

散度偏大计数数据分布的贝叶斯估计 *

陈 雪 东

(云南大学数学统计学院, 昆明, 650091; 湖州师范学院理学院, 湖州, 313000)

摘 要

本文针对索赔次数数据的特点, 讨论了两类可导致散度偏大特征数据的分布类型: 零点膨胀分布与膨胀参数分布, 并根据 Bayes 理论与 MCMC 方法, 利用 WinBUGS 对其进行建模和抽样. 经过比较, 给出了实现分布拟合的途径, 最后通过两个数值例子加以展示.

关键词: 散度偏大, 零点膨胀分布, 膨胀参数分布, MCMC 方法.

学科分类号: O212.8.

§1. 引 言

1992 年, Lambert^[1] 提出了零点膨胀 Poisson 回归 (Zero-inflated Poisson Regression, ZIP) 模型, 用来对具有较多零次记录的数据进行统计分析. 与一般服从 Poisson 分布的计数数据相比, 上述数据的方差大于均值, 这种特征被称为散度偏大 (Over-dispersion)^[2]. 通常情况下, 对于散度偏大数据的研究较多采用负二项分布与广义 Poisson 分布^[3] 来刻画其统计特征. 近年来, 对该问题的研究, 在生物学^[4], 医学^[5] 以及保险精算领域^[6] 受到很大的重视, 针对不同的数据特点和背景, 一些新的模型和方法被提出, 比较有代表性的如广义线性混合模型 (GLMM)^[3], 广义 Poisson 回归模型^[7], 随机效应模型^[6], 条件自回归模型以及 EM 算法^[12] 等等.

在保险精算领域, 索赔次数分布的研究具有非常重要的地位. 在保险实务中, 一方面, 由于保险公司和投保人风险意识的增强以及保险公司采用了如免赔制度、无赔款折扣 (NCD) 制度等风险回避机制, 使得索赔次数在零点处有更大的概率, 产生散度偏大现象; 另一方面, 保单组合由于受诸多因素的影响, 其风险往往具有非同质性, 这又使得索赔次数分布更为复杂. 鉴于上述情况, 在研究散度偏大索赔次数时, 除了某些由常用分布生成的零点膨胀分布 (Zero-inflated Distribution, ZID) 外, 本文引入另一类具有膨胀参数的分布 (Inflated-parameter Distribution, IPD) 并从生成过程进行实质性比较. 在统计推断时, 通常方法处理上述分布很困难且复杂, 与之相比, 贝叶斯推断中的 MCMC 方法显得行之有效, 特别是专门的 MCMC 抽样工具 WinBUGS 的使用, 使得完全条件分布的选择和模型的建立更加方便灵活. 为了避免多参数与复杂分布计算的困难, 本文借鉴了 Ioannis, N., Petros, D.^[8] 和 Scollnik, P.M.^[10] 的思想, 采用 MCMC 方法在 WinBUGS 上进行抽样和参数估计, 经过模型间的比较, 给出实现分布拟合的途径. 最后, 以两类常见索赔次数分布的数值例子作为例证.

* 国家自然科学基金 (10561008), 云南省自然科学基金 (2004A0002M) 和浙江省自然科学基金 (Y606667) 资助项目.
本文 2005 年 5 月 23 日收到, 2006 年 1 月 12 日收到修改稿.

§ 2. 两类可导致散度偏大计数数据的分布

2.1 由已知分布 $\pi(x|\theta)$ 生成的零点膨胀分布 (Zero-inflated Distributions)

设 $\pi(x|\theta)$, $x = 0, 1, 2, \dots$ 是已知的具有未知参数 θ 的某类计数数据的概率分布, 通常的零点膨胀 (Zero-inflated) 分布定义为

$$P(Y = y|\rho, \theta) = \begin{cases} \rho + (1 - \rho)\pi(y = 0|\theta), & y = 0; \\ (1 - \rho)\pi(y|\theta), & y > 0. \end{cases} \quad (1)$$

显然, 该分布由已知分布 $\pi(x|\theta)$ 和 ρ 确定, 其中 $\pi(x|\theta)$ 称为生成分布, $\rho \in [0, 1]$ 为零点膨胀参数, 对应的均值和方差可以表示为

$$E(Y|\rho, \theta) = (1 - \rho)E_{\pi}(Y|\theta) \quad (2)$$

和

$$\text{Var}(Y|\rho, \theta) = \rho(1 - \rho)[E_{\pi}(Y|\theta)]^2 + (1 - \rho)\text{Var}_{\pi}(Y|\theta). \quad (3)$$

于是, 当生成分布分别为 Poisson 分布, 二项分布, 混合 Poisson 分布, 广义 Poisson 分布, 负二项分布时, 由 (1) 式就可以得到相应的 Zero-inflated 分布, 分别简记为 ZIP, ZIB, ZIMP, ZIGP, ZINB 分布, 它们具有的散度偏大性可通过 (2), (3) 式得到, 为简略起见, 仅给出后三类分布的情形, 具体为

- ZIMP 分布: 均值, 方差分别为

$$\begin{cases} E(Y|\rho, \mu, \lambda_1, \lambda_2) = (1 - \rho)[\mu\lambda_1 + (1 - \mu)\lambda_2], \\ \text{Var}(Y|\rho, \mu, \lambda_1, \lambda_2) = (1 - \rho)[\mu\lambda_1 + (1 - \mu)\lambda_2] + \rho(1 - \rho)[\mu\lambda_1 + (1 - \mu)\lambda_2]^2 \\ \quad + (1 - \rho)\mu(1 - \mu)(\lambda_1 - \lambda_2)^2. \end{cases}$$

- ZIGP 分布: 均值, 方差分别为

$$E(Y|\rho, \lambda, \mu) = (1 - \rho)\lambda, \quad \text{Var}(Y|\rho, \lambda, \mu) = (1 - \rho)\lambda\left[\rho\lambda + \frac{1}{(1 - \mu)^2}\right].$$

- ZINB 分布: 均值, 方差分别为

$$E(Y|\rho, r, \beta) = (1 - \rho)\beta r, \quad \text{Var}(Y|\rho, r, \beta) = (1 - \rho)\beta r(1 + \beta).$$

2.2 由膨胀参数几何分布生成的膨胀参数分布 (Inflated-parameter Distributions)

根据 Zero-inflated 分布的定义, 设 X 为服从 $\pi(x|\theta)$ 的非负随机变量, 记 $\pi(0|\theta) = \pi_0$, $G_X(t) = E(t^X)$ 为 X 的母函数, 则由 X 生成的具有参数 $\rho \in (0, 1)$ 的 Zero-inflated 分布可表示为

$$P(Y = 0) = \rho + (1 - \rho)\pi_0; \quad P(Y = j) = (1 - \rho)\pi(j|\theta), \quad j = 1, 2, \dots \quad (4)$$

容易知道,

$$G_Y(t) = \rho + (1 - \rho)G_X(t).$$

现对 Y 的分布做进一步的分析, 由 (4) 式可知 $P(Y \geq 1) = (1 - \rho)(1 - \pi_0)$, 如果撇开 (4) 式中的 $\pi(j|\theta)$, 假设 $Y = j, j \geq 1$ 的概率随 j 呈几何递减, 则可由

$$P(Y \geq 1) = (1 - \rho)(1 - \pi_0) = \sum_{j=1}^{\infty} (1 - \pi_0)(1 - \rho)^2 \rho^{j-1}$$

得到

$$P(Y = 0) = \pi_0 + \rho(1 - \pi_0); \quad P(Y = j) = (1 - \pi_0)(1 - \rho)^2 \rho^{j-1}, \quad j = 1, 2, \dots$$

为了简单起见, 令 $(1 - \pi_0)(1 - \rho) = \mu$, 上式表示为

$$P(Y = 0) = 1 - \mu; \quad P(Y = j) = \mu(1 - \rho)\rho^{j-1}, \quad j = 1, 2, \dots$$

该分布称为膨胀参数几何 (Inflated-parameter Geometric, IPG) 分布^[7], 膨胀参数为 ρ , 且母函数为

$$G(t) = \sum_{j=0}^{\infty} t^j P(Y = j) = 1 - \frac{\mu(1-t)}{1-\rho t}.$$

设 $X_j (j = 1, 2, \dots, r)$ 独立同分布且服从上述 IPG 分布, 则 $S_r = \sum_{j=1}^r X_j$ 的母函数为

$$G_{S_r}(t) = \left(1 - \frac{\mu(1-t)}{1-\rho t}\right)^r. \quad (5)$$

称 S_r 服从的分布为由 IPG 分布生成的膨胀参数负二项 (Inflated-parameter Negative Binomial, IPNB) 分布, 分布列为

$$\begin{cases} P(S_r = 0|r, \mu, \rho) = (1 - \mu)^r, \\ P(S_r = j|r, \mu, \rho) = (1 - \mu)^r \sum_{i=1}^j \binom{j-1}{i-1} \binom{r+i-1}{i} [\mu(1-\rho)]^i \rho^{j-i}, \quad j = 1, 2, \dots \end{cases} \quad (6)$$

其散度偏大现象表现为

$$E(S_r|r, \mu, \rho) = \frac{r\mu}{(1-\mu)(1-\rho)}; \quad \text{Var}(S_r|r, \mu, \rho) = \frac{r\mu[1 + (1-\mu)\rho]}{(1-\mu)^2(1-\rho)^2}.$$

当 $r \rightarrow \infty$ 时, 记 $S = \sum_{j=1}^{\infty} X_j$, 若 $r \rightarrow \infty, \mu \rightarrow 0, r\mu \rightarrow \lambda$, 则由 (5) 式有

$$G_S(t) = \exp\left(\frac{\lambda(t-1)}{1-\rho t}\right).$$

该分布称为由 IPG 分布生成的膨胀参数 Poisson (Inflated-parameter Poisson, IPP) 分布, 分布列为

$$\begin{cases} P(S = 0|\rho, \lambda) = e^{-\lambda}, \\ P(S = j|\rho, \lambda) = e^{-\lambda} \sum_{i=1}^j \frac{1}{i!} \binom{j-1}{i-1} [\lambda(1-\rho)]^i \rho^{j-i}, \quad j = 1, 2, \dots \end{cases} \quad (7)$$

同样, 其散度偏大现象表现为

$$E(S|\rho, \lambda) = \frac{\lambda}{(1-\rho)}; \quad \text{Var}(S|\rho, \lambda) = \frac{\lambda(1+\rho)}{(1-\rho)^2}.$$

§ 3. 基于 ZID 和 IPD 的索赔数据的 MCMC 模型

3.1 一般索赔次数数据

设随机变量 X 表示索赔次数 (Number of Claims), 其分布为 $p(x|\theta)$, θ 为未知参数. 一般的索赔次数的频数分布表为

Number of Claims	0	1	2	3	...	$\geq k-1$
Number of Policiers	n_1	n_2	n_3	n_4	...	n_k

设上述数据存在散度偏大特征, 考虑用零点膨胀分布中的 ZIMP, ZIGP, ZINB 和膨胀参数分布中的 IPNB, IPP 来对这些数据进行拟合.

记 $N = \sum_{i=1}^k n_i$, 一般的, 若概率

$$P(X = i - 1) = p_i, \quad i = 1, 2, \dots, k - 1; \quad P(X \geq k - 1) = p_k$$

已知, 由于 $p_k = 1 - \sum_{i=1}^{k-1} p_i$, 则 (n_1, n_2, \dots, n_k) 服从多项分布, 样本的似然函数为

$$L(n_1, \dots, n_k | p_1, \dots, p_k, N) = \frac{N!}{\prod_{i=1}^k n_i!} \prod_{i=1}^k p_i^{n_i} \propto \prod_{i=1}^k p_i^{n_i}. \quad (8)$$

于是给出观察值 (n_1, n_2, \dots, n_k) 后, 关于未知参数 θ 的后验分布为

$$\begin{aligned} f(\theta | n_1, \dots, n_k) &\propto f(n_1, \dots, n_k | p_1, \dots, p_k) f(p_1, \dots, p_k | \theta) f(\theta) \\ &\propto \frac{N!}{\prod_{i=1}^k n_i!} \prod_{i=1}^k p_i^{n_i} \prod_{i=1}^k f(p_i | \theta) f(\theta). \end{aligned}$$

由以上完全条件分布即可给出 MCMC 方法在 WinBUGS 中的抽样方案:

- $(n_1, \dots, n_k) \sim \frac{N!}{\prod_{i=1}^k n_i!} \prod_{i=1}^k p_i^{n_i}$
- $p_i \sim f(p_i | \theta), j = 1, 2, \dots, k$
- $\theta \sim f(\theta)$

其中 $f(\theta)$ 为适当选择的未知参数的先验分布, $f(p_i | \theta)$ 由假设的几类 ZID 和 IPD 的分布确定. 由于 $\log p_i < 0$, 根据 (8) 式引入负似然对数函数 $NLL = -\sum_{i=1}^k n_i \log p_i$ 作为拟合分布选择的依据.

3.2 包含分组因素的索赔次数数据

在保险实务中, 另一类较常见的索赔次数还包含有分组因素, 设因素 A 的水平为 $i = 1, 2, \dots, n$, 因素 B 的水平为 $j = 1, 2, \dots, m$, 而记号 $n_i^{(i,j)}$, $p_i^{(i,j)}$ 分别表示 A 在 i 水平, B 在 j 水平下的索赔次数和概率. 其频数分布可以表示为

A	A ₁			...	A _n		
B	B ₁	...	B _m	B ₁	...	B _m
0	n ₁ ^(1,1)	...	n ₁ ^(1,m)	n ₁ ^(n,1)	...	n ₁ ^(n,m)
1	n ₂ ^(1,1)	...	n ₂ ^(1,m)	n ₂ ^(n,1)	...	n ₂ ^(n,m)
2	n ₃ ^(1,1)	...	n ₃ ^(1,m)	n ₃ ^(n,1)	...	n ₃ ^(n,m)
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
≥ k - 1	n _k ^(1,1)	...	n _k ^(1,m)	n _k ^(n,1)	...	n _k ^(n,m)

与讨论过的情况相比，这些数据均对应着新的分组因素，在拟合时应考虑其影响。

记 $N = \sum_{i=1}^n \sum_{j=1}^m \sum_{l=1}^k n_l^{(i,j)}$ ，类似地有

$$\begin{aligned}
 f(\theta|n_l^{(i,j)}) &\propto f(n_l^{(i,j)}|\theta)f(\theta) \\
 &\propto \left(\frac{N!}{\prod \prod \prod n_l^{(i,j)}!} \prod_{i=1}^n \prod_{j=1}^m \prod_{l=1}^k [p_l^{(i,j)}]^{n_l^{(i,j)}} \right) f(p_l^{(i,j)}|\theta)f(\theta) \\
 &\propto (\prod \prod \prod [p_l^{(i,j)}]^{n_l^{(i,j)}}) (\prod \prod \prod f(p_l^{(i,j)}|\theta_{ij})) (\prod \prod f(\theta_{ij}|\alpha_i, \beta_j)) f(\alpha_i)f(\beta_j).
 \end{aligned}$$

上式最后的 $f(p_l^{(i,j)}|\theta_{ij})$ 由适当的 ZID 和 IPD 分布确定， $f(\alpha_i)$ ， $f(\beta_j)$ 分别为参数 α_i 和 β_j 的先验分布，其中参数 α_i 表示 A 因素 i 水平的效应， β_j 表示 B 因素 j 水平的效应，而 $f(\theta_{ij}|\alpha_i, \beta_j)$ 则通过如下的对数线性模型确定

$$\begin{cases} \log(\theta_{ij}) = \alpha_i + \beta_j + \mu_{ij}, \\ \mu_{ij} \sim N(0, \tau^2), \tau \sim \text{Beta}(0.001, 0.001). \end{cases}$$

相应的 MCMC 方法在 WinBUGS 中的抽样方案：

- $n_l^{(i,j)} \sim \frac{N!}{\prod \prod \prod n_l^{(i,j)}!} \prod_{i=1}^n \prod_{j=1}^m \prod_{l=1}^k [p_l^{(i,j)}]^{n_l^{(i,j)}}$
- $p_l^{(i,j)} \sim f(p_l^{(i,j)}|\theta_{ij})$
- $\theta_{ij} \sim \begin{cases} \log(\theta_{ij}) = \alpha_i + \beta_j + \mu_{ij} \\ \mu_{ij} \sim N(0, \tau^2), \tau \sim \text{Beta}(0.001, 0.001) \end{cases}$
- $\alpha_i \sim f(\alpha_i), \beta_j \sim f(\beta_j)$

§ 4. 数值例子

首先看一个一般索赔次数数据的例子。

数据来源于 Klugman, Panjer 和 Willmot 讨论的例子 ([9])，25000 份保单中索赔次数分别为 0, 1, 2, 3, 4, 5 及 6 次以上的保单数目依次为 23148, 1639, 173, 35, 3, 2 和 0。

经过简单的计算和统计检验可知，索赔次数的均值小于方差 (Mean = 0.08, Variance = 0.103) 且差异显著，故认为存在散度过大特征。按照前述方法，分别利用 ZID 和 IPD 类型分布进行拟合，结果如下 (其中 P[1], P[2], ..., P[7] 分别表示索赔次数为 0, 1, ..., ≥ 6 的概率)

《应用概率统计》版权所有

表 1 例 1 各分布拟合数据

类型	分布	NLL	P[1]	P[2]	P[3]	P[4]	P[5]	P[6]	P[7]
ZID	ZIMP	7387.0	0.9259	0.06538	0.007214	0.001214	2.114E-4	3.269E-5	5.142E-6
	ZIGP	7389.0	0.926	0.06489	0.007778	0.001083	1.651E-4	2.683E-5	5.565E-6
	ZINB	7389.0	0.926	0.06484	0.0079	0.001052	1.465E-4	2.104E-5	3.642E-6
IPD	IPNB	7390.0	0.926	0.06485	0.00798	0.001023	1.346E-4	1.81E-5	2.895E-6
	IPP	7391.0	0.926	0.06481	0.008084	9.852E-4	1.18E-4	1.396E-5	1.85E-6

根据表 1 中的结果, 当采用 ZIMP 分布时, NLL = 7387 为最小值, 说明其拟合效果最好. 进一步的, 确定上述分布的参数值由下表 2 给出

表 2 例 1 各分布参数值

类型	分布	ρ	μ	λ	λ_1	λ_2	r	β
ZID	ZIMP	0.9391	0.2974	-	0.09124	0.74	-	-
	ZIGP	0.08292	0.09479	0.1632	-	-	-	-
	ZINB	0.177	-	-	-	-	0.567	0.1833
IPD	IPNB	0.03957	0.1024	-	-	-	1.131	-
	IPP	0.0898	-	0.0769	-	-	-	-

在应用 WinBUGS 进行抽样时, 设定了四个初始值, 对应于 ZIMP 模型, 其 NLL 值的抽样过程, 统计量计算的采样部分和收敛过程可分别由下列图 1 与图 2 展示.

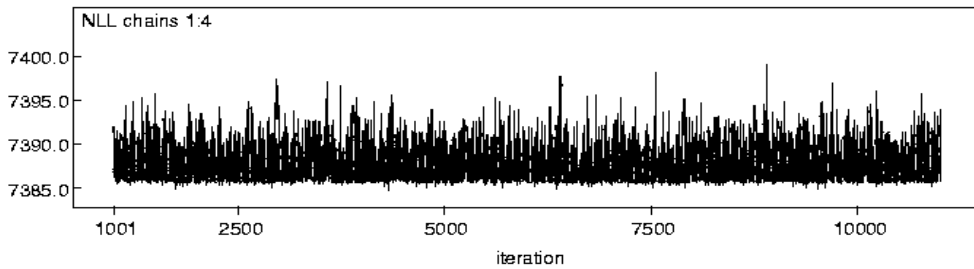


图 1 NLL 值抽样过程

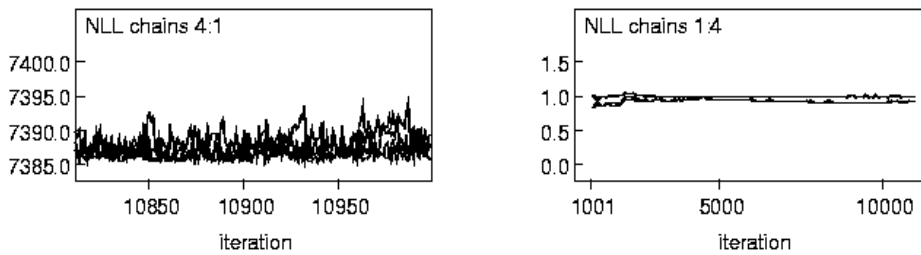


图 2 NLL 值采样部分与收敛过程

再来看一个包含分组因素的索赔次数数据.

数据来源于 ([11]), 某保险公司 12299 辆投保机动车辆车身险的保单按投保人的年龄 (因素 A) (是否大于 25 岁) 和车辆类型 (因素 B) (家用和高性能) 分成四组, 索赔次数分别为 0, 1, 2, 3 及 4 次以上的数量, 对于 25 岁以上家用组为 5019, 738, 65, 4 和 0; 对于 25 岁以上高性能组为 1068, 182, 27, 4 和 0; 对于 25 岁以下家用组为 2907, 592, 66, 5 和 0; 对于 25 岁以下高性能组为 1232, 334, 50, 6 和 0.

《应用概率统计》版权所有

为简化起见，利用前述模型及 ZIP 和 IPP 分布进行拟合，结果如下 (其中 $P[1], P[2], \dots, P[5]$ 分别表示索赔次数为 $0, 1, \dots, \geq 4$ 的概率)

表 3 例 2 各分布参数值

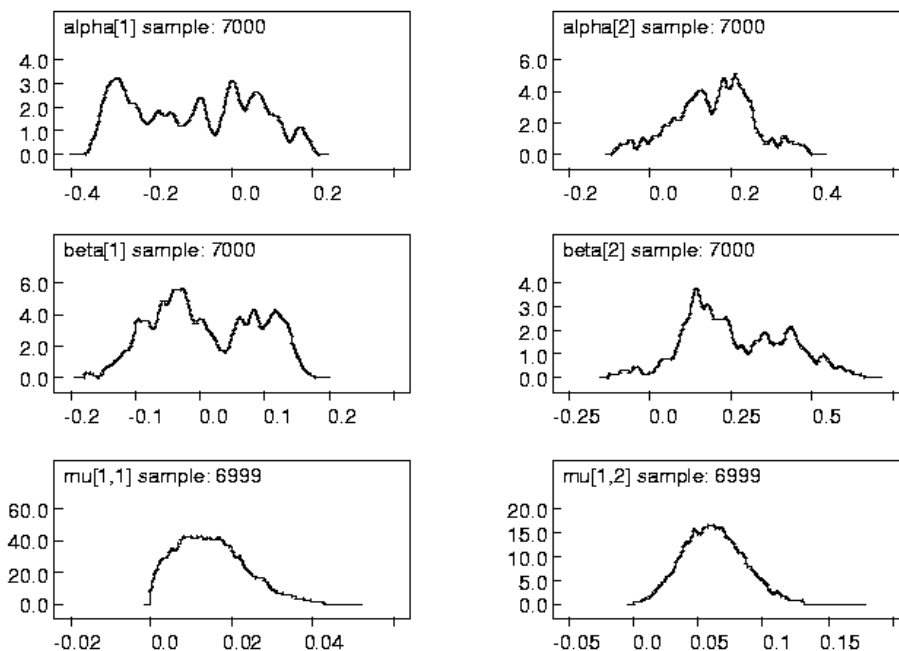
分布	NLL	因素 A	因素 B	P[1]	P[2]	P[3]	P[4]	P[5]
ZIP	6307.0	1	1	0.8616	0.1263	0.01141	6.939E-4	3.322E-5
			2	0.8333	0.1424	0.02187	0.002296	1.983E-4
		2	1	0.8151	0.1644	0.01894	0.001464	8.96E-5
			2	0.7606	0.2032	0.03244	0.003484	3.027E-4
IPP	6306.0	1	1	0.8615	0.1265	0.01118	7.689E-4	4.869E-5
			2	0.832	0.1435	0.02124	0.002845	4.193E-4
		2	1	0.8145	0.1649	0.01887	0.001606	1.224E-4
			2	0.7603	0.2039	0.03164	0.003707	4.061E-4

同样，IPP 的 NLL 优于 ZIP，其拟合效果更好。模型的参数列入下表 4

表 4 例 2 各分布参数值

分布	因素 A	因素 B	α	β	μ	λ
ZIP	1	1	0.1092	0.04853	-	0.1557
		2	0.1092	0.1734	-	0.3575
	2	1	0.1084	0.04853	-	0.09687
		2	0.1084	0.1734	-	0.1193
IPP	1	1	-0.09454	0.01068	0.01504	-
		2	-0.09454	0.2638	0.06213	-
	2	1	0.1583	0.01068	0.01321	-
		2	0.1583	0.2638	0.02096	-

对于 IPP 模型，其参数的抽样分布如图 3 所示



《应用概率统计》版权所有

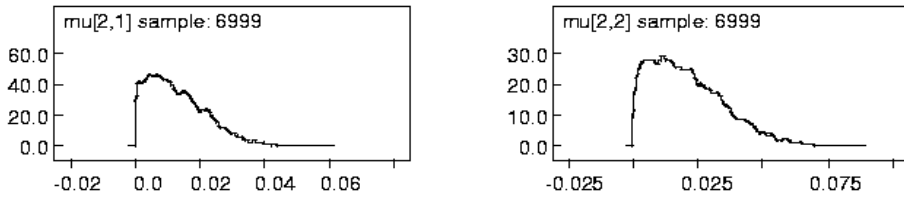


图3 模型参数的分布曲线

参 考 文 献

- [1] Lambert, D., Zero-inflated Poisson regression with an application to defects in manufacturing, *Technometrics*, **34**(1992), 1–14.
- [2] McCullagh, P. and Nelder, J.A., *Generalized Linear Models (Second edition)*, Chapman and Hall, London, 1989.
- [3] Consul, P.C. and Famoye, F., Generalized Poisson regression model, *Communications in Statistics-Theory and Methods*, **21**(1992), 89–109.
- [4] Daniel, B.H., Zero-inflated Poisson and binomial regression with random effects: a case study, *Biometrics*, **56**(2000), 1030–1039.
- [5] Stephane, R., A compound Poisson model for word occurrences in DNA sequences, *Journal of the Royal Statistical Society: Series C*, **51**(2002), 437–453.
- [6] Dagne, A., Hierarchical Bayesian analysis of correlated zero-inflated count data, *Biometrical Journal*, **46**(2004), 653–663.
- [7] Leda, D.M., A generalization of the classical discrete distributions, *Communications in Statistics: Theory and Methods*, **31**(2002), 871–888.
- [8] Ioannis, N. and Petros, D., Bayesian modelling of outstanding liabilities incorporating claim count uncertainty, *North American Actuarial Journal*, **6**(2002), 113–128.
- [9] Klugman, S.A., Panjer, H.H. and Willmot, G.E., *Loss Models: From Data to Decisions*, John Wiley and Sons, New York, 1998.
- [10] Scollnik, P.M., Actuarial modeling with MCMC and BUGS, *North American Actuarial Journal*, **5**(2002), 96–124.
- [11] 王静龙等, 非寿险精算学, 中国人民大学出版社, 2004.
- [12] Stephens, D.A., Crowder, M.J. and Dellaportas, P., Quantification of automobile insurance liability: a Bayesian failure time approach, *Insurance: Mathematics and Economics*, **34**(2004), 1–21.

Bayesian Estimates of Distribution for Count Data with Overdispersion

CHEN XUEDONG

(School of Mathematics and Statistics, Yunnan University, Kunming, 650091)

(Faculty of Science, Huzhou Teachers' College, Huzhou, 313000)

This paper deals with two classes distribution of count data with overdispersion: Zero-inflated Distribution and Inflated-parameter Distribution, which are accordance with data of claims. We consider several model formulations of those distributions by using Bayesian theory and MCMC methods in WinBUGS. By comparison, a approach of modelling data is obtained and two illustrations with real data are provided.