# FEATURE EXTRACTION COMPARISON OF IMAGE ANALYSIS SYSTEMS AND GEOGRAPHIC INFORMATION SYSTEMS

J Gairns, Intera Information Technologies, Canada
T Taylor, DIPIX Technologies Incorporated, Canada

## ABSTRACT

Today, user demands and improvements in Information Systems are bringing together data formats that were not previously physically linked. With the advent of Geographic Information Systems (GIS), it was realized that data from a wide variety of sources could be used in a complimentary fashion. One example is remotely sensed imagery in a GIS environment. A GIS coverage has a unique feature that remotely sensed imagery does not. A database. A database is an integral component of a GIS, but it requires extensive management in order that useful information may be stored and manipulated. Remotely sensed data offers an elegant supplement to such a database, in the form of extensive spatial information. Information from a raster image can be extracted automatically, with or without operator supervision. By using an image as a data source and the GIS vectors as boundaries to delimit that data, the two environments can offer more functionality than either alone.

There are concerns regarding the marriage of these two data sources. If the GIS data, once brought together with the raster data is not properly registered, problems can arise. For example, data that is not adequately geo-referenced is hardly useful at all, and results derived from such can easily be highly deceiving.

The raster and vector environments can offer a great deal of information exchange to each other. In fact, the existence of either data type enhances the information content of the other. Given this perfect marriage, one must consider the consequences of bringing one type of data into the realm of the other. What happens to raster data when it becomes vectorized and conversely, what happens to rasterized vector data. Is there perfect registration, or is the registration of these complimentary data not as straightforward ? An ARC/INFO vector polygon coverage of water boundaries was integrated with an ARIES format NOAA AVHRR Local Area Coverage (LAC) image of the corresponding geographic area. The GIS polygons were rasterized and visually compared to the existing NOAA sub-scene. A Feature Extraction technique was performed on the NOAA sub-scene and the resulting vectors were compared against the original GIS coverage using a simple visual comparison.

Keywords: Feature Extraction, GIS/IAS Integration, Accuracy

## 1.0 INTRODUCTION

Traditional Image Analysis Systems (IAS) offer an ideal compliment to GIS data extraction, manipulation and archiving functionality. The extraction of image statistics using a GIS overlay is an obvious benefit. For example, an operator can automatically select training areas by using the functionality of the GIS and querying the pixels that fall within a polygon. This information could easily be stored in a database, and subsequently manipulated as a database attribute.

Given that spatial data has an extremely high information content for a relatively low cost, it is desirable to integrate spatial data with a topological database, such as is inherent in a GIS. Spatial data offers vast quantities of information, but one must consider what happens when spatial data is brought together with other data types. The purpose of this paper is to explore the implications of merging traditionally detached sources of information via automated, or semi-automated procedures, in particular, feature extraction.

There is a desirable effect when data from a GIS is merged with remotely sensed imagery. This serendipitous effect is information synergism. Information synergism is the overall increase in information content of a system, exceeding that of the individual data components. By modelling the various data types in a single environment, information that was not previously obvious becomes evident.

Most existing research into comparative or relative feature extraction deals specifically with the use of some sort of interactive component to the procedure (Schowengerdt and Pries (1988), Zelek (1990), O'Brien (1991)). These are perfectly valid approaches, but a user may not always have such intimate knowledge of a study site which reduces the potential efficiency of the extraction technique. What does one do in this case ? The answer points to an automatic approach. Work in this area is still very much in the research phase, although it is approaching an operational level.

A particularly important concept in the field of feature extraction is how an algorithm actually recognizes an edge or boundary. The procedure for locating linear features is very similar to that of locating spatial features. Both have edges, which can be thought of as a "contrast amongst distinct features in the image" (Zelek, 1990). Characteristics of an edge such as pattern, size, shape and colour are important elements in the recognition of the contrast that delineates the edge of a particular feature. It is difficult to quantify these characteristics, however a qualitative approach can prove useful as a tool for comparison in this case.

## 2.0 DATA DESCRIPTION

Practically speaking, most remotely sensed imagery could be used in a study such as this, however NOAA AVHRR was chosen for the task. The choice to use this data was based on the availability of

data within the framework of an on-going project at the Canada Centre for Remote Sensing (CCRS) known as the Crop Information System (Manore et al. (1989)). Water bodies were chosen for the comparison because they are simple to recognize visually. The data stratification scheme performed was a simple maximum likelihood classification on the original imagery. The result of the classification was an 11 class land cover theme image. In particular, the water-body theme was used for processing into a resultant vector coverage.

A manually digitized coverage of water bodies for the study area was available. This vector coverage was the control vector coverage against which the output from the feature extraction were compared.

## 3.0 METHODOLOGY

The following is a generic methodology for feature extraction from raster imagery. The assumption made is that the imagery contains some useful geographic information, but that this information is not in the proper format, which in turn provides the impetus for the operator to extract these features from the imagery. Feature extraction gives a user a vector product that can easily be integrated into the GIS environment with as little operator interaction as possible. It is recognized, however, that given the current state of software development, this fully-automated methodology, while promising, is not yet feasible for operational use. The approach above was followed in this paper as much as possible with exceptions discussed.

### 3.1 Raster to Vector Conversion

#### Classification of Raster Image

It is understood that it is possible to extract features from an unclassified image, likewise it is also understood that the job of feature extraction would proceed much more simply if the data were stratified. There are procedures that can be used to stratify the data (e.g., density slicing or supervised/unsupervised classification). The purpose of stratifying the data is to make the analysis procedures more practical, in terms of processing time and disk storage.

#### Feature Extraction

The purpose of the feature extraction procedure is to identify homogeneous clusters of pixels. In this case, only a single theme class (water bodies) was used from the original 11 theme classification. Sub-pixel elements were not considered to be significant in this comparison.

#### Boundary Extraction

Once the homogeneous clusters of pixels have been identified by the feature extraction process the next step is to delineate the edge of these clusters. This is the process of boundary detection, which is also known as image segmentation. The output from the boundary extraction procedure is, presumably, a vector representation of the original polygonal structure or feature.

#### Generalization of extracted vectors

If the extracted boundaries were examined at this point, they would appear to be 'step-like'. That is, they would follow the exact contours of the pixel edges. In order to smooth out these 'steps' and create a more realistic representation of the feature, a smoothing-filter needs to be passed over the edge. The larger the filter size, the greater the effects of the smoothing. Thus, while a 3 x 3 filter might smooth the step-edge slightly, a 9 x 9 filter might distort the edge and even larger filters might shift the X and Y coordinates.

By generalizing the data, the data volume is also decreased. The benefit of this is decreased data storage requirements and increased processing speed. The potential deficiency of this is that the data may become too generalized, and not very well registered.

#### Export of vectors to GIS

Up to this point, the work done has been entirely in the image processing domain. The extracted features are now ready to be exported to the GIS. This procedure is a straightforward translation of the extracted vectors to a format compatible with the GIS, such as the Digital Line graph (DLG) format.

## 3.2 Vector to Raster Conversion

**Convert the vector strings into a raster representation.**

This data conversion step is relatively simple. Since both the raster image and the GIS linework are, presumably, geographically referenced, the task at hand is to determine whether a given vector falls within a specific pixel. Since we are generally not interested in sub-pixel features (ie. features smaller than the spatial resolution of the image), such trivial elements should be removed. These could be deleted by filtering out or deleting elements less than a user-specified threshold.

**Export resulting raster to IAS**

This step is similar to bringing data from an IAS to a GIS. The data is translated to an intermediate format, such as DLG, and subsequently exported to the IAS.

The rasterized vectors can then be displayed as an overlay on the raw imagery to assess the relative accuracy of the linework. Using water bodies, for instance, allows an operator to visually inspect whether the linework is geographically accurate with respect to the image. In some cases, a 'live-link' to the GIS database can be maintained, but a discussion of this is beyond the scope of this paper.

## 3.3 Comparison of the two data conversion routes

In this paper, the accuracy of a feature extraction technique using data from a land cover classification was qualitatively compared against a rasterization of a GIS vector coverage. Specifically, water boundaries were used to reference the two data sets. A scheme of scoring both of the procedures based on 4 of the 9 elements of image interpretation (Bowden and Pruitt (1974)) was adopted. The 4 criteria chosen were size, shape, resolution (scale) and geometric accuracy of the end products of the processing. If the size and shape of each of the elements were similar, a high score was given. If the resolution of the elements were closely matched, a high score was given. If the elements overlapped well, a high score was given for accuracy. The values assigned to each criteria were ranked from 1 (poor) to 10 (excellent). The results of

the qualitative comparison are tabulated in Tables 1 and 2.

**Feature Extraction Technique**

In this case, a control dataset of classified NOAA AVHRR imagery that had undergone the feature extraction procedure was used. The resulting vector data were imported into the GIS and displayed with the manually digitized water body coverage.

**Rasterization of Vector Coverage**

In this case, a control dataset of manually digitized water bodies that had been rasterized was used. The rasterized vector data was exported to the IAS environment and displayed as an image overlay on the unclassified image. The proximity of the raster water body theme to known water features was observed and then scored.

### Table 1.

Qualitative Evaluation of the Accuracy
of a Rasterized Vector Coverage against
a Georeferenced NOAA Image Composite

| Criteria | Performance Score (1-10) |
|---|---|
| Size | 5 |
| Shape | 8 |
| Resolutio | 7 |
| Accuracy | 9 |
| | 29 |

### Table 2.

Qualitative Evaluation of the Accuracy
of Vectors Extracted from a Land Cover
Classification against a Digitized Coverage
of Water Bodies

| Criteria | Performance Score (1-10) |
|---|---|
| Size | 5 |
| Shape | 5 |
| Resolutio | 6 |
| Accuracy | 7 |
| | 23 |

367

## 4.0 DISCUSSION

The results profile the implications of importing raster data into the vector domain of a GIS and importing vector data into the raster environment of an IAS. The effects of the various data transformations with respect to geographic accuracy is addressed.

There is a obstacle with some GISs, in that, there are limitations in the software. Most GISs have a practical limit in terms of the number of elements that can be addressed in a single coverage. While this limit is generally large, it is still a limit. This leads one to consider the complexity of an image. That is, is the coverage going to exceed the limits of the software ? More and more, this is becoming a bottleneck for analysis. Many researchers must devise creative solutions to deal with these inherent software limitations. Occasionally, these limits are practical, rather than physical. That is, they reflect the hardware limits more than the software restrictions. Hardware limitations include disk storage capacity and processing speed. By increasing either of these, the user is faced with increasingly cost-ineffective solutions to their problems.

Since feature extraction is concerned with the simplification of highly complex information, it follows that the actual process involved is likewise highly complex. The simplest solution to this problem is to first stratify the complex data and then perform the feature extraction procedure. Image classification is just such a stratification scheme, albeit a complex one. In a simple case, we observe a standard 8 bit image channel to have 256 possible digital values. Imagine the increased complexity by adding further channels. Conversely, consider the case where the original image can be stratified, through a supervised classification technique, to a mere 11 classes. The task of feature extraction becomes considerably easier.

Because we have stratified the data into a number of desired classes, we have some control over how the individual pixels become classified. This implies that the number of potential artifacts that could result from the feature extraction process is minimized. In a homogeneous field, the number of misclassified pixels is minimal. However, in a heterogeneous field, the potential number of misclassified pixels increases. Thus, obviously, one could expect quite a few artifacts from a heterogeneous field, and little, if any, from a homogeneous field.

The procedure of extracting vector features from a raster theme image is currently still in the research phase (O'Brien (1991), Taylor et al. (1991)). Experiments are going on that are scene and situation specific, and thus can not be applied to a general case. Progress is being made in this field, but results must be taken with a grain of salt (Taylor (1991)). There are semi-automated procedures for feature extraction that are being used, but they require considerable operator input (Zelek (1990), Van Cleynenbreugel et al. (1990)). For example, a feature is identified by an operator and at a certain point the extraction or recognition algorithm takes over. This sort of procedure generally produces more reliable results than the fully automated procedures, but at the cost of greater operator interaction. Currently rule-based feature extraction techniques use operator expertise and knowledge of a specific site to aid in the extraction process (e.g. Van Cleynenbreugel, 1990). Although one would expect this approach to yield the most robust results, it is not always possible to have such in-depth knowledge of a study site. In cases where knowledge of a specific site is not known, the extraction algorithm must proceed without the benefit of any additional information.

A common feature extraction application is image segmentation. Image segmentation refers to the selection of linear features from an image. Typically, segmentation is used to select road features from an image. It is understood that pixel resolution has a profound effect on the ability of a feature extraction algorithm to pick out specific elements (Van Cleynenbreugel et al. (1990)). By increasing the resolution of the pixel, the feature being sensed is more truly represented, and is, therefore, more easily recognized. Feature extraction algorithms seek out regions of homogeneity. There is far more information in a digital image than can be seen with the naked eye. Image segmentation algorithms are designed to seek out a specific element and identify it as such. The end result is that, often, image artifacts or noise are extracted in addition to the desired elements. This noise can be dealt with through spatial filtering techniques or by selecting elements that meet a certain criteria and subsequently deleting them. In images where regions of homogeneity are fuzzy, a data stratification approach must be adopted.

Abstracting vector data to a raster representation is a different matter, and is generally more straightforward. The procedure in

this case, is to determine whether a portion of a specific vector falls within a particular pixel in a raster grid. Generally, when one needs to integrate vector data into an image analysis environment, the output image size and coordinates are set. That is, the size and complexity of the output is fixed to within a certain number of pixels and lines.

A visual comparison of the results are shown in Table 1 and Table 2. Overall, it was found that bringing a rasterized vector coverage into the raster domain was preferable, in terms of general appearance. That is, the rasterized vectors had the greatest visual appeal when overlayed on an image. In the raster domain, one is bound by the fact that sub-pixel registration is not considered, and that where a vector lies within a pixel is academic. The performance of this method is summarized in Table 1.

When one looks at a vector representation of extracted features, it becomes evident that the edges of the vectors often do not match up. This is largely dependent on the quality of the geographic referencing of the original raster image and the effects of any edge smoothing that was performed on the boundaries. The overall performance of this method is summarized in Table 2.

## 5.0 CONCLUSION

Automated feature extraction techniques can not replace manual digitizing, as of yet. The potential for image segmentation or feature extraction to supplement the job of an operator is certainly there. Automated and semi-automated techniques are desirable to enhance operational turn-around time for getting data through a system. The logical end of this process is a more efficient system for decision-making. As with most things, there are strong elements of give and take, in this case, with respect to image analysis and Geographic Information Systems. The more complicated the data, the greater the demands on the operator to manage the data. Data abstraction is an important aspect of decision making tools, but the user must always keep in mind the accuracy of the information, and hence utility and value of the decisions made.

## REFERENCES

Estes, J.E. and D.S. Simonett, eds. Manual of Remote Sensing: Volume II (Falls Church: American Society of Photogrammetry, 1974) pp. 869-887.

Manore, M., R.J. Brown, K. Korporal and H. Press (1989) "GIS in the Analysis of Satellite Data for Vegetation Monitoring: The Crop Information System" Proceedings of the National Conference on Geographic Information Systems - GIS 89 pp. 118-124.

O'Brien, D (1991) "Computer Assisted Feature Extraction (INTEREX)" Proceedings of the 14th Canadian Symposium on Remote Sensing pp. 423-429.

Schowengerdt, R.A. and R.A. Pries (1988) "interactive Image Feature Compilation for Geographic Information Systems" SPIE - Recent Advances is Sensors, Radiometry, and Data Processing for Remote Sensing Vol. 924. pp. 305-311.

Taylor, T., P. Landriau and A. Scott (1991) "Context Classification for Road Network Extraction from Landsat and SPOT Imagery" Proceedings of the 14th Canadian Symposium on Remote Sensing pp. 413-419.

Van Cleynenbreugel, J., F. Fierens, P. Suetens and A. Oosterlinck (1990) "Delineating Road Structures on Satellite Imagery by a GIS-Guided Technique" Photogrammetric Engineering and Remote Sensing Vol. 56, No. 6. pp. 893-898.

Walsh, Stephen J., D.R. Lightfoot and D.R. Butler (1987) "Recognition and Assessment of Error in Geographic Information Systems" Photogrammetric Engineering and Remote Sensing Vol. 53, No. 10, pp. 1423-1430.

Wang, M. and P.J. Howarth (1991) "Some Generic Issues in Spatial Data Integration" Proceedings of the 14th Canadian Symposium on Remote Sensing pp. 444-447.

Zelek, J.S. (1990) "Computer-Aided Linear Planimetric Feature Extraction" (1990) IEEE Transactions on Geoscience and Remote Sensing Vol. 28, No. 4. pp. 567-572.