# 17 VISUALIZING FUZZY MAPS

M. GOODCHILD, L. CHIH-CHANG & Y. LEUNG

## ...INTRODUCTION

Increasing emphasis on analysis, modelling and decision support within the GIS applications community in recent years has led to a general concern for issues of data quality. If the purpose of spatial data handling is to make maps, then perhaps it is sufficient to require merely that the output map product be as accurate as the input and to present visualizations of the data as if they were perfectly accurate and certain. But the detailed analytic and modelling applications that underlie much of the recent literature of GIS (Tomlin, 1991; Laurini and Thompson, 1992) demand much more stringent and robust approaches. If the input is known to be inaccurate, uncertain or error-prone, then it is important that the effects of such inaccuracies on the output also be known. Without such knowledge, the apparent value of GIS in supporting spatial decision-making may be illusory. Similarly, it is important that the user be fully informed about data quality, and the most effective way to do this may be by some method of visual communication. In this sense, the arguments for visualization of data quality in GIS seem much stronger than in cartography, and perhaps explain why there are so few techniques for data quality display in the cartographic tradition.

As noted in Chapter 15, uncertainty in spatial databases has many sources, not all of which are captured by the terms 'error' and 'inaccuracy'. The contents of the database can differ from the truth, or from some source of higher accuracy, because of the effects of abstraction in the map-making process, through generalization, abstraction, exaggeration, simplification or classification. All of these create forms of uncertainty for the analyst, and in principle it is desirable to be able to measure, or at least perceive, their effects.

Although inaccuracy is pervasive in spatial data, some types of data are clearly less accurate than others. A GPS (Global Positioning System) survey provides known levels of positional accuracy, potentially down to the nearest centimetre. We focus in this chapter on a class of data known to be subject to relatively high levels of uncertainty, and for which there are no such straightforward measures of accuracy. In this class, every point on the plane is characterised by a single value measured on a nominal or multinomial scale; examples include soil class, land cover class and land use. We refer to this as a multinomial field and such fields are often displayed by chorochromatic, or k-colour, maps. Two

data models are commonly used to build digital representations of such fields. The first, the raster model, is used when the field is obtained by remote sensing, by making use of one of a number of standard procedures for classification. In this model, spatial variation is represented through an array of rectangular cells or pixels, and all information on within-pixel variability is lost. The second, or polygon model, partitions the plane into a number of polygons of arbitrary shape but homogeneous class, thus losing all variability within polygons. The polygon model is also commonly used in making maps of multinomial fields, although the boundary lines on such maps are drawn as continuous curves, and will be discretised as polygons, with straight line segments between vertices, only when the lines are digitised.

Both models are clearly approximations. In the raster model, all variation over distances less than the cell size is lost, and cell size is therefore a convenient descriptor of positional accuracy. The accuracy of the map as a whole is also determined in part by the selection of classes, since in reality, soils or land cover variation is only approximately described by a fixed, finite number of classes, and there will always be within-class variation in the real world which the map purports to represent. In the polygon model, variation within polygons is lost, but polygons are not fixed in size and so there is no simple measure of positional accuracy to compare to the cell size of the raster model. In practice, positional accuracy is often quoted as boundary width, typically 0.5mm, but although boundaries may be drawn with thin lines, the actual width of transition zones in reality may be much larger, and such measures fail to capture the positional uncertainty attributable to a lack of homogeneity within polygons. A better alternative is to derive a measure of positional accuracy from the area of the smallest mapped polygon, or 'minimum mapping unit', since patches smaller than this must have been ignored. A suitable measure is the diameter of a circle of area equal to the minimum mapping unit.

By itself, a measure of positional accuracy does little to capture the inherent uncertainty in raster or polygon representations of multinomial fields. It fails to deal with the problem of within-class variability, or variation over distances less than the positional accuracy of the data. Recently, GIS researchers have begun to turn to concepts of fuzzy classification to provide a more versatile approach to characterising uncertainty, and fuzzy classifiers have become popular sources for error descriptors and error models. The objective of this chapter is to illustrate a number of methods for visualizing this particular approach to uncertainty in multinomial fields. In the raster case, fuzzy classifiers provide a means of describing uncertainty, by associating each pixel not with a single class, but with a vector of class memberships, each one interpreted as a measure of belonging. Thus pixel $x$'s degree of belonging in class $i$ might be denoted by $\pi_i(x)$, and the vector of class memberships might be written:

$$\{\pi_1(x), \pi_2(x), ..., \pi_n(x)\}$$

where $n$ is the number of classes.

In the polygon case, inaccuracy occurs in the form of variation within polygons, perhaps at the edges where boundaries are merely approximations to zones of transition (Mark and Csillag, 1989), or perhaps centrally where small

inclusions and islands of different classes have not been mapped because they fall below the minimum mapping unit area. Neither of these issues is dealt with effectively by giving the polygon model a fuzzy class membership. Instead, it is necessary to abandon the polygon model because it is fundamentally unable to serve as an adequate basis for representing within-polygon variation. Instead, we see the geometry of the polygon model as an artifact of the mapping process, having little value in an effective approach to data quality, and transform to the raster model. Thus both heterogeneity of polygon class and transition near the boundary are represented through the use of pixel class memberships.

While the concept of fuzzy pixel classification is a familiar feature of the remote sensing literature, there has been very little research on its visualization, or on the processing of such data within GIS. In part this may be because of concerns over data volume, since $n$ memberships must be stored for each pixel, rather than one integer between 1 and $n$. In practice, however, it is rare for more than two class memberships to be significantly greater than zero in any one pixel. Fuzzy-classified scenes are difficult to visualize for similar reasons, because the objective in principle would be to communicate $n$ memberships to the user per pixel. Moreover it is not clear how measurements such as class area can be made from such data. Thus despite the availability of fuzzy classifiers, and the greater information content of fuzzy-classified scenes, it is tempting to convert such data to a simple maximum likelihood classification on the grounds that the latter are much easier to handle. As a result, estimates of the area of a given class are biased, and the user viewing a maximum likelihood display is given a falsely optimistic impression of successful classification.

In cartography, a polygon with fuzzy or mixed classification is sometimes shown filled with bars of alternating colours, corresponding to the two classes that are mixed in the polygon. The use of bars of constant width and straight parallel sides allows the eye to distinguish correctly between true boundaries, which have generally complex shapes, and the artificial boundaries formed by the bars. However the method is effectively limited to mixtures of two classes, and becomes confusing if more than a small proportion of polygons are mixed, and there is still the risk that the uninformed user will misunderstand the convention.

The purpose of this chapter is to discuss methods of visualization and processing for fuzzy-classified maps and scenes within GIS. We include with this term not only the results of fuzzy classification in remote sensing, but also derivatives of the polygon model where each pixel is associated with a mixture of classes, or with probabilities of class membership. The next section discusses the meaning of such data from a statistical perspective, and introduces the chapter's error model. This is followed by a description of the environment for visualization of fuzzy-classified scenes developed by the authors. The final summary discusses directions for future research. Sections of the chapter draw on Leung, Goodchild and Lin (1992).

## ... PROBABILISTIC PERSPECTIVE

Consider a raster in which each pixel is associated with a vector of class memberships. The various possible sources of this data were discussed in the previ-

---

ous section. To provide a probabilistic interpretation, we assume that the memberships are normalized by pixel:

$$P_i(x) = \frac{\pi_i(x)}{\Sigma_i \pi_i(x)}$$

Thus $P(x)$ is interpreted as the probability that pixel $x$ belongs to class i out of the $n$ classes. This might be understood in a mixed pixel context as the proportion of pixel $x$'s area that is of class i; or the proportion of interpreters who would have assigned the pixel's area to class i; or the proportion of pixels with the same spectral response as $x$ that are truly i; or in a seasonal sense as the proportion of the year during which the pixel should be assigned to class i rather than some other class. Numerous other interpretations are possible.

We define the term multinomial probability field (MPF) as a vector field whose value at any point is a normalised vector of class membership probabilities of length $n$. In other words, the database can be queried to determine, for any point $(x,y)$, the probabilities that point belongs to any of classes 1 through $n$. A raster provides a suitable way of creating an acceptable approximation of such a field in a digital database, with a spatial resolution defined by the cell size.

Although a display of pixels showing the membership in each class is informative, it nevertheless fails to convey an impression of uncertainty, suggesting that memberships are expressions of deterministic knowledge, rather than of lack of knowledge, or of fuzziness. Instead, we focus on class memberships as descriptions of uncertainty, and on the range of possible maps that might therefore exist. In other words, and using the terms defined in the introduction to this section, the fuzzy class memberships is defined by realisations of that model, and the range of possibilities is defined by realisations of that model. Goodchild, Sun and Yang (1992) define an error model in the context of spatial databases as 'a stochastic process capable of generating distorted versions of the same reality'. The best known error model is the Gaussian, used to describe uncertainty in measurements of a simple scalar quantity like the elevation at a point. Each of the outcomes of such an error model provides one possible version of the truth, as it might be interpreted by one soil scientist, or as it might be digitised by one operator. Thus we see a map as a collection of interdependent measurements, and an error model for maps as a means of generating simulations of alternative maps within the error model's inherent range of uncertainty in those measurements. One realisation of a map error model might be one scientist's interpretation of variation in land cover over a given area, or one digitiser operator's effort to capture the contents of a given map. Goodchild, Sun and Yang (1992) describe an error model for an MPF. Each realisation is a map in which each pixel is assigned to a single class. The model's two essential properties are:

(a) between realisations, the proportion of times pixel $x$ is assigned to class i approaches $P_i(x)$ as the number of realisations becomes large; and
(b) within realisations, the outcomes in neighbouring pixels are correlated, the degree of correlation being controlled by a spatial dependence parameter $\rho$.

When the spatial dependence parameter is zero, outcomes are independent in

each pixel. This is the case illustrated by Fisher (1991b). However, this is almost certainly unrealistic since few if any real processes are likely to create such independent outcomes. As the parameter increases, outcomes are correlated over longer and longer distances. One suitable interpretation of this is that larger and larger inclusions within polygons are ignored, or fall below the minimum mapping unit area.

For example, consider an agricultural field of 100ha, captured as a raster database with a cell size of 0.01ha (10m square). A soil map of the field might include four classes, with significant uncertainty because of fuzzy boundaries between classes, inclusions, etc. Because the pixel size of 1ha has no connection with any process of soil formation or development, it is next to impossible that the true class will be independent in neighbouring pixels. Instead, we expect neighbouring pixels to have strongly correlated classes, and inclusions to extend over distances substantially more than 10m. On the other hand, while there may be uncertainty over the crop grown in the field, it is certain that the field has only one crop. Spatial dependence is so strong in this case that in any realisation of the error model, the entire agricultural field can have only one class.

Many commonly used descriptions of map error fail to meet the requirements of an error model, since they fall short of the complete specification of a stochastic process. Such descriptions include the width of an epsilon band, the measures mandated by many map accuracy standards, the statistics of the misclassification matrix used in remote sensing, and the reliability diagram found on many topographic maps. All of these are useful error descriptors, but fall short of being useful error models. Neither is there a useful connection between many such descriptors and the necessary parameters of error models. For example, it is not possible to connect the parameters in the model described above with such measures as positional accuracy of polygon boundaries, or per-polygon misclassification of attributes. Visualization is perhaps the only way of overcoming this conceptual barrier, since it allows us to connect the model and its parameters with their implications in the form of simple pictures.

## ... TOOLS FOR VISUALIZATION

In this section we describe the tools we have developed for visualization of MPFs. These techniques include some designed to display the descriptors of error (parameters of the error model) and others designed to generate and display realisations of the error model. As we argued in the first section, classification procedures are important for remotely sensed imagery, but it is also desirable to be able to visualize MPFs from sources such as land cover maps in which the classification is performed by other means. For this reason, the system is modular in design, and includes a classification module, display module, and modules for data manipulation. Only the display modules are discussed here (for additional details see Leung, Goodchild and Lin, 1992). The system is interactive and uses a graphic user interface, all instructions and operations being triggered by selecting appropriate screen buttons. Windows are opened and closed as appropriate. It has been developed in C and X Windows for the IBM RS/6000 under the AIX operating system.

In the display module, images can be displayed directly by associating colours with spectral bands without classification, in order to support direct visualization of the preclassified scene. However the most important component of the module supports the display of classified images. In general, techniques of dithering and bit-mapping can be used to display uncertainty in terms of levels of class membership, to expose the spatial variation in membership within regions or across region boundaries. In addition the system provides several other measures and methods for conveying information about an MPF to the user. The following sections briefly describe and illustrate the principal tools.

The illustrations, Plates 17 to 20, are derived from a Landsat TM scene for an area to the west of Santa Barbara, California. For the purpose of illustration, the scene was classified using only three classes, identified as water, vegetation and soil. Since these clearly do not characterise the full range of spectral responses, pixel fuzzy memberships tend to be relatively low and mixed. This somewhat artificial classification example produces artifacts which are the subject of comment below.

### UNCLASSIFIED IMAGE

Colours can be assigned to spectral bands to create conventional false-colour representations of the unclassified scene. This allows the user to see the raw data before classification.

### CLASSIFIED IMAGE

The RGB (red/green/blue) colour model is used to display the results generated by the fuzzy classifier, or input from some other source. Each class is associated with a point in RGB space, and each vector of class memberships is mapped to an intermediate point in the colour space by linear interpolation. This method is successful for two classes ($n=2$) provided the pure-class colours are chosen carefully, but it is difficult for the eye to decode the results for $n=3$, and for $n>3$ the mapping from class membership vector to colour space is no longer unique. Moreover mapping is non-unique for $n=3$ if the class memberships have not been normalized to sum to 1 (see previous section).

Plate 17 shows a display of the image using this approach assigning water to blue, vegetation to green and soil to red. Although all of the information contained in the fuzzy memberships is displayed in this rendering, it is virtually impossible for the eye and brain to deconvolute the linear mixing of colours. On the other hand, this may be a useful display if the user merely wishes to identify the memberships associated with a given pixel; the bars at the top of the image display the RGB components of the pixel currently selected by the cursor.

To deal with the difficulty of visualizing membership in many classes, it is possible to display each class's memberships separately using a grey scale. By using multiple windows one can display the general distribution of each class for up to four or even six classes simultaneously. Plate 18 shows the memberships for the three classes water, vegetation and soil in the upper left, lower left

and upper right windows respectively. While these are easier to interpret than Plate 17, the displays convey no collective sense of a pixel's memberships.

Sometimes it is desirable to have a non-fuzzy image of a fuzzy scene. A simple defuzzing mechanism is maximum likelihood, where the displayed class $f(x) = i$ if $\pi_i(x) > \pi_j(x)$ for all i, j, i not equal to j; that is, a pixel is assigned to class i (and displayed with class i's colour) if its degree of membership in class i is highest of all its memberships. The user has control over the colours assigned to each class. Frequency distributions of the entire image can be displayed, and the user can zoom into a selected area, or display the contents of any pixel.

## AREA

Calculation of the area occupied by each class is a common GIS function. For conventionally classified scenes or other forms of raster data it is calculated by counting the pixels assigned to each class and multiplying by pixel area. However the solution is less clear in the case of fuzzy-classified scenes. If $P_i(x)$ is interpreted as the proportion of pixel x that is truly class i, as in a mixed pixel interpretation of fuzziness, then the area of class i will be the sum of such fractions added over the scene. On the other hand if $P_i(x)$ is interpreted as the expected probabilistically, the same estimate must be interpreted as the expected area of class i. Similar approaches are appropriate if $P_i(x)$ is given other probabilistic interpretations. Thus the calculation of area on a fuzzy-classified scene seems adequately addressed by calculating:

$$A_i = b\sum_x P_i(x)$$

where b is the area of each raster cell. Note that this estimate may be very different from the conventional one based on maximum likelihood, that is,

$$A_i^* = b\sum_x f(x)$$

More difficult is the estimation of error variance, standard error, or the uncertainty associated with such estimates. In the mixed pixel interpretation A, is deterministic, with zero uncertainty. In a probabilistic interpretation, and assuming that outcomes in each pixel are independent of outcomes in neighbouring pixels (zero spatial dependence) then the uncertainty associated with area estimates can be determined from the statistics of the binomial distribution in the form of a standard error:

$$e_i = b(\sum_x P_i(x)[1 - P_i(x)])^{0.5}$$

where $e_i$ is the root mean square uncertainty in estimate $A_i$. Fisher (1991) used Monte Carlo simulation to estimate this standard error. When spatial dependence is present, as it almost always is, and outcomes in neighbouring pixels are correlated, it is necessary to resort to the methods described by Goodchild, Sun and Yang (1992).

## ENTROPY

The degree of certainty in a pixel's classification can be measured in various ways, but one that expresses the degree to which membership is concentrated in a particular class, rather than spread over a number of classes, is the information statistic or entropy measure:

$$H(x) = \frac{1}{\log_e n}\sum_i P_i(x)\log_e P_i(x)$$

where $H(x)$ is the entropy associated with pixel x. $H(x)$ varies from 0 (one class has probability 1, all others have probability 0) to 1 (all classes have probability equal to $1/n$). The system allows a map of H to be displayed using a grey scale; light areas have high certainty (probability concentrated in one class) while dark areas have low certainty.

The distribution of the entropy statistic for the Santa Barbara scene is shown in Plate 19. Clearly evident is the ocean turbidity off Point Arguello, where mixing of currents produces substantial reduction in the membership of the water class, and thus an increasing level of entropy in these pixels. In general the ocean is dark because of its high membership in the water class, but the land pixels are much more mixed and thus lighter in this rendering. Some of the greatest levels of uncertainty are in pixels located in the Town of Lompoc, because urban is not a recognised class.

The degree of fuzziness associated with membership in each class can be assessed by another form of the entropy measure:

$$H_i = \frac{1}{N\log_e 2}\sum_x \{P_i(x)\log_e P_i(x) + [1 - P_i(x)]\log_e[1 - P_i(x)]\}$$

where the sum is now over the pixels and N is the number of pixels. Entropy is zero if the probability of membership in class i is 0 or 1 in all pixels, and 1 if probability is 0.5 in all pixels. The overall entropy H of the entire fuzzy scene can be obtained by adding these measures over all classes.

## REALISATIONS

As noted earlier, an important aspect of visualizing uncertainty is the ability to view individual realisations of an error model, rather than its parameters. All of the previously noted methods display some aspect of the probability vectors, which are the parameters of the error model's stochastic process, rather than its outcomes. Viewing a display of probability vectors necessarily diverts attention from the variation between realisations, and focuses more on the average or expected case.

The system includes the ability to display realisations of the error model, using user-determined levels of spatial dependence. Goodchild, Sun and Yang (1992) discuss possible methods for determining appropriate levels, as attributes of the entire map, or of individual classes, or of geographic regions. A display of four or six different realisations in different windows on the screen provides

graphic illustration of the implications of uncertainty in spatial data, and draws attention to its influence on analysis, modelling and decision-making.

Plate 20 shows a series of realisations of the Santa Barbara scene, using probabilities computed by normalising class memberships. Each of the four images shows a different value of spatial dependence $\rho$, ranging from 0.20 in the upper left to 0.25 in the lower right. In reality, spatial dependence is certainly a function of class, and probably also varies regionally. However each illustration represents a realisation under a uniform value of $\rho$.

The most obvious visual effect of the spatial dependence parameter is shown in the sizes of inclusions. As $\rho$ approaches 0.25, inclusions become larger, as illustrated by the large blobs in the ocean in the lower right illustration. The lower class membership for water around Point Arguello have produced a higher density of inclusions. When compared to Plate 17, these illustrations demonstrate the sharp difference between the two approaches to display of uncertainty: the parameters of the error model in the form of class memberships (Plate 17) and realisations of the model as a stochastic process (Plate 20).

## ...SUMMARY AND FUTURE DIRECTIONS

It is often argued in the GIS community that while uncertainty is endemic to spatial data and undoubtedly affects the outcomes of spatial data processing, it is best not to draw attention to it because of its complexity and potentially damaging effects on decision-making. The user 'does not want to know'. Analogous software systems, such as the statistical packages and database management systems, do not include techniques for capturing, storing and manipulating explicit information on uncertainty, so why should GIS? We believe that this argument is both intellectually unsound and disastrously short-sighted. Most spatial decisions, particularly important ones, are made in an environment of conflict and controversy. As GIS matures and becomes available to more and more parties to a debate, the naive view that the party with the GIS somehow carries greater weight will become less and less realistic, and easier and easier to attack. Pressures for better quality assurance and control are already emerging from instances of GIS-related litigation.

In the context provided by a stochastic error model, the user of an error-handling GIS has three alternatives. First, displays can omit all reference to uncertainty by showing the most likely class for each pixel. This is the most commonly encountered option, and has pervaded GIS to date, particularly in its applications to the spatial distribution of environmental parameters such as soil class or land cover class. Second, displays can inform the user about uncertainty through error descriptors, or through the parameters of error models. This option is represented here by the display of class membership probabilities. Finally, the user can be shown samples from the range of possible maps in the form of simulations under the error model. The differences between these three options are striking, and immediately suggestive. Realisations draw attention to the effects of error, and may be much more meaningful to a user who lacks a full understanding of the concepts of statistics. On the other hand, the arguments that have sustained a lack of attention to uncertainty in cartographic tradition – the desire not to confuse the process of communication, and a will-

ingness to portray the world as simpler than it really is – presumably apply here also. We hope these illustrations, and those provided in the next chapter, will stimulate debate on these important issues. What impact can the use of these options have on the process of spatial decision-making, and where is each option most useful?

Spatial statistics is a complex and difficult field, and few GIS practitioners have more than an elementary understanding of its techniques and concepts. Moreover visual techniques are inherently convincing and communicative. Thus it seems that visualization will have to be a fundamental part of any concerted effort to handle uncertainty within GIS. Goodchild, Sun and Yang (1992) have argued that visualization is the key to user participation in the determination of the key spatial dependence parameters in spatial statistical models of uncertainty. In the case of the error model used in this chapter, visualization holds the key to involving mapping specialists in the determination of appropriate values of the spatial dependence parameter $\rho$.

An MPF is inherently multidimensional, and this chapter has presented a number of techniques for improving the user's ability to understand this particular form of spatial variation. However any communication system must satisfy the requirements of the user as much as it exploits the capabilities of the system, and it seems clear to us that an ideal design can only come from the experience of working with these tools in a real analytic environment.

## ...ACKNOWLEDGEMENT

## ...REFERENCES

Fisher P F (1991b) Modelling soil map-unit inclusions by Monte Carlo simulation. *International Journal of Geographical Information Systems* 5(2): 193–208.

Goodchild M F, Sun G, Yang S (1992) Development and test of an error model for categorical data. *International Journal of Geographical Information Systems* 6(2): 87–104

Laurini R, Thompson D (1992) *Fundamentals of Spatial Information Systems*, Academic Press, San Diego, Calif.

Leung Y, Goodchild M F, Lin Chih-Chang (1992) Visualization of fuzzy scenes and probability fields. *Proceedings, Fifth International Symposium on Spatial Data Handling*, Charleston, South Carolina, International Geographical Union, Columbus, Ohio

Mark D M, Csillag F (1989) The nature of boundaries in 'area-class' maps. *Cartographica* 26(1): 65–78

Tomlin C D (1991) *Geographic Information Systems and Cartographic Modeling*, Prentice-Hall, Englewood Cliffs, NJ