

# 非参数回归模型样条估计量的分布

詹陆丽\*, 武新乾

河南科技大学数学与统计学院, 河南 洛阳

收稿日期: 2022年9月24日; 录用日期: 2022年10月17日; 发布日期: 2022年10月26日

## 摘要

为探究非参数回归模型中非参数函数估计量的分布, 本文在标准正态误差情形下首先得到了均值函数样条估计量的正态分布, 然后得到了方差函数基于残差的样条估计量的渐近分布, 并采用单个卡方变量线性函数来近似方差函数估计量的渐近分布。通过数值模拟验证了均值函数估计量的分布和方差函数估计量的渐近分布。

## 关键词

非参数回归模型, 样条估计, 渐近分布, 卡方分布线性组合

# Distribution of Spline Estimators for Nonparametric Regression Models

Luli Zhan\*, Xinqian Wu

School of Mathematics and Statistics, Henan University of Science and Technology, Luoyang Henan

Received: Sep. 24<sup>th</sup>, 2022; accepted: Oct. 17<sup>th</sup>, 2022; published: Oct. 26<sup>th</sup>, 2022

## Abstract

To explore the distribution of nonparametric function estimator in the nonparametric regression model, we first obtain the normal distribution of the spline estimator of the mean function in the standard normal error case in this paper. Then the asymptotic distribution of the spline estimator of the variance function based on the residuals is obtained. And the linear function of individual chi-square variable is used to approximate the asymptotic distribution of the variance function estimator. The distribution of the mean function estimator and the asymptotic distribution of the variance function estimator are illustrated by numerical simulations.

\*通讯作者。

## Keywords

Nonparametric Regression Model, Spline Estimation, Asymptotic Distribution, Linear Combination of Chi-Square Distribution

Copyright © 2022 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

考虑如下非参数回归模型

$$y_t = m(x_t) + \sigma(x_t)\varepsilon_t, \quad t = 1, 2, \dots, n, \quad (1)$$

其中  $x_t = t/n$  为固定设计点,  $y_t$  为被解释变量或响应变量,  $m(\cdot)$  与  $\sigma(\cdot)$  分别为未知的均值函数和标准差函数,  $\{\varepsilon_t\}$  为独立同分布的随机误差序列, 且  $E\varepsilon_t = 0$ ,  $E\varepsilon_t^2 = 1$ ,  $t = 1, 2, \dots, n$ 。

近些年来, 国内外学者对非参数回归模型做了大量的研究[1]-[7]。渐近分布作为估计量大样本性质中的一个重要方面, 也引起了一些学者的兴趣。秦永松(1991) [8]基于加权核估计得到了均值函数的估计量, 并研究了其导数的渐近分布; Liang 与 Jing (2004) [9]采用加权核估计的方法, 基于负相关序列研究了非参数回归模型中未知函数估计量的逐点一致收敛性与渐近正态性; Jin 等(2014) [10]研究了一步 Newton-Raphson 估计和局部轮廓似然估计, 并给出了基于两种方法的估计量的分布; Alharbi 与 Patili (2018) [11]提出对响应变量和残差的乘积进行平滑处理, 并研究了基于此方法得到的方差函数估计量的渐近分布; Li 和 Lin (2020) [12]在没有独立性假设的情形下, 推导出了误差方差  $\sigma^2 = \text{var}(\varepsilon)$  的最佳半参数效率约束, 并建立了基于残差的  $\sigma^2$  有效估计量的渐近正态性。

本文对非参数回归模型中基于样条方法的均值函数估计量和方差函数估计量的分布问题进行研究, 并通过数值模拟验证效果。

## 2. 估计量及主要结论

根据文献[7], 将区间  $D = [0, 1]$  进行包括两端点的  $k+1$  等距分割, 结点序列为

$$0 = t_0 < t_1 < \dots < t_{k+1} = 1,$$

构造相应的  $\nu$  次样条空间  $S_{k,\nu}$ , 其基函数记作  $B_s(x)$  ( $s = 1, 2, \dots, k + \nu$ )。设  $K = k + \nu$ , 又令

$$B_s(x) = (B_1(x), B_2(x), \dots, B_K(x))^T,$$

则均值函数  $m(x)$  的样条估计为

$$\hat{m}(x) = \mathbf{B}^T(x) \hat{\boldsymbol{\phi}} = \sum_{s=1}^K \hat{\phi}_s B_s(x), \quad (2)$$

这里  $\hat{\boldsymbol{\phi}} = (\hat{\phi}_1, \dots, \hat{\phi}_K)^T$  为能使

$$l_m(\boldsymbol{\phi}) = \sum_{t=1}^n \left\{ y_t - \sum_{s=1}^K \phi_s B_s(x_t) \right\}^2$$

最小化的参数向量,  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_K)^T$ , 它的最小二乘估计为  $\hat{\boldsymbol{\phi}} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{y}$ , 其中

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^T,$$

$$\mathbf{W} = \begin{pmatrix} B_{k1}(x_1) & B_{k2}(x_1) & \dots & B_{kk}(x_1) \\ B_{k1}(x_2) & B_{k2}(x_2) & \dots & B_{kk}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ B_{k1}(x_n) & B_{k2}(x_n) & \dots & B_{kk}(x_n) \end{pmatrix}.$$

记残差  $\hat{\zeta}_t = y_t - \hat{m}(x_t), t = 1, \dots, n$ , 令  $z_t = \hat{\zeta}_t^2, \mathbf{z} = (z_1, \dots, z_n)^T$ , 则方差函数  $\sigma^2(x)$  的基于残差的样条估计为

$$\hat{\sigma}^2(x) = \mathbf{B}^T(x) \hat{\boldsymbol{\theta}} = \sum_{s=1}^K \hat{\theta}_s B_s(x), \tag{3}$$

其中  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_K)^T$  为使得

$$l_{\sigma^2}(\boldsymbol{\theta}) = \sum_{t=1}^n \left\{ z_t - \sum_{s=1}^K \theta_s B_s(x_t) \right\}^2$$

最小化的参数向量,  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^T$ , 它的最小二乘估计为  $\hat{\boldsymbol{\theta}} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{z}$ .

### 2.1. 均值函数估计量的分布

本文假定误差序列  $\{\varepsilon_t\}$  来自标准正态总体, 即  $\varepsilon_t \sim N(0,1), t = 1, 2, \dots, n$ . 不妨记

$$\mathbf{M} = (m(x_1), m(x_2), \dots, m(x_n))^T,$$

$$\Sigma = \begin{pmatrix} \sigma(x_1) & 0 & \dots & 0 \\ 0 & \sigma(x_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma(x_n) \end{pmatrix},$$

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T,$$

则模型(1)可简记为

$$\mathbf{y} = \mathbf{M} + \Sigma \boldsymbol{\varepsilon}, \tag{4}$$

其中

$$\boldsymbol{\varepsilon} \sim N \left( \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \right).$$

因为  $\mathbf{M}$  与  $\Sigma$  分别为常数向量和常数矩阵, 又根据期望与方差的性质可知,  $\mathbf{y} \sim N(\mathbf{M}, \Sigma^2)$ .

**定理 1** 均值函数  $m(x)$  的样条估计量  $\hat{m}(x)$  服从均值为  $\mathbf{B}^T(x) \mathbf{A} \mathbf{M}$ , 方差为  $\mathbf{B}^T(x) \mathbf{A} \Sigma^2 \mathbf{A}^T \mathbf{B}(x)$  的正态分布, 即

$$\hat{m}(x) \sim N(\mathbf{B}^T(x) \mathbf{A} \mathbf{M}, \mathbf{B}^T(x) \mathbf{A} \Sigma^2 \mathbf{A}^T \mathbf{B}(x)),$$

其中  $\mathbf{A} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$ .

**证明:** 根据(2)式知, 要证  $\hat{m}(x)$  的分布, 只需证  $\hat{\boldsymbol{\phi}}$  的分布, 又因为  $\hat{\boldsymbol{\phi}} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{y}$ , 那么有

$$\hat{\boldsymbol{\phi}} \sim N \left( (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{M}, (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \Sigma^2 \mathbf{W} (\mathbf{W}^T \mathbf{W})^{-1} \right),$$

故有

$$\hat{m}(x) = \mathbf{B}^T(x) \hat{\boldsymbol{\phi}} \sim N(\mathbf{B}^T(x) \mathbf{A} \mathbf{M}, \mathbf{B}^T(x) \mathbf{A} \Sigma^2 \mathbf{A}^T \mathbf{B}(x))。$$

## 2.2. 方差函数估计量的分布

本文记  $X \leftrightarrow Y$  表示随机变量序列  $X$  与  $Y$  渐近同分布。接下来讨论方差函数估计量  $\hat{\sigma}^2(x)$  的渐近分布。

**定理 2** 方差函数  $\sigma^2(x)$  基于残差的样条估计量  $\hat{\sigma}^2(x)$  与  $\sum_{t=1}^n c_t \varepsilon_t^2$  渐近同分布, 即

$$\hat{\sigma}^2(x) \leftrightarrow \sum_{t=1}^n c_t \varepsilon_t^2,$$

这里  $\mathbf{c} = (c_1, c_2, \dots, c_n) = \mathbf{B}^T(x) \mathbf{A} \Sigma^2$ 。

**证明:** 根据(3)式, 要证方差函数估计量  $\hat{\sigma}^2(x)$  的分布, 只需证  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$  的分布。

因为

$$\hat{\zeta}_t = y_t - \hat{m}(x_t) = y_t - m(x_t) + m(x_t) - \hat{m}(x_t) = \sigma(x_t) \varepsilon_t - [\hat{m}(x_t) - m(x_t)],$$

则有

$$z_t = \hat{\zeta}_t^2 = \left\{ \sigma(x_t) \varepsilon_t - [\hat{m}(x_t) - m(x_t)] \right\}^2,$$

即

$$z_t = \left[ \sigma(x_t) \varepsilon_t \right]^2 - 2\sigma(x_t) \varepsilon_t [\hat{m}(x_t) - m(x_t)] + [\hat{m}(x_t) - m(x_t)]^2。$$

由  $\hat{m}(x)$  一致收敛到  $m(x)$  (见文献[7]中定理 1) 且  $\varepsilon_t$  是正态随机变量, 可知

$$\hat{\sigma}^2(x) \leftrightarrow \mathbf{B}^T(x) \cdot (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \cdot \begin{pmatrix} \sigma^2(x_1) \varepsilon_1^2 \\ \sigma^2(x_2) \varepsilon_2^2 \\ \vdots \\ \sigma^2(x_n) \varepsilon_n^2 \end{pmatrix} = \mathbf{B}^T(x) \mathbf{A} \Sigma^2 \begin{pmatrix} \varepsilon_1^2 \\ \varepsilon_2^2 \\ \vdots \\ \varepsilon_n^2 \end{pmatrix},$$

即

$$\hat{\sigma}^2(x) \leftrightarrow \sum_{t=1}^n c_t \varepsilon_t^2。$$

因为  $\{\varepsilon_t\}$  为相互独立且服从标准正态分布的随机误差序列, 故有

$$\varepsilon_t^2 \sim \chi^2(1), \quad \forall t = 1, 2, \dots, n,$$

因此, 要求方差函数估计量  $\hat{\sigma}^2(x)$  的渐近分布就是求服从卡方分布的独立随机变量的线性组合的分布。根据文献[13][14]的结果, 尝试用单个  $\chi^2$  变量的线性函数近似  $n$  个相互独立的  $\chi^2$  变量的线性组合。

首先, 考虑用  $\tilde{Y} = a\chi^2(d)$  作为  $\sum_{t=1}^n c_t \varepsilon_t^2$  的近似分布。采用一、二阶矩拟合的原则选取  $a, d$ , 即由方程

$$\begin{cases} E(\tilde{Y}) = E\left(\sum_{t=1}^n c_t \varepsilon_t^2\right) \\ D(\tilde{Y}) = D\left(\sum_{t=1}^n c_t \varepsilon_t^2\right) \end{cases}$$

确定。从而  $a, d$  应满足方程

$$\begin{cases} ad = \sum_{i=1}^n c_i \triangleq P \\ a^2 d = \sum_{i=1}^n c_i^2 \triangleq Q \end{cases}$$

解得

$$a = Q/P, \quad d = P^2/Q。$$

考虑到  $\chi^2$  变量的自由度为正整数, 故将  $d$  修正为

$$d^* = [P^2/Q + 0.5],$$

这里  $[x]$  表示不超过  $x$  的最大整数, 若  $P^2/Q < 0.5$ , 则取  $d^* = 1$ 。

再用  $Y \sim a\chi^2(d^*) + e_1$  来近似  $\sum_{i=1}^n c_i \varepsilon_i^2$  的分布, 采用上述方法得到

$$e_1 = P - \sqrt{Qd^*}, \quad a = (P - e_1)/d^*,$$

于是, 可用  $a\chi^2(d^*) + e_1$  作为  $\sum_{i=1}^n c_i \varepsilon_i^2$  的近似。

### 3. 数值模拟

考虑模型(1), 其中

$$\begin{cases} m(x) = 50x^3(1-x)^3 \\ \sigma(x) = 0.2 + 0.4\sin(\pi x), \\ \varepsilon_t \sim N(0,1), t = 1, 2, \dots, n \end{cases} \quad (5)$$

这里  $x \in [0,1]$ 。

对模型(5)进行蒙特卡罗模拟, 利用三次 B-样条基函数估计未知均值函数  $m(x)$ , 基于残差估计方差函数  $\sigma^2(x)$ , 并在 AIC 准则下自动选取等距结点数。

#### 3.1. 均值函数估计量的模拟

为验证理论分布效果, 使用 MATLAB 软件进行模拟运算, 选取显著性水平为  $\alpha = 0.05$ , 具体步骤如下:

第一步, 根据模型(5), 计算出  $\mathbf{B}^T(x)\mathbf{A}\mathbf{M}$  与  $\mathbf{B}^T(x)\mathbf{A}\Sigma^2\mathbf{A}^T\mathbf{B}(x)$ ;

第二步, 生成  $N$  个服从  $N(\mathbf{B}^T(x)\mathbf{A}\mathbf{M}, \mathbf{B}^T(x)\mathbf{A}\Sigma^2\mathbf{A}^T\mathbf{B}(x))$  的随机数;

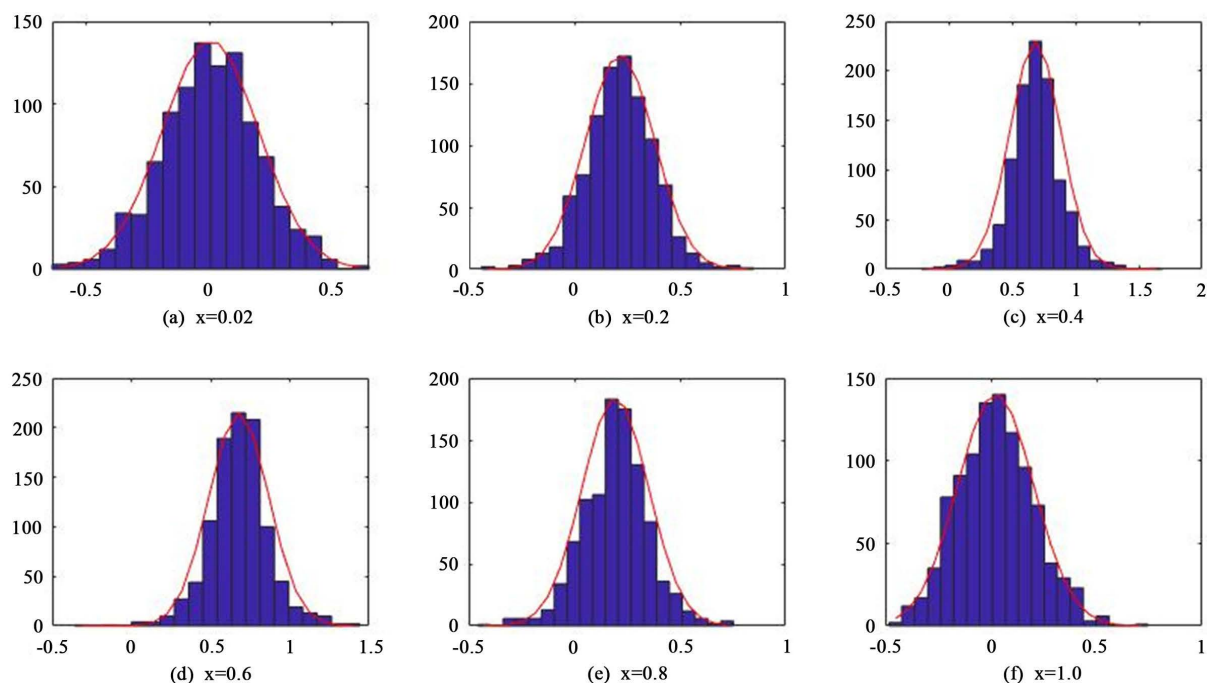
第三步, 基于样条方法, 通过蒙特卡罗模拟产生  $N$  个均值函数的估计值  $\hat{m}(x)$ ;

第四步, 在给定置信水平下, 检验第二步与第三步产生的随机数与均值函数估计值是否来自于同一分布;

第五步, 对上述四步进行多次重复模拟, 分析所得结果。

图 1 绘制了  $x = 0.02, 0.2, 0.4, 0.6, 0.8, 1.0$  处在  $N = 1000$  时的直方图和概率密度函数曲线图。

由图 1 可见, 各点处拟合的均值函数的直方图与概率密度函数曲线呈倒 U 型, 直观地可以认为各点处拟合的均值函数估计值来自于正态分布。进一步地, 对各点处的均值函数估计量的分布与正态分布进行 Two-sample t-test 检验, 并分别循环模拟  $N = 10, 50, 100, 500, 1000$  次, 检验的 P 值如表 1 所示。



**Figure 1.** Histogram of the estimates at each point of the mean function and its probability density function curve  
**图 1.** 均值函数各点处估计值的直方图及其概率密度函数曲线

**Table 1.** The P-values of the Two-sample t-test-test for the mean function

**表 1.** 均值函数 Two-sample t-test 检验的 P 值

$x/N$	10	50	100	500	1000
0.02	0.0007	0.4628	0.4062	0.9438	0.7338
0.20	0.0586	0.4463	0.5546	0.8140	0.7879
0.40	0.0472	0.7216	0.5068	0.6059	0.9070
0.60	0.1672	0.0525	0.3608	0.5405	0.6887
0.80	0.3916	0.7595	0.9055	0.6088	0.4992
1.00	0.0060	0.8371	0.7298	0.2669	0.9159

由表 1 可知, 当  $x = 0.2, 0.4, 0.8$  时, P 值均大于 0.05; 当  $N$  较大时, 各点处的 P 值均大于 0.05, 说明在给定的显著性水平 0.05 下, 应该接受原假设, 即认为检验数据服从正态分布。

### 3.2. 方差函数估计量的模拟

检验步骤如下:

第一步: 依据模型(5)给定  $x$  值, 计算  $a$ ,  $d^*$  与  $e_1$ ;

第二步: 随机生成  $N$  个服从  $a\chi^2(d^*) + e_1$  的随机数;

第三步: 基于残差样条方法, 通过蒙特卡罗模拟生成  $N$  个方差函数  $\hat{\sigma}^2(x)$  的估计值;

第四步: 在给定置信水平下, 检验第二步与第三步产生的随机数与方差函数估计值是否来自于同一分布;

第五步: 对上述四步进行多次重复模拟, 分析所得结果。

图 2 绘制了  $x = 0.02, 0.2, 0.4, 0.6, 0.8, 1.0$  处在  $N = 1000$  时的直方图和概率密度函数曲线图。

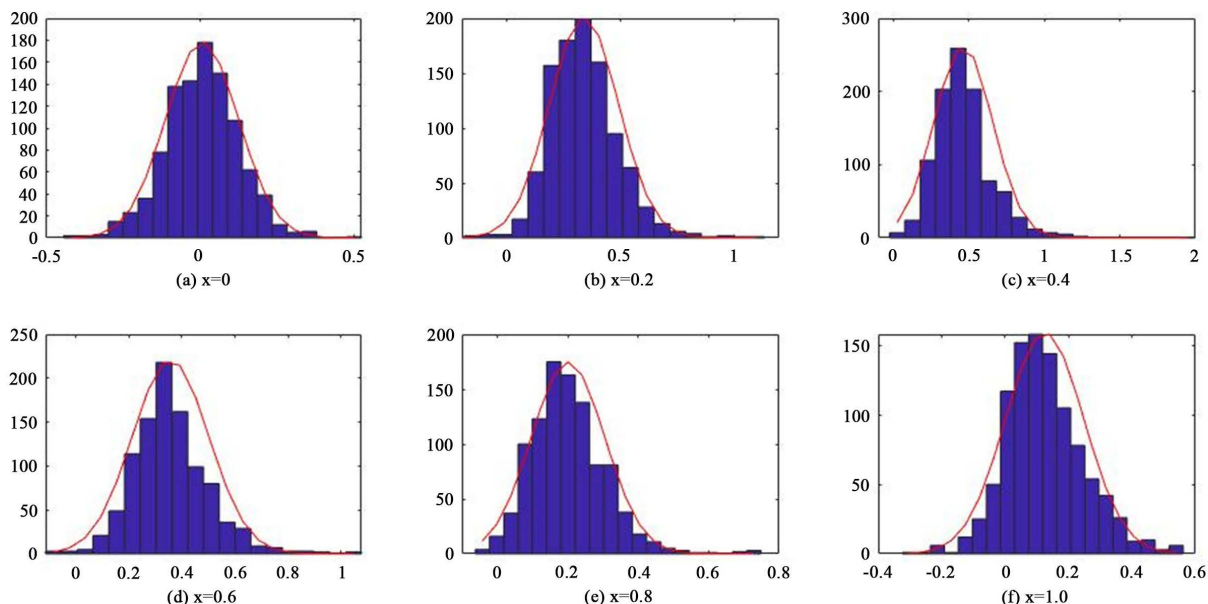


Figure 2. Histogram of the estimates at each point and its probability density function curve

图 2. 方差函数各点处估计值的直方图及其概率密度函数曲线

由图 2 可见, 各点处拟合的方差函数估计量的直方图与概率密度函数曲线呈不对称的倒 U 型, 且整体上右偏, 直观地可以认为各点处拟合的方差函数估计量近似服从卡方分布。进一步地, 对各点处的方差函数估计量的渐近分布与  $a\chi^2(d^*) + e_1$  进行 Two-sample t-test 检验, 并分别循环模拟  $N = 10, 50, 100, 500, 1000$  次, 结果如表 2 所示。

由表 2 可知, 各点处的 P 值均大于给定的显著性水平 0.05, 说明可以接受原假设, 即认为检验数据近似服从  $a\chi^2(d^*) + e_1$ 。

Table 2. The P-values of the Two-sample t-test-test for the variance function

表 2. 方差函数 Two-sample t-test 检验的 P 值

$x/N$	10	50	100	500	1000
0.02	0.3875	0.2297	0.0620	0.1405	0.0955
0.20	0.3715	0.0989	0.9234	0.9621	0.6536
0.40	0.8930	0.3863	0.9431	0.3844	0.9430
0.60	0.2782	0.6464	0.7323	0.1275	0.7215
0.80	0.3610	0.9049	0.9524	0.7373	0.9535
1.00	0.2025	0.4128	0.8573	0.9954	0.9361

#### 4. 结论

本文基于样条方法研究了固定设计下异方差非参数回归模型的均值函数与方差函数估计量的分布, 均值函数的估计量服从正态分布, 方差函数估计量的渐近分布可由单个  $\chi^2$  变量的线性函数来近似。模拟结果显示: 在给定显著性水平 0.05 下, 分布拟合效果较优。

本文所研究的固定设计下异方差非参数回归模型的均值函数与方差函数估计量的近似分布为生物、医学、地质、经济等领域的研究带来了便利。

## 基金项目

国家自然科学基金项目(11601126); 河南省重点攻关项目(182102210286)。

## 参考文献

- [1] Chown, J. (2016) Efficient Estimation of the Error Distribution Function in Heteroskedastic Nonparametric Regression with Missing Data. *Statistics & Probability Letters*, **117**, 31-39. <https://doi.org/10.1016/j.spl.2016.04.009>
- [2] 齐培艳, 田铮, 段西发, 袁芳. 异方差非参数回归模型均值与方差变点的小波估计与应用[J]. 系统工程理论与实践, 2013, 33(4): 988-995.
- [3] Burman, P. (1991) Regression Function Estimation from Dependent Observations. *Journal of Multivariate Analysis*, **36**, 263-279. [https://doi.org/10.1016/0047-259X\(91\)90061-6](https://doi.org/10.1016/0047-259X(91)90061-6)
- [4] Song, Q. and Yang, L. (2009) Spline Confidence Bands for Variance Functions. *Journal of Nonparametric Statistics*, **21**, 589-609. <https://doi.org/10.1080/10485250902811151>
- [5] 武新乾, 张刚. 非参数回归模型中误差方差的样条估计[J]. 郑州大学学报(理学版), 2015, 47(3): 17-20.
- [6] 郑美洁, 田波平. 基于两步样条光滑法的非参数回归模型研究[J]. 统计与决策, 2020, 36(3): 14-20.
- [7] 马晓跃, 武新乾. 非参数回归模型基于残差的样条估计[J]. 河南科技大学学报(自然科学版), 2021, 42(4): 91-96+10.
- [8] 秦永松. 一类非参数回归函数导数估计的渐近分布[J]. 工程数学学报, 1991, 8(1): 67-74.
- [9] Liang, H. and Jing, B. (2004) Asymptotic Properties for Estimates of Nonparametric Regression Models Based on Negatively Associated Sequences. *Journal of Multivariate Analysis*, **95**, 227-245. <https://doi.org/10.1016/j.jmva.2004.06.004>
- [10] Jin, S., Su, L. and Xiao, Z. (2014) Adaptive Nonparametric Regression with Conditional Heteroskedasticity. *Economic Theory*, **31**, 1153-1191. <https://doi.org/10.1017/S0266466614000450>
- [11] Alharbi, Y. and Patili, P. (2018) Error Variance Function Estimation in Nonparametric Regression Models. *Communications in Statistics-Simulation and Computation*, **47**, 1479-1491. <https://doi.org/10.1080/03610918.2017.1315774>
- [12] Li, Z. and Lin, W. (2020) Efficient Error Variance Estimation in Non-Parametric Regression. *Australian & New Zealand Journal of Statistics*, **62**, 467-484. <https://doi.org/10.1111/anzs.12311>
- [13] 范大茵, 冯云. 独立  $\chi^2$  变量线性组合的近似分布[J]. 高校应用数学学报 A 辑(中文版), 1993, 8(3): 335-338.
- [14] Zhang, J. (2011) Approximate and Asymptotic Distributions of Chi-Squared-Type Mixtures with Applications. *Journal of the American Statistical Association*, **100**, 273-285. <https://doi.org/10.1198/016214504000000575>