

Support Vector Machine for Interval Regression

Dug Hun Hong¹⁾, Changha Hwang²⁾

Abstract

Support vector machine (SVM) has been very successful in pattern recognition and function estimation problems for crisp data. This paper proposes a new method to evaluate interval linear and nonlinear regression models combining the possibility and necessity estimation formulation with the principle of SVM. For data sets with crisp inputs and interval outputs, the possibility and necessity models have been recently utilized, which are based on quadratic programming approach giving more diverse spread coefficients than a linear programming one. SVM also uses quadratic programming approach whose another advantage in interval regression analysis is to be able to integrate both the property of central tendency in least squares and the possibilistic property in fuzzy regression. However this is not a computationally expensive way. SVM allows us to perform interval nonlinear regression analysis by constructing an interval linear regression function in a high dimensional feature space. In particular, SVM is a very attractive approach to model nonlinear interval data. The proposed algorithm here is model-free method in the sense that we do not have to assume the underlying model function for interval nonlinear regression model with crisp inputs and interval output. Experimental results are then presented which indicate the performance of this algorithm.

Key words : Interval Regression Analysis, Possibility and Necessity Models, Quadratic Programming, Support Verctor Machine.

1. Introduction

In many real applications, information is often uncertain, imprecise and incomplete, which can be represented by fuzzy data or a generalization of interval data. For handling interval data, fuzzy regression analysis becomes an important tool and successfully applied to different applications such as market forecasting, system identification, and so on. Fuzzy regression analysis can be simplified to interval regression analysis, where interval regression models are implemented. In interval linear regression, possibility and necessity models have been employed under given interval data. Coefficients in interval regression model are assumed to be interval. In fact, interval regression is regarded as the simplest version of possibilistic regression analysis. Possibilistic regression analysis has been first proposed by Tanaka et al.(1992) where a fuzzy linear system has been used as a regression

1) Dept. of Mathematics, Myungji University

2) Dept. of Statistical Information, Catholic University of Daegu

model. To determine the interval parameters of interval regression models, a basic linear programming (LP) or quadratic programming (QP) problem should be solved.

In previous LP-based approaches, some coefficients become crisp because of the characteristic of LP, since these regression analyses have been reduced to LP problems. To overcome this crisp characteristic of LP, Tanaka and Lee(1998) propose interval regression analysis based on QP. They adopt the dual method for solving their QP formulations. First, they introduce a basic formulation where the sum of squared residuals of the estimated interval outputs is considered as an objective function. This QP based approach is more flexible than the LP based one in the sense that the noncrisp coefficients obtained by QP are more desirable than those by LP. Next, they introduce another QP formulation that combines the property of central tendency in least squares and the possibilistic property in fuzzy regression. This approach has more central tendency than their previous approaches.

In this paper, we propose a new method to evaluate interval linear and nonlinear regression models combining the possibility and necessity estimation formulation with the principle of support vector machine (SVM). This proposed SVM is also based on QP approach giving fairly diverse spread coefficients and integrates well both the property of central tendency in least squares and the possibilistic property in fuzzy regression. However, this is not a computationally expensive way. This SVM allows us to perform interval nonlinear regression analysis by constructing an interval linear regression function in a high dimensional feature space. In particular, SVM is a very attractive approach to model nonlinear interval data. The proposed algorithm here is model-free method in the sense that we do not have to assume the underlying model function for interval nonlinear regression model with crisp inputs and interval output. This model-free method turns out to be a promising method which has been attempted to treat interval nonlinear regression model.

2. Interval Regression for Crisp Input-Interval Output

In this section, we need to briefly look at how to get solutions for interval regression models using QP approach proposed by Tanaka and Lee(1998). For a data set with crisp inputs and interval outputs, we can consider two interval regression models, i.e., the possibility and necessity models. In this section, we review the unified QP approach to obtain the possibility and necessity models simultaneously. In this unified approach, we assume for simplicity that the center coefficients of the possibility regression model and necessity regression model are same.

Suppose that we are given training data $\{(\mathbf{x}_i, Y_i), i = 1, \dots, n\}$, where $\mathbf{x}_i = (1, x_{i1}, \dots, x_{im})^t$ is the i th input vector, $Y_i = (y_i, e_i)$ is the corresponding interval output that consists of a center y_i and a radius e_i . For this data set, the possibility and necessity estimation models are denoted as

$$Y^*(\mathbf{x}_i) = A_0^* + A_1^*x_{i1} + \dots + A_m^*x_{im}, i = 1, \dots, n$$

$$Y_*(\mathbf{x}_i) = A_{0*} + A_{1*}x_{i1} + \dots + A_{m*}x_{im}, i = 1, \dots, n$$

where the interval coefficients A_i^* and A_{i*} are denoted as $A_i^* = (a_i^*, c_i^*)$ and $A_{i*} = (a_{i*}, c_{i*})$, respectively. The estimated interval $Y^*(\mathbf{x}_i)$ by the possibility model always includes the observed interval Y_i , whereas the estimated interval $Y_*(\mathbf{x}_i)$ by the

necessity model should be included in the observed interval Y_i . In fact, we can denote the interval coefficients A_i^* and A_{*i} as

$$A_i^* = (a_i, c_i + d_i), A_{*i} = (a_i, c_i)$$

which satisfies the condition $A_{*i} \subseteq A_i^*, i = 0, 1, \dots, m$ since c_i and d_i are assumed to be positive. Therefore, by interval arithmetic the possibility model $Y^*(x_i)$ and the necessity model $Y_*(x_i)$ can be written as

$$Y^*(x_i) = (a^t x_i, c^t |x_i| + d^t |x_i|), Y_*(x_i) = (a^t x_i, c^t |x_i|).$$

Then, the objective function in the unified approach by QP can be assumed as the following quadratic function:

$$\min_{a,c,d} J = d^t \left(\sum_{i=1}^n |x_i| |x_i|^t \right) d + \xi (a^t b + c^t c).$$

3. Quadratic Loss SVM for Interval Regression

In this section, we propose a new method to evaluate interval linear and nonlinear regression models combining the possibility estimation formulation integrating the property of central tendency with the principle of SVM. We first need to look at how to get solutions for interval linear regression models by implementing the SVM approach. We follow the way of constructing objective function in SVM regression. Then, the objective function can be assumed as the following quadratic function:

$$\min_{a,c} \frac{1}{2} (\|a\|^2 + \|c\|^2 + \|d\|^2) + \frac{\chi}{2} \left(\sum_{i=1}^n \xi_{1i}^2 + \sum_{i=1}^n (\xi_{2i}^2 + \xi_{2i}^{*2}) \right) \quad (1)$$

subject to

$$d^t |x_i| \leq \xi_{1i}$$

$$y_i - a^t x_i \leq \xi_{2i}, \quad a^t x_i - y_i \leq \xi_{2i}^*, \quad i = 1, \dots, n$$

$$a^t x_i + c^t |x_i| + d^t |x_i| \geq y_i + e_i, \quad a^t x_i - c^t |x_i| - d^t |x_i| \leq y_i - e_i, \quad i = 1, \dots, n.$$

$$a^t x_i + c^t |x_i| \leq y_i + e_i, \quad a^t x_i - c^t |x_i| \geq y_i - e_i, \quad i = 1, \dots, n.$$

Although it is possible to use two weight coefficients like Tanaka and Lee(1998), we use one weight coefficient. Here, ξ_{1i} represents spreads of the estimated outputs, and ξ_{2i}, ξ_{2i}^* are slack variables representing upper and lower constraints on the outputs of the model. Hence, we can construct a Lagrange function as follows:

$$\begin{aligned} L = & \frac{1}{2} (\|a\|^2 + \|c\|^2 + \|d\|^2) + \frac{\chi}{2} \left(\sum_{i=1}^n \xi_{1i}^2 + \sum_{i=1}^n (\xi_{2i}^2 + \xi_{2i}^{*2}) \right) \\ & - \sum_{i=1}^n \alpha_{1i} (\xi_{1i} - d^t |x_i|) \\ & - \sum_{i=1}^n \alpha_{2i} (\xi_{2i} - y_i + a^t x_i) - \sum_{i=1}^n \alpha_{2i}^* (\xi_{2i} - a^t x_i + y_i) \\ & - \sum_{i=1}^n \alpha_{3i} (a^t x_i + c^t |x_i| + d^t |x_i| - y_i - e_i) - \sum_{i=1}^n \alpha_{3i}^* (y_i - e_i - a^t x_i + c^t |x_i| + d^t |x_i|) \\ & - \sum_{i=1}^n \alpha_{4i} (y_i + e_i - a^t x_i - c^t |x_i|) - \sum_{i=1}^n \alpha_{4i}^* (a^t x_i - c^t |x_i| - y_i + e_i) \end{aligned} \quad (2)$$

Here, $\alpha_{1i}, \alpha_{2i}, \alpha_{2i}^*, \alpha_{3i}, \alpha_{3i}^*, \alpha_{4i}, \alpha_{4i}^*$ are Lagrange multipliers. It follows from the saddle point condition that the partial derivatives of L with respect to the primal variables

$(\mathbf{a}, \mathbf{c}, \mathbf{d}, \xi_{1i}, \xi_{2i}, \xi_{2i}^*)$ have to vanish for optimality.

$$\frac{\partial L}{\partial \mathbf{a}} = 0 \rightarrow \mathbf{a} = \sum_{i=1}^n (\alpha_{2i} - \alpha_{2i}^*) \mathbf{x}_i + \sum_{i=1}^n (\alpha_{3i} - \alpha_{3i}^*) \mathbf{x}_i - \sum_{i=1}^n (\alpha_{4i} - \alpha_{4i}^*) \mathbf{x}_i \quad (3)$$

$$\frac{\partial L}{\partial \mathbf{c}} = 0 \rightarrow \mathbf{c} = \sum_{i=1}^n (\alpha_{3i} + \alpha_{3i}^*) |\mathbf{x}_i| - \sum_{i=1}^n (\alpha_{4i} + \alpha_{4i}^*) |\mathbf{x}_i| \quad (4)$$

$$\frac{\partial L}{\partial \mathbf{d}} = 0 \rightarrow \mathbf{c} = - \sum_{i=1}^n \alpha_{1i} |\mathbf{x}_i| + \sum_{i=1}^n (\alpha_{3i} + \alpha_{3i}^*) |\mathbf{x}_i| \quad (5)$$

$$\frac{\partial L}{\partial \xi_{1i}} = 0 \rightarrow \xi_{1i} = \frac{1}{8} \alpha_{1i} \quad (6)$$

$$\frac{\partial L}{\partial \xi_{2i}^*} = 0 \rightarrow \xi_{2i}^* = \frac{1}{8} \alpha_{2i}^* \quad (7)$$

Substituting (3)-(7) into (2) yields the dual optimization problem.

$$\begin{aligned} & \text{maximize} \left\{ -\frac{1}{2} \left(\sum_{i,j=1}^n (\alpha_{2i} - \alpha_{2i}^*) (\alpha_{2j} - \alpha_{2j}^*) \mathbf{x}_i^t \mathbf{x}_j \right. \right. \\ & + \sum_{i,j=1}^n (\alpha_{3i} - \alpha_{3i}^*) (\alpha_{3j} - \alpha_{3j}^*) \mathbf{x}_i^t \mathbf{x}_j + \sum_{i,j=1}^n (\alpha_{4i} - \alpha_{4i}^*) (\alpha_{4j} - \alpha_{4j}^*) \mathbf{x}_i^t \mathbf{x}_j \\ & + 2 \sum_{i,j=1}^n (\alpha_{2i} - \alpha_{2i}^*) (\alpha_{3j} - \alpha_{3j}^*) \mathbf{x}_i^t \mathbf{x}_j - 2 \sum_{i,j=1}^n (\alpha_{2i} - \alpha_{2i}^*) (\alpha_{4j} - \alpha_{4j}^*) \mathbf{x}_i^t \mathbf{x}_j \\ & - 2 \sum_{i,j=1}^n (\alpha_{3i} - \alpha_{3i}^*) (\alpha_{4j} - \alpha_{4j}^*) \mathbf{x}_i^t \mathbf{x}_j + \sum_{i,j=1}^n (\alpha_{3i} + \alpha_{3i}^*) (\alpha_{3j} + \alpha_{3j}^*) |\mathbf{x}_i|^t |\mathbf{x}_j| \\ & + \sum_{i,j=1}^n (\alpha_{4i} + \alpha_{4i}^*) (\alpha_{4j} + \alpha_{4j}^*) |\mathbf{x}_i|^t |\mathbf{x}_j| - 2 \sum_{i,j=1}^n (\alpha_{3i} + \alpha_{3i}^*) (\alpha_{4j} + \alpha_{4j}^*) |\mathbf{x}_i|^t |\mathbf{x}_j| \\ & + \sum_{i,j=1}^n \alpha_{1i} \alpha_{1j} |\mathbf{x}_i|^t |\mathbf{x}_j| - 2 \sum_{i,j=1}^n \alpha_{1i} (\alpha_{3j} + \alpha_{3j}^*) |\mathbf{x}_i|^t |\mathbf{x}_j| \left. \right) \\ & + \sum_{i,j=1}^n (\alpha_{3i} + \alpha_{3i}^*) (\alpha_{3j} + \alpha_{3j}^*) |\mathbf{x}_i|^t |\mathbf{x}_j| \Big) - \frac{1}{2C} \sum_{i=1}^n \alpha_{1i}^2 \\ & - \frac{1}{2C} \sum_{i=1}^n \alpha_{1i}^2 - \frac{1}{2C} \sum_{i=1}^n (\alpha_{2i}^2 + \alpha_{2i}^{*2}) \\ & + \sum_{i=1}^n (\alpha_{2i} - \alpha_{2i}^*) y_i + \sum_{i=1}^n (\alpha_{3i} - \alpha_{3i}^*) y_i - \sum_{i=1}^n (\alpha_{4i} - \alpha_{4i}^*) y_i \\ & + \sum_{i=1}^n (\alpha_{3i} + \alpha_{3i}^*) e_i - \sum_{i=1}^n (\alpha_{4i} + \alpha_{4i}^*) e_i \left. \right\} \end{aligned} \quad (8)$$

subject to

$$\alpha_{1i}, \alpha_{ki}, \alpha_{ki}^* \geq 0, \quad k=2,3,4.$$

Solving (8) with above constraints determines the Lagrange multipliers, $\alpha_{1i}, \alpha_{ki}, \alpha_{ki}^*$.

Hence, if $\mathbf{c}'|\mathbf{x}| \geq 0$ and $\mathbf{d}'|\mathbf{x}| \geq 0$, then the linear interval regression function is as follows:

$$Y^*(\mathbf{x}) = (\mathbf{a}'\mathbf{x}, \mathbf{c}'|\mathbf{x}| + \mathbf{d}'|\mathbf{x}|) \quad (9)$$

$$Y_*(\mathbf{x}) = (\mathbf{a}'\mathbf{x}, \mathbf{c}'|\mathbf{x}|) \quad (10)$$

Next, we will consider nonlinear interval regression model. In contrast to linear interval regression, there have been no articles on nonlinear interval regression. In this paper we treat nonlinear interval regression, without assuming the underlying model function. In the case where a linear regression function is inappropriate SVM makes algorithm nonlinear. This could be achieved by simply preprocessing input patterns \mathbf{x}_i by a map $\Phi: R^d \rightarrow E$ into some feature space E and then applying SVM regression algorithm. This is an astonishingly straightforward way.

First notice that the only way in which the data appears in (8) is in the form of inner products $\mathbf{x}_i^t \mathbf{x}_j, |\mathbf{x}_i|^t |\mathbf{x}_j|$. The algorithm would only depend on the data through dot

products in E , i.e. on functions of the form $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^t \Phi(\mathbf{x}_j)$, $K(|\mathbf{x}_i|, |\mathbf{x}_j|) = \Phi(|\mathbf{x}_i|)^t \Phi(|\mathbf{x}_j|)$. The well used kernels for regression problem are given below.

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^t \mathbf{y} + 1)^p, \quad K(\mathbf{x}, \mathbf{y}) = e^{-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}}.$$

Here, p and σ^2 are kernel parameters. In final, the nonlinear interval regression solution is given by

$$\begin{aligned} & \text{maximize} \left\{ -\frac{1}{2} \left(\sum_{i,j=1}^n (a_{2i} - a_{2i}^*)(a_{2j} - a_{2j}^*) K(\mathbf{x}_i, \mathbf{x}_j) \right. \right. \\ & + \sum_{i,j=1}^n (a_{3i} - a_{3i}^*)(a_{3j} - a_{3j}^*) K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i,j=1}^n (a_{4i} - a_{4i}^*)(a_{4j} - a_{4j}^*) K(\mathbf{x}_i, \mathbf{x}_j) \\ & + 2 \sum_{i,j=1}^n (a_{2i} - a_{2i}^*)(a_{3j} - a_{3j}^*) K(\mathbf{x}_i, \mathbf{x}_j) - 2 \sum_{i,j=1}^n (a_{2i} - a_{2i}^*)(a_{4j} - a_{4j}^*) K(\mathbf{x}_i, \mathbf{x}_j) \\ & - 2 \sum_{i,j=1}^n (a_{3i} - a_{3i}^*)(a_{4j} - a_{4j}^*) K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i,j=1}^n (a_{3i} + a_{3i}^*)(a_{3j} + a_{3j}^*) K(|\mathbf{x}_i|, |\mathbf{x}_j|) \\ & + \sum_{i,j=1}^n (a_{4i} + a_{4i}^*)(a_{4j} + a_{4j}^*) K(|\mathbf{x}_i|, |\mathbf{x}_j|) - 2 \sum_{i,j=1}^n (a_{3i} + a_{3i}^*)(a_{4j} + a_{4j}^*) K(|\mathbf{x}_i|, |\mathbf{x}_j|) \\ & + \sum_{i,j=1}^n a_{1i} a_{1j} K(|\mathbf{x}_i|, |\mathbf{x}_j|) - 2 \sum_{i,j=1}^n a_{1i} (a_{3j} + a_{3j}^*) K(|\mathbf{x}_i|, |\mathbf{x}_j|) \left. \right\} \quad (12) \\ & + \sum_{i,j=1}^n (a_{3i} + a_{3i}^*)(a_{3j} + a_{3j}^*) K(|\mathbf{x}_i|, |\mathbf{x}_j|) \left. \right) - \frac{1}{2C} \sum_{i=1}^n a_{1i}^2 \\ & - \frac{1}{2C} \sum_{i=1}^n a_{1i}^2 - \frac{1}{2C} \sum_{i=1}^n (a_{2i}^2 + a_{2i}^{*2}) \\ & + \sum_{i=1}^n (a_{2i} - a_{2i}^*) y_i + \sum_{i=1}^n (a_{3i} - a_{3i}^*) y_i - \sum_{i=1}^n (a_{4i} - a_{4i}^*) y_i \\ & + \sum_{i=1}^n (a_{3i} + a_{3i}^*) e_i - \sum_{i=1}^n (a_{4i} + a_{4i}^*) e_i \left. \right\} \end{aligned}$$

subject to

$$a_{1i}, a_{ki}, a_{ki}^* \geq 0, \quad k = 2, 3, 4.$$

Solving (12) with the above constraints determines the Lagrange multipliers, a_{1i}, a_{ki}, a_{ki}^* .

Therefore, the interval nonlinear regression function is given as follows:

$$Y^*(\mathbf{x}) = \left(\sum_{i=1}^n [(a_{2i} - a_{2i}^*) + (a_{3i} - a_{3i}^*) - (a_{4i} - a_{4i}^*)] K(\mathbf{x}_i, \mathbf{x}), \right. \\ \left. \sum_{i=1}^n [-a_{1i} + 2(a_{3i} + a_{3i}^*) - (a_{4i} + a_{4i}^*)] K(|\mathbf{x}_i|, |\mathbf{x}|), \right) \quad (13)$$

$$Y_*(\mathbf{x}) = \left(\sum_{i=1}^n [(a_{2i} - a_{2i}^*) + (a_{3i} - a_{3i}^*) - (a_{4i} - a_{4i}^*)] K(\mathbf{x}_i, \mathbf{x}), \right. \\ \left. \sum_{i=1}^n [(a_{3i} + a_{3i}^*) - (a_{4i} + a_{4i}^*)] K(|\mathbf{x}_i|, |\mathbf{x}|), \right) \quad (14)$$

References

- P. Diamond, H. Tanaka(1998), Fuzzy regression analysis, in: R. Slowinski (Ed.), Fuzzy Sets in Decision Analysis, Operations Research and Statistics, Kluwer, Boston, 349-387.
- D.H Hong, C. Hwang(2003), Support vector fuzzy regression machines, Fuzzy Sets and

Support Vector Machine for Interval Regression

Systems, 138, 271-281.

H. Tanaka(1987), Fuzzy data analysis by possibilistic linear models, *Fuzzy Sets and Systems*, 24, 363-375.

H. Tanaka, I. Hayashi, K. Nagasaka(1988), Interval regression analysis by possibilistic measures (in Japanese), *The Japanese Journal of Behaviormetrics*, 16, 1-7.

H. Tanaka, H. Lee(1998), Interval regression analysis by quadratic programming approach, *IEEE Trans. Fuzzy Syst.*, 6, 473-481.

H. Tanaka, S. Uejima, K. Asia(1992), Linear regression analysis with fuzzy model, *IEEE Trans. Man. Cybernet.* 12, 903-907.

V. Vapnik(1998), *Statistical Learning Theory*, Springer, Germany.