

Kansas State University Libraries

New Prairie Press

---

Conference on Applied Statistics in Agriculture

1997 - 9th Annual Conference Proceedings

---

## THE USE OF INVERSE THEORY ON AN ILL-POSED ENVIRONMENTAL COMPOSITE SAMPLING PROBLEM

V. A. Lancaster

S. Keller-McNulty

Follow this and additional works at: <https://newprairiepress.org/agstatconference>



Part of the [Agriculture Commons](#), and the [Applied Statistics Commons](#)



This work is licensed under a [Creative Commons Attribution-Noncommercial-No Derivative Works 4.0 License](#).

---

### Recommended Citation

Lancaster, V. A. and Keller-McNulty, S. (1997). "THE USE OF INVERSE THEORY ON AN ILL-POSED ENVIRONMENTAL COMPOSITE SAMPLING PROBLEM," *Conference on Applied Statistics in Agriculture*. <https://doi.org/10.4148/2475-7772.1308>

This is brought to you for free and open access by the Conferences at New Prairie Press. It has been accepted for inclusion in Conference on Applied Statistics in Agriculture by an authorized administrator of New Prairie Press. For more information, please contact [cads@k-state.edu](mailto:cads@k-state.edu).

## The Use of Inverse Theory on an Ill-Posed Environmental Composite Sampling Problem

V. A. Lancaster and S. Keller-McNulty, Department of Statistics, Kansas State University, Manhattan, KS 66506

**Abstract:** *As an alternative to retesting, the use of inverse theory techniques is proposed to resolve the lack of information inherent in composite sampling methods. This paper evaluates the feasibility of combining composite sampling with the inverse theory technique of linear regularization on an environmental site characterization investigation. Federal legislation has mandated the cleanup of hazardous waste sites, creating the need to characterize these sites for various chemical constituents. An abundance of samples, high measurement costs, and limited budgets create the appeal of compositing samples. We propose that the number of costly laboratory analyses can be reduced by combining composite sampling and inverse theory techniques. The goal of the paper is to estimate the constituent concentration for the sample units used to construct the composite.*

*A novel application of linear regularization is illustrated with a data set from an environmental investigation into mercury contamination at a waste disposal site in New Mexico. The disposal site has measurements at both the composite and sample unit level. This allows for a rare opportunity to evaluate the assumptions made about the compositing process and to compare estimates based on composite measurements with estimates based on sample unit measurements.*

### 1. Introduction

In recent years there has been renewed interest in composite sampling which has been prompted by demands from within the environmental community. Federal legislation has mandated the monitoring and cleanup of hazardous waste sites across the country. Before a remediation decision can be made, a waste site must be characterized for the various chemical constituents that are suspected to be present. Once remediation is complete, additional sampling and analyses are required to verify that the clean-up criteria have been met. What these investigations have in common are an abundance of samples, high measurement costs, and limited budgets. These characteristics create the appeal of composite sampling.

At times the objective of an environmental investigation is to identify all sample units above a criterion level or to estimate the constituent concentration for each sample unit. Traditionally, additional measurements have to be taken to get at the sample unit information. This involves employing some sort of combinatorial scheme after the composite sampling results which are aimed at identifying the anomalous sample unit or units. In some environmental investigations the entire sample unit is used to form the composite so retesting schemes are not a viable option.

A composite sample is a sample that is formed from two or more sample units. The subsequent measurement is taken on the entire composite or an aliquot removed from the composite. The constituent information on the individual sample units is lost; it is observed only indirectly through the composite measurement. An alternative methodology to retesting, that employs the use of combinatorial schemes prior to forming the composites and inverse theory techniques for estimation at the level of the sample unit is described here. For additional reading on composite sampling the reader is referred to an annotated bibliography by Boswell *et al.* (1992a) and two reviews of composite sampling methods, Lancaster and Keller-McNulty (1996) and Lovison *et al.*

(1993).

Inverse theory addresses problems in making inferences about a phenomenon from partial or incomplete information. These problems are often categorized as well- or ill-posed. The notion of a well-posed problem was first introduced by Hadamard in his book *Lectures on Cauchy's Problem in Linear Partial Differential Equations* (1923). It represented a significant step forward in the classification of the multitude of problems associated with differential equations. He singled out those problems with the sufficiently general properties of existence and uniqueness. That is, a solution exists for arbitrary data, there is only one solution and it depends continuously on the data. When one or more of these criteria are violated the problem is categorized as ill-posed. Although not traditionally labeled as such, well- and ill-posed problems arise in statistics when solving systems of linear equations of the type used in analysis of variance and regression. They can be represented notationally as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ , where  $\mathbf{e}$  is the random error associated with  $\mathbf{y}$ . When solving these problems the ability to get an inverse on the design matrix  $\mathbf{X}$  or on  $\mathbf{X}'\mathbf{X}$  will determine if the solution is unique. When an inverse does not exist there are many solutions to the problem. There are more parameters to be estimated than can be uniquely determined from the available data; hence, the problem is ill-posed. The traditional least squares techniques for coping with this lack of information are various reparameterization techniques such as the use of generalized or conditional inverse theory.

Various methods of dealing with ill-posed inverse problems have evolved outside the discipline of statistics in astronomy, geophysics, and mathematics (Craig and Brown, 1986; Tarantola, 1994; and Delves and Mohamed, 1985). These methodologies are referred to in the literature as inverse theory techniques. Craig and Brown (1986) classify these methodologies as classical and non-classical inversion techniques. The difference between the two lies in the philosophical approach to the problem. The term classical is used to denote any technique not explicitly utilizing prior information or assumptions about the solution, such as least squares techniques. Non-classical techniques augment the given information with some prior knowledge of the nature of the solution. Unlike the techniques used to find solutions to well-posed problems the usefulness of a particular technique for an ill-posed problem depends on the nature of the data being inverted (Loredo and Epstein, 1989). No technique is appropriate for all problems.

Section 2.0 describes the compositing process and introduces the notation. A brief history of the development of inverse theory techniques is given in Section 3.0. Section 4.0 describes an application from an environmental investigation and illustrates the use of composite sampling and linear regularization for sample unit estimation. Finally, the conclusions and areas for future research are discussed in Section 5.0.

## 2. Composite Sampling Process

It is useful to view composite sampling as a process which consists of the *statement of the objectives*, the *sampling design*, the *compositing design*, the *measurement process*, and the *data analysis process*. Although, it is the data analysis process that is the focus of this paper, the steps preceding it are described for completeness.

The *statistical objectives* of an investigation involving composite sampling form a dichotomy. For some objectives, the information from the composite samples provides a unique and consistent estimator. An example would be a global estimate of a single parameter such as the

prevalence or the mean concentration of a constituent. Other objectives involve estimation or classification at the level of the sample unit such as classifying each individual sample as to the presence or absence of a constituent. Typically, with these later objectives the number of samples is far greater than the number of observations taken on a few composite samples. If *all* the composites indicated that the trait was absent, all the individual samples could be classified. However, if one or more of the composites indicate the presence of the trait, the samples of the positive composites cannot be classified without additional information.

The *sampling design* is the method by which the  $n$  samples are selected from the target site. Attached to each sample is the continuous random variable  $x_j$  with finite first and second moments, where  $j = 1, 2, \dots, n$ . This random variable represents the value of the constituent of interest to the investigator. In environmental site characterization the constituent is most often the concentration of some organic or inorganic substance or radionuclide believed to pose a risk to the environment.

The *compositing design* is a combinatorial scheme for physically manipulating the sampling units. It specifies which sample units go into which composite, the number of composites, the number of sample units per composite, and the amount of material from each sample unit to use in forming the composite. It takes the selected sample units to the measurement stage by subsampling, mixing, pooling, and reducing. These manipulations are done with the goal of decreasing the number of samples to a few *representative* units on which subsequent measurements are taken.

The *measurement process* is often what motivates the use of a composite sampling method. In many investigations the cost of analyzing a sample is far greater than the costs of sampling, thus generating the appeal of composite sampling. These measurement costs are a function of the trait or traits being measured, the material of the sampling units (i.e. water, rock, soil, blood), the analytic method used to make the measurement, and the desired precision and accuracy of the measurement.

The *data analysis process* is dependent on the objectives of the investigation, the sampling design, the compositing design, and the measurement process. How the information from the composite measurements is used to accomplish the objective of the investigation, will depend on the sources of variability introduced in the processes of sampling, compositing, and measurement and the assumptions that are made along the way. The challenge lies in the fact that the trait of interest is only observed indirectly through measurements on the composites. This creates the need for an estimation methodology such as inverse theory that can incorporate auxiliary information and provide unique solutions.

Notation is needed to represent the result of the compositing process where the  $x_j$ 's are not measured directly but are grouped to form a composite and the measurement taken on the composite. Let  $\mathbf{C}$  represent a  $(cn \times n)$  composite design matrix, where each row  $\mathbf{c}_i$ , represents the formation of a composite for  $i = 1, 2, \dots, c$ , the number of composites. The elements in the rows,  $c_{ij}$ , are the weights given to the sample units forming the composite. The weights are based on the assumptions one makes about the relationship between the physical process of compositing and the composite measurement. In the environmental applications studied by Boswell and Patil (1987), where the sample unit material is soil and water, they argue that the physical manipulation of the sample units to form a composite is equivalent to the numerical manipulation of the data to calculate a mean. In this case, the weights are one over the number of sample units that are combined into the composite.

Notationally, a composite sample will be represented by  $\mathbf{c}_i \mathbf{x}$ , where  $\mathbf{x}$  is a  $(n \times 1)$  vector of unknown sample unit measurements. The problem can then be represented as  $\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{e}$ , where  $\mathbf{y}$

is the  $(cxI)$  data vector of the composite measurements and  $\mathbf{e}$  is the  $(cxI)$  vector of measurement error where  $\mathbf{e} \sim (\mathbf{0}, \Sigma_e)$ . Since  $c$  is usually less than  $n$ ,  $\mathbf{C}'\mathbf{C}$  is rank deficient and the problem is ill-posed.

### 3. Inverse Theory: Linear Regularization

Inverse theory techniques have their origin in the physical sciences. The characteristics common to physical science investigations are variables of interest that can only be measured indirectly, data that are a consequence of observation rather than a designed experiment, existing laws of physics that describe the functional relationship between the variable of interest and the indirect measurements, the use of linear integral equations to describe the relationship between the dependent and independent variables (i.e. Fredholm and Volterra equations), and a deterministic tradition. Vast numbers of physical systems can be described using the Fredholm integral equation of the first kind (Craig and Brown, 1986; Tikhonov and Arsenin, 1977; Twomey, 1977). This integral equation is extremely ill-conditioned. The methodologies that have been developed to bring about a well-posed extension for this equation form the foundation of inverse theory techniques. Fredholm equations involve definite integrals with fixed upper and lower limits. An inhomogeneous Fredholm equation of the first kind has the form

$$g(x) = \int_a^b k(x,y)f(y)dy,$$

where  $f(y)$  is the unknown function to be determined,  $g(x)$  is the known data function, and  $k(x,y)$  is the known kernel of the equation. Fredholm discovered that these integral equations behave very much like a system of linear equations. In the numerical solution of these equations they are first discretized to a system of linear equations and then augmented with prior information on the nature of the solution to create a system of equations that can be solved directly.

The central idea in inverse theory is the prescription

$$\text{minimize}(\mathbf{x}): A[\mathbf{x}] + \lambda B[\mathbf{x}] \quad (1)$$

where  $A[\mathbf{x}]$  and  $B[\mathbf{x}]$  are two positive functionals and  $0 < \lambda < \infty$ . In the integral equation application,  $A[\mathbf{x}]$  would represent the discretized linear system of equations and  $B[\mathbf{x}]$  the prior information. The prescription in (1) follows from the quadratic minimization principle; when a quadratic form is degenerate, adding a constant times a nondegenerate form and minimizing the weighted sum, will lead to a unique and stable solution. The first functional,  $A[\mathbf{x}]$ , is commonly some measure of fit that is degenerate, and the second functional,  $B[\mathbf{x}]$ , is constructed to introduce a measure of smoothness into the solution. If  $A[\mathbf{x}]$  by itself is minimized, the fit becomes very good but the solution is unstable. Small perturbations in the data result in wild oscillations in the solution vector  $\mathbf{x}$ . On the other hand, minimizing  $B[\mathbf{x}]$  by itself gives a smooth solution that has nothing to do with the data. Press *et al.* (1994) make an interesting observation regarding (1). The minimization of the weighted sum,  $A[\mathbf{x}] + \lambda B[\mathbf{x}]$ , can also be interpreted as the minimization of  $A[\mathbf{x}]$  for some constrained value of  $B[\mathbf{x}]$  or the minimization of  $B[\mathbf{x}]$  for some constrained value of  $A[\mathbf{x}]$ ; all three formulations lead

to the same solution vector  $\mathbf{x}$ .

The general principle for dealing with the numerical instability of  $A[\mathbf{x}]$  is referred to as *regularization*, or the term more commonly seen in the statistical literature, *smoothing*. For the inverse technique investigated in this paper,  $A[\mathbf{x}]$  will have the form,

$$A[\mathbf{x}] = \|\mathbf{C}\mathbf{x} - \mathbf{y}\|^2, \quad (2)$$

where  $\mathbf{C}$  is the composite design matrix,  $\mathbf{y}$  the vector of known composite measurements, and  $\mathbf{x}$  the vector of unknown measurements on the sample units. There are two choices to make regarding smoothing, the nature of  $B[\mathbf{x}]$  and the value of  $\lambda$ . The choice of  $B[\mathbf{x}]$  is qualitative, determining the manner of smoothing in the solution vector. Prior or auxiliary information are used to construct a  $B[\mathbf{x}]$  that describes the relationship between members of the solution vector.  $B[\mathbf{x}]$  is referred to as the stabilizing functional or regularizing operator. In contrast, the choice of  $\lambda$  is quantitative and reflects how much to smooth. As  $\lambda$  moves toward 0, the solution tends to the unstable least squares estimates; how far  $\lambda$  moves away from 0 is based upon ones degree of belief in the prior information. The constant  $\lambda$ , adjudicates a delicate balance between remaining faithful to the data and achieving a smooth solution. It is referred to as the smoothing or regularization parameter.

Definitions for the prescription (1) differ as to the choice of  $A[\mathbf{x}]$  and  $B[\mathbf{x}]$ , the computational burdens of these choices, the recommendations for selecting  $\lambda$ , how variability is introduced into the problem formulation, and the philosophical motivations. An intuitively appealing formulation of (1) has been described in the literature by Phillips (1962), Twomey (1977), Tikhonov and Arsenin (1977), and Craig and Brown (1986), under the names of the Phillips-Twomey method, constrained linear regularization, Tikhonov-Miller regularization, and method of regularization. For simplicity it is referred to here as linear regularization. Its appeal is created by the fact it has a closed form solution. An excellent review of linear regularization methods from a numerical analysis perspective is given in the book by Press *et al.* (1994).

In linear regularization techniques most measures of non-smoothness are simple quadratic forms of  $\mathbf{x}$ . For example, a common form for  $B[\mathbf{x}]$  is  $\mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x}$ , where  $\mathbf{B}'\mathbf{B}$  is a positive definite matrix. The resulting minimization principle in (1) is

$$\text{minimize}(\mathbf{x}): \|\mathbf{C}\mathbf{x} - \mathbf{y}\|^2 + \lambda\mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x}, \quad (3)$$

which can easily be reduced to a linear set of normal equations.

The construction of  $\mathbf{B}$  relies on prior information or knowledge about the solution vector. If no prior information is available, the identity matrix can be used. In this case, the minimization of (3) would set about selecting a solution vector with the smallest length  $\|\mathbf{x}\|^2$ . This solution is often referred to as the solution of the inverse problem with zeroth-order regularization. It has a direct link to the numerical analysis technique of singular value decomposition (SVD). In the face of singularities, the SVD algorithm selects from the set of possible solution vectors the one that has minimum length.

If prior knowledge leads one to assume a one-dimensional ordering within  $\mathbf{x}$ , possible measures of smoothness include:

$$\begin{aligned}
 \text{i.} \quad & \mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x} = \sum (x_i - x_{i+1})^2, \\
 \text{ii.} \quad & \mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x} = \sum (-x_i + 2x_{i+1} - x_{i+2})^2, \text{ or} \\
 \text{iii.} \quad & \mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x} = \sum (-x_i + 3x_{i+1} - 3x_{i+2} + x_{i+3})^2.
 \end{aligned} \tag{4}$$

Functional (i) corresponds to a priori belief that a credible solution is not too different from a constant. Functionals (ii) and (iii) correspond to one-dimensional linear and quadratic trends respectively. All of these can easily be extended to the two-dimensional case.

Environmental investigations into site characterization often include historical information about the site; at a minimum they provide auxiliary information on the sample units in the form of spatial coordinates. Historical information can include information on the geology and topology of the site, the source and extent of the contamination, and the constituents that are likely to be present. This information can include less accurate field measurements as well as more accurate measurements from a prior investigation. This information can be used to aid in the construction of  $\mathbf{B}$ .

#### 4. Compositing Sampling Example From an Environmental Investigation

The data are from a Resource Conservation and Recovery Act Facility Investigation at Los Alamos National Laboratory (LANL) in Los Alamos, New Mexico. The investigation targeted a disposal site, approximately 40 feet long and 15 feet wide, located on a moderately steep hillside. The disposal site contained waste originating from a vacuum pump repair shop from 1950 to 1957. The expended vacuum pump oil contained radionuclides and inorganics. One objective of the investigation was to define the extent of the mercury contamination and to remediate those areas with mercury concentrations greater than or equal to 24 mg/kg. At the time of this investigation the environmental scientists at LANL were interested in the feasibility of composite sampling as a methodology for site characterization. Therefore, sampling and compositing designs were implemented to study the application of composite sampling on mercury contaminated soil.

In 1993, a 10x5 sampling grid was constructed over the disposal site and soil samples were taken from 42 of the 50 grid squares at a depth of 0 to 10 centimeters. Within a grid square the sample was taken at the centrally aligned location. All 42 sample units were analyzed for mercury using cold vapor atomic absorption spectrometry. In addition, fifteen composite samples were constructed from the sample units within the ten rows and five columns. All fifteen composite samples were analyzed for mercury with the same analytic technique. What is unusual about these data is the existence of measurements at both the composite and sample unit levels. This allows for a rare opportunity to evaluate the assumptions one makes about the physical process of compositing and to compare estimates based on composite measurements to those same estimates based on sample unit measurements.

Figure 1 illustrates the systematic sampling design and composite design for this investigation. The 42 sampling locations are identified by the squares, with the solid squares indicating which sample units were in a particular row or column composite. Figure 1(a) illustrates the column composite design. Five column composites (referred to as Y11 - Y15) were constructed, with the number of sample units making up the column composites ranging from six to ten. The measurements of the column composites are plotted in the panel above the column composite design. Figure 1(b) illustrates the row composite design. Ten row composites (Y1 - Y10) were constructed,

with the number of sample units making up the row composites ranging from two to five. The ten row composite measurements are plotted in the panel to the right the row design. Loess lines are superimposed on the composite measurement plots to aid in interpretation.

Figure 2 displays the 42 sample unit measurements on the original and log ten scale. The sample units measurements range from 0.867 mg/kg to 14566.667 mg/kg. The histogram and boxplot of the original data displayed in Figure 2(a) identify a positive skewness in the mercury data, a feature which is often seen in environmental data. The histogram and boxplot for the transformed data are displayed in Figure 2(b). Even with this transformation the data still show some positive skewness. The transformed data are also displayed in an image plot in Figure 2(c). The darker the shade in the image plot the higher the mercury concentration; white squares indicate no information was collected from that grid square. The image plot shows the area of highest concentration is located along the middle column and there is a decreasing trend. Of the 42 grid squares, (i.e., remediation units), 19 have mercury concentrations at or above the remediation level of 24 mg/kg, and these remediation units are located in the three middle columns. The log ten transformation is used in constructing image plots in order to enhance visual interpretation. All analyses are done on the untransformed data. The log ten was not used in the analyses since the additive relationship between the composite measurement and the individual sample unit measurements cannot be expressed using logarithms.

The assumption made about the compositing process is explored in Figure 3. Figure 3 plots the mean of the sample units which were composited versus the composite measurement with a 45° line for reference. Also included are plots of the differences between the composite measurement and the mean of the sample units making up the composite. Figure 3(a) contains information on all fifteen composites and draws attention to Y13 (middle column composite) which has a concentration thirteen times larger than any of the other composite measurements. Since Y13 conceals the relationship of the remaining composites, it was removed and the data replotted in Figure 3(b). One can now see the anomalous results for Y8 which has a composite measurement of 164.00 mg/kg and a sample unit mean of 55.78 mg/kg. Figure 3(b) shows a slight trend in the relationship between the composite measurement and mean of the sample unit measurements. The composite measurement is under estimating at lower concentrations and over estimating at the higher concentrations. With the exception of composite Y8, the correspondence between the sample unit and the composite measurements was better than expected by the environmental scientists.

Linear regularization is used on the mercury data to illustrate how one can adapt this theory to get estimates at the sample unit level. Recall for this example,  $\mathbf{x}$  is the vector of unknown mercury concentrations for the 42 sample units. A common form in linear regularization for  $B[\mathbf{x}]$ , and the one used for this illustration is  $\mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x}$ .  $\mathbf{B}$  is used to describe the relationship between sample units on the 10x5 sampling grid. In this investigation information on the location of the contaminant source and the topology of the site suggests a possible linear relationship within the columns of the sampling grid. Therefore, the functional 4(ii) was adapted for the two-dimensional case and also to take into account the missing sample units and the edge effects. This lead to a  $\mathbf{B}$  matrix (42 x 42) for the mercury example with the structure,



$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -3/2 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -3/2 & 0 & 0 & 0 & 1 & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -3/2 & 0 & 0 & 0 & 1 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & -2 & 0 & 0 & 1 & \dots & 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & & & & & & & & & \vdots & & & & & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \dots & 1 & 0 & 0 & 0 & -2/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \dots & 0 & 1 & 0 & 0 & 0 & -2/3 \end{bmatrix}.$$

Each row of the **B** defines the nearest neighbor relationship for a particular sample unit. Samples units are numbered 1 to 42 by row (see Figure 1). As an example, the sixth row of **B** defines the nearest neighbor relationship for the sample unit 6. This row of **B** implies that the constituent concentration of sample unit 6 is equal to the sum of the concentrations from sample units 2 and 9.

For illustration purposes, the value of  $\lambda$  is calculated as  $\text{Tr}(\mathbf{C}'\mathbf{C})/\text{Tr}(\mathbf{B}'\mathbf{B})$ . This choice gives comparable weight to both of the functionals in (3).

The results of this straight forward application of linear regularization are displayed using image plots in Figures 4 and 5. Figure 4(b) displays the log ten transformation of the observed sample unit concentrations and 4(a) the log ten transformation of the estimates from the linear regularization. The difference between the observed concentrations and the estimates for the untransformed data are displayed in Figure 4(c). The residuals used in Figure 4(c) have been grouped into four categories based on the quartiles. The darkest shading represents the largest positive residuals (under estimates of the observed concentrations) and the lightest shading represents the largest negative residuals (over estimates of the observed concentrations). Even with this very naive construction of the **B** matrix the linear regularization estimates can identify the area of high mercury contamination.

The linear regularization estimates can also be used to identify the grid squares (i.e., remediation units) at or above the remediation level of 24 mg/kg. This methodology identifies 34 of the sample units correctly with 6 false positives and 2 false negatives. These results are displayed in Figure 5. This information gives an estimate of the prevalence of mercury as 0.50 compared to the observed rate of 0.45.

Clearly there are many sources of variability inherent in this compositing process that have not been isolated. Two of these are within sample and within composite unit variability. An *ideal* investigation would have evaluated these sources of variability as a preliminary first step by removing and analyzing numerous aliquots from the same sample and composite units. This would have provided valuable information on the homogeneity of the material. Other valuable information such as sample support size, how the composites were constructed, and measurement variability are not available for these data. Even with these limitations the data provide valuable information on the use of composite sampling methodologies for site characterization where the sample units are soil and the contaminant is mercury.

## 5. Conclusions

Although the idea of compositing samples has generated much interest (Garner *et al.*, 1988), applications of composite sampling methods in the environmental science literature are nearly nonexistent. The articles that have appeared are retrospective or theoretical. The authors suggest that the reason for the lack of transfer from theory to application is in part due to the numerous physical forms a composite can take, including soil, rock, sludge, and air. This, combined with the numerous constituents and analytic measurement techniques, prohibits generalization between investigations. More controlled empirical research needs to be performed in the laboratory and in the field before any widespread application of composite sampling methodologies can occur in the environmental sciences.

This brief study into the use of linear regularization on an ill-posed composite sampling problem is encouraging for two reasons. First, the relationship between the composite measurement and the mean of the sampling units is surprisingly good. This, despite the fact a soil sample contaminated with elemental mercury, as in the investigation cited here, probably represents the furthest extreme in sample heterogeneity (Bloom, 1992). Second, despite the “nonideal” conditions of the investigation, which include the missing sample units in row composites Y1-Y5, the absence of a sample units from 8 of the 50 sampling grids, and the anomalous information on row composite Y8, the estimates follow the trend in the observed measurements.

The authors are currently investigating alternative constructions for **B**, the construction of optimal composite designs, and optimal choices for  $\lambda$ . Other inverse theory methodologies are also being explored.

## Acknowledgments

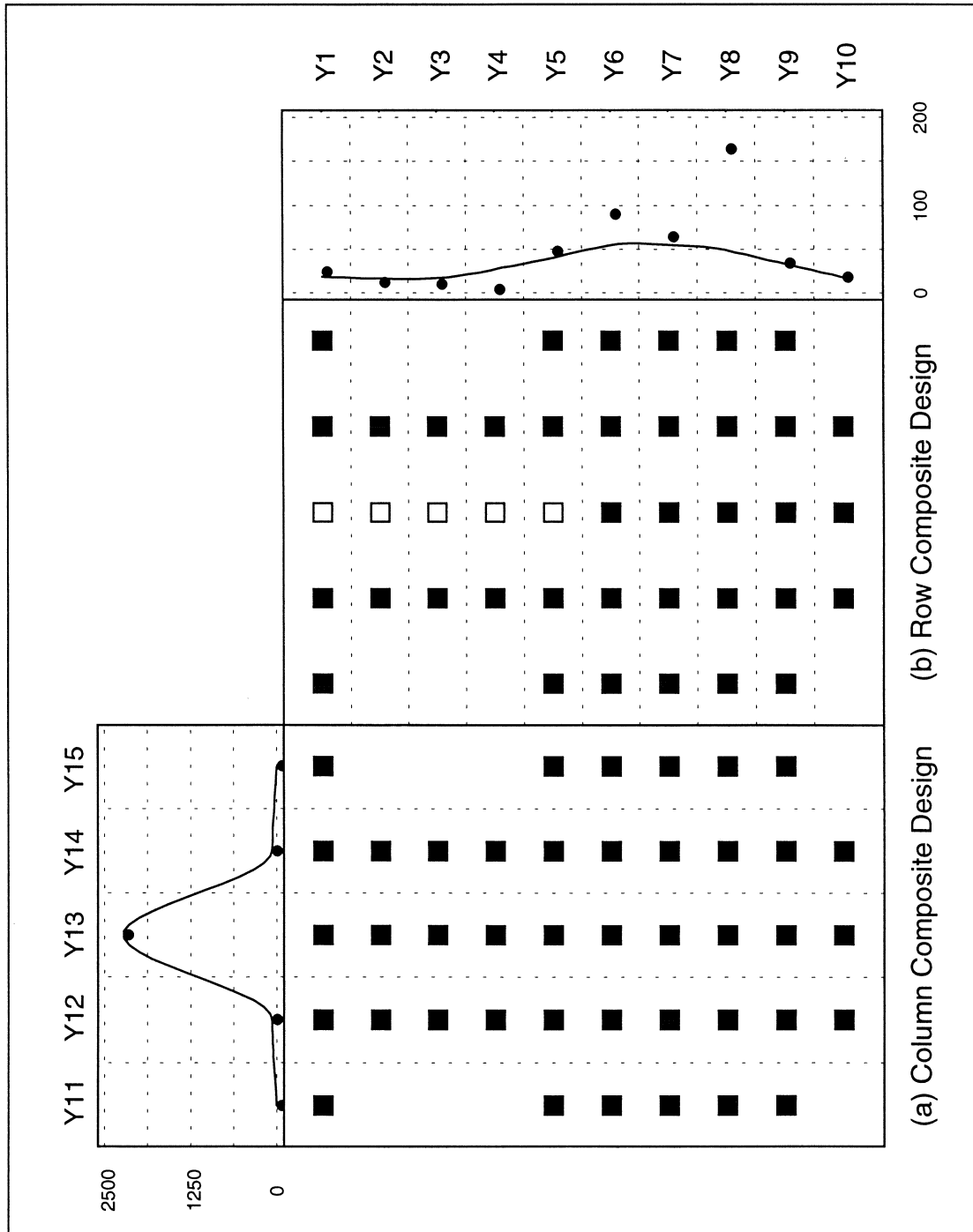
The authors thank the environmental scientists at Los Alamos National Laboratory along with Randy Ryti and Sylvia Hammerdinger at Neptune & Co.

## References

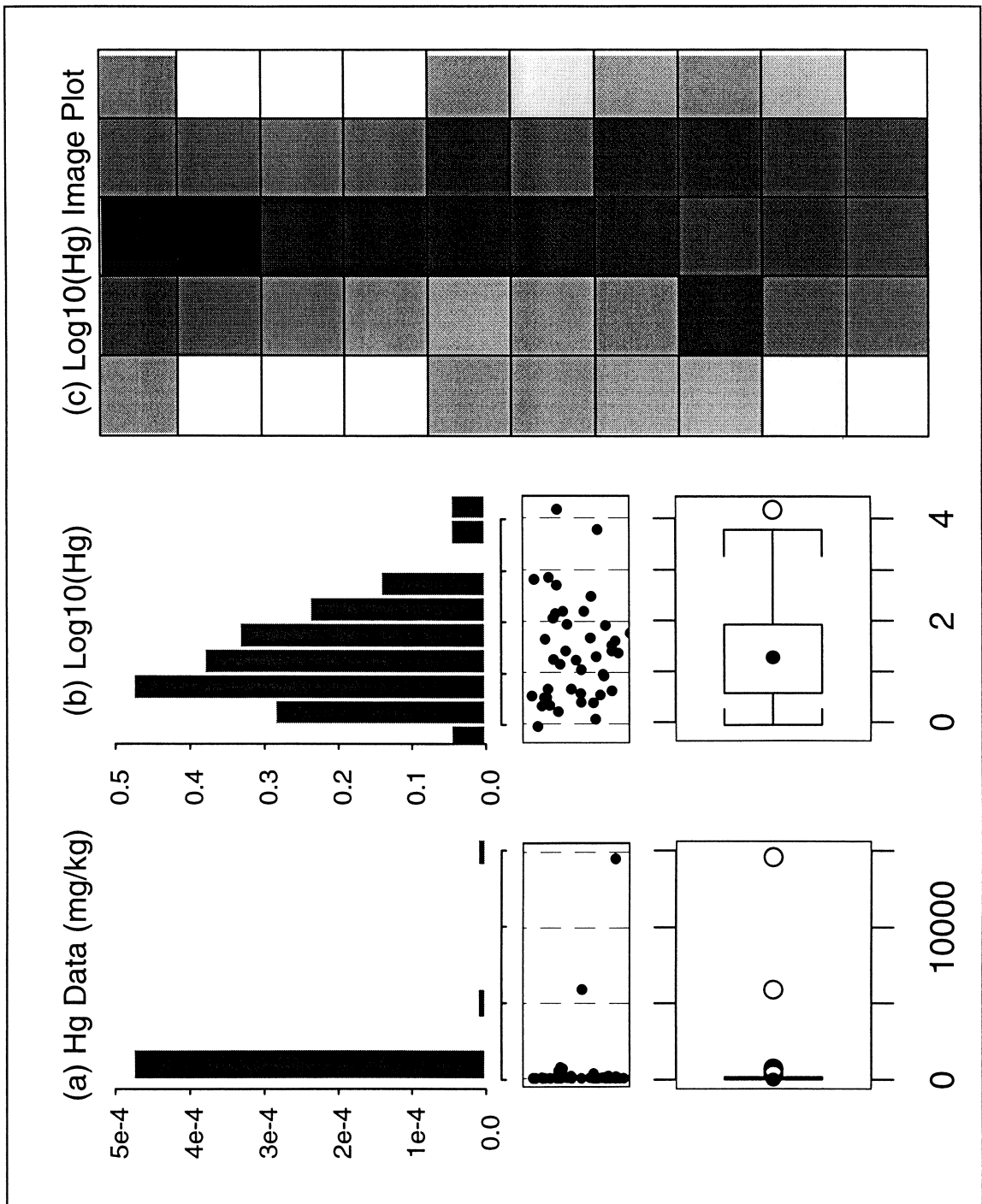
- Bloom, N. (1992), “Considerations in Sampling for and Analysis of Mercury at Uncharacterized Spill Sites”, Workshop on Mercury Contamination at Natural Gas Industry Sites, Chicago, IL.
- Boswell, M. T., Gore, S. D., Lovison, G. and Patil, G. P. (1992a), “Annotated Bibliography of Composite Sampling”, Technical Report No. 92-0802, Center for Statistical Ecology and Environmental Statistics, Dept. of Statistics, Penn. St. Univ., Univ. Park, PA 16802.
- Boswell, M. T., Gore, S. D. and Patil, G. P. (1992b), “Efficiency of Various Composite Retesting Schemes to Classify Samples with Presence/Absence Measurements”, Technical Report No. 90-0901, Center for Statistical Ecology and Environmental Statistics, Dept. of Statistics, Penn. St. Univ., Univ. Park, PA 16802.
- Boswell, M. T. and Patil, G. P. (1987), “A Perspective on Composite Sampling”, *Communications in Statistics, Theory and Methods*, 16, 3069-3093.

- Craig, I. J. D., and Brown, J. C. (1986), *Inverse Problems in Astronomy: A Guide to Inversion Strategies for Remotely Sensed Data*, England: Adam Hilger Ltd.
- Delves, L.M. and Mohamed, J.L. (1985), *Computational Methods for Integral Equations*, Cambridge, U.K.: Cambridge University Press.
- Garner, F. C., Stapanian, M. A. and Williams, L. R. (1988), "Composite Sampling for Environmental Monitoring", *Principles of Environmental Sampling*, American Chemical Society, Washington, DC, 363-374.
- Hadamard, J. (1932), *Le Problem de Cauchy et les équations aux dérivées partielles linéaires hypéboliques*, Hermann, Paris.
- Lancaster, V. A. and Keller-McNulty, S. (1996), "A Review of Composite Sampling Methods", Technical Report No. 96-01. Dept. of Statistics, Kansas St. Univ., Manhattan, KS 66506.
- Loredo, T. J. And Epstein, R. I. (1989), "Analyzing Gama-Ray Burst Spectral Data," *The Astrophysical Journal*, 336, 896-919.
- Lovison, G., Gore, S. D. and Patil, G. P. (1993), "Design and Analysis of Composite Sampling Procedures", Technical Report No. 92-1007, Center for Statistical Ecology and Environmental Statistics, Dept. of Statistics, Penn. St. Univ., Univ. Park, PA 16802.
- Phillips, D.L. (1962), "A technique for the numerical solution of certain integral equations of the first kind", *Journal of the Association for Computing Machinery*, 9, 84-97.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (1994), *Numerical Recipes in C: The Art of Scientific Computing*, Cambridge: Cambridge University Press, pp. 788-888.
- Tarantola, A. (1994), *Inverse Problem Theory: Methods for Data Fitting and Model Parameter Estimation*, The Netherlands: Elsevier Science B. V.
- Tikhonov, A. N., and Arsenin, V. Y. (1977), *Solutions of Ill-Posed Problems*, Washington, D. C.: V. H. Winston & Sons.
- Twomey, S. (1977), *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurements*, The Netherlands: Elsevier Scientific Publishing Company.

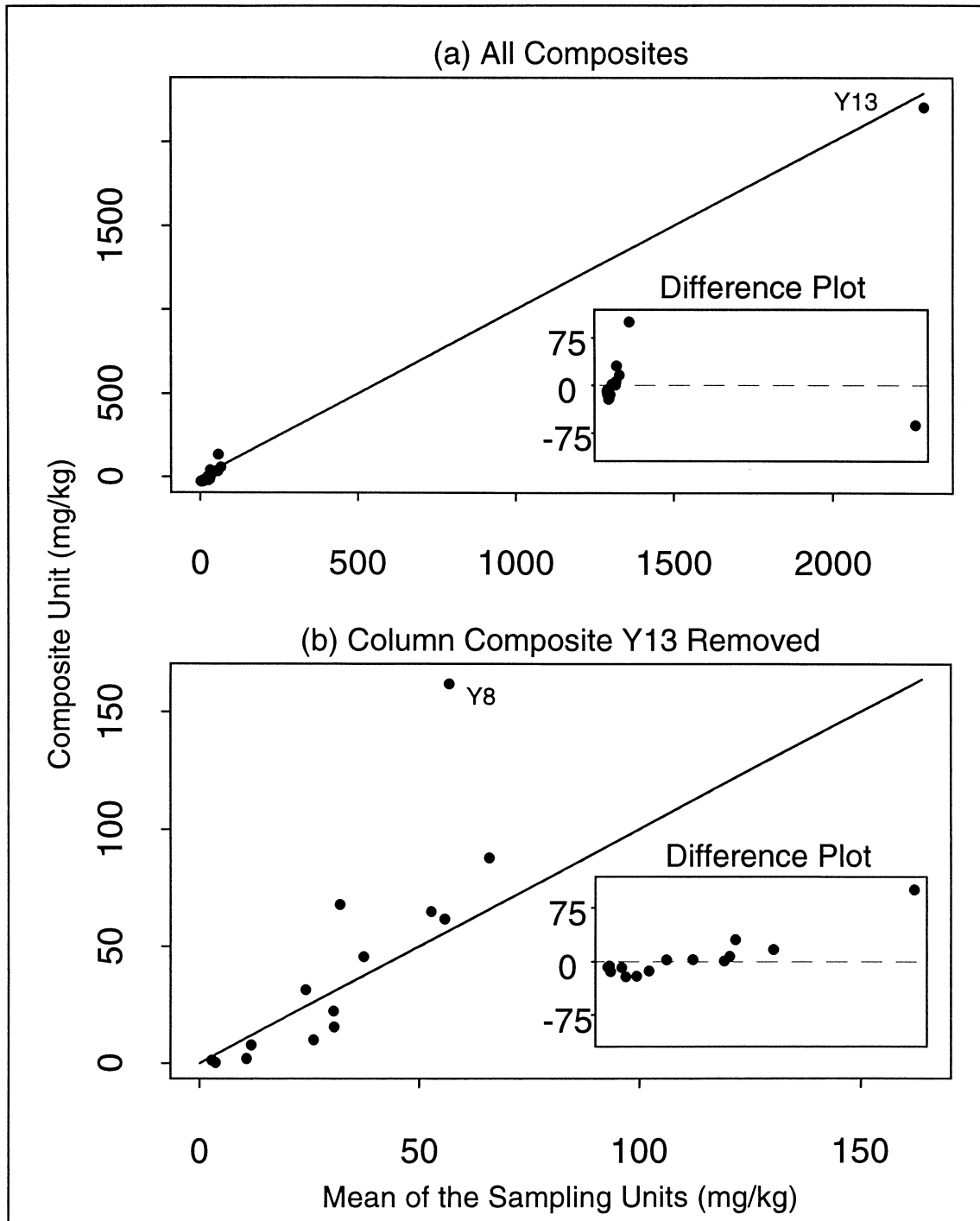
**Figure 1.** Sampling and composite designs for the LANL mercury investigation. The 42 sampling locations are identified by squares. (a) and (b) illustrate the column and row composite designs, where a solid square indicates a sampling unit from that grid point was included in a particular row or column composite. The panels on the top and the far right plot the column and row composite measurements for mercury in mg/kg, with loess lines superimposed to aid in interpretation.



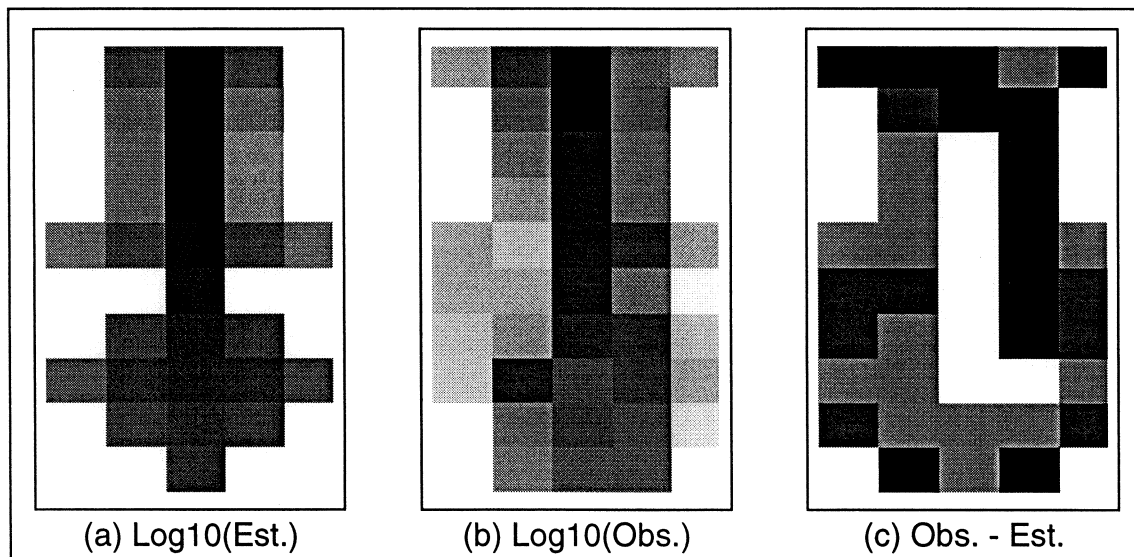
**Figure 2.** Descriptive plots for the mercury (mg/kg) measurements of the 42 sampling units from the LANL investigation. (a) is a histogram, one dimensional scatter plot, and boxplot for the data in its original scale. (b) is a histogram, one dimensional scatter plot, and boxplot for the data in log ten scale. (c) is an image plot of the log ten transformed data, the darkest shade represents 4.16, the lightest shade represents -0.06, and white areas indicate no sample unit was collected from that grid square.



**Figure 3.** Mean of the sampling units vs. the composite measurement. The solid line in the large panels are 45° reference lines. The small panels are plots of the differences between the composite measurement and the mean of the sampling units making of the composite with a dashed reference line at zero. (a) contains the information for all fifteen composites. (b) is a plot without the data from composite Y13.



**Figure 4.** Image plots for mercury concentration. (a) displays the linear regularization estimates for the 42 remediation units on the log ten scale (darkest shade represents 3.39 and lightest shade represents -0.01). (b) displays the observed mercury concentrations for the 42 remediation units on the log ten scale (darkest shade represents 4.14 and lightest shade represents -0.06). (c) displays the difference between the observed and estimated values for the untransformed data (mg/kg) grouped into four categories based on quartiles. The darkest shading represents the largest positive residuals (under estimates) and the lightest shading represents the largest negative residuals (over estimates). White areas indicate no sample unit was collected from that grid square.



**Figure 5.** Image plots for the classification of remediation units. (a) displays classification based on the linear regularization estimates for the 42 remediation units. (b) displays classification based on the 42 observed mercury concentrations. Black areas represent mercury concentrations  $\geq 24$  mg/kg, grey areas represent mercury concentrations  $< 24$  mg/kg, and white areas indicate no sample unit was collected from that grid square.

