

分位数模型回归分析

胡良平^{1 2*}

(1. 军事科学院研究生院, 北京 100850;

2. 世界中医药学会联合会临床科研统计学专业委员会, 北京 100029

* 通信作者: 胡良平, E-mail: lphu812@sina.com)

【摘要】 本文目的是介绍分位数模型回归分析的概念、作用以及用软件实现计算的方法。先介绍有关的基本概念, 再介绍基本原理, 最后通过一个实例并基于 SAS 软件演示如何实施此分析。结果表明: 分位数模型回归分析能够很好地解决因变量不服从正态分布并具有异方差且资料中存在一定比例异常点的回归分析问题。

【关键词】 分位数; 分位数模型回归分析; 简单算法; 内部点算法; 光滑算法

中图分类号: R195.1

文献标识码: A

doi: 10.11886/j.issn.1007-3256.2018.04.002

Quantile model regression analysis

Hu Liangping^{1 2*}

(1. Graduate School, Academy of Military Sciences PLA China, Beijing 100850, China;

2. Specialty Committee of Clinical Scientific Research Statistics of World Federation of Chinese Medicine Societies, Beijing 100029, China

* Corresponding author: Hu Liangping, E-mail: lphu812@sina.com)

【Abstract】 The purpose of this paper was to introduce the concepts, functions and applying calculation methods by using the statistical software of the quantile model regression analysis. Firstly, the basic concepts of the regression analysis were introduced. Secondly, the current paper illustrated the basic principle of the regression analysis. Finally, the quantile model regression analysis was demonstrated through one example by using the SAS software. The method mentioned above could perfectly solve the problem of the regression analysis of the data with the following situations: ①the dependent variable was not the normal distribution and had the heterogeneity of variance, ②there were some outliers in the data.

【Keywords】 Quantile; Quantile model regression analysis; Simplex algorithm; Interior point algorithm; Smoothing algorithm

1 概 述^[1-2]

1.1 分位数回归模型

1.1.1 分位数

分位数是一种位置指标, 一个特定的分位数将任何一个频数曲线下的面积(其数值为 1)分为两部分, 若小于等于此分位数的观测值个数占全部观测值个数的比例为 $1/4$, 则称该分位数为第 1 四分位数, 记作 Q_1 , 同理, 还有第 2、第 3 四分位数, 分别记作 Q_2 、 Q_3 ; 若小于等于此分位数的观测值个数占全部观测值个数的比例为 $1/10$, 则称该分位数为第 1 十分位数, 记作 D_1 , 同理, 还有第 2、第 3、……、第 9 十分位数, 分别记作 D_2 、 D_3 、……、 D_9 ; 若小于等于此分位数的观测值个数占全部观测值个数的比例为 $1/100$, 则称该分位数为第 1 百分位数, 记作 P_1 , 同理, 还有第 2、第 3、……、第 99 百分位数, 分别记作 P_2 、 P_3 、……、 P_{99} 。

显然, 第 1 四分位数 = 第 25 百分位数, 即 $Q_1 = P_{25}$; 第 2 四分位数 = 第 5 十分位数 = 第 50 百分位数 = 中位

数, 即 $Q_2 = D_5 = P_{50} = M$ (代表中位数); 第 3 四分位数 = 第 75 百分位数, 即 $Q_3 = P_{75}$ 。如此, 常用百分位数代替四分位数和十分位数。通过给出一组资料的若干个分位数, 可初步描述该组资料的离散程度和分布概况, 故在实际工作中, 常用百分位数法确定服从偏态分布资料的医学指标的正常值范围。

1.1.2 均值回归模型与均值回归方程

在仅有一个定量自变量和一个定量因变量的简单直线回归分析中, 经典统计学中需要事先给出如下几个假定。

其一, 正态假定。即当自变量 x 在其取值区间内取定任何一个特定值 x_i 时, 因变量 y 可以有一组取值与其对应。例如很多身高 x 为 165cm 的成年人其体重 y 是不等的。如此多的 y 值就会有一个概率分布, 粗略地说, 该分布可能是正偏态的、对称的(最理想的为正态的)或负偏态的, 简单记作 $P(y|x_i)$, 经典统计学假定 $P(y|x_i)$ 为“正态分布”, 对所有的 x_i 都成立。

其二, 同方差假定。即当 x_i 取任何值时, 其对应的上述正态分布 $P(y|x_i)$ 的方差均相等, 称为“同

项目基金: 国家高技术研究发展计划课题资助(2015AA020102)

方差”或“方差齐性”。否则,就称为“异方差”或“方差不齐”。

其三,均值假定。即当 x_i 取任何值时,一律采用 y 的算术平均值 \bar{y}_i 取代任何一个 y_j 值,记作 $\hat{y} = \bar{y}_i | x_i$ 。于是,所有的数据点 (x, y) 可表示为 (x_i, \hat{y}_i) 。

其四,独立误差假定。即当 x_i 取任何值时,观测值 $(y_j | x_i)$ 与其算术均值 $\bar{y}_i | x_i$ (或 \hat{y}_i) 之间的偏差为“误差 ε_i ”,假定各误差 ε_i 之间互相独立。

基于以上假定且按“普通最小二乘原理或最小平方原理”^[3] 构建的简单直线回归模型和多重线性回归模型都被称为“均值回归模型”,分别见式(1)和式(2):

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i = 1, 2, \dots, n \quad (1)$$

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_m x_{mi} + \varepsilon_i \quad i = 1, 2, \dots, n \quad (2)$$

在模型(1)和(2)中,基于样本数据并按最小二乘原理估计出其参数的数值,然后忽略模型中的误差项,就可得到“均值回归方程”,见式(3)和式(4):

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad i = 1, 2, \dots, n \quad (3)$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_m x_{mi} \quad i = 1, 2, \dots, n \quad (4)$$

在式(3)和式(4)中,等号左边的“ \hat{y}_i ”是在单个自变量 x_i (在只有一个自变量的场合下)或多个自变量构成的向量 (X_i) (在有多个自变量的场合下)取一个或一组特定数值的前提下,重复 k 次观测或试验得到了因变量 y 的 k 个观测值的“算术平均值”,简称为“均值”。严格地说,它应该被分别表示成如下的形式,见式(5)和式(6)。

对式(3)而言 \hat{y}_i 的真实含义应该是式(5):

$$\hat{y}_i = \bar{y}_i | (x = x_i) = \frac{1}{k} \sum_{j=1}^k y_{ij} | (x = x_i) \quad i = 1, 2, \dots, n \quad (5)$$

对式(4)而言 \hat{y}_i 的真实含义应该是式(6):

$$\hat{y}_i = \bar{y}_i | (X = X_i) = \frac{1}{k} \sum_{j=1}^k y_{ij} | (X = X_i) \quad i = 1, 2, \dots, n \quad (6)$$

在式(6)中,大写的字母“ X ”代表由多个自变量构成的向量。

也就是说,常用的简单直线回归模型或方程和多重线性回归模型或方程都属于“均值回归模型或方程”。“最小二乘原理”可用语言表述如下:

在各 x_i 或 X_i 条件下,求因变量 y 的任何一个观测值与其模型对应的估计值 \hat{y}_i 之间偏差的平方,再将所有各偏差平方求和,称此“和”为“目标函数 Q ”,设法求其最小值,与此同时,获得求解模型中参

数估计值的计算公式。

值得一提的是:对基于最小二乘原理构造出来的目标函数 Q 求其最小值时,需要将 Q 中的待估参数视为“变量”,而将各观测点上变量的取值 (x_i, y_i) (一个自变量的情况下)或 $(x_{1i}, x_{2i}, \dots, x_{mi}, y_i)$ (m 个自变量的情况下)视为“常数”,利用高等数学中求函数“极大值”和“极小值”的方法,即求函数 Q 关于未知参数的偏导数等多个计算步骤来实现目的。事实上,求得的偏导数又是关于研究问题中“变量 (x, y) 或 $(x_1, x_2, \dots, x_m, y)$ ”的函数,故称其为“导函数”。令求得的导函数为零,就得到若干个方程,其个数为所构建的回归模型中参数的个数。在直线回归模型中,有两个参数,即截距和斜率;在含有 m 个自变量的多重线性回归模型中,有 $(m + 1)$ 个参数。将这些线性方程联立在一起,就构成了“线性方程组”。求出该方程组的解,就得到了求解待估计参数的计算公式。

1.1.3 中位数回归模型与中位数回归方程

在前面构建“均值回归模型或方程”的“均值假定”中,若用“中位数”取代“算术平均值”就得到了“中位数回归模型或回归方程”,在有 m 个自变量的场合下,中位数回归模型见式(7)、中位数回归方程见式(8):

$$Y_{(0.5)i} | (X = X_i) = \beta_{0(0.5)} + \beta_{1(0.5)i} x_{1i} + \beta_{2(0.5)i} x_{2i} + \dots + \beta_{m(0.5)i} x_{mi} + \varepsilon_{(0.5)i} \quad (7)$$

$$\hat{y}_{(0.5)i} | (X = X_i) = \hat{\beta}_{0(0.5)} + \hat{\beta}_{1(0.5)i} x_{1i} + \hat{\beta}_{2(0.5)i} x_{2i} + \dots + \hat{\beta}_{m(0.5)i} x_{mi} \quad (8)$$

在式(7)和式(8)中,下角标中的 i 代表观测的编号 $i = 1, 2, \dots, n$ 。“0”代表“常数项”;“1、2、……、 m ”分别代表“第 1 个”“第 2 个”……“第 m 个”自变量。“0.5”代表“第 50 百分位数”或“第 2 四分位数”或“中位数”所对应的累计概率,也就是说,“ $y_{(0.5)}$ ”就是与其对应的“分位数”,即“中位数”。

1.1.4 分位数回归模型与分位数回归方程

由前面关于“分位数”的概念可知“中位数”是一个特定的“分位数”。若将上面的“中位数回归模型或方程”中的“中位数”替换成任意一个分位数“ $y_\tau | X_i$ ”,则所得到的回归模型就被称为“分位数回归模型或方程”了。前述的“ $y_\tau | X_i$ ”的含义是:在自变量 x 取特定值 x_i 的条件下,求出全部因变量 y 的第 τ 分位数 $0 < \tau < 1$ 。当 $\tau = 0.5$ 时,就是“第 0.5 分位数”或“中位数”;当 $\tau = 0.25$ 时,就是“第 0.25

分位数”或“第 1/4 分位数”,也叫做“第 1 四分位数”;当 $\tau = 0.75$ 时,就是“第 0.75 分位数”或“第 3/4 分位数”,也叫做“第 3 四分位数”。

在式(7)和式(8)中,若将“0.5”改换成“ $\tau \in (0, 1)$ ”(其含义是 $0 < \tau < 1$),就得到了与第 τ 分位数对应的回归模型或回归方程,见式(9)和式(10):

$$y_{(\tau)i} | (X = X_i) = \beta_{0(\tau)} + \beta_{1(\tau)i} X_{1i} + \beta_{2(\tau)i} X_{2i} + \dots + \beta_{m(\tau)i} X_{mi} + \varepsilon_{(\tau)i} \quad (9)$$

$$\hat{y}_{(\tau)i} | (X = X_i) = \hat{\beta}_{0(\tau)} + \hat{\beta}_{1(\tau)i} X_{1i} + \hat{\beta}_{2(\tau)i} X_{2i} + \dots + \hat{\beta}_{m(\tau)i} X_{mi} \quad (10)$$

在分位数回归分析中,每给定一个“ τ 值”,就可求出一个相应的“分位数回归模型或方程”。故对于一个给定的资料来说,因 τ 在开区间 $(0, 1)$ 范围内有无穷多个取值,就有无穷多个“分位数回归模型或方程”。其中,最常见的分位数回归模型或方程为“第 1 四分位数回归模型或方程”“第 2 四分位数回归模型(即中位数回归模型)或方程”和“第 3 四分位数回归模型或方程”。

1.2 分位数模型回归分析应用的场合

当拟做多重线性回归分析的原始数据中的定量因变量不服从正态分布、有时还存在异方差、资料中存在一定比例的“异常点”且自变量间不存在严重多重共线性时,采用此方法构建多重线性回归模型,可以最大限度地消除资料中违反经典回归分析的“部分或全部假定”对建模结果造成的影响。

1.3 分位数模型回归分析的计算原理

1.3.1 概述

值得注意的是:在求解“均值回归模型”中参数时,是基于“最小二乘原理”推导出来的正规方程组并求解得到的;而当采用“分位数”取代“均值”时,就不适合采取最小二乘原理来构建“目标函数 Q”,

而需要求因变量 y 的任何一个观测值与其模型对应的估计值 \hat{y}_i 之间偏差的“绝对值”,再将所有各点上的偏差绝对值求和,称此“和”为“目标函数 G”。令式(10)等号左边为“W”,则为求解式(10)中参数估计值而构建的目标函数 G 如式(11)所示:

$$G = - \sum_{i: y_{(\tau)i} < W} (1 - \tau) (y_i - W) + \sum_{i: y_{(\tau)i} \geq W} \tau (y_i - W) \quad (11)$$

设法求出式(11)的最小值,与此同时,获得模型中参数的估计值,这样求出的参数估计值被称为“分位数回归参数估计值”。

1.3.2 式(11)最小值的计算方法

由 SAS 软件的帮助信息 [从 SAS/STAT 的 QUANTREG 过程的“details”(即详细情况)菜单进入,其第 2 行为“Optimization Algorithms”(优化算法)]可知,求式(11)最小值的计算方法有三种,分别为“简单算法”“内部点算法”和“光滑算法”,在样本含量 $n < 5\,000$ 且自变量个数 $m < 50$ 时,三种算法的参数估计结果是基本相同的。

事实上,上述提及的三种算法都相当复杂,涉及到复杂的矩阵运算和反复迭代计算过程。换言之,由式(11)无法给出与各种分位数回归模型或方程中各参数的解析式,即计算公式。故具体分析时,建议采用统计软件来实现计算。因篇幅所限,详细的计算方法此处不再赘述。

2 基于分位数模型回归分析解决实际问题^[1]

2.1 问题与数据结构

【例 1】假定有一个总样本含量 $n = 1\,000$ 的数据集中包含 5% 异常点的资料,每组数据(即每个个体)包含三个变量(x_1, x_2, y)的观测值。见表 1。

表 1 某资料中首尾各 10 组数据($n = 1000$)

num	x_1	x_2	y	num	x_1	x_2	y
1	1.42151	1.13105	21.201	991	-1.15836	-1.04066	101.788
2	2.00893	0.79083	21.192	992	0.71023	-0.61811	100.557
3	1.82697	-0.02043	18.602	993	1.04106	0.26261	103.892
4	-1.49036	-0.93768	-1.025	994	-0.26348	-0.05647	101.371
5	-0.45912	-3.42593	-2.552	995	-1.09514	-1.76723	99.935
6	0.39892	-2.02172	5.997	996	1.65162	-0.83433	87.594
7	1.08865	-0.98391	13.049	997	-0.26465	-0.73262	104.739

续表 1:

8	-0.48139	1.33821	11.208	998	0.55630	-1.52099	105.340
9	-0.23384	-0.85954	6.060	999	0.05137	0.24319	95.097
10	0.00900	-0.24376	9.661	1000	-1.53994	0.96521	105.054

注:表 1 省略编号为 11~990 的 980 行数据;在全部 1000 行数据中,最后 50 行数据为“异常点”,占 5%

【特别说明】例 1 是人为构造的,它来自 SAS 9.3 的 QUANTREG 过程中的“样例”。三个变量“ x_1 、 x_2 和 y ”没有实际的专业含义,仅为了造出一个样本含量为 1 000 且含 5% 异常点的数据集。

表 1 中数据构造的方法如下:设定 x_1 和 x_2 及测量误差 e 都是服从标准正态分布的随机变量(其均值为 0、方差为 1),前 950 个 y 的数值按下面的模型(12)计算出来;后 50 个 y 的数值按下面的模型(13)计算出来:

$$y = 10 + 5 * x_1 + 3 * x_2 + 0.5 * e \quad (12)$$

$$y = 100 + e \quad (13)$$

比较式(12)与式(13)可知: y 的前 950 个数据中的每一个都在基数“10”的基础上再加上三项并不大的数值,其平均值约为“ $10 + 5 + 3 = 18$ ”;而 y 的后 50 个数据中的每一个都在基数“100”的基础上再加上一个随机误差,其平均值约为 100。由此可知:表 1 的 1000 行数据中,对因变量 y 而言,后 50 个 y 值明显大于前 950 个 y 值,故属于“异常值”,它们所对应的那 50 行数据点就属于“异常点”了。

【问题】试拟合表 1 中 y 依赖 x_1 、 x_2 变化的二重线性回归模型。

2.2 所需要的 SAS 程序

2.2.1 产生数据集

先用下面的一段 SAS 数据步程序产生表 1 中

的 1000 行 3 列数据,创建数据集 a。

```
data a ( drop = i );
do i = 1 to 1000;
    x1 = rannor( 1234 );
    x2 = rannor( 1234 );
    e = rannor( 1234 );
    if i > 950 then y = 100 + 10 * e;
    else y = 10 + 5 * x1 + 3 * x2 + 0.5 * e;
output;
end;
run;
/* 以上程序产生 1000 行数据( x1 ,x2 ,y) ,其中 ,有
5% 的是异常值 * /
```

以上程序产生的数据见表 1。

2.2.2 对资料中因变量 y 的分布情况进行探索性分析

采用下面的 SAS 程序对因变量 y 进行正态性检验并绘制其频数直方图,直观了解 y 的频数分布情况。

```
proc univariate data = a;
var y;
histogram y;
run;
```

以上 SAS 程序输出结果见表 2、图 1。

表 2 对表 1 中 1000 个 y 值是否服从正态分布检验结果

检验	统计量		p 值
Shapiro - Wilk	W	0.500517	Pr < W < 0.0001
Kolmogorov - Smirnov	D	0.308976	Pr > D < 0.0100
Cramer - von Mises	W - Sq	29.62101	Pr > W - Sq < 0.0050
Anderson - Darling	A - Sq	162.54	Pr > A - Sq < 0.0050

由表 2 第 1 行“W 检验”可知,本例中 y 值不服从正态分布。

图 1 显示:1 000 个 y 值中的 950 个在图中左边部分,呈正偏态分布;中间出现了较大一段空缺,还有 50 个 y 值位于图中右边部分,其数值波动约为

72 ~ 120。

2.2.3 对资料进行分位数模型回归分析

下面基于“ $\tau = 0.25$ 、 0.50 和 0.75 ”分别用“分位数回归分析法”来创建二重线性回归模型:

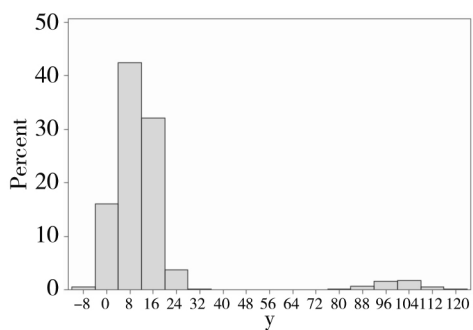


图1 表1中全部1000个y值的频数分布直方图

```
proc quantreg algorithm = smooth( ratio = .5) data = a;
model y = x1 x2/quantile = 0.25 0.50 0.75;
run;
```

【SAS 程序说明】model 语句中的选项“quantile = 0.25 0.50 0.75”是要求 SAS 系统分别创建分位数为 0.25、0.50(中位数)和 0.75 的 3 个二重线性回归模型。

【SAS 主要输出结果】

Quantile and Objective Function	
Quantile	0.25
Objective Function	1282.8452
Predicted Value at Mean	9.6857

Parameter Estimates

Parameter	DF	Estimate	Standard Error	95% Confidence Limits		t Value	Pr > t
Intercept	1	9.6957	0.0198	9.6569	9.7345	490.27	<.0001
x1	1	5.0083	0.0195	4.9701	5.0466	257.09	<.0001
x2	1	3.0170	0.0148	2.9880	3.0460	204.22	<.0001

以上是“ $\tau = 0.25$ ”对应的“第 1/4 四分位数回归模型”的参数估计及检验结果

Quantile and Objective Function	
Quantile	0.5
Objective Function	2441.1927
Predicted Value at Mean	10.0259

Parameter Estimates

Parameter	DF	Estimate	Standard Error	95% Confidence Limits		t Value	Pr > t
Intercept	1	10.0364	0.0219	9.9934	10.0794	457.96	<.0001
x1	1	5.0106	0.0194	4.9725	5.0487	257.90	<.0001
x2	1	3.0294	0.0186	2.9930	3.0658	163.31	<.0001

以上是“ $\tau = 0.50$ ”时对应的“第 2/4 四分位数回归模型”(即中位数回归模型)的参数估计及检验结果,结果表明:截距项 = 10.0364、两个自变量的斜率分别为 5.0106 和 3.0294,参考它们的理论值[见式(12)]分别为 10.5 和 3,说明在因变量不服从正态分布且资料中含有 5% 异常点时,采用“中位数回归分析”创建的二重线性回归模型与其“原模型”之

间的吻合程度非常高。

Quantile and Objective Function	
Quantile	0.75
Objective Function	3512.2464
Predicted Value at Mean	10.4106

Parameter Estimates

Parameter	DF	Estimate	Standard Error	95% Confidence Limits		t Value	Pr > t
Intercept	1	10.4205	0.0247	10.3720	10.4690	421.73	<.0001
x1	1	4.9993	0.0275	4.9453	5.0532	181.90	<.0001
x2	1	3.0093	0.0280	2.9544	3.0641	107.66	<.0001

以上呈现的是“ $\tau = 0.75$ ”时对应的“第 3/4 四分位数回归模型”的参数估计及检验结果。

2.3 小结^[3-4]

2.3.1 采用普通最小二乘法构建二重线性回归模型

```
proc reg data = a;
    model y = x1 x2;
```

变量	自由度	参数估计值	标准误差	t 值	Pr > t
Intercept	1	14.48953	0.62584	23.15	<.0001
x1	1	4.39030	0.62997	6.97	<.0001
x2	1	2.50293	0.60204	4.16	<.0001

以上结果表明:截距项 = 14.48953,两个自变量的斜率分别为 4.39030、2.50293,参考它们的理论值[见式(12)]分别为 10、5 和 3,说明在因变量不服从正态分布且资料中含有 5% 异常点时,直接基于普通最小二乘原理创建的二重线性回归模型很不理想。

2.3.2 采用“稳健回归分析”创建二重线性回归模型

```
proc robustreg data = a method = lts seed = 100;
    model y = x1 x2;
run;
```

以上 SAS 程序的主要输出结果如下:

复相关系数的平方 $R^2 = 0.9933$ 非常大,表明此二重线性回归模型具有很高的实用价值,其参数估计如下:

LTS Parameter Estimates

Parameter	DF	Estimate
Intercept	1	10.0054
x1	1	5.0240
x2	1	3.0598

以上结果表明:截距项 = 10.0054,两个自变量的斜率分别为 5.0240 和 3.0598,参考它们的理论值[见式(12)]分别为 10、5 和 3,说明在因变量不服从正态分布且资料中含有 5% 异常点时,基于最小截平方和(Least trimmed squares, LTS)法的“稳健回归分析”创建的二重线性回归模型与其“原模型”非常吻合。

2.3.3 总结三种建模方法所得结果的差异

在因变量不服从正态分布且资料中存在异常点

run;

以上 SAS 程序的主要输出结果如下:

总模型的 $F = 33.76$, $P < 0.0001$,表明所创建的二重线性回归模型具有统计学意义;但复相关系数的平方 $R^2 = 0.0634$ 非常小,表明此二重线性模型的实用价值并不高;模型中三个参数的估计与假设检验结果如下:

变量	自由度	参数估计值	标准误差	t 值	Pr > t
Intercept	1	14.48953	0.62584	23.15	<.0001
x1	1	4.39030	0.62997	6.97	<.0001
x2	1	2.50293	0.60204	4.16	<.0001

时,采用 SAS 中 REG 过程并基于普通最小二乘原理直接创建多重线性回归模型的效果很差,而基于“稳健回归分析法”和“分位数回归分析法”得到的结果非常接近。事实上,本例中的数据是基于上面的二重线性回归模型(12)产生的。这个模型意味着:截距项为 10、 x_1 前的回归系数为 5、 x_2 前的回归系数为 3,在此基础上,加一个随机误差的二分之一。此处的“随机误差”是服从均值为 0、方差为 1 的正态分布的随机误差。基于上面的计算结果,可以写出三个二重线性回归模型的具体表达式如下,见式(14)、式(15)、式(16):

$$\hat{y} = 14.48953 + 4.39030x_1 + 2.50293x_2 \quad (14) \text{ (OLS 估计法)}$$

$$\hat{y} = 10.0054 + 5.0240x_1 + 3.0598x_2 \quad (15) \text{ (LTS 估计法)}$$

$$\hat{y} = 10.0364 + 5.0106x_1 + 3.0294x_2 \quad (16) \text{ (分位数估计法)}$$

【结论】以模型(12)为“金标准”模型(14)偏离很远;模型(15)和(16)的质量都很高。相比之下,中位数回归分析的结果稍优于稳健回归分析法(以“回归系数与其真值之间的偏差最小”为评价依据)。

参考文献

- [1] SAS Institute Inc. STAT SAS 9.3 User's Guide[M]. Cary, NC: SAS Institute Inc, 2011: 6261-6338.
- [2] 方匡南,朱建平,姜叶飞. R 数据分析方法与案例详解[M]. 北京:电子工业出版社, 2015: 330-345.
- [3] 胡良平,胡纯严,鲍晓蕾. 应用数理统计[M]. 北京:电子工业出版社, 2015: 145-152.
- [4] 胡良平. 稳健回归分析[J]. 四川精神卫生, 2018, 31(3): 201-204.

(收稿日期:2018-08-10)

(本文编辑:陈霞)