

1 **Electronic supplement**

2

3 **APPENDIX: DETAILED DESCRIPTION OF METHODS**

4 **Environmental data**

5 I used 45 environmental variables to explain differences in species richness among grid cells
6 with a resolution (i.e. cell size) of 10' longitude × 6' latitude, i.e. approx. 130 km², following
7 the 1:25,000 ordnance survey maps for Germany. These parameters were transformed from
8 digital maps with polygon-topology to the grid by intersection and exported into the database
9 format that I used.

10 *Spatial Coordinates*

11 I used South and East enumeration of the 1:25,000 ordnance survey maps as southing and
12 easting of the grid cells for incorporating spatial autocorrelation and defining neighbourhood
13 matrices.

14 *Land Cover*

15 Corine Land Cover (CLC) data, provided by the “Statistisches Bundesamt” (1997), was used
16 to calculate the following variables per grid cell: number of patches, average patch size,
17 variation coefficient of patch size, number of different types, number of aggregated types
18 (aggregated types are urban area, agricultural area, forests and near-nature area, wetland area,
19 water surface area).

20 *Soil data*

21 Soil data was taken from the soil survey map ('Bodenübersichtskarte') 1:1,000,000 of the
22 “Bundesanstalt für Geowissenschaften und Rohstoffe (BGR)” (<http://www.bgr.de>). I
23 calculated number of patches, average patch size, variation coefficient of patch size, number

24 of types, and number of aggregated types per grid cell. The soil types were aggregated to the
25 following classes: soils of coasts and bogs, soils of floodplains and valleys, soils of lowlands,
26 soils of loess landscapes, soils of low mountain, soils of high mountain, anthropogenic soils.

27 *Geological data*

28 For geological data, I used the geological survey map 1:1,000,000 ('Geologische
29 Übersichtskarte') of the "Bundesanstalt für Geowissenschaften und Rohstoffe (BGR)"
30 (<http://www.bgr.de>). For each grid cell, I calculated number of patches, average patch size,
31 variation coefficient of patch size, number of geological types, number of aggregated types
32 (aggregated types are lime, sand, loess, clay, others).

33 *Climate data*

34 Climate data on a 1 km² grid scale were provided by the "Deutscher Wetterdienst, Dept.
35 Klima und Umwelt". The recording periods were 1951-1980 for the temperature data and
36 1961-1990 for the precipitation data. I calculated grid cells averages and coefficients of
37 variation for mean January temperature, mean July temperature, mean annual temperature,
38 mean annual precipitation, and for the difference between mean July temperature and mean
39 January temperature.

40 *Altitude*

41 Averages and coefficients of variation of altitude per grid cell were calculated after the
42 ARCDDeutschland500 dataset, scale 1:500.000, provided by ESRI.

43

44 The scales of the maps were different among these data but grain and extent for this analysis
45 were the same for all data classes. The grain of this analysis (i.e. 10' longitude x 6' latitude)
46 was much larger than the resolution of even the coarsest scale map used, therefore the
47 information per grid cell should be sufficiently detailed. As the different parameters are

48 transformed into principal components, differences in mapping scale should matter even less.
49 However, finer resolution information on environmental variables might lead to reduced noise
50 in the models.

51

52 **Plant distribution data**

53 Data on plant distributions came from the FLORKART database of the German Centre for
54 Phytodiversity at the German Federal Agency for Nature Conservation (see
55 www.floraweb.de). This database was collated from current and former regional floristic
56 mapping projects, mainly based on the field work of thousands of volunteers. I used data with
57 a resolution of 10' longitude and 6' latitude.

58

59 As mapping was organised decentrally, mapping intensity proved to be heterogeneous
60 throughout Germany. To reduce this bias, I considered mapping intensity by designating 50
61 control species, all of which had to be present in order to include a grid cell in the analysis.
62 These control species are ubiquitous and assumed to occur in every grid cell. Grid cells that
63 lack any of the control species were regarded as not sufficiently sampled and excluded from
64 the data set. 45 of the control species were the most ubiquitous species in Germany according
65 to Krause (1998) and five are additional ubiquitous species which are either inconspicuous or
66 difficult to determine (see Kühn *et al.* 2004). This left 1928 of the 2995 grid cells of Germany
67 for analysis. Native plant species were identified after Kühn & Klotz (2002) yielding 2411
68 species for analysis. Species numbers were log-transformed to achieve normality.

69

70 **Statistical Analysis**

71 The 45 environmental variables were transformed in a principal component analysis on the
72 correlation matrix to condense most of the environmental variance on the first few principal

73 components (PCs). I only used the first four PCs, which explained roughly 41% of the
74 variation in the dataset.

75

76 I used Ordinary Least Squares (OLS) regression as non-spatial model to relate log-
77 transformed species richness to predictors. OLS regression can be written in matrix notation
78 as $y = \beta X + \varepsilon$ where y is the vector of observation, β is the vector of regression coefficients
79 (including the intercept) and X is the matrix of explanatory variables. ε is the normally
80 distributed error (i.e. the residuals) with $\varepsilon \sim N(0, V)$ where the variance-covariance matrix
81 $V = \sigma^2 I$ with σ^2 as the variance and I as the identity matrix (i.e. ones as diagonal elements and
82 zeros for all off-diagonal elements).

83

84 In a conditionally autoregressive (CAR) model, the variance-covariance matrix is defined as
85 $V = (I - \rho W)^{-1} M$ where ρ is a spatial autocorrelation coefficient, W as a matrix of $n \times n$ spatial
86 weights derived from a defined neighbourhood and M as an $n \times n$ matrix with the conditional
87 variances ($\sigma_1^2, \dots, \sigma_n^2$) of y (i.e., the variances of y given the realized values of the spatial
88 neighbours) on the diagonal and zeros in the off-diagonal positions (Lichstein *et al.* 2002,
89 Haining 2003). In the present analysis, I assumed that the conditional variances of y were
90 constant, i.e. $M = I\sigma^2$.

91 In Simultaneous autoregressive (SAR) models we have $V = \sigma^2 [(I - \rho W)'(I - \rho W)]^{-1}$. SAR models
92 have several ways to incorporate spatial autocorrelation (SAC): The error model (ESAR)
93 corrects for SAC in the error term ($y = \beta X + \rho W \xi + \varepsilon$, where ξ is the autocorrelated error and
94 ε is the uncorrelated error), the lag model corrects for SAC in the response variable
95 ($y = \rho W y + \beta X + \varepsilon$) and the spatial Durbin model (or mixed autoregressive model,
96 $y = \rho W y + \beta X - \rho W \beta X + \varepsilon$) combines both error and lag model (Anselin 1988). To my
97 knowledge there is currently no sufficient and consistent ecological theory to provide a basis

98 for decided which of these autoregressive models to (but some hints are given in econometric
99 literature, Anselin & Bera 1998). I chose a data driven approach (Haining 2003) to find the
100 model which provided the best fit to the data (as measured by AIC) and which was most
101 effective in the removal of spatial autocorrelation (as measured by Moran's *I*).

102

103 Crucial in all these models is the specification of the neighbourhood which defines the local
104 zone of influence. I defined several neighbourhoods which include all grid cells within a
105 distance of 1.5, 2 and 2.9 cells Euclidean distance (i.e. all 8 adjacent neighbours, the 12
106 neighbours including the second nearest one, and 24 nearest neighbours including the closest
107 third order neighbours around a focus grid cell) respectively. I chose these distance classes as
108 they stepwise included larger neighbourhoods. Neighbourhood weights were row
109 standardized for SAR models as recommended by Anselin (1988) and binary to provide a
110 symmetric neighbourhood matrix as necessary for CAR models.

111

112 I included all four PCs and their second and third order polynomials as predictors. However,
113 in the spatial model only PC4² remained significant. Since this was not relevant for the main
114 message of this article, the result is not shown.

115

116 Legendre et al. (2002) distinguish between spatial dependence (or spatial structure) and
117 spatial autocorrelation. The former implies that a response variable is structured because it
118 depends upon explanatory variables that are themselves spatially structured. However, as the
119 failure to include important (spatially structured) may also result in spatial autocorrelation
120 (Cliff & Ord 1981, Haining 2003), both concepts can sometimes be linked. The removal of
121 spatial structure, however, does not necessarily remove spatial autocorrelation. I used this
122 concept in the analysis to remove large-scale trends from my dataset. Therefore, I used a third
123 order polynomial trend surface regression (Legendre & Legendre 1998) i.e., I used an OLS

124 regression to explain $\log(\text{species richness})$ as a function of *southing* + *easting* +
125 *southing***easting* + *southing*² + *easting*² + *southing*²**easting* + *southing***easting*² + *southing*³
126 + *easting*³. All higher-order parameters except *southing*² were significant, which was then
127 removed. The residuals of this model were used as response in a regression on the four PCs
128 (results see table S2). Although originally intended as a method to partial out the spatial
129 structure of a pattern, it is appropriate in this context: At the spatial scale of this analysis, it is
130 highly unlikely that local dispersal processes etc. will lead to a spatially structured species
131 richness pattern per se. Therefore, large-scale environmental parameters will be most
132 important for plant distribution patterns. A spatial gradient will thus integrate across several
133 large-scale environmental gradients not represented in the variables which I used for the
134 principal component analysis. It is hence a suitable way to account for unknown large-scale,
135 spatially structured environmental gradients.

136

137 The fit of different spatial and the non-spatial model to the four PCs was compared using AIC
138 (Akaike's Information Criterion, $AIC = -2LL + 2n$ where *LL* is the log-likelihood of the model
139 and *n* is the numbers of parameters in the fitted model) (Quinn & Keough 2002). Since the
140 four PCs are orthogonal to each other (thus avoiding collinearity problems often associated
141 with model selection procedures), I used error-probabilities (p-values) within a method to
142 assess the importance of covariates, which was much faster than model simplification and
143 calculation of AIC. I also calculated a pseudo- $R^2 = 1 - D_A/D_0$ with D_A as the deviance of the
144 model of interest and D_0 as the deviance of the non-spatial intercept-only model. For
145 Gaussian distributed errors, the deviance is the sum of squared residuals (therefore, for OLS
146 with variance=deviance, pseudo- $R^2 = R^2$). I used Moran's *I* correlograms (Legendre &
147 Legendre 1998) to evaluate the amount of spatial autocorrelation of the residuals. Moran's *I* is
148 an autocorrelation coefficient and could be regarded as spatial equivalent to Pearson's
149 correlation coefficient. Significance was assessed after 1000 permutations.

150

151 The error-model (ESAR) with a neighbourhood distance of 2 grid cells yielded the best fit

152 (AIC) and was the only one that reduced spatial autocorrelation to a non-significant amount.

153

154 I did all calculations in R (R Development Core Team 2005) using functions 'spautolm' for

155 CAR models and 'errrorsarlm' and 'lagssarlm' for SAR models from package SPDEP (Bivand

156 *et al.* 2005) and ncf (Bjørnstad 2004) for spatial correlograms.

157

158

REFERENCES

- Anselin, L. (1988) *Spatial econometrics: methods and models*. Kluwer, Dordrecht.
- Anselin, L. & Bera, A.K. (1998) Spatial dependence in linear regression models with an introduction to spatial econometrics. *Handbook of applied economic statistics* (ed. by A. Ullah and D.E.A. Giles), pp. 237-289. Marcel Dekker, New York.
- Bivand, R., Anselin, L., Bernat, A., Carvalho, M. M., Chun, Y., Dormann, C., Dray, S., Halbersma, R., Lewin-Koh, N., Ono, H., Tiefelsdorf, M. & Yu, D. (2005) *spdep: Spatial dependence: weighting schemes, statistics and models*. R package version 0.3-17.
- Bjørnstad, O. N. (2004) *ncf: spatial nonparametric covariance functions*. R package version 1.0-6. <http://onb.ent.psu.edu/onb1/R>.
- Cliff, A.D. & Ord, J.K. (1981) *Spatial Processes: Models and Applications*. Pion, London.
- Dormann, C.F. (2006) Effects of incorporating spatial autocorrelation into the analysis of species distribution data. *Global Change Biology* **in press**.
- Haining, R.P. (2003) *Spatial data analysis: Theory and practice*. Cambridge University Press, Cambridge.
- Krause, A. (1998) Floras Alltagskleid oder Deutschlands 100 häufigste Pflanzenarten. *Natur und Landschaft* **73**, 486-491.
- Kühn, I., Brandl, R., & Klotz, S. (2004) The flora of German cities is naturally species rich. *Evolutionary Ecology Research* **6**, 749-764.
- Kühn, I. & Klotz, S. (2002) Floristischer Status und gebietsfremde Arten. *BIOLFLOR – Eine Datenbank zu biologisch-ökologischen Merkmalen der Gefäßpflanzen in Deutschland* (ed. by S. Klotz, I. Kühn and W. Durka). *Schriftenreihe für Vegetationskunde* **38**, pp. 47-56. Bundesamt für Naturschutz, Bonn.
- Legendre, P., Dale, M.R.T., Fortin, M.J., Gurevitch, J., Hohn, M., & Myers, D. (2002) The consequences of spatial structure for the design and analysis of ecological field surveys. *Ecography* **25**, 601-615.
- Legendre, P. & Legendre, L. (1998) *Numerical Ecology*. Developments in environmental modelling **20**. Elsevier, Amsterdam.
- Lichstein, J.W., Simons, T.R., Shiner, S.A., & Franzreb, K.E. (2002) Spatial autocorrelation and autoregressive models in ecology. *Ecological Monographs* **72**, 445-463.
- Quinn, G.P. & Keough, M.J. (2002) *Experimental design and data analysis for biologists*. Cambridge University Press, Cambridge.
- R Development Core Team (2005) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.

Short communication:

Incorporating spatial autocorrelation may invert observed patterns (electronic supplement)

Statistisches Bundesamt (1997). *Daten zur Bodenbedeckung für die Bundesrepublik Deutschland 1:100.000*. Statistisches Bundesamt, Wiesbaden.

Table S1 Loadings of environmental variables on the first four dimensions of a Principal Component analysis on 45 environmental variables across Germany.

	PC1	PC2	PC3	PC4
avg July temperature	-0.64	0.01	0.39	-0.11
cv of July temperature	0.63	0.51	0.02	0.02
avg January temperature	-0.54	-0.06	0.57	0.26
cv of January temperature	-0.05	-0.05	0.02	0.06
temperature difference July-January	0.09	0.07	-0.31	-0.35
avg mean annual temperature	-0.66	-0.04	0.57	0.10
cv of mean annual temperature	0.64	0.39	0.00	-0.02
avg annual p	0.65	0.23	0.19	0.07
cv of annual precipitation	0.48	0.49	0.00	0.07
avg altitude	0.81	0.35	-0.06	-0.01
cv of altitude	-0.02	0.13	-0.03	0.06
number of CLC patches	0.31	0.33	0.42	-0.30
number of CLC types	-0.35	0.20	0.45	-0.38
avg size of CLC patches	-0.16	-0.30	-0.46	0.31
cv of size of CLC patches	-0.19	-0.20	-0.32	0.38
number of aggregated CLC types	-0.31	-0.25	-0.13	-0.42
agricultural area	-0.28	-0.18	-0.30	0.38
wetlands area	-0.14	-0.20	-0.13	0.04
urbanised area	-0.25	0.26	0.59	-0.07
forest or near-nature area	0.58	0.19	-0.01	-0.42
watersurface area	-0.12	-0.19	-0.07	-0.28
number of geological patches	-0.20	0.45	-0.26	-0.24
number of geological types	-0.26	0.67	-0.19	-0.01
avg size of geological patches	0.30	-0.34	0.30	0.21
cv of size of geological patches	0.17	-0.05	-0.09	-0.23
number of aggregated geological types	-0.33	0.50	0.11	0.34
area of lime subsoil	0.14	0.33	0.02	0.18
area of loess subsoil	-0.19	0.26	0.20	0.49
area of sandy subsoil	-0.32	-0.11	0.24	-0.41
area of clay subsoil	-0.15	-0.05	0.04	0.24
area of other subsoil	0.44	-0.04	-0.30	0.06
number of soil patches	-0.54	0.49	-0.30	-0.08
number of soil types	-0.56	0.58	-0.30	-0.08
number of aggregated soil types	-0.65	0.42	-0.13	-0.05
number of natural soil types	-0.52	0.56	-0.37	-0.07
number of natural aggregated soil types	-0.60	0.39	-0.26	-0.03
avg size of soil patches	0.55	-0.33	0.26	0.09
cv of size of soil patches	-0.11	0.08	-0.12	-0.12
area of anthropogenic soils	-0.22	0.12	0.40	-0.12
area of soils of high mountains	0.07	0.07	-0.04	0.00
area of coastal soils	-0.29	-0.36	-0.29	0.07
area of soils of loess landscape	-0.13	0.38	0.13	0.59
area of soils of low mountains	0.77	0.35	0.05	-0.01
area of soils of planes	-0.26	-0.47	-0.33	-0.32
area of soils of valleys and floodplains	-0.48	-0.05	0.28	-0.21
Eigenvalues	7.88	4.59	3.43	2.60
Percentage variance	17.50	10.20	7.62	5.78
Cumulative percentage variance	17.50	27.71	35.33	41.11

Short communication:

Incorporating spatial autocorrelation may invert observed patterns (electronic supplement)

Table S2 Results of an Ordinary Least Square regression of log-transformed species richness in Germany on four environmental Principal Components after statistically controlling for large-scale spatial gradients (i.e., using residuals of a third order polynomial trend surface regression).

	Estimate	Std. Error	t value	p
Intercept	-0.00	0.001	0.00	1.000
PC1, Altitude	-0.16	0.023	-7.12	<0.001
PC2, Geodiversity	0.39	0.03	13.12	<0.001
PC3, Urbanization	-0.04	0.034	-1.3	0.2
PC4, Loess	-0.22	0.039	-5.66	<0.001
