPER version 1.6, Precision Agriculture  
18 Laboratory, Sydney, Australia.



1 Minasny, B., McBratney, A.B., 2002. FuzME version 3.0, Precision Agriculture Laboratory, The  
2 University of Sydney, Australia.

3 Palmer, R.C., 2007. The soils of Nottingham Trent University farm at Brackenhurst, Nottinghamshire.  
4 NSRI unpublished research, report no. YE20006E for the British Geological Survey.

5 Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D., The R Development Core Team, 2013. nlme: Linear  
6 and Nonlinear Mixed Effects Models. R package version 3.1-110.

7 Priori, S., Martini, E., Andrenelli, M. C., Magini, S., Agnelli, A. E., Bucelli, P., Biagi, M., Pellegrini, S.,  
8 Costantini, E. A. C. 2013. Improving wine quality through harvest zoning and combined use of  
9 remote and soil proximal sensing. *Soil Science Society of America Journal*, 77, 1338-1348.

10 R Development Core Team 2010. R: A language and environment for statistical computing. R  
11 Foundation for Statistical Computing, Vienna, Austria.

12 Robinson, D.A. Binley A., Crook N., Day-Lewis F.D., Ferre T.P.A., Grauch V.J.S. Knight R., Knoll  
13 M., Lakshmi V., Miller R., Nyquist J., Pellerin L., Singha K., and Slater L., 2008. Advancing  
14 process-based watershed hydrological research using near-surface geophysics: A vision for, and  
15 review of, electrical and magnetic geophysical methods. *Hydrological Processes*, 20, 3604-3635.

16 Robinson, D. A., Abdu, H., Lebron, I., Jones, S. B., 2012. Imaging of hill-slope soil moisture wetting  
17 patterns in a semi-arid oak savanna catchment using time-lapse electromagnetic induction. *Journal of*  
18 *Hydrology* 416, 39-49.

1 Triantafilis, J., Lesch, S. M., La Lau, K. & Buchanan, S. M. 2009a. Field level digital soil mapping of  
2 cation exchange capacity using electromagnetic induction and a hierarchical spatial regression model.  
3 *Australian Journal of Soil Research*, **47**, 651-663.

4 Triantafilis, J., Kerridge, B., Buchanan, S.M., 2009b. Digital soil-class mapping from proximal and  
5 remotely sensed data at the field level. *Agronomy Journal* 101, 841-853.

6 Triantafilis, J., Gibbs, I.D., Earl, N.Y., 2013. Digital soil pattern recognition in the lower Namoi valley  
7 using numerical clustering of gamma-ray spectrometry data. *Geoderma* 192, 407-421.

8 Tye, A. M., Kessler, H., Ambrose, K., Williams, J. D. O., Tragheim, D., Scheib, A., Kuras, O., 2011.  
9 Using integrated near-surface geophysical surveys to aid mapping and interpretation of geology in  
10 an alluvial landscape within a 3D soil-geology framework. *Near Surface Geophysics* 9, 15-31.

11 Van Meirvenne, M., Islam, M. M., De Smedt, P., Meerschman, E., Van De Vijver, E., & Saey, T. 2013.  
12 Key variables for the identification of soil management classes in the aeolian landscapes of north-  
13 west Europe. *Geoderma*, 199, 99-105.

14 Verbeke, G., Molenberghs, G. 2000. *Linear Mixed Models for Longitudinal Data*. Springer, New York.

15 Webster, R., Oliver, M. A., 1992. Sample adequately to estimate variograms of soil properties. *Journal*  
16 *of Soil Science* 43, 177-192.

17 Woodforth, A., Triantafilis, J., Cupitt J., Malik, R.S., Geering, H., 2012. Mapping estimated deep  
18 drainage in the lower Namoi Valley using a chloride mass balance model and EM34 data.  
19 *Geophysics* 77, WB245-256.



1 **Table 1** Euclidean centroid values of proximally sensed ancillary data clustered using FKM and for classes  $k = 4, 5$  and  $7$ .  
 2 Note: centroids shown for K (%), U (ppm), Th (ppm), TC (cps-counts per second), 1mPcon and 1mHcon, respectively.  
 3 Note: Number of soil samples which are members of the classes also shown.  
 4

	$k = 4$		$k = 5$		$k = 7$	
	<i>Centroid values</i>	<i>members</i>	<i>Centroid values</i>	<i>members</i>	<i>Centroid values</i>	<i>members</i>
	0.86, 1.28, 3.74, 30.51, 8, 5	57	0.83, 1.26, 3.67, 29.79, 7, 5	49	0.81, 1.25, 3.64, 29.46, 6, 4	41
	1.90, 1.93, 5.76, 50.95, 15, 9	45	1.95, 1.96, 5.86, 51.86, 15, 9	36	1.95, 1.96, 5.85, 51.84, 15, 9	21
	2.09, 2.19, 6.18, 55.38, 35, 19	45	2.12, 2.20, 6.23, 55.91, 35, 20	47	2.07, 2.18, 6.15, 54.98, 34, 19	37
	2.36, 2.23, 6.61, 60.05, 42, 27	58	2.37, 2.23, 6.62, 60.16, 42, 26	54	2.43, 2.23, 6.72, 61.18, 44, 29	23
			1.30, 1.53, 4.55, 38.96, 19, 13	19	1.20, 1.48, 4.41, 37.23, 16, 11	23
					2.31, 2.22, 6.51, 59.06, 40, 24	42
					2.02, 1.99, 5.98, 53.27, 17.36, 11.17	15

5  
 6  
 7

1 **Table 2 Reduced models for different soil properties.**

Properties	Selected predictors	$\sigma^2$	Covariance function	Nugget fraction	Distance parameter
Clay content	1mHcon, 1mPcon, K	31.360	Spherical	0	33 m
pH	1mHcon, 1mPcon, U	0.276	Exponential	0.3	65 m

2 Note:  $\sigma^2$  - residual variance of the fitted regression; nugget fraction - nugget/sill; and, distance parameter - distance  
3 parameter of the exponential correlation function or range of the spherical correlation function.

4

1 **Figure Captions**

2

3 **Figure 1** Location of study area east of Nottingham and River Trent and the Soil Series Map.

4

5 **Figure 2** a) Air-photo of the pasture and arable study fields, b) spatial location of the DUELEM-1S and  
6 gamma-ray ( $\gamma$ -ray) spectrometry survey transects, and c) soil sample locations.

7

8 **Figure 3** Spatial distribution of gamma-ray ( $\gamma$ -ray) spectrometry data including; a) Total count (TC –  
9 counts per second), and b) potassium (K - %), and DUALEM-1S electrical conductivity ( $EC_a$  – mS/m)  
10 of c) 1mHcon (i.e. Deep).

11

12 **Figure 4** Plot of; (a) fuzziness performance index (FPI), (b) normalized classification entropy (NCE)  
13 versus classes ( $k = 2$  to 10) and (c) fuzziness exponent ( $\phi$ ) versus  $-dJ(M,C)/d\phi$ .

14

15 **Figure 5** Spatial distribution of FKM classes for  $k =$  a) 4, b) 5, and c) 7.

16

17 **Figure 6** Plot of mean squared prediction error for FKM classes (i.e.  $\sigma_{p,C}^2$ , solid lines) and linear  
18 regression models (i.e.  $\sigma_{p,R}^2$ , dashed lines) of (a) clay (%) and (b) pH. Note: Sample size are labeled  
19 (i.e.  $N = 20, 40, 80, \dots, 205$ ).

20

21 **Figure 7** Plot of mean and standard deviation of soil (a) clay (%) and (b) pH.

1  
2  
3  
4

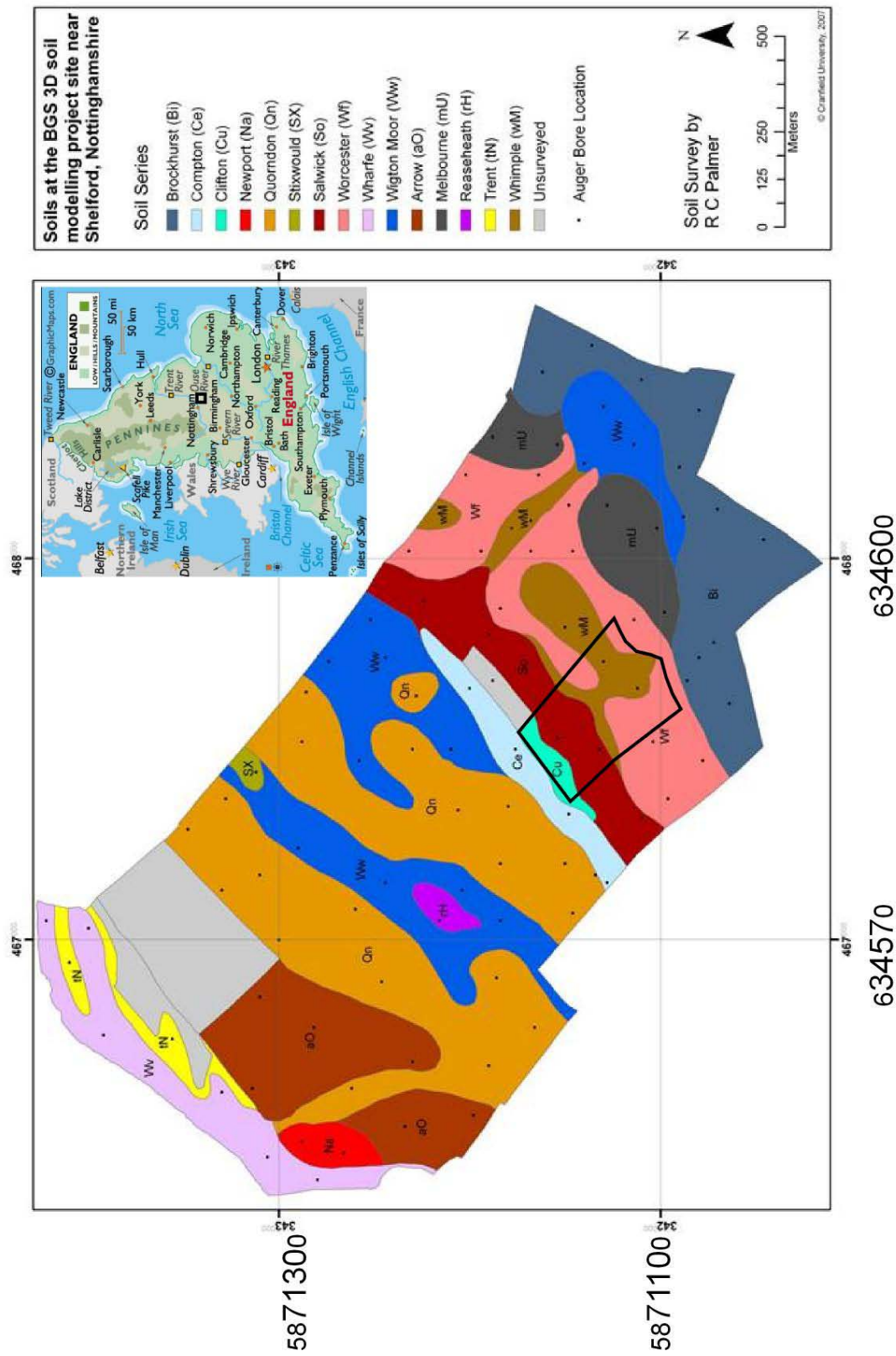
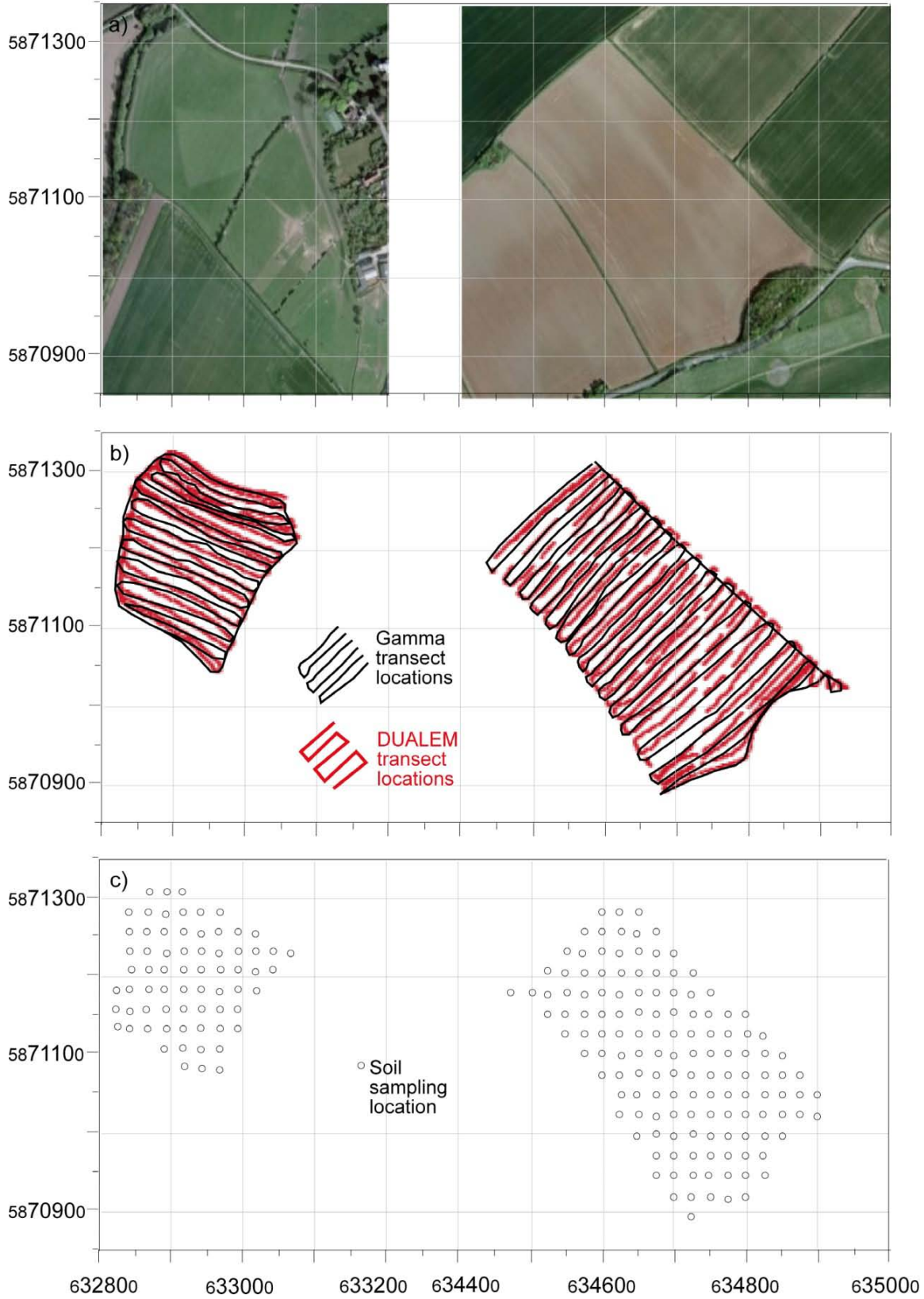


Figure 1

Northings (m)



1  
2  
3

**Figure 2**



Northings (m)

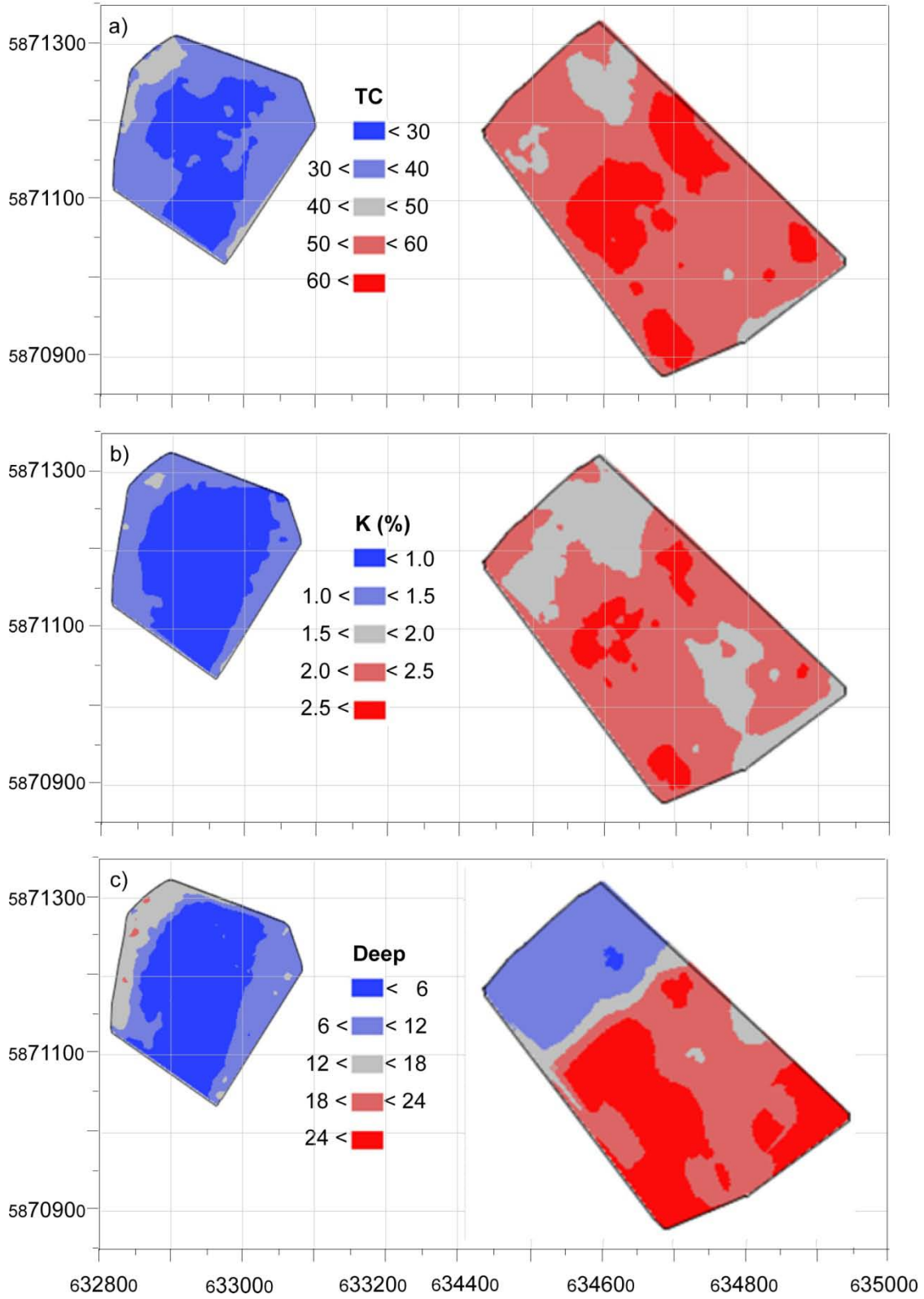


Figure 3

1  
2  
3  
4

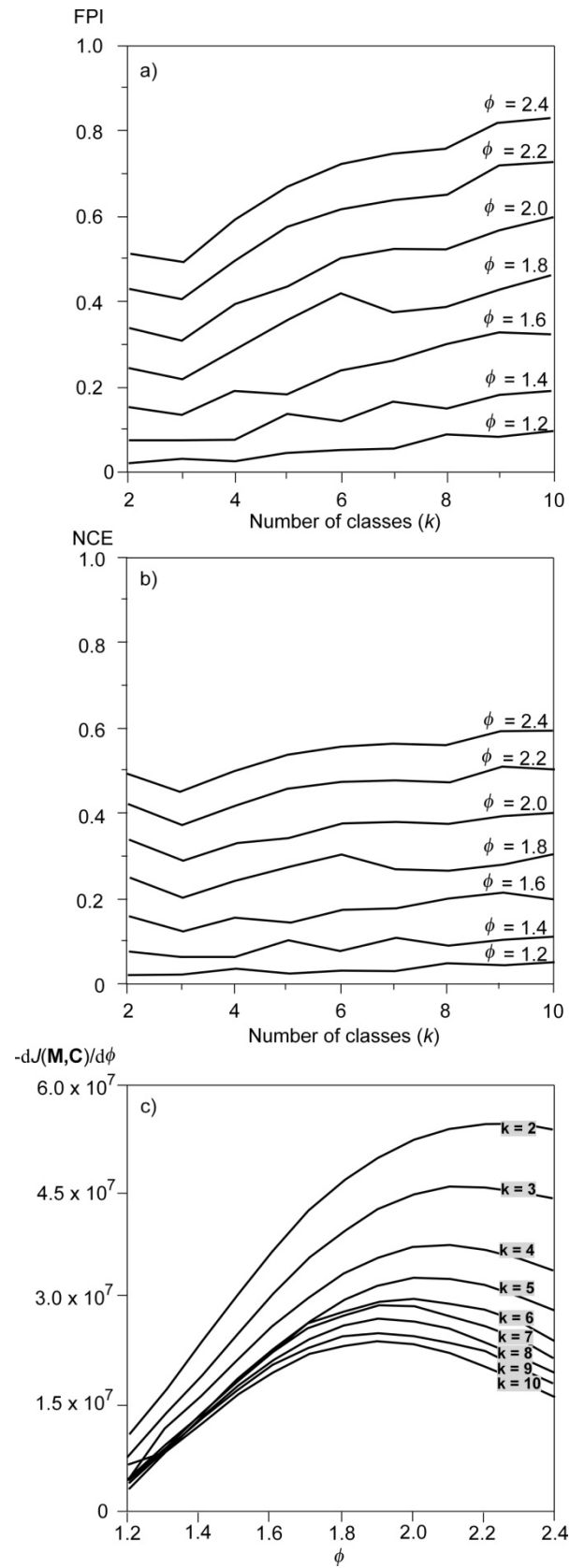


Figure 4

1  
2  
3

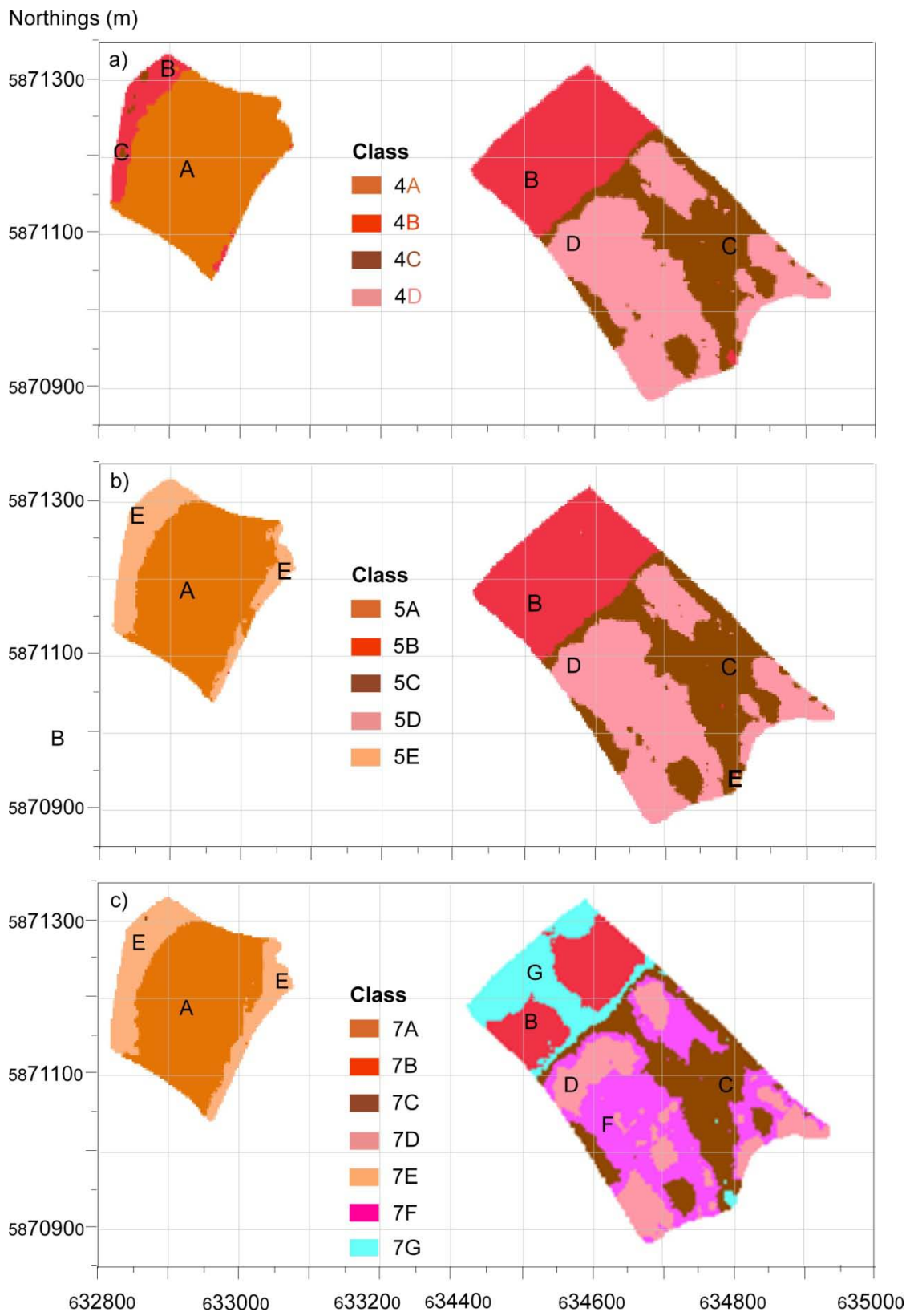


Figure 5

1  
2  
3

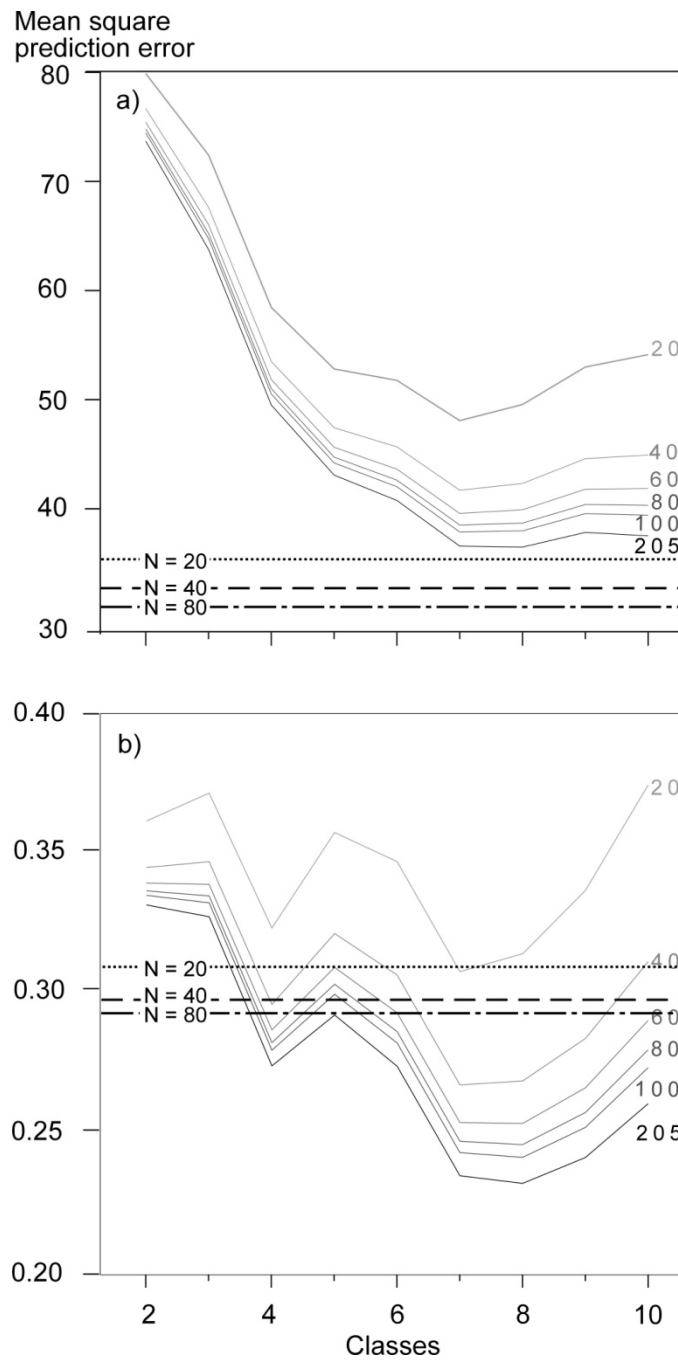
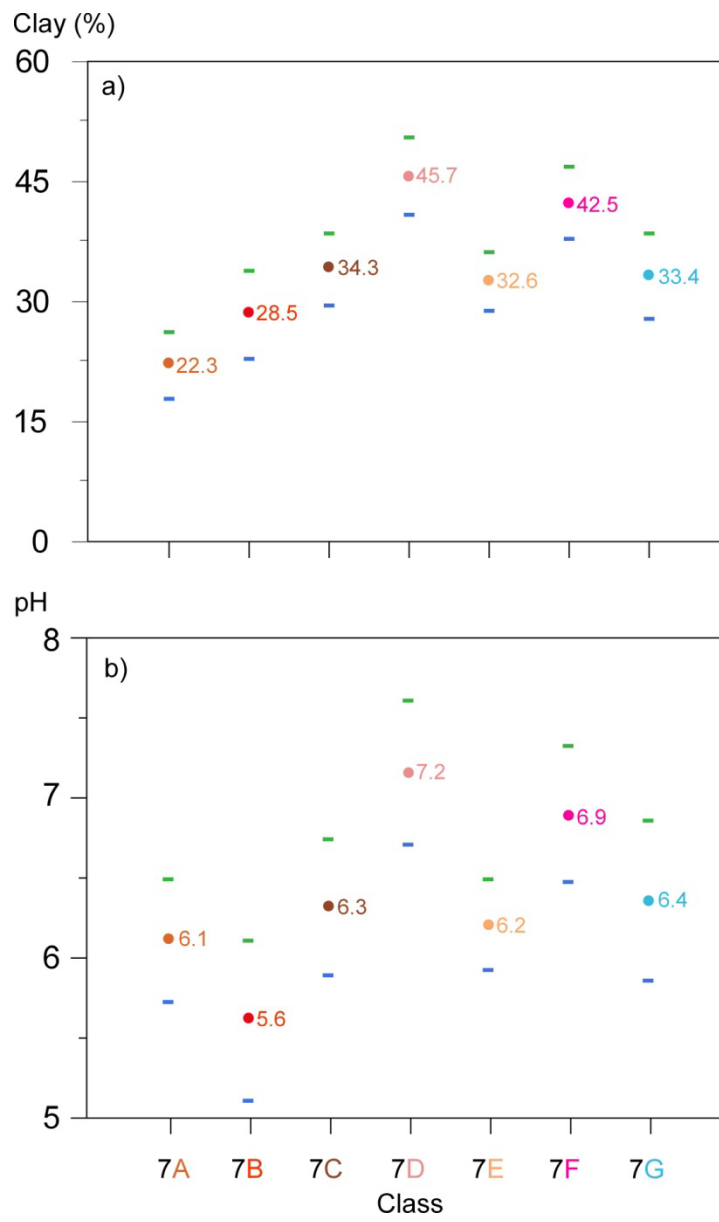


Figure 6

1  
2  
3



**Figure 7**

1  
2  
3