

# 基于 XGBoost 算法的房价预测模型

王冬雪, 郭秀娟\*

(吉林建筑大学电气与计算机学院, 吉林 长春 130118)

**摘要:**随着人工智能技术的快速发展,机器学习被广泛应用于各行各业,并在多个领域都取得了较好的成果。房价是一个影响因素复杂的热点问题,难以对其做出全面准确的预测。因此,本文尝试在一个相对稳定条件下将相关的机器学习算法应用于房价预测中。本文首先对数据进行缺失值处理、样本的因变量相关分析及标准化处理等一系列预处理,然后结合互联网数据和机器学习中的 XGBoost 算法对数据集进行建模和训练,最终实现对房价的精准预测。

**关键词:**房价预测;机器学习;XGBoost

中图分类号:F832.4

文献标识码:A

文章编号:2096-2118(2021)03-0079-04

## Housing Price Prediction Model Based on XGBoost

WANG Dongxue, GUO Xiujuan\*

(School of Electrical and Computer Engineering, Jilin Jianzhu University, Changchun, Jilin 130118, China)

**Abstract:** With the rapid development of artificial intelligence technology, machine learning has been widely used in all walks of life and has achieved good results in many fields. Housing price is a hot issue with complicated influencing factors, so it is difficult to make a comprehensive and accurate forecast. Therefore, this paper attempts to apply the relevant machine learning algorithm to the housing price prediction under a relatively stable condition. In this paper, a series of preprocessing such as missing value processing, dependent variable correlation analysis of samples and standardization processing are carried out on the data at first. Then, the data set is modeled and trained by combining internet data and XGBoost algorithm in machine learning, and the accurate prediction of housing price is finally realized.

**Keywords:** housing price prediction; machine learning; XGBoost

### 0 引言

近年来,人们对房价的关注一直居高不下。由于城市化的加剧,对房屋租赁和购房的需求也持续增加,而房价问题不仅关系到人民生活水平,更是与国民经济发展息息相关。因此,对房价进行预测不仅对人们买卖房屋具有参考意义,而且对于政府进行房价调控也有积极作用<sup>[1]</sup>。因而确定一种可以精准反映房价走势的算法具有重要意义。

本文通过使用 XGBoost 算法来预测房价。通过对数据分析、预处理及基于 XGBoost 模型来构建房

价预测模型。影响房价的因素多且复杂,如房屋面积、房屋地理位置、房屋户型等,本文仅选取对于房价影响较大的 79 个特征对房价预测模型进行评估,并选择 RMSLE 算法作为预测房价的评估算法。

### 1 数据预处理

在实际情况下,由于环境复杂等因素,我们获取的数据往往是存在缺失和异常的,因此,在建模前要对数据进行预处理。

#### 1.1 数据集来源

收稿日期:2021-03-11

作者简介:王冬雪(1995-),女,吉林省长春市人,在读硕士研究生,研究方向:电气工程。

\* 通讯作者:郭秀娟(1961-),女,吉林省长春市人,教授,博士,研究方向:计算机图形图像处理, E-mail:779523836@qq.com

(C)1994-2021 China Academic Journal Electronic Publishing House. All rights reserved. <http://www.cnki.net>

## 1.2 缺失值处理

— 正态分布( $\mu=12.02$  且  $\sigma=0.40$ )

Term (词)	Frequency (频次)
汶川地震	100
抗震质量	96
人口特征	94
人口类型	92
围护质量	81
震害程度	49
房屋倒塌面积	17
车辆状况	5
车质重量	5
车辆超载日期	5
毁损程度	5
车辆位置	4
地下室情况	3
吃回锅冷饭	2
地下室高度	2
成品面积	2
施工面积	1
施工单模型	1
装饰材料	1
城市分区类型	1
全浴套数	1
地域类型	1
销售方式	1
销售条件	1
未竣工面积	1
门庭建筑面积	1
顺风走廊面积	1
已出租年份	1
已售年份	1
美元价格	1

图 1 缺失数据统计图

### 1.3 样本因自变量相关分析

(a) 单价分布情况

	出售价格	房屋整体材料	高档居住面积	车库容量大小	车库面积	地下室总面积	1 型成品面积	地下全浴室	总房间	建成年份
出售价格	1.00	0.79	0.71	0.64	0.62	0.61	0.61	0.56	0.53	0.52
房屋整体材料	0.79	1.00	0.59	0.60	0.56	0.54	0.45	0.55	0.43	0.57
高档居住面积	0.71	0.59	1.00	0.47	0.47	0.45	0.57	0.63	0.83	0.20
车库容量大小	0.64	0.60	0.47	1.00	0.88	0.43	0.44	0.47	0.36	0.54
车库面积	0.62	0.56	0.47	0.88	1.00	0.49	0.49	0.41	0.34	0.48
地下室总面积	0.61	0.54	0.45	0.43	0.49	1.00	0.82	0.32	0.29	0.39
1 型成品面积	0.61	0.48	0.57	0.44	0.49	0.82	1.00	0.38	0.41	0.28
地下全浴室	0.56	0.55	0.63	0.47	0.41	0.32	0.38	1.00	0.55	0.47
总房间	0.53	0.43	0.83	0.38	0.34	0.29	0.41	0.55	1.00	0.10
建成年份	0.52	0.57	0.20	0.54	0.48	0.39	0.28	0.47	0.10	1.00

图 2 相关矩阵热力图

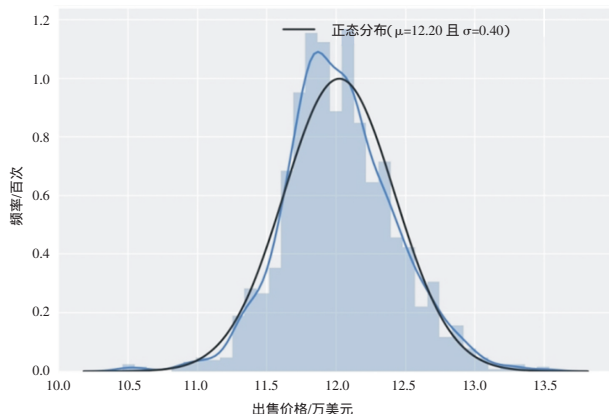
### 1.4 数据标准化处理

Box-Cox 变换的一般形式为：

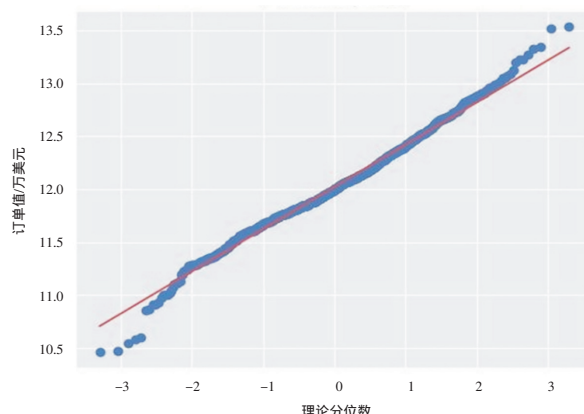
$$y(\lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln y & \lambda = 0 \end{cases} \quad (1)$$

(b) 单价分布的 P-P 图

图3 原始数据分布



(a) 变换后的单价分布情况



(b) 变换后的 P-P 图

图4 变换后数据分布

## 2 XGBoost 模型

XGBoost 又称极端梯度上升,它是大规模并行 Boosted Tree,是 Gradient Boosting Machine 的扩展,在相同的环境和条件下,XGBoost 比同类算法快 10 倍以上<sup>[5]</sup>。XGBoost 还可以通过分布式运算,进一步提高训练速度<sup>[6]</sup>。

### 2.1 基本模型

XGBoost 是由  $k$  个基模型组成的加法运算式:

$$\hat{y}_i = \sum_{t=1}^k f_t(x_i) \quad (2)$$

其中  $f_t$  为  $k$  个基模型  $\hat{y}_i$  为第  $i$  个样本的预测值。

损失函数可由预测值  $\hat{y}_i$  与真实值  $y_i$  进行表示:

$$loss = \sum_{i=1}^n l(y_i, \hat{y}_i) \quad (3)$$

其中  $n$  为样本数量。

模型的偏差和方差共同决定了该模型的预测精度,模型的偏差具体表现为损失函数,模型越简单则其方差越小,所以目标函数由模型的损失函数  $loss$  与抑制模型复杂度的正则项  $\Omega$  组成,所以目标函数可表示为:

$$obj = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{t=1}^k \Omega(f_t) \quad (4)$$

其中  $\Omega$  为模型的正则项。

以第  $t$  步的模型为例,模型对第  $i$  个样本  $x_i$  的预测为:

$$\hat{y}_i = \hat{y}_i^{t-1} + f_t(x_i) \quad (5)$$

其中  $\hat{y}_i^{t-1}$  是由第  $t-1$  步的模型给出的预测值  $f_t(x_i)$  是需要加入的新模型的预测值,则目标函数可写为:

$$obj^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1}) + \sum_{i=1}^n \Omega(f_t) = \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1} +$$

$$f_t(x_i)) + \sum_{i=1}^n \Omega(f_t) \quad (6)$$

而根据泰勒公式,可以把上述目标函数写为:

$$obj^{(t)} = \sum_{i=1}^n [l(y_i, \hat{y}_i^{t-1}) + g_t f_t(x_i) + \frac{1}{2} h_t f_t^2(x_i)] +$$

$$\sum_{i=1}^n \Omega(f_t) \quad (7)$$

其中  $g_t$  为损失函数的一阶导  $h_t$  为损失函数的二阶导。

将决策树定义为  $f_t(x) = w_{q(x)}$   $x$  为某一样本  $q(x)$  代表该样本所在的叶子结点,而  $w_q$  则代表叶子结点取值  $w$ ,所以  $w_{q(x)}$  代表每个样本的取值  $w$  (即预测值),则目标函数的正则项可以定义为<sup>[7]</sup>:

$$\Omega(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \quad (8)$$

其中  $\lambda$  和  $\gamma$  是 XGBoost 定义的,其值可设定,值越大,表示越希望获得结构简单的树,  $T$  为叶子数。

再定义  $G_j = \sum_{i \in I_j} g_i$   $H_j = \sum_{i \in I_j} h_i$ , 则目标函数可记

为:

$$obj^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T \quad (9)$$

而其中叶子结点  $j$  对应的权值可表示为:

$$w_j^* = -\frac{G_j}{H_j + \lambda} \quad (10)$$

所以目标函数可简化为:

$$obj^{(t)} = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T \quad (11)$$

记  $I_L, I_R$  分别是数据集的左右结点, 其中  $I = I_L \cup I_R$ , 则分裂后增益为:

$$L_{split} = Gain = \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} + \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \lambda \quad (12)$$

其中,  $\frac{G_L^2}{H_L + \lambda}$ ,  $\frac{G_R^2}{H_R + \lambda}$  分别为左右树的分数,

$\frac{(G_L + G_R)^2}{H_L + H_R + \lambda}$  为分裂前分数,  $\lambda$  为惩罚项<sup>[8]</sup>。

XGBoost 在构建树的节点时, 为每个节点添加了一个缺省方向, 当样本缺失对应特征时, 就会被归类到缺省方向上。如果样本存在特征缺失的情况, 则只需分配到左右节点而无需遍历, 故算法所需遍历的样本量大大减少。稀疏感知算法比 basic 算法速度快了超过 50 倍<sup>[9-10]</sup>。

## 2.2 模型评价

本文采用均方根对数误差(RMSLE)来作为模型评价的标准。其公式如下:

$$RMSLE = \sqrt{\frac{1}{n} \sum_{i=1}^n [\log(\hat{y}_i + 1) - \log(y_i + 1)]^2} \quad (13)$$

其中预测值  $\hat{y}_i$  为预测值,  $y_i$  为真实值,  $n$  为样本数量。

对训练集训练 100 次后的 RMSLE 为 0.041 646 875 398 8, 如图 5 所示。XGBoost 能更好地适应不平衡的数据集, 同时也更不容易过拟合, 泛化能力较好, 应用范围广泛, 因此该模型基本可以实现对房价的精准预测。对于该预测模型, 可应用到以下场景。

1) 链家、安居客等二手房的交易。该模型更加有利于买卖双方看清房价接下来的走势, 及时把握住期望成交价格。

2) 房产投资的应用。近年来, 从《新中产白皮书》中可以看到, 新中产人群, 除去自住房, 投资性房地产占比是最多的。因此, 该模型对于投资者有一定的指导性作用。

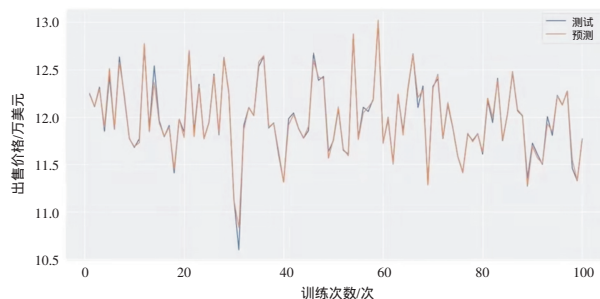


图5 预测结果图

## 3 结论

房价预测问题本质上来说就是典型的回归问题。本文基于 XGBoost 算法进行房价预测, 首先对数据进行缺失值处理、相关分析及标准化处理等一系列预处理, 然后使用 XGBoost 算法对数据集进行建模和训练, 最终实现对房价的精准预测。

## 参考文献

- [1] 张家棋, 杜金. 基于 XGBoost 与多种机器学习方法的房价预测模型[J]. 现代信息科技, 2020, 4(4): 15-18.
- [2] 汪佳琦. 基于机器学习方法的美国房屋价格影响因素分析[D]. 昆明: 云南大学, 2019.
- [3] 陈绵旺. 基于 RS-SVM 的商品住宅价格预测研究[D]. 南昌: 华东交通大学, 2016.
- [4] 高建, 周丽萍. 基于 Box-Cox 变换的住宅特征价格理论研究[J]. 河北科技大学学报, 2007(3): 247-250.
- [5] 曾婷婷. 基于机器学习的房价预测模型研究[D]. 绵阳: 西南科技大学, 2020.
- [6] Byeonghwa Park, Jae Kwon Bae. Using machine learning algorithms for housing price prediction: The case of Fairfax County, Virginia housing data[J]. 2015, 42(6): 2928-2934.
- [7] 龚洪亮. 基于 XGBoost 算法的武汉市二手房价格预测模型的实证研究[D]. 武汉: 华中师范大学, 2018.
- [8] FLOREZ-LOPEZ R, RAMON-JERONIMO J M. Enhancing accuracy and interpretability of ensemble strategies in credit risk assessment. A correlated-adjusted decision forest proposal [J]. Expert Systems with Applications, 2015, 42(13): 5737-5753.
- [9] 叶倩怡. 基于 XGBoost 方法的实体零售业销售额预测研究[D]. 南昌: 南昌大学, 2016.
- [10] Tianqi Chen, Carlos Guestrin. XGBoost: A Scalable Tree Boosting System[J]. IJCE Transactions on Fundamentals of Electronics, Communications and Computer Sciences, 2016, 53(7): 4903-4919.

编辑: 杨洋