

基于0-1膨胀二项分布的客观贝叶斯分析

吴懿祺¹ 肖翔¹ 古晞²

¹(上海工程技术大学数理与统计学院 上海 201620)

²(同济大学数学科学学院 上海 200092)

摘要 在医疗卫生、金融证券等应用领域,经常会同时出现零观测值、一观测值较多的情况。为更好地拟合这类数据,提出一种0-1膨胀二项分布模型并进行客观贝叶斯分析。采用数据扩充策略,基于完全似然函数,得到Jeffreys先验和reference先验。采用WinBUGS软件和R软件进行数值模拟,设定不同的样本量和参数真值,对不同的无信息先验进行评估。对2020年1月28日与2月22日COVID-19死亡人数进行分析,结果表明,在小样本情形下基于客观贝叶斯先验 π_{R3} 下的拟合效果比 π_{R1} 和 π_{R2} 要好。

关键词 0-1膨胀二项分布 客观贝叶斯 Jeffreys先验 reference先验 数据扩充

中图分类号 TP3

文献标志码 A

DOI:10.3969/j.issn.1000-386x.2024.04.007

OBJECTIVE BAYESIAN ANALYSIS BASED ON ZERO-AND-ONE-INFLATED BINOMIAL DISTRIBUTION

Wu Yiqi¹ Xiao Xiang¹ Gu Xi²

¹(School of Mathematics, Physics and Statistics, Shanghai University of Engineering Science, Shanghai 201620, China)

²(School of Mathematical Science, Tongji University, Shanghai 200092, China)

Abstract Count data with excess zeros and ones arise frequently in various fields such as medical health, finance and securities. To better fit such data, a zero-and one-inflated binomial distribution model is proposed and the objective Bayesian analysis is carried. Based on data augmentation strategy and the complete likelihood function, the Jeffreys prior and the reference priors were derived for this model. For different sample sizes and different true values of the parameters, simulations were adopted to assess the performance of the different uninformed priors through WinBUGS and R software. We analyze the death toll of COVID-19 on January 28 and February 22, 2020. The results show that the fitting effect based on objective Bayesian prior π_{R3} is better than π_{R1} and π_{R2} in the case of small sample.

Keywords Zero-and-one-inflated binomial distribution Objective Bayesian Jeffreys prior Reference prior Data augmentation

0 引言

计数数据一直是统计学研究的热点,在公共卫生、道路安全、保险精算和工业生产等众多领域存在着大量的计数数据。泊松模型和0膨胀泊松分布(ZIP)是用来处理计数数据常用的方法^[1-4]。然而在实际应用中,经常会遇到0和1过多(称之为“0-1膨胀”)的数据样本,其中0代表未发生,1代表发生1次。例如,

在COVID-19大流行中,大部分个体在感染过一次新冠肺炎病毒后,自身产生了抗体,使得其感染的次数可能最多为一次。又如,由于交通安全法规意识的增强,或者一些微小事故未申报理赔,所以保险公司申报索赔的频数中0和1的数据较高。这时如果仍然用传统的模型进行拟合,拟合效果不尽如人意。

近年来,国内外学者把对ZIP的研究推广到0-1膨胀泊松分布(ZOIP),并取得了丰富的研究成果。Maria^[5]在研究瑞典成年人看牙医的次数时,发现数据

中出现了较多 0 和 1,首次提出了 0-1 膨胀泊松模型。田震^[6]研究了 0-1 膨胀 ZOIP 回归模型及其参数估计,并基于数据删失和数据加权扰动对模型进行统计诊断。Tang 等^[7]通过引入隐变量,构造了 ZOIP 模型新的结构形式,采用极大似然估计与贝叶斯方法对模型进行参数估计,并对新加坡军团菌感染数据和美国底特律城市交通事故死亡数据进行研究,取得了较好的拟合效果。Liu 等^[8]通过参数变化的方法,计算了 ZOIP 模型中参数的 Jeffreys 先验和 reference 先验,并进行了客观贝叶斯分析,拟合效果比使用 naive flat 先验要更好。夏丽丽等^[9]采用了局部多项式核回归法对 ZOIP 模型进行参数估计,结合 EM 算法和 Newton-Raphson 迭代法对参数进行近似求解,通过对糖尿病患者数据的实例分析,验证了方法的有效性。

目前,国内外学者的研究主要集中在 0-1 膨胀泊松分布和 0-1 膨胀负二项分布,对 0-1 膨胀二项分布的研究几乎是一片空白。众所周知,泊松分布是二项分布的极限分布,但在实际应用和观测中,观测数据值往往是有限的,如果采用 0-1 膨胀二项分布进行拟合,对于不同的观测数据集,可以选择不同的独立重复实验次数,这样就比选择 0-1 膨胀泊松分布进行拟合更加具有灵活性。鉴于以上特点,本文提出了 0-1 膨胀二项分布模型,基于数据扩充策略,运用客观贝叶斯方法对模型参数进行估计。

1 0-1 膨胀二项分布模型

本节提出 0-1 膨胀二项分布(ZOIB)模型,即一个非负的 ZOIB 的随机变量 Y 可以表示为 $Y = V(1 - B_1) + B_1(1 - B_2)$,其中 B_1, B_2, V 相互独立, B_1 服从于实验成功概率为 p_1 的伯努利分布; B_2 服从于实验成功概率为 p_2 的伯努利分布; V 服从于 m 次独立重复实验且成功概率为 θ 的二项分布,即 $P(V = k) = C_m^k \theta^k (1 - \theta)^{m-k}$, $k = 0, 1, \dots, m$ 。 Y 与 (B_1, B_2, V) 之间的关系如下:

$$\begin{cases} (Y = 0) \Leftrightarrow (V = 0, B_1 = 0) \cup (B_1 = 1, B_2 = 1) \\ (Y = 1) \Leftrightarrow (V = 1, B_1 = 0) \cup (B_1 = 1, B_2 = 0) \\ (Y = k) \Leftrightarrow (V = k, B = 0), k = 2, 3, \dots, m \end{cases} \quad (1)$$

则随机变量 Y 的分布律为:

$$P(Y = k) = \begin{cases} p_1 p_2 + (1 - p_1)(1 - \theta)^m & k = 0 \\ p_1(1 - p_2) + (1 - p_1)m\theta(1 - \theta)^{m-1} & k = 1 \\ (1 - p_1)C_m^k \theta^k (1 - \theta)^{m-k} & k = 2, 3, \dots, m \end{cases} \quad (2)$$

式中: $0 \leq p_1 \leq 1, 0 \leq p_2 \leq 1, 0 \leq \theta \leq 1$, 记 $Y \sim ZOIB(p_1, p_2, \theta)$, 可以看出 p_1 和 $1 - p_1$ 分别是一个伯努利分布与一个二项分布的混合比例。当 $p_2 = 1$ 时, ZOIB 变成零

膨胀二项分布;当 $p_1 = 0$, ZOIB 退化二项分布。

设 $Y = (Y_1, Y_2, \dots, Y_n)$ 是来自 0-1 膨胀二项分布容量为 n 的观测值,则 (p_1, p_2, θ) 的似然函数为:

$$L(p_1, p_2, \theta | Y) \propto [p_1 p_2 + (1 - p_1)(1 - \theta)^m]^{S_0} \times [p_1(1 - p_2) + (1 - p_1)m\theta(1 - \theta)^{m-1}]^{S_1} \times (1 - p_1)^{n - S_0 - S_1} \theta^S (1 - \theta)^{m(n - S_0 - S_1) - S} \quad (3)$$

式中: S_0 表示 $\{i: Y_i = 0\}$ 中元素的个数; S_1 表示 $\{i: Y_i = 1\}$ 中元素的个数; $S = \sum_{Y_i \geq 2} Y_i$ 。

由于式(3)很复杂,很难获得参数 (p_1, p_2, θ) 的 Jeffreys 先验和 reference 先验。因此本节基于隐变量 $B_1 = (B_{11}, B_{12}, \dots, B_{1n}), B_2 = (B_{21}, B_{22}, \dots, B_{2n}), V = (V_1, V_2, \dots, V_n)$, 构建基于数据扩充的完全似然函数,其表达式为:

$$L(p_1, p_2, \theta | Y, B_1, B_2, V) = \prod_{i=1}^n [p_1 p_2^{B_{2i}} (1 - p_2)^{1 - B_{2i}}]^{B_{1i}} [(1 - p_1) C_m^{V_i} \theta^{V_i} (1 - \theta)^{m - V_i}]^{1 - B_{1i}} = \prod_{i=1}^n p_1^{B_{1i}} (1 - p_1)^{1 - B_{1i}} p_2^{B_{1i} B_{2i}} (1 - p_2)^{B_{1i}(1 - B_{2i})} [C_m^{V_i} \theta^{V_i} (1 - \theta)^{m - V_i}]^{1 - B_{1i}} \quad (4)$$

相应的对数完全似然函数表示为:

$$l(p_1, p_2, \theta | Y, B_1, B_2, V) = \sum_{i=1}^n [B_{1i} \ln p_1 + (1 - B_{1i}) \ln(1 - p_1)] + \sum_{i=1}^n [B_{1i} B_{2i} \ln p_2 + B_{1i}(1 - B_{2i}) \ln(1 - p_2)] + \sum_{i=1}^n (1 - B_{1i}) [V_i \ln \theta + (m - V_i) \ln(1 - \theta) + \ln C_m^{V_i}] \quad (5)$$

2 客观贝叶斯分析

2.1 Fisher 信息矩阵

由于原始的对数似然函数式(3)非常复杂,不便于推导无信息先验,本节使用对数完全似然函数式(5)进行推导,得到 Fisher 信息矩阵是对角矩阵,极大地简化了客观先验的计算。

定理 1 对于 ZOIB 模型式(2),基于隐变量 (B_1, B_2, V) , 参数 $\Phi = (p_1, p_2, V)$ 的 Fisher 信息矩阵为:

$$I(\Phi) = \begin{bmatrix} \frac{n}{p_1(1 - p_1)} & 0 & 0 \\ 0 & \frac{np_1}{p_2(1 - p_2)} & 0 \\ 0 & 0 & \frac{mn(1 - p_1)}{\theta(1 - \theta)} \end{bmatrix} \quad (6)$$

证明 对数完全似然函数式(5)的二阶偏导数如下:

$$\begin{cases} \frac{\partial^2 l}{\partial p_1^2} = - \sum_{i=1}^n \left(\frac{B_{1i}}{p_1^2} + \frac{1-B_{1i}}{(1-p_1)^2} \right) \\ \frac{\partial^2 l}{\partial p_2^2} = - \sum_{i=1}^n \left(\frac{B_{2i} B_{2i}}{p_2^2} + \frac{B_{2i}(1-B_{2i})}{(1-p_2)^2} \right) \\ \frac{\partial^2 l}{\partial \theta^2} = - \sum_{i=1}^n (1-B_{1i}) \left(\frac{V_i}{\theta^2} + \frac{m-V_i}{(1-\theta)^2} \right) \end{cases}, \begin{cases} \frac{\partial^2 l}{\partial p_1 \partial p_2} = \frac{\partial^2 l}{\partial p_2 \partial p_1} = 0 \\ \frac{\partial^2 l}{\partial p_1 \partial \theta} = \frac{\partial^2 l}{\partial \theta \partial p_1} = 0 \\ \frac{\partial^2 l}{\partial p_2 \partial \theta} = \frac{\partial^2 l}{\partial \theta \partial p_2} = 0 \end{cases} \quad (7)$$

由于 $B_{1i} \sim \text{Bernoulli}(p_1)$, $B_{2i} \sim \text{Bernoulli}(p_2)$, $V_i \sim B(m, \theta)$, 则 $E(B_{1i}) = p_1$, $E(B_{2i}) = p_2$, $E(V_i) = m\theta$, $i = 1, 2, \dots, n$ 。因此, Fisher 信息矩阵中各元素计算如下:

$$\begin{cases} I_{11} = -E\left(\frac{\partial^2 l}{\partial p_1^2}\right) = \frac{n}{p_1(1-p_1)} \\ I_{22} = -E\left(\frac{\partial^2 l}{\partial p_2^2}\right) = \frac{np_1}{p_2(1-p_2)} \\ I_{33} = -E\left(\frac{\partial^2 l}{\partial \theta^2}\right) = \frac{mn(1-p_1)}{\theta(1-\theta)} \\ I_{12} = I_{21} = I_{13} = I_{31} = I_{23} = I_{32} = 0 \end{cases} \quad (8)$$

最终得到参数 $\Phi = (p_1, p_2, \theta)$ 的 Fisher 信息矩阵式(6)。

2.2 客观先验

与 Laplace 先验比较, Jeffreys 先验能够在参数变换下保持不变性,比 Laplace 先验具有更广泛的应用场合。

定理 2 参数 $\Phi = (p_1, p_2, \theta)$ 的 Jeffreys 先验为:

$$\pi_J \propto p_2^{-1/2} (1-p_2)^{-1/2} \theta^{-1/2} (1-\theta)^{-1/2} \quad (9)$$

证明 参数 $\Phi = (p_1, p_2, \theta)$ 的 Jeffreys 先验与 Fisher 信息矩阵行列式的平方根成正比:

$$\pi_J \propto \det(\mathbf{I})^{1/2} = (I_{11} I_{22} I_{33})^{1/2} \propto p_2^{-1/2} (1-p_2)^{-1/2} \theta^{-1/2} (1-\theta)^{-1/2} \quad (10)$$

由于在多维情形下,由 Jeffreys 先验得到的后验估计有时并不相合,因此, Bernardo^[10] 将参数分为感兴趣参数和讨厌参数,研究了在多维情形下对 Jeffreys 先验修正的方法,即基于信息量准则使参数先验分布和后验分布之间的 Kullback-Liebler 距离最大。这种从信息量准则出发推导出先验分布称为 reference 先验,当参数是一维时,reference 先验就变成了 Jeffreys 先验。

在统计推断中,模型的参数往往不止一个,总是有些参数比另外一些参数更重要,更值得人们去关注。根据不同的重要性,就会得到不同的 reference 先验。例如,记号 $\{(p_1, p_2), \theta\}$ 表示 2 组参数, p_1 与 p_2 之间的重要性相同,它们比 θ 更重要。再如,记号 $\{p_1, p_2, \theta\}$ 表示 3 组参数,重要性程度排序为 p_1, p_2, θ 依次递减。本节研究了 7 种不同的参数组合,并推导出相应的 reference 先验。

定理 3 对于参数组合 $\{(p_1, p_2), \theta\}$ 和 $\{\theta, (p_1, p_2)\}$, Φ 的 reference 先验为:

$$\pi_{R1} \propto (1-p_1)^{-1/2} p_2^{-1/2} (1-p_2)^{-1/2} \theta^{-1/2} (1-\theta)^{-1/2} \quad (11)$$

证明 对于参数组合 $\{(p_1, p_2), \theta\}$, (p_1, p_2) 是感兴趣的参数, Fisher 信息矩阵 $\mathbf{I}(\Phi)$ 可写为:

$$\Sigma_1 = \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & I_{33} \end{bmatrix} \quad (12)$$

$$\text{式中: } \Sigma_{11} = \begin{bmatrix} I_{11} & 0 \\ 0 & I_{22} \end{bmatrix}。$$

根据 Berger^[11], reference 先验求解过程中,需要找到关于参数 (p_1, p_2, θ) 的两个函数 h_1 和 h_2 , 再完成以下四个步骤。

$$\begin{aligned} h_1 &= \frac{\det(\Sigma_1)}{I_{33}} = \frac{n^2}{(1-p_1)p_2(1-p_2)} \\ h_2 &= I_{33} = \frac{mn(1-p_1)}{\theta(1-\theta)} \end{aligned} \quad (13)$$

步骤 1 选取参数空间的一组紧子集:

$$\begin{aligned} \Omega_i &= \Omega_{12i} \times \Omega_{3i} = \\ &\{(p_1, p_2) \mid a_{1i} < p_1 < b_{1i}, a_{2i} < p_2 < b_{2i}\} \times \\ &\{\theta \mid a_{3i} < \theta < b_{3i}\} \end{aligned} \quad (14)$$

使得 $a_{1i}, a_{2i}, a_{3i} \rightarrow 0, b_{1i}, b_{2i}, b_{3i} \rightarrow 1$ 。

步骤 2 当 (p_1, p_2) 给定时, θ 的条件先验为:

$$\pi_{R1}^i(\theta \mid p_1, p_2) = \frac{\sqrt{h_2} \Omega_{3i}}{\int_{\Omega_{3i}} \sqrt{h_2} d\theta} \propto \theta^{-1/2} (1-\theta)^{-1/2} \Omega_{3i} \quad (15)$$

步骤 3 (p_1, p_2) 的边缘先验为:

$$\pi_{R1}^i(p_1, p_2) = \frac{\exp\left\{\frac{1}{2} \int_{\Omega_{3i}} \pi_{R1}^i(\theta \mid p_1, p_2) \log(h_1) d\theta\right\} \Omega_{12i}}{\iint_{\Omega_{12i}} \exp\left\{\frac{1}{2} \int_{\Omega_{3i}} \pi_{R1}^i(\theta \mid p_1, p_2) \log(h_1) d\theta\right\} dp_1 dp_2} \propto (1-p_1)^{-1/2} p_2^{-1/2} (1-p_2)^{-1/2} \Omega_{12i} \quad (16)$$

步骤 4 Φ 的 reference 先验为:

$$\pi_{R1} = \lim_{i \rightarrow \infty} \frac{\pi_{R1}^i(p_1, p_2) \pi_{R1}^i(\theta \mid p_1, p_2)}{\pi_{R1}^i(p_1^*, p_2^*) \pi_{R1}^i(\theta^* \mid p_1^*, p_2^*)} \propto (1-p_1)^{-1/2} p_2^{-1/2} (1-p_2)^{-1/2} \theta^{-1/2} (1-\theta)^{-1/2} \quad (17)$$

式中: p_1^*, p_2^* 和 θ^* 是参数空间中事先给定的值。

对于参数组合 $\{\theta, (p_1, p_2)\}$, θ 是感兴趣的参数, Fisher 信息矩阵 $\mathbf{I}(\Phi)$ 可写为:

$$\Sigma_2 = \begin{bmatrix} I_{33} & 0 \\ 0 & \Sigma_{11} \end{bmatrix} \quad (18)$$

$$\text{得到 } h_1 = \frac{mn(1-p_1)}{\theta(1-\theta)}, h_2 = \frac{n^2}{(1-p_1)p_2(1-p_2)}。$$

步骤 1 选取与 $\{(p_1, p_2), \theta\}$ 参数空间中相同的一组紧子集。

步骤 2 当 θ 给定时, (p_1, p_2) 的条件先验为:

$$\pi_{R1}^i(p_1, p_2 | \theta) = \frac{\sqrt{h_2} \Omega_{12i}}{\iint_{\Omega_{12i}} \sqrt{h_2} dp_1 dp_2} \propto (1 - p_1)^{-1/2} p_2^{-1/2} (1 - p_2)^{-1/2} \Omega_{12i} \quad (19)$$

步骤 3 θ 的边缘先验为:

$$\pi_{R1}^i(\theta) = \frac{\exp\left\{\frac{1}{2} \iint_{\Omega_{12i}} \pi_{R1}^i(p_1, p_2 | \theta) \log(h_1) dp_1 dp_2\right\} \Omega_{3i}}{\int_{\Omega_{3i}} \exp\left\{\frac{1}{2} \iint_{\Omega_{12i}} \pi_{R1}^i(p_1, p_2 | \theta) \log(h_1) dp_1 dp_2\right\} d\theta} \propto \theta^{-1/2} (1 - \theta)^{-1/2} \Omega_{3i} \quad (20)$$

步骤 4 Φ 的 reference 先验为:

$$\pi_{R1} = \lim_{i \rightarrow \infty} \frac{\pi_{R1}^i(\theta) \pi_{R1}^i(p_1, p_2 | \theta)}{\pi_{R1}^i(\theta^*) \pi_{R1}^i(p_1^*, p_2^* | \theta^*)} \propto (1 - p_1)^{-1/2} p_2^{-1/2} (1 - p_2)^{-1/2} \theta^{-1/2} (1 - \theta)^{-1/2} \quad (21)$$

式中: p_1^* , p_2^* 和 θ^* 是参数空间中事先给定的值。

定理 4 对于参数组合 $\{(p_1, \theta), p_2\}$ 和 $\{p_2, (p_1, \theta)\}$, Φ 的 reference 先验为:

$$\pi_{R2} \propto p_1^{-1/2} p_2^{-1/2} (1 - p_2)^{-1/2} \theta^{-1/2} (1 - \theta)^{-1/2} \quad (22)$$

证明 方法与定理 3 相同,故省略。

定理 5 对于参数组合 $\{(p_2, \theta), p_1\}$, $\{p_1, (p_2, \theta)\}$ 和 $\{p_1, p_2, \theta\}$, Φ 的 reference 先验为:

$$\pi_{R3} \propto p_1^{-1/2} (1 - p_1)^{-1/2} p_2^{-1/2} (1 - p_2)^{-1/2} \theta^{-1/2} (1 - \theta)^{-1/2} \quad (23)$$

证明 参数组合 $\{(p_2, \theta), p_1\}$, $\{p_1, (p_2, \theta)\}$ reference 先验的证明方法与定理 3 和定理 4 相同,这里只对第 3 个参数组合 $\{p_1, p_2, \theta\}$ 的 reference 先验进行推导。将 Σ_1 的逆矩阵记为:

$$\Sigma_1^{-1} = \begin{bmatrix} s_{11} & 0 & 0 \\ 0 & s_{22} & 0 \\ 0 & 0 & s_{33} \end{bmatrix} \quad (24)$$

令 $s_1 = s_{11}$, $s_2 = \begin{bmatrix} s_{11} & 0 \\ 0 & s_{22} \end{bmatrix}$, $s_3 = s_{33}$, 根据 Berger

和 Bernardo^[12], 令 $H_j = s_j^{-1}$, h_j 是 H_j 右下角的元素, $j = 1, 2, 3$, 得到:

$$h_1 = \frac{n}{p_1(1-p_1)}, h_2 = \frac{np_1}{p_2(1-p_2)}, h_3 = \frac{mn(1-p_1)}{\theta(1-\theta)}$$

选取参数空间的一组交集 $\Omega_i = \Omega_{1i} \times \Omega_{2i} \times \Omega_{3i}$, 其中: $\Omega_{1i} = \{p_1 | a_{1i} < p_1 < b_{1i}\}$, $\Omega_{2i} = \{p_2 | a_{2i} < p_2 < b_{2i}\}$, $\Omega_{3i} = \{\theta | a_{3i} < \theta < b_{3i}\}$, 使得 $a_{1i}, a_{2i}, a_{3i} \rightarrow 0, b_{1i}, b_{2i}, b_{3i} \rightarrow 1$, 令 $\theta_1 = p_1, \theta_2 = p_2, \theta_3 = \theta$, 对于 $j = 1, 2, 3$, 令 $\theta_{[j]} = (\theta_1, \dots, \theta_j), \theta_{[-j]} = (\theta_{j+1}, \dots, \theta_3)$, 特别地, $\theta_{[-0]} = (\theta_1, \theta_2, \theta_3), \theta_{[0]}$ 为空。令:

$$\Theta^i(\theta_{[j]}) = \{\theta_{j+1} : (\theta_{[j]}, \theta_{j+1}, \theta_{[-j+1]}) \in \Theta^i\} \quad (25)$$

对于某 $\theta_{[-j+1]}$, 从 h_j 的表达式看出, 每一个 h_j 只依赖于 $\theta_{[j]}, j = 1, 2, 3$, 这样就比较容易得到:

$$\pi_{R3}^i(\theta_1, \theta_2, \theta_3) = \left\{ \prod_{j=1}^3 \frac{|h_j|^{1/2}}{\int_{\Theta^i(\theta_{[j-1]})} |h_j|^{1/2} d\theta_j} \right\} \Omega_i \quad (26)$$

对于 $j = 1, 2, 3$, 分别计算 $\int_{\Theta^i(\theta_{[j-1]})} |h_j|^{1/2} d\theta_j$, 结果如下:

$$\int_{\Theta^i(\theta_{[0]})} |h_1|^{1/2} d\theta_1 = \int_{a_{1i}}^{b_{1i}} \sqrt{\frac{n}{p_1(1-p_1)}} dp_1 = c_1 \quad (27)$$

$$\int_{\Theta^i(\theta_{[1]})} |h_2|^{1/2} d\theta_2 = \int_{a_{2i}}^{b_{2i}} \sqrt{\frac{np_1}{p_2(1-p_2)}} dp_2 = p_1^{1/2} c_2 \quad (28)$$

$$\int_{\Theta^i(\theta_{[2]})} |h_3|^{1/2} d\theta_3 = \int_{a_{3i}}^{b_{3i}} \sqrt{\frac{mn(1-p_1)}{\theta(1-\theta)}} d\theta = (1-p_1)^{1/2} c_3 \quad (29)$$

式中: c_1, c_2 和 c_3 都是常数。因此, Φ 的 reference 先验为: $\pi_{R3} \propto p_1^{-1/2} (1 - p_1)^{-1/2} p_2^{-1/2} (1 - p_2)^{-1/2} \theta^{-1/2} (1 - \theta)^{-1/2}$ (30)

2.3 后验分析

本节先证明基于先验分布 $\pi_J, \pi_{R1}, \pi_{R2}, \pi_{R3}$ 得到 (p_1, p_2, θ) 的后验分布是恰当的, 再通过隐变量的条件分布, 设计 Gibbs 的抽样机制, 获取高效的后验样本。

定理 6 基于先验分布 $\pi_J, \pi_{R1}, \pi_{R2}, \pi_{R3}$ 得到 (p_1, p_2, θ) 的后验分布都是恰当的。

证明 (p_1, p_2, θ) 的后验分布为:

$$\pi(p_1, p_2, \theta, B_1, B_2, V | Y) \propto L(p_1, p_2, \theta | Y, B_1, B_2, V) \pi(p_1, p_2, \theta) \quad (31)$$

式中: $Y = (Y_1, Y_2, \dots, Y_n)$ 为观测数据; $B_1 = (B_{11}, B_{12}, \dots, B_{1n}), B_2 = (B_{21}, B_{22}, \dots, B_{2n}), V = (V_1, V_2, \dots, V_n)$ 为隐变量。结合式(4), 并以 π_{R3} 为例:

$$\begin{aligned} \pi_{R3}(p_1, p_2, \theta | Y, B_1, B_2, V) \propto & p_1^{\sum_{i=1}^n B_{1i} - \frac{1}{2}} (1 - p_1)^{\sum_{i=1}^n (1 - B_{1i}) - \frac{1}{2}} \times \\ & p_2^{\sum_{i=1}^n B_{1i} B_{2i} - \frac{1}{2}} (1 - p_1)^{\sum_{i=1}^n B_{1i} (1 - B_{2i}) - \frac{1}{2}} \times \\ & \theta^{\sum_{i=1}^n V_i (1 - B_{1i}) - \frac{1}{2}} (1 - \theta)^{\sum_{i=1}^n (m - V_i) (1 - B_{1i}) - \frac{1}{2}} \end{aligned} \quad (32)$$

由于:

$$\int_0^1 p_1^{\sum_{i=1}^n B_{1i} - \frac{1}{2}} (1 - p_1)^{\sum_{i=1}^n (1 - B_{1i}) - \frac{1}{2}} dp_1 = \frac{\Gamma\left(\sum_{i=1}^n B_{1i} + \frac{1}{2}\right) \Gamma\left(\sum_{i=1}^n (1 - B_{1i}) + \frac{1}{2}\right)}{\Gamma(n + 1)} \quad (33)$$

$$\int_0^1 p_2^{\sum_{i=1}^n B_{1i} B_{2i} - \frac{1}{2}} (1 - p_2)^{\sum_{i=1}^n B_{1i} (1 - B_{2i}) - \frac{1}{2}} dp_2 = \frac{\Gamma\left(\sum_{i=1}^n B_{1i} B_{2i} + \frac{1}{2}\right) \Gamma\left(\sum_{i=1}^n B_{1i} (1 - B_{2i}) + \frac{1}{2}\right)}{\Gamma\left(\sum_{i=1}^n B_{1i} + 1\right)} \quad (34)$$

$$\int_0^1 \theta \sum_{i=1}^n V_i (1 - B_{1i})^{-\frac{1}{2}} (1 - \theta) \sum_{i=1}^n (m - V_i) (1 - B_{1i})^{-\frac{1}{2}} d\theta =$$

$$\frac{\Gamma\left(\sum_{i=1}^n V_i (1 - B_{1i}) + \frac{1}{2}\right) \Gamma\left(\sum_{i=1}^n (m - V_i) (1 - B_{1i}) + \frac{1}{2}\right)}{\Gamma\left(\sum_{i=1}^n m (1 - B_{1i}) + \frac{1}{2}\right)} \quad (35)$$

因此,得到 $\int_0^1 \int_0^1 \int_0^1 \pi_{R3}(p_1, p_2, \theta | Y, B_1, B_2, V) dp_1 dp_2 d\theta < \infty$, 从而证明了基于先验分布 π_{R3} 得到的 (p_1, p_2, θ) 的后验分布是恰当的。类似的方法,基于先验分布 π_J, π_{R1} 和 π_{R2} 得到 (p_1, p_2, θ) 的后验分布也都是恰当的。

2.4 抽样算法设计

引理 1^[7] 设 $Y = (Y_1, Y_2, \dots, Y_n)$ 是来自 0-1 膨胀二项分布容量为 n 的观测值, $B_1 = (B_{11}, B_{12}, \dots, B_{1n})$, $B_2 = (B_{21}, B_{22}, \dots, B_{2n})$, $V = (V_1, V_2, \dots, V_n)$ 为隐变量, (B_1, B_2, V) 的条件分布 $\pi(B_1, B_2, V | Y, p_1, p_2, \theta)$ 为:

$$P(B_1 = i, B_2 = j, V = v | Y = 0, p_1, p_2, \theta) =$$

$$\begin{cases} \frac{(1-p_1)(1-p_2)P(V=0)}{p_1 p_2 + (1-p_1)P(V=0)} & i=0, j=0, v=0 \\ \frac{(1-p_1)p_2 P(V=0)}{p_1 p_2 + (1-p_1)P(V=0)} & i=0, j=1, v=0 \\ \frac{p_1 p_2 P(V=v)}{p_1 p_2 + (1-p_1)P(V=0)} & i=1, j=1, v=0, 1, \dots, m \\ 0 & \text{其他} \end{cases} \quad (36)$$

$$P(B_1 = i, B_2 = j, V = v | Y = 1, p_1, p_2, \theta) =$$

$$\begin{cases} \frac{(1-p_1)(1-p_2)P(V=1)}{p_1(1-p_2) + (1-p_1)P(V=1)} & i=0, j=0, v=1 \\ \frac{(1-p_1)p_2 P(V=1)}{p_1(1-p_2) + (1-p_1)P(V=1)} & i=0, j=1, v=1 \\ \frac{p_1(1-p_2)P(V=v)}{p_1(1-p_2) + (1-p_1)P(V=1)} & i=1, j=0, v=0, 1, \dots, m \\ 0 & \text{其他} \end{cases} \quad (37)$$

$$P(B_1 = i, B_2 = j, V = v | 2 \leq Y \leq m, p_1, p_2, \theta) =$$

$$\begin{cases} 1 - p_2 & i=0, j=0, v=2, 3, \dots, m \\ p_2 & i=0, j=1, v=2, 3, \dots, m \\ 0 & \text{其他} \end{cases} \quad (38)$$

首先,由引理 1,当观测数据 $Y = (Y_1, Y_2, \dots, Y_n)$ 给定时,得到 (B_1, B_2, V) 的样本。

(1) 当 $Y_i = 0$ 时,由式(36),抛掷一枚硬币,其正面朝上的概率为 $\frac{p_1 p_2}{p_1 p_2 + (1-p_1)(1-\theta)^m}$ 。当正面朝上

时,令 $B_{1i} = 1, B_{2i} = 1, V_i$ 通过二项分布抽样得到;当反面朝上时,令 $B_{1i} = 0, V_i = 0$ 。此时,再抛掷另一枚硬币,其正面朝上概率为 p_2 ,当正面朝上时,令 $B_{2i} = 1$,否则,令 $B_{2i} = 0$ 。

(2) 当 $Y_i = 1$ 时,由式(37),抛掷一枚硬币,其正面朝上的概率为 $\frac{p_1(1-p_2)}{p_1 p_2 + (1-p_1)m\theta(1-\theta)^{m-1}}$ 。当正面朝上时,令 $B_{1i} = 1, B_{2i} = 0, V_i$ 通过二项分布抽样得到;当反面朝上时,令 $B_{1i} = 0, V_i = 1$ 。此时,再抛掷另一枚硬币,其正面朝上概率为 p_2 ,当正面朝上时,令 $B_{2i} = 1$,否则,令 $B_{2i} = 0$ 。

(3) 当 $Y_i = v$ 时, $v = 2, 3, \dots, m$,由式(38),抛掷一枚硬币,其正面朝上的概率为 p_2 ,当正面朝上时,令 $B_{1i} = 0, B_{2i} = 1, V_i = v$;当反面朝上时,令 $B_{1i} = 0, B_{2i} = 0, V_i = v$ 。

其次,由式(32),分别得到参数 p_1, p_2, θ 的满条件概率分布,均为常见的概率分布,从而更加容易实现对后验分布的抽样,即:

$$p_1 | p_2, \theta, Y, B_1, B_2, V \sim \text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n B_{1i}, n + \frac{1}{2} - \sum_{i=1}^n B_{1i}\right) \quad (39)$$

$$p_2 | p_1, \theta, Y, B_1, B_2, V \sim \text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n B_{1i} B_{2i}, \frac{1}{2} + \sum_{i=1}^n B_{1i} (1 - B_{2i})\right) \quad (40)$$

$$\theta | p_1, p_2, Y, B_1, B_2, V \sim \text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n V_i (1 - B_{1i}), \frac{1}{2} + \sum_{i=1}^n (m - V_i) (1 - B_{1i})\right) \quad (41)$$

最后,实施 Gibbs 抽样,具体步骤如下:

(1) 设定参数初始值 $p_1^{(0)}, p_2^{(0)}, \theta^{(0)}$ 。

(2) 对 $t = 1, 2, \dots, n$,进行以下迭代:

① 利用参数估计值 $(p_1^{(t-1)}, p_2^{(t-1)}, \theta^{(t-1)})$ 得到样本 $(B_{1i}^{(t)}, B_{2i}^{(t)}, V_i^{(t)})$, $i = 1, 2, \dots, n$;

② 由 $\text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n B_{1i}, n + \frac{1}{2} - \sum_{i=1}^n B_{1i}\right)$ 抽样得到 $p_1^{(t)}$;

③ 由 $\text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n B_{1i} B_{2i}, \frac{1}{2} + \sum_{i=1}^n B_{1i} (1 - B_{2i})\right)$ 抽样得到 $p_2^{(t)}$;

④ 由 $\text{Beta}\left(\frac{1}{2} + \sum_{i=1}^n V_i (1 - B_{1i}), \frac{1}{2} + \sum_{i=1}^n (m - V_i) (1 - B_{1i})\right)$ 抽样得到 $\theta^{(t)}$ 。

3 数值模拟

本节基于三种 reference 先验,通过数值模拟对

ZOIB 分布的参数进行估计。样本容量分别设为 $n = 20$ 和 $n = 50$,二项分布中 m 的值设为 10, θ 的值设为 0.8, p_1 的值分别设为 0.3 和 0.7, p_2 的值分别设为 0.4 和 0.6,所有模拟重复 5 000 次,分别计算了参数估计量的均值和均方误差,如表 1 和表 2 所示,可以看出,随着样本容量的增加,三种客观贝叶斯先验下的估计值越来越接近真值,均方误差也越来越小。三种方法对 θ 和 p_2 的估计效果相当,但对于 p_1 的估计, π_{R3} 比 π_{R1} 和 π_{R2} 表现要更好,这是因为在 π_{R3} 中包含 p_1 的信息更加丰富。

表 1 和下参数估计量的均值

p_1	p_2	n	θ_{R1}	p_{1R1}	p_{2R1}	θ_{R2}	p_{1R2}
0.3	0.4	20	0.768	0.283	0.362	0.783	0.279
		50	0.779	0.288	0.374	0.787	0.284
	0.6	20	0.776	0.282	0.566	0.785	0.331
		50	0.781	0.287	0.584	0.791	0.319
0.7	0.4	20	0.772	0.672	0.354	0.786	0.669
		50	0.787	0.685	0.376	0.795	0.681
	0.6	20	0.775	0.734	0.645	0.784	0.724
		50	0.791	0.721	0.626	0.793	0.716

p_1	p_2	n	p_{2R2}	θ_{R3}	p_{1R3}	p_{2R3}
0.3	0.4	20	0.357	0.786	0.282	0.381
		50	0.388	0.791	0.294	0.392
	0.6	20	0.504	0.769	0.287	0.578
		50	0.557	0.793	0.296	0.581
0.7	0.4	20	0.346	0.742	0.686	0.374
		50	0.378	0.788	0.694	0.389
	0.6	20	0.645	0.768	0.672	0.633
		50	0.621	0.793	0.695	0.612

表 2 和下参数估计量的均方误差

p_1	p_2	n	θ_{R1}	p_{1R1}	p_{2R1}	θ_{R2}	p_{1R2}
0.3	0.4	20	0.083	0.074	0.098	0.081	0.077
		50	0.065	0.037	0.075	0.067	0.035
	0.6	20	0.076	0.072	0.093	0.078	0.082
		50	0.061	0.036	0.062	0.058	0.044
0.7	0.4	20	0.087	0.045	0.096	0.086	0.056
		50	0.056	0.038	0.074	0.066	0.042
	0.6	20	0.065	0.037	0.083	0.075	0.041
		50	0.057	0.025	0.068	0.066	0.035

续表 2

p_1	p_2	n	p_{2R2}	θ_{R3}	p_{1R3}	p_{2R3}
0.3	0.4	20	0.088	0.078	0.057	0.068
		50	0.085	0.052	0.026	0.055
	0.6	20	0.092	0.076	0.068	0.072
		50	0.073	0.043	0.034	0.053
0.7	0.4	20	0.086	0.082	0.053	0.066
		50	0.073	0.063	0.041	0.043
	0.6	20	0.082	0.065	0.048	0.072
		50	0.073	0.046	0.025	0.053

4 应用实例

COVID-19 威胁了全世界人民的健康与安全。因此,掌握 COVID-19 病例的分布规律及相关影响因素是传染病溯源追踪的重要环节,也是积极实施防控策略的有效途径。

本文采用湖北省 17 个城市(武汉、黄石、十堰、宜昌、襄阳、鄂州、荆门、孝感、荆州、黄冈、咸宁、随州、恩施、利川、仙桃、潜江和天门) COVID-19 疫情数据进行实例分析。样本数据来源于湖北省卫健委。基于流行病学机理:在大面积切断传播途径的 5 至 7 天后,感染人数与死亡人数会出现好转。因为武汉市于 2020 年 1 月 23 日封城,所以选取封城后的第 5 天,即 1 月 28 日为样本。2 月 22 日雷神山、火神山医院以及各方舱医院建设完毕,并且各地派往湖北省的医疗团队基本全部到位,所以再选取 2 月 22 日为样本。综上,选取具有代表性的这两天,即 1 月 28 日和 2 月 22 日的死亡人数为研究对象。

4.1 样本数据描述

如图 1 所示,1 月 28 日武汉市死亡人数最多,为 19 人;鄂州、荆门、黄冈、天门,死亡人数为 1 人;黄石、十堰、宜昌、襄阳、荆州、咸宁、随州、恩施、利川、仙桃和潜江死亡人数最少,为 0 人。2 月 22 日武汉市死亡人数依旧最多,为 82 人;黄石、鄂州、荆门、随州和天门五个城市,死亡人数为 1 人;十堰、宜昌、荆州、咸宁、恩施、利川、仙桃和潜江,死亡人数最少,为 0 人。由此可见,在上述时空条件下,湖北省各地区的 COVID-19 死亡人数具有显著的 0-1 膨胀现象。

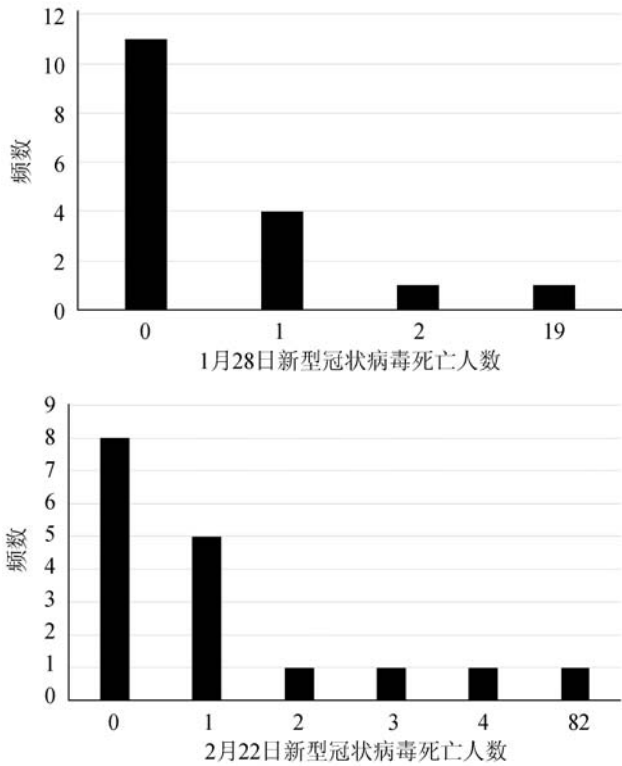


图1 样本数据的分布情况

4.2 拟合结果分析

根据 COVID-19 死亡人数的膨胀现象,本文采用 0-1 膨胀二项分布模型进行数据拟合,分别指定 $m = 19$ 和 $m = 82$ 。根据不同数据集指定 m ,比使用 0-1 膨胀泊松分布模型或者 0-1 膨胀负二项分布模型更加灵活。

三种不同客观先验下参数的点估计与区间估计如表 3 和表 4 所示,可以看出三种方法对参数的点估计效果相当,对于置信水平为 0.95 的区间估计,基于先验分布 π_{R3} 的客观贝叶斯分析更为精确。无论是从观测值与拟合值的接近程度(如图 2 和图 3 所示),还是 AIC(Akaike Information Criterion)的角度(如表 3 和表 4 所示),基于 π_{R3} 先验分布下的贝叶斯分析都能够实现更好的拟合效果。

表3 1月28日三种不同客观先验下参数的点估计与区间估计

	π_{R1}	π_{R2}	π_{R3}
p_1	0.548	0.542	0.543
	(0.237, 0.877)	(0.217, 0.962)	(0.221, 0.895)
p_2	0.435	0.435	0.435
	(0.219, 0.752)	(0.227, 0.740)	(0.223, 0.697)
θ	0.748	0.742	0.745
	(0.493, 0.880)	(0.492, 0.933)	(0.521, 0.862)
AIC	38.24	38.69	35.58

表4 2月22日三种不同客观先验下参数的点估计与区间估计

	π_{R1}	π_{R2}	π_{R3}
p_1	0.527	0.529	0.525
	(0.225, 0.867)	(0.215, 0.879)	(0.211, 0.886)
p_2	0.425	0.424	0.426
	(0.232, 0.741)	(0.225, 0.739)	(0.221, 0.688)
θ	0.746	0.743	0.747
	(0.494, 0.899)	(0.491, 0.913)	(0.511, 0.858)
AIC	37.95	37.52	36.87

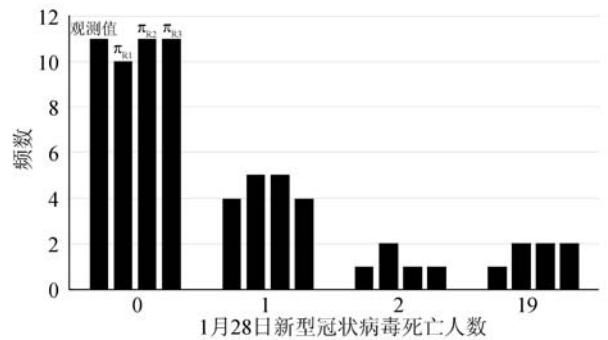


图2 1月28日三种不同客观先验下的拟合值

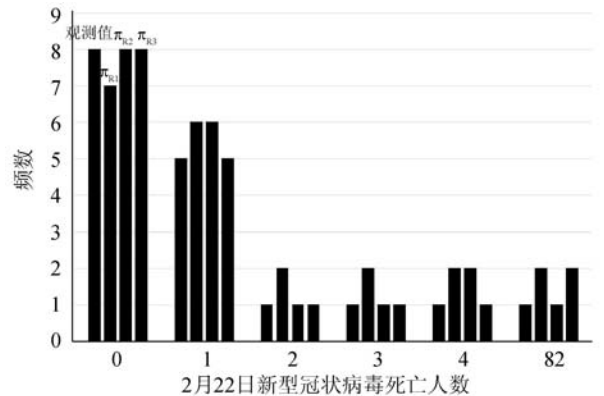


图3 2月22日三种不同客观先验下的拟合值

5 结语

本文对 0-1 膨胀二项分布模型进行了客观贝叶斯分析,巧妙地利用基于隐变量的完全似然函数写出了参数的 Fisher 信息矩阵。由于 Fisher 信息矩阵是对角矩阵,很容易推导出参数的 Jeffreys 先验和 reference 先验。采用 WinBUGS 软件和 R 软件进行数值模拟,设定不同的样本量和参数真值,对不同的无信息先验进行评估。最后分别对 2020 年 1 月 28 日与 2 月 22 日湖北省 COVID-19 死亡人数进行了实例分析,结果表明,在小样本情形下基于客观先验 π_{R3} 的拟合效果比 π_{R1} 和 π_{R2} 要好。在实际应用中,当数据集服从其他形式的 0-1 膨胀分布模型时,也可以采用本文所提出的

(下转第 59 页)

综上所述,使用 MATLAB 系统辨识工具箱准确辨识了变频循环泵的系统模型,并对系统进行了控制优化,为进一步结合预测结果与经济背压进行控制打下了坚实的基础。

4 结 语

本文通过分析背压影响因素,使用两种数据处理方法,对数据进行降维,在 Python 平台上使用循环神经网络进行预测,得到了精度较高的背压预测模型。预测结果均优于^[1-9]的预测结果。本文还分析了间接空冷的控制手段,并且选取变频循环泵作为其控制对象,依靠数据对变频循环泵进行建模。并使用模糊 PID 进行控制优化,优化了其上升时间和超调量。对未来将经济背压引入,精确计算实际条件下循环泵出力具有重要的意义。

参 考 文 献

- [1] 葛晓霞,赵舒莹,肖洪闯,等. 基于改进果蝇算法优化 SVM 的凝汽器真空预测[J]. 热能动力工程,2020,35(11):39-45.
- [2] 王建国,林乐平. 粒子群算法与径向神经网络相结合的凝汽器真空预测模型[J]. 热力发电,2015,44(10):72-76.
- [3] 张利平,陈浩天,王伟锋,等. 应用 PSO 算法改进 Elman 神经网络的双压凝汽器真空预测[J]. 热力发电,2015,44(3):53-57.
- [4] 田松峰,吴昭延,王子光,等. 基于神经网络的凝汽器污垢热阻预测模型[J]. 热力发电,2019,48(2):78-82.
- [5] 夏琳琳,台金娟,刘惠敏,等. 权重提取与 Dempster 多重融合的凝汽器真空预测[J]. 沈阳工业大学学报,2015,37(3):329-334.
- [6] 陈婷,孟娜,王建国. 凝汽器真空的监测模型[J]. 东北电力大学学报,2012,32(3):54-58.
- [7] 夏琳琳,台金娟,文磊. 基于组合证据冲突修正的凝汽器真空度多网络融合预测[J]. 化工自动化及仪表,2015,42(12):1331-1335.
- [8] 王国涛. 循环水泵变频改造节能分析[J]. 能源与节能,2020(7):51-52.
- [9] 靖长财,张冬青. 某 1000MW 汽轮机循环水泵变频存在问题及运行优化探讨[J]. 能源科技,2020,18(2):44-46.
- [10] 弗朗索瓦肖莱. Python 深度学习[M]. 张亮,译. 北京:人民邮电出版社,2018.
- [11] 伊德里斯. Python 数据分析[M]. 韩波,译. 北京:人民邮电出版社,2016.
- [12] Stepanek H. Thinking in pandas[M]. Berkeley: Apress,2020.
- [13] Naik G R. Advances in principal component analysis[M]. Singapore: Springer,2018.
- [14] 武彬,张栾英. 模糊自整定 PID 控制在主汽温控制中的应用[J]. 计算机仿真,2015,32(2):387-390.
- [15] 吴凡,李伟雄. 基于 MATLAB 系统辨识工具的系统辨识[J]. 河北农机,2016(11):59-60.

(上接第 52 页)

方法,便捷地写出不同的客观先验,并找到最优先验,从而得到较为精确的研究结果。此外,这种思想和方法还可以推广到混合治愈模型中去,这正是本文今后研究的方向。

参 考 文 献

- [1] Lambert D. Zero-inflated Poisson regression, with an application to defects in manufacturing[J]. Technometrics,1992,34(1):1-14.
- [2] Xie F C, Wei B C, Lin J G. Score tests for zero-inflated generalized Poisson mixed regression models[J]. Computational Statistics & Data Analysis,2009,53(9):3478-3489.
- [3] 罗付岩,赵佳星. 基于零膨胀模型的我国银企关系规模影响因素研究[J]. 数理统计与管理,2017,36(2):351-360.
- [4] 李二倩,田茂再. 零膨胀泊松分布参数的固定宽度置信区间构造方法[J]. 应用概率统计,2018,34(1):49-74.
- [5] Maria M, Christina O. Is visiting the dentist a good habit? Analyzing count data with excess zeros and excess ones[M]// Working Paper, Umea Economic Studies. Västerbotten; Umeå universitet,1999.
- [6] 田震. 零一膨胀回归模型及其统计诊断[D]. 昆明:云南大学,2016.
- [7] Tang Y C, Liu W C, Xu A C. Statistical inference for zero-and-one-inflated Poisson models[J]. Statistics Theory and Related Fields,2017,1(2):216-226.
- [8] Liu W C, Tang Y C, Xu A C. A zero-and-one inflated Poisson model and its application[J]. Statistics and its Interface,2018,11(2):339-351.
- [9] 夏丽丽,田茂再. 零一膨胀泊松回归模型的非参数统计分析及其应用[J]. 数理统计与管理,2019,38(2):235-246.
- [10] Bernardo J M. Reference posterior distributions for Bayesian inference[J]. Journal of the Royal Statistical Society,1979,41(2):113-128.
- [11] Berger J O, Bernardo J M. On the development of the reference prior method[J]. Bayesian Statistics,1992,4(4):35-60.
- [12] Berger J O, Bernardo J M. Reference priors in a variance components problem[J]. Bayesian Analysis in Statistics and Econometrics,1992,177-194.