

## Article (refereed) - postprint

---

Henrys, P.A.; Bee, E.J.; Watkins, J.W.; Smith, N.A.; Griffiths, R.I. 2015.  
**Mapping natural capital: optimising the use of national scale datasets.**  
*Ecography*, 38 (6). [10.1111/ecog.00402](https://doi.org/10.1111/ecog.00402)

© 2014 The Authors. *Ecography* © 2014 Nordic Society Oikos

This version available <http://nora.nerc.ac.uk/508614/>

NERC has developed NORA to enable users to access research outputs wholly or partially funded by NERC. Copyright and other rights for material on this site are retained by the rights owners. Users should read the terms and conditions of use of this material at <http://nora.nerc.ac.uk/policies.html#access>

**This document is the author's final manuscript version of the journal article, incorporating any revisions agreed during the peer review process. There may be differences between this and the publisher's version. You are advised to consult the publisher's version if you wish to cite from this article.**

The definitive version is available at <http://onlinelibrary.wiley.com/>

Contact CEH NORA team at  
[noraceh@ceh.ac.uk](mailto:noraceh@ceh.ac.uk)

1 **Mapping Natural Capital: Optimising the use of national scale datasets**

2 Henrys, P. A.<sup>1</sup>, Bee, E. J.<sup>2</sup>, Watkins, J. W.<sup>1</sup>, Smith, N.A.<sup>3</sup> and Griffiths, R. I.<sup>4</sup>

3 <sup>1</sup> Centre for Ecology and Hydrology, Lancaster Environment Centre, Library Avenue,  
4 Bailrigg, Lancaster, LA1 4AP, UK.

5 <sup>2</sup> British Geological Survey, Environmental Science Centre, Nicker Hill, Keyworth,  
6 Nottingham, NG12 5GG, UK.

7 <sup>3</sup> British Geological Survey, Murchison House, West Mains Road, Edinburgh, EH9 3LA, UK.

8 <sup>4</sup> Centre for Ecology and Hydrology, Maclean Building, Benson Lane, Crowmarsh Gifford,  
9 Wallingford, Oxfordshire, OX10 8BB, UK.

10 *Correspondence:* Peter Henrys, Centre for Ecology and Hydrology, Lancaster Environment  
11 Centre, Library Avenue, Bailrigg, Lancaster, LA1 4AP, UK.

12 Email: pehn@ceh.ac.uk

13 Tel: +44(0)1524 595839

**14 Abstract**

15 Understanding the spatial distribution of specific environmental variables and the  
16 interdependencies of these variables is crucial for managing the environment in a sustainable  
17 way. Here we discuss two methods of mapping – a Geographical Information System  
18 classification-based approach and a statistical model-based approach. If detailed, spatially  
19 comprehensive covariate datasets exist to complement the ecological-response data, then  
20 using a statistical model-based analysis provides the potential for greater understanding of  
21 underlying relationships, as well as the uncertainty in the spatial predictions. Further, the  
22 model-based approach facilitates scenario testing. Although similar methods are already  
23 adopted in species distribution modeling, the flexibility of the model framework used is  
24 rarely exploited to go beyond modeling occupancy or suitability for a single species, into  
25 modeling complex derived metrics such as community composition and indicators of natural  
26 capital. As an example, we assess the potential benefits of the statistical model-based  
27 approach to mapping natural capital through the use of two national survey datasets; The  
28 Centre for Ecology and Hydrology (CEH) Land Cover Map (LCM) and the British  
29 Geological Survey's (BGS) Parent Material Model (PMM), to predict national soil microbial  
30 community distributions based on data from a sample of > 1000 soils covering Great Britain.  
31 The results are mapped and compared against a more traditional, land classification-based  
32 approach. The comparison shows that, although the maps look broadly similar, the model-  
33 based approach provides better overall spatial prediction, and the contribution of individual  
34 model terms (along with their uncertainty) are far easier to understand and interpret, whilst  
35 also facilitating any scenario testing. We therefore both recommend the use of spatial  
36 statistical modelling techniques to map natural capital and anticipate that they will become  
37 more prominent over the forthcoming years.

## 38 **Introduction**

39 The Millennium Ecosystem Assessment (2005) and more recently in the UK, the National  
40 Ecosystem Assessment (2011), stress the importance of ecosystems and understanding the  
41 interdependencies between their underlying drivers of change (Carpenter et al. 2009; Feld et  
42 al. 2009; Norgaard 2009). Ecosystem ‘natural capital’ can be identified, according to  
43 Costanza and Daly (1992), as the “assets” or “stock that yields a flow of valuable goods or  
44 services into the future”. This concept of “natural capital” and “flow of goods” has gained  
45 traction in recent years and has been used as a way of bridging the scientific-economic-  
46 policymaking divide, enabling the potential impact of ecosystem modification to be better  
47 evaluated, and more meaningfully incorporated, into decisions affecting society (National  
48 Research Council, 2005; Millennium Ecosystem Assessment, 2005). Knowledge regarding  
49 the spatial distribution of ecological systems and the natural capital stocks that they produce  
50 is of crucial importance for managing the effects of human pressure and environmental  
51 change on natural resources (Swetnam et al. 2011; Naidoo and Ricketts 2006).

52 In order to investigate spatial distribution and variation in natural capital, it is crucial to make  
53 use of all available data, both on the natural capital indicator itself and on complementary  
54 datasets that are *a priori* thought to drive changes in this response – it is important from the  
55 outset that ecological understanding of the system and any synthesis of it are clearly thought  
56 about (Austin, 2002). This is to provide unbiased estimates of stocks of natural capital and  
57 related ecosystems, enabling planners and policy makers to identify the most economically or  
58 environmentally desirable trade-offs (Turner et al. 2010; Nengwang et al. 2009). For  
59 example, the availability of suitable habitat for wild pollinator populations may vary  
60 depending on the relative strength of the different abiotic and biotic environmental drivers  
61 present, such as climate, soil, geology or types of habitats. One approach to investigate the  
62 spatial distribution of natural capital may be based on a geographical stratification of the

63 region of interest according to environmental conditions. However, a simple environmental  
64 stratification or categorical classification does not provide the flexibility to analyse different  
65 drivers, measure their relative strength in determining how stocks are currently distributed, or  
66 predict how these may change under future management or environmental change scenarios.  
67 All of these require a more flexible approach capable of making best use of a range of source  
68 data.

69 Two examples of the Geographical Information System (GIS) classification-based approach  
70 illustrate its shortcomings. For example, the US Geological Society (USGS) generated a map  
71 of standardised terrestrial ecosystems across the US that could be useful for studies of the  
72 production and value of ecosystem goods and services and indicators thereof (Sayre et al.  
73 2009). The map is derived by classifying areas according to a set of environmental covariates  
74 that describe features such as climate and geology. The Institute of Terrestrial Ecology's  
75 (ITE) land classification of Great Britain provides a similar map of environmental classes,  
76 defined according to a clustering technique imposed on a multivariate ordination, and was  
77 based on multiple covariate data sets such as geology, topography and climate (Bunce et al.  
78 1996). The assumption made is that all important covariate effects are accounted for in the  
79 classification. These classification maps of environmental or ecosystem strata can provide a  
80 basis on which one can overlay, and hence map, specific indicators of natural capital based  
81 on the spatial pattern of the strata. However, any further inference, uncertainty analysis, or  
82 testing of assumptions and hypotheses, is not possible as the classes are fixed and we cannot  
83 disaggregate which drivers are most important for understanding the regional variation or  
84 extent of the natural capital indicator in question. Furthermore, one can only make inferences  
85 regarding change and association within the existing classification structure, and they cannot  
86 be used to predict the outcome after environmental changes (such as climate change or  
87 different land use regimens). Such GIS classification-based approaches are commonly used to

88 map natural capital and ecosystem service indicators (eg. Norton et al., 2012; Troy and  
89 Wilson, 2006; Raymond et al., 2009; Costanza et al., 2006), but any uncertainty analysis or  
90 understanding of spatial dependence is rarely explored as the classification approach does not  
91 lend itself to this.

92 In contrast, using spatial statistical models with an ability to compensate for or make use of  
93 spatial autocorrelation, the high quality, geographically widespread spatial data used in the  
94 aforementioned GIS classifications can be further exploited to enable both predictive  
95 geographic infilling across space and estimation of specific covariate effects. Such  
96 approaches, however, rely on good spatial coverage of both the observation data and the  
97 predictor variables used to build the models. As many different forms of spatial  
98 environmental data (such as rasters) are becoming more accessible, and GIS tools become  
99 more ubiquitous, the development of methods which make best use of these data for  
100 environmental research is timely and of increasing importance in dealing with environmental  
101 change scenarios and providing appropriate advice to policy makers and environmental  
102 stakeholders.

103 The use of similar statistical regression modeling techniques, such as standard GLMs  
104 (McCullagh and Nelder, 1989), GAMs (Hastie and Tibshirani, 1990) and MARS (Friedman,  
105 1991), has been common in both epidemiology and in species distribution / ecological niche  
106 modeling for some time. In the epidemiology literature such approaches are commonly used  
107 to map disease risk, incidence and spread (Vieira et al., 2005, Nguyen et al., 2012, French  
108 and Wand, 2004). In the ecology literature attention has been more focused on predictive  
109 modeling and understanding environmental effects rather than purely spatial analysis (e.g.  
110 Kriging or GIS classification). The mapping approach presented here demonstrates the use of  
111 a species distribution modeling regression approach with the inclusion of a spatial correlation  
112 structure (as we are ultimately interested in the spatial distribution). Although sometimes

113 included when modeling and mapping individual species' distributions, this approach has  
114 rarely been applied specifically to the concept of mapping natural capital and indicators  
115 thereof.

116 In this paper, we present the application of a spatial statistical regression model using two  
117 national-scale data sets to explore the benefits of this approach against the use of simple  
118 environmental stratification. We reflect on how such approaches could be used to gain  
119 information on the distribution and extent of natural capital, and multiple environmental  
120 indicators across Great Britain.

## 121 **Materials and Methods**

### 122 National scale environmental data

123 The mapping of environmental indicators, either by GIS classification or statistical  
124 modelling, requires high quality observation data and covariate data with good spatial  
125 coverage (no obviously sparse areas) over the region of interest, preferably at high resolution  
126 with sufficient sample size. The Centre for Ecology and Hydrology (CEH) and the British  
127 Geological Survey (BGS) provide spatial information across Great Britain at 25 m and 50 m  
128 resolution on land-cover and parent material, respectively. Having national coverage of two  
129 key land-surface influences is important in determining the potential location of natural  
130 capital. Hence the two covariates can provide a solid basis for modelling and mapping natural  
131 capital and ecological responses to changes in land cover and parent material at the national  
132 scale. In the future, other covariates could be incorporated into the methodology, but for  
133 simplicity and as an example only two have been used in this paper.

134 The Land Cover Map 2007 (LCM2007) provides information about physical materials on the  
135 Earth's surface over the UK (Morton et al. 2011). Such physical materials may be manmade

136 urbanised areas consisting of roads or buildings, or natural materials such as vegetation,  
137 exposed rock on inland water. The LCM2007, derived from satellite imagery, was produced  
138 as part of the Countryside Survey of the UK as a snapshot audit (Morton et al. 2011). Ground  
139 truthing and knowledge-based enhancements are also used to derive the physical coverage  
140 from the satellite images that make up the dataset, which is a continuous parcel-based  
141 (polygon) dataset accompanied by a suite of derived raster products with 25 m and 1 km  
142 resolution.

143 The Parent Material Model (PMM) is a spatial database representing below ground material  
144 from which the topsoil develops (Lawley, 2008). The PMM enables the distribution of  
145 physiochemical properties of the weathered and un-weathered parent materials to be mapped.  
146 It details over 30 rock and sediment characteristics adding simplified classifications of  
147 lithological properties. The attribute content includes a range of texture information, colour,  
148 structure, mineralogy, lithology, carbonate content and information about how the parent rock  
149 was formed (genetic origin) (British Geological Survey, 2013).

#### 150 Natural Capital Data

151 As an example assessment of the possible benefit gained by adopting a geostatistical  
152 modelling approach over classification methods, we consider data on soil microbial  
153 community structure obtained from Countryside Survey (CS) 2007 (Norton et al. 2012). This  
154 dataset represents information on bacterial biodiversity at a nationwide extent. Soil bacterial  
155 biodiversity can be considered a good indicator unifying various parameters pertaining to  
156 natural capital, in that it is a biodiversity measure responsive to both natural fixed  
157 environmental factors such as geology and also changes in climate and land use (Griffiths et  
158 al., 2011). In a previous study analyzing these data, Griffiths et al., 2011 used a molecular  
159 approach (Terminal Restriction Fragment Length Polymorphism) to characterise the bacterial

160 communities in soils from over 1000 cores sampled across Great Britain within the  
161 Countryside Survey sampling framework, which consisted of up to five randomly sampled  
162 soils taken from over 200 1-km<sup>2</sup> locations across GB. In their study, non-metric  
163 multidimensional scaling (NMDS) was used on the Bray-Curtis similarities of the community  
164 profiling results to define community composition in two dimensions. The first axis scores  
165 resulting from their ordination form the microbial community data used in the remainder of  
166 the work presented here.

167 The data were assessed by Griffiths et al., 2011 in relation to other environmental variables  
168 collected as part of the survey, including abiotic aspects of the environment as well as soil  
169 physical and chemical parameters. Those authors found that bacterial communities at this  
170 landscape scale were structured in similar manner to plants, and were highly correlated with a  
171 general gradient of soil parameters from acidic-organic soils to neutral soils of lower organic  
172 matter. This gradient was apparent in the first axis NMDS site scores, which generally  
173 increased with increased soil pH, and declining organic matter. These soil features are  
174 generally determined by the underlying geology and climate as well as associated human land  
175 usage. Therefore soil pH and plant biodiversity ordination scores were found to be amongst  
176 the best variables correlating with measures of bacterial biodiversity, but the aggregate  
177 vegetation classification (AVC) was also a strongly predictive factor.

178 To upscale the data from the discrete sampled locations and produce a GB scale map,  
179 Griffiths et al (2011) used the interpolation technique inverse distance weighting (Figure 1).  
180 Such a map is successful in illustrating the broad differences in communities between, for  
181 example, England and Scotland, but is unlikely to hold predictive power at smaller spatial  
182 scales. Here, we suggest that since vegetation cover and pH are strong predictors, and that the  
183 observed dataset has good spatial coverage due to the stratified sampling design of CS, we  
184 can use a more informative model-based approach to make more predictive spatial

185 extrapolations. In particular we seek to test whether a more predictive spatial mapping can be  
186 obtained by using the LCM and PPM national coverage maps, compared to making naive use  
187 of an existing classification.

### 188 Statistical analysis

189 Given data on a numerical indicator of natural capital with suitable spatial coverage over the  
190 region in question, statistical models can be used to model the relationship between the  
191 indicator and other environmental covariate data. The model framework adopted needs to be  
192 flexible enough to cover the potentially complex structure of the observational data, whilst at  
193 the same time taking care to avoid false assumptions of independence, normality and  
194 linearity. An example of such a framework is the Generalised Linear Geostatistical Model  
195 (GLGM) of Diggle et al. (1998). This framework can easily be extended to a more generic  
196 setting where the linearity assumption is relaxed to form a Generalised Additive  
197 Geostatistical Model (GAGM) following on from the Generalised Additive Model framework  
198 (Hastie and Tibshirani 1990), which is already commonly adopted in species distribution  
199 modeling. The underlying model framework of a GAGM consists of three parts: 1) a linear  
200 combination of potentially smoothly varying covariate functions; 2) a spatial random field,  
201 which we will define as a Stationary Gaussian Process (SGP); and 3) random effects  
202 representing underlying, potentially non-spatial, error structure.

203 Having modelled the relationship between stock estimates of particular indicators reflecting  
204 national capital (such as: soil carbon; water quality; plant species occurrence; and in this  
205 instance soil microbial community structure) and the environmental covariates, one can,  
206 within the bounds of the training data, interpolate across unsampled geographic regions using  
207 information on the covariates available over finer spatial scales. For prediction of this sort it  
208 is essential that the observed data demonstrate both good spatial coverage and good covariate

209 coverage such that predictions are not made beyond the range of this training data set—i.e. all  
210 geographic areas where we wish to make predictions are represented and the full range of  
211 covariate values are represented in the data set that the models were built on. In species  
212 distribution modeling, this is often referred to as the difference between analog and non-  
213 analog conditions (see for example Williams and Jackson, 2007; Veloz et al., 2012; Algar et  
214 al., 2009), where non-analog conditions are those unlike any previously observed in the  
215 study. Providing that the geographic and covariate space over which predictions are sought  
216 has a suitable analog in the observed data, substituting the wide coverage covariate data into  
217 the estimated model achieves predictions over the same spatial extent for the same snapshot  
218 in time as the observed response data. The geostatistical model-based approach of Diggle et  
219 al. (1998) has the clear advantage over simple kriging and GIS classification that both spatial  
220 correlation structure and covariate effects are taken into account. Furthermore, the model-  
221 based approach allows for simple extraction of the estimated error structure, and hence we  
222 can quantify the uncertainty in the predictions. Further details on the model framework  
223 including mathematical specification are provided in Supplementary Material Appendix 1.

224 In following this modeling procedure, we first carried out a GIS ‘points in polygon’  
225 procedure to concatenate the CS data on microbial communities with corresponding data on  
226 land cover and calcium carbonate content. The final dataset consisted of 1010 observations.  
227 The raw data on soil microbial community ordination scores were modeled against broad  
228 habitat and calcium carbonate content using a generalised additive mixed-model (Lin and  
229 Zhang, 1999) approach. This follows the same generic approach as the GAGM without the  
230 inclusion of a spatial random field, which was deemed redundant upon examination of model  
231 residuals using Moran’s *I*. The random components in the mixed model were needed to  
232 account for the apparent non-independence between any two soil cores taken from the same  
233 1km square. These were more likely to be similar than two cores taken from two different

234 squares. Alongside the random effects and fixed effects of habitat and calcium carbonate, an  
235 additional spatial surface was included to account for residual large scale spatial variation.

236 The model was fitted, including the smoothly varying spatial surface using tensor product  
237 smooth interactions, via the gamm function in the ‘mgcv’ library (Wood, 2011) in the R  
238 statistical environment (R Development Core Team, 2008). Estimates of the model  
239 parameters were obtained using restricted maximum. Full details of model specification and  
240 testing are provided in Supplementary Material Appendix 2, which also provides details on  
241 model fitting when the spatial random field is needed in the model formula.

242 For purposes of comparison, we then used the ITE land classification (Bunce et al., 1996) to  
243 produce a classification-based assessment. This was obtained by simply taking the mean  
244 microbial ordination axis score per land class. As the same land classification is used to  
245 classify the CS samples, sufficient sample size was guaranteed in each classification segment.  
246 What we are hence comparing is a model-based map versus the naive use of an existing  
247 geographic classification. Existing classification maps are often used in this way as it is not  
248 always feasible to develop a new classification for each purpose.

249 Examination of the mean square error of the predictions against the observed data provides a  
250 formal comparison of the goodness of fit of the model-based approach versus the  
251 classification-based approach. Mean square errors are obviously produced at an observation  
252 level, but here we wanted to map them to assess any spatial characteristics and areas where  
253 the model was and was not performing well. To do this the average mean square error in each  
254 habitat\*calcium carbonate category was calculated (or land class category) and this value  
255 mapped according to where that category is present over GB.

## 256 **Results**

257 In the model-based approach, parameter estimates and associated P values of the fixed effects  
258 show a high degree of dissimilarity amongst the factor levels of each of the category values  
259 (Table 1). High levels of calcium carbonate content are correlated with high values of the  
260 microbial community metric. This is consistent with the findings in Griffiths et al. (2011)  
261 who showed a positive relationship with the community metric and pH. Likewise, the acidic  
262 habitats, such as dwarf shrub heath, coniferous woodland and acid grassland, show low  
263 values for the community score, again consistent with findings of those authors.

264 After estimating all unknown parameters in the relationship between microbial community  
265 structure on one hand and land cover and calcium carbonate content on the other, and  
266 checking these parameters against expert knowledge gained from previously published  
267 results, predictions were obtained over Great Britain by substituting the full LCM and PMM  
268 data into the equation from the fitted model together with the spatial coordinates (Figure 2C).  
269 Similar models and maps were produced for the two sub-models which contain a single  
270 predictor variable each: land cover OR calcium carbonate content (Figures 2A-B). This  
271 separation enables a visual inspection of effects of each specific covariate and is a clear  
272 advantage over the classification-based approaches where it is fully unknown what is driving  
273 the spatial pattern and how. Although informative with regards to specific covariates, the  
274 model is a correlative assessment and any robust inference on drivers of change is  
275 confounded by the possible correlation between covariates included the model and missing  
276 ones. Care is therefore needed when interpreting the estimated relationships between the  
277 response and individual model terms.

278 As an interpretation of the maps presented in Figure 2, it appears that the land cover data  
279 enable separation of the response between the upland and lowland dominated habitats (Figure  
280 2A), a feature clearly visible in the Kriging-based map (Figure 1), whereas the calcium  
281 carbonate data allow separation of the lowland habitats into the alkaline and acidic soils

282 (Figure 2B). The maps produced also echo the findings in Griffiths et al., 2011 that both  
283 factors are required to adequately describe the spatial variation exhibited in microbial  
284 community structure (Figure 2C).

#### 285 Comparison with classification-based approach

286 The classification-based map, derived using the ITE land classification (Bunce et al., 1996),  
287 uses colours on the same gradient scale as the model-based results to indicate the estimated  
288 mean within each class (Figure 2D). Comparing the full model-based map (Figure 2C) to the  
289 map drawn using classification means (Figure 2D), shows that although the two maps look  
290 broadly similar, it is unknown what key components make up the soil microbial community  
291 structure and what drives the spatial segregation in the classification-based map.

292 The mean square errors from each of the mapping approaches are mapped with the darker  
293 colours representing a lower mean square error and hence better goodness of fit (Figure 3). It  
294 is clear that the modelled approach of using both land cover and parent material provides the  
295 best fit to the data. It also shows how the model-based approach is more informative, by  
296 examining the relative contribution of each variable as layers are included or discounted in  
297 the model. Integrating the mean square error over the whole area provides a simple single  
298 statistic assessment and shows that the model-based map using land cover and geology  
299 provides the best fit (lowest total mean square error of 10482.60 versus 22929.54 for the  
300 classification-based map). The classification-based map, however, still provides some  
301 information, indicating potential areas where it provides a better fit than the model-based  
302 approach. An example here would be around The Fens in East Anglia (highlighted by the red  
303 box). Thus it is clear the model-based approach may be missing an important driving variable  
304 (or any correlate of that missing driver) that represents the differing microbial community  
305 structure found in this area.

**306 Discussion**

307 Understanding spatial trends in natural capital indicators and their relationships with  
308 environmental conditions is vital in supporting evidence-based policy. The example  
309 presented demonstrated a procedure to facilitate this by modeling and subsequently mapping  
310 one particular indicator of natural capital known to have a significant impact on terrestrial  
311 ecosystem functioning. Though it is tempting to use these types of models to draw inference  
312 on drivers of change and the causal pathway behind the current state of natural capital, they  
313 can only identify potential environmental drivers and the variables that show a clear  
314 relationship to the response. This is because the models themselves represent a correlative  
315 assessment to establish relationships present in the observation data. To understand the role  
316 of mechanistic drivers, an assessment involving experiments and specifically designed long-  
317 term studies is necessary (Holland, 1986). However, if the sole purpose of the analysis is  
318 prediction, as spatial mapping is, rather than understanding drivers of change, then any  
319 confounding correlation between included covariates and missing covariates is not critical  
320 (Araújo and Guisan, 2006). The example used only two covariate datasets, however, it would  
321 be trivial to add further environmental variables such as climate or topographical features.  
322 This would increase the flexibility of the model-based approach and is likely to reduce the  
323 mean square error further across the geographic range.

324 Previous work in this area has often focused on the use of classification-based maps to  
325 provide a framework onto which one can express the value of natural capital. The results  
326 showed that the model-based map outperformed the classification approach. In our particular  
327 example this was perhaps not surprising - Griffiths et. al. 2011 had already demonstrated land  
328 cover was a key factor in microbial community response, and land cover is omitted in the  
329 classification of Bunce et al (1996). Classification maps are often developed without the  
330 inclusion of variables that may be subject to change over time. This is to ensure that the

331 geographic classification remains robust. Hence the exclusion of land cover occurs in many  
332 classifications, as it can be highly temporally variable. This example further highlights the  
333 issues surrounding naive use of existing classifications and why, given appropriate data, a  
334 model-based approach ought to be favoured.

335 The model-based approach to mapping natural capital presented, whilst extremely powerful  
336 and informative, relies heavily upon data with good spatial coverage, both in terms of the  
337 response one wishes to model and the variables with which to make prediction across a wider  
338 range of unsampled locations. It is therefore clear that coordinated, large scale, nationwide  
339 monitoring schemes such as the Countryside Survey (Norton et al., 2013), which play a  
340 pivotal role in providing source data on natural capital assets, should be maintained and  
341 exposed to inform policy decisions.

342 With increasing pressure on our natural assets from increasing human requirements and  
343 environmental change, there is an urgent need to provide better information for policy  
344 development and decision support. If we are to fully understand and value natural assets and  
345 ensure that they feed into decision-making, then it is important that we understand their  
346 distribution and trends in national extent and condition. Initiatives such as the Valuing Nature  
347 Network (VNN) and Natural Capital Committee (NCC) in the UK are government funded  
348 initiatives with the remit of ensuring that the national contribution of natural assets to a range  
349 of societal and economic benefits is well understood and helpfully informs decision making.  
350 This is done whilst balancing competing pressures and assessing the impact of different  
351 policy scenarios. Natural capital initiatives like the VNN and NCC also often seek to  
352 understand trade-offs and co-benefits across multiple environmental responses to help in  
353 conservation management, planning and resource distribution. We therefore anticipate that  
354 the powerful, information rich, model-based approach to understanding and mapping natural

355 capital will increase in use over the coming years as we seek to value our natural assets and  
356 predict landscape scale responses to change in environmental or policy drivers.

### 357 **Acknowledgements**

358 The authors gratefully acknowledge the support of Simon Smart, Russell Lawley and  
359 Matthew Harrison in the review and editing of this paper. We also express significant thanks  
360 in particular to the associate editor and also to three anonymous referees for their helpful  
361 comments and suggested edits that resulted in a significantly improved revision of this paper.  
362 This paper has been published as a result of research conducted by the CEH and BGS through  
363 their NERC directed research programmes. In addition, RIG was supported by the European  
364 Commission under the EcoFINDERS project (FP7-264465). The authors publish this article  
365 with the permission of the Executive Director, British Geological Survey (BGS), Natural  
366 Environment Research Council (NERC).

### 367 **References**

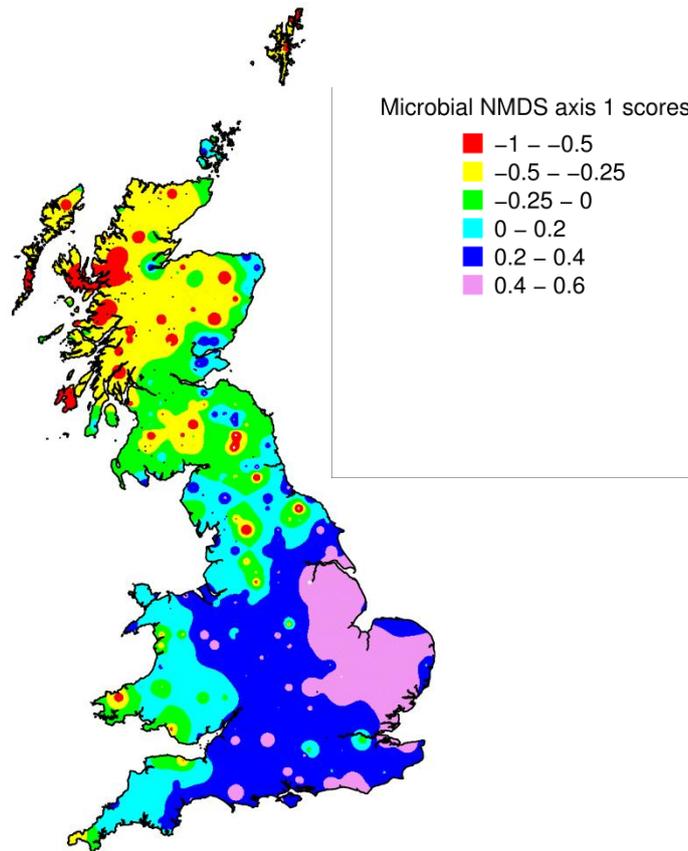
- 368 Algar, A. C., Kharouba, H. M., Young, E. R., & Kerr, J. T. (2009). Predicting the future of  
369 species diversity: macroecological theory, climate change, and direct tests of alternative  
370 forecasting methods. - *Ecography*, 32(1): 22-33.
- 371 Araújo, M. B., and Guisan, A. 2006. Five (or so) challenges for species distribution  
372 modelling. – *Journal of Biogeography*, 33(10): 1677-1688.
- 373 Bunce, R. G. H. et al. 1996. Land Classification for Strategic Ecological Survey. – *Journal of*  
374 *Environmental management* 47: 37-60.
- 375 British Geological Survey. 2013. Soil parent Material Model. Available online at:  
376 <http://www.bgs.ac.uk/products/onshore/soilPMM.html>. [Cited 07 February 2012].

- 377 Carpenter, S.R. et al. 2009. Science for managing ecosystem services: Beyond the  
378 Millennium Ecosystem Assessment. – Proceedings of the National Academy of Sciences of  
379 the United States of America. 106(5): 1305–1312.
- 380 Costanza, R., and Daly, H. E. 1992. Natural Capital and Sustainable Development. –  
381 Conservation Biology. 6 (1): 37-46.
- 382 Costanza, R., Wilson, M., Troy, A., Voinov, A., Liu, S., and D’Agostino, J. 2006. The value  
383 of New Jersey’s ecosystem services and natural capital. – Gund Institute for Ecological  
384 Economics, University of Vermont and New Jersey Department of Environmental Protection,  
385 Trenton, New Jersey, 13.
- 386 Diggle, P. J. et al. 1998. Model-based geostatistics (with discussion). – Journal of the Royal  
387 Statistical Society, Series C. 47( 3): 299-351.
- 388 Friedman, J. H. 1991. Multivariate adaptive regression splines. – The annals of statistics, 1-  
389 67.
- 390 Griffiths, R. I. et al. 2011. The bacterial biogeography of British soils. – Environmental  
391 Microbiology. 13 (6): 1642-1654.
- 392 Hastie, T. J. and Tibshirani, R. J. 1990. Generalized Additive Models. – Chapman and Hall,  
393 New York
- 394 Holland, P. W. 1986. Statistics and causal inference. – Journal of the American statistical  
395 Association, 81(396): 945-960.
- 396 Lawley, R. 2008. The soil-parent material database: A user guide. – British Geological  
397 Survey Internal Report .OR/08/034. 45pp.

- 398 Lin, X. and Zhang, D. 1999. Inference in generalized additive mixed models by using  
399 smoothing splines. – *Journal of the Royal Statistical Society, Series B.* 61: 381-400
- 400 McCullagh, P. and Nelder, J.A. 1989. *Generalized Linear Models*, 2nd edition. London:  
401 Chapman and Hall/CRC Press
- 402 Morton, D. et al. 2011. *Land Cover Map. 2007. Countryside Survey: Final Report for*  
403 *LCM2007 – the new UK Land Cover Map.* – Centre for Ecology and Hydrology. CS  
404 technical Report No 11/07.
- 405 Millennium Ecosystem Assessment. 2005. *Ecosystems and Human Well-being: Synthesis.* –  
406 Island Press. Washington DC.
- 407 Naidoo, R and Ricketts T. 2006. Mapping the economic costs and benefits of conservation .  
408 *PLoS Biology.* 4 (11): doi:10.1371/journal.pbio.0040360.
- 409 National Research Council. 2005. *Valuing ecosystem services: Toward better environmental*  
410 *decision-making.* National Academy Press, Washington, DC.
- 411 Nengwang, C. et al. 2009. A GIS-based approach for mapping direct use value of ecosystem  
412 services at a county scale: Management implications. – *Ecological Economics.* 68(11): 2768-  
413 2776.
- 414 Norgaard, R. B. 2009. Ecosystem services: From eye-opening metaphor to complexity  
415 blinder. – *Ecological Economics.* 69 (6): 1219-1227.
- 416 Norton, L. R., Inwood, H., Crowe, A., and Baker, A. 2012. Trialling a method to quantify the  
417 ‘cultural services’ of the English landscape using Countryside Survey data. – *Land Use*  
418 *Policy.* 29(2): 449-455.

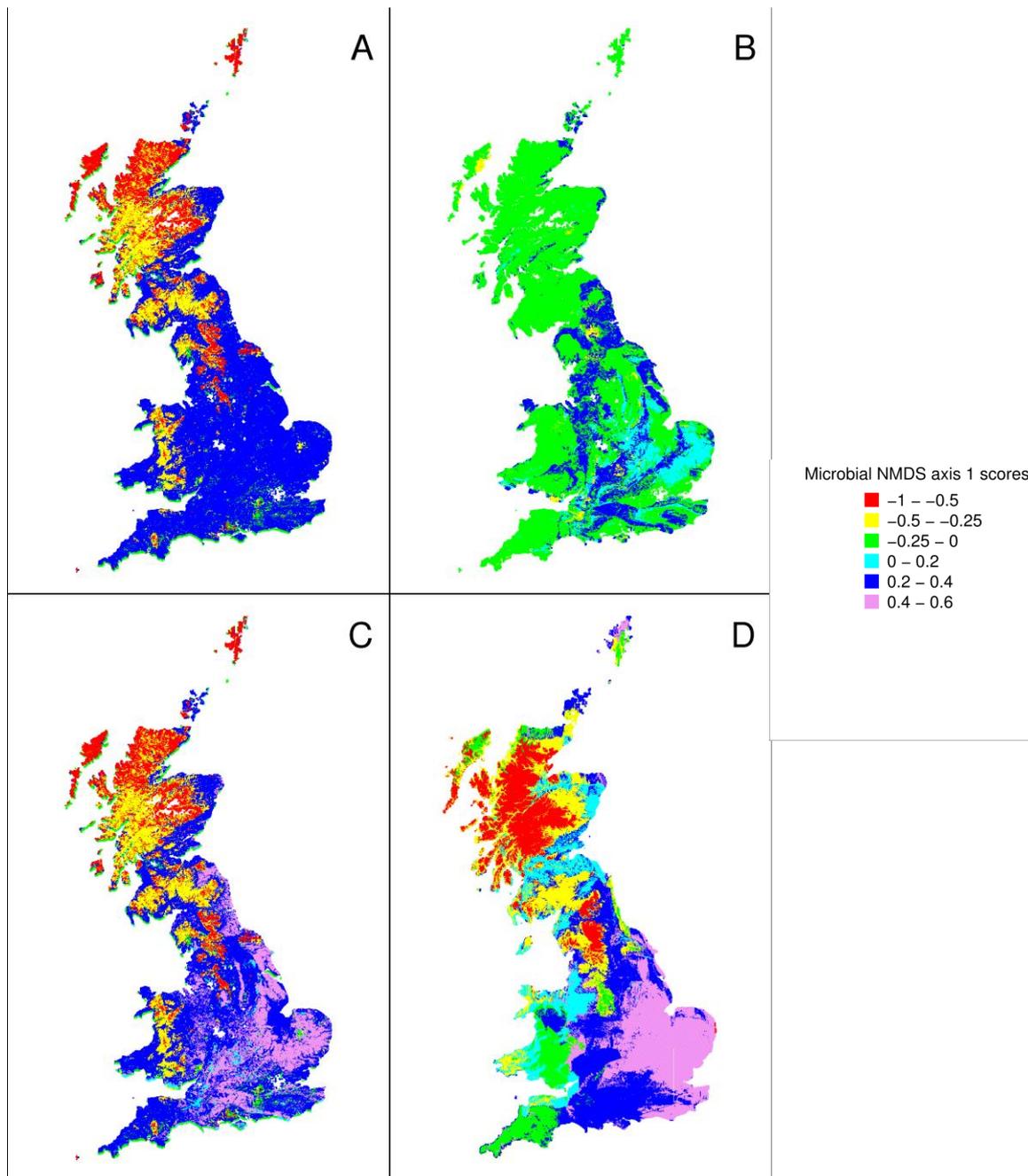
- 419 Norton, L. R. et al. 2013. Measuring stock and change in the GB countryside for policy - Key  
420 findings and developments from the Countryside Survey 2007 field survey. – Journal of  
421 Environmental Management. 113: 117-127.
- 422 R Development Core Team. 2008. R: A language and environment for statistical computing.  
423 – R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL  
424 <http://www.R-project.org>.
- 425 Raymond, C. M., Bryan, B. A., MacDonald, D. H., Cast, A., Strathearn, S., Grandgirard, A.,  
426 and Kalivas, T. 2009. Mapping community values for natural capital and ecosystem services.  
427 – Ecological economics. 68(5): 1301-1315.
- 428 Sayre, R. et al. 2009. A new map of standardized terrestrial ecosystems of the conterminous  
429 United States. – U.S. Geological Survey Professional Paper 1768, 17 p. (Also available  
430 online – [pubs.usgs.gov/pp/1768](http://pubs.usgs.gov/pp/1768).)
- 431 Swetnam, R. D. et al. 2011. Mapping socio-economic scenarios of land cover change: A GIS  
432 method to enable ecosystem service modelling. – Journal of Environmental Management. 92:  
433 563-574.
- 434 Troy, A. and Wilson, M. A. 2006. Mapping ecosystem services: practical challenges and  
435 opportunities in linking GIS and value transfer. – Ecological economics. 60(2): 435-449.
- 436 Turner, R. K., et al. 2010. An Introduction to socio-economic assessment with a marine  
437 strategy framework CSERGE, UEA – Published by the Department for Environment, Food  
438 and Rural Affairs.
- 439 The UK National Ecosystem Assessment. 2011. The UK National Ecosystem Assessment:  
440 Synthesis of the Key findings. – UNEP-WCMC, Cambridge.

- 441 Veloz, S. D., Williams, J. W., Blois, J. L., He, F., Otto-Bliesner, B., & Liu, Z. 2012.  
442 No-analog climates and shifting realized niches during the late quaternary: implications for  
443 21st-century predictions by species distribution models. – *Global Change Biology*, 18(5):  
444 1698-1713.
- 445 Williams, J. W. and Jackson, S. T. 2007. Novel Climates, No-Analog Plant Communities,  
446 and Ecological Surprises: Past and Future. – *Frontiers in Ecology and the Environment* 5:  
447 475-482
- 448 Wood, S.N. 2011. Fast stable restricted maximum likelihood and marginal likelihood  
449 estimation of semiparametric generalized linear models. – *Journal of the Royal Statistical*  
450 *Society (B)*. 73(1): 3-36
- 451



452

453 **Figure 1:** Map of soil microbial community structure (NMDS first axis scores) based on kriging of data  
454 obtained from the Countryside Survey – a stratified random sample of 591 1km survey squares located across  
455 the whole of Great Britain.



456

457 **Figure 2:** Maps of predictions in soil microbial community structure over Great Britain at 1km resolution,  
 458 showing comparisons among covariates of the model-based approach, and contrasting results of the  
 459 model-based and classification analyses. A - C using model-based approaches with covariates: A) land  
 460 cover only; B) calcium carbonate content only; and C) land cover and calcium carbonate content  
 461 combined. D estimating mean levels in each environmental stratum defined by the ITE land classification  
 462 of GB, then displaying on map using the spatial outline of each stratum.



463

464 **Figure 3:** Goodness of fit of the spatial statistical model used to derive the relationship between soil microbial  
465 community structure and environmental variables (land use and calcium carbonate content of the soil parent  
466 material). Map shows mean square error in each of the land use\*calcium carbonate classes. Darker shades  
467 indicate areas with low error.

468 **Table 1:** Estimated parameters and associated standard errors and p values resulting from the spatial statistical  
 469 model estimated defining the relationship between soil microbial community scores and environmental variables  
 470 (land use type and calcium carbonate class).

Parameter	Estimate	Standard Error	P Value
Intercept (CACO3 VARIABLE(LOW) * Bog)	-0.45	0.04	< 0.001
CACO3 HIGH	0.16	0.05	< 0.001
CACO3 LOW	-0.11	0.04	0.004
CACO3 MODERATE	0.04	0.14	0.786
CACO3 NONE	-0.10	0.03	< 0.001
CACO3 UNKNOWN	-0.25	0.14	0.078
CACO3 VARIABLE	-0.01	0.04	0.848
CACO3 VARIABLE(HIGH)	-0.31	0.10	0.002
Broadleaved, Mixed and Yew Woodland	0.46	0.04	< 0.001
Coniferous Woodland	0.10	0.04	0.011
Arable and Horticultural	0.87	0.03	< 0.001
Improved Grassland	0.77	0.03	< 0.001
Neutral Grassland	0.68	0.03	< 0.001
Calcareous Grassland	0.60	0.12	< 0.001
Acid Grassland	0.12	0.03	< 0.001
Bracken	0.24	0.07	0.001
Dwarf Shrub Heath	-0.01	0.04	0.818
Fen, Marsh, Swamp	0.46	0.06	< 0.001

471

472

## 473 **Supplementary Material**

### 474 **Appendix 1**

475 A model framework suitable for spatial modelling and mapping is the Generalised Linear  
 476 Geostatistical Model (GLGM) of Diggle et al. (1998). This framework can easily be extended  
 477 to a more generic setting where the linearity assumption is relaxed to form a Generalised  
 478 Additive Geostatistical Model (GAGM) following on from the Generalised Additive Model  
 479 framework (Hastie and Tibshirani 1990). Both model frameworks allow for the key  
 480 relationships to be estimated between the response of interest and the environmental  
 481 covariates, whilst at the same time controlling for additional spatial effects. This is because  
 482 observations close to one another are more likely to be similar than observations far away,  
 483 even after accounting for the environmental covariates in the model.

484 Spatial autocorrelation can be accounted for by including a purely spatial term in the model,  
 485 often a spatial random field, which captures any residual spatial variation in the data. This  
 486 ensures that parameter estimates and their associated standard errors are unaffected by any  
 487 residual spatial dependence. It also has the advantage that one can use the estimated spatial  
 488 correlation structure when making predictions, thus maximising the use of information, in an  
 489 approach similar to simple kriging. The underlying model framework of the GAGM  
 490 considered is presented below, where the geostatistical model consists of three parts: 1) a  
 491 linear combination of potentially smoothly varying covariate functions; 2) a spatial random  
 492 field, which we will define as a Stationary Gaussian Process (SGP); and 3) random effects  
 493 representing underlying, potentially non-spatial, error structure. Mathematically the model  
 494 framework is represented as

$$E[Y_i] = g\{\eta_i\}$$

$$495 \quad (1) \quad \eta_i = \alpha + \sum_{j=1}^k f_j(x_{ji}) + \mathbf{S}(u_i) + \mathbf{Z}_i \mathbf{b}$$

496 where  $Y$  is the response variable,  $f_j$  are smooth functions (generally cubic regression splines)  
497 of environmental covariates  $x_j$ ,  $g$  is the link function (as with standard GLMs),  $\alpha$  is the  
498 intercept term,  $\mathbf{Z}$  represents different grouping levels,  $b \sim N(0, \sigma)$  represents the differing  
499 variation assigned to each of the groups in  $\mathbf{Z}$  and  $S$  is a Stationary Gaussian Process at  
500 location  $u_i$  with zero mean and covariance structure given by  $\text{Cov}(u, u') = \sigma^2 \rho(\|u - u'\|)$ .

501

**502 Appendix 2**

503 As with all statistical modelling approaches it is more appropriate to start with a model  
504 consisting of a fixed effects formula dictated by scientific understanding and a simple error  
505 structure. Then, upon testing residuals and model assumptions, adapt the error structure as  
506 necessary. In this example we hence started with a simple GAM with land cover and calcium  
507 carbonate data as predictor variables together with a purely spatial interaction term of latitude  
508 and longitude to account for large scale spatial effects. Fitting a spatial trend surface is  
509 crucial to ensure adequate attribution of the response to the model covariates (Legendre and  
510 Fortin, 1989).

511 Upon examination of the residuals, it was clear that within square variance was not the same  
512 as the between square variation; hence the assumption of independence in the residuals was  
513 flawed. We therefore re-fitted the model with a random intercept effect to account for which  
514 CS 1km square the soil data were obtained from. This allowed for small scale random  
515 adjustments in the model. The residuals from the re-fitted model did not appear to imply any  
516 heteroscedasticity or any obvious key missing hierarchy in the error structure.

517 The residuals were then analysed for any small scale spatial autocorrelation. This was done  
518 using Moran's I, which showed no signs of small scale spatial autocorrelation apparent in the  
519 residuals. As this spatial autocorrelation was assessed on the residuals there was no need to  
520 include any disconnection when calculating Moran's I as any differences should have been  
521 accounted for in the main effects. Previous studies (eg Franklin and Mills, 2003) have shown  
522 spatial autocorrelation of soil microbial community data is evident at distances of up to 7  
523 metres. CS squares are separated by a minimum of 15 km and within square observations are  
524 separated by a minimum of 80 metres with an average separation distance of 558 metres.  
525 Given this, and the results of Franklin and Mills, the redundancy of fine scale spatial  
526 autocorrelation in the model is perhaps not surprising.

527 We therefore modelled the raw data on soil microbial communities against broad habitat and  
 528 calcium carbonate content using a generalised additive mixed-model based approach. This  
 529 follows the same generic approach as the GAGM without the inclusion of a spatial random  
 530 field. Generalised additive mixed models (Lin and Zhang, 1999) extend the framework of the  
 531 standard GAM by allowing both fixed and random affects to be present in the model. The  
 532 random components can account for unobserved affects that could influence the outcome of  
 533 the response variable and therefore ensure that estimated standard errors are accurate and any  
 534 inference is reliable. Extending the general GAM equation to include random effects gives us  
 535 a model of the following form:

$$g(E[Y_i|\mathbf{x}, \mathbf{b}]) = \alpha + \sum_{j=1}^k f_j(x_{ji}) + \mathbf{Z}_i \mathbf{b}$$

536 where  $y$  is the response variable,  $f_j$  are smooth functions (generally cubic regression splines),  
 537  $g$  is the link function (as with standard GLMs),  $\alpha$  is the intercept term,  $Z$  represents different  
 538 grouping levels and  $\mathbf{b} \sim N(0, \sigma)$  represents the differing variation assigned to each of the  
 539 groups in  $Z$ .

540 The random components are used here to allow us to account for the fact that any two soil  
 541 cores taken from the same 1km square are more likely to be similar than two cores taken  
 542 from two different squares. The non-linear smooth form allows fitting of an additional  
 543 smoothly varying spatial surface to soak up any residual large scale spatial variation and  
 544 hence captures the spatial structure present in the data that our covariates may not adequately  
 545 explain. This is akin to including time as a covariate in time series modelling – the user is  
 546 effectively de-trending the data. Even in the absence of small scale spatial autocorrelation,  
 547 Legendre and Fortin (1989) emphasised the importance of including this term. Including the  
 548 random effects, additional spatial surface and the habitat and calcium carbonate covariate  
 549 effects, the fitted model is thus represented by

$$g\{\mathbf{E}[smc_{i,s}]\} = \alpha + \beta_{h_i} + \eta_{c_i} + f(\text{Latitude}_i, \text{Longitude}_i) + \omega_s + \sigma_i$$

550  
551 where for each observation  $i$  in square  $s$ ,  $smc$  is the soil microbial community score,  $\beta_h$  is the  
552 estimated value of habitat  $h$  associated with observation  $i$ ,  $\eta_c$  represents the value for calcium  
553 carbonate category  $c$ ,  $\omega$  represents the error (normally distributed) associated specifically  
554 with square  $s$  and  $\sigma$  represents the residual model error also assumed to follow a normal  
555 distribution. The model was fitted, including the smoothly varying spatial surface using  
556 tensor product smooth interactions, using the `gamm` function in the ‘`mgcv`’ library (Wood,  
557 2011) in the R statistical environment (R Development Core Team, 2008).

558 Had the re-fitted model failed the independence assumptions and the Moran’s I test showed  
559 evidence for fine scale spatial autocorrelation, then the inclusion of the spatial random field  
560 term in the model would have been necessary. Practically, the Gaussian Random Field (GRF)  
561 is often estimated by making the assumption that it is adequately specified by a Markov  
562 Random Field (MRF) whereby each location only depends on its “neighbours” and is  
563 conditionally independent of all other locations. The neighbourhood structure of the MRF  
564 allows the spatial component of the model to be estimated by methods such as Conditional  
565 Autoregressive Models (CAR) or Simultaneous Autoregressive Models (SAR). Dormann et  
566 al (2007) provide an overview of methods for accounting for spatial autocorrelation including  
567 description of CAR and SAR models and how to fit them in practice with clearly referenced  
568 R packages.

569 It is worth noting that both CAR and SAR models can also be estimated in a Bayesian  
570 framework, where estimated parameters and standard errors are often more reliable than in  
571 likelihood approximation methods, though with an added computational cost. The advantage  
572 is the added flexibility that moving to the Bayesian paradigm brings. Specifically in this case  
573 the possible inclusion of smoothly varying penalised regression splines following the

574 approach taken by Crainiceanu et. al. (2005). This provides the full ability to fit the model  
575 specified in Eqn (A1). This type of model can also be easily fitted using Integrated Nested  
576 Laplace Approximation (Rue et. al., 2009), where robust parameter estimates can be obtained  
577 quickly and efficiently. The R package R-Inla ([www.r-inla.org](http://www.r-inla.org)) is a user friendly resource for  
578 fitting the model in Eqn (A1) using this approach.

579

580 **References**

- 581 Crainiceanu, C. M., Ruppert, D. And Wand, M. P. 2005. Bayesian Analysis for Penalized  
582 Spline Regression Using WinBUGS. – Journal of Statistical Software. 14
- 583 Diggle, P. J. et al. 1998. Model-based geostatistics (with discussion). – Journal of the Royal  
584 Statistical Society, Series C. 47( 3): 299-351.
- 585 Dormann, C. F. et al. 2007. Methods to account for spatial autocorrelation in the analysis of  
586 species distributional data: a review. – Ecography. 30: 609-628.
- 587 Franklin, R. B. and Mills, A. L. 2003. Multi-scale variation in spatial heterogeneity for  
588 microbial community structure in an eastern Virginia agricultural field. – FEMS  
589 Microbiology Ecology, 44: 335–346.
- 590 Hastie, T. J. and Tibshirani, R. J. 1990. Generalized Additive Models. – Chapman and Hall,  
591 New York
- 592 Fortin, M. J., Drapeau, P., & Legendre, P. 1989. Spatial autocorrelation and sampling design  
593 in plant ecology. – Vegetatio, 83(1-2): 209-222.
- 594 Lin, X. and Zhang, D. 1999. Inference in generalized additive mixed models by using  
595 smoothing splines. – Journal of the Royal Statistical Society, Series B. 61: 381-400
- 596 R Development Core Team. 2008. R: A language and environment for statistical computing.  
597 – R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL  
598 <http://www.R-project.org>.
- 599 Rue, H., Martino, S. & Chopin, N. 2009. Approximate Bayesian inference for latent Gaussian  
600 models by using integrated nested Laplace approximations. – Journal of the Royal Statistical  
601 Society, Series B. 71(2): 319-392.

602 Wood, S.N. 2011. Fast stable restricted maximum likelihood and marginal likelihood  
603 estimation of semiparametric generalized linear models. – Journal of the Royal Statistical  
604 Society (B). 73(1): 3-36

605

606

607

608