

Raster-Vector Conversion Methods for Automated Cartography  
With Applications in Polygon Maps and Feature Analysis

Shin-yi Hsu  
Department of Geography  
SUNY-Binghamton  
Binghamton, NY 13901  
U.S.A.

Xingyuan Huang  
Department of Geography  
Nanking University  
Nanking, China  
Visiting Scholar  
SUNY-Binghamton

ABSTRACT

In this paper, we discussed the concept of data base conversions between raster format and vector format. With a series of cartographic experiments, we have demonstrated that polygon maps can be generated from raster data without losing the visual quality provided that the editing process is performed to remove the grid cell effect. The methodologies discussed in this paper can therefore serve as a model for using raster data for automated cartography.

INTRODUCTION

Polygon maps constructed by plotters are usually based on vectorized data sets, each constituting a distinctive, labeled region enclosed by a series of line segments. With areal density or typed symbols, these maps are generally called choropleth maps in the cartographic literature. Using the most popular mapping program SYMAP, as an example, the array of (X,Y) coordinates used to bound a region, called A-Conformal Line in that program, is in fact a vectorized data set. The quantitative or qualitative measure for that region, called E-Values, is the basis for generating the statistical surface with the choropleth method.

Instead of using only the edge information, the cartographer can employ a data set that covers the entire study area with a matrix format for mapping purposes. Such data set is called raster data; digitized imagery is one of the most popular forms.

To use raster data for production of maps with line plotters, the raster data have to be processed, and structured in such a way that they contain only two types of information: edges and the interior, corresponding to the A-Conformal Line and E-Values of SYMAP, respectively.

This paper discusses general methodologies for producing edge and interior information using imagery data as examples, and illustrates the techniques for the generation of polygon maps based on raster data with a series of computer maps.

PROCESSING OF RASTER DATA FOR CARTOGRAPHIC APPLICATION

For a given study area, raster data characterizing the statistical surface can be of either univariate or multivariate nature. Using image data for example, digitized black and white imagery can be considered as univariate; whereas multispectral imagery is multivariate. This classification is based on whether merging of two or more data files is needed to create a single file for mapping purposes. Methods for processing these two types of data are discussed below.

## Processing of Single-Channel Data for Raster-Vector Data Base Conversion

The purpose of data processing is to generate distinctive regions and produce edge and interior information for each region. In the context of image data analysis, methodologies for such purposes belong to the general concept of supervised classification and scene segmentation or unsupervised classification.

Supervised classification utilizes calibration samples, called training sets in image processing literature, to classify the entire study area according to given categories plus a rejected class. The techniques for performing supervised classification have been discussed by many researchers and can be obtained from standard textbooks in remote sensing such as these by Sabins (1978) and Hall (1979). A simplified version was given by Hsu (1979).

To classify features or terrain types with one channel data, multiple measures for a given point are generally required to obtain a high rate of correct classification. This is because a single measure for the raster data usually does not give enough information for discrimination purposes. To increase the number of measures, spatial information is generally used. This type of approach is one of the forms of texture analysis; a cartographic approach was given by Hsu (1978).

Provided that there exists enough information in the spectral and spatial domains of the raster data, the training sets are properly selected and analyzed, and finally the classification logic is capable of handling the distributional characteristics of the data, a good classification map can be obtained with appropriate data processing techniques.

The final classification map is in fact, composed of two basic cartographic elements: edges enclosing the classified regions, and the interiors representing the characteristics of the regions. In terms of the SYMAP language, edges are A-Conformal Line, whereas interiors are E-Values. Therefore, when these edges and the information of the interior of each region are extracted and stored in a different data file from the classified map, we have in fact converted the raster data into vectorized data, which can be used by line plotters to generate polygon maps.

In addition to the above-discussed supervised classification method, a family of image processing techniques based on the concept of segmentation can be utilized to generate distinctive regions. In the remote sensing literature, it is generally called unsupervised classification method.

Scene segmentation can be approached from either edge detection, or region growing point of view. The former technique discovers the boundaries of distinctive region using local statistics from adjacent points, whereas the latter delineates distinctive areas by clustering "homogeneous" data points until the growing process touches the edges where another region begins spatially. A more detailed discussion on these topics can be obtained from Hall (1979).

Similar to the classification map generated by a supervised method, the segmented scene can be coded in terms of the edge and interior information using a vector format. Thus a conversion from raster data to vector data can also be achieved based on scene segmentation techniques.

## Processing of Multi-Channel Data for Raster-Vector Data Base Conversion

Similar to single-channel data, multi-channel data in raster format like LANDSAT imagery which has four spectral bands, can be processed by means of both supervised and unsupervised classification methods for mapping purposes.

The methodologies for using multi-channel data to classify a scene are essentially the same as those used in the processing of a single-channel data except that the number of features variables increases by a factor of equal to or larger than the number of channels. For instance, if there are three variables (one tone plus two texture measures) from each band, the number of variables available for analysis in a 4-channel system is at least twelve (12) because additional variables can be derived from ratio bands between any of two channels.

To segment scenes with multi-channel data, certain types of single-band data must be generated by merging these multi-bands. Above-mentioned ratioing technique is one of the commonly-used methods for merging two frames into one.

Another useful technique is the principal component analysis. As it is well-known in the multivariate statistical analysis literature, the number of components is equal to the number of variables; however, only the first few components would provide meaningful information. Using the LANDSAT MSS data for example, usually the first and the second components provide meaningful information.

For segmentation analysis, the component scores map is used as the base representing a combination of the multiple-band data. The meaning of the component has to be interpreted from the relationship between the component and the original variables. Using the LANDSAT data as an example again, the first component usually represents a linear combination of four bands, which is equivalent to panchromatic imagery.

Once the classification or segmentation maps are generated using the above discussed methods with multi-variable data, the same edge and interior extraction algorithm for the analysis of a single-band data can be used to generate vectorized data for plotting polygon maps. The following sections discuss the methods for generating polygon maps using a series of experiments to illustrate the concept of raster-vector data conversion methods as discussed.

### EXPERIMENTATION ON RASTER -VECTOR DATA BASE CONVERSION

#### The Original Data Base

Our experiments begin with a polygon map constructed by a line plotter showing different soil regions as in Figure 1. Note that each region is composed of a series of line segments enclosing that region. Some line segments are shared by two adjacent regions. This map is therefore based upon vectorized data.

To show that raster data can be utilized to generate polygon maps via a data base conversion method, Figure 1 was first converted to Figure 2 showing the interior information instead of the edge information, and then to Figure 3 conveying the same idea but with a raster data set composed of (65 x 58) data points.

## From Raster Data Back to Vector Data

Figure 4 was generated from Figure 3 to depict the edges and the interior of each region using the following scheme (Figure 4a) and algorithm. To extract the edges and the interior simultaneously, we need to identify three types of boundaries:

- (1) exterior boundary separating two regions like that between region 5 and region 6 of Figure 4a;
- (2) interior boundary identifying the inner region like region 6 using (-1 sign) as the region ID Code;
- (3) common boundary between grid cells of the inner region.

The purpose of using these boundaries is to create the necessary edge information (4 edges) for each control point or raster data point so that the exterior boundary can be determined by identifying and subsequently eliminating the interior and the common boundaries that identify a given region.

For example, in Figure 4b IX<sub>Y</sub> (i,j) identifies the center point of a given grid; and (i, j-1), (i+ 1,j), (i,j+1) and (i-1,j) are the ID codes for the four edges of the control point. In addition, (i+1,j), the second edge, is (-1) according to the information given from the adjacent control points, thus it is the interior boundary. Using the same principle, all of the exterior boundaries can be determined, displayed and stored.

The data structure of the stored data is given in Figure 4c. Data Set 1 (Column 1 of Figure 4c) is composed of all the distinctive regions. Each region is subdivided into three sections:

- (1) left-hand side is the mother region;
- (2) right-hand side indicates adjacent regions;
- (3) center part are the (X,Y) coordinates for all of the boundary points; and the boundary points are further identified by the condition whether they belong to one or two adjacent regions. For instance, A are boundaries points without adjacent regions, whereas B, C and D are points shared by region 1 (mother region) and adjacent regions 6, 3 and 2 respectively. Such information is extracted and given in Data Set 2 (Column 2 of Figure 4c). Data set 3 (Column 3 of Figure 4c) identifies line segments that are not shared by two or more adjacent regions.

Figure 5 is constructed from Data Set 1 of Figure 4c, indicating that boundaries separating adjacent regions are plotted (repeated) twice because of the raw data of the boundary points are used.

To eliminate such double-plotting of boundaries, and to be able to extract individual regions, the data structure composed of Data Set 2 and Data Set 3 has to be utilized based on the fact that:

- (1) no overlapping line segments exist in Data Set 2;
- (2) boundaries either belonging to a single region or shared by two regions are identified in Data Set 3;

- (3) data points shared by two or more line segments are identified.

#### Refinement of the Polygon Map

As can be noted on Figure 5, the edges between adjacent regions are darker than the border lines because they are drawn twice by the plotters based on the fact that these line segments are used by two adjacent regions by the plotter.

The polygon map in Figure 5 can be refined first by removing the second plotting on the border line segments as shown in Figure 6 where line weight is even throughout the entire map.

Since Figure 6 is based on raster data, the grid-cell effect exists between nodes of line segments as compared to the original, vector-based map of Figure 1. To produce a visually-pleasing map, Figure 6 was edited by "chiseling" corners produced by grid cells according to the method illustrated in Figure 6a. To eliminate corner points, every 3-point series is examined to determine the existence of corner points. Intermediate points are eliminated by examining the slopes between every two adjacent line segments (each segment is defined by two points). If slopes are equal the center point is eliminated, and vice versa. The final polygon map is shown in Figure 7, and it is almost identical to Figure 1, the original map based on vector data. This proves that our methodologies for raster-vector data base conversion are very effective.

#### APPLICATIONS IN FEATURE ANALYSIS

In image processing, polygons determined by distinctive edges usually represent unique features, such as soil, vegetation, cultivated fields, etc. Once the edges of these features are determined, the texture and tone information of these features can be extracted from the pixels in the interior enclosed by the polygons.

For feature analysis, we in fact utilize the raster data and the vectorized data simultaneously; the former is the interior and the latter is the edge. This model's analysis is applicable to multi-channel and multi-temporal image data sets.

#### REFERENCES

Sabins, F. F., Remote Sensing, Principles and Interpretation, San Francisco; W. H. Freeman and Company, 1978.

Hall, E., Computer Image Processing and Recognition, New York: Academic Press, 1979.

Hsu, S., "Texture Analysis: A Cartographic Approach and Its Applications in Pattern Recognition," The Canadian Cartographer, Vol. 15, No. 2, 1978. pp. 151-166.

Hsu, S., "Automation in Cartography with Remote Sensing Methodologies and Technologies," The Canadian Cartographer, Vol. 16, No. 2, 1979, pp. 183-194.



Figure 1: The Original Soil Region Map with Vector Data

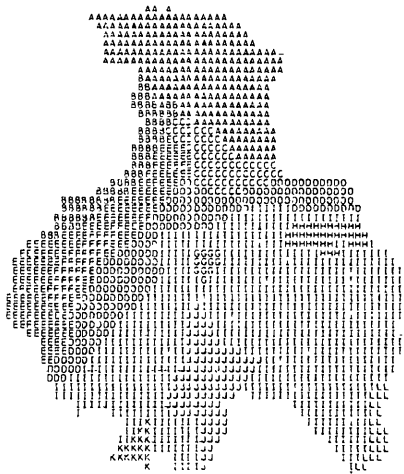


Figure 2: A Raster Data Map of Figure 1



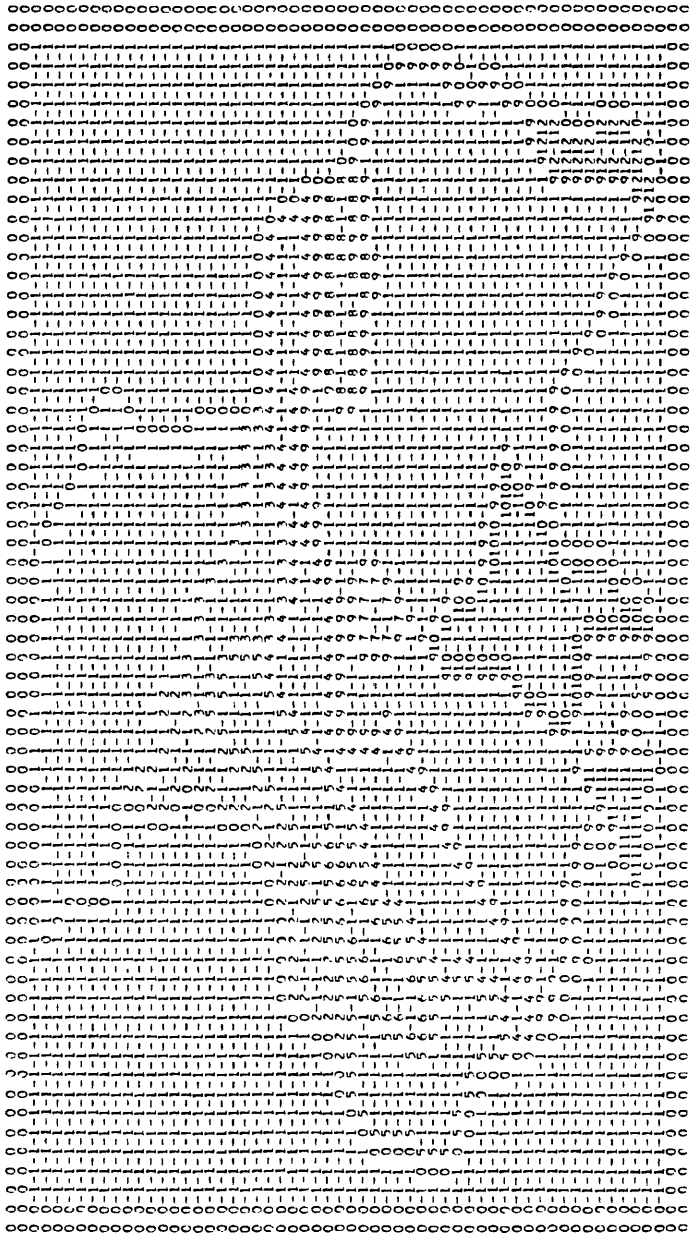


Figure 4: Edge and Interior Map of Figure 3



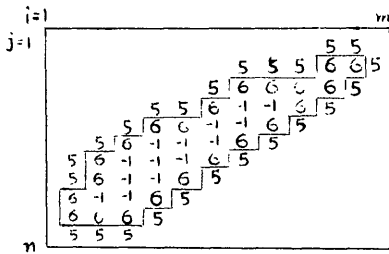


Figure 4a: Three types of boundaries identified

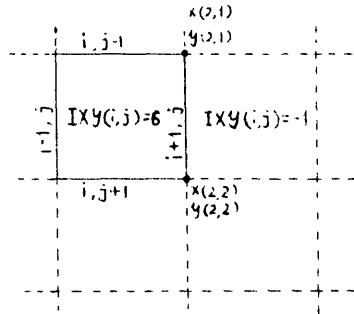


Figure 4b: The ID codes for the center point and four edges of a given grid

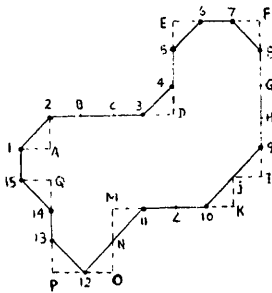


Figure 6a: The points edited and stored

Corner points eliminated:  
A, D, E, F, I, K, M, O, P, Q

Intermediate points eliminated:  
B, C, G, H, J, L, N

Points stored:  
1, 2, 3, 4, 5, 6, 7, 8, 9, 10,  
11, 12, 13, 14, 15

DATA SET 1			DATA SET 2			DATA SET 3		
Mother regions	Boundaries points	Adjacent regions	Boundaries points	Boundaries points	Boundaries points	Boundaries points	Boundaries points	Boundaries points
1	$\frac{a_1}{n_1}$ A	6 3 2	$\frac{c_{10}}{n_{10}}$ B	$\frac{c_{10}}{n_{10}}$ B $\frac{c_{11}}{n_{11}}$ C $\frac{c_{12}}{n_{12}}$ D	$\frac{c_{10}}{n_{10}}$ $\frac{c_{11}}{n_{11}}$ $\frac{c_{12}}{n_{12}}$ A	$\frac{c_{10}}{n_{10}}$ $\frac{c_{11}}{n_{11}}$ $\frac{c_{12}}{n_{12}}$ E $\frac{c_{13}}{n_{13}}$ $\frac{c_{14}}{n_{14}}$ $\frac{c_{15}}{n_{15}}$ F $\frac{c_{16}}{n_{16}}$ $\frac{c_{17}}{n_{17}}$ $\frac{c_{18}}{n_{18}}$ H $\frac{c_{19}}{n_{19}}$ $\frac{c_{20}}{n_{20}}$ $\frac{c_{21}}{n_{21}}$ I		
	B		$\frac{c_{16}}{n_{16}}$ $\frac{c_{17}}{n_{17}}$ $\frac{c_{18}}{n_{18}}$ G		$\frac{c_{16}}{n_{16}}$ $\frac{c_{17}}{n_{17}}$ $\frac{c_{18}}{n_{18}}$ H			
	C		$\frac{c_{19}}{n_{19}}$ $\frac{c_{20}}{n_{20}}$ $\frac{c_{21}}{n_{21}}$ I					
	D							
	E							
2	$\frac{a_2}{n_2}$ F	1 3	$\frac{c_{20}}{n_{20}}$ G	$\frac{c_{20}}{n_{20}}$ G	$\frac{c_{20}}{n_{20}}$ $\frac{c_{21}}{n_{21}}$ $\frac{c_{22}}{n_{22}}$ J	$\frac{c_{20}}{n_{20}}$ $\frac{c_{21}}{n_{21}}$ $\frac{c_{22}}{n_{22}}$ M		
	D							
	G							
	H							
3	$\frac{a_3}{n_3}$ I	2 1 6 4 6	$\frac{c_{30}}{n_{30}}$ J	$\frac{c_{30}}{n_{30}}$ J $\frac{c_{31}}{n_{31}}$ K $\frac{c_{32}}{n_{32}}$ L	$\frac{c_{30}}{n_{30}}$ $\frac{c_{31}}{n_{31}}$ $\frac{c_{32}}{n_{32}}$ M	$\frac{c_{30}}{n_{30}}$ $\frac{c_{31}}{n_{31}}$ $\frac{c_{32}}{n_{32}}$ M		
	G							
	C							
	J							
	K							
	L							
M								
7	$\frac{a_7}{n_7}$ R	6			$\frac{c_{70}}{n_{70}}$ $\frac{c_{71}}{n_{71}}$ $\frac{c_{72}}{n_{72}}$ S	$\frac{c_{70}}{n_{70}}$ $\frac{c_{71}}{n_{71}}$ $\frac{c_{72}}{n_{72}}$ S		
	S							
	T							
	• • •		• • •		• • •			

Figure 4c: The Data Structure

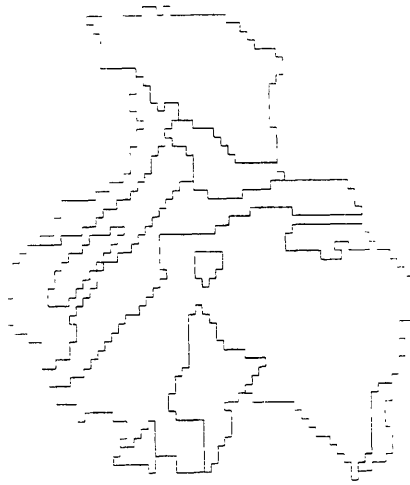


Figure 5. Raw Vector-Plotter Map  
from Figure 4

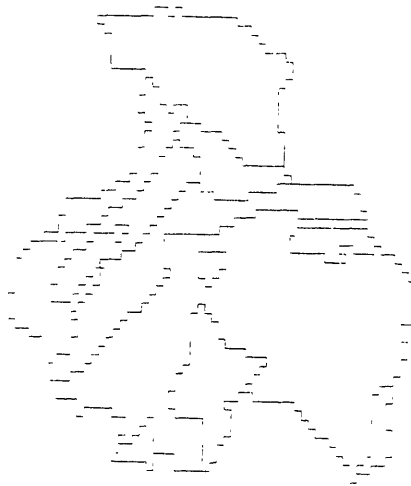


Figure 6 Refined Vector-Plotter  
Map from Figure 5

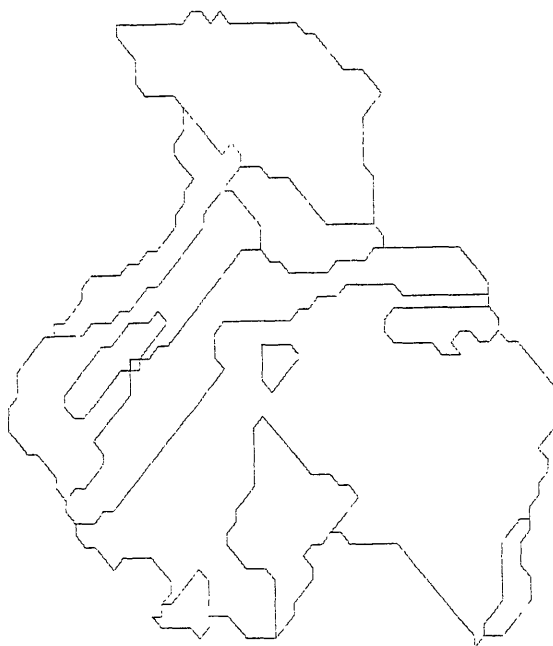


Figure 7: The Final, Edited Map  
from Figure 6