

GIS, SPATIAL STATISTICAL GRAPHICS, AND FOREST HEALTH.

James J. Majure, Noel Cressie, Dianne Cook, and Jürgen Symanzik

ABSTRACT

This paper discusses the use of a geographic information systems (GIS), Arcview 2.1, linked with a dynamic graphics program, XGobi, in the statistical analysis of spatial data. The link allows multivariate data, collected at geographic locations and stored in Arcview, to be passed into XGobi and analyzed dynamically. The connection between the points in XGobi and the spatial locations from which they were collected is maintained so that points in either Arcview or XGobi can be brushed and the corresponding points in the other application identified immediately. Spatial cumulative distribution functions (SCDFs), spatially lagged scatter plots, and variogram cloud plots can be displayed in XGobi using the link. In each type of plot, the connection to the spatial sampling location is maintained and user interaction can take place in either application.

The link is used to predict and analyze SCDFs of forest crown health in the northeastern United States. The SCDFs are predicted from field data collected as part of the U.S. Environmental Protection Agency's (USEPA) Environmental Monitoring and Assessment Program (EMAP). The field data are augmented with concomitant geographic information, including Landsat Thematic Mapper images, digital elevation models, and population information, which are used to improve the SCDF prediction.

INTRODUCTION

This paper discusses the integration of a dynamic graphics program, XGobi, into a geographic information system (GIS), Arcview 2.1 (ESRI 1995), and its use in the statistical analysis of spatial data. The link between XGobi and Arcview allows multivariate data, collected at geographic locations and stored in Arcview to be passed into XGobi and viewed. The connection between the points in XGobi and the spatial locations from which they were collected is maintained so that points in either XGobi or Arcview can be *brushed* (see Note 1 at the end of the paper), resulting in simultaneous brushing of the corresponding points in the other application. The link also has the ability to use XGobi to display spatial cumulative distribution functions (SCDFs), spatially lagged scatter plots, and variogram-cloud plots. The connection to the spatial sampling locations is maintained in each type of plot and user interaction can take place in either application.

The particular problem to which these tools are applied involves the spatial prediction and analysis of SCDFs for forest crown health in the northeastern United States. The SCDFs are predicted from field data collected as part of the U.S. Environmental Protection Agency's (USEPA) Environmental Monitoring and Assessment Program (EMAP). In addition to the field data, concomitant geographic information, including Landsat Thematic Mapper images, digital elevation models, and population information are included in the analysis. This additional information has the potential to improve the SCDF prediction.

In this paper, we will first give an overview of the linking technology between Arcview and XGobi. We will then discuss the use of the link in the prediction of SCDFs.

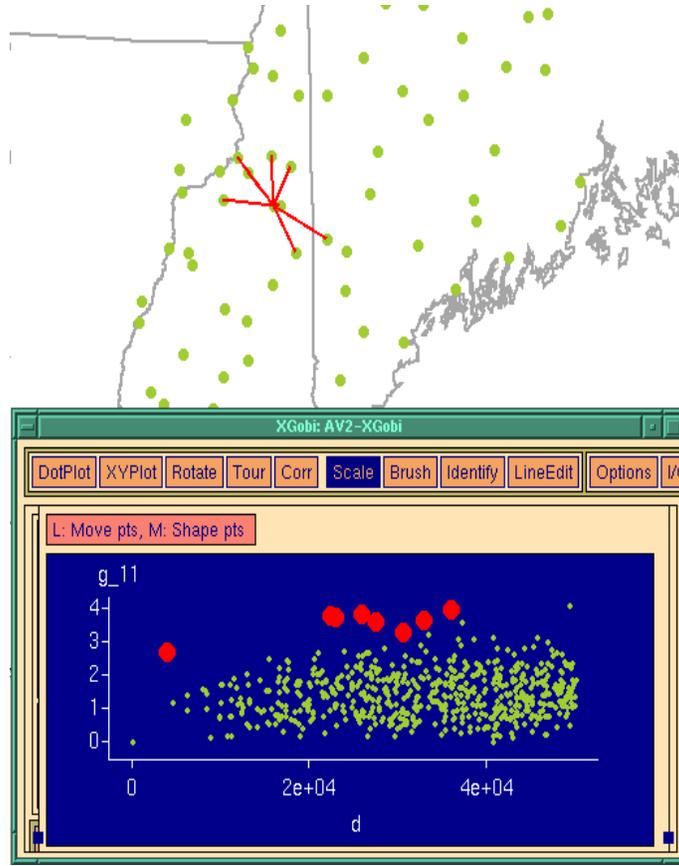


Figure 1. A variogram-cloud plot (bottom of figure) with large values brushed. The map view (top of figure) indicates that all brushed points have a common sampling location.

INTEGRATION OF DYNAMIC GRAPHICS TOOLS INTO A GIS

Interactive and dynamic graphics programs are very useful in the exploration of high-dimensional data. With data collected at spatial locations, it is important to include the locations as part of the analysis. This leads very naturally to the integration of a GIS with a dynamic graphics program; the GIS is used for displaying spatial locations and concomitant geographic variables, and the dynamic graphics program is used for visualizing and exploring the corresponding data space. This type of link has been constructed between Arcview 2.1 and XGobi (Swayne et al. 1991), an interactive dynamic graphics program in the X Window SystemTM environment. Technical details of the link can be found in Symanzik et al. (1995) and Majure et al. (1995).

The link between Arcview and XGobi is intended to provide functionality that is not provided by either the GIS or the dynamic graphics program alone. While GISs provide sophisticated capabilities for the input of spatial data, its management, and the display of maps, graphics and tables, their capability for statistical analysis is generally limited and dynamic graphical analysis is non-existent. Although most dynamic graphics programs can plot the coordinates of spatial locations, they do not have the capabilities of producing high quality maps that provide a geographic frame of reference. Together, then, Arcview and XGobi share their strengths and produce a product that is more than the sum of the parts.

The specific tools made available by the link include the resident capabilities of both Arcview and XGobi, as well as the ability to do *linked brushing* (see Note 1 at the end of the paper) between the two systems. The capabilities of Arcview 2.1 include the display and manipulation of sample locations and other geographic information. XGobi provides an array of graphic options through the manipulation of scatter plots. The types of plots available include univariate and bivariate plots, three-dimensional point rotation, and higher-dimensional rotation with the grand tour (Asimov 1985, Buja and Asimov 1986) and the correlation tour (Buja et al. 1988). Both the grand tour and correlation tour allow rotation toward “interesting” projections of the data through projection pursuit (Cook et al. 1993). The link between the two programs allows the analyst to brush points, in either Arcview or XGobi, with a color/size/glyph and to see where the corresponding points are located in the other application. Thus, outliers in an XGobi plot can be brushed to see (in Arcview) where they were collected, or a spatial region in Arcview can be brushed to see (in XGobi) where the corresponding attribute measurements fall in the data space. Together, these tools provide a powerful and flexible environment for the graphical analysis of spatial data.

In addition to these basic capabilities, the link has been extended to include the display and analysis of SCDFs, spatially lagged scatter plots (Cressie 1993) and variogram clouds (Haslett et al. 1991, Bradley and Haslett 1992). In these cases, the data being passed from the GIS is processed before being displayed in XGobi. An explanation and examples of the SCDF link are given in the next section. The variogram-cloud link is used when exploring the spatial dependence in a data set and when looking for spatial outliers. In this option, the points displayed in XGobi represent all possible pairs of sampling locations. For each pair of locations, XGobi plots the square-root of the absolute difference between attribute values at the locations versus the Euclidean distance between the locations. In data sets exhibiting strong spatial dependence, the variance in the attribute differences will increase with increasing distance between locations. Locations that are near to one another, but with large attribute differences, might indicate a spatial outlier, even though the values at both locations may appear to be reasonable when examining the data set non-spatially.

Figure 1 shows a variogram-cloud plot for precipitation sampling stations in which several potentially outlying points have been brushed. Because each point in the XGobi window corresponds to a pair of sampling locations, when the points in XGobi are brushed the Arcview window shows each pair of sampling locations connected by a line. This is also shown in Figure 1. Notice that all of the outlying points have a single sampling location in common. When the Arcview window is displayed with elevation contours, it is immediately obvious that the location in question is located on top of a mountain, which accounts for the large difference in precipitation.

PREDICTION OF THE SCDF FOR TREE CROWN HEALTH

In this section, the link described previously will be applied to the spatial prediction and visualization of the SCDF for the crown defoliation index (CDI) (Anderson et al. 1992), calculated from data collected in the northeastern United States. The CDI represents the nature of tree crown health as a response to stressors. In this analysis, the SCDF for the CDI process is predicted from data collected from a probability-based sample. Further-

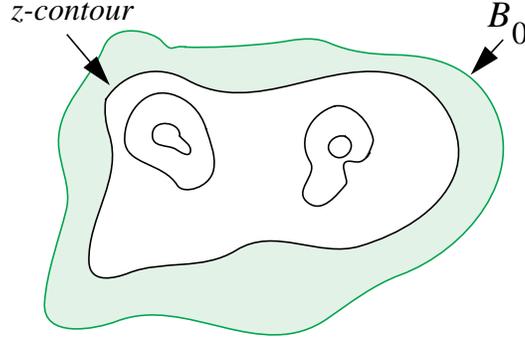


Figure 2. A graphical representation of the SCDF.

more, we will use concomitant information, such as remotely sensed images, digital elevation models, and population densities, to improve the power of SCDF prediction for small areas. From the SCDF, it is possible to predict the area of forested land that falls in health classes (e.g., poor, marginal, good) as defined by the CDI. Using the link, SCDFs can be compared between regions or between the entire spatial domain and a subset of that domain.

Definition of the SCDF

Before we proceed, some background is necessary. Consider the spatial process

$$\{Z(s) : s \in D\}, \quad (1)$$

where D represents the region of interest. Because we are interested in tree crown health, there is a scaling issue of when individual trees, after aggregation, begin to look like a forest. After suitable aggregation, one can represent the ecological index as a random field with continuous spatial index.

Because the field data are taken over a small study site, which we shall denote as Δ , we chose this as our standard area. Henceforth, we shall define Δ as the spatial support unit (SSU). Thus, at location s , we have SSU Δ and $Z(s)$ defined over $\Delta(s)$.

The SCDF for this process is defined as follows:

$$F_{\infty}(z; B_0) \equiv \int_{B_0} I(Z(u) \leq z) du / |B_0|; z \in \mathfrak{R}, \quad (2)$$

where $B_0 \in D$ is the forested portion of D , $|B_0|$ denotes the area of B_0 , and $I(A)$ denotes the indicator function equal to one if A is true and equal to zero otherwise. Then the SCDF is the fraction of area in the region B_0 for which the value of the spatial process Z is less than a cutoff value z . This is depicted graphically in Figure 2.

Because the information that we have is at a countable number of sampling locations and because we will use satellite data and other concomitant information to predict the SCDF, we shall tessell-

late the region B_0 into “tiles” made up of the image pixels. Let

$$B_0 = \bigcup_{i=1}^{N(B_0)} \{A(\mathbf{u}_i)\}, \quad (3)$$

where $A(\mathbf{u}_i)$ represents the image pixel defined at center point \mathbf{u}_i . There are $N(B_0)$ such pixels that make up B_0 . For this analysis, then, we will use (3) and replace (2) with

$$F_\infty(z; B_0) \equiv \sum_{i=1}^{N(B_0)} I(Z(\mathbf{u}_i) \leq z) / N(B_0), \quad (4)$$

where $Z(\mathbf{u}_i)$ refers to the crown index defined over $\Delta(\mathbf{u}_i)$ located at the point \mathbf{u}_i ; $i = 1, \dots, N(B_0)$.

Notice that we have effectively replaced the process $\{Z(s) : s \in D\}$, with a discrete process

$$\{Z(s_i) : i = 1, \dots, N(D)\}, \quad (5)$$

where $N(D)$ is the number of pixels that tessellate D in a manner analogous to (3). This discretization is essential for making progress but does introduce an approximation, the effect of which deserves further study.

Available to the researcher are data from the field,

$$\mathbf{Z} \equiv (Z(s_1), \dots, Z(s_n))', \quad (6)$$

obtained at sampling locations $\{s_1, \dots, s_n\}$. Given these data, a basic predictor of (4) is

$$\hat{F}_n(z; B_0) = \sum_{i=1}^n u(s_i) I(Z(s_i) \leq z) / \sum_{i=1}^n u(s_i), \quad (7)$$

where $\{u(s_1), \dots, u(s_n)\}$ is a set of known weights, for example, the reciprocals of the inclusion probabilities in a sampling design. This is the form of the predictor that is used in this analysis.

Data

SCDF prediction will be examined for the CDI of deciduous trees in the northeast United States. The data were collected as part of the Forest Health Monitoring program within the USEPA's EMAP. The CDI is the weighted average of two variables: crown dieback (CDB) and foliage transparency (FTR). The CDI for SSU $\Delta(s)$ is defined as:

$$Z(s) \equiv \frac{\sum_{j=1}^{n(s)} DBH_j \cdot (CDB_j + FTR_j) / 2}{\sum_{j=1}^{n(s)} DBH_j}; \quad s \in B_0, \quad (8)$$

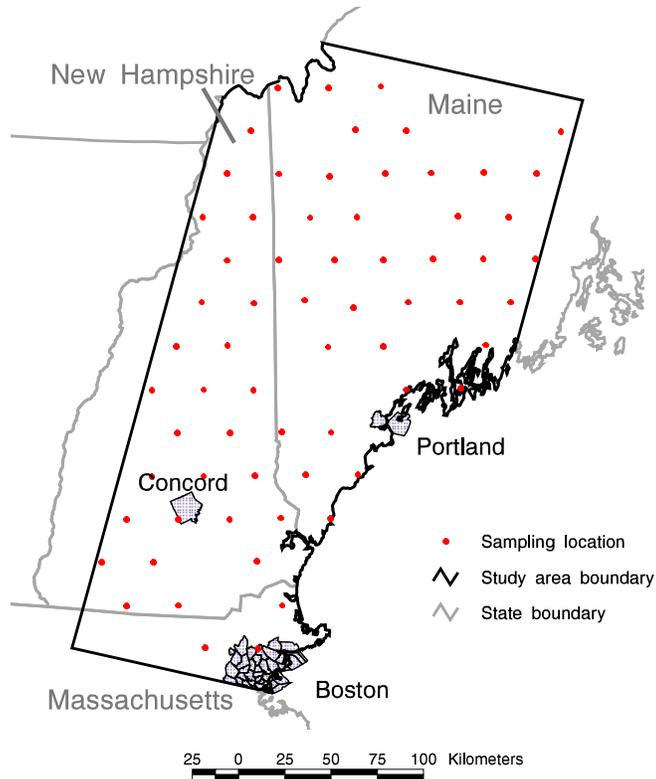


Figure 3. The study area.

where $n(s)$ is the number of trees at sampling location s , and DBH_j is the diameter at breast height of tree j ; $j = 1, \dots, n(s)$.

Crown dieback refers to the percentage of dead branches in the upper, sunlight-exposed parts of the tree crown. The assumption is that these branches have died from stressors in the environment other than lack of light. It is measured as a percentage in increments of 5 from 0% to 100%. Foliage transparency refers to the amount of light penetrating foliated branches. It ignores “holes” in the tree due to bare branches and is measured on the same scale as crown dieback.

The data were collected at sampling sites on the EMAP hexagonal sampling grid (White et al. 1992). The samples analyzed here were collected in the summer of 1992. In the study area, there are 66 sampling sites with deciduous trees.

The region under consideration is in the northeastern United States and includes portions of Maine, Massachusetts, and New Hampshire. This region, which is shown in Figure 3, corresponds to the area of two Landsat satellite scenes.

Methodology

Our goal is to be able to predict the SCDF for small areas. In order to do this, we will exploit associations between sample data and data for which we have complete coverage, for example, remotely sensed data and digital elevation models. Observed associations will be used to predict values for the spatial process being studied at additional locations in the spatial domain. These points will then be used to predict the SCDF of the process for small areas.

The association between sampled data and the concomitant information is assumed to follow a simple linear model. Express the log of the CDI as the linear combination of concomitant variables plus a small-scale stochastic term:

$$Y(s) \equiv \log(Z(s)) = X(s)\beta + \varepsilon(s). \quad (9)$$

This model is fitted using weighted-least-squares regression, with the weights $\{w(s_i)\}$ being equal to the sum of the DBH of trees at each location. The small-scale term is estimated from the residuals of the weighted regression model:

$$\hat{\varepsilon}(s_i) = w(s_i)^{1/2} \left(Y(s_i) - X(s_i)\hat{\beta}_{wls} \right); i = 1, \dots, n. \quad (10)$$

This term is assumed to be intrinsically stationary and can be predicted at any location, s_0 , using optimal spatial prediction (kriging).

After both models (large-scale and small-scale) have been fitted, the spatial process $Z(s_0)$ can be predicted for any location, s_0 , in the spatial domain by:

$$\hat{Z}(s_0) = \exp \left(X(s_0)\hat{\beta}_{wls} + w(s_0)^{-1/2}(\hat{\varepsilon}(s_0)) \right), \quad (11)$$

where $\hat{\beta}_{wls}$ is the fitted regression coefficients from the large-scale model, $w(s_0)$ is the weight for location s_0 , and $\hat{\varepsilon}(s_0)$ is the predicted value for the small-scale term at location s_0 . Details are given below, including the determination of $\{w(s_i)\}$, in equation (14).

The Large-scale Model

The large-scale model is used to exploit associations between sample data and concomitant geographic information. This model was fitted using weighted least squares to express the log of the CDI for deciduous trees as a linear combination of regressor variables. The observations were weighted by the sum of the diameter at breast height for all deciduous trees at each location. The regressors that were considered include:

- *X and Y coordinates*: the coordinates of the sample locations (indicates a spatial trend)
- *precipitation*: the amount of precipitation in each of the four quarters prior to the sample date (predicted from NOAA precipitation values using optimal spatial prediction (i.e., a form of kriging))
- *greenness*: the greenness index of the tasselled cap transformation of Landsat remotely sensed imagery (Crist and Cicone 1984). This variable was calculated for the Landsat scene after a 3x3 average was applied to each pixel. The Landsat images used were acquired with the Landsat 5 sensor on June 12, 1993 (one year after the sample data were collected).
- *topography*: calculated from USGS 3-arc-second digital elevation models
 - *elevation*: the transformation, $\log(\text{elevation})$, was used
 - *slope*: expressed in percent

- *aspect*: the transformation, $\sin((1/2)*\text{aspect})$, was used
- *population density*: the population density was derived from the 1990 U.S. census block groups; the transformation, $\log(\text{population density})$, was used

Model selection

All possible models using the eleven regressor variables were fitted using weighted least squares. The final model was selected using four criteria:

1. low colinearity of regressor variables;
2. low residual sum of squares;
3. high value of R-square; and
4. significance of coefficients.

The colinearity of the regressor variables was evaluated using the condition index (Belsey et al. 1980). Any models with a condition index greater than 500 were not considered. Of the remaining models, the one with the lowest residual sum of squares and highest R-squared was evaluated based on the significance of coefficients. The goal is to find a model for which all coefficients are significantly different from zero at the 95% confidence level. This criteria was applied somewhat loosely, and the final model, which has a coefficient (the coefficient of the variable, *sinaspect*) that doesn't meet the criteria, is deemed acceptable. The largest condition index was 429.

The selected model is given below:

Residual Standard Error = 2.2517, Multiple R-Square = 0.3395
 N = 66, F-statistic = 7.8396 on 4 and 61 df, p-value = 0

	coef	std.err	t.stat	p.value
Intercept	-4.7035	1.9921	-2.3611	0.0214
y	1.2783e-6	0.0000	3.4788	0.0009
<i>sinaspect</i>	0.1551	0.0844	1.8381	0.0709
<i>greenness</i>	-0.0047	0.0022	-2.1381	0.0365
<i>p91q3</i>	0.0534	0.0141	3.7868	0.0004

where *y* is the y coordinate, *p91q3* is the precipitation in the 3rd quarter of 1991, *greenness* is the Landsat greenness index, and *sinaspect* is the transformed aspect variable.

During the large-scale model fitting process, XGobi and the link between Arcview and XGobi were useful for several purposes. First, they helped in the exploratory spatial data analysis and the detection of the spatial outlier in the precipitation data set (see Figure 1). This data set was used to estimate the precipitation at each forest health sampling location. Second, through the use of the correlation tour, XGobi allowed us to check visually to see if there were associations between the explanatory and dependent variables and to check for collinearity among the explanatory variables. Finally, XGobi helped to assess visually regression diagnostics and outliers among the residuals.

Small-scale Model

The small-scale term of the linear model is estimated from the weighted residuals from the fitted

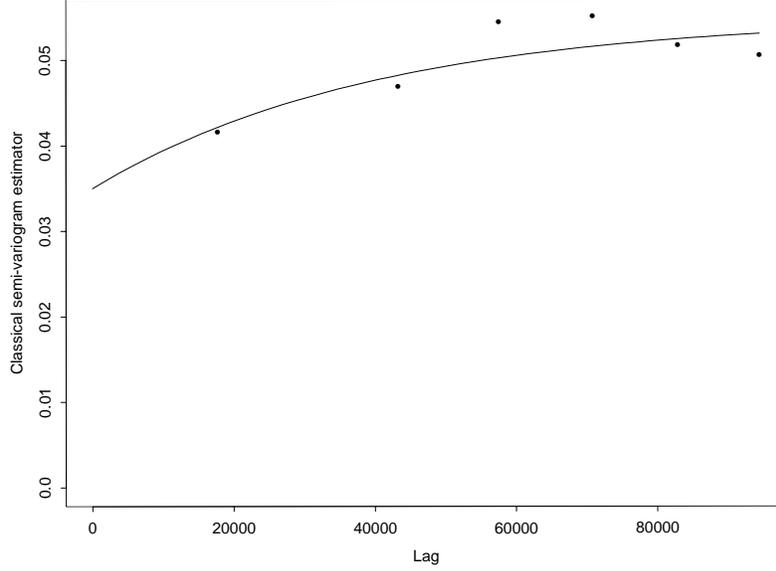


Figure 4. Variogram model fitted to the residuals of the fitted large-scale model.

linear model; see (10). Variogram analysis on these residuals indicate that there is clear spatial structure. The variogram estimates, along with a fitted exponential variogram model, are shown in Figure 4.

When predicting the small-scale term for a location s , constrained kriging (Cressie 1993) was used. Ordinary kriging involves the constraint

$$E(Z(s_0)) = E(p(\mathbf{Z}, s_0)), \quad (12)$$

where $p(\mathbf{Z}, s_0)$ is the kriging predictor. It has been shown (written communication, Aldworth and Cressie) that ordinary kriging produces a process that is too smooth to be used for SCDF prediction. Constrained kriging adds the additional constraint

$$\text{var}(Z(s_0)) = \text{var}(p(\mathbf{Z}, s_0)). \quad (13)$$

Together (12) and (13) match the first two moments of the predictor with the first two moments of the process. If we are to use the predicted values as if they were real data, as we do for SCDF prediction, the additional constraint (13) becomes very important.

Determination of the Spatial Domain

Before SCDF prediction can be carried out, the spatial domain of interest, B_0 , must be determined. In this case, B_0 is the portion of the study area that contains deciduous forests. For our analysis, we estimated this area by using the naturalized difference vegetation index (NDVI) and B_0 is defined as those areas for which the NDVI is greater than 0.5. This area is shown in Figure 5.

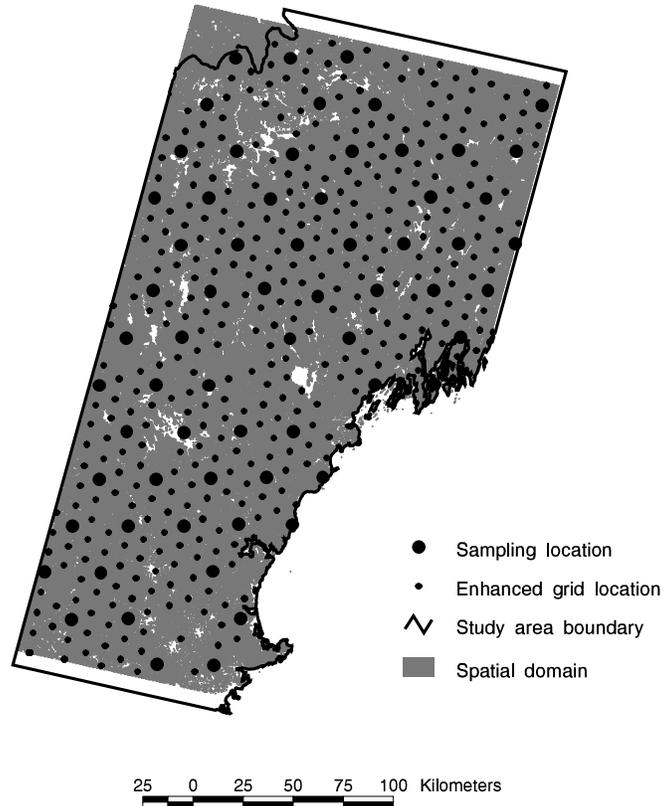


Figure 5. The approximated spatial domain of the CDI process in the study area. Also shown are enhanced grid locations.

Prediction of the Spatial Process at Additional Locations

After the preliminary work of model fitting and determination of the spatial domain has been completed, prediction and visualization of the SCDF can proceed. The first step is to predict the spatial process at additional points within the spatial domain. Points were added that correspond to a 7-factor enhancement of the original hexagonal sampling grid (White et al. 1992). This added 6 points for every point in the original grid (see Figure 5). Using (11) the CDI can be predicted for each new point that falls within the spatial domain. Because the weights used in (11), which are the sum of DBH for all trees at each location, are not known, they must be predicted. In this case, the weights were modeled as a function of the tasselled cap transformation (greenness) of the landsat image. The relationship between the two was determined by simple linear regression. The regression had an R-squared of approximately .25 and resulted in the following model to obtain the weights,

$$w(s_0) = greennes(s_0) \bullet 3.815.$$

Prediction and Visualization of the SCDF

The Arcview 2.1-XGobi SCDF link, introduced in the first section, can be used to predict, view

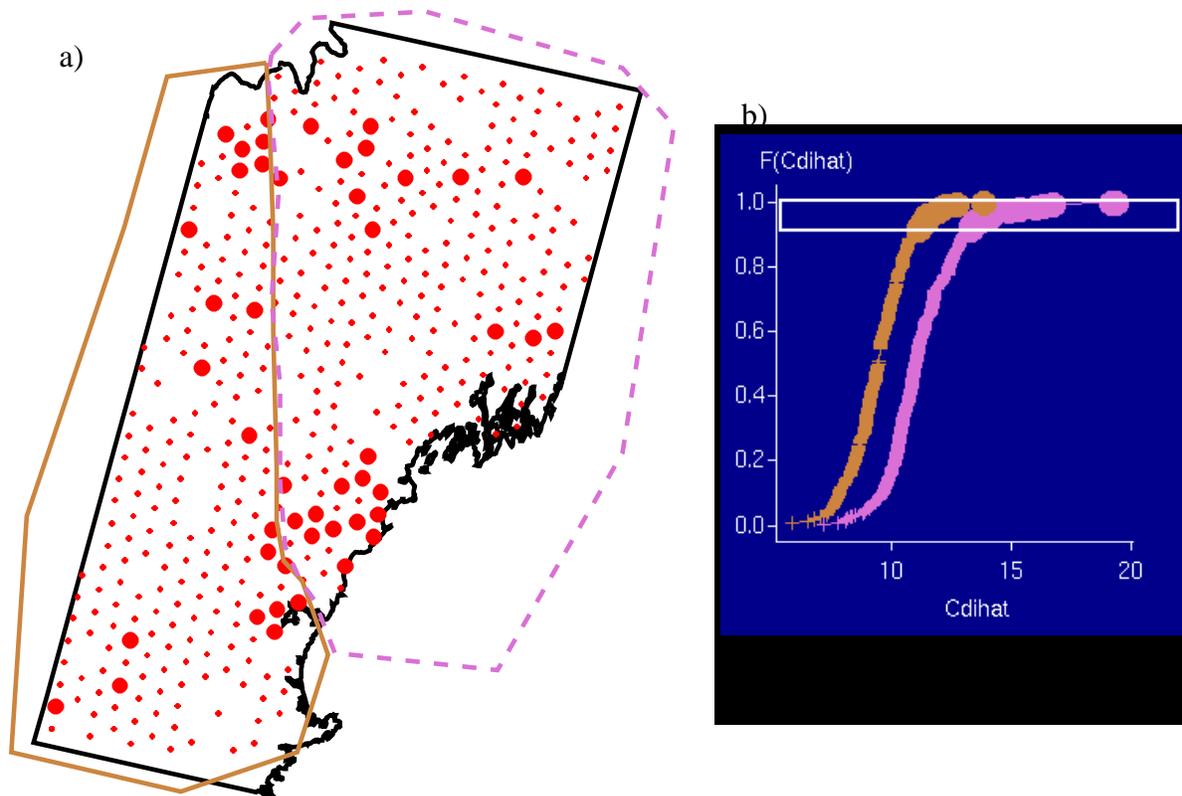


Figure 6. a) Map of study area showing the regions that have been defined for SCDF estimation. b) Predicted SCDFs for regions shown in Figure 6a. The horizontal, white box is being used to brush approximately the top 10% of the values in both SCDFs. Brushed points are shown as large filled circles.

and interactively query the SCDFs. This link provides several capabilities, including: (1) the definition of subregions of the spatial domain over which the SCDF will be calculated (up to 10 regions can be specified); and (2) linked brushing, in both directions, between the Arcview map window and the XGobi SCDF plot.

An example of an analysis using this link is shown in Figure 6. This figure shows the SCDFs calculated for the CDI in two regions: that portion of the study area that falls in the state of Maine (the dashed polygon), and that portion that falls in New Hampshire and Massachusetts (the solid polygon). Figure 6b shows the predicted SCDFs for these regions; the SCDF on the left is for the New Hampshire/Massachusetts region and the SCDF on the right is for the Maine region. Figure 6b indicates that there is a difference in the CDI for the two regions.

Figure 6 also gives an example of the brushing capabilities of the link. In this case, a horizontally shaped brush has been used to brush approximately the highest 10% of the values in both regions (Figure 6b). These points are shown in the map view as large filled circles, indicating the sampling locations containing high values. By moving the brush up and down, various quantiles of the data can be explored. Alternatively, a vertically shaped brush could be used to brush specific ranges of values in the SCDF. This might be done, for example, if *a priori* cutoff values for the index were known that divide the resource into levels. In the current example, these cutoff values might correspond to health classes.

Acknowledgments

Research related to this article was supported by an EPA EMAP grant under cooperative agreement #CR822919. The article has not been subjected to the review of the EPA and thus does not necessarily reflect the view of the agency and no official endorsement should be inferred.

REFERENCES

- Anderson, R. L., Burkman, W.G., Millers, I., and Hoffard, W.H. (1992) Visual crown rating model for upper canopy trees in the eastern United States. USDA Forest Service, Southeastern Region, Forest Pest Management. 15 pp.
- Asimov, D. (1985) The grand tour: A tool for viewing multidimensional data. *SIAM Journal on Scientific and Statistical Computing*, 6(1): 128-143.
- Belsey, D.A., Kuh, E., and Welsch, R.E. (1980). *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. Wiley, New York
- Bradley, R. and Haslett J. (1992) Interactive graphics for the exploratory analysis of spatial data - the interactive variogram cloud. *2nd CODATA Conference on Geomathematics and Geostatistics. Sci. de la Terre, Sér. Inf., Nancy, 1992, 31: 373-386.*
- Buja, A. and Asimov, D. (1986) Grand tour methods: an outline. *Computing Science and Statistics*, 17:63-67.
- Cook, D., Buja, A., Calorera, J., and Hurley, C. (1995) Grand tour and projection pursuit. *Journal of Computational and Graphical Statistics*, 4(3), pp. 155-172
- Cressie, N. (1993) Aggregation in geostatistical problems. In *Geostatistics Tróia '92*, Soares, A. ed, Kluwer, Dordrecht, Vol. 1, pp. 25-36.
- Cressie, N. (1993) *Statistics for Spatial Data*. Wiley, New York.
- Crist, E. P., and Cicone, R. C. (1984) A physically-based transformation of Thematic Mapper data-the TM tasseled cap. *IEEE Transactions on Geoscience and Remote Sensing*, 22(3): 256-263.
- Haslett, et al. (1991) Dynamic graphics for exploring spatial data with application to locating global and local anomalies. *The American Statistician*, 45: 234-242.
- Majure, J. J., Cook, D., Cressie, N., Kaiser, M., Lahiri, S., and Symanzik, J. (1995). Spatial CDF Estimation and Visualization with Applications to Forest Health Monitoring, *Computing Science and Statistics*, Vol. 27, to appear.
- Swayne, D. F., Cook, D., and Buja, A. (1991) XGobi: Interactive dynamic graphics in the X window systems with a link to S. In *ASA Proceedings of the Section on Statistical Graphics*, American Statistical Association, Alexandria, VA, pp. 1-8.
- Symanzik, J., Majure, J. J., Cook, D. (1995). Dynamic graphics in a GIS: a bidirectional link between Arcview 2.0 and XGobi, *Computing Science and Statistics*, Vol. 27, to appear.
- White, D., Kimerling, J., and Overton, S. (1992) Cartographic and geometric components of a global sampling design for environmental monitoring. *Cartography and Geographic Information Systems*, 19(1): 5-21.

NOTES

1. Brushing refers to the ability to change the color/size/glyph of points in the graphics window. Linked brushing means that brushing conducted in any of the linked applications is immediately displayed in all the other applications.