

CONTRIBUTIONS TO COMPUTER ARITHMETIC

and Self-Validating numerical methods

C. Ullrich (editor)

J. C. Baltzer AG, Scientific Publishing Co. © IMACS 1990

pp. 133–147

## LEAST-SQUARE APPROXIMATIONS UNDER INTERVAL INPUT DATA

Svetoslav M. MARKOV

Institute of Mathematics and Informatics, Bulgarian Academy of Sciences

“Acad. G. Bonchev” st., block 8, 1113 Sofia, Bulgaria, smarkov@bio.bas.bg

**Abstract.** The method of least-square approximation is considered in the situation when the input data for the dependent variable are given in the form of intervals. For the multi-dimensional linear and for the polynomial cases algorithms with result verification are reported that can also deal with interval input data. Some open problems related to parameter estimation and curve fitting under interval input data are formulated.

**1. Introduction.** The origin of the errors in the observations often lies in the imprecision of the instruments for the experimental measurements. It may be more convenient for an experimental scientist to read-off intervals that contain with guarantee the true values of the measured quantities, than to read-off single numeric values. The new developments in computer arithmetic [4], interval analysis [1], [6] and numerical methods with result verification [5] should encourage experimental scientists to read-off interval-valued experimental data. In our presentation we assume that the experimental measurements for the dependent variables are provided in the form of intervals, which we shall shortly express by saying that we are given interval input data. The case of more general set-valued input data will not be considered here; we shall also not consider the situation involving interval data for the independent variables.

Even if the bounds for the input data are small, it may happen that they cause large deviations in the final results (especially by ill-conditioned problems, and such are often least-square approximation problems). A possible way to treat such problems is to consider corresponding “set-valued problems”. A set-valued problem involves set-valued input data  $X$ ,  $X$  being usually a vector or

a matrix, whose elements are compact sets like  $n$ -dimensional intervals, disks, ellipsoids, polyhedrons, etc. In the special case when these elements are intervals we talk about interval problems. A set-valued problem involving set-valued input data  $X$  (e. g.  $X$  may be a set-valued vector or a set-valued vector function) is considered as the set of problems with all possible numeric data  $x$ , such that  $x \in X$ . In particular, let  $X$  be an interval (vector, function) and  $x$  be a number (numeric vector, single-valued function), such that  $x \in X$ . If  $P(x)$  is the solution of a "numeric" problem that involves numeric input data  $x$ , then the solution of the corresponding interval problem is (by definition) the set of solutions  $P(x)$ , whenever  $x \in X$ , i. e.  $P(X) = \{P(x) : x \in X\}$ . The mathematical analysis of problems involving interval input data (that is of interval problems) is known under the name interval analysis [1], [6]. A typical simple result from interval analysis that we shall exploit in this paper is the following. If  $Y_i, i = 1, 2, \dots, N$ , are intervals, then  $\{\sum_{i=1}^N \alpha_i y_i : y_i \in Y_i\} = \sum_{i=1}^N \alpha_i Y_i$ , where in the right hand-side well-known interval arithmetic operations for addition and multiplication are employed [1], [6].

In the same manner we can treat errors due to the necessarily finite representation of numeric input data (e. g.  $1/3$  should be represented in a base 10 floating-point system by an interval of the form  $[0.33\dots33, 0.33\dots34]$ ). In such cases we replace the numeric input data by interval data but now we may also interfere to make these bounds as tight as we wish (using e. g. extended precision formats, such as STC-formatting technique [8]).

Consider a situation when we are given guaranteed intervals for the observations of a stochastic variable. Such intervals may take into account some systematic error in the experimental data due e. g. to the imprecision of the measuring instruments. In such a situation it may be useful to apply the least-square approximation method directly to the interval input data, obtaining thus (as usually in interval analysis) the set of all approximations corresponding to numeric data varying in the given intervals. In what follows we shall be concerned with the following two aspects of the least square approximation method: i) the treatment of interval input data for the dependent variable (making use of interval-arithmetic), and ii) the treatment of round-off errors (by means of computer arithmetic).

We first recall some well-known facts related to the least-square approximation method under numeric data, considering the most simple linear one-dimensional case. We shall make use of these facts in section 3, where the interval case will be considered.

**2. The one-dimensional linear regression model for numeric input data.** As we know the coefficients  $a$  and  $b$  of the line

$$(1) \quad l : \eta = a\xi + b$$

that fits to the (numeric) input data  $(x, y)$ ,  $x = (x_1, x_2, \dots, x_N) \in \mathbb{R}^N$ ,  $y = (y_1, y_2, \dots, y_N) \in \mathbb{R}^N$ , in such a way that  $\sum (ax_i + b - y_i)^2$  is minimal, are determined by the normal equations

$$\left(\sum x_i^2\right)a + \left(\sum x_i\right)b = \sum x_i y_i, \quad \left(\sum x_i\right)a + Nb = \sum y_i,$$

wherein  $\sum$  means summation from  $i$  to  $N$ . Denoting  $\bar{x} = (\sum x_i)/N$ ,  $\bar{y} = (\sum y_i)/N$ , and dividing the second equation by  $N$  we obtain

$$(2) \quad \left(\sum x_i^2\right)a + N\bar{x}b = \sum x_i y_i, \quad \bar{x}a + b = \bar{y}.$$

The determinant of (2) will be further denoted by

$$S_{xx} = \sum x_i^2 - N\bar{x}^2 = \sum (x_i - \bar{x})^2 > 0.$$

The slope  $a$  of the line  $l$  is

$$a = (1/S_{xx})\left(\sum x_i y_i - N\bar{x}\bar{y}\right),$$

which can also be written

$$a = \left(\sum x_i y_i - \bar{x} \sum y_i\right)/S_{xx} = \left(\sum (x_i - \bar{x})y_i\right)/S_{xx} \stackrel{def}{=} S_{xy}/S_{xx}.$$

For  $b$  we compute  $b = \bar{y} - a\bar{x} = \bar{y} - (S_{xy}/S_{xx})\bar{x}$ , so that  $l$  obtains the form

$$(3) \quad l : \eta = (S_{xy}/S_{xx})(\xi - \bar{x}) + \bar{y},$$

showing that  $l$  passes through the point  $(\bar{x}, \bar{y})$ .

The expression in the right hand-side of (3) can be rewritten in the form:

$$\begin{aligned} L : \eta &= (S_{xy}/S_{xx})(\xi - \bar{x}) + \bar{y} \\ &= (1/S_{xx})\left(\sum (x_i - \bar{x})y_i\right)(\xi - \bar{x}) + \left(\sum y_i\right)/N \\ &= \sum \left((x_i - \bar{x})(\xi - \bar{x})/S_{xx} + 1/N\right) y_i. \end{aligned}$$

Thus the line (3) can be represented in the form

$$(4) \quad l : \eta = \sum \gamma_i(\xi) y_i,$$

wherein the functions

$$(5) \quad \gamma_i(\xi) = \gamma_i(x; \xi) = (x_i - \bar{x})(\xi - \bar{x})/S_{xx} + 1/N, \quad i = 1, 2, \dots, N,$$

depend only on  $x$  and not on  $y$ .

Since  $\gamma_i$  is linear, it may have at most one zero. If  $x_i = \bar{x}$ , then  $\gamma_i = 1/N > 0$ . If  $x \neq \bar{x}$ , then  $\gamma_i(\xi)$  has a slope  $(x_i - \bar{x})/S_{xx}$ , such that

$$(x_i - \bar{x})/S_{xx} \begin{cases} < 0, & x_i < \bar{x}, \\ > 0, & x_i > \bar{x}. \end{cases}$$

Let  $x_i < \bar{x}$  for  $i = 1, 2, \dots, j$  and  $x_i > \bar{x}$  for  $i = j + 1, \dots, N$ . Denoting by  $\xi_i$  the zero of the linear function  $\gamma_i(\xi)$ , i. e.

$$\gamma_i(\xi_i) = (x_i - \bar{x})(\xi_i - \bar{x})/S_{xx} + 1/N = 0, \quad i = 1, \dots, N,$$

we have

$$(6) \quad \xi_i = \bar{x} + S_{xx}/(N(\bar{x} - x_i)), \quad i = 1, \dots, N.$$

The ordering of the  $x_i$ 's with respect to  $\bar{x}$  imply corresponding ordering of  $\xi_i$ 's. Namely, if  $\bar{x}$  lies between the knots  $x_j$  and  $x_{j+1}$ , then we have

$$\begin{aligned} x_1 < x_2 < \dots < x_j < \bar{x} < x_{j+1} < \dots < x_N \\ \implies \xi_{j+1} < \xi_{j+2} < \dots < \xi_N < \bar{x} < \xi_1 < \xi_2 < \dots < \xi_{j-1} < \xi_j. \end{aligned}$$

These relations remain true also for  $x_j = \bar{x}$ , providing that in this case  $\xi_j$  is understood as  $\infty$ , so that we can write

$$(7) \quad \begin{aligned} x_1 < x_2 < x_3 < \dots < x_j \leq \bar{x} < x_{j+1} < \dots < x_N \\ \implies \xi_{j+1} < \xi_{j+2} < \dots < \xi_N < \bar{x} < \xi_1 < \xi_2 < \dots < \xi_{j-1} < \xi_j. \end{aligned}$$

We shall adopt the notations  $D_k$ ,  $k = 0, 1, \dots, N$ , for the intervals with end-points  $\xi_i$  as follows:

$$\begin{aligned} D_{j+1} &= (-\infty, \xi_{j+1}], \quad D_{j+2} = [\xi_{j+1}, \xi_{j+2}], \dots, \quad D_N = [\xi_{N-1}, \xi_N], \\ D_0 &= [\xi_N, \xi_1], \quad D_1 = [\xi_1, \xi_2], \dots, \quad D_{j-1} = [\xi_{j-1}, \xi_j], \quad D_j = [\xi_j, \infty]. \end{aligned}$$

Let us now compute the sign of  $\gamma_i(\xi)$  in the interval  $D_k$ . In the "central interval"  $D_0 = [\xi_N, \xi_1]$  all  $\gamma_i(\xi)$ ,  $i = 1, 2, \dots, N$ , have positive signs. In the remaining intervals we have

i) to the right of  $D_0$ , that is for  $\xi \in D_k$ ,  $1 \leq k \leq j$ :

$$(8) \quad \text{sign } \gamma_i(\xi) = \left\{ \begin{array}{l} -, \quad i = 1, \dots, k, \\ +, \quad i = k + 1, \dots, N \end{array} \right\} = \text{sign } (i - k - 1/2), \quad i = 1, \dots, N;$$

ii) to the left of  $D_0$ , that is for  $\xi \in D_k$ ,  $j + 1 \leq k \leq N$ :

$$(9) \quad \text{sign } \gamma_i(\xi) = \left\{ \begin{array}{l} +, \quad i = 1, \dots, k, \\ -, \quad i = k + 1, \dots, N \end{array} \right\} = \text{sign } (k - i - 1/2), \quad i = 1, \dots, N.$$

**3. Least square approximation under interval input data.** Let us now consider the least square approximation method in the situation when interval-valued experimental data are provided for the true values  $y_i$ ,  $i = 1, 2, \dots, N$ , of the dependent variable  $y$ .

Suppose we are given an  $N$ -dimensional vector  $(x_1, x_2, \dots, x_N) = x \in \mathbb{R}^N$ , such that  $x_1 < x_2 < \dots < x_N$  and an  $N$ -dimensional interval vector  $(Y_1, Y_2, \dots, Y_N) = Y \in I\mathbb{R}^N$  ( $I\mathbb{R}^N$  stands for the set of all  $N$ -dimensional interval vectors). Let  $y \in \mathbb{R}^N$  be such that  $y \in Y$ , and  $l(x, y)$  be the regression line (3) generated by the input data  $(x, y)$ . Denote by  $L$  the family of all regression lines  $l(x, y)$  generated by the input data  $x, y$ , whenever the numeric vector  $y = (y_1, y_2, \dots, y_N)$  varies in the interval vector  $Y = (Y_1, Y_2, \dots, Y_N)$ , that is the set

$$(10) \quad L = L(x, Y) = \{l(x, y) : y \in Y\}.$$

Denote by  $L(\xi) = L(x, Y; \xi)$  the intersection of the set  $L$  by the vertical line through  $\xi$ . We thus define a set-valued (interval-valued) function which we shall denote again by  $L$ . We shall thus use the same notation  $L$  both for the set of regression lines (10) and for the corresponding set-valued function; we hope that no confusion occurs because of this.

**PROBLEM.** Compute a (best possible) inclusion for the interval-valued function  $L$ .

**SOLUTION.** We shall first compute a (rough) inclusion for  $L$ .

According to (3) the line  $l(x, y)$  generated by the data  $x, y$  is the line passing through the point  $m = (\bar{x}, \bar{y})$  and having as slope  $a = a(x, y) = S_{xy}/S_{xx} = (\sum(x_i - \bar{x})y_i)/S_{xx}$ . As  $y_i$  vary in  $Y_i$ ,  $i = 1, 2, \dots, N$ , the point  $m = (\bar{x}, N^{-1} \sum y_i)$  varies in some segment  $(\bar{x}, \bar{Y})$  and the slope  $a$  varies in some interval  $A$ . Using interval arithmetic we obtain

$$\bar{Y} = (1/N) \sum Y_i, \quad A = A(x, Y) = (\sum(x_i - \bar{x})Y_i)/S_{xx} = S_{xY}/S_{xx}.$$

The sets  $A$  and  $\bar{Y}$  are both obtained from the variation of the  $y$ 's and are therefore inter-related. If we consider these sets as independent, we may construct the following simple interval “linear” function

$$\hat{L} = \{a(\xi - \bar{x}) + \bar{y} : a \in A, \bar{y} \in \bar{Y}\} = A(\xi - \bar{x}) + \bar{Y},$$

which contains  $L$ , e. g.  $L(\xi) \subseteq \hat{L}(\xi)$ ,  $\xi \in \mathbb{R}$ . However, since the parameters  $a$  and  $\bar{y}$  are inter-related  $\hat{L}$  will only provide (rough) bounds for  $L$ .

We shall now compute the exact interval value of  $L$  at the point  $\xi$ . To this end we shall make use of representation (4) for  $l$ , that is  $\eta = \sum \gamma_i(\xi)y_i$ . Using again interval arithmetic we obtain for a fixed  $\xi$

$$\begin{aligned} L(\xi) &= \left\{ \sum \gamma_i(\xi)y_i : y \in Y \right\} = \sum \gamma_i(\xi)Y_i \\ (11) \quad &= \left[ \sum \gamma_i(\xi)y_i^{-\text{sign } \gamma_i(\xi)}, \sum \gamma_i(\xi)y_i^{\text{sign } \gamma_i(\xi)} \right] \\ &= [l^-(\xi), l^+(\xi)], \end{aligned}$$

where, according to (5)

$$\gamma_i(\xi) = (x_i - \bar{x})(\xi - \bar{x})/S_{xx} + 1/N, \quad i = 1, 2, \dots, N.$$

In formula (11) the end-points of the intervals  $Y$  are denoted by  $y_i^- \leq y_i^+$ , so that  $Y = [y_i^-, y_i^+]$ ,  $i = 1, 2, \dots, N$ ; also  $\text{sign } \gamma_i(\xi)$  means “+”, if  $\gamma_i(\xi) \geq 0$ , and “-”, if  $\gamma_i(\xi) < 0$ ;  $y_i^{--}$  means  $y_i^+$  (right end-point) and  $y_i^{-+}$  means  $y_i^-$  (left end-point). Formula (11) shows that for every  $\xi$ , the set  $L(\xi)$  is an interval, so that the set-valued function  $L$  may be considered as an interval function.

Denoting  $\gamma(\xi) = (\gamma_1(\xi), \gamma_2(\xi), \dots, \gamma_N(\xi))$  we may rewrite (11) in the form of an interval vector inner product

$$(12) \quad L(\xi) = L(x, Y; \xi) = \{l(x, Y; \xi) : y \in Y\} = \sum \gamma_i(\xi)Y_i = \gamma(\xi)Y.$$

We thus see that by means of interval arithmetic the interval function  $L(\xi)$  is expressed in the simple form  $L(\xi) = \sum \gamma_i(\xi)Y_i = \gamma(\xi)Y$ . This can be easily programmed within a program system providing interval and computer arithmetic operations, which means that the boundary functions  $l^-(\xi), l^+(\xi)$  of the interval function  $L(\xi)$  can be easily computed. However, it is interesting to know what these functions look like, that is what is the geometrical meaning of the expressions (11)–(12) for  $L(\xi)$ .

**4. On the geometrical meaning of the expression  $L(\xi) = \gamma(\xi)Y$  and its computation.** To see the geometrical meaning of the expression for  $L(\xi)$  we

have to know the signs of  $\gamma_i(\xi)$ . As we saw in section 2 these signs are constant in each of the intervals  $D_{j+1}, D_{j+2}, \dots, D_N, D_0, D_1, \dots, D_{j-1}, D_j$ , but may be different for different intervals  $D_k$  according to formulas (8)–(9).

Using (8)–(9) we see that in every fixed interval  $D_k$  the boundaries  $l^-, l^+$  of the set  $L$  are segments of regression lines corresponding to certain end-points of the input intervals  $Y_1, Y_2, \dots, Y_n$ . Thus, in the “central” interval  $D_0$  the boundary regression line  $l^+$  is generated by the set of all right end-points of  $Y_1, Y_2, \dots, Y_n$  and the boundary line  $l^-$  is generated by the set of all left end-points of  $Y_1, Y_2, \dots, Y_n$ , that is

$$L(\xi) = [l^-(\xi), l^+(\xi)] = \left[ \sum \gamma_i(\xi) y_i^-, \sum \gamma_i(\xi) y_i^+ \right].$$

We may also compute the width  $w$  of  $L$  in  $D_0$ . We have

$$w(L(\xi)) = \sum |\gamma_i(\xi)| w(Y_i) = \sum \gamma_i(\xi) (y_i^+ - y_i^-) = \sum \gamma_i(\xi) y_i^+ - \sum \gamma_i(\xi) y_i^-.$$

Let us compute the width of  $L$  on the whole real line under the additional assumption that the intervals  $Y_i$  have a constant width  $W$ . Then we have

$$\begin{aligned} w(L(\xi)) &= \sum |\gamma_i(\xi)| w(Y_i) = W \sum |\gamma_i(\xi)| \\ &= W \sum |(1/S_{xx})(x_i - \bar{x})(\xi - \bar{x}) + 1/N| \\ &\leq W(1/S_{xx}) |\xi - \bar{x}| (1 + \sum |x_i - \bar{x}|) = W + W(1/S_{xx}) |\xi - \bar{x}| \sum |x_i - \bar{x}|. \end{aligned}$$

From the above formula we see that the equality relation is reached in the interval  $D_0$ . Indeed, using that all  $\gamma_i$  are nonnegative in  $D_0$  and the relation  $\sum (x_i - \bar{x}) = 0$ , we obtain

$$\sum |\gamma_i(\xi)| = \sum \gamma_i(\xi) = \sum (x_i - \bar{x})(\xi - \bar{x})/S_{xx} + 1/N = \sum 1/N = 1, \quad \xi \in D_0.$$

It is easy to see that the width of  $L(\xi)$  increases as we move  $\xi$  away from  $\bar{x}$ .

For the midpoint  $\mu(L)$  of the interval  $L(\xi)$  we have

$$\mu(L(\xi)) = \mu(L(x, Y; \xi)) = \sum \gamma_i(x; \xi) \mu(Y_i),$$

showing that the midpoint always lies on the regression line generated by the midpoints of the interval input data  $Y_i$ .

Let us compute the slope of  $L(\xi)$  in the most outer intervals  $D_j$  and  $D_{j+1}$ . In the right-most interval  $D_j$  we have

$$\begin{aligned} L(\xi) &= \left[ \sum \gamma_i(\xi) y_i^- \operatorname{sign} \gamma_i(\xi), \sum \gamma_i(\xi) y_i^+ \operatorname{sign} \gamma_i(\xi) \right] \\ &= \left[ \sum \gamma_i(\xi) y_i^- \operatorname{sign}^{(i-j-1/2)}, \sum \gamma_i(\xi) y_i^+ \operatorname{sign}^{(i-j-1/2)} \right]. \end{aligned}$$

Replacing the expression for  $\gamma_i$  we obtain that the slope of  $L(\xi)$  in  $D_j$  is

$$\left[ \sum (1/S_{xx})(x_i - \bar{x})y_i^{-\text{sign}(i-j-1/2)}, \sum (1/S_{xx})(x_i - \bar{x})y_i^{\text{sign}(i-j-1/2)} \right].$$

For the left-most interval  $D_{j+1}$  we obtain the same expression. On the other side, the slope  $A(x, Y)$  of the interval line

$$\hat{L}(\xi) = A(x, Y)(\xi - \bar{x}) + \bar{Y}$$

is given by

$$\begin{aligned} A(x, Y) &= \{a(x, y) : y \in Y\} \\ &= \left\{ \sum (1/S_{xx})(x_i - \bar{x})y_i : y_i \in Y_i \right\} \\ &= \sum (1/S_{xx})(x_i - \bar{x})Y_i \\ &= \left[ \sum (1/S_{xx})(x_i - \bar{x})y_i^{-\text{sign}(x_i - \bar{x})}, \sum (1/S_{xx})(x_i - \bar{x})y_i^{\text{sign}(x_i - \bar{x})} \right]. \end{aligned}$$

This shows that the slope of  $\hat{L}$  coincides with the slope of  $L$  in the most outer intervals. Taking into account that both  $L$  and  $\hat{L}$  contain the segment  $(\bar{x}, \bar{Y})$  we obtain sufficient information about the geometrical disposition of  $\hat{L}$  with respect to  $L$ .

An algorithm with result verification [5] for the evaluation of the interval-valued function (12) at a fixed point  $\xi$  is straightforward. We first compute highly accurate interval inclusions  $\diamond\gamma_i(\xi)$  for the true values of  $\gamma_i(\xi)$ . Due to conversion errors the input intervals  $Y_i$  may also have to be expanded to machine intervals  $\diamond Y_i$ . Then the interval inner product  $\sum \diamond\gamma_i(\xi) \diamond Y_i$  should be computed by means of a computer arithmetic operation for highly accurate interval inner product [4]. Programs based on such an algorithm has been written both in PASCAL-SC [3] and in FORTRAN within the program package HIFICOMP [2].

**5. Multiple and polynomial regression under interval input data.** A generalization of the above considerations for the situation of a linear dependence of a variable  $\eta$  on  $m$  variables  $\xi_i$ ,  $i = 1, 2, \dots, m$ , of the form

$$(13) \quad \eta = \theta_0 + \theta_1\xi_1 + \theta_2\xi_2 + \dots + \theta_m\xi_m$$

is straightforward. Assume that for the variables  $\xi_1, \xi_2, \dots, \xi_m$  we are given numeric data  $x_{1i}, x_{2i}, \dots, x_{mi}$ ,  $i = 1, \dots, N$ ,  $N > m$  and for the dependent variable  $\eta$  the interval observations  $Y_i$ ,  $i = 1, 2, \dots, N$  are given. Denote



$$X = \begin{pmatrix} 1 & x_{11} & x_{21} & \dots & x_{m1} \\ 1 & x_{12} & x_{22} & \dots & x_{m2} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{1N} & x_{2N} & \dots & x_{mN} \end{pmatrix}, \quad Y = \begin{pmatrix} Y_{m1} \\ Y_{m2} \\ \dots \\ Y_{mN} \end{pmatrix}, \quad y = \begin{pmatrix} y_{m1} \\ y_{m2} \\ \dots \\ y_{mN} \end{pmatrix}$$

and assume that  $y$  is a numeric vector, such that  $y \in Y$ .

Denote  $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_m)^T$ ,  $\xi = (1, \xi_1, \xi_2, \dots, \xi_m)$ . As we know the parameter  $\theta$  of the hyperplane

$$(14) \quad \eta = \xi\theta$$

of best least-square approximation of the points  $(x_i, y_i)$ ,  $x_i = (x_{1i}, x_{2i}, \dots, x_{mi})$ ,  $i = 1, 2, \dots, N$ , satisfies the linear system  $(X^T X)\theta = X^T y$ , and if  $X^T X$  is non-singular, we may write  $\theta = (X^T X)^{-1}(X^T y)$ , so that the hyperplane (14) can be written in the form

$$(15) \quad l : \eta = \xi(X^T X)^{-1}(X^T y).$$

Following the idea of section 3 let us rewrite the above expression in the form

$$(16) \quad l : \eta = \xi((X^T X)^{-1} X^T)y = \Gamma(\xi)y,$$

where  $\Gamma(\xi) = \xi((X^T X)^{-1} X^T)$  is an  $N$ -dimensional vector function of the form  $\Gamma(\xi) = (\gamma_1(\xi), \gamma_2(\xi), \dots, \gamma_N(\xi))$ , such that  $\gamma_i(\xi)$  are linear functions of  $\xi$ . From (16) we can easily obtain in interval arithmetic

$$(17) \quad L = \{l : y \in Y\} = \Gamma(\xi)Y,$$

showing that the set of the hyperplanes corresponding to all  $y \in Y$  can be represented by means of the simple interval arithmetic function  $\Gamma(\xi)Y$ .

Similarly, the parameter  $\theta = (\theta_0, \theta_1, \dots, \theta_m)^T$  of the regression polynomial

$$(18) \quad \eta = \theta_0 + \theta_1\xi + \theta_2\xi^2 + \dots + \theta_m\xi^m$$

can be obtained by formally substituting in the above formulas  $\xi_j$  by  $\xi^j$ ,  $j = 1, \dots, m$ , and the entries  $x_{ji}$  in the matrix  $X$  by the numbers  $x_i^j$ ,  $j = 1, \dots, m$ ,  $i = 1, \dots, N$ . Then again formulas (15)–(17) hold true but now  $\Gamma(\xi) = \xi((X^T X)^{-1} X^T)$  is an  $N$ -dimensional vector function  $\Gamma(\xi) = (\gamma_1(\xi), \gamma_2(\xi), \dots, \gamma_N(\xi))$ , such that  $\gamma_i(\xi)$  are polynomial functions of  $\xi$  of degree  $m$ .

On the base of formula (17) we can construct algorithms with result verification for the safe and accurate computation both of the set of all regression hyper

planes (13) and of the set of all regression polynomials (18) in the same manner as described at the end of the previous section for the computation of the set of all regression lines (10) by means of (12).

Another approach to the polynomial regression problem under interval input data is proposed in [7]. J. Rokne considers an interval-arithmetic setting of the problem, using thereby interval orthogonal polynomials. A modification of Rokne's algorithm for the more general case of polynomials of many variables has been recently proposed [9].

## 6. Some open problems related to curve fitting and parameter identification under interval input data. Relation to interval interpolation.

In what follows we shall formulate some problems that might be of certain practical interest.

**Problem 1** (Robust parameter identification). We saw that the interval linear function  $\hat{L}(\xi) = A(x, Y)(\xi - \bar{x}) + \bar{Y}$  presents an outer approximation of the set  $L$ . However, from practical point of view it may be more interesting to find intervals  $A_1$  and  $Y_1$  for  $a$  and  $\bar{y}$  so that the interval function  $A_1(\xi - \bar{x}) + Y_1$  presents an inner approximation of  $L$  in certain interval for  $\xi$ . Various criteria for such an approximation can be used.

Least-square approximation problems under interval data can be considered in relation to some "interpolation" properties. Consider below an interval vector function  $Y(t) = (Y_0(t), Y_1(t), \dots, Y_N(t))$ , defined for  $t \in [0, T]$ , such that  $\mu(Y(t)) = \text{const}$  and  $w(Y(t))$  is an increasing function on  $t$ , such that  $w(Y(0)) = 0$ . A simple example of such an interval vector function is a function of the form  $Y(t) = Y(0) + [-t, t]$ , for which we have  $w(Y(t)) = 2t$ . It seems to be of practical interest to consider the following problems.

**Problem 2.** Find the smallest  $t$  such that the set-valued function  $L_t = L(x, Y(t); \xi)$  generated by  $Y(t)$  has a nonempty intersection with the intervals  $Y(t)$ . Find the smallest  $t$  such that  $L_t$  contains an element  $l$  (linear, polynomial function), that interpolates the intervals  $Y_i$ , (that is  $l$  passes through the intervals  $Y_i$ ).

We shall further denote by  $p(x, Y)$  the set of all interpolating polynomials  $\{p(x, y) : y \in Y\}$ , where  $p(x, y)$  is the interpolating polynomial for the data  $(x, y)$ ; as before  $p(x, Y)$  may also denote the corresponding interval-valued function.

**Problem 3.** Let the single-valued interpolation polynomial  $p(x, Y(0))$  be a polynomial of  $N$ -th degree but not a polynomial of  $(N - 1)$ -st degree. Compute the smallest  $t$  such that the family of polynomials  $P_t = p(x, Y(t))$  contains a

single-valued polynomial  $p^*$  of degree less than  $N$ . Compute the approximation by  $p^*$  to the single-valued vector  $(x, \mu(Y(0)))$ .

A generalization of this formulation can be considered for interval input data whose centers are not fixed.

**Problem 4.** Let the set  $p(x, Y(t))$  contain a polynomial of approximations  $(N - 1)$ -st degree. Find the largest  $\tau < t$  such that  $p(x, Y(\tau))$  does not contain polynomials of  $(N - 1)$ -st degree.

**Problem 5.** Given the set of interpolating polynomials of  $N$ -th degree  $P = P(x, Y) = \{p(x, y) : y \in Y\}$  for  $x = (x_0, x_1, \dots, x_N)$  and  $Y = (Y_0, Y_1, \dots, Y_N)$ , find the subset of all interpolating polynomials of  $(N - 1)$ -st degree that belong to  $P$ .

**Acknowledgments.** The present research is partially supported by the Committee of Science according to contract No. 755/1988 and by IIASA in the frames of a contracted study agreement “Mathematical modelling of dynamical processes”. Programs for interval least-square approximation and interval interpolation and tools for dynamic precision computation, written by E. Popova, N. Dimitrova and P. Bochev, were used in the present study.

## References

- [1] G. Alefeld, J. Herzberger. Introduction to Interval Computations. Academic Press, 1981.
- [2] HIFICOMP - Subroutine Library for Highly Efficient and Accurate Computations (I.A066.02112-01 13). Program Description and User's Guide, Sofia, 1987.
- [3] U. Kulisch (Ed.). PASCAL-SC: A PASCAL Extension for Scientific Computation, Information Manual and Floppy Disks, Version IBM PC/AT (DOS). Teubner, Stuttgart, 1987.
- [4] U. Kulisch, W. L. Miranker. Computer Arithmetic in Theory and Practice, Academic Press, New York, 1981.
- [5] U. Kulisch, H. J. Stetter (Eds.). A New Approach to Scientific Computation. Academic Press, New York, 1983.
- [6] R. Moore. Methods and Applications of Interval Analysis. SIAM, Philadelphia, 1979.

- [7] J. Rokne. Polynomial Least Square Interval Approximation. *Computing* 20, 165–176 (1978).
- [8] H. J. Stetter. Sequential Defect Correction for High-accuracy Floating Point Algorithms. *Lecture notes in mathematics* 1006, 1984, 186–202.
- [9] R. Trifonov. Estimation of Interval Models Using Orthogonal Multinomials (Lecture delivered at the 2nd Seminar on Dynamical models involving set-valued parameters, 21-23 Oct. 1989, Varna).