

基于 GM(1,1)-SVM 组合模型的中长期人口预测研究

徐翔燕 侯瑞环

塔里木大学信息工程学院 新疆 阿拉尔 843300

(1036269601@qq.com)

摘要 准确预测未来人口数量,对制定相关经济政策具有现实意义。文中针对人口中长期预测影响因素较复杂、可用历史数据较少、单一模型局限性等特点,构建了灰色预测和支持向量机的组合预测模型。该模型将灰色预测模型和支持向量机模型进行组合,利用标准差法确定模型的权值信息,将模型应用于一师阿拉尔市人口的中长期预测,选取一师阿拉尔市 1997—2017 年的人口数据进行分析,对 2018—2022 年的数据进行预测。结果表明:与单一模型相比较,组合模型预测精度更高,相对误差低,且预测结果比较稳定,结果更符合实际。

关键词:灰色预测;支持向量机;组合模型

中图法分类号 C924

Medium and Long-term Population Prediction Based on GM(1,1)-SVM Combination Model

XU Xiang-yan and HOU Rui-huan

College of Information Engineering, Tarim University, Alar, Xinjiang 843300, China

Abstract Accurate prediction of future population is of practical significance for the formulation of relevant economic policies. In this paper, a combined prediction model of grey and support vector machine is constructed according to the characteristics of complicated influencing factors of medium and long-term prediction, less available historical data, and the limitations of single model. The model combines the grey prediction model with the support vector machine model and uses the standard deviation method to determine the weight information. The model is applied to the medium and long-term prediction of the population of Alar City, and the population data of the first division of Alar City from 1997 to 2017 is selected for analysis, to predict the data 2018 to 2022. The result shows that, compared with the single model, the combined model has higher prediction accuracy and lower relative error, and the prediction result is relatively stable and more realistic.

Keywords Grey prediction, Support vector machine, Combined model

科学合理的人口预测对区域经济发展和制定相关的社会政策具有重要的现实意义,因此,如何提高人口预测的准确性和精度一直是研究者关注的问题。目前建立的人口预测模型主要有线性回归模型、马尔萨斯模型、Logistic 模型、GM(1,1)模型和 BP 神经网络模型等。从应用角度上来说,线性回归模型需要人口数据变动平稳、直线趋势较明显,而随着经济社会的发展,人口数量很难呈现直线变动的情况;马尔萨斯模型对于人口基数小、增长速度快的区域预测效果较好,这就使得模型的适用区域有限;Logistic 模型实质是马尔萨斯模型的改进,对于时间较长人口数据的预测误差较大且不稳定;BP 神经网络模型会由于信息量的缺少导致预测精度不高;而 GM(1,1)模型从自身的时间序列中寻找有用信息建立模型,在缺少大量数据的时间序列时有很好的预测效果,然而,GM(1,1)模型过多依赖于历史数据,在数据离散程度较大时,预测精度不高。

针对 GM(1,1)模型预测的局限性,不同学者提出了不同的改进措施。郭雪峰等^[1]针对传统灰色模型在预测流动人口方面存在精度不高的缺陷,基于自适应滤波法对传统灰色模

型进行残差修正,预测结果表明改进后的模型预测精度更高,适用性与可行性更强;蒋若凡等^[2]利用 PSO-BP 神经网络组合算法优化灰色模型,建立多指标灰色 PSO-BP 神经网络人口预测模型,实证分析表明模型预测和外推精度高,具有较好的实用价值;龙会典等^[3]以灰色系统理论的 GM(1,1)模型和随机过程理论的 Markov 链模型为基础构建了一个动态 GM(1,1)-Markov 链组合预测模型,并将改模型用于预测广东省单位 GDP 能耗,预测效果较好;李凯等^[4]利用辛普森 3/8 公式和牛顿插值公式的组合插值方法改进传统的 GM(1,1)模型,实验结果表明改进后的模型具有较好的预测稳定性;吴文泽等^[5]在已有研究的基础上进一步改进 GM(1,1)模型,并将改进后的模型应用到旅游客流量的预测中,结果表明改进后的模型具有较好的预测性能;徐丽丽等^[6]针对单一模型的局限性构建了灰色预测和径向基网络的组合预测模型预测山东省人口总量,验证结果表明组合模型的预测精度较高。大量的文献表明,改进后的模型预测精度较传统灰色预测模型好。而支持向量机预测方法由于其结构简单、泛化能力好、预测精度高优点,在众多领域得到了广泛的应用。如解伟等^[7]基

基金项目:塔里木大学校长基金青年创新资金项目(TDZKQN201824)

This work was supported by Tarim University President Fund Young Innovation Fund Project(TDZKQN201824).

通信作者:侯瑞环(1073293432@qq.com)

于支持向量机预测省级电网投资规模;贾娜等^[8]利用支持向量机算法对金刚石锯片锯切木材表面的粗糙度进行预测;徐路路等^[9]基于支持向量机和改进粒子群算法预测科学前沿;宋晓华等^[10]基于改进 GM(1,1)和 SVM 的组合模型预测长期电量等。近年来,组合预测模型^[11-13]成为预测领域中的重要研究方向,组合模型能够克服单一模型局限性,提高模型预测精度。

因此,本文结合 GM(1,1)预测模型和支持向量机模型的特点,针对人口的中长期预测,提供了一种新的建模思路,将 GM(1,1)预测模型和支持向量机模型相结合,利用标准差法确定单一模型的权重系数,形成组合预测模型,并对人口进行中长期预测。

1 模型介绍

1.1 GM(1,1)模型

GM(1,1)模型是最常用的一种灰色预测模型,主要通过先对原始数据作累加生成具有指数增长规律的数据列,然后对生成的数据列构建微分方程模型,求微分方程的时间响应函数,最后累减还原得到预测模型。具体建模过程如下。

设原始非负序列为 $x^{(0)} = \{x^{(0)}(1), x^{(0)}(2), \dots, x^{(0)}(n)\}$, 对原始数据进行累加,得到生成序列 $x^{(1)} = \{x^{(1)}(1), x^{(1)}(2), \dots, x^{(1)}(n)\}$, 其中 $x^{(1)}(k) = \sum_{i=1}^k x^{(0)}(i)$ 。

累加后的生成序列与原始序列之间满足灰微分方程:

$$x^{(0)}(k) + az^{(1)}(k) = \mu \quad (1)$$

其中, a 为发展灰数, μ 为内生控制灰数, $z^{(1)}(k) = 0.5(x^{(1)}(k) + x^{(1)}(k-1))$ 。

利用最小二乘法对参数进行估计,得到:

$$[a, \mu]^T = (B^T B)^{-1} B^T Y \quad (2)$$

其中,

$$B = \begin{bmatrix} -z^{(1)}(2) & 1 \\ -z^{(1)}(3) & 1 \\ \vdots & \vdots \\ -z^{(1)}(n) & 1 \end{bmatrix}, Y = \begin{bmatrix} x^{(0)}(2) \\ x^{(0)}(3) \\ \vdots \\ x^{(0)}(n) \end{bmatrix}$$

生成序列 $x^{(1)}$ 满足 GM(1,1)模型的白化微分方程:

$$\frac{dx^{(1)}}{dt} + ax^{(1)} = \mu \quad (3)$$

求解上述白化微分方程得到 GM(1,1)模型的时间响应函数:

$$\hat{x}^{(1)}(k) = (x^{(0)}(1) - \frac{\mu}{a})e^{-a(k-1)} + \frac{\mu}{a} \quad (4)$$

还原得预测模型为:

$$\begin{cases} \hat{x}^{(0)}(1) = x^{(0)}(1) \\ \hat{x}^{(0)}(k) = (1 - e^{-a})(x^{(0)}(1) - \frac{\mu}{a})e^{-a(k-1)} + \frac{\mu}{a}, k=2, 3, \dots, n \end{cases} \quad (5)$$

1.2 SVM 模型

支持向量机是由 Vapnik 于 1995 年首次提出,被广泛应用于多个领域来解决分类和预测问题的机器学习方法。支持向量回归(SVR)的主要思路是将低纬度空间函数通过非线性投影转化为高纬空间的线性投影,较好地解决了低纬度空间使用线性回归的困难。

假设训练样本 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ 。

$x_m \in R^n$ 为对应的输入向量, $y_m \in R$ 为输出向量, $m=1, 2, \dots, n$ 为观测值总数。回归方程为: $f(x) = \omega^T x + b$, ω 为超平面权重向量, b 为偏差向量。于是 SVR 问题可转化为:

$$\min_{\omega, b} \left\{ \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \phi_\epsilon(f(x_i) - y_i) \right\} \quad (6)$$

其中, C 为正则化常数, $\phi_\epsilon(\cdot)$ 为 ϵ 不敏感损失函数, 表达式为:

$$\phi_\epsilon(v) = \begin{cases} 0, & |v| \leq \epsilon \\ |v| - \epsilon, & |v| > \epsilon \end{cases}$$

引入松弛变量 ξ_i 与 $\hat{\xi}_i$, 重构式(6)最小化问题,即:

$$\begin{aligned} \min_{\omega, b, \xi_i, \hat{\xi}_i} & \left\{ \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m (\hat{\xi}_i - \xi_i) \right\} \\ \text{s. t.} & \begin{cases} f(x_i) - y_i \leq \epsilon + \xi_i \\ y_i - f(x_i) \leq \epsilon + \hat{\xi}_i \\ \hat{\xi}_i \geq 0, \xi_i \geq 0 \end{cases} \end{aligned} \quad (7)$$

引入 4 个非负拉格朗日乘子 $\mu_i, \hat{\mu}_i, \alpha_i, \hat{\alpha}_i$, 由拉朗日乘子法可以得到式(7)的拉格朗日函数为:

$$\begin{aligned} L = & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m (\hat{\xi}_i - \xi_i) - \sum_{i=1}^m \mu_i \xi_i - \\ & \sum_{i=1}^m \hat{\mu}_i \hat{\xi}_i + \sum_{i=1}^m \alpha_i (f(x_i) - y_i - \epsilon - \xi_i) + \sum_{i=1}^m \hat{\alpha}_i (y_i - \epsilon - \hat{\xi}_i) \end{aligned} \quad (8)$$

通过对式(8)求关于参数 $\omega, b, \xi_i, \hat{\xi}_i$ 的偏导数,令其为零,可得 SVR 的对偶问题:

$$\begin{aligned} \max_{\alpha_i, \hat{\alpha}_i} & \left\{ \sum_{i=1}^m y_i (\hat{\alpha}_i - \alpha_i) - \omega^T (\hat{\alpha}_i - \alpha_i) - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\hat{\alpha}_i - \alpha_i) \right. \\ & \left. (\hat{\alpha}_j - \alpha_j) x_i^T x_j \right\} \\ \text{s. t.} & \begin{cases} \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) = 0 \\ 0 \leq \hat{\alpha}_i, \alpha_i \leq C \end{cases} \end{aligned} \quad (9)$$

所得回归函数为:

$$f(x) = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) x_i^T x + b \quad (10)$$

将式(9)、式(10)中的 x_i^T 改成核函数 $\kappa(x, x_i)$, 则可得到回归函数:

$$f(x) = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \kappa(x, x_i) + b \quad (11)$$

2 组合预测模型构建

组合预测模型是一种将不同预测模型的预测结果,选取适当的权数进行加权组合的一种预测模型。其主要目的是最大限度地综合各个模型的预测结果,提高预测精度。组合模型能更系统、更科学地克服单一模型的局限性,从而提高预测精度。本文选取 GM(1,1)和 SVM 两种预测模型进行组合,其基本原理如下。

设 x_i 为第 i 年的实际人口数 ($i=1, 2, \dots, n$, n 为预测人口的年数), 则由 n 年人口的实际值得时间序列 $(x_i)_{1 \times n}$ 。设 \tilde{x}_{ik} 为第 k 种方法第 i 年的人口预测值 ($k=1, 2, \dots, K$), $e_{ik} = x_i - \tilde{x}_{ik}$ 为第 k 种方法第 i 年的人口预测值的预测误差, ω_k 为第 k 种方法权系数的估计值, \hat{x}_i 为组合预测值, 则有:

$$\hat{x}_i = \sum_{k=1}^K \omega_k \tilde{x}_{ik} \quad (12)$$

组合预测模型权重确立是以预测误差平方和达到最小为准则,具体公式如下:

$$\min E = \sum_{i=1}^n (x_i - \sum_{k=1}^K \omega_k \tilde{x}_{ik})^2$$

$$\text{s. t. } \sum_{k=1}^K \omega_k = 1, \omega_k \geq 0, k=1, 2, \dots, K$$

目前国内外对组合预测模型的研究较多,比较常用的有平均值法、方差-协方差优选组合预测法、最小二乘法、标准差法等。本文选取简单而有效的标准差法得到组合模型的权重,从而利用组合模型对人口进行预测。即采用标准差法得到预测模型的权重 ω_k ,然后将权系数 ω_k 代入预测模型求得预测结果。具体流程如图 1 所示。

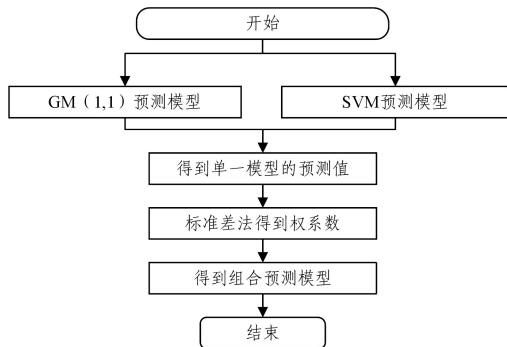


图 1 组合预测模型流程图

Fig. 1 Flow chart of combined prediction model

3 人口预测

《新疆生产建设兵团十三五规划纲要》明确指出,阿拉尔市在十三五期间需要增加人口十万,任务重,责任大。如今十三五规划执行时间已经过半,阿拉尔市若要完成这一指标,就需要较为精确地预测未来人口数量,从而制定相应的人口聚集政策。所以本文采用所构建的 GM(1,1) 和 SVM 组合预测模型,以 1997—2017 年阿拉尔市总人口为样本进行人口预测分析,原始数据(单位:万人)来源于《阿拉尔统计年鉴》与《新疆生产建设兵团统计年鉴》,具体数据如表 1 所列。

表 1 阿拉尔人口数量
Table 1 Population of Alar

年份	1997	1998	1999	2000	2001	2002	2003
人口数量	25.2434	26.4798	25.9934	26.7648	26.6207	27.5494	27.9396
年份	2004	2005	2006	2007	2008	2009	2010
人口数量	28.3249	28.6920	28.9359	29.1388	29.1193	29.2109	29.2305
年份	2011	2012	2013	2014	2015	2016	2017
人口数量	29.6946	29.8113	31.0018	31.0130	31.8009	32.6845	35.7961

采用相对误差检验模型的有效性,其具体公式如下:

$$\alpha_i = \frac{|x_i - \tilde{x}_i|}{x_i}$$

其中, x_i 为第 i 年的实际人口数, \tilde{x}_i 为第 i 年的预测人口数。通过标准误差计算可以得出 GM(1,1) 模型与 SVM 模型的组合权重分别为: 0.4808 和 0.5192, 至此, 可以得到具体模型为:

$$\hat{x}_i = 0.4808 \tilde{x}_{GM(1,1)} + 0.5192 \tilde{x}_{SVM} \quad (13)$$

分别使用传统 GM(1,1) 模型、SVM 模型以及组合模型式(13)对阿拉尔近 21 年人口数据做预测分析, 预测结果如表 2 所列。

表 2 预测结果分析表

Table 2 Analysis of forecast results

年份	GM(1,1) 模型		SVM 模型		组合模型	
	预测值	相对误差/%	预测值	相对误差/%	预测值	相对误差/%
1997	25.24	0.00	25.45	0.80	25.35	0.42
1998	25.87	2.30	26.62	0.53	26.26	0.83
1999	26.20	0.80	26.19	0.74	26.19	0.77
2000	26.54	0.84	26.92	0.59	26.74	0.10
2001	26.88	0.99	26.84	0.83	26.86	0.90
2002	27.23	1.16	27.72	0.63	27.49	0.23
2003	27.58	1.29	27.90	0.13	27.75	0.68
2004	27.94	1.38	28.30	0.07	28.13	0.70
2005	28.29	1.38	28.50	0.67	28.40	1.01
2006	28.66	0.96	28.81	0.44	28.74	0.69
2007	29.03	0.38	29.36	0.75	29.20	0.21
2008	29.40	0.97	29.29	0.60	29.35	0.78
2009	29.78	1.95	29.01	0.69	29.38	0.58
2010	30.17	3.20	29.61	1.29	29.09	0.48
2011	30.55	2.89	29.50	0.65	30.01	1.06
2012	30.95	3.81	29.65	0.53	30.28	1.56
2013	31.35	1.11	30.76	0.78	31.04	0.13
2014	31.75	2.38	30.46	1.78	31.08	0.22
2015	32.16	1.13	31.55	0.79	31.84	0.13
2016	32.57	0.34	32.38	0.93	32.47	0.65
2017	32.99	7.83	37.17	3.86	35.16	1.76
平均						
相对误差/%		1.85		0.91		0.69

通过表 2 的结果可知, 传统 GM(1,1) 模型预测相对误差最大值为 7.83%, 其平均相对误差高于 SVM 模型, 相比较而言, SVM 有更强的预测能力。组合模型的最大相对误差为 1.74%, 均小于单一的 GM(1,1) 模型的 7.83% 和 SVM 模型的 4.67%, 组合模型的平均相对误差最小, 为 0.69%。由此可见, 组合模型的预测精度比单一模型的精度高, 同时相对误差最小。

分别使用传统 GM(1,1) 模型、SVM 模型和组合模型预测阿拉尔市 2018—2022 年人口数量, 预测结果表 3 所列。

表 3 阿拉尔市 2018—2022 年人口预测结果

Table 3 Forecast results of population in Alar in 2018—2022

年份	GM(1,1) 模型 预测值	SVM 模型 预测值	组合模型
2018	33.42	37.19	35.38
2019	33.85	37.57	35.78
2020	34.28	38.16	36.29
2021	34.73	38.44	36.66
2022	35.17	38.87	37.09

已公布的 2018 年人口数据为 37.2 万人, 考虑到 2018 年阿拉尔市在兵团向南发展战略需求下, 实行向西部贫困山区招工落户政策, 使得迁入人口在 2018 年后半年增加 2 万余人, 自然增长下的人口数量为 35 万余人, 可以看出, 组合模型的实际预测精度更高, 更加符合实际情况。

结束语 人口数量问题在社会各项事业发展过程中起到关键作用, 人口政策的制定将会对社会发展影响持续几十年, 甚至上百年, 精确预测未来人口数量意义重大。人口数量变化受众多影响因素作用, 无法找到任何一种预测模型进行精确预测, 本文有效地利用了 GM(1,1) 模型在预测中计算简单、使用容易的优点, 结合 SVM 模型在中长期预测中比较稳健的特点, 根据误差标准构造相应权重因子, 得到抗干扰能力更强、预测精度较高的组合模型, 为中长期人口预测提供一种可行性较高的预测方法。

- dering imbalanced data distribution[J]. IEEE Transactions on Knowledge & Data Engineering, 2005(6):786-795.
- [13] BREIMAN L. Bagging predictors [J]. Machine Learning, 1996, 24(2):123-140.
- [14] ZAREAPOOR M, SHAMSOLMOALI P. Application of credit card fraud detection: Based on bagging ensemble classifier[J]. Procedia computer science, 2015, 48(2015):679-685.
- [15] WITTEN I H, FRANK E, HALL M A, et al. Data Mining: Practical machine learning tools and techniques [M]. Morgan Kaufmann, 2016:70-71.
- [16] 韩家伟, 坎伯. 数据挖掘: 概念与技术 [M]. 北京: 机械工业出版社, 2012:158-159.
- [17] DEORA C S, ARORA S, MAKANI Z. Comparison of Interestingness Measures: Support-Confidence Framework versus Lift-Rule Framework[J]. International Journal of Engineering Research & Applications, 2014, 3(2):208-215.
- [18] ALCALÁ-FDEZ J, FERNÁNDEZ A, LUENGO J, et al. KEEL Data-Mining Software Tool: Data Set Repository, Integration of Algorithms and Experimental Analysis Framework[J]. Journal of Multiple-Valued Logic & Soft Computing, 2011, 17:255-287.
- [19] PATIL T R, SHEREKAR S. Performance analysis of Naive Bayes and J48 classification algorithm for data classification[J]. International Journal of Computer Science and Applications, 2013, 6(2):256-261.
- [20] QUINLAN J R. Bagging, boosting, and C4.5 [C]// AAAI/IAAI. 1996:725-730.
- [21] LOBO J M, JIMÉNEZ-VALVERDE A, REAL R. AUC: a misleading measure of the performance of predictive distribution models[J]. Global Ecology and Biogeography, 2008, 17(2):145-151.
- [22] DAVIS J, GOADRICH M. The relationship between Precision-Recall and ROC curves [C]// Proceedings of the 23rd International Conference on Machine Learning. ACM, 2006:233-240.
- [23] POWERS D M. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation[J]. Journal of Machine Learning Technology, 2011, 2(1):37-63.
- [24] WILCOXON F, KATTI S, WILCOX R A. Critical values and probability levels for the Wil-coxon rank sum test and the Wil-coxon signed rank test[J]. Selected Tables in Mathematical Statistics, 1970, 1:171-259.



CUI Wei, born in 1994, postgraduate. His main research interests include machine learning and data mining.



ZHU Xiao-yan, born in 1983, Ph.D, associate professor, is a member of China Computer Federation. Her main research interests include machine learning and data mining.

(上接第 487 页)

将 GM(1,1)与 SVM 组合模型运用到阿拉尔市人口数量预测研究中,得到阿拉尔市未来五年的人口数量分别为 35.38 万人、35.78 万人、36.29 万人、36.66 万人和 37.09 万人。由于《新疆生产建设兵团十三五规划》中关于各师市人口聚集的有相关政策的要求,阿拉尔市在十三五期间需向西部贫困山区招工落户 10 万人,2018 年招工人数为 2 万余人,从而导致 2018 年预测结果与阿拉尔市公布的人口数据之间存在差距。因此,阿拉尔市实际人口数量与预测值之间相差 2 万余人,5 年之后的 2022 年阿拉尔市实际人口数量约为 40 万余人,达到了预期的人口聚集目标。

参考文献

- [1] 郭雪峰,黄健元,王欢.改进的灰色模型在流动人口预测中的应用[J].统计与决策,2018(8):76-79.
- [2] 蒋若凡,姜玉梅,李菲雅.基于灰色 PSO-BP 人口预测模型的研究与应用[J].西北人口,2011,32(3):23-26.
- [3] 龙会典,严广乐.基于改进的 GM(1,1)—Markov 链组合模型广东省单位 GDP 能耗预测[J].数理统计与管理,2017,36(2):200-207.
- [4] 李凯,张涛.基于组合插值的 GM(1,1)模型背景值的改进[J].计算机应用研究,2018,35(10):2994-2999.
- [5] 吴文泽,张涛. GM(1,1)模型的改进及应用[J].统计与决策,2019(9):15-18.
- [6] 徐丽丽,李洪,李劲.基于灰色预测和径向基网络的人口预测研究[J].计算机科学,2019,46(Z1):431-435.
- [7] 解伟,潘文明,王成化,等.基于支持向量机的省级电网中长期投资规模预测模型研究[J].工业技术经济,2019(8):154-160.
- [8] 贾娜,郭佳欣,花军,等.采用支持向量机算法对金刚石锯片锯切木材表面粗糙度的预测[J].东北林业大学学报,2019,47(10):85-89.
- [9] 徐路路,王芳.基于支持向量机和改进粒子群算法的科学前沿预测模型研究[J].情报科学,2019,37(8):22-28.
- [10] 宋晓华,祖丕娥,伊静,等.基于改进 GM(1,1)和 SVM 的长期电量优化组合预测模型[J].中南大学学报(自然科学版),2012,43(5):1803-1807.
- [11] BATES J M, GRANGER C W. Combination of Forecasts[J]. Operational Res-Ouart, 1969, 20(4):451-468.
- [12] LIU S L, HU Z Q, CHI X K. The research of power load forecasting method on combination forecasting model[J]. Information Science and Engineering, 2010(26).
- [13] 吴静敏,左洪福,陈勇.基于免疫粒子群算法的组合预测方法[J].系统工程理论方法应用,2006,15(3):229-233.



XU Xiang-yan, born in 1990, master, lecturer. Her main research interests include intelligent optimization algorithm and so on.



HOU Rui-huan, born in 1986, master, lecturer. His main research interests include nonparametric statistics and so on.