

# **Computing Ordinary Least-Squares Parameter Estimates for the National Descriptive Model of Mercury in Fish**

Chapter 10 of  
Section C, Computer Programs  
**Book 7, Annotated Data Processing and Computations**

Techniques and Methods 7–C10



# **Computing Ordinary Least-Squares Parameter Estimates for the National Descriptive Model of Mercury in Fish**

By David I. Donato

Chapter 10 of  
Section C, Computer Programs  
**Book 7, Annotated Data Processing and Computations**

Techniques and Methods 7–C10

**U.S. Department of the Interior  
U.S. Geological Survey**

**U.S. Department of the Interior**  
KEN SALAZAR, Secretary

**U.S. Geological Survey**  
Marcia K. McNutt, Director

U.S. Geological Survey, Reston, Virginia: 2013

For more information on the USGS—the Federal source for science about the Earth, its natural and living resources, natural hazards, and the environment, visit <http://www.usgs.gov> or call 1–888–ASK–USGS.

For an overview of USGS information products, including maps, imagery, and publications, visit <http://www.usgs.gov/pubprod>

To order this and other USGS information products, visit <http://store.usgs.gov>

Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

Although this information product for the most part is in the public domain, it also may contain copyrighted materials as noted in the text. Permission to reproduce copyrighted items must be secured from the copyright owner.

Suggested citation:

Donato, D.I., 2013, Computing ordinary least-squares parameter estimates for the National Descriptive Model of Mercury in Fish: U.S. Geological Survey Techniques and Methods, book 7, chap. C10, 9 p., <http://pubs.usgs.gov/tm/07/c10>.

# Contents

Abstract.....	1
Introduction .....	1
The National Descriptive Model of Mercury in Fish as a Formal Statistical Model.....	2
The Normal Equations for Ordinary Least-Squares Estimation .....	2
Matrix Form of Data for the National Descriptive Model of Mercury in Fish .....	3
Computing the Normal-Equation Matrix Products During Data Input.....	4
Computing With Weighted Observations.....	5
Computational Methods for Solving the Normal Equations .....	6
Reference Software .....	6
Summary.....	7
References Cited.....	8
Appendix.....	9

This page intentionally left blank.

# Computing Ordinary Least-Squares Parameter Estimates for the National Descriptive Model of Mercury in Fish

By David I. Donato

## Abstract

A specialized technique is used to compute weighted ordinary least-squares (OLS) estimates of the parameters of the National Descriptive Model of Mercury in Fish (NDMMF) in less time using less computer memory than general methods. The characteristics of the NDMMF allow the two products  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$  in the normal equations to be filled out in a second or two of computer time during a single pass through the  $N$  data observations. As a result, the matrix  $\mathbf{X}$  does not have to be stored in computer memory and the computationally expensive matrix multiplications generally required to produce  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$  do not have to be carried out. The normal equations may then be solved to determine the best-fit parameters in the OLS sense. The computational solution based on this specialized technique requires  $O(8p^2+16p)$  bytes of computer memory for  $p$  parameters on a machine with 8-byte double-precision numbers. This publication includes a reference implementation of this technique and a Gaussian-elimination solver in preliminary custom software.

## Introduction

The National Descriptive Model of Mercury in Fish (NDMMF) is a statistical model used to predict the concentration of methylmercury in fish tissue (Wente, 2004). This model is of interest in current research at the U.S. Geological Survey (USGS) because of its ability to explain much of the variation in fish-tissue methylmercury concentrations as variation by geographic location, variation over time, and variation by fish species and length.

Before the NDMMF can be used to predict fish-tissue methylmercury concentrations, its parameters must be fitted to a collection of observations of fish-tissue methylmercury

concentrations. The statistical procedure used to fit the parameters of the NDMMF to observed data is that of **maximum-likelihood estimation** (MLE). The parameters of the NDMMF are fitted in the maximum-likelihood sense because the available national database of fish-tissue methylmercury concentrations includes a substantial number of **left-censored** observations<sup>1</sup>, and the MLE method makes better use of the information contained in censored observations than other common statistical methods of parameter estimation, including the method of **ordinary least-squares estimation** (Helsel, 2004). Ordinary least-squares (OLS) parameter estimation is, however, a preliminary step of choice in computing maximum-likelihood estimates and, therefore, an essential part of the complete procedure for estimating the parameters of the NDMMF.

This publication presents a technique, along with a reference (baseline) implementation in custom computer software, for computing OLS parameter estimates for the NDMMF based on weighted observations. By exploiting the specific characteristics of the NDMMF, this technique enables faster computation of OLS parameter estimates for the NDMMF using less computer random-access memory (RAM) than is possible with generalized statistical computer software. Faster computation and use of less RAM are major improvements in fitting parameters to the NDMMF because some generalized software procedures may fail to run because they require more RAM than is available, and those that do run will require many hours or days of computation. The reference software computes the best-fit parameters for the NDMMF in the ordinary least-squares sense using Gaussian elimination with back substitution (Noble, 1969; Draper and Smith, 1966); a revised version of the software using an alternative computational method, such as LU or Cholesky decomposition, is planned (Press and others, 1992).

The preliminary custom software included with this publication is not intended for general use; rather, it is intended for research and development use in conjunction with other custom software developed at the U.S. Geological Survey (USGS) for computing maximum-likelihood estimates of the parameters of the NDMMF more quickly, and using less RAM, than is possible with generalized statistical software (Donato, 2012). A full understanding of this report requires knowledge of linear algebra and familiarity with statistical models.

---

<sup>1</sup>A left-censored observation is a value determined to be below a particular detection limit but otherwise unspecified. For example, if a laboratory procedure cannot detect methylmercury concentration values below 0.020 part per million, then a sample with an undetected concentration would be recorded as a nondetected and, thus, as a left-censored value with a detection limit of 0.020 part per million (Helsel, 2004).

## The National Descriptive Model of Mercury in Fish as a Formal Statistical Model

The NDMMF is expressed formally and compactly as:

$$y_{ijk} = \alpha_j L_{ijk} + \beta_k + \varepsilon_{ijk} \quad (1)$$

where

- $y_{ijk}$  =  $\ln(C_{ijk} + 1)$  and  $C_{ijk}$  is the  $i$ th observed concentration value for species/cut  $j$  and sampling event  $k$ ;
- $\alpha_j$  = the parameter for the  $j$ th species/cut;
- $L_{ijk}$  =  $\ln(\text{length}_{ijk} + 1)$  and  $\text{length}_{ijk}$  is the  $i$ th observed fish length in inches for species/cut  $j$  and sampling event  $k$ ;
- $\beta_k$  = the parameter for the  $k$ th sampling event; and
- $\varepsilon_{ijk}$  = the random error for the  $i$ th observation for species/cut  $j$  and sampling event  $k$  and  $\varepsilon_{ijk} \sim N(0, \sigma)$ .

In the seminal paper on the NDMMF (Wente, 2004), the terms **species/cut** and **sampling event** are defined to have specialized meaning as applied to the NDMMF. Briefly, the species/cut is a combination of a species of fish and a method of preparing the tissue of the fish for laboratory analysis; a sampling event is a collection of samples from a particular geographic location within the same year. Wente's 2004 publication provides a full explanation of these terms.

Equation (1) implies that the number of parameters of the NDMMF equals the sum of the number of species/cut combinations and the number of sampling events. This is more readily apparent when the NDMMF is expressed less compactly as:

$$y_{ijk} = \alpha_1 SPC_1 L_{ijk} + \alpha_2 SPC_2 L_{ijk} + \cdots + \alpha_n SPC_n L_{ijk} + \beta_1 Event_1 + \beta_2 Event_2 + \cdots + \beta_m Event_m + \varepsilon_{ijk} \quad (2)$$

where

- $y_{ijk}, L_{ijk},$   
and  $\varepsilon_{ijk}$  are defined as for equation (1);
- $n$  = the number of species/cut combinations and parameters;
- $m$  = the number of sampling events and parameters;

$$\alpha_t SPC_t = \begin{cases} 0 & \text{if } t \neq j \\ \alpha_j & \text{if } t = j \end{cases}, \quad \text{and}$$

$$\beta_t Event_t = \begin{cases} 0 & \text{if } t \neq k \\ \beta_k & \text{if } t = k \end{cases}.$$

In equation (2), the sets of variables  $\{SPC_t\}_{t=1}^n$  and  $\{Event_t\}_{t=1}^m$  are **indicator variables** (Box, 1978). These indicator variables allow the NDMMF to be expressed in the form of the general linear model, a model linear with respect to its parameters (Mendenhall, 1968).

## The Normal Equations for Ordinary Least-Squares Estimation

An advantage of expressing the NDMMF in the form of the general linear model is that an expression for its best-fit parameters in the ordinary least-squares sense is known in a general form. The system of equations that determines the ordinary least-squares (OLS) parameters for a general linear model is called the **normal equations** (Monahan, 2001; Press and others, 1992). In matrix form, the general linear model is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (3)$$

and the normal equations for this model are:

$$(\mathbf{X}'\mathbf{X})\mathbf{b} = \mathbf{X}'\mathbf{y} \quad (4)$$

where specifically for the NDMMF

- $p = n + m$  is the total number of parameters (the number of species/cut parameters plus the number of sampling-event parameters);
- $N$  = the total number of observations (and generally  $N \gg p$ );
- $\mathbf{y}$  = the  $(N \times 1)$  vector of response observations;
- $\mathbf{X}$  = an  $(N \times p)$  matrix of values of observed data;
- $\mathbf{X}'$  = the  $(p \times N)$  matrix transpose of  $\mathbf{X}$ ; and
- $\boldsymbol{\beta}$  = a  $(p \times 1)$  vector of parameters and  $\mathbf{b}$  is the corresponding  $(p \times 1)$  vector of least-squares estimates of the parameters.

The symbol for the vector  $\boldsymbol{\beta}$  in the normal equations should not be confused with the symbol for the sampling-event parameters  $\{\beta_k\}$  in the NDMMF. The vector  $\boldsymbol{\beta}$  is a vector of all model parameters, so for the NDMMF,  $\boldsymbol{\beta}$  includes all  $n$  of the  $\{\alpha_j\}$  and all  $m$  of the  $\{\beta_k\}$  parameters. The matrix  $\mathbf{X}$  of data observations is sparse, containing only two nonzero values per row.

## Matrix Form of Data for the National Descriptive Model of Mercury in Fish

Equation (4) includes all of the matrices involved in the computation of OLS parameter estimates:  $\mathbf{y}$ ,  $\mathbf{b}$ ,  $\mathbf{X}$ , and  $\mathbf{X}'$ . Since the internal structure of these matrices depends on features of parameters and data specific to the NDMMF, each of these matrices must be specialized to the NDMMF before the details of computation of OLS parameter estimates for the NDMMF can be finalized.

$$\mathbf{y} = \begin{bmatrix} y_{1,1,1} \\ y_{2,1,1} \\ \vdots \\ y_{1,52,104} \\ y_{2,52,104} \\ \vdots \\ y_{15,n,m} \end{bmatrix} \quad (5)$$

The response vector  $\mathbf{y}$ , which contains one row for each of the  $N$  observations, may and usually does contain multiple observations (rows) for each combination of species/cut and sampling event. Therefore,  $N$  is typically much larger than the total number of parameters,  $p = n + m$ . Each row in  $\mathbf{y}$  corresponds to a row in  $\mathbf{X}$ .

$$\mathbf{b} = \begin{bmatrix} \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_n \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_m \end{bmatrix} \quad (6)$$

The vector of least-squares parameter estimates  $\mathbf{b}$  must be put into a sequence according to some convention in order to establish a correspondence between each parameter's estimate in  $\mathbf{b}$  and its relevant data observations in the matrices  $\mathbf{X}$  and  $\mathbf{X}'$ . Let the convention be that the rows of  $\mathbf{b}$  and  $\mathbf{X}'$  and the columns of  $\mathbf{X}$  must correspond to the  $n$  species/cut parameters  $\{\alpha_j\}$  in sequence by species/cut parameter number, followed by the  $m$  sampling-event parameters  $\{\beta_k\}$  in sequence by sampling-event parameter number.

$$\mathbf{X} = \begin{bmatrix} 0 & \dots & \dots & \dots & 0 & | & 0 & \dots & \dots & \dots & 0 \\ \vdots & \dots & \vdots & \dots & \vdots & | & \vdots & \dots & \vdots & \dots & \vdots \\ 0 & \dots & SPC_j L_{ijk} & \dots & 0 & | & 0 & \dots & \dots & Event_k & \dots & 0 \\ \vdots & \dots & \vdots & \dots & \vdots & | & \vdots & \dots & \vdots & \dots & \vdots \\ 0 & \dots & \dots & \dots & 0 & | & 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix} \quad (7)$$

Equation (7) illustrates the general appearance of the sparse matrix  $\mathbf{X}$  symbolically and equation (8) shows a numerical example of  $\mathbf{X}$  in outline.

$$\mathbf{X} = \begin{bmatrix} 0 & \dots & 1.4 & \dots & 0 & \dots & 0 & | & 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \dots & \vdots & \dots & \vdots & \dots & \vdots & | & \vdots & \dots & \vdots & \dots & \vdots & \dots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1.2 & | & 0 & \dots & 0 & \dots & 1 & \dots & 0 \\ \vdots & \dots & \vdots & \dots & \vdots & \dots & \vdots & | & \vdots & \dots & \vdots & \dots & \vdots & \dots & \vdots \\ 0 & \dots & 0 & \dots & 2.7 & \dots & 0 & | & 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \begin{matrix} 1 \\ \downarrow \\ N \end{matrix} \quad (8)$$

$\alpha_1 \qquad \rightarrow \qquad \alpha_n \quad | \quad \beta_1 \qquad \rightarrow \qquad \beta_m$

The matrix  $\mathbf{X}$  of observations of the independent variables contains one row for each of the  $N$  observations. Each row of  $\mathbf{X}$  is an ordered list of the coefficients of the parameters  $\{\alpha_t\}_{t=1}^n$  and then  $\{\beta_t\}_{t=1}^m$  for one equation in the system of  $N$  equations expressed equivalently, but in different forms, by equations (2) and (3). The correspondence between rows and the  $N$  observations, and between columns and parameters, is illustrated in equation (8); and the transposed correspondences are illustrated in equation (9).

Each row of  $\mathbf{X}$  provides an observation for exactly one species/cut  $j$  and exactly one sampling event  $k$ , **so each row contains exactly two nonzero values** determined by the set of observations: **first**, the value of  $L_{ijk}$  for the  $i^{\text{th}}$  observation for species/cut  $j$  and sampling event  $k$ , **then** a value of 1 for the single indicator variable for sampling event  $k$ . All other values in each row are zero (0). Thus in each row, all of the indicator variables  $\{SPC_t\}_{t=1}^n$  and  $\{Event_t\}_{t=1}^m$  are equal to zero, except for the one  $SPC_j$  indicator for the species/cut  $j$  associated with the observation for the row, and except for the one  $Event_k$  indicator for the sampling event  $k$  associated with the observation for the row.

$$\mathbf{X}' = \begin{array}{ccc|ccc}
 & 1 & \rightarrow & N & & \\
 \left[ \begin{array}{ccc}
 0 & \cdots & 0 \\
 \vdots & \vdots & \vdots \\
 1.4 & \cdots & 0 \\
 \vdots & \vdots & \vdots \\
 0 & \cdots & 2.7 \\
 \vdots & \vdots & \vdots \\
 0 & \cdots & 0 \\
 \hline
 1 & \cdots & 0 \\
 \vdots & \vdots & \vdots \\
 0 & \cdots & 0 \\
 \vdots & \vdots & \vdots \\
 0 & \cdots & 0 \\
 \vdots & \vdots & \vdots \\
 0 & \cdots & 1
 \end{array} \right] & \begin{array}{l} \alpha_1 \\ \\ \downarrow \\ \\ \alpha_n \\ \hline \beta_1 \\ \\ \\ \beta_m \end{array} & (9) & 
 \end{array}$$

The matrix transpose  $\mathbf{X}'$  contains one column for each observation, as shown above in equation (9), with the same illustrative values shown for  $\mathbf{X}$  in equation (8). Here, each **column** contains all zero values except for the two values determined by the observation it represents. (Although the mathematical form of the NDMMF permits  $L_{ijk}$  to equal zero, in practice  $L_{ijk}$  is always strictly greater than zero.) The first  $n$  rows of each column in the transpose correspond to the species/cut parameters  $\{\alpha_t\}_{t=1}^n$  in sequence, and the following  $m$  rows correspond to the sampling-event parameters  $\{\beta_t\}_{t=1}^m$  in sequence.

## Computing the Normal-Equation Matrix Products During Data Input

The normal equations in matrix form,  $(\mathbf{X}'\mathbf{X})\mathbf{b} = \mathbf{X}'\mathbf{y}$ , contain two matrix products:  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$ . Because of the sparseness and particular configuration of the matrix  $\mathbf{X}$  and its transpose  $\mathbf{X}'$  for the NDMMF, the two matrix products can be created directly by software during the input of the  $N$  data records, thus making it unnecessary to store  $\mathbf{X}$  or  $\mathbf{X}'$  in computer memory and unnecessary to perform any computationally expensive matrix multiplications. **This is the essence of the technique described in this publication: that a computation that might take hours or days with generalized computational methods is avoided by using an alternative, fast input process that requires only a few seconds of computer processing time.**

**The benefits of the alternative computation of  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$  during data input are substantial.** Generalized computational methods must store the  $(N \times p)$  matrix  $\mathbf{X}$  in computer memory, and because  $N \gg p$  this requires an amount of additional memory that is several times the amount

of memory required for storing just the  $(p \times p)$  product matrix  $\mathbf{X}'\mathbf{X}$ . Although some generalized software performs the multiplication of the sparse matrices  $\mathbf{X}'$  and  $\mathbf{X}$  with fewer computations or with less expensive computational operations than a full multiplication, the multiplication of  $\mathbf{X}'\mathbf{X}$  still requires a substantial proportion of the  $O(Np^2)$  floating-point multiplications and additions required in the general case. Thus generalized computational methods may require more time just to set up the normal equations in computer memory (by computing  $\mathbf{X}'\mathbf{X}$ ) than is required to solve them.

To understand how the product matrix  $\mathbf{X}'\mathbf{X}$  is created during data input, begin by applying the definition of matrix multiplication. Each element of the  $(p \times p)$  product matrix  $\mathbf{X}'\mathbf{X}$  is the sum of the  $N$  products of each element of a row multiplied by the corresponding element of a column of  $\mathbf{X}$ . Now, rather than considering what would be involved in multiplying these two matrices in full, consider instead how the values for a single data observation affect the final values of the elements of the product matrix  $\mathbf{X}'\mathbf{X}$ . Because of the characteristics of the NDMMF, a particular data observation  $(y_{ijk}, L_{ijk})$  for species/cut  $j$  and sampling-event  $k$  will affect the values in only two of the rows of  $\mathbf{X}'$ : the row for the species/cut  $j$  and the row for the sampling-event  $k$ . Thus, during matrix multiplication, only these two rows from  $\mathbf{X}'$  will have any role in the particular data observation's effect on the final values in the product matrix.

Now consider what happens as each of these two rows from  $\mathbf{X}'$  is multiplied by columns from  $\mathbf{X}$ . There are four cases to consider:

1. **Species/cut row from  $\mathbf{X}'$  and species/cut column from  $\mathbf{X}$ :** The row-column product will be zero for every species/cut column in  $\mathbf{X}$ , except the column for the same species/cut  $j$ . Note that the row from  $\mathbf{X}'$  for species/cut  $j$  is the transpose of the column from  $\mathbf{X}$  for species/cut  $j$ . Therefore, the element of the product matrix  $\mathbf{X}'\mathbf{X}$  at row  $j$  and column  $j$  will be affected by the addition of  $(L_{ijk})^2$ .
2. **Species/cut row from  $\mathbf{X}'$  and sampling-event column from  $\mathbf{X}$ :** The row-column product will be zero for every sampling-event column, except the column for sampling-event  $k$ . Therefore, the element of the product matrix  $\mathbf{X}'\mathbf{X}$  at row  $j$  and column  $n + k$  will be affected by the addition of  $L_{ijk} \times 1$ . (The column in  $\mathbf{X}$  for sampling-event  $k$  is column  $n + k$  because the  $n$  columns for species/cut parameters come first, so counting for sampling-event columns begins at  $n + 1$ .)
3. **Sampling-event row from  $\mathbf{X}'$  and species/cut column from  $\mathbf{X}$ :** The row-column product will be zero for every species/cut column, except the column for species/cut  $j$ . Therefore, the element of the product matrix  $\mathbf{X}'\mathbf{X}$  at row  $n + k$  and column  $j$  will be affected by the addition of  $1 \times L_{ijk}$ .

4. **Sampling-event row from  $X'$  and sampling-event column from  $X$ :** The row-column product will be zero for every sampling-event column, except the column for sampling-event  $k$ . Therefore, the element of the product matrix  $X'X$  at row  $n+k$  and column  $n+k$  will be affected by the addition of  $1 \times 1 = 1$ .

Thus, as a result of the fact that each row of  $X$  contains exactly two nonzero values, any particular data observation ( $y_{ijk}, L_{ijk}$ ) only influences four elements of the product matrix  $X'X$ . Let  $X'X = P$  and let  $p_{a,b}$  denote the element of  $P$  in the  $a$ th row and  $b$ th column. Then any particular data observation ( $y_{ijk}, L_{ijk}$ ) affects only four values in the product matrix  $P$  as follows:

1.  $p_{jj}$  : increased by the addition of  $(L_{ijk})^2$
2.  $p_{j,n+k}$  : increased by the addition of  $L_{ijk}$
3.  $p_{n+k,j}$  : increased by the addition of  $L_{ijk}$
4.  $p_{n+k,n+k}$  : increased by the addition of 1

This result allows the  $(p \times p)$  product matrix  $X'X = P$  to be computed by adding in values as each data observation is read in from an input file. When data input is complete, so is the product matrix  $X'X = P$ . Since  $p_{j,n+k}$  and  $p_{n+k,j}$  always receive the same additive contributions,  $X'X = P$  is a symmetric matrix.

The other product matrix,  $X'y$ , can also be computed during data input. Let  $X'y = q$  and let  $q_a$  denote the element of  $q$  in the  $a$ th row. Note that the product matrix  $X'y = q$  is a  $(p \times 1)$  column matrix (also called a **column vector**) and that there are exactly two rows in  $X'$  with values for any particular data observation:

1. **Species/cut row from  $X'$ :** The row-column product will be zero for every species/cut row in  $X'y$  except the row for species/cut  $j$ . The element of the product matrix  $X'y = q$  at row  $j$  will be affected by the addition of  $L_{ijk} \times y_{ijk}$ .
2. **Sampling-event row from  $X'$ :** The row-column product will be zero for every sampling-event row in  $X'$  except the row for sampling-event  $k$ . The element of the product matrix  $X'y = q$  at row  $n+k$  will be affected by the addition of  $1 \times y_{ijk}$ .

Consequently,  $X'y = q$  can be computed during input by observing that any particular data observation affects only two rows of  $q$  as follows:

1.  $q_j$  : increased by the addition of  $L_{ijk} \times y_{ijk}$
2.  $q_{n+k}$  : increased by the addition of  $y_{ijk}$

## Computing With Weighted Observations

The data used for fitting parameters to the NDMMF (the “NDMMF data”) include weights. The essential idea of a weight for a data observation is that an observation with a weight of 2.0 should have the same effect on the fitting of parameters as it would have if it were replaced by two identical observations, each with a weight of 1.0. The majority of observations in the NDMMF data are for the processing of tissue from a single fish; each such observation has a weight of 1.0.

Although the weights for the NDMMF data are sometimes adjusted nonlinearly into nonintegral values such as 1.65 or 2.8, all weights remain greater than or equal to 1.0 (and as previously mentioned, a majority of them equal 1.0 exactly). For the purpose of making use of the weights associated with NDMMF data observations when computing OLS parameter estimates, the essential idea of weights as simple integral multiples will be generalized to allow for nonintegral weights. In other words, an observation weighted by an integer should have the same effect that the corresponding number of identical observations would have, and nonintegral weights should have effects intermediate between the effects of the bracketing integral weights. For example, an observation with a weight of 2.4 should have an effect between the effect of a weight of 2 and a weight of 3. The previous section of this publication showed that any particular data observation (ignoring its weight for the moment) will contribute additively to exactly four elements of  $X'X = P$  and to exactly two elements of  $X'y = q$ . In principle, because two identical observations would make the same additive contributions twice, and three identical observations would make the same additive contributions three times, the additive contributions for a weighted observation are equal to each of the contributions shown in the preceding section multiplied by the weight. If  $w_{ijk}$  is a real number represented in computation as a floating-point number, then when the weight  $w_{ijk}$  is associated with an observation ( $y_{ijk}, L_{ijk}$ ) the weighted effects of this observation on  $X'X = P$  and  $X'y = q$  are as follows:

1.  $p_{jj}$  : increased by the addition of  $(L_{ijk})^2 \times w_{ijk}$
2.  $p_{j,n+k}$  : increased by the addition of  $L_{ijk} \times w_{ijk}$
3.  $p_{n+k,j}$  : increased by the addition of  $L_{ijk} \times w_{ijk}$
4.  $p_{n+k,n+k}$  : increased by the addition of  $1 \times w_{ijk}$
5.  $q_j$  : increased by the addition of  $L_{ijk} \times y_{ijk} \times w_{ijk}$
6.  $q_{n+k}$  : increased by the addition of  $y_{ijk} \times w_{ijk}$

## Computational Methods for Solving the Normal Equations

There are several different, well-documented computational methods available for solving the normal equations for  $\mathbf{b}$  once the  $(p \times p)$  matrix  $\mathbf{X}'\mathbf{X}$  and the  $(p \times 1)$  column vector  $\mathbf{X}'\mathbf{y}$  from the normal equations  $(\mathbf{X}'\mathbf{X})\mathbf{b} = \mathbf{X}'\mathbf{y}$  are available in computer RAM. By analogy with standard algebra, the solution of the normal equations may be expressed as

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \quad (10)$$

in terms of the left multiplicative inverse of the  $(p \times p)$  matrix  $\mathbf{X}'\mathbf{X}$ ; and one way to compute the solution of the normal equations  $\mathbf{b}$  is to invert the matrix  $\mathbf{X}'\mathbf{X}$  and left-multiply this inverse with  $\mathbf{X}'\mathbf{y}$ . In general, however, it is not necessary in computations to invert the matrix  $\mathbf{X}'\mathbf{X}$  in order to solve for  $\mathbf{b}$ , and when the matrix inverse is not required for additional computations, the inversion is often undesirable because it entails more computation than other methods and is, therefore, relatively slow (Monahan, 2001).

Well-known and widely used methods for solving the normal equations and other systems of linear equations computationally (Monahan, 2001; Press and others, 1992) include the following:

- Gaussian elimination with back substitution;
- Gauss-Jordan elimination with back substitution;
- Cholesky decomposition;
- LU decomposition with back substitution;
- Singular value decomposition; and
- QR decomposition.

Of these methods, no one stands out as the best for all problems though in general LU decomposition comes close to being the presumptive method of choice (Press and others, 1992). Each method has advantages and disadvantages that make it more suitable for some sets of normal equations and less suitable for others. Gauss-Jordan elimination is acceptably efficient when the matrix inverse is required for other operations, but it is three times slower than alternative methods for simply solving a system for a single solution vector  $\mathbf{b}$ . Although Gaussian elimination

is faster than Gauss-Jordan elimination and comparable in speed to LU decomposition, it is seldom used in statistical packages because it is subject to round-off errors and other computational inaccuracies. Cholesky decomposition is the fastest of the commonly used methods (approximately twice as fast as Gaussian elimination or LU decomposition), but it can only be used when the matrix  $\mathbf{X}'\mathbf{X}$  is symmetric and positive definite. Singular value decomposition is a “can’t fail” method that effectively overcomes ill-conditioning and near-zero values in the matrix, but it is slow and, consequently, it is not suited for solving large systems of equations such as those encountered in computing OLS parameters for the NDMMF. Finally, QR decomposition requires about twice as much computation as LU decomposition and is not generally chosen, except in special cases (Press and others, 1992).

Detailed description of the derivation of the normal equations for OLS parameter estimation is beyond the scope of this publication, as are detailed descriptions of the various methods of solving systems of linear equations. These topics are covered thoroughly in a number of texts and other books (Mendenhall, 1968; Draper and Smith, 1966; Monahan, 2001; Press and others, 1992).

## Reference Software

The reference custom software included with this publication is a single program module of structured and commented source code written in the C programming language. This code has been compiled using Version 4.4.4 of GCC and has been executed successfully on a computer workstation running under Version 2.6.33 of the Linux kernel. This reference software uses the method of Gaussian elimination with back substitution to solve the normal equations and find the best-fit parameter estimates for the NDMMF. For a computation involving about 15,400 parameters, the real run time (elapsed wall-clock time) on an otherwise lightly loaded workstation was about 16 hours.

A run time of 16 hours may seem long, especially to those who are accustomed to computing solutions for systems of linear equations with 100 or fewer parameters. The long run time is understandable, however, in view of the fact that the time required to solve a system of linear equations is approximately proportional to the cube of the number of parameters. More precisely, the computation required

for Gaussian elimination or LU decomposition is  $O(\frac{5}{6}p^3)$  additions and multiplications where  $p$  is the number of parameters. Back substitution requires the comparatively modest, additional computation of  $\frac{1}{2}p^2$  multiplications and additions (Press and others, 1992). Doubling  $p$  increases the computation required for Gaussian elimination or LU decomposition by a factor of approximately eight (8). Computing a solution for a model with 15,400 parameters will take about 3,650,000 times longer than computing a solution for a model with 100 parameters.

The reference software uses relatively little memory, except for what is required for the  $(p \times p)$  matrix  $\mathbf{X}'\mathbf{X}$  and the two  $(p \times 1)$  column vectors  $\mathbf{b}$  and  $\mathbf{X}'\mathbf{y}$  where  $p$  is the number of parameters. Thus, the software requires  $O(8p^2 \times 16p)$  bytes of computer memory for  $p$  parameters. For example, when the number of parameters is 15,400, then memory required is 1.767 gigabytes. The memory saved by not storing the  $(N \times p)$  matrix of observations (with  $N \cong 101,000$ ) is 11.589 gigabytes.

Also provided with the C source code are:

- a compilation script,
- a sample data input file, and
- samples of the program's three output files.

The software requires that there be no gaps in the species/cut and sampling-event parameter numbers. Each set of parameter numbers must begin with 1 and end with a number that does not exceed the value of its respective manifest constant—NUMSPC or NUMEVENTS. The values of these two manifest constants, along with the value of the manifest constant MATRIXDIMENSION, must be set correctly before compiling and running the program code. The value of MATRIXDIMENSION must equal or exceed the sum of the values of NUMSPC and NUMEVENTS.

This software does not perform checks on the validity of input data. The software will fail to produce valid results if inputs are not valid. Input concentrations are assumed to be in units of parts per million. Input lengths are assumed to be in inches.

This software is provided as a research tool, not as production code. Use of this code requires basic proficiency in reading, understanding, modifying, and compiling C program code. Setting up the input data file requires careful attention to its contents.

## Summary

This publication describes a technique for efficient computation of ordinary least-squares (OLS) parameter estimates for the National Descriptive Model of Mercury in Fish (NDMMF). Included with this report is a preliminary reference implementation of the technique in software, using the method of Gaussian elimination. The technique enables rapid setup of the normal equations for the NDMMF so that the equations can then be solved for the OLS parameter estimates by any of several available methods.

The characteristic of the NDMMF that enables the rapid and direct setup of the normal equations is its restriction to exactly two parameters for each observation of methylmercury concentration. This characteristic allows the matrix products in the normal equations to be filled out quickly during data input so that the matrix  $\mathbf{X}$  of data values for the independent variables does not need to be stored in computer memory and the full matrix multiplications implied by the normal equations do not actually have to be carried out.

It is not essential that all computation be performed using the provided reference software. The technique is also potentially usable with generalized statistical software. The reference software can be modified by a qualified programmer to omit the computation for Gaussian elimination and just write the two matrices created on data input ( $\mathbf{X}'\mathbf{X}$  and  $\mathbf{b}$ ) to a file or files for subsequent input into any generalized statistical package or program capable of accepting the normal equations in matrix form as file input to a linear-system solver.

Computing the OLS parameter estimates for the NDMMF using the specialized technique described in this publication requires  $O(8p^2 \times 16p)$  bytes of computer memory on a system with 8-byte double-precision numbers, where  $p$  is the number of parameters of the model. In recent usage, for illustration, the custom software, with  $p \cong 15,400$ , used about 1.767 gigabytes of computer memory and required about 16 hours for computation. The computer memory saved by not storing the  $\mathbf{X}$  matrix in simplest form (with  $N \cong 101,000$ ) was 11.589 gigabytes. The computational time saved by avoiding a full matrix multiplication to create  $\mathbf{X}'\mathbf{X}$  would vary by generalized statistical software but could range up to the time required for  $O(Np^2)$  multiplications and additions; so the computational time saved could be a multiple of the  $O(\frac{5}{6}p^3)$  multiplications and additions used to compute the ordinary least-squares parameter estimates by Gaussian elimination. With  $p \cong 15,400$  and  $N \cong 101,000$  the computational time saved could be up to  $\sim 7.87$  times what is used for Gaussian elimination.

## References Cited

- Box, George E.P., and others, 1978, *Statistics for experimenters; an introduction to design, data analysis, and model building*: John Wiley & Sons, Inc., 653 p.
- Donato, D.I., 2012, Computing maximum-likelihood estimates for parameters of the National Descriptive Model of Mercury in Fish: U.S. Geological Survey Open-File Report 2012-1181, 15 p. (Available only at <http://pubs.usgs.gov/of/2012/1181>.)
- Draper, N.R. and Smith, Harry, 1966, *Applied regression analysis*: John Wiley & Sons, Inc., 407 p.
- Helsel, D.R., 2004, *Nondetects and data analysis; statistics for censored environmental data*: Wiley-Interscience, 268 p.
- Mendenhall, William, 1968, *Introduction to linear models and the design and analysis of experiments*: Duxbury Press, 465 p.
- Monahan, J.F., 2001, *Numerical methods of statistics*: Cambridge University Press, 428 p.
- Noble, Ben, 1969, *Applied linear algebra*: Prentice-Hall, Inc., 523 p.
- Press, W.H., and others, 1992, *Numerical recipes in C; the art of scientific computing (2d ed.)*: Cambridge University Press, 994 p.
- Wente, S.P., 2004, A statistical model and national dataset for partitioning fish-tissue mercury concentration variation between spatiotemporal and sample characteristic effects: U.S. Geological Survey Scientific Investigations Report 2004-5199, 15 p. (Available only at <http://pubs.usgs.gov/sir/2004/5199>.)

# Appendix

This appendix contains the C language source computer code (`NdmO1s82.c`) for the reference software that implements and illustrates the technique described in this publication. Also included in this appendix are:

- a compilation script—`comp`,
- a sample data input file—`Hgdata.srt`, and
- samples of the three output files produced by the reference software—`EVENTparameters.txt`, `SPCparameters.txt`, and `log`.

Please take notice of the disclaimers included in the C language source-code file.

Appendix files are available for download at <http://pubs.usgs.gov/tm/07/c10>.

This page intentionally left blank.

Publishing support provided by:  
U.S. Geological Survey Science Publishing Network,  
Reston and Tacoma Publishing Service Centers

For more information concerning the research in this report, contact the  
Eastern Geographic Science Center  
U.S. Geological Survey  
521 National Center  
12201 Sunrise Valley Drive  
Reston, VA 20192  
<http://egsc.usgs.gov>

