

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.Doi Number

An Efficient ν -minimum Absolute Deviation Distribution Regression Machine

Yan Wang¹, Yao Wang¹, Yingying Song¹, Xuping Xie¹, Lan Huang^{1,2*}, Wei Pang^{3,4*}, George M. Coghill⁴

¹ Key Laboratory of Symbol Computation and Knowledge Engineering, Ministry of Education, College of Computer Science and Technology, Jilin University, Changchun, 130012, P. R. China

² Zhuhai Laboratory of Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Department of Computer Science and Technology, Zhuhai College of Jilin University, Zhuhai, 519041, China

³ School of Mathematical and Computer Sciences, Heriot-Watt University, EH14 4AS, United Kingdom

⁴ Department of Computing Science, University of Aberdeen, AB24 3UE, United Kingdom

Corresponding authors: Lan Huang(huanglan@jlu.edu.cn); Wei Pang(w.pang@hw.ac.uk)

This research was funded by the National Natural Science Foundation of China (Nos. 61702214, 61772227), the Development Project of Jilin Province of China (Nos 20180414012GH, 20180201045GX, 2020C003), Guangdong Key Project for Applied Fundamental Research (Grant 2018KZDXM076). This work was also supported by Jilin Provincial Key Laboratory of Big Data Intelligent Computing (No. 20180622002JC).

ABSTRACT Support Vector Regression (SVR) and its variants are widely used regression algorithms, and they have demonstrated high generalization ability. This research proposes a new SVR-based regressor: ν -minimum absolute deviation distribution regression (ν -MADR) machine. Instead of merely minimizing structural risk, as with ν -SVR, ν -MADR aims to achieve better generalization performance by minimizing both the absolute regression deviation mean and the absolute regression deviation variance, which takes into account the positive and negative values of the regression deviation of sample points. For optimization, we propose a dual coordinate descent (DCD) algorithm for small sample problems, and we also propose an averaged stochastic gradient descent (ASGD) algorithm for large-scale problems. Furthermore, we study the statistical property of ν -MADR that leads to a bound on the expectation of error. The experimental results on both artificial and real datasets indicate that our ν -MADR has significant improvement in generalization performance with less training time compared to the widely used ν -SVR, LS-SVR, ε -TSVR, and linear ε -SVR. Finally, we open source the code of ν -MADR at <https://github.com/AsunaYY/v-MADR> for wider dissemination.

INDEX TERMS ν -Support vector regression, absolute regression deviation mean, absolute regression deviation variance, dual coordinate descent algorithm.

I. INTRODUCTION

Support vector regression (SVR) [1-3] has been widely used in machine learning, since it can achieve better structural risk minimization. SVR realizes linear regression mainly by constructing linear decision functions in high dimensional space. Compared with other regression methods, such as least square regression [4], Neural Networks (NN) regression [5], logistic regression [6], and ridge regression [7], SVR has better generalization ability for regression problems [8-10]. In recent years, there have been many studies about SVR-based algorithms. Several SVR approaches have been developed, such as ε -support vector regression (ε -SVR) [1, 11], ν -support vector regression (ν -SVR) [12], and least square support vector regression (LS-SVR) [13, 14]. The basic idea of these methods is to find the decision function by

maximizing the boundaries of two parallel hyperplanes. Different from ε -SVR, ν -SVR introduces another parameter, ν , to control the number of support vectors and adjust the parameter ε automatically. The parameter ν has a certain range of values, that is, (0,1]. When solving the quadratic specification problem (QPP), ν -SVR reduces the number of computational parameters by half, which greatly reduces the computational complexity. Besides, some researchers have proposed the non-parallel planar regressors, such as twin support vector regression (TSVR) [15], ε -twin support vector regression (ε -TSVR) [16], parametric-insensitive nonparallel support vector regression (PINSVR) [17], lagrangian support vector regression [18], and lagrangian twin support vector

regression (LTSVR) [19]. These algorithms demonstrate good ability to capture data structure and boundary information.

Support vector (SV) theory indicates that maximizing the minimum margin is not the only way to construct the separating hyperplane for SVM. Zhang and Zhou [20-23] proposed the large margin distribution machine (LDM), which was designed to maximize the margin mean and minimize the margin variance simultaneously. Gao and Zhou [23] proved that the margin distribution including the margin mean and the margin variance was crucial for generalization compared to a single margin, and optimizing the margin distribution can also naturally accommodate class imbalance and unequal misclassification costs [21]. Inspired by the idea of LDM, Liu *et al.* proposed a minimum deviation distribution regression (MDR) [24], which introduced the statistics of regression deviation into ε -SVR. More specifically, MDR minimizes the regression deviation mean and the regression deviation variance while optimizing the minimum margin. In addition, Reshma and Pritam were also inspired by the idea of LDM, and they proposed a large-margin distribution machine-based regression model (LDMR) and a new loss function [25, 26]. However, the definition of the deviation mean in MDR is not very appropriate for positive and negative samples, and the speed of ε -SVR strategy that MDR used can be further improved.

Considering the above advances in SVR, in this research, we introduce the statistical information into ν -SVR and propose an ν -minimum absolute deviation distribution regression (ν -MADR). We give the definition of regression deviation mean which takes into account both the positive and negative values of the regression deviation of sample points. Inspired by recent theoretical results [20-24], ν -MADR simultaneously minimizes the absolute regression deviation mean and the absolute regression deviation variance based on the ν -SVR strategy, thereby greatly improving the generalization performance [21, 23]. To solve the optimization problem, we propose a dual coordinate descent (DCD) algorithm for small sample problems, and we also propose an averaged stochastic gradient descent (ASGD) algorithm for large-scale problems. Furthermore, the boundary on error expectation of ν -MADR is studied. The performance of ν -MADR is assessed on both artificial and real datasets in comparison with other typical regression algorithms, such as ν -SVR, LS-SVR, ε -TSVR, and linear ε -SVR. According to previous research, SVR-based algorithms show better generalization ability for regression problems [8-10]. In conclusion, our experimental results demonstrate that the proposed ν -MADR can lead to better performance than other algorithms for regression problems. The main contributions of this paper are as follows:

- 1) We propose a new regression algorithm that minimizes both the absolute regression deviation mean and the absolute regression deviation variance, and this new algorithm takes into account the positive and negative values of the regression deviation of sample points.

- 2) We propose two optimization algorithms, i.e., the dual coordinate descent (DCD) algorithm for small samples problems and the averaged stochastic gradient descent (ASGD) algorithm for large-scale problems.
- 3) We theoretically prove the upper bound on the generalization error of ν -MADR and analyze the computational complexity of our optimization algorithms.

As SVR-based algorithms are widely used for regression problems, ν -MADR has great application potential.

The rest of this paper is organized as follows: Section 2 introduces the notations used in this paper and presents a brief review of SVR as well as the recent progress in SV theory. Section 3 introduces the proposed ν -MADR, including the kernel and the bound on the expectation of error. Experimental results are reported in Section 4, and finally, the conclusions are drawn in Section 5.

II. BACKGROUND

Suppose $\mathbf{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ is a training set of n samples, where $\mathbf{x}_i \in \mathcal{X}$ is the input sample in the form of d -dimensional vectors and $y_i \in \mathcal{R}$ is the corresponding target value. The objective function is $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$, where $\mathbf{w} \in \mathcal{R}^m$ is the weight vector, $b \in \mathcal{R}$ is the bias term, and $\phi(\mathbf{x})$ is the mapping function induced by a kernel κ , i.e., $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$. To reduce the complexity brought by b , we enlarge the dimension of \mathbf{w} and $\phi(\mathbf{x}_i)$ as in [27], i.e., $\mathbf{w} \leftarrow [\mathbf{w}, b]^T$, $\phi(\mathbf{x}_i) \leftarrow [\phi(\mathbf{x}_i), 1]$. Thus, the function $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ becomes the following form:
 $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$.

In what follows, we only consider problems in the form of the above function.

Formally, we denote \mathbf{X} as the matrix whose i -th column is $\phi(\mathbf{x}_i)$, i.e., $\mathbf{X} = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]$, and $\mathbf{y} = [y_1, \dots, y_n]^T$ is a column vector.

A. THE SVR ALGORITHMS

There are two traditional methods for solving support vector regression (SVR) algorithms, namely ε -SVR [1, 11] and ν -SVR [12]. In order to find the best fitting surface, ε -SVR maximizes the minimum margin containing the data in the so-called ε -tube, in which the distances of the data to the fitting hyperplane are not larger than ε . Therefore, ε -SVR with soft-margin can be expressed as follows:

$$\begin{aligned} & \min_{\mathbf{w}, \xi, \xi^*} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C(\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \\ & s. t. \mathbf{y} - \mathbf{X}^T \mathbf{w} \leq \varepsilon \mathbf{e} + \xi, \\ & \mathbf{X}^T \mathbf{w} - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \xi, \xi^* \geq \mathbf{0}, \end{aligned}$$

where parameter C is used for the tradeoff between the flatness of $f(\mathbf{x})$ and the tolerance of the deviation larger than ε ; $\xi = [\xi_1, \xi_2, \dots, \xi_n]$ and $\xi^* = [\xi_1^*, \xi_2^*, \dots, \xi_n^*]$ are the slack variables measuring the distances of the training samples

outside the ε -tube from the ε -tube itself as soft-margin; \mathbf{e} stands for the all-one vector of appropriate dimensions.

The dual problem of ε -SVR is formulated as

$$\begin{aligned} \min_{\alpha, \alpha^*} & \frac{1}{2}(\alpha - \alpha^*)^T \mathbf{Q}(\alpha - \alpha^*) + \varepsilon(\alpha + \alpha^*) + \mathbf{y}^T(\alpha - \alpha^*) \\ \text{s. t. } & \mathbf{e}^T(\alpha - \alpha^*) = 0, \quad 0 \leq \alpha_i, \alpha_i^* \leq C, \quad i = 1, 2, \dots, n, \end{aligned} \quad (1)$$

where α and α^* are the Lagrange multipliers; $Q_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.

In order to facilitate the calculation, Formula (1) can be transformed as follows:

$$\begin{aligned} \min_{\tilde{\alpha}} & \frac{1}{2} \tilde{\alpha}^T \begin{bmatrix} \mathbf{Q} & -\mathbf{Q} \\ -\mathbf{Q} & \mathbf{Q} \end{bmatrix} \tilde{\alpha} + \begin{bmatrix} \varepsilon \mathbf{e} + \mathbf{y} \\ \varepsilon \mathbf{e} - \mathbf{y} \end{bmatrix}^T \tilde{\alpha} \\ \text{s. t. } & \begin{bmatrix} \mathbf{e} \\ -\mathbf{e} \end{bmatrix}^T \tilde{\alpha} = 0, \quad 0 \leq \tilde{\alpha}_i \leq C, \quad i = 1, 2, \dots, 2n, \end{aligned} \quad (2)$$

where $\tilde{\alpha} = [\alpha^T, \alpha^{*T}]^T$.

v -SVR [12] is another commonly used algorithm for solving SVR. Compared with ε -SVR, v -SVR uses a new parameter $v \in (0, 1]$ to control the number of support vectors and training errors and adjust parameter ε automatically. According to Gu *et al.*, the objective function $f(\mathbf{x})$ in v -SVR is represented by the following constrained minimization problem with soft-margin [28-30]:

$$\begin{aligned} \min_{\mathbf{w}, \varepsilon, \xi, \xi^*} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \left(v\varepsilon + \frac{1}{n}(\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \right) \\ \text{s. t. } & \mathbf{y} - \mathbf{X}^T \mathbf{w} \leq \varepsilon \mathbf{e} + \xi, \\ & \mathbf{X}^T \mathbf{w} - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \xi, \xi^* \geq \mathbf{0}, \quad \varepsilon \geq 0. \end{aligned}$$

The dual problem of v -SVR is

$$\begin{aligned} \min_{\alpha, \alpha^*} & \frac{1}{2}(\alpha - \alpha^*)^T \mathbf{Q}(\alpha - \alpha^*) + \mathbf{y}^T(\alpha - \alpha^*) \\ \text{s. t. } & \mathbf{e}^T(\alpha - \alpha^*) = 0, \quad \mathbf{e}^T(\alpha + \alpha^*) \leq Cv, \\ & 0 \leq \alpha_i, \alpha_i^* \leq \frac{C}{n}, \quad i = 1, 2, \dots, n. \end{aligned}$$

According to Chang *et al.* and Crisp *et al.*, the inequality $\mathbf{e}^T(\alpha + \alpha^*) \leq Cv$ in v -SVR can be replaced by the equality form of $\mathbf{e}^T(\alpha + \alpha^*) = Cv$ with the constraint $0 < v \leq 1$ [11, 31], so we have

$$\begin{aligned} \min_{\alpha, \alpha^*} & \frac{1}{2}(\alpha - \alpha^*)^T \mathbf{Q}(\alpha - \alpha^*) + \mathbf{y}^T(\alpha - \alpha^*) \\ \text{s. t. } & \mathbf{e}^T(\alpha - \alpha^*) = 0, \quad \mathbf{e}^T(\alpha + \alpha^*) = Cv, \\ & 0 \leq \alpha_i, \alpha_i^* \leq \frac{C}{n}, \quad i = 1, 2, \dots, n. \end{aligned} \quad (3)$$

We substitute the equation $\alpha^* = C\mathbf{v}\mathbf{e} - \alpha$ into Formula (3), and Formula (3) can be written as follows:

$$\begin{aligned} \min_{\alpha} & \frac{1}{2}(2\alpha - C\mathbf{v}\mathbf{e})^T \mathbf{Q}(2\alpha - C\mathbf{v}\mathbf{e}) + \mathbf{y}^T(2\alpha - C\mathbf{v}\mathbf{e}) \\ \text{s. t. } & \mathbf{e}^T(2\alpha - C\mathbf{v}\mathbf{e}) = 0, \end{aligned} \quad (4)$$

$$0 \leq \alpha_i \leq \frac{C}{n}, \quad i = 1, 2, \dots, n.$$

As one can see from Formula (2) and (4), by substituting the equation $\alpha^* = C\mathbf{v}\mathbf{e} - \alpha$ into the dual problem, the number of computational parameters of the v -SVR has been reduced by half compared to ε -SVR when solving the QPP. The difference in both time complexity and spatial complexity between ε -SVR and v -SVR can be expressed as follows:

$$\frac{O(\varepsilon - \text{SVR})}{O(v - \text{SVR})} = O\left(\frac{\text{Formula(2)}}{\text{Formula(4)}}\right) = O\left(\frac{2n * 2n}{n * n}\right) = 4.$$

B. RECENT PROGRESS IN SV THEORY

Recent SV theory indicates that maximizing the minimum margin is not the only way to construct the separating hyperplane for SVR, because it does not necessarily lead to better generalization performance [20]. There may exist the so-called data piling problem in SVR [32], that is, the separating hyperplane produced by SVR tends to maximize data piling, which makes the data pile together when they are projected onto the hyperplane. If the distribution of the boundary data is different from that of the internal data, the hyperplane constructed by SVR will be inconsistent with the actual data distribution, which reduces the performance of SVR.

Fortunately, Gao and Zhou have demonstrated that marginal distribution was critical to the generalization performance [23]. By using the margin mean and the margin variance, the model is robust to different distributions of boundary data and noise. Inspired by the above research, MDR [24] introduced the statistics of deviation into ε -SVR and this allows more data to have impact on the construction of the hyperplane.

In MDR, the regression deviation of sample (\mathbf{x}_i, y_i) is formulated as

$$y_i = y_i - f(\mathbf{x}_i), \quad \forall i = 1, \dots, n. \quad (5)$$

So, the regression deviation mean is

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i)) = \frac{1}{n} \mathbf{e}^T(\mathbf{y} - \mathbf{X}^T \mathbf{w}),$$

and the regression deviation variance is defined as

$$\begin{aligned} \hat{y}^2 &= \left(\frac{1}{n} \sqrt{\sum_{i=1}^n \sum_{j=1}^n [y_i - f(\mathbf{x}_i) - y_j + f(\mathbf{x}_j)]^2} \right)^2 \\ &= \frac{1}{n^2} \{2[\mathbf{w}^T \mathbf{X}(\mathbf{n}\mathbf{I} - \mathbf{e}\mathbf{e}^T)\mathbf{X}^T \mathbf{w} - 2\mathbf{y}^T(\mathbf{n}\mathbf{I} - \mathbf{e}\mathbf{e}^T)\mathbf{X}^T \mathbf{w} \\ &\quad + \mathbf{y}^T(\mathbf{n}\mathbf{I} - \mathbf{e}\mathbf{e}^T)\mathbf{y}]\}. \end{aligned}$$

MDR minimizes the regression deviation mean and the regression deviation variance simultaneously, so we have the following primal problem of soft-margin MDR:

$$\begin{aligned} & \min_{\mathbf{w}, \xi, \xi^*} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \lambda_1 \hat{\gamma}^2 + \lambda_2 \bar{\gamma}^2 + C(\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \\ & \text{s. t. } \mathbf{y} - \mathbf{X}^T \mathbf{w} \leq \varepsilon \mathbf{e} + \xi, \\ & \quad \mathbf{w} - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \quad \xi, \xi^* \geq \mathbf{0}, \end{aligned}$$

where λ_1 and λ_2 are the parameters for trading-off the regression deviation variance, the regression deviation mean and the model complexity.

Here, we can see from Equation (5) that the regression deviation, γ_i , is positive when the sample (\mathbf{x}_i, y_i) lies above the regressor and negative when the sample (\mathbf{x}_i, y_i) lies under the regressor. But in fact, for regression, the regression deviation of the sample (\mathbf{x}_i, y_i) is the distance between the actual value and the estimated one, that is, $\gamma_i = |y_i - f(\mathbf{x}_i)|, \forall i = 1, \dots, n$. Therefore, the definition of the deviation mean in MDR here is not very appropriate.

On the other hand, when solving QPP, MDR uses the ε -SVR strategy, and it needs to calculate $2n$ (n is the number of training samples) parameters. Calculating a large number of parameters will increase the computational complexity and reduce the speed of the algorithm. Considering this, in the remainder of this paper, we will introduce our latest advances in SV theory and address the limitations of ε -SVR strategy.

III. v -MINIMUM ABSOLUTE DEVIATION DISTRIBUTION REGRESSION

In this section, we first formulate the absolute deviation distribution which takes into account the positive and negative values of the regression deviation of samples. Then we give the optimization algorithms and the theoretical proof.

A. FORMULATION OF v -MADR

The two most straightforward statistics for characterizing the absolute deviation distribution are the mean and the variance of absolute deviation. In regression problems, the absolute regression deviation of sample (\mathbf{x}_i, y_i) is formulated as

$$\varphi_i = |y_i - f(\mathbf{x}_i)|, \forall i = 1, \dots, n. \quad (6)$$

φ_i is actually the distance between the actual value of the sample (\mathbf{x}_i, y_i) and the estimated one. According to the definition in Equation (6), we give the definitions of statistics of absolute deviation in regression.

Definition 1. Absolute regression deviation mean is defined as follows:

$$\begin{aligned} \bar{\varphi} &= \frac{1}{n} \sum_{i=1}^n \varphi_i^2 = \frac{1}{n} \sum_{i=1}^n |y_i - f(\mathbf{x}_i)|^2 \\ &= \frac{1}{n} (\mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w} - 2\mathbf{y}^T \mathbf{X}^T \mathbf{w} + \mathbf{y}^T \mathbf{y}). \end{aligned} \quad (7)$$

The absolute regression deviation mean actually represents the expected value of difference between the actual values of data and the estimated ones. In order to facilitate the calculation, we have done a square process in this definition. In fact, we can view the absolute regression deviation mean as the adjusted distances of data to their fitting hyperplane. Next,

we give the concept of the absolute regression deviation variance as follows:

Definition 2. Absolute regression deviation variance is defined as follows:

$$\begin{aligned} \hat{\varphi} &= \left(\frac{1}{n} \sqrt{\sum_{i=1}^n \sum_{j=1}^n |y_i - f(\mathbf{x}_i) - y_j + f(\mathbf{x}_j)|^2} \right)^2 \\ &= \frac{2}{n^2} [\mathbf{w}^T \mathbf{X} (\mathbf{nI} - \mathbf{e} \mathbf{e}^T) \mathbf{X}^T \mathbf{w} - 2\mathbf{y}^T (\mathbf{nI} - \mathbf{e} \mathbf{e}^T) \mathbf{X}^T \mathbf{w} + \mathbf{y}^T (\mathbf{nI} - \mathbf{e} \mathbf{e}^T) \mathbf{y}]. \end{aligned} \quad (8)$$

We can see that the absolute regression deviation variance quantifies the scatter of regression.

Existing SVR's loss is calculated only if the absolute value of the difference between the actual data and the estimated values is greater than a threshold. The fitting hyperplane constructed by SVR is only affected by the distribution of the boundary data. If the distribution of the boundary data largely deviates from that of the internal data, the hyperplane constructed will be inconsistent with the actual overall data distribution. To overcome this issue, v -MADR aims to obtain a tradeoff between the distribution of the boundary data and that of the internal data. This means that the fitting hyperplane constructed by v -MADR is not only determined by the distribution of the boundary data, but also measures the influence of the overall data distribution on the fitting hyperplane by simultaneously minimizing the absolute regression deviation mean and the absolute regression deviation variance, which is closer to the real distribution for many datasets and is more robust to noise.

Therefore, similar to the soft-margin of v -SVR [28], the final optimization problem considering the soft-margin has the following form:

$$\begin{aligned} & \min_{\mathbf{w}, \varepsilon, \xi, \xi^*} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \lambda_1 \hat{\varphi} + \lambda_2 \bar{\varphi} + C \left(v\varepsilon + \frac{1}{n} (\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \right) \\ & \text{s. t. } \mathbf{y} - \mathbf{X}^T \mathbf{w} \leq \varepsilon \mathbf{e} + \xi, \\ & \quad \mathbf{X}^T \mathbf{w} - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \quad \xi, \xi^* \geq \mathbf{0}, \\ & \quad \varepsilon \geq 0, \end{aligned} \quad (9)$$

where parameters λ_1 and λ_2 are aimed at achieving the tradeoff among the absolute regression deviation mean, the absolute regression deviation variance and the model complexity. It is evident that the soft-margin v -MADR subsumes the soft-margin v -SVR when λ_1 and λ_2 both equal 0. The meanings of the other variables have been introduced in previous formula.

B. ALGORITHMS FOR v -MADR

Solving Formula (9) is a key point for v -MADR in practical use. In this section, we first design a dual coordinate descent (DCD) algorithm for kernel v -MADR, and then present an average stochastic gradient descent (ASGD) algorithm for large-scale linear kernel v -MADR.

1) KERNEL v -MADR

By substituting the absolute regression deviation mean $\hat{\varphi}$ (Definition 1) and the absolute regression deviation variance $\hat{\varphi}$ (Definition 2) into Formula (9), we obtain Formula (10) as follows:

$$\begin{aligned} \min_{\mathbf{w}, \varepsilon, \xi, \xi^*} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \mathbf{w}^T \mathbf{X} \left(\frac{2\lambda_1 + \lambda_2}{n} \mathbf{I} - \frac{2\lambda_1}{n^2} \mathbf{e} \mathbf{e}^T \right) \mathbf{X}^T \mathbf{w} \\ & - \left(\frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{y}^T - \frac{4\lambda_1}{n^2} \mathbf{y}^T \mathbf{e} \mathbf{e}^T \right) \mathbf{X}^T \mathbf{w} \\ & + C \left(v\varepsilon + \frac{1}{n} (\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \right) \end{aligned} \quad (10)$$

$$\begin{aligned} \text{s. t. } & \mathbf{y} - \mathbf{X}^T \mathbf{w} \leq \varepsilon \mathbf{e} + \xi, \\ & \mathbf{X}^T \mathbf{w} - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \xi, \xi^* \geq \mathbf{0}, \\ & \varepsilon \geq \mathbf{0}, \end{aligned}$$

The \mathbf{y}^T and $\mathbf{y}^T \mathbf{e} \mathbf{e}^T \mathbf{y}$ terms in $\hat{\varphi}$ (Definition 1) and $\hat{\varphi}$ (Definition 2) are constants in an optimization problem, so we omit this term. However, Formula (10) is still intractable because of the high dimensionality of $\phi(\mathbf{x})$ and its complicated form. Inspired by [20, 33], we give the following theorem to state the optimal solution \mathbf{w} for Formula (10).

Theorem 1. The optimal solution \mathbf{w} for Formula (10) can be represented by the following form:

$$\mathbf{w} = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i) = \mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*), \quad (11)$$

where $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_n]^T$ and $\boldsymbol{\alpha}^* = [\alpha_1^*, \alpha_2^*, \dots, \alpha_n^*]^T$ are the parameters of v -MADR.

Proof. Suppose that \mathbf{w} can be decomposed into the span of $\phi(\mathbf{x}_i)$ and an orthogonal vector, that is,

$$\mathbf{w} = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i) + \mathbf{z} = \mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{z},$$

where \mathbf{z} satisfies $(\phi(\mathbf{x}_j)^T \cdot \mathbf{z}) = 0$ for all j , that is, $\mathbf{X}^T \mathbf{z} = \mathbf{0}$.

Then we obtain the following equation:

$$\mathbf{X}^T \mathbf{w} = \mathbf{X}^T (\mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{z}) = \mathbf{X}^T \mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*). \quad (12)$$

According to Equation (12), the second and the third terms and the constraints of Formula (10) are independent of \mathbf{z} . Besides, the last term of Formula (10) can also be considered as being independent of \mathbf{z} . To simplify the first term of Formula (10), and consider $\mathbf{X}^T \mathbf{z} = \mathbf{0}$, we get

$$\begin{aligned} \mathbf{w}^T \mathbf{w} &= (\mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{z})^T (\mathbf{X}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{z}) \\ &= (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{z}^T \mathbf{z} \\ &\geq (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*), \end{aligned}$$

where the equal relationship in the above “ \geq ” is valid if and only if $\mathbf{z} = \mathbf{0}$. Thus, setting $\mathbf{z} = \mathbf{0}$ does not affect the rest of the terms and strictly reduces the first term of Formula (10). Based on all above, \mathbf{w} in Formula (10) can be represented as the form of Equation (11). **Q.E.D.**

Based on Theorem 1, we have

$$\mathbf{X}^T \mathbf{w} = \mathbf{X}^T \mathbf{X} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) = \mathbf{Q}(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*),$$

$$\mathbf{w}^T \mathbf{w} = (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) = (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)^T \mathbf{Q} (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*),$$

where $\mathbf{Q} = \mathbf{X}^T \mathbf{X}$ is the kernel matrix. Let $\boldsymbol{\alpha}' = (\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)$, thus Formula (9) leads to

$$\begin{aligned} \min_{\boldsymbol{\alpha}', \varepsilon, \xi, \xi^*} & \frac{1}{2} (\boldsymbol{\alpha}')^T \mathbf{G} (\boldsymbol{\alpha}') + \mathbf{H}^T \boldsymbol{\alpha}' + C \left(v\varepsilon + \frac{1}{n} (\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \right) \\ \text{s. t. } & \mathbf{y} - \mathbf{Q} \boldsymbol{\alpha}' \leq \varepsilon \mathbf{e} + \xi, \end{aligned} \quad (13)$$

$$\begin{aligned} & \mathbf{Q} \boldsymbol{\alpha}' - \mathbf{y} \leq \varepsilon \mathbf{e} + \xi^*, \\ & \xi, \xi^* \geq \mathbf{0}, \\ & \varepsilon \geq \mathbf{0}, \end{aligned}$$

where $\mathbf{G} = \mathbf{Q} + \frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{Q}^T \mathbf{Q} - \frac{4\lambda_1}{n^2} \mathbf{Q}^T \mathbf{e} \mathbf{e}^T \mathbf{Q}$ and $\mathbf{H} = -\frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{Q}^T \mathbf{y} + \frac{4\lambda_1}{n^2} \mathbf{Q}^T \mathbf{e} \mathbf{e}^T \mathbf{y}$. By introducing the Lagrange multipliers $\boldsymbol{\eta}$, $\boldsymbol{\eta}^*$, $\boldsymbol{\beta}$, $\boldsymbol{\beta}^*$ and $\boldsymbol{\gamma}$, the Lagrange function of Formula (13) is given as follows:

$$\begin{aligned} L(\boldsymbol{\alpha}', \xi, \xi^*, \varepsilon, \boldsymbol{\beta}, \boldsymbol{\beta}^*, \boldsymbol{\eta}, \boldsymbol{\eta}^*, \boldsymbol{\gamma}) &= \frac{1}{2} (\boldsymbol{\alpha}')^T \mathbf{G} (\boldsymbol{\alpha}') + \mathbf{H}^T \boldsymbol{\alpha}' + \\ & C \left(v\varepsilon + \frac{1}{n} (\mathbf{e}^T \xi + \mathbf{e}^T \xi^*) \right) - \boldsymbol{\beta}^T (\varepsilon \mathbf{e} + \xi - \mathbf{y} + \mathbf{Q} \boldsymbol{\alpha}') - \\ & \boldsymbol{\beta}^{*T} (\varepsilon \mathbf{e} + \xi^* + \mathbf{y} - \mathbf{Q} \boldsymbol{\alpha}') - \boldsymbol{\eta}^T \xi - \boldsymbol{\eta}^{*T} \xi^* - \boldsymbol{\gamma}^T \varepsilon \mathbf{e}, \end{aligned} \quad (14)$$

where $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_n]^T$, $\boldsymbol{\beta}^* = [\beta_1^*, \beta_2^*, \dots, \beta_n^*]^T$, $\boldsymbol{\eta} = [\eta_1, \eta_2, \dots, \eta_n]^T$, $\boldsymbol{\eta}^* = [\eta_1^*, \eta_2^*, \dots, \eta_n^*]^T$, and $\boldsymbol{\gamma} = [\gamma_1, \gamma_2, \dots, \gamma_n]^T$. By setting the partial derivatives $\{\boldsymbol{\alpha}', \xi, \xi^*, \varepsilon\}$ to zero for satisfying the KKT conditions [34], we can get the following equations:

$$\frac{\partial L}{\partial \boldsymbol{\alpha}'} = \mathbf{G} \boldsymbol{\alpha}' + \mathbf{H} - \mathbf{Q}^T \boldsymbol{\beta} + \mathbf{Q}^T \boldsymbol{\beta}^* = \mathbf{0}, \quad (15)$$

$$\frac{\partial L}{\partial \xi} = \frac{C}{n} \mathbf{e} - \boldsymbol{\beta} - \boldsymbol{\eta} = \mathbf{0}, \quad (16)$$

$$\frac{\partial L}{\partial \xi^*} = \frac{C}{n} \mathbf{e} - \boldsymbol{\beta}^* - \boldsymbol{\eta}^* = \mathbf{0}, \quad (17)$$

$$\frac{\partial L}{\partial \varepsilon} = C v - \mathbf{e}^T \boldsymbol{\beta} - \mathbf{e}^T \boldsymbol{\beta}^* - \mathbf{e}^T \boldsymbol{\gamma} = 0. \quad (18)$$

By substituting Equations (15), (16), (17) and (18) into Equation (14), Equation (14) is re-written as:

$$\begin{aligned} \min_{\boldsymbol{\beta}, \boldsymbol{\beta}^*} f(\boldsymbol{\beta}, \boldsymbol{\beta}^*) &= \frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^T \mathbf{P} (\boldsymbol{\beta} - \boldsymbol{\beta}^*) + s^T (\boldsymbol{\beta} - \boldsymbol{\beta}^*) \\ \text{s. t. } & \mathbf{e}^T (\boldsymbol{\beta} + \boldsymbol{\beta}^*) \leq C v, \\ & 0 \leq \beta_i, \beta_i^* \leq \frac{C}{n}, i = 1, 2, \dots, n, \end{aligned} \quad (19)$$

where $\mathbf{P} = \mathbf{Q} \mathbf{G}^{-1} \mathbf{Q}^T$ and $s = -\mathbf{Q} \mathbf{G}^{-1} \mathbf{H} - \mathbf{y}$, \mathbf{G}^{-1} stands for the inverse matrix of \mathbf{G} .

According to Chang and Lin, the inequality $\mathbf{e}^T (\boldsymbol{\alpha} + \boldsymbol{\alpha}^*) \leq C v$ in v -SVR can be replaced by the equality form of $\mathbf{e}^T (\boldsymbol{\alpha} + \boldsymbol{\alpha}^*) = C v$ with the constraint $0 < v \leq 1$, and there always exists the optimal solution [11]. Based on this conclusion, we can attain the equation for the following form:

$$\mathbf{e}^T (\boldsymbol{\beta} + \boldsymbol{\beta}^*) = C v. \quad (20)$$

We thus substitute the equation $\boldsymbol{\beta}^* = C v \mathbf{e} - \boldsymbol{\beta}$ into Formula (19), and Formula (19) can be obtained as follows:

$$\begin{aligned} \min_{\beta} f(\beta) &= \frac{1}{2} (2\beta - Cve)^T P (2\beta - Cve) + s^T (2\beta - Cve) \\ \text{s. t. } 0 &\leq \beta_i \leq \frac{C}{n}, i = 1, 2, \dots, n. \end{aligned} \quad (21)$$

As one can see from Formula (21), by substituting the equation $\beta^* = Cve - \beta$ into Formula (19), the number of computational parameters of the v -MADR has been halved.

Due to the simple box constraint and the convex quadratic objective function, there exist many methods to solve the optimization problem [35-38]. To solve Formula (21), we use the DCD algorithm [39], which continuously selects one of the variables for minimization and keeps others as constants, thus a closed-form solution can be achieved at each iteration. In our situation, we minimize the variation of $f(\beta)$ by adjusting the value of $\beta_i \in \beta$ with a step size of t while keeping other $\beta_{k \neq i}$ as constants, then we need to solve the following sub-problem:

$$\begin{aligned} \min_t f(\beta + t\mathbf{d}_i) \\ \text{s. t. } 0 &\leq \beta_i + t \leq \frac{C}{n}, i = 1, 2, \dots, n, \end{aligned}$$

where \mathbf{d}_i denotes the vector with 1 in the i -th element and 0's elsewhere. Thus, we have

$$f(\beta + t\mathbf{d}_i) = f(\beta) + [\nabla f(\beta)]_i t + 2p_{ii} t^2, \quad (22)$$

where p_{ii} is the diagonal entry of P . Then we calculate the gradient $\nabla f(\beta)_i$ in Equation (22) as follows:

$$[\nabla f(\beta)]_i = 2\mathbf{d}_i^T P (2\beta - Cve) + 2s^T \mathbf{d}_i.$$

As $f(\beta)$ is independent of t , it can be omitted from Equation (22). Hence $f(\beta + t\mathbf{d}_i)$ can be transformed into a simple quadratic function. If we denote β_i^{iter} as the value of β_i at the $iter$ -th iteration, $\beta_i^{iter+1} = \beta_i^{iter} + t$ is the value at the $(iter + 1)$ -th iteration. To solve Equation (22), we can have the minimization of t which satisfies Equation (22) for the following form:

$$t = -\frac{[\nabla f(\beta)]_i}{4p_{ii}}.$$

Thus, the value of β_i^{iter+1} is obtained as

$$\beta_i^{iter+1} = \beta_i^{iter} - \frac{[\nabla f(\beta)]_i}{4p_{ii}}.$$

Furthermore, considering the box constraint $0 \leq \beta_i \leq \frac{C}{n}$, we have the minimization for β_i^{iter+1} as follows:

$$\beta_i^{iter+1} \leftarrow \min(\max(\beta_i^{iter} - \frac{[\nabla f(\beta^{iter})]_i}{4p_{ii}}, 0), C/n).$$

After β converges, we can obtain α' according to Equation (15) and Equation (20) as follows:

$$\alpha' = G^{-1}(Q^T(\beta - \beta^*) - H) = G^{-1}(Q^T(2\beta - Cve) - H).$$

Thus, the final function is

$$f(x) = \sum_{i=1}^n \alpha'_i \kappa(x_i, x),$$

where $\alpha'_i = (\alpha_i - \alpha_i^*)$.

Algorithm 1 summarizes the procedure of v -MADR with the kernel functions. The initial value of β is $Cve/2$, which simplifies the calculation procedure of v -MADR and satisfies Equation (20). Parameter v is controllable and its range is $(0,1]$.

Algorithm 1 Dual coordinate descent solver for kernel v -MADR.

Input: Dataset X , λ_1 , λ_2 , C , v ;

Output: α' ;

Initialization: $\beta = \frac{Cve}{2}$, $\alpha' = \frac{4\lambda_1 + 2\lambda_2}{n} G^{-1} Q^T y -$

$\frac{4\lambda_1}{n^2} G^{-1} Q^T e e^T y$, $A = G^{-1} Q^T$, $p_{ii} = \mathbf{d}_i^T Q G^{-1} Q^T \mathbf{d}_i$;

1: **for** $iter = 1, 2, \dots, maxIter$ **do**

2: Randomly disturb β and then get the random index;

3: **for** $i = 1, 2, \dots, n$ **do**

4: $[\nabla f(\beta)]_i \leftarrow 2(\mathbf{d}_i^T Q \alpha' - y_i)$;

5: $\beta_i^{tmp} \leftarrow \beta_i$;

6: $\beta_i \leftarrow \min(\max(\beta_i - \frac{[\nabla f(\beta)]_i}{4p_{ii}}, 0), C/n)$;

7: $\alpha' \leftarrow \alpha' + 2(\beta_i - \beta_i^{tmp}) A \mathbf{d}_i$;

8: **end for**

9: **if** β converges **then**

10: **break**;

11: **end if**

12: **end for**

We now analyze the computational complexity of Algorithm 1 as follows:

The parameters initialization is shown in Table 1, where n represents the number of the examples and m represents the number of features.

TABLE 1
TIME COMPLEXITY OF THE FORMULA

Formula being calculated	Time complexity of the formula
$Q = X^T X$	$n * m * n$
$Q_e = Q^T e$	n^2
$G = Q + \frac{4\lambda_1 + 2\lambda_2}{n} Q^T Q - \frac{4\lambda_1}{n^2} Q^T e e^T Q$	$1 + n^3 + n^2$
$invG = G^{-1}$	n^3
$A = G^{-1} Q$	n^3
$sumY = e^T y$	n
$H = -(\frac{4\lambda_1 + 2\lambda_2}{n} Q^T y - \frac{4\lambda_1}{n^2} Q^T e e^T y)$	n^2
$\alpha' = \frac{4\lambda_1 + 2\lambda_2}{n} G^{-1} Q^T y - \frac{4\lambda_1}{n^2} G^{-1} Q^T e e^T y$ $= -G^{-1} H$	n^2
$P = Q G^{-1} Q^T = Q A$	n^3

The time complexity for the dual coordinate descent (DCD) algorithm is $maxIter * n * n$, where $maxIter$ is 1000.

We can infer the time complexity of the DCD algorithm is the sum of the above time complexity. In summary, the time

complexity of the DCD algorithm is $O(n^3)$ and it has the space complexity of $O(n^2)$.

2) LARGE-SCALE LINEAR KERNEL v -MADR

In regression analysis, processing larger datasets may increase the time complexity. Although the DCD algorithm could solve kernel v -MADR efficiently for small sample problems, it is not the best strategy for larger problems. Considering computational time cost, we adopt an averaged stochastic gradient descent (ASGD) algorithm [40] to linear kernel v -MADR to improve the scalability of v -MADR, and ASGD solves the optimization problem by computing a noisy unbiased estimate of the gradient, and it randomly samples a subset of the training instances rather than all data.

We reformulate Formula (10) into a linear kernel v -MADR as follows:

$$\begin{aligned} \min_{\mathbf{w}} g(\mathbf{w}) = & \frac{1}{2} \mathbf{w}^T \left[\mathbf{I} + \frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X}\mathbf{X}^T - \frac{4\lambda_1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{X}^T \right] \mathbf{w} + \\ & \left[-\frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X}\mathbf{y} + \frac{4\lambda_1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{y} \right]^T \mathbf{w} + \frac{C}{n} \left(\sum_{i=1}^n \max(0, y_i - \right. \\ & \left. \mathbf{w}^T \mathbf{x}_i - \varepsilon) + \sum_{i=1}^n \max(0, \mathbf{w}^T \mathbf{x}_i - y_i - \varepsilon) \right), \end{aligned} \quad (23)$$

where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ and $\mathbf{y} = [y_1, y_2, \dots, y_n]^T$. The term $Cv\varepsilon$ in Formula (10) is constant in an optimization problem, so we omit this term.

For large-scale problems, it is expensive to compute the gradient of Formula (23) because we need all the training samples for computation. Stochastic gradient descent (SGD) [41, 42] works by computing a noisy unbiased estimation of the gradient via sampling a subset of the training samples. When the objective is convex, the SGD is expected to converge to the global optimal solution. In recent years, SGD has been successfully used in various machine learning problems with powerful computation efficiency [43-46].

In order to obtain an unbiased estimation of the gradient $\nabla g(\mathbf{w})$, we first present the following theorem which can be proved by computing $\nabla g(\mathbf{w})$.

Theorem 2. If two samples (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) are sampled from the training data set randomly, then

$$\begin{aligned} \nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j) = & (4\lambda_1 + 2\lambda_2) \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} - 4\lambda_1 e_i \mathbf{x}_i e_j \mathbf{x}_j^T \mathbf{w} + \\ & \mathbf{w} - (4\lambda_1 + 2\lambda_2) y_i \mathbf{x}_i + 4\lambda_1 e_i \mathbf{x}_i e_j y_j - \\ & C \begin{cases} \mathbf{x}_i & i \in I_1, \\ -\mathbf{x}_i & i \in I_2, \\ \mathbf{0} & \text{otherwise,} \end{cases} \end{aligned} \quad (24)$$

is an unbiased estimation of $\nabla g(\mathbf{w})$. Here $I_1 = \{i | y_i - \mathbf{w}^T \mathbf{x}_i \geq \varepsilon\}$, $I_2 = \{i | \mathbf{w}^T \mathbf{x}_i - y_i \geq \varepsilon\}$.

Proof. Note that the gradient of $g(\mathbf{w})$ is

$$\nabla g(\mathbf{w}) = \mathbf{G}\mathbf{w} + \mathbf{H} - \frac{C}{n} \begin{cases} \sum_{i=1}^n \mathbf{x}_i & i \in I_1, \\ \sum_{i=1}^n -\mathbf{x}_i & i \in I_2, \\ \mathbf{0} & \text{otherwise,} \end{cases}$$

where $\mathbf{G} = \mathbf{I} + \frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X}\mathbf{X}^T - \frac{4\lambda_1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{X}^T$ and $\mathbf{H} = -\frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X}\mathbf{y} + \frac{4\lambda_1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{y}$. Further note that

$$\begin{aligned} E_{\mathbf{x}_i} [\mathbf{x}_i \mathbf{x}_i^T] &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{n} \mathbf{X}\mathbf{X}^T, \\ E_{\mathbf{x}_i} [y_i \mathbf{x}_i] &= \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i = \frac{1}{n} \mathbf{X}\mathbf{y}, \\ E_{\mathbf{x}_i} [e_i \mathbf{x}_i] &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} \mathbf{X}\mathbf{e}, \\ E_{\mathbf{x}_i} [e_i y_i] &= \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \mathbf{y}^T \mathbf{e}. \end{aligned} \quad (25)$$

According to the linearity of expectation, the independence between \mathbf{x}_i and \mathbf{x}_j , and with the set of equations (25), we have

$$\begin{aligned} E_{\mathbf{x}_i \mathbf{x}_j} [\nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j)] &= (4\lambda_1 + 2\lambda_2) E_{\mathbf{x}_i} [\mathbf{x}_i \mathbf{x}_i^T] \mathbf{w} - 4\lambda_1 E_{\mathbf{x}_i} [e_i \mathbf{x}_i] E_{\mathbf{x}_j} [e_j \mathbf{x}_j] \\ & \quad + \mathbf{w} - (4\lambda_1 + 2\lambda_2) E_{\mathbf{x}_i} [y_i \mathbf{x}_i] + 4\lambda_1 E_{\mathbf{x}_i} [e_i \mathbf{x}_i] E_{\mathbf{x}_j} [e_j y_j] \\ & \quad - C \begin{cases} E_{\mathbf{x}_i} [\mathbf{x}_i | i \in I_1], \\ E_{\mathbf{x}_i} [-\mathbf{x}_i | i \in I_2], \\ \mathbf{0} \text{ otherwise,} \end{cases} \\ &= (4\lambda_1 + 2\lambda_2) \frac{1}{n} \mathbf{X}\mathbf{X}^T \mathbf{w} - 4\lambda_1 \frac{1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{X}^T \mathbf{w} + \mathbf{w} \\ & \quad - (4\lambda_1 + 2\lambda_2) \frac{1}{n} \mathbf{X}\mathbf{y} + 4\lambda_1 \frac{1}{n^2} \mathbf{X}\mathbf{e}\mathbf{e}^T \mathbf{y} \\ & \quad - C \frac{1}{n} \begin{cases} \sum_{i=1}^n \mathbf{x}_i, i \in I_1, \\ \sum_{i=1}^n -\mathbf{x}_i, i \in I_2, \\ \mathbf{0} \text{ otherwise,} \end{cases} \\ &= \mathbf{G}\mathbf{w} + \mathbf{H} - \frac{C}{n} \begin{cases} \sum_{i=1}^n \mathbf{x}_i & i \in I_1, \\ \sum_{i=1}^n -\mathbf{x}_i & i \in I_2, \\ \mathbf{0} & \text{otherwise,} \end{cases} \\ &= \nabla g(\mathbf{w}). \end{aligned}$$

It is shown that $\nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j)$ is a noisy unbiased gradient of $g(\mathbf{w})$. **Q.E.D.**

Based on Theorem 2, the stochastic gradient can be updated as follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \varphi_t \nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j), \quad (26)$$

where φ_t is the learning rate at the t -th iteration.

Since the ASGD algorithm is more robust than the SGD algorithm [47], we actually adopt the ASGD algorithm to solve the optimization problem in Formula (23). At each iteration, in addition to updating the normal stochastic gradient in Equation (26), we also compute

$$\bar{\mathbf{w}}_t = \frac{1}{t - t_0} \sum_{i=t_0+1}^t \mathbf{w}_i,$$

where t_0 decides when to take the averaging operation. This average can also be calculated in a recursive formula as follows:

$$\bar{\mathbf{w}}_{t+1} = \bar{\mathbf{w}}_t + \delta_t(\mathbf{w}_{t+1} - \bar{\mathbf{w}}_t),$$

where $\delta_t = 1/\max\{1, t - t_0\}$.

Algorithm 2 summarizes the procedure of large-scale linear kernel ν -MADR.

Algorithm 2 Averaged stochastic gradient descent solver for linear kernel ν -MADR.

Input: Dataset \mathbf{X} , λ_1 , λ_2 , C , ε ;

Output: $\bar{\mathbf{w}}$;

Initialization: $\mathbf{u} = \mathbf{0}$, $t = 1$, $T = 5$;

1: **While** $t \leq T \cdot n$ **do**

2: Randomly select the training instances $(\mathbf{x}_i, \mathbf{y}_i)$ and $(\mathbf{x}_j, \mathbf{y}_j)$;

3: Compute $\nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j)$ as in Equation (24);

4: $\mathbf{w} \leftarrow \varphi_t \nabla g(\mathbf{w}, \mathbf{x}_i, \mathbf{x}_j)$;

5: $\bar{\mathbf{w}} \leftarrow \bar{\mathbf{w}} + \delta_t(\mathbf{w} - \bar{\mathbf{w}})$;

6: end while

The time complexity of the averaged stochastic gradient descent (ASGD) algorithm is $O(T * n * m)$ and its space complexity is $O(n * m)$.

3) PROPERTIES OF ν -MADR

We study the statistical property of ν -MADR that leads to a bound on the expectation of error for ν -MADR according to the leave-one-out cross-validation estimate, which is an unbiased estimate of the probability of test error. For the sake of simplicity, we only discuss the linear case as shown Formula (10) here, in which \mathbf{w} can be represented by the following form:

$$\mathbf{w} = \sum_{i=1}^n (\alpha_i - \alpha_i^*) = \boldsymbol{\alpha} - \boldsymbol{\alpha}^*,$$

while the result is also used in kernel mapping situations ϕ . Then we can get the dual problem of Formula (10) using the same steps as in Section III.B.1, i.e.

$$\begin{aligned} \min_{\boldsymbol{\beta}} f(\boldsymbol{\beta}) &= \frac{1}{2} (2\boldsymbol{\beta} - C\mathbf{v}\mathbf{e})^T \mathbf{P} (2\boldsymbol{\beta} - C\mathbf{v}\mathbf{e}) + \mathbf{s}^T (2\boldsymbol{\beta} - C\mathbf{v}\mathbf{e}) \\ \text{s. t. } 0 &\leq \beta_i \leq \frac{C}{n}, i = 1, 2, \dots, n, \end{aligned} \quad (27)$$

where $\mathbf{P} = \mathbf{X}^T \mathbf{G}^{-1} \mathbf{X}$, $\mathbf{s} = -\mathbf{X}^T \mathbf{G}^{-1} \mathbf{H} - \mathbf{y}$, $\mathbf{G} = \frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X} \mathbf{X}^T - \frac{4\lambda_1}{n^2} \mathbf{X} \mathbf{e} \mathbf{e}^T \mathbf{X}^T + \mathbf{I}$ and $\mathbf{H} = -\frac{4\lambda_1 + 2\lambda_2}{n} \mathbf{X} \mathbf{y} + \frac{4\lambda_1}{n^2} \mathbf{X} \mathbf{e} \mathbf{e}^T \mathbf{y}$.

Definition 3. Regression error is defined as follows:

$$\pi(\mathbf{x}, \mathbf{y}) = |\mathbf{y} - f(\mathbf{x})|.$$

We give the following theorem to state the expectation of the probability of test error.

Theorem 3. Let $\boldsymbol{\beta}$ be the optimal solution of (27), and $E[R(\boldsymbol{\beta})]$ be the expectation of the probability of test error, then we have

$$E[R(\boldsymbol{\theta})] \leq \frac{E[\varepsilon I_1 + 2p \sum_{i \in I_2} \beta_i + \sum_{i \in I_3} (\varepsilon + \xi_i)]}{n} \quad (28)$$

where $I_1 \equiv \{i | (\beta_i = 0) \cap (\beta_i^* = 0)\}$, $I_2 \equiv \{i | ((0 < \beta_i < C/n) \cap (\beta_i^* = 0)) \cup ((0 < \beta_i^* < C/n) \cap (\beta_i = 0))\}$, $I_3 \equiv \{i | ((\beta_i = C/n) \cap (\beta_i^* = 0)) \cup ((\beta_i^* = C/n) \cap (\beta_i = 0))\}$, $\xi_i = \max\{\xi_i, \xi_i^*\}$, $p = \max\{p_{ii}, i = 0, 1, \dots, n\}$, $\beta_i^* = C\nu - \beta_i$ and p_{ii} is the diagonal entry of \mathbf{P} .

Proof. Suppose

$$\boldsymbol{\beta}^* = \underset{0 \leq \beta_i \leq \frac{C}{n}}{\operatorname{argmin}} f(\boldsymbol{\beta}),$$

$$\boldsymbol{\beta}^i = \underset{0 \leq \beta_i \leq \frac{C}{n}, \beta_i = 0}{\operatorname{argmin}} f(\boldsymbol{\beta}), i = 1, 2, \dots, n$$

and the corresponding solution for the linear kernel ν -MADR are \mathbf{w}' and \mathbf{w}^i , respectively.

According to [48],

$$E[R(\boldsymbol{\theta})] \leq \frac{E[L((\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n))]}{n} \quad (29)$$

where $L((\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n))$ is the number of errors in the leave-one-out procedure.

In the process of solving Formula (27) using the Lagrange multipliers, every sample must meet the following KKT conditions:

$$\beta_i (\varepsilon + \xi_i - \mathbf{y}_i + \mathbf{x}_i^T \boldsymbol{\alpha}') = 0,$$

$$\beta_i^* (\varepsilon + \xi_i^* + \mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\alpha}') = 0,$$

$$\left(\frac{C}{n} - \beta_i\right) \xi_i = 0,$$

$$\left(\frac{C}{n} - \beta_i^*\right) \xi_i^* = 0,$$

$$\begin{aligned} \xi_i, \xi_i^* &\geq 0, i = 1, 2, \dots, n, \\ \varepsilon &\geq 0. \end{aligned}$$

According to the KKT conditions, we have that if and only if $\varepsilon + \xi_i - \mathbf{y}_i + \mathbf{x}_i^T \boldsymbol{\alpha}' = 0$, β_i can take a non-zero value, and if and only if $\varepsilon + \xi_i^* + \mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\alpha}' = 0$, β_i^* can take a non-zero value. In other words, if the sample $(\mathbf{x}_i, \mathbf{y}_i)$ is not in the ε -tube in the leave-one-out procedure, β_i and β_i^* can take a non-zero value. In addition, $\varepsilon + \xi_i - \mathbf{y}_i + \mathbf{x}_i^T \boldsymbol{\alpha}' = 0$ and $\varepsilon + \xi_i^* + \mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\alpha}' = 0$ cannot be established at the same time, so we get that at least one of β_i and β_i^* is zero. The specific breakdown is as follows:

i) If the sample $(\mathbf{x}_i, \mathbf{y}_i)$ is in the ε -tube in the leave-one-out procedure, then $\varepsilon + \xi_i - \mathbf{y}_i + \mathbf{x}_i^T \boldsymbol{\alpha}' \neq 0$ and $\varepsilon + \xi_i^* + \mathbf{y}_i - \mathbf{x}_i^T \boldsymbol{\alpha}' \neq 0$, so we have $\beta_i = 0$ and $\beta_i^* = 0$;

ii) If the sample $(\mathbf{x}_i, \mathbf{y}_i)$ is out of the ε -tube in the leave-one-out procedure, we have the following two situations:

- a) if the sample is above the ε -tube, then $\xi_i \neq 0$ and $\varepsilon + \xi_i^* + y_i - \mathbf{x}_i^T \boldsymbol{\alpha}' \neq 0$. So we have $\beta_i = C/n$ and $\beta_i^* = 0$;
- b) if the sample is under the ε -tube, then $\xi_i^* \neq 0$ and $\varepsilon + \xi_i - y_i + \mathbf{x}_i^T \boldsymbol{\alpha}' \neq 0$. So we have $\beta_i^* = C/n$ and $\beta_i = 0$;
- iii) If the sample (\mathbf{x}_i, y_i) is on the gap of the ε -tube in the leave-one-out procedure, we have the following two situations:
- a) if the sample is on the upper gap of the ε -tube, then $\xi_i = 0$, and we have $0 < \beta_i \leq C/n$ and $\beta_i^* = 0$;
- b) if the sample is on the lower gap of the ε -tube, then $\xi_i^* = 0$, and we have $0 < \beta_i^* \leq C/n$ and $\beta_i = 0$.

Based on the discussion above, we consider the following three cases to calculate the test error:

i) If both $\beta_i = 0$ and $\beta_i^* = 0$, we have that the sample (\mathbf{x}_i, y_i) is in the ε -tube in the leave-one-out procedure, and $\pi(\mathbf{x}_i, y_i) \leq \varepsilon$.

ii) If $(0 < \beta_i < C/n) \cap (\beta_i^* = 0)$ or $(0 < \beta_i^* < C/n) \cap (\beta_i = 0)$, we have that

$$f(\boldsymbol{\beta}^i) - \min_{\mathbf{t}} f(\boldsymbol{\beta}^i + \mathbf{t} \mathbf{d}_i) \leq f(\boldsymbol{\beta}^i) - f(\boldsymbol{\beta}'), \quad (30)$$

$$f(\boldsymbol{\beta}^i) - f(\boldsymbol{\beta}') \leq f(\boldsymbol{\beta}' - \beta_i' \mathbf{d}_i) - f(\boldsymbol{\beta}'), \quad (31)$$

where \mathbf{d}_i denotes the vector with 1 in the i -th element and 0's elsewhere. We can discover that the left-hand side of formula (30) is equal to $[\nabla f(\boldsymbol{\beta}^i)]_i^2 / (8p_{ii}) = (\mathbf{x}_i^T \mathbf{w}^i - y_i)^2 / (2p_{ii})$ and the right-hand side of formula (31) is equal to $2p_{ii} \beta_i'^2$. So by combining formula (30) and (31), we have $\pi(\mathbf{x}_i, y_i)^2 / (2p_{ii}) = (\mathbf{x}_i^T \mathbf{w}^i - y_i)^2 / (2p_{ii}) \leq 2p_{ii} \beta_i'^2$. Further, we can obtain $\pi(\mathbf{x}_i, y_i) \leq 2p_{ii} \beta_i'$.

iii) If $(\beta_i = C/n) \cap (\beta_i^* = 0)$ or $(\beta_i^* = C/n) \cap (\beta_i = 0)$, we have that the sample (\mathbf{x}_i, y_i) is not in the ε -tube in the leave-one-out procedure. So we can get $\pi(\mathbf{x}_i, y_i) = \varepsilon + \bar{\xi}_i'$, where $\bar{\xi}_i' = \max\{\xi_i', \xi_i^{*'}\}$.

So we have

$$L((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)) \leq \varepsilon |I_1| + 2p \sum_{i \in I_2} \beta_i' + \sum_{i \in I_3} (\varepsilon + \bar{\xi}_i'),$$

where $I_1 \equiv \{i | (\beta_i' = 0) \cap (\beta_i^{*'} = 0)\}$, $I_2 \equiv \{i | ((0 < \beta_i' < C/n) \cap (\beta_i^{*'} = 0)) \cup ((0 < \beta_i^{*'} < C/n) \cap (\beta_i' = 0))\}$, $I_3 \equiv \{i | ((\beta_i' = C/n) \cap (\beta_i^{*'} = 0)) \cup ((\beta_i^{*'} = C/n) \cap (\beta_i' = 0))\}$, $\bar{\xi}_i' = \max\{\xi_i', \xi_i^{*'}\}$, $p = \max\{p_{ii}, i = 0, 1, \dots, n\}$ and $\beta_i^{*'} = Cv - \beta_i'$. Take expectation on both side and with formula (29), we reach the conclusion that formula (28) holds. **Q.E.D.**

IV. EXPERIMENTAL RESULTS

Since SVR-based algorithms are now widely used for regression problems and demonstrate better generalization ability [8-10] than many existing algorithms, such as least square regression [4], Neural Networks (NN) regression [5], logistic regression [6], and ridge regression [7], we will not

repeat these comparisons. In this section, we empirically evaluate the performance of our v -MADR compared with other SVR-based algorithms, including v -SVR, LS-SVR, ε -TSVR, and linear ε -SVR on several datasets, including two artificial datasets, eight medium-scale datasets, and six large-scale datasets. All algorithms are implemented with MATLAB R2014a on a PC with a 2.00GHz CPU and 32 GB memory. v -SVR is solved by LIBSVM [49]; ε -SVR is solved by LIBLINEAR [50]; LS-SVR is solved by LSSVMLab [51]; and ε -TSVR is solved by the SOR technique [52, 53]. RBF kernel $\kappa(\mathbf{x}_i^T, \mathbf{x}_j^T) = \exp(-\|\mathbf{x}_i^T - \mathbf{x}_j^T\|^2 / \sigma^2)$ and polynomial kernel $\kappa(\mathbf{x}_i^T, \mathbf{x}_j^T) = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^d$ are employed for nonlinear regression. The values of the parameters are obtained by means of a grid-search method [54]. For brevity, we set $c_1 = c_2, c_3 = c_4$ and $\varepsilon_1 = \varepsilon_2$ for ε -TSVR and $\lambda_1 = \lambda_2$ for our nonlinear v -MADR. The parameter v in v -MADR is selected from the set $\{2^{-9}, 2^{-8}, \dots, 2^0\}$, and the remaining parameters in the five methods and the parameters in the Gaussian kernel are selected from the set $\{2^{-9}, 2^{-8}, \dots, 2^9\}$ by 10-fold cross-validation. Specifically, the parameter d in polynomial kernel is selected from $\{2, 3, 4, 5, 6\}$.

In order to evaluate the performance of the proposed algorithm, the performance metrics are specified before presenting the experimental results. Without loss of generality, let n be the number of training samples and m be the number of testing samples, denote \hat{y}_i as the prediction value of y_i , and $\bar{y} = (\sum_{i=1}^m y_i) / m$ as the average value of y_1, y_2, \dots, y_m . Then the details of the metrics used for assessing the performance of all regression algorithms are stated in Table 2. To demonstrate the overall performance of a method, a performance metric referred to average rank of each method is defined as

$$\text{average rank}(R) = \frac{1}{s} \sum_{i=1}^s \text{rank}(R)_i,$$

where $R \in \{v\text{-SVR, LS-SVR, } \varepsilon\text{-TSVR, LIBLINE-AR, } v\text{-MADR}\}$ is the regression method, s is the number of datasets, and $\text{rank}(R)_i$ means the performance rank of method R on the i -th dataset among all regression methods.

In our experiments, we test the performance of the above methods on two artificial datasets, eight medium-scale datasets and six large-scale data sets. The basic information of these datasets is given in Table 3. All real-world datasets are taken from UCI (<http://archive.ics.uci.edu/ml>) and StatLib (<http://lib.stat.cmu.edu/>), and more detailed information can be accessed from those websites. Before regression analysis, all of these real datasets are normalized to zero mean and unit deviation. For medium-scale datasets, RBF kernel and polynomial kernel are used, and for large-scale datasets, only the linear kernel v -MADR is used considering the computational complexity. Each experiment is repeated for 30 trials with 10-fold cross validation and the mean evaluation of R^2 , NMSE, MAPE and their standard deviations were recorded. Particularly, the two datasets "Diabetes" and

“Motorcycle” have smaller numbers of samples and features, so we use the leave-one-out cross validation instead.

TABLE 2
PERFORMANCE METRICS

Metrics	Definition
SSE	$SSE = \sum_{i=1}^m (y_i - \hat{y}_i)^2$
SST	$SST = \sum_{i=1}^m (y_i - \bar{y})^2$
SSR	$SSR = \sum_{i=1}^m (\hat{y}_i - \bar{y})^2$
NMSE	$NMSE = SSE/SST$
R^2	$R^2 = \frac{\left(\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - E[\hat{y}_i])(y_i - E[y_i])\right)^2}{\sigma_y^2 \sigma_{\hat{y}}^2}$
MAPE	$MAPE = \frac{1}{m} \sum_{i=1}^m \left \frac{y_i - \hat{y}_i}{y_i} \right $

TABLE 3
THE REAL-WORLD DATASETS USED FOR EXPERIMENTS

Scale	Dataset	Samples	Features	Dataset	Samples	Features
medium	Diabetes	43	2	Motorcycle	133	1
	Autoprice	159	15	Servo	167	4
	Wisconsin	194	32	MachineCPU	209	31
	AutoMpg	398	7	WDBC	569	30
large	ConcreteCS	1030	8	Abalone	4177	8
	CPUsmall	8192	12	Bike	10886	9
	Driftdataset	13910	128	Cadate	20640	8

A. ARTIFICIAL DATASETS

In order to compare our v -MADR with v -SVR, LS-SVR, and ε -TSVR, we choose two artificial datasets with different distributions. Firstly, we consider the function: $y = x^{\frac{2}{3}}$. In order to fully assess the performance of the methods, the training samples are added with Gaussian noises with zero means and 0.5 standard deviation, that is, we have the following training samples (x_i, y_i) :

$$y_i = x_i^{\frac{2}{3}} + \xi_i, x_i \sim U[-2, 2], \xi_i \sim N(0, 0.5^2), \quad (32)$$

where $U[a, b]$ represents the uniformly random variable in $[a, b]$ and $N(\mu, \sigma^2)$ represents the Gaussian random variable with means μ and standard deviation σ , respectively. To avoid biased comparisons, ten independent groups of noisy samples are randomly generated, including 200 training samples and 400 none noise test samples.

The estimated functions obtained by these four methods are shown in Figure 1. Obviously, all four methods have obtained good fitted values, but our v -MADR has the best approximation compared to the rest of the methods. Table 4

shows the corresponding performance metrics and training time. Compared with the other methods, our v -MADR has the highest R^2 , lowest NMSE and MAPE, which indicates that our v -MADR achieves good fitting performance and a good presentation of the statistical information in the training dataset. In addition, the CPU time of our v -MADR is not much different from other methods, and equivalent to v -SVR.

The second artificial example is the regression estimation on the *Sinc* function: $y = \sin(x)/x$. The training samples are added with Gaussian noise with zero means and 0.5 standard deviation. Therefore, we have the following training samples (x_i, y_i) :

$$y_i = \frac{\sin(x_i)}{x_i} + \left(0.5 - \frac{|x_i|}{8\pi}\right) \xi_i, \quad (33)$$

$$x_i \sim U[-4\pi, 4\pi], \xi_i \sim N(0, 0.5^2).$$

The dataset consists of 200 training samples and 400 test samples. Figure 2 illustrates the estimated functions obtained by these four methods and Table 4 shows the corresponding performance. These results also demonstrate the superiority of our method. At the bottom of Table 4, we list the average ranks of all four methods on the artificial datasets for different

performance metrics. It can be seen that our v -MADR is superior to other three methods on R^2 and NMSE, and is comparable to LS-SVR and ε -TSVR in terms of MAPE.

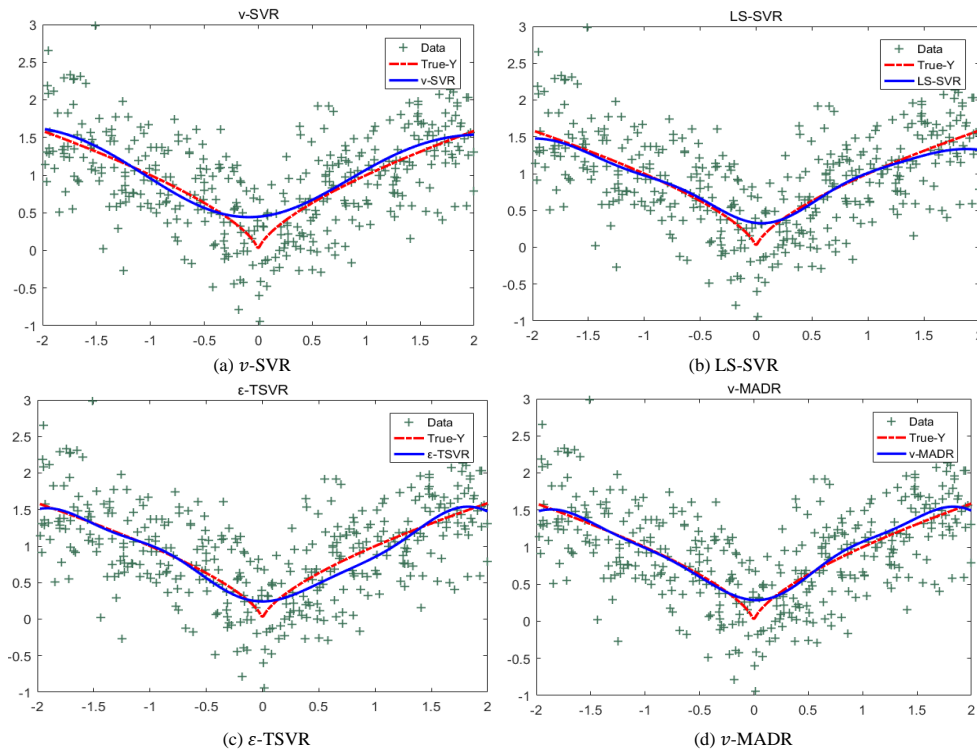


FIGURE 1. The predictions of v -SVR, LS-SVR, ε -TSVR and our v -MADR on function $y = x^2$.

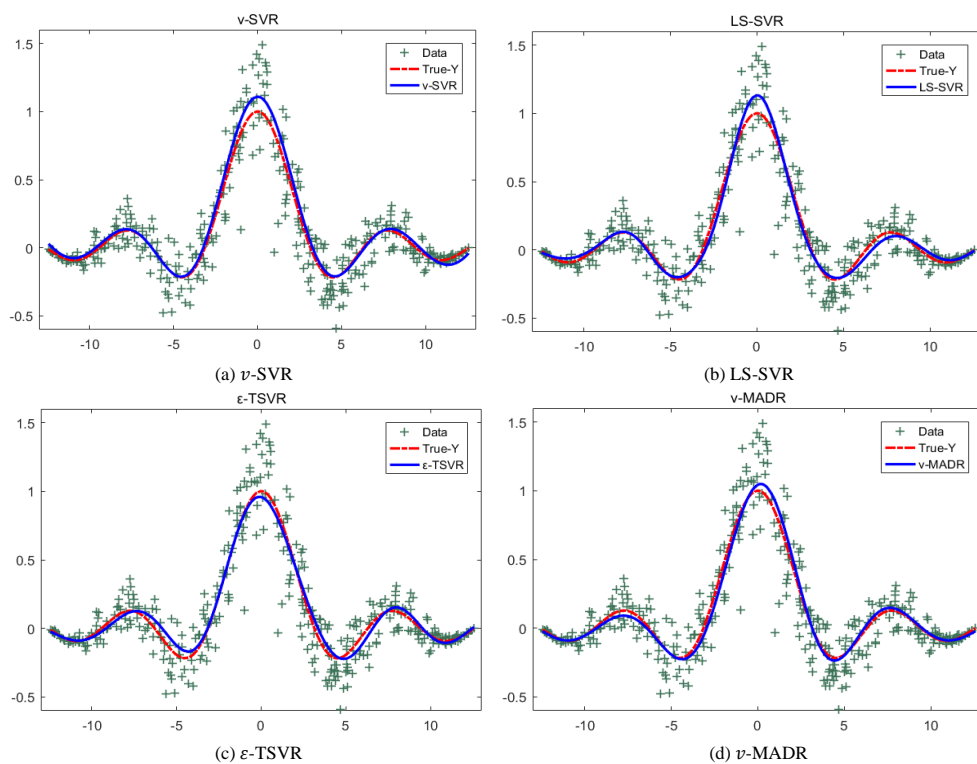


FIGURE 2. The predictions of v -SVR, LS-SVR, ε -TSVR and our v -MADR on the *sinc* function.

TABLE 4
THE RESULT COMPARISONS OF v -SVR, LS-SVR, ϵ -TSVR AND OUR v -MADR ON ARTIFICIAL DATASETS.

Dataset	regressor	R^2 (rank)	NMSE (rank)	MAPE (rank)	CPU(sec)
$x^{\frac{2}{3}}$	v -SVR	0.9319(4)	0.0856(4)	0.2479(4)	0.0333
	LS-SVR	0.9446(3)	0.0698(3)	0.2152(2)	0.0188
	ϵ -TSVR	0.9500(2)	0.0593(2)	0.2091(1)	0.0196
	v-MADR	0.9584(1)	0.0529(1)	0.2165(3)	0.0471
$Sinc(x)$	v -SVR	0.9889(2)	0.0183(2)	1.1792(4)	0.0837
	LS-SVR	0.9844(3)	0.0190(3)	0.8389(2)	0.0172
	ϵ -TSVR	0.9823(4)	0.0200(4)	0.9118(3)	0.0202
	v-MADR	0.9940(1)	0.0083(1)	0.7333(1)	0.0474
average rank	v -SVR	3	3	4	-
	LS-SVR	3	3	2	-
	ϵ -TSVR	3	3	2	-
	v-MADR	1	1	2	-

B. MEDIUM-SCALE DATASETS

Table 5 and Table 6 list the experimental results on the eight medium-scale datasets from UCI and StatLib with RBF and polynomial kernels, respectively. From the average rank at the bottom of Table 5 and Table 6, our v -MADR is superior to the other three methods. In detail, on most datasets, our v -MADR has the highest R^2 , lowest NMSE and MAPE. Although on several datasets, such as ‘‘MachineCPU’’, our v -MADR does not achieve the best experimental results compared with other methods, it is not the worst. Our v -MADR also has good performance in terms of CPU running time. The above experimental results indicate that v -MADR is an efficient and promising algorithm for regression. Table 7 and Table 8 list the optimal parameters with RBF and polynomial kernels, respectively. Figure 3(a) and Figure 3(b) show the comparisons of CPU time among our v -MADR, v -SVR, LS-SVR and ϵ -TSVR on each medium-scale dataset with RBF kernel and polynomial kernel.

For further evaluation, we investigate the absolute regression deviation mean and variance of our v -MADR with RBF kernel, v -SVR, LS-SVR and ϵ -TSVR on medium-scale datasets as shown in Figure 4. From Figure 4, our v -MADR has the smallest absolute regression deviation mean and

variance on most datasets. In addition, v -MADR also has the most compact mean and variance distribution, which demonstrates its robustness. From the above results, it is obvious that our v -MADR outperforms other three methods.

The change of parameter values may have a great effect on the results of regression analysis. For our RBF kernel v -MADR, there are mainly three trade-off parameters, i.e., λ_1 , λ_2 , C and one kernel parameter σ . Figure 5(a) and Figure 5(b) shows the influence of λ_1 on NMSE and CPU time by varying it from 2^{-9} to 2^9 while fixing λ_2 , C and σ as the optimal ones by cross validation. Figures 5(c)~5(h) show the influence of λ_2 , C and σ on NMSE and CPU time, respectively. As one can see from Figure 5(a), Figure 5(c) and Figure 5(e), the NMSE values on medium-scale datasets do not change significantly when the values of the three parameters λ_1 , λ_2 , and C are changed. Figure 5(g) shows that σ has more obvious influence on NMSE. On most datasets, as σ becomes larger, NMSE will become smaller and smaller until it converges at a fixed value. Figure 5(b), Figure 5(d), Figure 5(f) and Figure 5(h) show the influence of parameters λ_1 , λ_2 , C and σ on CPU time. Experimental results indicate that the performance of v -MADR is not sensitive to parameter changes, which further demonstrates the robustness of v -MADR.

TABLE 5
THE RESULT COMPARISONS OF v -SVR, LS-SVR, ϵ -TSVR AND v -MADR ON MEDIUM-SCALE DATASETS WITH RBF KERNEL

Dataset	regressor	R^2 (rank)	NMSE (rank)	MAPE (rank)	CPU(sec)
Diabetes	v -SVR	0.5343 \pm 0.0028(2)	0.4768 \pm 0.0024(2)	1.5971 \pm 0.0197(4)	0.009
	LS-SVR	0.5151 \pm 0.0381(4)	0.4986 \pm 0.0504(4)	1.5742 \pm 0.0477(2)	0.014
	ϵ -TSVR	0.5281 \pm 0.1182(3)	0.4805 \pm 0.1185(3)	1.5909 \pm 0.1188(3)	0.004
	v-MADR	0.5891 \pm 0.0137(1)	0.4127 \pm 0.0141(1)	1.4750 \pm 0.0493(1)	0.002
Motorcycle	v -SVR	0.7938 \pm 0.0020(3)	0.2081 \pm 0.0029(4)	1.2636 \pm 0.0352(4)	0.007
	LS-SVR	0.7975 \pm 0.0011(2)	0.2027 \pm 0.0009(3)	1.2444 \pm 0.0154(3)	0.013
	ϵ -TSVR	0.7680 \pm 0.0003(4)	0.2021 \pm 0.0003(2)	1.2441 \pm 0.0055(2)	0.012
	v-MADR	0.7984 \pm 0.0006(1)	0.2017 \pm 0.0007(1)	1.2316 \pm 0.0063(1)	0.016
Autoprice	v -SVR	0.9481 \pm 0.0680(2)	0.0530 \pm 0.0698(2)	0.4187 \pm 0.1070(1)	0.004
	LS-SVR	0.9465 \pm 0.0674(3)	0.0541 \pm 0.0682(3)	0.5039 \pm 0.1392(3)	0.049
	ϵ -TSVR	0.9338 \pm 0.0744(4)	0.0665 \pm 0.0745(4)	0.5355 \pm 0.1524(4)	0.007
	v-MADR	0.9549 \pm 0.0040(1)	0.0467 \pm 0.0035(1)	0.4889 \pm 0.0220(2)	0.027
Servo	v -SVR	0.9337 \pm 0.0565(4)	0.0686 \pm 0.0582(4)	0.2491 \pm 0.0420(2)	0.291
	LS-SVR	0.9630 \pm 0.0403(2)	0.0372 \pm 0.0406(2)	0.3374 \pm 0.1835(3)	0.019

	ϵ -TSVR	0.9507 \pm 0.0445(3)	0.0495 \pm 0.0446(3)	0.3931 \pm 0.2082(4)	0.016
	v-MADR	0.9755 \pm 0.0187(1)	0.0253 \pm 0.0320(1)	0.2473 \pm 0.1280(1)	0.045
Wisconsin	v-SVR	0.2420 \pm 0.0354(3)	0.7712 \pm 0.0371(3)	1.8255 \pm 0.1265(4)	0.008
	LS-SVR	0.2546 \pm 0.0146(2)	0.7574 \pm 0.0202(2)	1.5667 \pm 0.0462(1)	0.086
	ϵ -TSVR	0.2285 \pm 0.0130(4)	0.7752 \pm 0.0137(4)	1.7028 \pm 0.1055(3)	0.027
	v-MADR	0.2641 \pm 0.0016(1)	0.7486 \pm 0.0012(1)	1.5840 \pm 0.0125(2)	0.046
MachineCPU	v-SVR	0.9994 \pm 0.0005(1)	0.0006 \pm 0.0006(1)	0.0888 \pm 0.1165(1)	0.006
	LS-SVR	0.9978 \pm 0.0020(2)	0.0022 \pm 0.0020(2)	0.1775 \pm 0.1113(2)	0.027
	ϵ -TSVR	0.9921 \pm 0.0048(4)	0.0080 \pm 0.0048(4)	0.3580 \pm 0.1641(4)	0.034
	v-MADR	0.9942 \pm 0.0014(3)	0.0063 \pm 0.0017(3)	0.2134 \pm 0.0779(3)	0.047
AutoMpg	v-SVR	0.9196 \pm 0.0108(4)	0.0807 \pm 0.0121(4)	1.0111 \pm 0.2842(2)	0.118
	LS-SVR	0.9262 \pm 0.0072(2)	0.0741 \pm 0.0079(2)	1.0196 \pm 0.0977(3)	0.054
	ϵ -TSVR	0.9228 \pm 0.0034(3)	0.0773 \pm 0.0035(3)	1.0549 \pm 0.0421(4)	0.105
	v-MADR	0.9267 \pm 0.0017(1)	0.0736 \pm 0.0015(1)	1.0080 \pm 0.0182(1)	0.326
WDBC	v-SVR	0.9382 \pm 0.0106(3)	0.0632 \pm 0.0086(3)	0.0988 \pm 0.0041(3)	0.216
	LS-SVR	0.9520 \pm 0.0109(2)	0.0489 \pm 0.0113(2)	0.0182 \pm 0.0052(1)	0.196
	ϵ -TSVR	0.9344 \pm 0.0066(4)	0.0659 \pm 0.0047(4)	0.1726 \pm 0.0085(4)	0.615
	v-MADR	0.9710 \pm 0.0258(1)	0.0298 \pm 0.0198(1)	0.0714 \pm 0.0043(2)	0.913
average rank	v-SVR	2.7500	2.8750	2.6250	-
	LS-SVR	2.3750	2.5000	2.2500	-
	ϵ -TSVR	3.6250	3.3750	3.5000	-
	v-MADR	1.2500	1.2500	1.6250	-

TABLE 6

THE RESULT COMPARISONS OF v-SVR, LS-SVR, ϵ -TSVR AND v-MADR ON MEDIUM-SCALE DATASETS WITH POLYNOMIAL KERNEL

Dataset	regressor	R^2 (rank)	NMSE (rank)	MAPE (rank)	CPU(sec)
Diabetes	v-SVR	0.371 \pm 0.0931(4)	0.6424 \pm 0.1007(4)	1.2135 \pm 0.5799(2)	0.001
	LS-SVR	0.526 \pm 0.0096(2)	0.4901 \pm 0.0313(3)	1.4929 \pm 0.1274(3)	0.012
	ϵ -TSVR	0.525 \pm 0.0063(3)	0.4846 \pm 0.0214(2)	1.6795 \pm 0.1074(4)	0.004
	v-MADR	0.549 \pm 0.0324(1)	0.4834 \pm 0.0317(1)	1.1254 \pm 0.0465(1)	0.002
Motorcycle	v-SVR	0.115 \pm 0.0106(4)	0.8898 \pm 0.0103(4)	1.2730 \pm 0.1945(1)	0.001
	LS-SVR	0.547 \pm 0.0024(3)	0.4544 \pm 0.0055(3)	1.6929 \pm 0.0356 (4)	0.014
	ϵ -TSVR	0.548 \pm 0.0010(2)	0.4540 \pm 0.0042(2)	1.6771 \pm 0.0540 (3)	0.032
	v-MADR	0.549 \pm 0.0006(1)	0.4516 \pm 0.0028(1)	1.6456 \pm 0.0315 (2)	0.017
Autoprice	v-SVR	0.881 \pm 0.0433(4)	0.1217 \pm 0.0551(4)	0.6501 \pm 0.1275(4)	0.004
	LS-SVR	0.965 \pm 0.0109(3)	0.0349 \pm 0.0113(3)	0.4184 \pm 0.0446(2)	0.043
	ϵ -TSVR	0.973 \pm 0.0082(2)	0.0315 \pm 0.0152(2)	0.4750 \pm 0.1303(3)	0.017
	v-MADR	0.976 \pm 0.0108(1)	0.0256 \pm 0.0115(1)	0.4056 \pm 0.0772(1)	0.025
Servo	v-SVR	0.545 \pm 0.0034(4)	0.4905 \pm 0.0079(4)	0.8068 \pm 0.0122(4)	0.016
	LS-SVR	0.935 \pm 0.0385(3)	0.0645 \pm 0.0384(3)	0.4870 \pm 0.1473(3)	0.020
	ϵ -TSVR	0.942 \pm 0.0012(2)	0.0576 \pm 0.0017(1)	0.4856 \pm 0.0112(2)	0.019
	v-MADR	0.945 \pm 0.0047(1)	0.0631 \pm 0.0066(2)	0.4425 \pm 0.1647(1)	0.083
Wisconsin	v-SVR	0.203 \pm 0.0227(4)	0.8326 \pm 0.0357(4)	1.3011 \pm 0.0160(2)	0.007
	LS-SVR	0.469 \pm 0.0021(3)	0.5692 \pm 0.0092(3)	1.3614 \pm 0.0036(3)	0.095
	ϵ -TSVR	0.777 \pm 0.0020(1)	0.2439 \pm 0.0336(1)	1.3793 \pm 0.1926(4)	0.039
	v-MADR	0.585 \pm 0.0013(2)	0.4577 \pm 0.0066(2)	0.8205 \pm 0.0081(1)	0.041
MachineCPU	v-SVR	0.933 \pm 0.0175(4)	0.0682 \pm 0.0171(4)	1.3525 \pm 0.2995(4)	0.015
	LS-SVR	0.999 \pm 0.0008(3)	0.0006 \pm 0.0009(3)	0.1068 \pm 0.0696(2)	0.029
	ϵ -TSVR	0.999 \pm 0.0004(2)	0.0005 \pm 0.0005(2)	0.1159 \pm 0.0503(3)	0.021
	v-MADR	0.999 \pm 0.0002(1)	0.0004 \pm 0.0003(1)	0.0916 \pm 0.0251(1)	0.045
AutoMpg	v-SVR	0.802 \pm 0.0127(4)	0.1997 \pm 0.0136(4)	1.2530 \pm 0.2663(2)	0.080
	LS-SVR	0.895 \pm 0.0307(2)	0.1047 \pm 0.0309(2)	1.2613 \pm 0.2146(3)	0.047
	ϵ -TSVR	0.892 \pm 0.0396(3)	0.1079 \pm 0.0356(3)	1.3668 \pm 0.0156(4)	0.143
	v-MADR	0.922 \pm 0.0131(1)	0.0815 \pm 0.0154(1)	1.1219 \pm 0.2950(1)	0.298

WDBC	v -SVR	$0.471 \pm 0.0198(4)$	$0.5494 \pm 0.0537(4)$	$0.7195 \pm 0.0166(4)$	0.061
	LS-SVR	$0.921 \pm 0.0202(2)$	$0.0812 \pm 0.0211(2)$	$0.0318 \pm 0.0094(1)$	0.177
	ε -TSVR	$0.902 \pm 0.0132(3)$	$0.0977 \pm 0.0124(3)$	$0.2362 \pm 0.0198(3)$	0.360
	v-MADR	$0.980 \pm 0.0157(1)$	$0.0202 \pm 0.0062(1)$	$0.0510 \pm 0.0251(2)$	1.112
average rank	v -SVR	4.0000	4.0000	2.8750	-
	LS-SVR	2.6250	2.7500	2.6250	-
	ε -TSVR	2.2500	2.0000	3.2500	-
	v-MADR	1.1250	1.2500	1.2500	-

TABLE 7

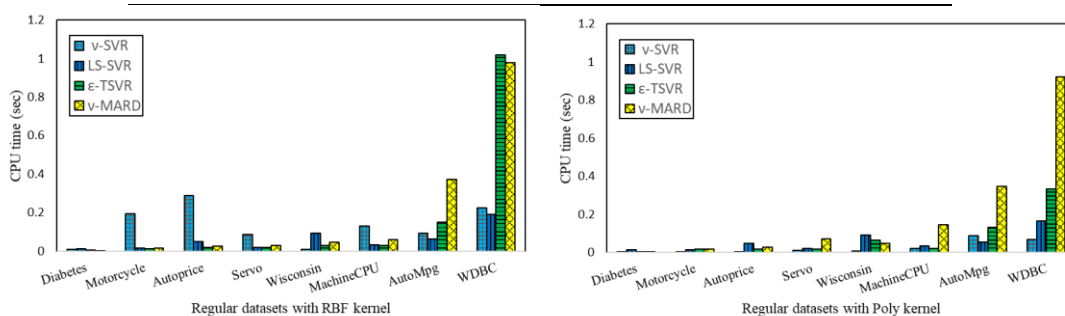
THE OPTIMAL PARAMETERS ON MEDIUM-SCALE DATASETS WITH RBF KERNEL

Dataset	v -SVR		LS-SVR		ε -TSVR			v -MADR		
	C	σ	C	σ	$c_1 = c_2$	$c_3 = c_4$	σ	C	$\lambda_1 = \lambda_2$	σ
Diabetes	2^9	2^{-4}	2^7	2^4	2^{-4}	2^{-3}	2^{-2}	2^5	2^7	2^{-3}
Motorcycle	2^9	2^0	2^4	2^{-2}	2^{-8}	2^{-6}	2^1	2^9	2^9	2^1
Autoprice	2^9	2^{-8}	2^8	2^{-6}	2^{-6}	2^{-9}	2^{-6}	2^{-9}	2^9	2^{-6}
Servo	2^6	2^{-1}	2^9	2^1	2^{-9}	2^{-9}	2^{-2}	2^{-7}	2^9	2^0
Wisconsin	2^2	2^{-9}	2^2	2^9	2^{-5}	2^2	2^{-7}	2^{-1}	2^6	2^{-9}
MachineCPU	2^8	2^{-9}	2^9	2^8	2^{-7}	2^{-9}	2^{-8}	2^{-9}	2^9	2^{-7}
AutoMpg	2^3	2^{-3}	2^4	2^2	2^{-5}	2^{-7}	2^{-3}	2^{-5}	2^9	2^{-3}
WDBC	2^2	2^{-5}	2^3	2^4	2^{-5}	2^{-6}	2^{-5}	2^{-9}	2^9	2^{-4}

TABLE 8

THE OPTIMAL PARAMETERS ON MEDIUM-SCALE DATASETS WITH POLYNOMIAL KERNEL

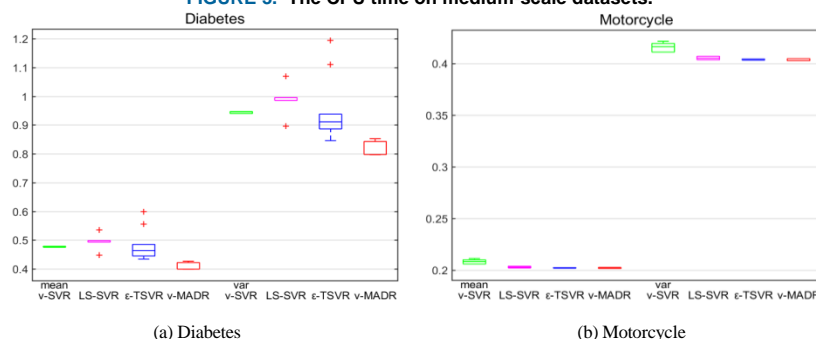
Dataset	v -SVR		LS-SVR		ε -TSVR			v -MADR			
	C	d	C	d	$c_1 = c_2$	$c_3 = c_4$	d	C	$\lambda_1 = \lambda_2$	d	
Diabetes	2^{-1}	3	2^{-4}	2	2^{-3}	2^8	2	2^9	2^{-5}	4	
Motorcycle	2^{-3}	2	2^{-3}	6	2^{-5}	2^7	6	2^{-4}	2^4	6	
Autoprice	2^{-1}	3	2^{-5}	2	2^{-3}	2^3	2	2^{-1}	2^3	2	
Servo	2^2	3	2^{-6}	5	2^{-9}	2^5	4	2^{-3}	2^{-4}	5	
Wisconsin	2^{-3}	3	2^{-9}	2	2^{-1}	2^9	2	2^9	2^{-9}	3	
MachineCPU	2^0	2	2^{-2}	2	2^{-8}	2^{-1}	2	2^{-3}	2^{-9}	2	
AutoMpg	2^2	3	2^{-4}	3	2^{-2}	2^9	2	2^{-3}	2^2	3	
WDBC	2^{-3}	3	2^{-7}	2	2^{-5}	2^4	2	2^9	2^{-3}	4	



(a) CPU time on medium-scale datasets with RBF kernel

(b) CPU time on medium-scale datasets with polynomial kernel

FIGURE 3. The CPU time on medium-scale datasets.



(a) Diabetes

(b) Motorcycle

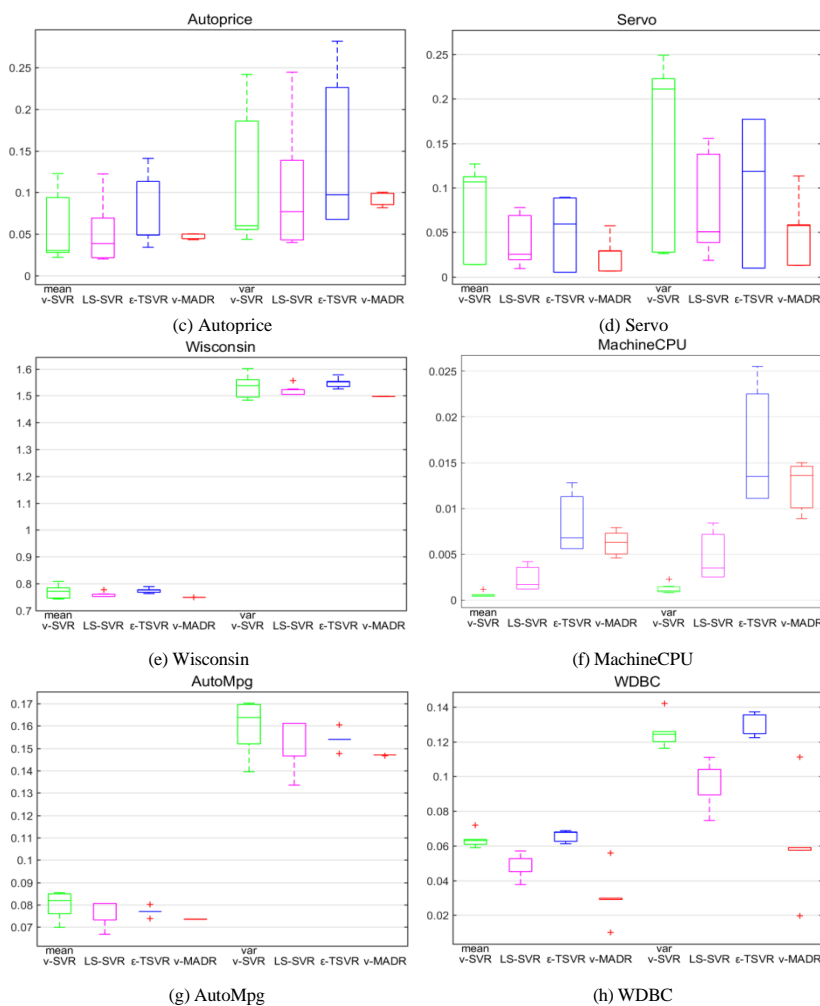
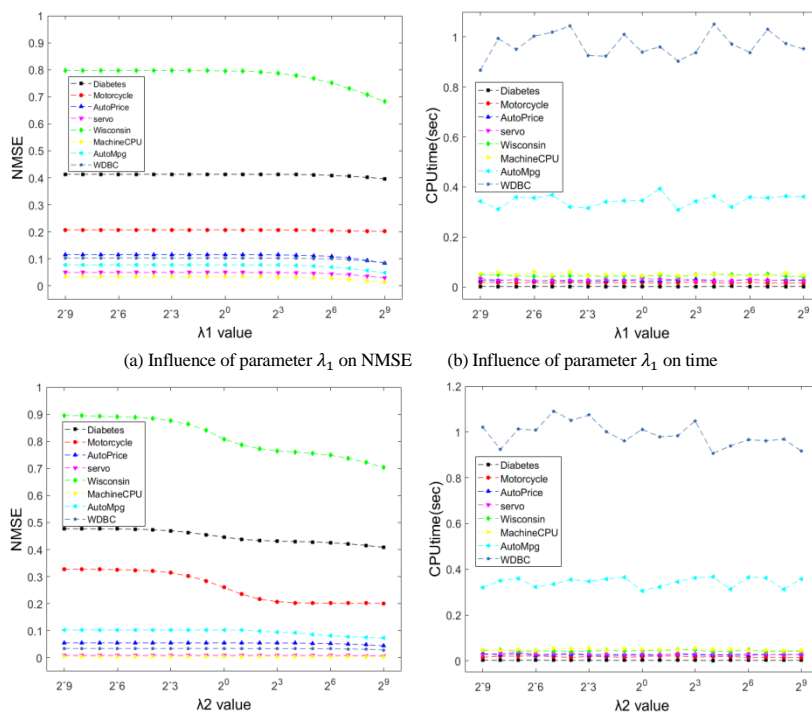


FIGURE 4. The absolute regression deviation mean and variance of v -MADR with RBF kernel on medium-scale datasets.



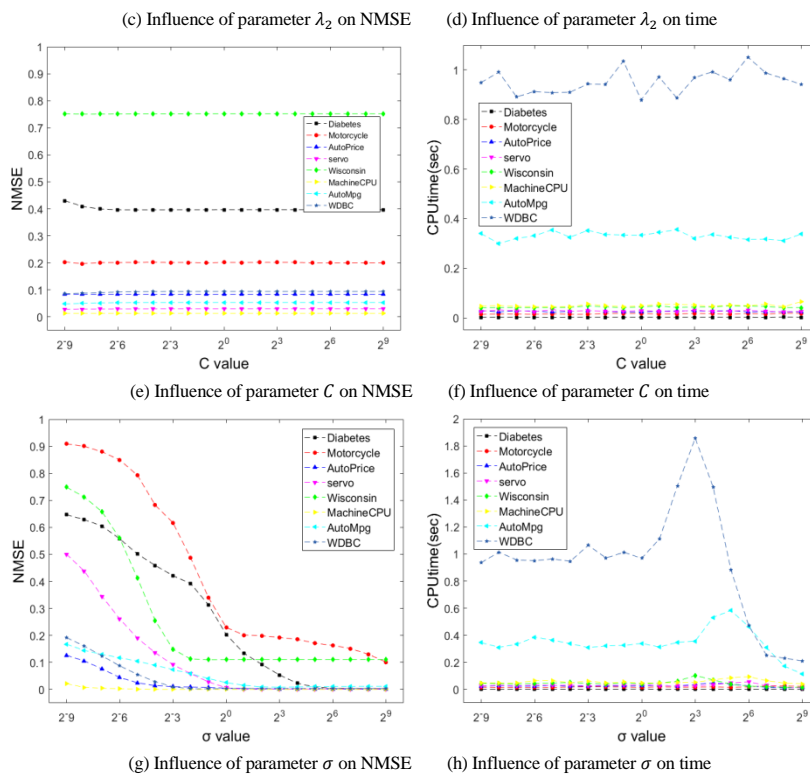


FIGURE 5. Parameter influence on NMSE and CPU time on medium-scale datasets with RBF kernel.

C. LARGE-SCALE DATASETS

Table 9 lists the experimental results on six large-scale datasets with linear kernel. We have additionally added a comparison of linear ϵ -SVR, which was solved by LIBLINEAR [50] that can handle large-scale datasets. In this experiment, because the datasets are too large, for each dataset, 2/3 of the dataset is randomly selected as the training set for feature selection, and the rest 1/3 of the dataset is used as the

test set for evaluation. From the average rank at the bottom of Table 9, the overall performance of ν -MADR is better than other compared methods or is highly competitive. The optimal parameters are listed in Table 10. Figure 6 shows the comparisons of CPU time. From Figure 6, linear kernel ν -MADR is the fastest learning method. In particular, the CPU time of linear kernel ν -MADR is far superior to ν -SVR, LS-SVR and ϵ -TSVR.

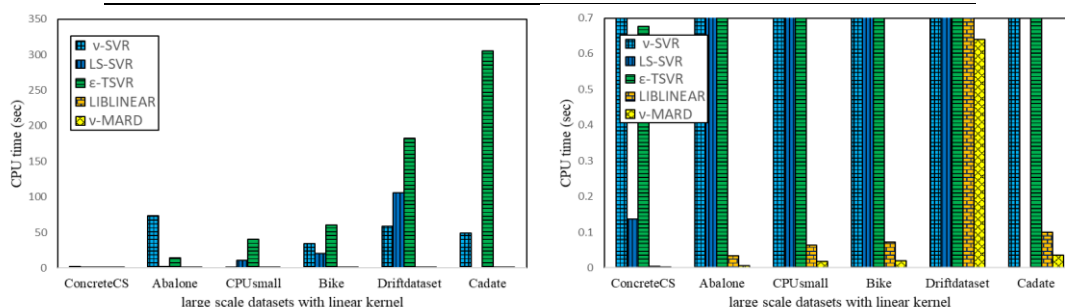
TABLE 9
THE RESULT COMPARISONS OF ν -SVR, LS-SVR, ϵ -TSVR, ϵ -SVR AND ν -MADR ON LARGE-SCALE DATASETS WITH LINEAR KERNEL

Dataset	regressor	R^2 (rank)	NMSE (rank)	MAPE (rank)	CPU(sec)
ConcreteCS	ν -SVR	0.5849 \pm 0.0511(4)	0.4180 \pm 0.0535(4)	2.1691 \pm 0.4274(4)	2.1470
	LS-SVR	0.6030 \pm 0.0440(2)	0.4033 \pm 0.0504(2)	2.0463 \pm 0.6259(2)	0.1371
	ϵ -TSVR	0.5934 \pm 0.0483(3)	0.4076 \pm 0.0476(3)	2.0746 \pm 0.5095(3)	0.6757
	LIBLINEAR	0.3945 \pm 0.0493(5)	0.7194 \pm 0.1146(5)	3.9093 \pm 2.1885(5)	0.0040
	ν -MADR	0.6124 \pm 0.0251(1)	0.3925 \pm 0.0278(1)	1.8771 \pm 0.6436(1)	0.0024
Abalone	ν -SVR	0.5237 \pm 0.0399(3)	0.4838 \pm 0.0403(3)	3.3437 \pm 0.3106(2)	73.7317
	LS-SVR	0.5103 \pm 0.0345(4)	0.4921 \pm 0.0351(4)	3.4757 \pm 0.4091(5)	1.7781
	ϵ -TSVR	0.5310 \pm 0.0267(2)	0.4712 \pm 0.0286(2)	3.4732 \pm 0.2775(4)	14.2850
	LIBLINEAR	0.3605 \pm 0.0295(5)	0.6708 \pm 0.0586(5)	3.3933 \pm 0.6411(3)	0.0337
	ν -MADR	0.5491 \pm 0.0322(1)	0.4569 \pm 0.0370(1)	3.2690 \pm 0.2572(1)	0.0065
CPUsmall	ν -SVR	0.6918 \pm 0.0371(5)	0.3120 \pm 0.0398(5)	4.8450 \pm 0.8492(2)	0.7953
	LS-SVR	0.7140 \pm 0.0185(3)	0.2882 \pm 0.0192(2)	6.0237 \pm 1.3026(5)	10.8428
	ϵ -TSVR	0.7160 \pm 0.0285(2)	0.2994 \pm 0.0627(4)	5.6246 \pm 1.0662(3)	40.5285
	LIBLINEAR	0.8579 \pm 0.0320(1)	0.2138 \pm 0.0367(1)	1.0918 \pm 0.3100(1)	0.0626
	ν -MADR	0.7107 \pm 0.0251(4)	0.2928 \pm 0.0281(3)	5.8786 \pm 0.5713(4)	0.0183
Bike	ν -SVR	0.3054 \pm 0.0237(4)	0.7102 \pm 0.0418(4)	2.0073 \pm 0.2880(1)	33.9155
	LS-SVR	0.3116 \pm 0.0108(2)	0.6887 \pm 0.0109(2)	2.1121 \pm 0.2953(3)	20.2850

	ϵ -TSVR	0.3088 \pm 0.0160(3)	0.6916 \pm 0.0167(3)	2.1389 \pm 0.3706(4)	60.5827
	LIBLINEAR	0.1456 \pm 0.0249(5)	1.2910 \pm 0.0973(5)	4.5741 \pm 0.4089(5)	0.0725
	<i>v</i>-MADR	0.3124 \pm 0.0101(1)	0.6879 \pm 0.0112(1)	2.0562 \pm 0.3886(2)	0.0200
Driftdataset	<i>v</i> -SVR	0.5838 \pm 0.0294(4)	0.4223 \pm 0.0341(3)	0.6057 \pm 0.0099(4)	58.8105
	LS-SVR	0.4904 \pm 0.2335(5)	0.6566 \pm 0.4629(5)	0.5328 \pm 0.0617(1)	105.9514
	ϵ -TSVR	0.6018 \pm 0.0916(3)	0.4035 \pm 0.0959(2)	0.5932 \pm 0.0931(2)	182.7503
	LIBLINEAR	0.6219 \pm 0.0285(1)	0.4270 \pm 0.0174(4)	0.7941 \pm 0.0161(5)	1.5930
	<i>v</i>-MADR	0.6145 \pm 0.0946(2)	0.3980 \pm 0.0456(1)	0.5966 \pm 0.0461(3)	0.6395
Cadate	<i>v</i> -SVR	0.6179 \pm 0.0138(3)	0.3845 \pm 0.0162(3)	2.3677 \pm 0.3229(3)	49.4465
	LS-SVR	N/A(5)	N/A(5)	N/A(5)	N/A
	ϵ -TSVR	0.6204 \pm 0.0061(2)	0.3799 \pm 0.0062(2)	2.3034 \pm 0.1976(2)	304.7820
	LIBLINEAR	0.4728 \pm 0.0132(4)	0.8760 \pm 0.0356(4)	5.0565 \pm 1.6076(4)	0.1007
	<i>v</i>-MADR	0.6207 \pm 0.0158(1)	0.3795 \pm 0.0169(1)	2.1349 \pm 0.3066(1)	0.0356
average rank	<i>v</i> -SVR	3.8333	3.6667	2.6667	-
	LS-SVR	3.5000	3.3333	3.5000	-
	ϵ -TSVR	2.5000	2.6667	3.0000	-
	LIBLINEAR	3.5000	4.0000	3.8333	-
	<i>v</i>-MADR	1.6667	1.3333	2.0000	-

TABLE 10
THE OPTIMAL PARAMETERS ON LARGE-SCALE DATASETS WITH LINEAR KERNEL.

Dataset	<i>v</i> -SVR	LS-SVR	ϵ -TSVR		LIBLINEAR		<i>v</i> -MADR	
	C	C	$c_1 = c_2$	$c_3 = c_4$	C	C	λ_1	λ_2
ConcreteCS	2 ⁷	2 ¹	2 ⁻³	2 ¹	2 ⁻⁹	2 ⁹	2 ⁻²	2 ⁻⁵
Abalone	2 ⁹	2 ⁹	2 ⁻³	2 ⁻¹	2 ⁻⁵	2 ⁹	2 ⁻⁸	2 ³
CPUsmall	2 ⁻⁸	2 ⁻⁶	2 ⁻⁹	2 ⁷	2 ⁻⁴	2 ⁹	2 ⁻⁸	2 ⁵
Bike	2 ⁵	2 ⁻⁵	2 ⁻⁹	2 ⁴	2 ⁻⁹	2 ⁹	2 ⁴	2 ³
Driftdataset	2 ⁻⁹	2 ⁻⁶	2 ⁻⁷	2 ⁷	2 ⁸	2 ³	2 ⁻⁵	2 ⁻⁹
Cadate	2 ³	N/A	2 ⁻⁷	2 ⁵	2 ¹	2 ⁹	2 ⁻¹	2 ⁻⁷



(a) CPU time on large datasets with linear kernel

(b) A closer look at the CPU time in the range of [0,0.7] sec on large datasets

FIGURE 6. The CPU time on large datasets with linear kernel.

V. CONCLUSIONS

In this research, we introduce statistical information into *v*-SVR and propose a novel SVR method called *v*-MADR. *v*-MADR improves the performance of SVR and overcomes the limitations of existing SVR algorithms by minimizing both the absolute regression deviation mean and the absolute regression deviation variance, which takes into account both the positive and negative values of the regression deviation of sample points. *v*-MADR proposes a dual coordinate descent (DCD) algorithm for small sample problems, and we also propose an averaged stochastic gradient descent (ASGD) algorithm for large-scale problems, which greatly reduces the

computational complexity and thus improves the algorithm speed. We provide a theoretical analysis on the boundary of the expectation of error for *v*-MADR. Experimental results have shown that *v*-MADR outperforms several regression methods and demonstrates great application potential. Our *v*-MADR Matlab codes can be accessed from: <https://github.com/AsunaYY/v-MADR>.

In the near future, we will further investigate the potential of *v*-MADR for big data problems, e.g., predictive analysis for bioinformatics and systems biology problems, and problems in finance. We envision a great application potential in these problems.

REFERENCES

- [1] A. J. Smola and B. Scholkopf, "A tutorial on support vector regression," *Statistics and Computing*, vol. 14, no. 3, pp. 199-222, 2004.
- [2] A. J. Smola and B. Scholkopf, *Learning with Kernels. (1st ed.)*, 1998.
- [3] J. D. Brown, M. F. Summers, and B. A. Johnson, "Prediction of hydrogen and carbon chemical shifts from RNA using database mining and support vector regression," *Journal of Biomolecular Nmr*, vol. 63, no. 1, pp. 39-52, 2015.
- [4] P. K. Rajaraman, T. A. Manteuffel, M. Belohlavek, E. McMahon, and J. J. Heys, "Echocardiographic particle imaging velocimetry data assimilation with least square finite element methods," *Computers & Mathematics with Applications*, vol. 68, no. 11, pp. 1569-1580, 2014.
- [5] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929-1958, 2014.
- [6] Y. Ke, B. Fu, and W. Zhang, "Semi-varying coefficient multinomial logistic regression for disease progression risk prediction," *Statistics in Medicine*, vol. 35, no. 26, pp. 4764-4778, 2016.
- [7] L. H. Dicker, "Ridge regression and asymptotic minimax estimation over spheres of growing dimension," *Bernoulli*, vol. 22, no. 1, pp. 1-37, 2016.
- [8] R. M. Balabin and E. I. Lomakina, "Support vector machine regression (SVR/LS-SVM)—an alternative to neural networks (ANN) for analytical chemistry? Comparison of nonlinear methods on near infrared (NIR) spectroscopy data," *Analyst*, vol. 136, no. 8, pp. 1703-1712, 2011.
- [9] I. Khosravi, Y. Jouybari-Moghaddam, and M. R. Sarajian, "The comparison of NN, SVR, LSSVR and ANFIS at modeling meteorological and remotely sensed drought indices over the eastern district of Isfahan, Iran," *Natural Hazards*, vol. 87, no. 3, pp. 1507-1522, 2017.
- [10] N. Suyaraj, N. Theera-Umpon, and S. Auephanwiriyakul, "A comparison of NN-based and SVR-based power prediction for mobile DS/CDMA systems," in *2008 International Symposium on Intelligent Signal Processing and Communications Systems*, pp. 1-4, 2009.
- [11] C.-C. Chang and C.-J. Lin, "LIBSVM: A Library for Support Vector Machines," *Acm Transactions on Intelligent Systems and Technology*, vol. 2, no. 3, 2011.
- [12] B. Schölkopf, A. J. Smola, R. C. Williamson, and P. L. Bartlett, "New Support Vector Algorithms," *Neural Computation*, vol. 12, no. 5, pp. 1207-1245, 2000.
- [13] J. A. K. Suykens, L. Lukas, P. Van Dooren, B. De Moor, and J. Vandewalle, "Least squares support vector machine classifiers: a large scale algorithm," *European Conference on Circuit Theory and Design*, vol. 99, 1999.
- [14] J. A. K. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural Processing Letters*, vol. 9, no. 3, pp. 293-300, 1999.
- [15] X. Peng, "TSVR: An efficient Twin Support Vector Machine for regression," *Neural Networks*, vol. 23, no. 3, pp. 365-372, 2010.
- [16] Y. H. Shao, C. H. Zhang, Z. M. Yang, L. Jing, and N. Y. Deng, "An ϵ -twin support vector machine for regression," *Neural Computing & Applications*, vol. 23, no. 1, pp. 175-185, 2013.
- [17] Z.-M. Yang, X.-Y. Hua, Y.-H. Shao, and Y.-F. Ye, "A novel parametric-insensitive nonparallel support vector machine for regression," *Neurocomputing*, vol. 171, pp. 649-663, 2016.
- [18] S. Balasundaram and G. Benipal, "On a new approach for Lagrangian support vector regression," *Neural Computing and Applications*, vol. 29, no. 9, pp. 533-551, 2018.
- [19] S. Balasundaram and D. Gupta, "On implicit Lagrangian twin support vector regression by Newton method," *International Journal of Computational Intelligence Systems*, vol. 7, no. 1, pp. 50-64, 2014.
- [20] T. Zhang and Z.-H. Zhou, "Large margin distribution machine," in *20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 313-322, 2014.
- [21] Z.-H. Zhou, "Large margin distribution learning," in *6th IAPR TC3 International Workshop on Artificial Neural Networks for Pattern Recognition*, vol. 8774, pp. 1-11, 2014.
- [22] T. Zhang and Z.-H. Zhou, "Optimal Margin Distribution Machine," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1-1, 2019.
- [23] W. Gao and Z.-H. Zhou, "On the doubt about margin explanation of boosting," *Artificial Intelligence*, vol. 203, pp. 1-18, 2013.
- [24] M.-Z. Liu, Y.-H. Shao, Z. Wang, C.-N. Li, and W.-J. Chen, "Minimum deviation distribution machine for large scale regression," *Knowledge-Based Systems*, vol. 146, pp. 167-180, 2018.
- [25] R. Rastogi, P. Anand, and S. Chandra, "Large-margin Distribution Machine-based regression," *Neural Computing and Applications*, vol. 32, no. 8, pp. 3633-3648, 2020.
- [26] P. Anand, R. Rastogi, and S. Chandra, "Generalized ϵ -Loss Function-Based Regression," in *Machine Intelligence and Signal Analysis*, pp. 395-409, 2019.
- [27] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A library for large linear classification," *Journal of Machine Learning Research*, vol. 9, pp. 1871-1874, 2008.
- [28] B. Gu, V. S. Sheng, Z. Wang, D. Ho, S. Osman, and S. Li, "Incremental learning for v-Support Vector Regression," *Neural Networks*, vol. 67, pp. 140-150, 2015.
- [29] B. Gu, V. S. Sheng, K. Y. Tay, W. Romano, and S. Li, "Incremental Support Vector Learning for Ordinal Regression," *Ieee Transactions on Neural Networks and Learning Systems*, vol. 26, no. 7, pp. 1403-1416, 2015.
- [30] B. Gu and V. S. Sheng, "A Robust Regularization Path Algorithm for v-Support Vector Classification," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 5, pp. 1241-1248, 2017.
- [31] D. J. Crisp and C. J. C. Burges, "A geometric interpretation of v-SVM classifiers," in *13th Annual Neural Information Processing Systems Conference*, pp. 244-250, 2000.
- [32] J. S. Marron, M. J. Todd, and J. Ahn, "Distance-weighted discrimination," *Journal of the American Statistical Association*, vol. 102, no. 480, pp. 1267-1271, 2007.
- [33] B. Schölkopf, "Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond," *Publications of the American Statistical Association*, vol. 98, no. 462, pp. 489-489, 2003.
- [34] A. F. Izmailov and M. V. Solodov, "Karush-Kuhn-Tucker systems: regularity conditions, error bounds and a class of Newton-type methods," *Mathematical Programming*, vol. 95, no. 3, pp. 631-650, 2003.
- [35] D. Anguita, A. Ghio, S. Pisciutta, and S. Ridella, "A support vector machine with integer parameters," *Neurocomputing*, vol. 72, no. 1-3, pp. 480-489, 2008.
- [36] L. Oneto, A. Ghio, S. Ridella, and D. Anguita, "Learning Resource-Aware Classifiers for Mobile Devices: From Regularization to Energy Efficiency," *Neurocomputing*, vol. 169, pp. 225-235, 2015.
- [37] L. Oneto, J. L. R. Ortiz, and D. Anguita, "Constraint-Aware Data Analysis on Mobile Devices: An Application to Human Activity Recognition on Smartphones," in *Adaptive Mobile Computing*, 2017.
- [38] C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and S. Sundararajan, "A dual coordinate descent method for large-scale linear SVM," in *25th International Conference on Machine Learning*, pp. 408-415, 2008.
- [39] G.-X. Yuan, C.-H. Ho, and C.-J. Lin, "Recent Advances of Large-Scale Linear Classification," *Proceedings of the Ieee*, vol. 100, no. 9, pp. 2584-2603, 2012.
- [40] B. T. Polyak and A. B. Juditsky, "Acceleration of stochastic approximation by averaging," *SIAM Journal on Control and Optimization*, vol. 30, no. 4, pp. 838-855, 1992.
- [41] H. Kushner and G. G. Yin, *Stochastic approximation and recursive algorithms and applications*. Springer Science & Business Media, 2003.
- [42] L. Bottou, *Large-Scale Machine Learning with Stochastic Gradient Descent*. 2010.
- [43] T. Zhang, "Solving large scale linear prediction problems using stochastic gradient descent algorithms," in *Proceedings, Twenty-First International Conference on Machine Learning*, pp. 919-926, 2004.
- [44] S. Shalev-Shwartz, Y. Singer, and N. Srebro, "Pegasos: Primal estimated sub-Gradient solver for SVM," in *24th International Conference on Machine Learning*, vol. 227, pp. 807-814, 2007.
- [45] A. Bordes, L. Bottou, and P. Gallinari, "SGD-QN: Careful Quasi-Newton Stochastic Gradient Descent," *Journal of Machine Learning Research*, vol. 10, pp. 1737-1754, 2009.
- [46] O. Shamir and T. Zhang, "Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes," in *30th International Conference on Machine Learning*, pp. 71-79, 2013.
- [47] W. Xu, "Towards Optimal One Pass Large Scale Learning with Averaged Stochastic Gradient Descent," *arXiv e-prints*, 2011.
- [48] A. Luntz, "On estimation of characters obtained in statistical procedure of recognition," *Technicheskaya Kibernetika*, vol. 3, 1969.

- [49] C.-C. Chang and C.-J. Lin, "LIBSVM -- A Library for Support Vector Machines," ed. <https://www.csie.ntu.edu.tw/~cjlin/libsvm>, 2018.
- [50] Machine Learning Group at National Taiwan University, "LIBLINEAR -- A Library for Large Linear Classification," ed. <https://www.csie.ntu.edu.tw/~cjlin/liblinear>, 2019.
- [51] K. Pelckmans *et al.*, "LSSVmlab: A Matlab/c toolbox for least square support vector machines," ed. <http://www.esat.kuleuven.ac.be/sista/lssvmlab>, 2002.
- [52] Y.-H. Shao, C.-H. Zhang, X.-B. Wang, and N.-Y. Deng, "Improvements on Twin Support Vector Machines," *Ieee Transactions on Neural Networks*, vol. 22, no. 6, pp. 962-968, 2011.
- [53] O. L. Mangasarian and D. R. Musicant, "Successive overrelaxation for support vector machines," *Ieee Transactions on Neural Networks*, vol. 10, no. 5, pp. 1032-1037, 1999.
- [54] C. W. Hsu and C. J. Lin, "A comparison of methods for multiclass support vector machines," *Ieee Transactions on Neural Networks*, vol. 13, no. 2, pp. 415-425, 2002.



Yan Wang, Ph.D., graduated from the College of Computer Science and Technology of Jilin University, China, in 2007. He is currently a professor and doctoral supervisor of the College of Computer Science and Technology of Jilin University. He was doing research on a bioinformatics collaborative project as a visiting scholar and postdoctoral at University of Georgia, USA from 2007 to 2011. During 2012 to 2013, he got a postdoctoral position in University of Trento, Italy. His research interests focus on Computational Intelligence, Bioinformatics and E-Business, such as machine learning and deep learning, big data intelligent computing, gene expression analysis, cancer biomarker prediction and overlapping community detection. He has had over 70 research papers published including more than 30 indexed by SCI, such as *Nucleic Acids Res*, *Briefing in Bioinformatics*, *Computational and Structural Biotechnology Journal*, *Frontiers in Genetics*, *Scientific Reports*, *Applied Soft Computing*, and *Bioinformatics*, *etc.* And he has also presided over or taken part in several projects funded by the National Natural Science Foundation and "863" project in China, respectively.



image processing.



Yingying Song, Master, is currently studying at the College of Computer Science and Technology of Jilin University, China. She received her bachelor's degree in the College of Computer Science and Technology of Jilin University, in 2017. She has published two papers and won two school scholarships. Her research interests are mainly in the field of machine learning, regression analysis, feature selection, and deep learning.



machine learning, and deep learning.

Xuping Xie is currently pursuing a master's degree at the College of Computer Science and Technology of Jilin University, China. She is also a student in the Jilin Provincial Key Laboratory of Big Data Intelligent Computing. She studied at the College of Information Science and Technology of Northeast Normal University from 2015 to 2019 and obtained a bachelor's degree. Her research interests are mainly in the field of



LAN HUANG is currently a professor, supervisor of PhD candidates at Jilin University, Changchun. She received her B.S., M.S. and Ph.D. degrees in computer science and technology from Jilin University, Changchun, in 1994, 1999 and 2003 respectively. She is mainly engaged in intelligent computing, data mining theory and application research, and high performance computing. She is Outstanding Member of CCF. She won the middle-aged and young experts with outstanding contributions in Jilin province. She is the director of key laboratory of big data intelligent computing in Jilin province. As PI and Co-PI, she has been undertaking or accomplished more than 10 teaching and scientific research projects, granted by the National 863 Hi-tech Research and Development Program, the National Science Foundation China, provincial/ministerial foundations, and other sources. The projects that she participated as main investigator, were awarded the first prizes for the Jilin Province Scientific and Technological Progress Award (first prizewinner in 2017, second prizewinner in 2019). In recent five years, she has published more than 80 academic papers, won 1 second prize for teaching achievements of Jilin Province (2018), published 4 textbooks as chief editor.



WEI PANG received his Ph.D. degree in computing science from the University of Aberdeen, in 2009. He is currently an Associate Professor with Heriot-Watt University, Edinburgh, U.K. He is also an Honorary Senior Lecturer at University of Aberdeen, UK. He has authored over 100 articles, including more than 40 journal articles. His research interests include bio-inspired computing, data mining, machine learning, and qualitative reasoning.

Dr. Pang was a recipient of the Best Paper Award in the 19th Annual UK Workshop on Computational Intelligence (UKCI 2019) and the Best Paper Runner Up Award in the 12th International Conference on Advanced Data Mining and Applications (ADMA 2016).



George M. Coghill, Ph.D., graduated from the Intelligent Systems Laboratory of Heriot-Watt University on the topic of Fuzzy Qualitative Reasoning. He was promoted to Reader of Aberdeen in Computing Science in 2009, and to a Chair in 2011. At Aberdeen he has undertaken a number of roles. He was Head of Computing Science from 2010 to 2014; Director of Research and a Deputy Head of the School of Natural and Computing Sciences from 2015 to 2018. He has served on a several committees in the institution including the Academic Standards Committee (Post-graduate), the Quality Assurance Committee, and the Research Policy Committee. His research interests are mainly in the field of Model-based Systems and Qualitative reasoning.