

## 密度加权孪生支持向量回归机

程昊翔, 王 坚

(同济大学 a. 电子与信息工程学院, b. CIMS 研究中心, 上海 201804)

**摘要:** 为了使数据集的内在分布更好地影响训练模型, 提出一种密度加权孪生支持向量回归机算法. 该算法通过  $k$  近邻算法计算获得每个数据点基于数据密度分布的密度加权值, 并将密度加权值引入到标准孪生支持向量回归机算法中. 算法能够很好地反映训练数据集的内在分布, 使数据点准确影响训练模型. 通过 6 个 UCI 数据集上的实验结果分析验证了所提出算法的有效性.

**关键词:** 密度加权;  $k$  近邻法; 内在分布

**中图分类号:** TP273

**文献标志码:** A

## Density-weighted twin support vector regression

CHENG Hao-xiang, WANG Jian

(a. College of Electronics and Information Engineering, b. CIMS Research Center, Tongji University, Shanghai 201804, China. Correspondent: CHENG Hao-xiang, E-mail: 373668304@qq.com)

**Abstract:** To better impact the training model with the inherent distribution of the training dataset, a twin support vector regression called density-weighted twin support vector regression is proposed. Firstly, the density-weighted value is computed based on the  $k$ -nearest neighbor algorithm. Then, the values of density-weighted are introduced to the standard twin support vector regression. It is found that the proposed algorithm can well reflect the inherent distribution of the training dataset and lead to a more accurate impact on the training model. Experimental results on six UCI datasets show the effectiveness of the proposed algorithm.

**Keywords:** density-weighted;  $k$ -nearest neighbor; inherent distribution

## 0 引 言

近几年, 支持向量机<sup>[1]</sup>(SVM)已成功应用于分类和回归问题. SVM 是以 VC 维理论和结构风险最小化为理论基础<sup>[2]</sup>提出的, 相比于人工神经网络等其他机器学习算法, SVM 有如下优势: 1) SVM 解决一个带约束的二次规划问题, 能够得到唯一的最优解; 2) SVM 基于结构风险最小化理论原则, 有更好的泛化性. 然而, SVM 对于大数据量的模型训练问题, 存在时间复杂度较高的问题. SVM 模型训练的时间复杂度为  $O(m^3)$ ,  $m$  为训练样本的规模.

最近, Jayadeva 等<sup>[3]</sup>提出了一种用于解决分类问题的新算法——孪生支持向量机 (TWSVM). TWSVM 算法的动机来源于 GEPSVM<sup>[4]</sup>, 与 SVM 训练获得两个平行的超平面不同, TWSVM 通过解决两个规模较小的二次规划问题获得两个不平行的超平面. 通过计

算分析, 在数据集规模相同的情况下, TWSVM 的训练时间是标准 SVM 的 1/4. Peng<sup>[5]</sup>将 TWSVM 从解决分类问题扩展到回归问题, 提出了孪生支持向量回归机 (TSVR). TSVR 经过近几年的研究有了不同的扩展<sup>[6-10]</sup>, 但是这些研究所考虑的每一个训练数据对于生成超平面的贡献和作用是不同的, 没有反映出数据在真实空间中的内在分布情况.

针对上述问题, 本文提出一种基于密度加权的孪生支持向量回归机 (DW-TSVR), 在目标函数中对不同数据所对应的经验误差进行密度加权, 密度加权值在真实空间中采用 kNN 算法进行计算, 并通过计算规则将密度加权值的范围缩小在  $[0,1]$ . 对标准 TSVR 的目标函数进行正则化, 将其扩展到结构风险最小化误差. 在 6 个 UCI 数据集上的实验比较结果表明, 所提出的 DW-TSVR 相比于 TSVR、SVR、Weighted

收稿日期: 2015-01-10; 修回日期: 2015-04-11.

基金项目: 国家自然科学基金面上项目(71273188); 国家自然科学基金重大项目(91024031).

作者简介: 程昊翔(1986—), 男, 博士, 从事系统工程、机器学习的研究; 王坚(1961—), 男, 教授, 博士生导师, 从事系统工程、智能优化算法等研究.

TSVR<sup>[6]</sup>具有更高的精确度.

## 1 TSVR

对 TSVR 进行简单介绍. 假设一个大小为  $m$  的数据集  $S$ , 表示为

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \in (\chi \times y)^m. \quad (1)$$

其中:  $\chi \in R^n$ ,  $y \in R$ ,  $x_i \in R^n$ ,  $i = 1, 2, \dots, m$ ,  $n$  为输入空间的特征数. 目标是通过构建两个优化问题获得两个函数来预测无输出的实例. 两个预测函数表示为

$$f_1(x) = w_1^T x + b_1, \quad f_2(x) = w_2^T x + b_2.$$

最后的预测函数为

$$f(x) = \frac{1}{2}(f_1(x) + f_2(x)).$$

TSVR 的几何解释如图 1 所示.

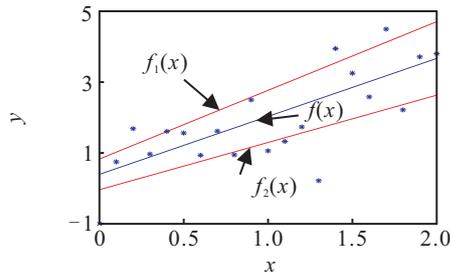


图 1 TSVR 几何解释

为了实现本文目标, TSVR 构建如下带约束的优化函数:

$$\begin{aligned} \min \quad & \frac{1}{2}(Y - e\epsilon_1 - (Aw_1 + eb_1))^T \times \\ & (Y - e\epsilon_1 - (Aw_1 + eb_1)) + c_1 e^T \xi, \\ \text{s.t.} \quad & Y - (Aw_1 + eb_1) \geq e\epsilon_1 - \xi, \quad \xi \geq 0; \end{aligned} \quad (2)$$

$$\begin{aligned} \min \quad & \frac{1}{2}(Y - e\epsilon_2 - (Aw_2 + eb_2))^T \times \\ & (Y - e\epsilon_2 - (Aw_2 + eb_2)) + c_2 e^T \eta, \\ \text{s.t.} \quad & (Aw_2 + eb_2) - Y \geq e\epsilon_2 - \eta, \quad \eta \geq 0. \end{aligned} \quad (3)$$

其中

$$c_1, c_2 > 0, \quad A = \{x_1^T, x_2^T, \dots, x_m^T\} \in R^{m \times n},$$

$$Y = \{y_1, y_2, \dots, y_m\} \in R^m, \quad e \in R^m, \quad \epsilon_1, \epsilon_2 \geq 0, \quad \eta, \xi \in R^m \text{ 为松弛变量.}$$

通过拉格朗日变换, 引入  $\alpha$  和  $\gamma$ , 将式 (2) 和 (3) 转化为如下对偶形式:

$$\begin{aligned} \max \quad & -\frac{1}{2}\alpha^T G(G^T G)^{-1} G^T \alpha + \\ & f^T G(G^T G)^{-1} G^T \alpha - f^T \alpha, \\ \text{s.t.} \quad & 0 \leq \alpha \leq c_1 e; \end{aligned} \quad (4)$$

$$\begin{aligned} \max \quad & -\frac{1}{2}\gamma^T G(G^T G)^{-1} G^T \gamma - \\ & h^T G(G^T G)^{-1} G^T \gamma + h^T \gamma, \\ \text{s.t.} \quad & 0 \leq \gamma \leq c_2 e. \end{aligned} \quad (5)$$

其中

$$G = [A \quad e], \quad f = Y - e\epsilon_1, \quad h = Y + e\epsilon_2,$$

$$u_1 = [w_1^T \quad b_1], \quad u_2 = [w_2^T \quad b_2].$$

在式 (5) 对偶转化过程中可以得到  $u_1$ 、 $u_2$  与  $\alpha$ 、 $\gamma$  之间的关系为

$$u_1 = (G^T G + \sigma I)^{-1} G^T (f - \alpha),$$

$$u_2 = (G^T G + \sigma I)^{-1} G^T (h + \gamma).$$

通过求解式 (4) 和 (5) 两个二次规划问题, 可以求得  $\alpha$  和  $\gamma$ , 进一步求解得到  $u_1$  和  $u_2$ , 从而获得回归预测模型  $f(x)$  为

$$\begin{aligned} f(x) &= \frac{1}{2}(f_1(x) + f_2(x)) = \\ & \frac{1}{2}(w_1 + w_2)^T x + \frac{1}{2}(b_1 + b_2). \end{aligned}$$

## 2 DW-TSVR

引入密度加权值  $\rho_i$  ( $i = 1, 2, \dots, m$ ), 定义为

$$\rho_i = 1 - \frac{d(x_i, x_i^k)}{\max_{j \in \text{train set}} d(x_j, x_j^k)}. \quad (6)$$

其中:  $x_i^k$  为通过 kNN 算法<sup>[11]</sup>在输入原空间计算获得的与数据点  $x_i$  第  $k$  个邻近的点, 本文取  $k = 10$ ;  $d(x_i, x_i^k)$  为  $x_i$  与  $x_i^k$  之间在原空间的距离.

密度加权值衡量了训练数据点的相对密度, 通过计算每个数据点与其对应的第  $k$  个近邻点之间的距离与训练数据中此计算距离的最大值进行比较获得. 密度加权值的范围为  $[0, 1]$ . 通过密度加权值的引入, 当数据点拥有较大密度加权值时, 其对于训练获得的模型有较大的影响, 能准确地反映训练数据集的内在分布.

通过上述计算, 训练数据集  $S$  中每个数据点  $x_i$  都有对应的密度加权值  $\rho_i$ , 将  $\rho_i$  组成  $m \times m$  大小的对角矩阵  $W = \text{diag}(\rho_1, \rho_2, \dots, \rho_m)$ , 其对角线值为对应的密度加权值.

DW-TSVR 将标准 TWSVR 的经验风险最小化加入正则项  $(w^2 + b^2)/2$ , 扩展到结构风险最小化, 然后将密度加权矩阵  $W$  引入到新的目标函数中, 形成新的二次规划问题, 表示为

$$\begin{aligned} \min \quad & \frac{1}{2}(Y - (Aw_1 + eb_1))^T W (Y - (Aw_1 + eb_1)) + \\ & \frac{1}{2}c_3(w_1^T w_1 + b_1^2) + c_1 e^T \xi, \\ \text{s.t.} \quad & Y - (Aw_1 + eb_1) \geq -e\epsilon_1 - \xi, \quad \xi \geq 0; \end{aligned} \quad (7)$$

$$\begin{aligned} \min \quad & \frac{1}{2}(Y - (Aw_2 + eb_2))^T W (Y - (Aw_2 + eb_2)) + \\ & \frac{1}{2}c_4(w_2^T w_2 + b_2^2) + c_2 e^T \eta, \\ \text{s.t.} \quad & (Aw_2 + eb_2) - Y \geq -e\epsilon_2 - \eta, \quad \eta \geq 0. \end{aligned} \quad (8)$$

其中:  $c_1, c_2, c_3, c_4 > 0$ ;  $\eta, \xi \in R^m$  为松弛变量.

引入  $\alpha$  和  $\gamma$ , 通过拉格朗日变换, 得到式(7)和(8)的对偶问题形式为

$$\begin{aligned} \max \quad & -\frac{1}{2}\alpha^T G(G^T W G + c_3 I)^{-1} G^T \alpha + \\ & Y^T G(G^T W G + c_3 I)^{-1} G^T \alpha - \\ & (e^T \epsilon_1 + Y^T W) \alpha, \\ \text{s.t.} \quad & 0 \leq \alpha \leq c_1 e; \end{aligned} \tag{9}$$

$$\begin{aligned} \max \quad & -\frac{1}{2}\gamma^T G(G^T W G + c_4 I)^{-1} G^T \gamma - \\ & Y^T G(G^T W G + c_4 I)^{-1} G^T \gamma + \\ & (Y^T W - e^T \epsilon_2) \gamma, \\ \text{s.t.} \quad & 0 \leq \gamma \leq c_2 e. \end{aligned} \tag{10}$$

在上述问题(9)和(10)的转化过程中得到

$$\begin{aligned} h_1 &= [w_1^T \ b_1]^T = \\ & (G^T W G + c_3 I)^{-1} G^T (W Y - \alpha), \\ h_2 &= [w_2^T \ b_2]^T = \\ & (G^T W G + c_4 I)^{-1} G^T (W Y + \gamma). \end{aligned}$$

由式(9)和(10)两个带约束的二次规划问题可以求解  $\alpha$  和  $\gamma$ , 从而得到  $h_1$  和  $h_2$ , 最终得到如下预测回归模型:

$$\begin{aligned} f(x) &= \frac{1}{2}(f_1(x) + f_2(x)) = \\ & \frac{1}{2}(w_1 + w_2)^T x + \frac{1}{2}(b_1 + b_2). \end{aligned}$$

### 3 逐次超松弛技术

采用逐次超松弛技术<sup>[12]</sup>, 对式(4)、(5)、(9)、(10)这4个二次规划问题进行求解. 首先, 将上述4个问题转化为如下统一形式:

$$\min \frac{1}{2} \alpha^T Q \alpha - d^T \alpha, \tag{11}$$

$$\text{s.t. } 0 \leq \alpha \leq c. \tag{12}$$

其中  $Q$  为正定矩阵. 以二次规划问题(9)为例, 有

$$\begin{aligned} Q &= G(G^T W G + c_3 I)^{-1} G^T, \\ d &= Y^T G(G^T W G + c_3 I)^{-1} G^T - (e^T \epsilon_1 + Y^T W), \\ c &= c_1 e. \end{aligned}$$

逐次超松弛技术描述如下. 已知  $\alpha^i$ , 按照下式计算  $\alpha^{i+1}$ :

$$\begin{aligned} \alpha^{i+1} &= \\ & (\alpha^i - tE^{-1}(Q\alpha^i - d + L(\alpha^{i+1} - \alpha^i))). \end{aligned} \tag{13}$$

其中:  $t \in (0, 2)$ ,  $L \in R^{m \times m}$  为  $Q$  的下三角部分, 对角矩阵  $E \in R^{m \times m}$  为  $Q$  的对角线. 逐次超松弛技术按式(13)迭代计算, 直到  $\|\alpha^{i+1} - \alpha^i\|$  小于某个规定的容错值, 本文容错值取  $10^{-3}$ .

逐次超松弛技术对解决大数据问题的孪生支持向量回归机非常有效, 而且已经证明该算法对于二次

规划问题求解线性收敛<sup>[13]</sup>. 文中标准 TWSVR 和 DW-TWSVR 算法均采用逐次超松弛技术进行求解.

### 4 实验分析

为了验证 DW-TSVR 算法的回归性能, 选择 UCI 标准数据集上的 6 个数据集, 使用高斯核将所提出的基于密度加权的孪生支持向量回归机与标准 SVR, 标准 TWSVR 算法和 Weighted TSVR 进行对比测试, 表 1 为 6 个数据集的规模和属性信息.

表 1 测试数据集

数据集名称	数据量	特征维度
Boston housing	506	13
Auto price	159	15
Auto-Mpg	398	7
Machine CPU	209	17
Diabetes	43	2
Pyrim	74	28

采用网格搜索算法 5 折交叉验证算法进行参数的选择. 5 折交叉验证即随机将数据集分成 5 个子集, 4 个子集用来训练模型, 剩下的子集用来测试训练得到的模型, 最后将 5 次的评价结果求平均值获得最终的训练模型. 高斯核函数为

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{r^2}}.$$

设定  $c_1 = c_2, c_3 = c_4, \epsilon_1 = \epsilon_2$ . 超参数  $c_1, c_2, c_3, c_4$  从  $\{2^i | i = -2, -1, \dots, 8\}$  中选取; 高斯核参数  $r$  从  $\{2^i | i = -4, -3, \dots, 8\}$  中选取;  $\epsilon_1, \epsilon_2$  从  $\{0.1, 0.2, \dots, 0.9, 1\}$  中选取. 性能评价标准采用均方根误差 (RMSE) 标准和规一化均方误差 (NMSE) 两个评价标准. RMSE 定义为

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{m}};$$

NMSE 定义为

$$\text{NMSE} = 10 \log_{10} \frac{\sum_{i=1}^m |y_i - \hat{y}_i|^2}{\sum_{i=1}^m |y_i|^2}.$$

其中:  $m$  为训练数据集数据的数量,  $y_i$  为实际的第  $i$  个输出值,  $\hat{y}_i$  为预测的第  $i$  个输出值. 表 2 为 SVR、TSVR、Weighted TSVR 和 DW-TSVR 四个算法通过 5 折交叉验证得到的各个数据集的 RMSE、NMSE 值.

从实验结果看, 在选取的 6 个数据集中, 本文所提出的 DW-TSVR 算法较 SVR、TSVR、Weighted TSVR 算法在回归精确度上有较明显的提高, 在 RMSE 和 NMSE 两个指标上都显示出了较好的性能, 验证了所提出的 DW-TSVR 算法的有效性. 本文所提出的密度加权孪生支持向量回归机算法通过引入密

表 2 算法测试结果比较

数据集	算法	RMSE	NMSE	Time/s
Boston housing	SVR	8.307	0.148	486.5
	TSVR	7.968	0.122	0.230
	Weighted TSVR	7.605	0.184	0.172
	DW-TSVR	7.462	0.118	0.170
Auto price	SVR	7824.6	0.389	8.677
	TSVR	7233.0	0.320	0.425
	Weighted TSVR	3845.0	0.303	0.422
	DW-TSVR	3432.6	0.288	0.410
Auto-Mpg	SVR	9.106	0.127	388.6
	TSVR	4.212	0.116	1.876
	Weighted TSVR	3.656	0.114	1.890
	DW-TSVR	3.356	0.112	1.825
Machine CPU	SVR	156.7	0.236	20.80
	TSVR	148.8	0.225	0.539
	Weighted TSVR	78.16	0.204	0.467
	DW-TSVR	69.45	0.210	0.425
Diabetes	SVR	0.572	0.216	0.150
	TSVR	0.552	0.205	0.072
	Weighted TSVR	0.533	0.198	0.063
	DW-TSVR	0.510	0.186	0.061
Pyrim	SVR	0.078	0.112	1.245
	TSVR	0.068	0.092	0.568
	Weighted TSVR	0.066	0.088	0.492
	DW-TSVR	0.066	0.089	0.488

度加权值,使模型能够更好地反映训练数据集的内在分布,从而获得了更加精确的训练模型.相比于标准SVR算法,在训练速度上有明显的提高,同时也验证了采用逐次超松弛技术解决TSVR算法求解问题的可行性.

## 5 结 论

本文提出了一种新的孪生支持向量回归机,称为密度加权孪生支持向量回归机(DW-TSVR).引入密度加权值的概念,将标准的孪生支持向量回归机优化目标从经验风险最小化扩展到结构风险最小化,并将密度加权值引入到新的目标函数中,所提出算法的对偶形式采用逐次超松弛算法进行求解. DW-TSVR算法与SVR、TSVR、Weighted TSVR在6个UCI数据集上的比较结果表明,所提出的算法在预测精确度上有一定的提高,验证了所提出算法的有效性.

## 参考文献(References)

- [1] Cortes C, Vapnik V. Support vector networks[J]. Machine Learning, 1995, 20(3): 273-297.
- [2] Vapnik V. The nature of statistical learning theory[M]. New York: Springer-Verlag, 1995: 93-110.
- [3] Jayadeva, Khemchandani R, Chandra S. Twin support vector machines for pattern classification[J]. IEEE Trans on Pattern Analysis and Machine Intelligence, 2007, 29(5): 905-910.
- [4] Mangasarian O L, Wild E W. Multisurface proximal support vector classification via generalized eigenvalues[J]. IEEE Trans on Pattern Analysis and Machine Intelligence, 2006, 28(1): 69-74.
- [5] Peng X. TSVR: An efficient twin support vector machine for regression[J]. Neural Networks, 2010, 23(3): 365-372.
- [6] Xu Yitian, Wang Laisheng. A weighted twin support vector regression[J]. Knowledge-Based Systems, 2012, 33: 92-101.
- [7] Shao Yuanhai, Zhang Chunhua, Wang Xiaobo, et al. Improvements on twin support vector machines[J]. IEEE Trans on Neural Networks, 2011, 20(6): 962-968.
- [8] Peng Xinjun, Xu Dong. A twin projection support vector machine for data regression[J]. Neurocomputing, 2014, 138: 131-141.
- [9] Ding Shifei, Wu Fulin, Shi Zhongzhi. Wavelet twin support vector machine[J]. Neural Computer & Application, 2014, 25(6): 1241-1247.
- [10] Shao Yuanhai, Zhang Chunhua. An  $\epsilon$ -twin support vector machine for regression[J]. Neural Computer & Application, 2013, 23(1): 175-185.
- [11] Li I, Chen J, Wu J. A fast prototype reduction method based on template reduction and visualization-induced self-organizing map for nearest neighbor algorithm[J]. Application Intelligence, 2013, 39(3): 564-582.
- [12] Olvi L Mangasarian, David R Musicant. Successive overrelaxation for support vector machines[J]. IEEE Trans on Neural Networks, 1999, 10(5): 1032-1037.
- [13] Luo Z Q, Tseng P. Error bounds and convergence analysis of feasible descent methods: A general approach[J]. Annals of Operations Research, 1993, 67(2): 157-178.

(责任编辑: 郑晓蕾)