

## A General Framework For Error Analysis In Measurement-based GIS — A Summary

Yee Leung

Jiang-Hong Ma

Michael F. Goodchild

Department of Geography and  
Resource Management, Center for  
Environmental Policy and Resource  
Management, and Joint Laboratory for  
Geoinformation Science,  
The Chinese University of Hong Kong,  
Hong Kong  
E-mail: yeeleung@cuhk.edu.hk

Faculty of Science,  
Xi'an Jiaotong University,  
Xi'an, and  
Chang'an University,  
Xi'an, P.R. China  
E-mail:  
jimath@pub.xaonline.com

Department of  
Geography,  
University of California,  
Santa Barbara, California,  
U.S.A.  
E-mail:  
good@ncgia.ucsb.edu

### Abstract

This paper gives a summary of the main results obtained from a series of research on the development of a general framework for error analysis in measurement-based geographic information systems (MBGIS). The study provides a rigorous statistical approach to measurement error analysis and error propagations throughout GIS and spatial operations. It first constructs a basic measurement error model from which relevant concepts such as the approximate law of error propagation, covariance-based error band and maximal allowable limits for positional error are developed. The research then proceeds to the point-in-polygon analysis under measurement errors. Quadratic forms in the joint coordinate vectors are introduced in order to identify whether a point is inside a polygon. They are not only used in point-in-polygon analysis, but are also utilized as the identification of intersection points as well as the computation of length of a line and area of a polygon. An algebra-based probability model is constructed for such a purpose. It is simple but rigorous and can circumvent the complexities surrounding the geometric relations between points and polygons. As a consequence, simple analytic expression and an approximate law for error propagation of intersection points are established, and they are applied to polygon-on-polygon overlay and its error propagation. The general framework facilitates a formal and practical error analysis in MBGIS. It renders a consistent and effective treatment to error propagation in a variety of interrelated GIS and spatial operations

### 1 Introduction

Since the conception of GIS, study of errors in GIS has been rather extensive and diverse (Goodchild and Gopal, 1989; Heurvelink, 1998; Leung and Yan, 1998; Mowrer and Congdon, 2000; Shi et al, 1999; Stanskiwski et al, 1996; Veregin, 1989; Wolf and Ghilani, 1997; Zhang and Goodchild, 2002). It is reckoned that different classes of spatial data exhibit different types of errors. Errors may be introduced during various stages of data manipulation and they may be propagated through GIS or spatial operations, ending up with an output consisting of certain level of inaccuracy.

Errors in spatial databases are complex and multivariate. They can generally be classified into the *inherent error* and the *operational error*. Inherent error is the error present in source documents, including the accumulated error in a map used as input to a GIS. Operational error, categorized as positional and identification errors, is produced through the data capture and manipulation functions of a GIS or is introduced during the process of data entry and occurs throughout data manipulation and spatial modeling. Since GIS databases, at the most general level, are based on a model of geographical data, errors can be classified into spatial, temporal and thematic error. From the modeling point of view, they can also be classified as either systematic or random. Systematic errors usually follow physical laws and they can be

predicted. It is however impossible to avoid entirely random errors in measurements (Wolf and Ghilani, 1997). Such random error is called measurement error (ME), which is one of the most important problems in the use of GIS data.

The epsilon band model of digitizing accuracy, for example, has been used to make estimates of the levels of positional uncertainty and ME that is due to digitizing polygon outlines (Dunn et al, 1990). Although the paper does not develop a comprehensive analytical approach to error for the conversion of map data into digital form, their empirical results on positional accuracy and ME suggest that the interaction of the scale and quality of source documents with the unique process of digitizing may introduce unexpectedly large amounts of error. The main form of this uncertainty is positional error, since it is the location of the polygon boundaries which is uncertain. But this in turn leads to ME, because the estimation of area will then be subjected to a large degree of uncertainty. There is thus an urgent need to develop methods for assessing the accuracy of vector-based GIS. Furthermore, these accuracy assessments need to be incorporated into the spatial analysis procedures, e.g., polygon overlay (the silver polygon problem) and point-in-polygon operations, used in GIS. Compared to raster-based GIS, error analysis in vector-based GIS is much more complex.

It should be noted that a feature of conventional GIS designs is the representation of position by derived coordinates, rather than by original measurements. In such coordinate-based GIS, it is impossible to apply traditional error analysis or estimate uncertainties in derived products. It is because the process of derivation is so complex, and frequently involves human interpretation, and the original measurements are almost always lost during the process. Furthermore, the complexity of the process induces strong spatial autocorrelations between objects in the result. To overcome the difficulties in managing uncertainty intrinsic to conventional GIS, the concept of measurement-based GIS (MBGIS) has been proposed in Goodchild (1999). The basic idea is to retain details of measurements such that error analysis can be made possible, and corrections to positions can be appropriately propagated through the database.

This paper is a summary of a series of research on the rigorous development of a general framework for error analysis in MBGIS with a sound statistical foundation. Figure 1 depicts in brief the main problems of the research series around which our discussion is centered.

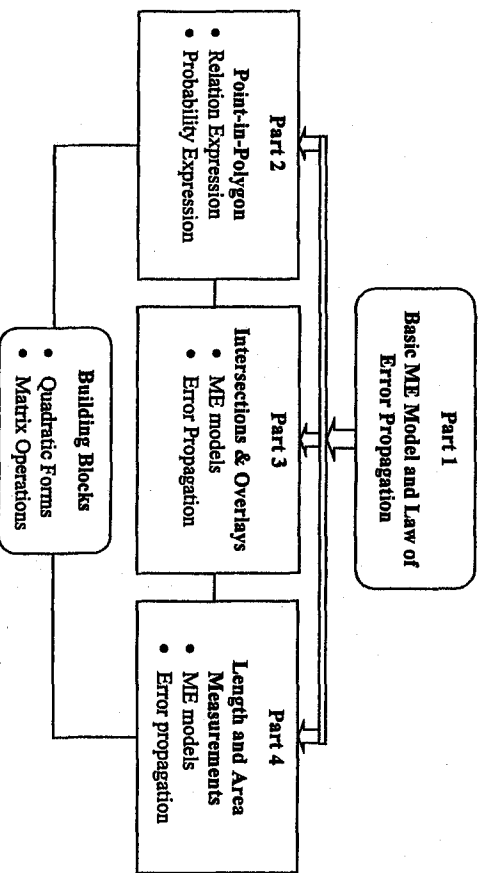


Figure 1. Project structure

## 2 Basic Measurement-Error Model and Related Concepts

This section summarizes the theoretical and empirical results in Leung et al (2003a).

### 2.1 The basic measurement-error model

Let  $\mathbf{x} \in R^p$  be a set of  $p$  measurements which can be directly observed,  $\mathbf{y} \in R^q$  be the quantities to be measured, and  $f$  (a well-defined geographical procedure, called the *operation function*, or more generally, the *transformation function*) be the function linking measurements  $\mathbf{x}$  to quantities  $\mathbf{y}$ , so that

$$\mathbf{y} = f(\mathbf{x}).$$

Once  $\mathbf{x}$  is measured and  $f$  is known,  $\mathbf{y}$  and the associated error distribution are uniquely and automatically determined. The basic ME model can simply be constructed as

$$(1): \begin{cases} \mathbf{Y} = f(\mathbf{X}), \\ \mathbf{X} = \mu_x + \epsilon_x, \quad \epsilon_x \sim (0, \Sigma_x), \end{cases} \quad (2.1)$$

where  $\mu_x$  is the true value vector,  $\mathbf{X}$  the random measurement value vector,  $\mathbf{Y}$  the indirect measurement value vector obtained through  $f$ , and  $\epsilon_x$  the random ME vector with zero mean  $\mathbf{0}$  and the variance-covariance matrix  $\Sigma_x$ . According to (2.1),  $\mathbf{Y}$  is random and its error variance-covariance matrix  $\Sigma_y$  is propagated by  $\Sigma_x$ . To simplify, the variance-covariance matrix is called the covariance matrix henceforth.

Furthermore, if a new vector  $\mathbf{z} \in R^r$  can be obtained through  $\mathbf{y}$  and a known function  $g(\cdot)$ , model (1) can again be utilized for ME analysis. As a result, error propagates along the process of operations via the basic ME model taking the general form of (1) (Figure 2).

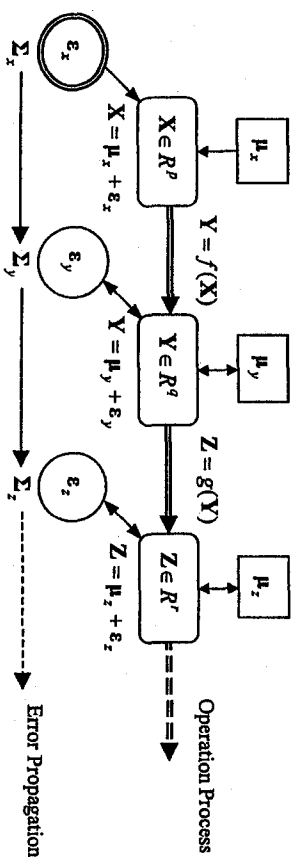


Figure 2. A flowchart for error propagation based on the basic ME model (1)

### 2.2 Approximate law of error propagation

In practice, the error distribution of  $\mathbf{Y}$  can seldom be derived exactly from the relationship  $\mathbf{Y} = f(\mathbf{X})$ . When the error distribution cannot be exactly obtained or is too complex to be derived because of nonlinearity of  $f$ , a common method is to employ a linear (or second-order) approximation to  $\mathbf{Y} = f(\mathbf{X})$ . Then, for nonlinear  $f$ , we have the following *approximate law of error propagation*:

$$\Sigma_y \approx \Sigma_x \mathbf{B}_\mu \mathbf{B}_\mu^T, \quad (2.2)$$

where  $\mathbf{B}_\mu$  is a *Jacobian matrix* of  $f$  at  $\mu_x$ . That is,  $\Sigma_y$  can be approximately expressed by  $\Sigma_x$  and  $\mathbf{B}_\mu$ . In general, we have the following law of error propagation:

$$(II): \Sigma_{\mu} \begin{cases} \approx \Sigma_{\mu} = B_{\mu} \Sigma_{\mu} B_{\mu}^T, & \text{if } f(\mathbf{x}) \text{ is nonlinear in } \mathbf{x}, \\ = B \Sigma_{\mu} B^T, & \text{if } f(\mathbf{x}) \text{ is linear in } \mathbf{x}, \text{ i.e., } f(\mathbf{x}) = \mathbf{a} + B\mathbf{x}, \end{cases}$$

where  $\mathbf{a}$  and  $B$  are constants.

Although the term "law of error propagation" has been commonly used in the literature, it often lacks a rigorous mathematical expression and its usage is often unclear or sometimes incorrect. We especially stress in here the approximation and local property of this law for a nonlinear transformation function  $f$ . The emphasis on "approximation" is necessary in order to apply correctly the first-order Taylor approximation to derive (2.2). In addition, what should be pointed out is its local property as such an approximation is effective only in the local neighborhood of  $\mu_x$ . Different from the notation adopted in Wolf and Ghilani (1997), Heuvelink (1998) and other publications, the Jacobian matrix  $B_{\mu}$  is appended with the subscript  $\mu$  to reflect its dependence on the local point  $\mu$ .

Based on the basic model (I) and the error propagation law (II), we have made three simple applications in geodesy: (1) measurement with a distance and a direction, (2) measurement with two distances, and (3) measurement with two angles, and derived their approximate laws for error propagation (see Leung et al (2003a) for details).

### 2.3 Covariance-based error bands

Uncertainty of a point can be derived from the covariance matrix associated with it and can be presented by an error ellipse. There are many uncertainty models for a line, such as the most commonly used epsilon band model (Perkal, 1956, 1966; Blakemore, 1984), the error band model (Shi, 1994, 1999), the  $\epsilon_{\sigma}$  and  $\epsilon_m$  error band models (Tong et al, 1999), and the positional uncertainty model (Alesheikh and Li, 1996; Alesheikh et al, 1999). Some of these models are either simple with no statistical rigor, or they are too complex to be lucid and effective. We formulate in this study a simple but statistically strict and unified error band model, named the "covariance-based error band", which contains error structures of the endpoints of a line segment. This turns out to be a natural and effective error band concept.

Since a line consists of points, uncertainty (i.e. the covariance matrix) of any point on the line should naturally reflect that of the line. For  $i=1, 2$ , let  $X_i = (X_{i1}, X_{i2})^T$ ,  $\mu_i = (\mu_{i1}, \mu_{i2})^T$ , and  $\epsilon_i = (\epsilon_{i1}, \epsilon_{i2})^T$  be respectively the random, true, and ME vectors of the endpoint coordinates of the segment  $V_i V_2$ . To generalize our discussion, we consider the joint ME vector  $\epsilon_{(2)} = (\epsilon_1^T, \epsilon_2^T)^T$  and let its covariance matrix be  $\Sigma_{(2)}$ , i.e.,

$$\Sigma_{(2)} \equiv \text{COV}(\epsilon_{(2)}) = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \quad \Sigma_{ij} \equiv \text{COV}(\epsilon_i, \epsilon_j), \quad i, j = 1, 2,$$

where the subscript (2) indicates that we have two endpoints.

Thus for any point  $V'$  in the line segment  $V_1 V_2$ , the covariance matrix of its coordinate vector  $X' = tX_1 + (1-t)X_2$  is derived as

$$\Sigma_{t'} \equiv \text{COV}(X') = t^2 \Sigma_{11} + (1-t)^2 \Sigma_{22} + t(1-t)(\Sigma_{12} + \Sigma_{21}), \quad 0 \leq t \leq 1.$$

If we assume that  $\epsilon_{(2)}$  is distributed as a normal distribution, then  $P[X_i \in R_i^{(\alpha)}] = 1 - \alpha$ , with

$$R_i^{(\alpha)} \equiv \{x: (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \leq \chi_{2,\alpha}^2\}, \quad i=1, 2,$$

where  $\chi_{2,\alpha}^2$  is the upper  $\alpha$ -quantile of the chi-squares distribution with 2 degrees of freedom. Based on  $R_1^{(\alpha)}$  and  $R_2^{(\alpha)}$ , we can construct a larger region:

$$R_{(1,2)}^{(\alpha)} \equiv \{x: \text{there is a real number } t (0 \leq t \leq 1) \text{ such that } (x - \mu_t)^T \Sigma_t^{-1} (x - \mu_t) \leq \chi_{2,\alpha}^2\},$$

where  $\mu_t \equiv t\mu_1 + (1-t)\mu_2$ . It is obvious that  $R_{(1,2)}^{(\alpha)}$  is the union of the confidence regions (ellipses) of all points on the line segment and is called the *covariance-based error band* for the line segment. It can be observed that when  $\Sigma_{11} = \Sigma_{22} = \Sigma_{12} = \Sigma_{21} = \epsilon^2 I_2$ ,  $R_{(1,2)}^{(\alpha)}$  becomes the epsilon band model; when  $\Sigma_{11} = \Sigma_{22} = \sigma^2 I_2$  and  $\Sigma_{12} = \Sigma_{21} = 0$ ,  $R_{(1,2)}^{(\alpha)}$  becomes the error band model. For different  $\Sigma_{(2)}$ , the covariance-based error band take on different shapes. Detailed theoretical analysis and simulation experiments can be found in Leung et al (2003a).

### 2.4 Maximal allowable limits for positional error

It is well known that one of the most important characteristics of spatial objects is their topological (or geometrical) property. Although we can tolerate a certain degree of ME in locations, ME must be restricted to a certain range so that it will not distort the original topology. That is, it is necessary to give a ME limit in order to guarantee that the original topology of a spatial object or topological relationship among spatial objects is invariant under ME. We call such a limit the *maximal allowable limit* (MAL). Large ME may change the topology and geometrical property of a spatial object and result in logical inconsistency or reality distortion. A simple original polygon, for example, can be distorted into a complex one if error goes above the MAL.

Let  $\sigma_{1i}^2$  and  $\sigma_{2i}^2$  be respectively the variances of the ME  $\epsilon_{1i}$  and  $\epsilon_{2i}$  of a vertex  $V_i$ ,  $i=1, \dots, n$ . Within our framework, the maximal allowable limits  $\bar{\sigma}_{1i}$  and  $\bar{\sigma}_{2i}$  for  $V_i$  are defined as the maximal allowable values of  $\sigma_{1i}$  and  $\sigma_{2i}$  so that the confidence ellipse of  $V_i$  does not intersect the covariance-based error bands of its disjoint edges. Since there is a confidence level attached to the covariance-based error band, it may also be treated as the confidence level of the MAL. Thus, as long as the variances  $\sigma_{1i}^2$  and  $\sigma_{2i}^2$  of ME in  $V_i$  are not greater than  $\bar{\sigma}_{1i}^2$  and  $\bar{\sigma}_{2i}^2$  respectively, the observed polygon in the presence of ME will not change the topological property of the true polygon with the confidence level, and their difference results only from errors in position of the vertices.

## 3 An Algebra-based Probability Model for Point-in-Polygon Analysis

Point-in-polygon, line-in-polygon, and polygon-on-polygon analyses are fundamental but important problems in vector-based GIS. They can arise in spatial queries or overlay operations. The basic starting point is point-in-polygon analysis (Leung and Yan, 1997, 1998; Rigaux et al, 2002). Because of the uncertainty about spatial objects such as points and polygons, the analysis becomes rather difficult and complex, especially from the geometric point of view.

How to perform point-in-polygon analysis under ME is thus a basic problem in MBGIS. Except for the conceptual discussion in Leung and Yan (1997), in-depth theoretical analysis of the issue when points and polygons are both in error has not been made. We construct in Leung et al (2003b) an algebra-based probability model for point-in-polygon analysis under ME and provide a concise and unified expression for computing the probability that a point is inside a polygon (concave or convex). Our basic idea is to first convert a point-in-polygon analysis into several point-in-triangle problems. For each triangle, the corresponding probability can then be concisely computed. The algebraic approach is not only simple and effective, it can also circumvent complications and difficulties encountered in geometric solutions of the problem. This section summarizes the theoretical and empirical results in Leung et al (2003b).

### 3.1 The triangle model

Under the effect of ME, let three vertex positions of the underlying triangle  $\Delta Y_1 Y_2 Y_3$  be  $Y_i(X_i)$ ,  $i=1,2,3$ , the corresponding coordinate column vectors are  $X_i = (X_{i1}, X_{i2})^T \in R^2$ ,  $i=1,2,3$ , and  $X_0 = (X_{01}, X_{02})^T \in R^2$  is the coordinate column vector of any point  $Y$ . Define  $X_{(4)}^T = (X_1^T X_2^T X_3^T X_0^T)$ , and  $H_1 = \Delta_1 \otimes H_0$ ,  $H_0 = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$ ,  $i=1,2,3,4$ ,

$$Z_i = X_{(4)}^T H_i X_{(4)},$$

$$\Delta_1 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 0 & 1 \end{pmatrix}, \Delta_2 = \begin{pmatrix} 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 \\ -1 & 0 & 1 & 0 \end{pmatrix}, \Delta_3 = \begin{pmatrix} 0 & 1 & 0 & -1 \\ -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{pmatrix}, \Delta_4 = \begin{pmatrix} 0 & 1 & -1 & 0 \\ -1 & 0 & 1 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

The probability that the point  $Y$  is inside the triangle region  $R(Y_1 Y_2 Y_3)$  is obtained as

$$P\{Y \in R(Y_1 Y_2 Y_3)\} = P[Z_{i_1}, i=1,2,3,4, \text{ have the same sign}], \quad (3.1)$$

or  $P\{Y \in R(Y_1 Y_2 Y_3)\} = P[Z_{\Delta} \geq 0]$ ,

where  $Z_{\Delta} = Z_{\min} Z_{\max}^{-1} \geq 0$ ,  $Z_{\min} = \min\{Z_i\}$ , and  $Z_{\max} = \max\{Z_i\}$ .

### 3.2 The polygon model

Since triangles are the most essential and simple geometric objects and any polygon (whether it is convex or concave) can always be decomposed into triangles (Berg et al, 2000), we can transform the polygon problem into a multiple triangle problem on the basis of the triangle model (3.1). Thus, unlike any geometric method, the algebraic model makes the convexity or non-convexity of a polygon an non-issue.

Let  $Y$  be an uncertain point with ME,  $R(Y_1 Y_2 \dots Y_n)$  be a point set representing the region of a polygon with vertices  $Y_1, Y_2, \dots, Y_n$ . Let a partition of a simple polygon region  $R(Y_1 Y_2 \dots Y_n)$  be

$$R(Y_1 Y_2 \dots Y_n) = R(Y_{11} Y_{12} Y_{13}) + R(Y_{21} Y_{22} Y_{23}) + \dots + R(Y_{m1} Y_{m2} Y_{m3}),$$

where  $Y_{11}, Y_{12}, Y_{13}$  are the vertices of the triangle corresponding to  $\Delta Y_{11} Y_{12} Y_{13}$ . Then we can derive the following probability formula for point-in polygon analysis:

$$P\{Y \in R(Y_1 Y_2 \dots Y_n)\} = P\{Y \in R(Y_{11} Y_{12} Y_{13})\} + P\{Y \in R(Y_{21} Y_{22} Y_{23})\} + \dots + P\{Y \in R(Y_{m1} Y_{m2} Y_{m3})\} \\ = P[Z_{\Delta_1} \geq 0] + P[Z_{\Delta_2} \geq 0] + \dots + P[Z_{\Delta_m} \geq 0], \quad (3.2)$$

where  $Z_{\Delta_i} = Z_{\min_i} Z_{\max_i}^{-1}$  is the random variable corresponding to  $\Delta Y_{i1} Y_{i2} Y_{i3}$ .

The results can be applied to polygons with holes or objects consisting of several non-connected polygons. It should be noted that although the number of triangles resulted from triangulation may be large, the number of triangles necessary to the point-in-polygon analysis can usually be reduced since triangles outside the allowable maximal range of a point do not contribute to the probability in (3.2).

## 4 Error Analysis for Intersections and Overlays

Besides point-in-polygon analysis, line-in-polygon and polygon-on-polygon operations are also important in GIS. When a polygon is convex, line-in-polygon operation can be reduced to point-in-polygon operation since a line segment is inside a polygon as long as its endpoints are inside the polygon. When a polygon is non-convex, however, it involves the problem of knowing whether a line segment intersects

with the polygon. Moreover, polygon-on-polygon operations are also closely related to the intersection problems. Therefore, error analysis for intersection points is a key problem in many GIS operations. This section summarizes the theoretical and empirical results in Leung et al (2003c).

### 4.1 Approximate law for error propagation in intersection coordinates

Under ME, denote by  $X_{(4)}^T = (X_1^T X_2^T X_3^T X_0^T)$  the joint coordinate vector of the endpoints for two line segments  $Y_1(X_1) Y_2(X_2)$  and  $Y_3(X_3) Y_4(X_4)$ ,  $\bar{H} = \bar{\Delta} \otimes H_0$ ,  $\bar{H}_i = \bar{\Delta}_i \otimes H_0$ ,  $\bar{\Delta}_i = \Delta_i$ ,  $i=1,2,4$ ,  $\bar{\Delta}_3 = -\Delta_3$ , and

$$\bar{\Delta} = \begin{pmatrix} 0 & 0 & -1 & 1 \\ 0 & 0 & 1 & -1 \\ 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \end{pmatrix}, \text{ and } \lambda_i(X_{(4)}) = \frac{X_{(4)}^T \bar{H}_i X_{(4)}}{X_{(4)}^T \bar{H} X_{(4)}}, \quad i=1,2,3,4.$$

Then the ME model for the position vector  $X_c$  of the intersection point  $Y_c$  of two line segments  $Y_1 Y_2$  and  $Y_3 Y_4$  is

$$\begin{cases} X_c = f(X_{(4)}) = [\lambda_1(X_{(4)}) D_1 + \lambda_2(X_{(4)}) D_2] X_{(4)}, \\ X_{(4)} = \mu_{(4)} + \varepsilon_{(4)}, \quad \varepsilon_{(4)} \sim (0, \Sigma_{\varepsilon}), \end{cases} \quad (4.1)$$

where  $D_i = e_{4i}^T \otimes I_2$  is a  $2 \times 8$  matrix,  $i=1,2$ ,  $e_{4i}$  is a  $n \times 1$  unit column vector whose the  $i$ th component is 1 and 0 elsewhere,  $\mu_{(4)} = (\mu_1^T \mu_2^T \mu_3^T \mu_4^T)$  and  $\varepsilon_{(4)} = (\varepsilon_1^T \varepsilon_2^T \varepsilon_3^T \varepsilon_4^T)$  are respectively the corresponding true and ME vectors.

Since the transformation function  $f$  in (4.1) is nonlinear in  $X_{(4)}$ , an approximate law of error propagation for  $X_c$  is obtained as

$$\Sigma_c \approx \bar{\Sigma}_c = B_{\mu} \Sigma_{\varepsilon} B_{\mu}^T,$$

where  $B_{\mu}$  is a  $2 \times 8$  matrix given by

$$B_{\mu} = \frac{1}{2} \begin{bmatrix} D_1 X_{(4)} X_{(4)}^T \bar{H}_1 + D_2 X_{(4)} X_{(4)}^T \bar{H}_2 \\ X_{(4)}^T \bar{H} X_{(4)} \\ + [\lambda_1(X_{(4)}) D_1 + \lambda_2(X_{(4)}) D_2] \left[ I_2 - \frac{2}{X_{(4)}^T \bar{H} X_{(4)}} X_{(4)} X_{(4)}^T \bar{H} \right] \end{bmatrix}.$$

### 4.2 Approximate law for error propagation in polygon-on-polygon analysis

Due to its methodological and technical complexities, developing error propagation models for overlay operations, especially on vector-based data, has seldom been attempted. Error propagation in polygon-on-polygon operation is a typical example. A formal model is constructed in Leung et al (2003c) and is briefly described as follows:

Consider directed polygons  $P$  with  $n_i$  vertices  $Y_j^{(i)}$  listed in a counter-clockwise manner:  $Y_1^{(i)}, \dots, Y_{n_i}^{(i)}, Y_{n_i+1}^{(i)}, \dots, Y_{n_i+1}^{(i)} = Y_1^{(i)}$ ,  $i=1,2$ , where in the superscript  $(i)$  denotes the  $i$ th polygon and the subscript  $j$  denotes the  $j$ th vertex. Assume that vertices  $Y_j^{(i)}$  have the coordinates vectors  $X_j^{(i)}$  under ME, their true coordinates vectors are  $\mu_j^{(i)}$ , and the corresponding ME vectors are  $\varepsilon_j^{(i)}$ ,  $j=1, \dots, n_i$ ,

$i=1, 2$ . Let  $\epsilon_{(1,2)}$ ,  $\epsilon_{(2,1)}$ ,  $\epsilon_{(1,2)}$ ,  $\epsilon_{(1,1)}$ , and  $\epsilon_{(2,2)}$  be respectively the joint ME vectors of vertices coordinates counter-clockwisely listed, of the overlaid polygons  $P_1 - P_2$ ,  $P_2 - P_1$ ,  $P_1 \cap P_2$ , and the original input polygons  $P_1$  and  $P_2$ . Denote the joint input ME vector by  $\epsilon = (\epsilon_{(1,1)}^T, \epsilon_{(2,2)}^T)^T$ .

Employing the results in the approximate law for error propagation in intersection points, we can express approximately the joint output ME vectors  $\epsilon_{(1,2)}$ ,  $\epsilon_{(2,1)}$ , and  $\epsilon_{(1,2)}$  as

$$\tilde{\epsilon}_{(1,2)} = \mathbf{D}_{(1,2)} \epsilon, \quad \tilde{\epsilon}_{(2,1)} = \mathbf{D}_{(2,1)} \epsilon, \quad \tilde{\epsilon}_{(1,2)} = \mathbf{D}_{(2,1)} \epsilon,$$

where the determinations of  $\mathbf{D}_{(1,2)}$ ,  $\mathbf{D}_{(2,1)}$  and  $\mathbf{D}_{(2,1)}$  depend on the intersection points. Thus, we can obtain the following approximate laws of error propagation in polygon-on-polygon operation:

$$\tilde{\Sigma}_{(1,2)} = \mathbf{D}_{(1,2)} \Sigma \mathbf{D}_{(1,2)}^T, \quad \tilde{\Sigma}_{(1,2)} = \mathbf{D}_{(1,2)} \Sigma \mathbf{D}_{(1,2)}^T, \quad \tilde{\Sigma}_{(2,1)} = \mathbf{D}_{(2,1)} \Sigma \mathbf{D}_{(2,1)}^T.$$

## 5 Error Analysis in Length and Area Measurements

In addition to the study of absolute errors for a single location, we, on the basis of the discussions in Kivert (1997) and Hunter and Goodchild (1996), and in light of the locational error models advanced in Leung and Yan (1998), have also performed detailed analysis of relative errors in length and area measurements in MBGIS. This section summarizes the theoretical and empirical results in Leung et al (2003d).

### 5.1 Error analysis in length measurements

Let  $Y_1(\mathbf{X}_1)$  and  $Y_2(\mathbf{X}_2)$  be the endpoints of a line segment under ME,  $\mathbf{X}_1^T = (\mathbf{X}_1^T, \mathbf{X}_2^T)$ , and  $L$  the length of the line segment. Let  $\mathbf{G}_{(2)} \equiv \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \otimes \mathbf{I}_2$ . Then the ME model for length measurement becomes

$$\begin{cases} L = f(\mathbf{X}_{(2)}) = (\mathbf{X}_{(2)}^T \mathbf{G}_{(2)} \mathbf{X}_{(2)})^{1/2} \\ \mathbf{X}_{(2)} = \mu_{(2)} + \epsilon_{(2)}, \quad \epsilon_{(2)} \sim (0, \Sigma_{(2)}). \end{cases}$$

It is shown in Leung et al (2003d) that  $L^2$  can be represented as

$$L^2 = \sum_{i=1}^m \lambda_i x_{i, \delta_i}^2,$$

where  $\lambda_i$  are the distinct non-zero eigenvalues of  $\mathbf{G}_{(2)} \Sigma_{(2)}$ ,  $P_i$  their respective orders of multiplicity, and  $\delta_i^2$  the quantities determined by  $\Sigma_{(2)}$  and  $\mu_{(2)}$ . Theoretically, the exact law for error propagation in  $L$  can be obtained by the statistical methods. For simplicity, we can give its approximate law of error propagation as:

$$\sigma_L^2 = \mathbf{B}_{\mu_{(2)}} \Sigma_{(2)} \mathbf{B}_{\mu_{(2)}}^T, \quad \text{where } \mathbf{B}_{\mu_{(2)}} = (\mu_{(2)}^T \mathbf{G}_{(2)} \mu_{(2)})^{-1/2} \mu_{(2)}^T \mathbf{G}_{(2)}.$$

We can also analyze the error of the perimeter of a polygon propagated from the ME of the vertices coordinates. Assume that the vertices of a simple  $n$ -sided polygon are  $Y_i(\mathbf{X}_i)$ ,  $i=1, \dots, n$ ,  $Y_{n+1} = Y_1$ . Let  $\mathbf{X}_{(n)}^T = (\mathbf{X}_1^T, \dots, \mathbf{X}_n^T)$  and  $\mathbf{G}_{(n)} \equiv [(\mathbf{e}_{n,i} - \mathbf{e}_{n,i+1})^T (\mathbf{e}_{n,i} - \mathbf{e}_{n,i+1})] \otimes \mathbf{I}_2$ .

The ME model for the perimeter measurement is

$$\begin{cases} L_{(n)} = f(\mathbf{X}_{(n)}) = \sum_{i=1}^n (\mathbf{X}_{(n)}^T \mathbf{G}_{(2)} \mathbf{X}_{(n)})^{1/2} \\ \mathbf{X}_{(n)} = \mu_{(n)} + \epsilon_{(n)}, \quad \epsilon_{(n)} \sim (0, \Sigma_{(n)}), \end{cases}$$

and the corresponding approximate law of error propagation is

$$\sigma_{L_{(n)}}^2 = \mathbf{B}_{\mu_{(n)}} \Sigma_{(n)} \mathbf{B}_{\mu_{(n)}}^T, \quad \text{where } \mathbf{B}_{\mu_{(n)}} = \sum_{i=1}^n [\mu_{(n)}^T \mathbf{G}_{(2)} \mu_{(n)}]^{-1/2} \mu_{(n)}^T \mathbf{G}_{(2)}.$$

### 5.2 Error analysis in area measurements

For area measurement, most of the research in the literature is on raster-based data (Lloyd, 1976). Under the assumption that error at each vertex is independently and identically distributed around the local mean, effects of point error on area calculations has been reported in Chrisman and Yandell (1988). For more results, see Goodchild and Gopal (1989).

We consider in here an  $n$ -sided polygon with the vertices  $Y_i(\mathbf{X}_i)$  under the ME vectors  $\epsilon_i \equiv (\epsilon_{i1}, \epsilon_{i2})^T$ , where  $\mathbf{X}_i \equiv (X_{i1}, X_{i2})^T$ ,  $i=1, \dots, n$ ,  $Y_{n+1} = Y_1$ . The signed area of this polygon is given by

$$A_{(n)} = \frac{1}{2} \sum_{i=1}^n (X_{i1} X_{i+1,2} - X_{i2} X_{i+1,1}).$$

Let  $\mathbf{X}_{(n)} \equiv (\mathbf{X}_1^T, \dots, \mathbf{X}_n^T)^T$ ,  $\mu_{(n)} \equiv (\mu_1^T, \dots, \mu_n^T)^T$ ,  $\epsilon_{(n)} \equiv (\epsilon_1^T, \dots, \epsilon_n^T)^T$ ,  $\mathbf{H}_{(n)} \equiv \Delta_{(n)} \otimes \mathbf{H}_0$ , and

$$\Delta_{(n)} = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & -1 \\ -1 & 0 & 1 & \dots & 0 & 0 \\ 0 & -1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \dots & 0 & 0 \\ 1 & 0 & \dots & \dots & 0 & -1 \end{pmatrix}, \quad n \geq 3$$

Accordingly, the ME model for the area measurement of a polygon becomes

$$\begin{cases} A_{(n)} = f(\mathbf{X}_{(n)}) = \frac{1}{2} \mathbf{X}_{(n)}^T \mathbf{H}_{(n)} \mathbf{X}_{(n)}, \\ \mathbf{X}_{(n)} = \mu_{(n)} + \epsilon_{(n)}, \quad \epsilon_{(n)} \sim (0, \Sigma_{(n)}). \end{cases} \quad (5.4)$$

Assume that  $\epsilon_{(n)} \sim N_{2n}(0, \Sigma_{(n)})$ . Then the area of the polygon  $A_{(n)}$  in (5.1) can be represented as

$$A_{(n)} = \frac{1}{4} \sum_{i=1}^m \lambda_i x_{i, \delta_i}^2,$$

where  $\lambda_i \equiv \lambda_i(\mathbf{H}_{(n)} \Sigma_{(n)})$  are the distinct non-zero eigenvalues of  $\mathbf{H}_{(n)} \Sigma_{(n)}$ ,  $P_i$  their respective orders of multiplicity, and  $\delta_i^2 \equiv \delta_i^2(\Sigma_{(n)}, \mu_{(n)})$  the quantities determined by  $\Sigma_{(n)}$  and  $\mu_{(n)}$ . The resulting exact law of error propagation is obtained as

$$\sigma_{A_{(n)}}^2 = \text{Var}(A_{(n)}) = F_n(\Sigma_{(n)}; \mu_{(n)}) \equiv \frac{1}{8} \sum_{i=1}^m \lambda_i^2 (\mathbf{H}_{(n)} \Sigma_{(n)}) [P_i + 2\delta_i^2(\Sigma_{(n)}, \mu_{(n)})].$$

This equation indicates that the variance  $\sigma_{A_{(n)}}^2$  of area measurement for polygons is not only related to the covariance matrix  $\Sigma_{(n)}$  of the locational error of the vertices, but is also related to the true locations  $\mu_{(n)}$  of the vertices. This conclusion may give a valuable insight on area measurements in general.

## 6 Conclusion

Corresponding to Figure 1, the main results in Leung et al (2003a-d) are summarized in Figure 3. The theoretical underpinnings and empirical analyses of the present study provide a sound and formal basis for error analysis in MBGIS. It will play an even more important role when future GIS are no longer digitized maps but maps constructed directly with multi-source data with varying measurement accuracies.

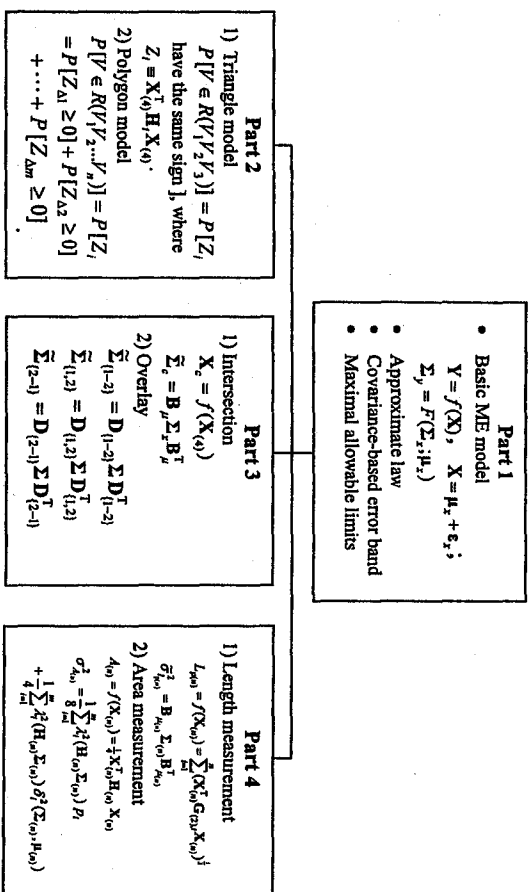


Figure 3. Basic results

## Acknowledgement

This project was supported by the earmarked grant CUHK 4362/00H of The Hong Kong Research Grants Council.

## References

- Aleskikh, A.A., J.A.R. Blais, M.A. Chapman, and H. Kartim. 1999. Rigorous geospatial data uncertainty models for GIS. In *Spatial Accuracy Assessment: Land Information Uncertainty in Natural Resources*, K. Lowell and A. Jaton (eds), pp. 195-202. Chelsea, Michigan: Ann Arbor Press.
- Aleskikh, A.A., and R. Li. 1996. Rigorous uncertainty models of line and polygon objects in GIS. *Proceedings of GIS/SLIS '96*, Denver, CO, pp. 906-920.
- Berg, M. de, M. van Kreveld, M. Overmars, and O. Schwartzkopf. 2000. *Computational Geometry: Algorithms and Applications* (2nd ed.). Berlin: Springer-Verlag.
- Blakemore, M. 1984. Generalization and error in spatial data bases. *Cartographica*, 21, 131-139.
- Chrisman, N.R., and B.S. Yandell. 1988. Effects of point error on area calculations: a statistical model. *Surveying and Mapping*, 48(4): 241-246.
- Cressie, N.A.C. 1993. *Statistics for Spatial Data, Revised Edition*. John Wiley & Sons, New York.

- Dunn, R., A.R. Harrison, and J.C. White. 1990. Positional accuracy and measurement error in digital databases of land use: an empirical study. *Int. J. Geographical Information Systems*, 4(4): 385-398.
- Goodchild, M.F. and S. Gopal. (Eds). 1989. *Accuracy of Spatial Databases*. London: Taylor & Francis.
- Goodchild, M.F. 1999. Measurement-based GIS. In Shi, W.Z., Goodchild, M.F. and Fisher, P.F. (Eds), *Proceedings of the International Symposium on Spatial Data Quality '99*, Hong Kong: Hong Kong Polytechnic University, 1-9.
- Heuvelink, G.B.M. 1998. *Error Propagation in Environmental Modelling with GIS*. London: Taylor & Francis.
- Hunter, G.J. and M.F. Goodchild. 1996. A new model for handling vector data uncertainty in geographical information systems. *Journal of the Urban and Regional Information Systems Association*, 8(1), 51-57.
- Klivert, H.T. 1997. Assessing, representing and transmitting positional uncertainty in maps. *Int. J. Geographical Information Science*, 11(1): 33-52.
- Leung, Y., and J.P. Yan. 1999. Point-in-polygon analysis under certainty and uncertainty. *Geographical Information Science*, 13(1): 93-114.
- Leung, Y., and J.P. Yan. 1998. A locational error model for spatial features. *Int. J. Geographical Information Science*, 12, 607-620.
- Leung, Y., J.H. Ma, and M.F. Goodchild. 2003a. A general framework for error analysis in measurement-based GIS—Part 1: the basic measurement-error model and related concepts. (unpublished paper)
- Leung, Y., J.H. Ma, and M.F. Goodchild. 2003b. A general framework for error analysis in measurement-based GIS—Part 2: the algebra-based probability model for point-in-polygon analysis. (unpublished paper)
- Leung, Y., J.H. Ma, and M.F. Goodchild. 2003c. A general framework for error analysis in measurement-based GIS—Part 3: error analysis for intersections and overlays. (unpublished paper)
- Leung, Y., J.H. Ma, and M.F. Goodchild. 2003d. A general framework for error analysis in measurement-based GIS—Part 4: error analysis for length and area measurements. (unpublished paper)
- Lloyd, P.R. 1976. Quantisation error in area measurement. *The Cartographic Journal*, 13(1), 22-25.
- Movwer, H. T., and R. G. Congelton. (eds). 2000. *Quantifying Spatial Uncertainty in Natural Resources: Theory and Applications for GIS and Remote Sensing*. Chelsea (Michigan): Ann Arbor Press.
- Perkal, J. 1956. On epsilon length. *Bulletin de l'Académie Polonaise des Sciences*, 4: 399-403.
- Perkal, J. 1966. On the length of empirical curves. Discussion Paper Number 10, Michigan Inter-University Community of Mathematical Geography.
- Rigaux, P., M. Scholl, and A. Voisard. 2002. *Spatial Databases with Application to GIS*. San Francisco: Morgan Kaufmann Publishers.
- Shi, W. 1994. *Modeling Positional and Thematic Uncertainties in Integration of Remote Sensing and GIS*. Ph.D. Thesis, ITC Publication No. 22. ITC, The Netherlands.
- Shi, W., M. Ehlers, and K. Tempfli. 1999. Analytical modeling positional and thematic uncertainties in the integration of remote sensing and geographical information systems. *Transactions in GIS*, 3(2): 119-136.
- Stanislowski, L.V., B.A. Dewitt, and R. S. Shrestha. 1996. Estimating positional accuracy of data layers within a GIS through error propagation. *Photogrammetric Engineering and Remote Sensing*, 62(4), 429-433.
- Tong, X., W. Shi, and D. Liu. 1999. Development of error models for transition curves in GIS. In Shi, W.Z., Goodchild, M.F. and Fisher, P.F. (Eds), *Proceedings of the International Symposium on Spatial Data Quality '99*, Hong Kong: Hong Kong Polytechnic University, 299-307.
- Vergin, H. 1989. A taxonomy of error in spatial databases. *Technical Paper 89-12, National Center for Geographic Information & Analysis*, Geography Department, University of California, Santa Barbara, California.
- Wolf, P.R., and C.D. Ghilani. 1997. *Adjustment Computations: Statistics and Least Squares in Surveying and GIS*. New York: John Wiley & Sons, Inc.
- Zhang, J., and M.F. Goodchild. 2002. *Uncertainty in Geographical Information*. New York: Taylor and Francis.