

第3章损失分布的统计推断方法

【内容要点】

1. 经验分布函数

(1) 经验分布函数的定义

以频率代替概率， $F_n(x)$ 实际上就是 n 个观测值落在区间 $(-\infty, x]$ 的频率，注意分布函数在 $x(k) \leq x < x(k+1)$ 这一左闭右开区间为一恒值。

变量观测值不多时，此时即为离散型变量的分布函数；变量观测值较多时，将其人为分组，此时即为连续型变量的分布函数，同时为了消除组距不同对频率的影响，可以计算频率密度：频率密度=频率组距。

(2) 频率直方图和曲线图

频率密度直方图以损失随机变量的观测值分组为横坐标，以各组的频率密度为纵坐标，将频率密度直方图的各个矩形上底中点用线段连接起来，就是频率密度折线图。把折线图修匀(光滑化)，就成为频率密度曲线图。

(3) 损失次数的经验分布

即以损失次数为离散随机变量的分布函数。

(4) 损失额的经验分布

即为连续型变量的经验分布，如果直接以损失额为因变量的修匀分布函数不够光滑，可试用损失额的其他函数作为因变量，如损失额对数形式等。

2. 损失分布参数的估计、假设检验和损失分布函数的拟合检验

完整的拟合过程为：

(1) 损失分布的选择

根据损失观测数据的频率密度的直方图或曲线图形态(如连续还是离散，单峰还是多峰，正态还是偏态，左偏还是右偏等)，从与其形似的理论分布族中，选择适当的分布律或概率密度函数曲线。

特别地，根据损失次数经验分布离散、右偏的特点，一般常选择的理论分布有：泊松(Poisson)分布、二项分布、负二项(Pascal)分布等，其中泊松分布是单参数的，二项分布和负二项分布都是双参数的分布。

根据损失额经验分布连续、非负、右偏和长尾的特点，一般常选择对数正态分布、帕雷托(pareto)分布、伽玛(Gamma)分布和韦伯

(Weibull)分布等作为理论分布，它们都是双参数的分布。有的还用三参数的广义帕雷托分布和布尔(Burr)分布作为损失额的理论分布。

(2) 利用观测数据对理论分布的参数进行估计(点估计含矩法估计、极大似然估计、分位数估计、最小二乘估计、区间估计)和假设检验。

① 点估计

a. 矩法估计

即以样本原点矩作为随机变量 x 的同类同阶原点矩，一般的，待估参数有几个就建立几个矩

等方程。样本原点矩是根据样本数据计算的，随机变量的原点矩是根据其分布函数计算的，是随机变量的数字特征，其中含有未知参数。

b. 极大似然估计

是使得联合概率函数（称为似然函数）的 $L(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$ 达到最大的参数 θ 。这一点很好理解，因为似然函数是概率的乘积，显然使得该乘积最大的 θ 是较好的估计。

至于其求解就是最大化 L 函数，一般求使其一阶倒数为 0 即可，有时为了简化计算，极大化 L 的对数 $\ln L$ 。

值得注意的是均匀分布的唯一参数 θ （均匀分布的区间长度）的极大似然估计为样本的最大值。

c. 分位数估计

就是应用样本或者总体的分布函数，找到使得 $P(X \leq x_i) =$ 分位百分比的 x_i ，一般的有几个未知参数就建立几个等式。

要切记的是，在样本中查找分位数对应的分位点注意是要使得 $P(X \leq x_i) =$ 分位百分比，而非 x_i 的下一个或上一个样本。另外，首先要将样本按从小到大排序。

d. 最小二乘估计

随机变量 X 的先验分布函数为 $F(x; \theta)$ ，根据观测值 x_1, x_2, \dots, x_n 建立的经验分布函数为 $F_n(x)$ ，所谓未知参数 θ 的最小二乘估计，就是使得先验分布函数与经验分布函数的距离 $D(\theta) = \sum_{i=1}^n [F(x_i; \theta) - F_n(x_i)]^2$ 达到最小的 θ 。

在实际应用时为了提高最小二乘估计的效果，常将 $D(\theta)$ 加权调整为

$D(\theta) = \sum_{i=1}^n W(x_i) [F(x_i; \theta) - F_n(x_i)]^2$ ，其中权 $W(x_i)$ 可以有不同的取法，一般可取

$$W(x_i) = F(x_i; \theta) [1 - F(x_i; \theta)]$$

根据概率统计理论，有三种性质常用来评价估计的好坏：无偏性、有效性和一致性。具体视实际问题的要求在这三个标准中权衡。

② 区间估计

是使得参数 θ 处于某一置信区间的概率等于置信度的估计，当然一般要求事先知道样本或者总体的分布。

③ 假设检验

先对随机变量的参数作出假设，然后选择检验统计量，确定检验规则和拒绝域，再根据观测到的数据计算检验统计量的值，从而作出拒绝还是接受原假设的判断。

(3) 对理论分布作拟合优度检验（注意：与上面的参数假设检验不同，参数假设检验是先已知样本或者总体分布，假设参数的某些数字特征，然后来检验假设的正确性与否；而这里的拟合优度检验是对上述前两项工作即损失分布的选择和参数估计或检验的评估）。

一般选用 χ^2 检验。 χ^2 拟合优度检验的原假设是 H_0 ：损失数据服从某个理论分布。为了检验 H_0 ，先把观测到的数据分成 n 组。选用检验统计量 $\chi^2 = \sum_{i=1}^n (Q_i - E_i)^2 / E_i$ ，式中的 Q_i 是第 i 组中观测到的实际频数， E_i 是根据假设的理论分布计算出来的理论频数。在假设 H_0 成立的前提下，统计量 χ^2 服从自由度为 $n-k-1$ 的 χ^2 分布，其中 k 为理论分布中用估计方法得到的参数的个数。

值得注意的是，在分组时，一定要包含样本的所有可能取值范围，切不可因为在某一区间没有样本，就丢掉这一区间，因为随机变量的分布在整个区间是完整的。比如，一索赔额分布如下，索赔额在 0~1000 有 200 次，1000~2000 有 300 次，2000 以上 0 次，一定要将索赔额分为三段，0~1000，1000~2000 和 2000 以上。

3. 部分保险数据损失分布的条件分布问题

主要研究免赔和再保险业务风险的有关问题（详见下面的重难点解析）。

4. 多风险损失的分布函数的拟合

已知两种风险（如火灾和盗窃）损失数据，如何据此求解总损失的分布函数。

第一种方法是把两类损失分布函数的加权平均作为总损失的分布函数。若 $F_1(x)$ 和 $F_2(x)$ 分别为两类损失的分布函数， p 是在已知发生损失的条件下发生的是盗窃损失的概率，是根据观测数据来估计的。那么总损失的分布可以确定为

$$F(x) = pF_1(x) + (1-p)F_2(x)$$

第二种方法是两类损失的观测值合并在一起，相当于一种风险进行拟合。

5. 贝叶斯估计

步骤：

(1) 选择先验分布，注意是关于待估参数 θ 的分布，而不是关于 x 的分布；

(2) 确定条件密度函数，是在假设 θ 已知的情况下，观测值 x 的密度，并求该密度的乘积即联合密度函数；

(3) 求出 θ 的后验分布，实际就是一个条件密度函数， $\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int f(x|\theta)\pi(\theta)d\theta}$ 。

由于其分母是关于 θ 的定积分，是关于 θ 的常数，与 θ 的分布无关，因此，可以不求解其结果。即 $\pi(\theta|x) \propto f(x|\theta)\pi(\theta)$ (\propto 表示成比例的意思)

常见的先验分布和后验分布

损失随机变量的分布观测值为 x_1, x_2, \dots, x_n

未知参数

参数的先验分布

参数的后验分布

泊松分布 $P(\lambda)$

伽玛分布 $\Gamma(\alpha, \lambda)$

$\lambda > 0$

$\Gamma(\beta, \gamma)$

$\Gamma(\sum x_i + \beta, n + \gamma)$

$\Gamma(n\alpha + \beta, \sum x_i + \gamma)$

泊松分布 $P(\lambda)$

伽玛分布 $\Gamma(\alpha, \lambda)$

$\lambda > 0$

$U(0, +\infty)$

$\Gamma \sum x_i + 1, n$

$\Gamma n + \alpha, \sum x_i$

二项分布 $B(m, p)$

负二项分布 $NB(k, p)$

$0 < p < 1$

$Beta(\alpha, \beta)$

$Beta \sum x_i + \alpha, nm - \sum x_i + \beta$

$Beta n + \alpha, \sum x_i + \beta$

正态分布 $N(\mu, \sigma^2)$

$-\infty < \mu < +\infty$

$N(\theta, \tau^2)$

$N \sum x_i \sigma^2 + \theta, \tau^2 n \sigma^2 + 1, 1 n \sigma^2 + 1 \tau^2$

正态分布 $N(\mu, \sigma^2)$

$-\infty < \mu < +\infty$

$U(-\infty, +\infty)$

$N \sum x_i, \sigma^2 n$

(4) 选择损失函数 $Loss(\hat{\theta}, \theta)$

所谓损失函数是指真实 θ 与估计值 $\hat{\theta}$ 的误差的评价函数，常用的有三种损失函数：平方损失函数 $(\hat{\theta} - \theta)^2$ 、绝对值损失函数 $|\hat{\theta} - \theta|$ 、0-1 损失函数 $L=1, |\hat{\theta} - \theta| > \varepsilon$

$0, |\hat{\theta} - \theta| \leq \varepsilon$ ($\varepsilon > 0, \varepsilon$ 极小)。

(5) 不同损失函数下的 θ 的估计

① 平方损失函数下的贝叶斯估计是后验分布的均值

$$E[Loss(\hat{\theta}, \theta)] = \int Loss(\hat{\theta}, \theta) \pi(\theta | \mathbf{x}) d\theta = \int (\hat{\theta} - \theta)^2 \pi(\theta | \mathbf{x}) d\theta$$

$$d\hat{\theta} \theta \wedge E[Loss(\hat{\theta}, \theta)] = -2 \int (\hat{\theta} - \theta) \pi(\theta | \mathbf{x}) d\theta = 0$$

所以 $\theta^{\wedge} \int \pi(\theta | x) d\theta = \int \theta \pi(\theta | x) d\theta$

而 $\int \pi(\theta | x) d\theta = 1$,

所以 $\theta^{\wedge} = E(\theta | x) = \int \theta \pi(\theta | x) d\theta$

② 绝对误差函数下的贝叶斯估计是后验分布的中位数

$$E[\text{Loss}(\theta^{\wedge}, \theta)] = \int \text{Loss}(\theta^{\wedge}, \theta) \pi(\theta | x) d\theta = \int \theta^{\wedge} - \theta \pi(\theta | x) d\theta$$

$$= \int_{-\infty}^{\theta^{\wedge}} (\theta^{\wedge} - \theta) \pi(\theta | x) d\theta + \int_{\theta^{\wedge}}^{+\infty} \theta^{\wedge} (\theta - \theta^{\wedge}) \pi(\theta | x) d\theta$$

$$d\theta^{\wedge} E(\text{Loss}(\theta^{\wedge}, \theta)) = \theta^{\wedge} \pi(\theta^{\wedge} | x) - \int_{-\infty}^{\theta^{\wedge}} \pi(\theta | x) d\theta - \theta^{\wedge} \pi(\theta^{\wedge} | x) - \theta^{\wedge} \pi(\theta^{\wedge} | x)$$

$$+ \int_{\theta^{\wedge}}^{+\infty} \theta^{\wedge} \pi(\theta | x) d\theta + \theta^{\wedge} \pi(\theta^{\wedge} | x)$$

$$= 0$$

$$\text{所以 } \int_{-\infty}^{\theta^{\wedge}} \pi(\theta | x) d\theta = \int_{\theta^{\wedge}}^{+\infty} \pi(\theta | x) d\theta$$

所以 θ^{\wedge} 为 $\pi(\theta | x)$ 的中位数。

③ 绝对误差函数下的贝叶斯估计是后验分布的众数

$$E[\text{Loss}(\theta^{\wedge}, \theta)] = \int \text{Loss}(\theta^{\wedge}, \theta) \pi(\theta | x) d\theta = 1 - \int_{\theta^{\wedge} - \varepsilon}^{\theta^{\wedge} + \varepsilon} \pi(\theta | x) d\theta$$

ε 极小, 故只有 $\pi(\theta | x)$ 最大时上式最小,

所以 θ^{\wedge} 为 $\pi(\theta | x)$ 的众数。

6. 信度理论

(1) 信度理论概述

非寿险保费的估算可以根据两类数据: 一类是通过观察得到的本险种一组保单的近期数据; 另一类是同险种保单早期损失数据或者类似险种保单的同期损失数据。根据前一类数据确定的保费称为经验保费

PM_o , 根据后一类数据得到的保费称为先验信息保费 PM_e 。所谓信度理论, 是如何合理将这两者加权 $((1-z) PM_o + z PM_e)$ 以作为后验保费的估计, 其中 z 称为信度或者信度因子。

(2) 有限扰动信度

① 原理和思想: 如果 C 的估计量 C^{\wedge} 是由先验信息数据 M 和近期观察值 T 加权获得, 即 $C^{\wedge} = (1-z) M + zT = (1-z) M + zET + z(T-ET)$, $z(T-ET)$ 就是 C^{\wedge} 和真值 C 距离的扰动, 有限扰动信度是要该扰动的相对值不超过一定限度 γ 的概率足够达到 $1-\alpha$, 有

$$P\{x - \mu \leq \gamma\} \geq 1 - \alpha$$

② 完全可信性

a. 完全可信性的条件

完全可信性即指完全信度经验， $z=1$ ，因此，

$$P\{x - \mu \leq \gamma\}$$

$\geq 1 - \alpha$ ，其逆定理同样成立，即完全可信性成立的条件是

$$P\{x - \mu \leq \gamma\}$$

$$\geq 1 - \alpha$$

b. 意义和应用：对样本量的要求

运用中心极限定理，假设 X 的方差为 $\text{Var}(X) = \sigma^2$ ，则 $(x - \mu) / \sigma$ 的渐近分布为标准正态分布 $N(0, 1)$ 。根据上式，我们有

$$P\left\{\frac{x - \mu}{\sigma} \leq \frac{\gamma}{\sigma}\right\}$$

$$\geq 1 - \alpha$$

因而

$$n \geq \frac{U_{1-\alpha/2}^2 \sigma^2}{\gamma^2}$$

式中的 $U_{1-\alpha/2}$ 是标准正态分布 $N(0, 1)$ 的 $1-\alpha/2$ 分位点。于是，完全可信性条件可以简化为

$$n \geq \frac{U_{1-\alpha/2}^2 \sigma^2}{\gamma^2}$$

$$\sigma^2 = U_{1-\alpha/2}^2 \gamma^2 \text{Var}(X) / [E(X)]^2$$

$$\Delta n \geq 0$$

③ 平方根法则

在上面所说的同样假设（运用中心极限定理，假设 X 的方差为 Var

$(X) = \sigma^2$ ，则 $(x - \mu) / \sigma$ 的渐近分布为标准正态分布 $N(0, 1)$ ）下

$$z = n\gamma\mu\sigma U_{1-\alpha/2} = n\lambda_0\sigma\mu$$

所以 $z = nn_0$ ，这就是所谓的平方根法则。

④ 值得注意的是， z 的最大值， $n \geq n_0$ 时，不可根据平方根法则 $z = nn_0$

$$z = \min nn_0, \quad 1 = \min \mu \sigma n \lambda_0, \quad 1 = \min E(X) \text{Var}(X) n \lambda_0, \quad 1$$

(3) Buhlmann 信度

① 背景知识

由于贝叶斯估计方法需要事先知道参数 θ 的先验分布 $\pi(\theta)$ 条件密度 $f(x|\theta)$ ，而且计算也较困难。因此，在非寿险精算中，贝叶斯信度并不实用。于是，瑞士精算学家 Buhlmann 提出了用观测值的线性函数作为损失随机变量 X 或者其参数的可信性估计，并由此来推算保费的信度模型。

② 模型结构参数及其无偏估计

该模型同样假设损失随机变量 X 的分布参数 θ 也是随机变量，非寿险精算中常把 θ 的不同取值理解为异质风险的不同风险水平，因而也叫风险参数。记 $\mu(\theta) = E(X|\theta)$ ， $v(\theta) = \text{Var}(X|\theta)$ ，并称 $\mu(\theta)$ 为假设均值，称 $v(\theta)$ 为过程方差 (process variance)，它们都是随机变量 θ 的函数。记

$$\mu = E(X) = E[E(X|\theta)] = E[\mu(\theta)]$$

(先对 θ_i 下的 x 求得期望值 $\mu(\theta_i)$ ，再对不同 θ_i 下的 $\mu(\theta_i)$ 求得均值，所以 μ 称为 X 的总均值。因此，可以引入估计量

$$\hat{\mu}(\theta_i) = \bar{x}_i = \sum_{j=1}^n x_{ij} / n_i$$

$$\hat{\mu} = \sum_{i=1}^r \bar{x}_i = \sum_{i=1}^r \sum_{j=1}^n x_{ij} / n_i$$

$$\hat{\mu} = \sum_{i=1}^r \sum_{j=1}^n x_{ij} / n_i$$

容易证得这些估计都是无偏估计。

$$\hat{v} = E[\text{Var}(X|\theta)] = E[v(\theta)]$$

(先对 θ_i 下的 x 求得方差 $v^*(\theta_i)$ ，是在同一风险水平下的内在差异，故称为同质方差，再对不同 θ_i 下的 $v^*(\theta_i)$ 求得均值，因此，可以引入估计量

$$v^*(\theta_i) = \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2 / (n-1), \quad v^* = \sum_{i=1}^r v^*(\theta_i) / r$$

$$v^*(\theta_i) = \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2 / (n-1)$$

容易证得这些估计都是无偏估计。

$$a = \text{Var} [E(X|\theta)] = \text{Var} [v^*(\theta)]$$

(先对 θ_i 下的 x 求得期望值 $\mu(\theta_i)$ ，再对这些期望值求得方差，描述不同风险水平下的 x 的均值之间的差异，故称为异质方差。据此考虑引入估计量 $\sum_{i=1}^r (x_i - \bar{x})^2 / (r-1)$ ，由于 $E \sum_{i=1}^r (x_i - \bar{x})^2 / (r-1)$

$= a + v^*$ ，即 $\sum_{i=1}^r (x_i - \bar{x})^2 / (r-1)$ 不是 a 的无偏估计，所以修改后的 $a^* = \sum_{i=1}^r (x_i - \bar{x})^2 / (r-1) - v^*$ 为 a 的无偏估计。)

③ 两个定理

$\text{Var}(X) = E[\text{Var}(X|\theta)] + \text{Var}[E(X|\theta)] = v^* + a$ (可以用来简便计算)

若两个损失随机变量 X_1 和 X_2 的风险参数相等，则 $\text{Cov}(X_1, X_2) = a$ 。

④ Buhlmann 方法

基本思想是：把损失随机变量 X 的观测值 x_1, x_2, \dots, x_n 的线性函数 $\beta + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$ 作为将来损失的可信性估计，并使均方损失达到最小。正由于这个原因，Buhlmann 信度也叫最小平方信度。

这里的 β_i 是观测值 x_i 在线性函数中的权重， $i=1, 2, \dots, n$ 。由于观测值 x_1, x_2, \dots, x_n 的风险特征相同，所以可令 $\beta_1 = \beta_2 = \dots = \beta_n$ ，这时 $\beta + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n = \beta + n \beta_1 \bar{x}$ ，若记 $z = n \beta_1$ ，则将来损失的可信性估计为 $\beta + z \bar{x}$ 。

式中的 β 和 z 可以通过求均方损失函数

$$L(\beta, z) = E\{[\beta + z\bar{x} - \mu(\theta)]^2\}$$

的极小值点得到。经过简单的计算可知，当

$$\beta = (1-z) \mu$$

$z = n v^* / (n v^* + a)$ 时， $L(\beta, z)$ 达到极小。

用这样计算得到 $\beta + z\bar{x} = (1-z) \mu + z\bar{x}$ 作为将来损失的估计而厘定的保费，称为 Buhlmann 信

度保费 (buhlmann credibility premium), 简称信度保费。

⑤ Buhlmann 信度保费与平方损失条件下的贝叶斯估计得到的保费是一致的。这一点非常有用, 可以通过平方损失条件下贝叶斯估计求得 Buhlmann 信度保费的信度, 同时可以通过 Buhlmann 信度保费的信度求得平方损失条件下贝叶斯估计。

(4) Buhlmann Straub 信度

实际就是加权的 Buhlmann 信度, 具体来说是对同质方差的加权, 可以参考韩天雄主编的教材, 此处不再赘述。

【重难点解析】

【例 3.1】某公司在 2002 年中共发生了 100 个理赔个案, 并按损失大小整理出下列频率密度直方图。但在整理过程中, 数字 x 和 y 遗失了, 仅知道有 60 个理赔在 x 和 y 之间发生。试根据直方图信息求 x 的值。

解频率密度直方图面积为 1, 即 $\int 3000100f(x) dx=1$

$$0.0006 \times (y-x) = 0.6$$

$$0.0009 \times (x-100) + 0.0001 \times (3000-y) = 0.4$$

得 $x=362.5$

【例 3.2】已知索赔记录如下表所示。

| 索 赔 次 数 | 平均索赔金额 |
|---------|--------|
| 2001 年 | 100 |
| 1000 | 2002 年 |
| 200 | 1250 |

每年的通货膨胀率为 10%, 现用 Pareto (3, λ) 分布作为 2003 年内平均索赔金额的模型。求 λ 的矩估计值。

解

$$E(X) = \lambda \alpha^{-1} = \lambda 3^{-1} =$$

$$x = 100300 \times 1000 \times 1.12 + 200300 \times 1250 \times 1.1 = 1320$$

$$\lambda = 2640$$

部分保险数据损失分布的条件分布问题，是容易让读者疑惑的问题，在这里给予重点说明如下。

存在免赔额时保险人承担的风险 Y 和保险标的风险 X 的关系为

$$Y = X - d, X > d$$

Y 未定义，其他

此时

$$F_Y(y) = P(Y \leq y) = 0, y < d$$

$$P(Y \leq y | X > d) = F_X(y) - F_X(d) / 1 - F_X(d), y \geq d$$

$$f_Y(y) = 0, y < d$$

$$f_X(y) / 1 - F_X(d), y \geq d$$

$$E(Y) = \int_0^{\infty} dx f_X(x) dx$$

超额分保再保险（自留额为 d ）的再保险人承担的保险责任与保险标的风险的关系为

$$Y = X - d, X > d$$

Y 未定义，其他

此时

$$F_Y(y) = P(Y \leq y) = 0, y < 0$$

$$P(Y \leq y | X > d) = F_X(y + d) - F_X(d) / 1 - F_X(d), y \geq 0$$

$$f_Y(y) = 0, y < 0$$

$$f_X(y+d) - F_X(d), y \geq 0$$

$$E(Y) = \int_0^{\infty} d(x-d) f_X(x) dx$$

【例 3.3】设原保险人与再保险人签订超额分保合同，自留额为 d ，无最高限额，设保险标的风险 x 服从参数为 (u, σ^2) 的对数正态分布，求此时再保险人承担的保险责任 y 的期望值。

解

$$\ln x \sim N(u, \sigma^2),$$

$$y = x - d, x > d$$

y 未定义，其他

$$E(Y) = \int_0^{\infty} d(x-d) f_X(x) dx = \int_0^{\infty} d(x-d) \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(\ln x - u)^2} dx$$

$$= e^{u - \frac{\sigma^2}{2}} \int_{\ln d - u}^{\infty} (x-d) \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(\ln x - u)^2} dx$$

$$= e^{u - \frac{\sigma^2}{2}} \int_{\ln d - u}^{\infty} (x-d) \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(\ln x - u)^2} dx$$

$$= e^{u - \frac{\sigma^2}{2}} \int_{\ln d - u}^{\infty} (x-d) \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(\ln x - u)^2} dx$$

(此积分较为复杂，本例宜记住)

$$\text{同时, } E(Y) = E(X) - \int_0^{\infty} d [1 - F_X(x)] dx$$

$$= E(X) - \int_0^{\infty} d [1 - F_X(x)] dx$$

$$E(Y) = E(X) - \int_0^{\infty} d [1 - F_X(x)] dx$$

【例 3.4】

$F(x) = 1 - qe^{-px} (x \geq 0)$ ，求其矩母函数。

解这是一个典型的混合分布函数，

由 $F(x) = 1 - qe^{-px} = p \times 1 + q \times (1 - e^{-px})$ 知其是

$p(x=0)=1$ 和参数为 p 的指数分布的混合分布，前者的概率是 p ，后者的概率是 $q(=1-p)$ ，

$F(x) = pqe^{-px} (x \geq 0)$ ，且 $F(x=0) = p(x=0) = p$ ，故矩母函数

$$MX(t) = E(etx) = p \times e^0 + \int_0^{\infty} 0etxpe^{-px} dx = p + qp e^{-t}$$

实际上就是 $p(x=0) = 1$ 和参数为 p 的指数分布的混合分布的矩母函数的线性组合。

贝叶斯估计是重点内容之一，该估计可应用于多个领域，并且程式化较强。

贝叶斯估计主观性的来源：贝叶斯方法主观性主要来源于先验分布的选择和对损失函数的选择。

【例 3.5】 被保险人每年的索赔次数服从泊松分布，50%的被保险人每年平均索赔次数为 2，另外 50%的被保险人每年平均索赔次数为 4，随机选择一位被保险人，发现其在开始两年的索赔次数都是 4 次，求该被保险人在第三年内索赔次数的贝叶斯估计值。

解

$$\text{先验分布 } p(\lambda = 2) = p(\lambda = 4) = 0.5$$

$$p(x_1 = x_2 = 4 | \lambda = 2) = 24e^{-24}!2 = 0.00814$$

$$\text{同理, } p(x_1 = x_2 = 4 | \lambda = 4) = 0.03817$$

$$p(\lambda = 2 | x_1 = x_2 = 4)$$

$$= p(x_1 = x_2 = 4 | \lambda = 2) p(\lambda = 2)$$

$$p(x_1 = x_2 = 4 | \lambda = 2) p(\lambda = 2) + p(x_1 = x_2 = 4 | \lambda = 4) p(\lambda = 4)$$

$$= 0.18213$$

$$\text{同理, } p(\lambda = 4 | x_1 = x_2 = 4) = 0.81787$$

$$E(\lambda | x_1 = x_2 = 4) = 2 \times 0.18213 + 4 \times 0.81787 = 3.6357$$

【例 3.6】 设

x_1, x_2, \dots, x_n 是来自总体分布的伯努利分布的样本。

伯努利分布的参数为 p ，其先验分布为均匀分布 $U(0, 1)$ 。设后验分布的均值为 15。试求 n ($n \leq 5$) 和 $\sum x_i$ 。

解伯努利分布函数为

$$f(x|p) = p^x q^{1-x}, L(x|p) = p^{x_i} q$$

$$(1-x_i) = p^{x_i} q^{n-x_i}, \text{ 其先验分布为均匀分布即 } \pi(p) = 1, \text{ 则后验分布 } \pi(p|x) = f(x|p) \cdot \pi(p)$$

$$\int 10f(x|p) \pi(p) dp$$

为 Beta 分布 $\text{Beta}(\sum x_i+1, n-\sum x_i+1)$, 均值为 $\frac{\sum x_i+1}{n+2} = 15$, 所以

$$n=3, \sum x_i=0$$

【例 3.7】 已知索赔额变量的先验分布是 Pareto ($\alpha, 10$), 参数 α 等可能地取 1, 2, 3。现观察到某风险的索赔额为 20。计算该风险下次索赔额大于 30 的概率。

Pareto (α, λ) 的函数密度为 $f(x) = \alpha \lambda^\alpha (\lambda+x)^{-\alpha-1}, x>0$

$$\text{解 } \lambda=10, f(x=20|\alpha=1)$$

$$= \alpha \lambda^\alpha (\lambda+x)^{-\alpha-1}$$

$$= 1 \times 10^\alpha (10+20)^{-\alpha-1} =$$

$$0.011111$$

同理可得

$$f(x=20|\alpha=2) = 0.007407, f(x=20|\alpha=3) = 0.003704$$

先验分布

$$P(\alpha=1) = P(\alpha=2) = P(\alpha=3) = \frac{1}{3}$$

后验分布

$$\pi(\theta|x) =$$

$$f(x|\theta) \pi(\theta) / \int f(x|\theta) \pi(\theta) d\theta$$

$$P(\alpha=1|x=20) = \frac{0.011111 \times \frac{1}{3}}{0.011111 \times \frac{1}{3} + 0.007407 \times \frac{1}{3} + 0.003704 \times \frac{1}{3}} = 0.5$$

$$\text{同理 } P(\alpha=2|x=20) = 0.33, P(\alpha=3|x=20) = 0.17$$

索赔额大于 30 的概率应该是在 α 取各可能值下索赔额大于 30 的概率乘以 α 取各可能值的后验概率, 其中 $P(x>30|\alpha=n) = \int_{30}^{\infty} \alpha \lambda^\alpha (\lambda+x)^{-\alpha-1} dx = \lambda^{-\alpha} (\lambda+\alpha)^{-\alpha} (\lambda=10)$, 具体的将 $\alpha=1, 2, 3$ 代入略, 因此, $p(x>30) = \sum_{n=1}^3 P(x>30|\alpha=n) P(\alpha=n|x=20) = 0.15$ 。

贝叶斯估计和极大似然估计的区别与联系是两者都需要求得联合密度函数。区别在于贝叶斯估计需要已知 θ 的先验分布和求解 θ 的后验分布, 极大似然估计不需要; 贝叶斯估计是以极小化估计值和真值的误差为目标, 极大似然估计以极大化联合密度函数为目标。

信度理论是一个难点, 其涉及的统计知识已经超出正常的数理统计范围, 我们具体例说如下。

(1) 有限扰动信度

【例 3.8】

保险公司有 8000 份保单, 当年总赔付成本是 450000 元, 而总风险保费是 520000 元。已知满足完全可信性条件是保单达到 14833 份。试问保险公司在下年应收取多少风险保费。

解

$$z = \frac{\sum_{i=1}^n x_i}{n} = \frac{450000}{8000} = 56.25, z_0 = 14833/n = 1.854125$$

$$p^* = 0.735 \times 450000 + (1-0.735) \times 520000 = 468550$$

因此，下一年每份保单的风险保费为 $4685508000=58.6$ 元

(2) Buhlmann 信度

【例 3.9】现有 4 类风险共 7 年损失记录。 X_{ij} 表示第 i 类风险在第 j 年内的损失。已知：

$$\sum_{i=1}^4 \sum_{j=1}^7 x_{ij} = 33.6$$

$$\sum_{i=1}^4 (x_i - \bar{x})^2 = 3.3$$

按最小平方信度方法，求信度因子。

解 $v = \frac{r}{r+1}$

$$v = \frac{\sum_{i=1}^4 r_i}{\sum_{i=1}^4 r_i + \sum_{j=1}^7 (x_{ij} - \bar{x}_i)^2} = \frac{2r}{2r + (n-1)}$$

$$a = \frac{\sum_{i=1}^4 r_i (x_i - \bar{x})^2}{2r - v} = 1.4$$

$$\therefore v = \frac{33.64 \times (7-1)}{77+1} = 0.82$$

$$a = 3.34 - 1 - 1.47 = 0.9$$

$$\therefore z = \frac{nn + v a}{v} = 77 + 1.40.9 = 0.82$$

(3) Buhlmann 信度与平方损失条件下的贝叶斯估计

【例 3.10】某类保险产品索赔次数 x 服从参数为 θ 的泊松分布，而参数 θ 服从参数为 α, β 的 Γ 分布，它的密度函数为

$$f(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} e^{-\beta\theta} \theta^{\alpha-1} \quad (\theta > 0)$$

先观察 n 份保单，在责任期内索赔次数分别为 x_i ，试求 θ 的贝叶斯估计。

解联合密度

$$f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!}$$

$$= e^{-n\theta} \theta^{\sum x_i} \prod_{i=1}^n \frac{1}{x_i!}$$

$$\propto e^{-(\beta+n)\theta} \theta^{\sum x_i + \alpha - 1}$$

$$\propto e^{-(\beta+n)\theta} \theta^{\sum x_i + \alpha - 1}$$

这是以 $\sum x_i + \alpha, \beta + n$ 为参数的 Γ 分布，在平方损失条件下 θ 的贝叶斯估计为

$$\frac{\sum x_i + \alpha}{\beta + n} = \frac{\beta + n}{\beta + n} \frac{\sum x_i + \alpha}{\beta + n}$$

【例 3.11】对上例采用 Buhlmann 信度求其信度及其信度估计。

解由于索赔次数 x 服从参数为 θ 的泊松分布，因此

$E(x | \theta) = \text{Var}(x | \theta) = \theta$ ，又由于参数 θ 服从参数为 α, β 的 Γ 分布所以

$$E(\text{Var}(x|\theta)) = E(\theta) = \alpha/\beta,$$

$$\text{Var}[(E(x|\theta))] = \text{Var}(\theta) = \alpha/\beta^2,$$

$$v_a = E(\text{Var}(x)) + \text{Var}(E(x)) = \beta$$

$$z = n\alpha + v_a = n\alpha + \beta$$

$$\mu = E(\theta) = \alpha/\beta$$

所以 $(1-z)\mu + zx = \beta/\beta + n\alpha/\beta + n\beta/\beta + n\sum x/n$, 和例 3.7 结果一致, 其中可以挖掘的其他知识不再赘述。

由于贝叶斯估计的重要性, 我们再给出几个例子来加深读者的认识。

【例 3.12】 设某运输车队每年的事故发生数服从泊松分布, 参数 λ 可取 1.0 或 1.5, 又设 λ 的先验分布为

$P(\lambda=1.0)=0.4$, $P(\lambda=1.5)=0.6$, 假如某一年该车队发生了三次事故, 求 λ 的后验分布, 并在平方损失函数下求参数的贝叶斯估计。

解 $P(x=3|\lambda) = \lambda^3/3!e^{-\lambda}$

$$P(x=3|\lambda=1) = 1^3/3! e^{-1}$$

$$P(x=3|\lambda=1.5) = 1.5^3/3! e^{-1.5}$$

$$P(\lambda=1|x=3) = e^{-1} \times 0.4/3!e^{-1} + 1.5^3/3!e^{-1.5} \times 0.6/3!e^{-1} \times 0.4 = 0.2457$$

$$P(\lambda=1.5|x=3) = e^{-1.5} \times 0.6/3!e^{-1.5} + 1.5^3/3!e^{-1.5} \times 0.6/3!e^{-1} \times 0.4 = 0.7543$$

$$\text{所以 } \lambda^{\wedge} = E(\lambda) = P(\lambda=1|x=3) \times 1 + P(\lambda=1.5|x=3) \times 1.5$$

$$= 0.2457 + 0.7543 \times 1.5 = 1.3772$$

【例 3.13】 设 θ 是某险种的索赔频率, 抽取 8 份保单, 发现在有效期内有多次索赔, 假如先验分布为:

- (1) $\theta \sim U(0, 1)$;
- (2) $\theta \sim f(\theta) = 2(1-\theta), 0 < \theta < 1$

0, 其他

分别求 θ 的后验分布, 并在平方损失函数下求参数的贝叶斯估计。

解 (1) $P(x=3|\theta) = C_3^5 \theta^3 (1-\theta)^2$
 $P(\theta|x=3) \propto C_3^5 \theta^3 (1-\theta)^2 \sim \text{Beta}(4, 6)$
 所以参数的贝叶斯估计为

$$\hat{\theta} = \frac{4}{4+6} = \frac{2}{5}$$

$$(2) \quad P(\theta|x=3) \propto C_3^5 \theta^3 (1-\theta)^2$$

$$\propto C_3^5 \theta^3 (1-\theta)^2$$

$$\sim \text{Beta}(4, 7)$$

所以参数的贝叶斯估计为

$$\hat{\theta} = \frac{4}{4+7} = \frac{4}{11}$$

【例 3.14】从一组有效保单中抽取 100 份, 发现有三个索赔, 假如该险种的索赔频率 θ 的先验分布为 $\text{Beta}(2, 200)$, 求 θ 的后验分布. 并在平方损失函数下求参数的贝叶斯估计。

解 $P(\theta|x=3) \propto C_3^{100} \theta^3 (1-\theta)^{97} C_2^{200} \theta^{2-1} (1-\theta)^{200-1}$
 $\propto C \times \theta^{5-1} (1-\theta)^{297-1} \sim \text{Beta}(5, 297)$

所以参数的贝叶斯估计为

$$\hat{\theta} = \frac{5}{5+297} = \frac{5}{302}$$

注: 上两题中, 在 $0 \leq \theta \leq 1$ 误差函数下的贝叶斯估计, 该是多少呢? 请读者思考。

【例 3.15】关于参数 θ 的贝叶斯估计, 下列选项哪一项是正确的?

- ① 在二次损失函数下, θ 的估计是后验分布的中位数。
 - ② 在二次损失函数下, θ 的估计是后验分布的众数。
 - ③ 在 $0 \leq \theta \leq 1$ 误差函数下, θ 的估计是后验分布的均值。
 - ④ 在 $0 \leq \theta \leq 1$ 误差函数下, θ 的估计是后验分布的众数。
- A. 仅①正确 B. 仅②正确 C. 仅③正确

D. 仅④正确 E. 全都不正确

解选 D。贝叶斯估计有三种损失函数，它们对应的参数的后验分布的估计要牢记。

【例 3.16】设 p 的先验分布为 $(0, 1)$ 上的均匀分布，已知 x_1, x_2, \dots, x_n 是来自总体分布为二点分布的样本，二点分布的参数为 p ，并且已知后验分布的均值为 14，问以下结论哪一个是正确的？

A. $\sum x_i=1, n=2$ B. $\sum x_i=1, n=4$ C. $\sum x_i=0, n=2$

D. $\sum x_i=0, n=4$ E. $\sum x_i=0, n=6$

解根据贝叶斯公式， p 的后验分布

$$f(p|x_1, x_2, \dots, x_n) \propto f(p) f(x_1, x_2, \dots, x_n)$$

$$\propto p^{x_1} (1-p)^{1-x_1} p^{x_2} (1-p)^{1-x_2} \cdots p^{x_n} (1-p)^{1-x_n}$$

$$\propto p^{n_1} (1-p)^{n-n_1}$$

$$n_1 = \sum_{i=1}^n x_i$$

$$\sim \text{Beta}(\sum x_i + 1, n - \sum x_i + 1)$$

所以其均值为

$$\frac{\sum x_i + 1}{\sum x_i + 1 + n - \sum x_i + 1} = \frac{\sum x_i + 1}{n + 2} = 14$$

选 C， $\sum x_i=0, n=2$ 可以满足要求。

【例 3.17】设定某种疾病发病次数服从泊松分布，大约一半的人每年的发病次数为 1 次，另一半的人每年发病次数为 2 次，随机选取一人，发现其在前两年的发病次数均为 1 次，求该人在第三年内的索赔 1 次的贝叶斯估计值。

A. $2.2+e$ B. $4.4+e$ C. $4+e$ D. $2.8+e$

E. $2+e$

解涉及贝叶斯估计，必须明确以下几个问题：

- (1) 总体分布是什么；
- (2) 参数是哪一个；
- (3) 参数的先验分布；
- (4) 样本值；
- (5) 损失函数是哪一个。

对本题而言，设 X 表示发病次数，则 X 为总体， $X \sim \text{Poisson}(\lambda)$ 。 λ 是参数，由“大约一半的人每年的发病次数为 1 次，另一半的人每年发病次数为 2 次”，则 λ 的先验分布为 $P(\lambda=1) = P(\lambda=2) = 0.5$ ；那样本值是多少？“随机选取一人，发现其在前两年的发病次数均为 1 次”，说明样本容量为 2，且 $x_1=x_2=1$ 。这里的损失函数没作说明，但这是离散型随机变量，就取平方损失函数即可。

λ 的后验分布

$$P(\lambda=1|x_1=1, x_2=1)$$

$$=P(\lambda=1)P(x_1=x_2=1|\lambda=1) + P(\lambda=2)P(x_1=x_2=1|\lambda=2)$$

$$=11+4e^{-2}$$

同理,

$$P(\lambda=2|x_1=1, x_2=1)$$

$$=P(\lambda=2)P(x_1=x_2=1|\lambda=2) + P(\lambda=1)P(x_1=x_2=1|\lambda=1)$$

$$=4e^{-2}+11$$

所以第三年内索赔次数估计的贝叶斯估计值为 (也就是后验分布的期望)

$$\hat{\lambda} = E(\lambda|x_1=1, x_2=1) = 1 \times 11 + 4e^{-2} + 2 \times 4e^{-2} + 11 = 8 + e^{-2} + 11$$

选 D。

【例 3.18】有关贝叶斯方法的陈述, 下列选项中正确的是哪一项?

- ① 在 $0 \sim 1$ 损失函数下, 贝叶斯方法得到的信度因子的估计与最小平方信度是一致的;
- ② 在估计非线性问题时, 贝叶斯方法比最小平方信度更有优越性;
- ③ 贝叶斯方法含有主观成分, 此主观成分主要表现在对先验分布及损失函数的选取上。

- A. 仅①正确 B. 仅②正确 C. 仅③正确
- D. ②③正确 E. 全部正确

解②③正确, 选 D。

在平方损失函数下, 贝叶斯估计值才与最小平方信度一致。

最小平方信度方法实际上更倾向于一个线性模型, 而贝叶斯方法没有此限制。

【例 3.19】若 X 服从参数为 p 的几何分布, p 为随机变量且 p 服从参数为 (α, β) 的贝塔分布, 那么 p 的后验分布是什么?

- A. 贝塔分布 B. 几何分布 C. 均匀分布
- D. 参数为 $(\alpha+1, x+\beta)$ 的贝塔分布
- E. 参数为 $(x+\beta, \alpha+1)$ 的贝塔分布

解 $f(p|x) \propto f(p)f(x|p)$

$$\propto \Gamma(\alpha+\beta) \Gamma(\alpha) \Gamma(\beta) p^{\alpha-1} (1-p)^{\beta-1} (1-p)^{\beta-1} (1-p)^x p$$

$$\propto p^{\alpha+1-1} (1-p)^{\beta+x-1} \sim \text{Beta}(\alpha+1, \beta+x)$$

选 AD。

【例 3.20】设 40 张同类保单，用 x_i 表示第 i 张保单的索赔次数，并设 $X_i \sim P(\lambda)$, $i=1, 2, \dots, 40$ ，又设参数 λ 为随机变量，且服从均值为 0.6，方差为 0.02 的 $\Gamma(\alpha, \beta)$ 分布，分布密度为 $f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$, $x > 0$ ，已知观察到 40 张保单共有 18 次索赔，试计算在平方损失函数下 λ 的贝叶斯估计。

解 λ 的后验分布

$$f(\lambda | x_1, x_2, \dots, x_n) \propto f(\lambda) \cdot f(x_1, x_2, \dots, x_n | \lambda)$$

$$\propto \lambda^{\alpha-1} e^{-\beta \lambda} \cdot e^{-\lambda} \lambda^{x_1} x_1! \cdot e^{-\lambda} \lambda^{x_2} x_2! \cdot \dots \cdot e^{-\lambda} \lambda^{x_{40}} x_{40}!$$

$$\propto \lambda^{\alpha-1} e^{-\beta \lambda} e^{-40 \lambda} \lambda^{\sum_{i=1}^{40} x_i}$$

$$\propto \lambda^{\alpha+\sum_{i=1}^{40} x_i - 1} e^{-(\beta+40)\lambda}$$

$$\sim \Gamma(\alpha + \sum_{i=1}^{40} x_i, \beta + 40)$$

由条件有 $\alpha/\beta = 0.6$, $\alpha/\beta^2 = 0.02$ ，可解得

$$\alpha = 18, \beta = 30$$

在平方损失函数下， λ 的贝叶斯估计为后验分布的均值。
所以 λ 的贝叶斯估计为

$$\hat{\lambda} = \frac{\alpha + \sum_{i=1}^{40} x_i}{\beta + 40} = \frac{18 + 1830 + 40}{30 + 40} = 0.51428$$

【例 3.21】评估损失分布的贝叶斯方法是一种主观方法，其主观性的主要表现有几个？

- A. 1 个 B. 2 个 C. 3 个
D. 4 个 E. 以上都不对

解选 B。

【例 3.22】(多选) 在用贝叶斯方法估计损失分布中，其主观性表现在何处？

- A. 选择先验分布 B. 确定似然函数 C. 确定参数 θ 的后验分布
D. 选择损失函数 E. 估计参数

解选 AD。本题是上一题的详细考察。

【例 3.23】(多选) 以下对未知参数的贝叶斯估计中，分别选择二次损失函数、绝对误差函数、0-1 误差函数，得到结果相同的有哪几项？

- A. 总体服从 $N(\mu, \sigma^2)$, 未知参数 μ 的先验分布为 $N(\mu', \sigma'^2)$, 求未知参数 μ 的贝叶斯估计
- B. 总体服从对数正态 $\log N(\mu, \sigma^2)$, 未知参数 μ 的先验分布为 $U(0, 1)$, 求未知参数 μ 的贝叶斯估计
- C. 总体服从 $\text{Exp}(\lambda)$, 未知参数 λ 的先验分布为 $U(0, \infty)$, 求未知参数 λ 的贝叶斯估计
- D. 总体服从 $\text{Exp}(\lambda)$, 未知参数 λ 的先验分布为 $\Gamma(\alpha', \lambda')$, 求未知参数 λ 的贝叶斯估计
- E. 总体服从 $\Gamma(\alpha, \lambda)$, 未知参数 λ 的先验分布为 $\Gamma(\alpha', \lambda')$, 求未知参数 λ 的贝叶斯估计

解在贝叶斯估计中, 二次损失函数下的贝叶斯估计为后验分布的均值, 绝对误差损失函数下为中位数, $0-1$ 误差损失函数下为众数。要想结果相同, 也就是要求后验分布的均值, 同时它也是众数, 也是中位数, 即均值=众数=中位数。这就要求后验分布为对称分布才行, 如正态分布, t 分布等。

- A 中, 参数 μ 的后验分布是正态分布;
 - B 中, 参数 μ 的后验分布是正态分布;
 - C 中, 参数 λ 的后验分布是伽玛分布;
 - D 中, 参数 λ 的后验分布也是伽玛分布;
 - E 中, 参数 λ 的后验分布也是伽玛分布;
- 故选 AB。

【习题解答】

1. 从一组有效保单中抽取 100 份, 发现有三个索赔, 假如该险种的索赔频率 θ 的先验分布为贝塔 $(2, 200)$, 求 θ 的后验分布均值。

$$\begin{aligned} \text{解 } P(\theta | X=3) &= C_{3100} \theta^3 (1-\theta)^{97} \\ &= C_{3100} \theta^3 (1-\theta)^{200-1} \\ &= C_{3100} \theta^3 (1-\theta)^{97} \\ &= C_{3100} \theta^3 (1-\theta)^{200-1} \\ &= C_{3100} \theta^3 (1-\theta)^{200-1} \\ &= C_{3100} \theta^3 (1-\theta)^{200-1} \\ &= C_{3100} \theta^3 (1-\theta)^{200-1} \end{aligned}$$

2. 设 40 张同类保单, 用 X_i 表示第 i 张保单的索赔次数, 并设 $X_i \sim P(\lambda)$, $i=1, 2, \dots, 40$, 又设参数 λ 为随机变量, 且服从均值为 0.6, 方差为 0.02 的伽玛 (α, β) 分布, 并且已知观察到 40 张保单共有 18 次索赔, 试计算在平方损失函数下 λ 的贝叶斯估计。

解由已知条件可知 40 张报单的联合概率密度函数为

$$\begin{aligned} f(x_1, x_2, \dots, x_{40} | \lambda) &= \prod_{i=1}^{40} \lambda x_i! e^{-\lambda} \\ &= \lambda^{40} e^{-40\lambda} x_1! x_2! \dots x_{40}! \end{aligned}$$

对于 $\Gamma(\alpha, \beta)$ 分布, 均值为 α/β , 方差为 α/β^2 。

由已知有 $\alpha/\beta=0.6$, $\alpha/\beta^2=0.02$

得 $\alpha=18$, $\beta=30$

$$\begin{aligned} f(\lambda | x_1, x_2, \dots, x_{40}) &= \prod_{i=1}^{40} \lambda x_i! e^{-\lambda} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \\ &= \prod_{i=1}^{40} \lambda x_i! e^{-\lambda} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha+\sum x_i-1} e^{-(40+\beta)\lambda} \end{aligned}$$

∴ 在平方损失函数下， λ 的估计为

$$\hat{\lambda} = \alpha + \sum x_i / 40 + \beta = 0.5143$$

3. 在对某工业产品的保险研究中，需要知道该产品在一年内损坏的概率。已知同为用过两年的该产品 500 件，在一年观察期中观察到 15 件损坏，求该产品在一年内损坏概率的极大似然估计。

解设该产品在一年内损坏的概率为 q ，那么一年内损坏的件数 N 服从参数 500 与 q 的二项分布。似然函数 $L(q) = P(N=k) = C_{500}^k q^k (1-q)^{500-k}$

$$\ln [L(q)] = \ln C_{500}^k + k \ln q + (500-k) \ln (1-q)$$

令

$$d [\ln L(q)] / dq = k/q + k - 500 = 0$$

得到 $\hat{q} = k/500 = 15/500 = 0.03$

4. 试求一枚匀称硬币投掷 4 次所得正面朝上次数的分布函数经验分布函数 $F(x)$ 。

解正面 0 次朝上的概率为： $1/16$ ；

正面 1 次朝上的概率为： $4/16$ ；

正面 2 次朝上的概率为： $6/16$ ；

正面 3 次朝上的概率为： $4/16$ ；

正面 4 次朝上的概率为： $1/16$ 。

∴ $F(x)$ 图形如下。

$$F(x) = 0, \quad x < 0$$

$$1/16, \quad 0 \leq x < 1$$

$$5/16, \quad 1 \leq x < 2$$

$$11/16, \quad 2 \leq x < 3$$

$$15/16, \quad 3 \leq x < 4$$

1, $x \geq 4$

5. 下表显示了索赔情况调查所揭示的索赔分布情况，画出观察到的分布函数的图形，并用它估计在 150~200 美元之间发生索赔的概率。

| 索赔额频数 | 索赔额频数 |
|----------|-------|
| 0~100 | 45 |
| 100~200 | 10 |
| 200~500 | 2 |
| 500~1000 | 0 |
| 1000 以上 | 0 |

解设 q 为所估计的索赔频率，于是似然函数

$$L(q) = (1-q)^{100} \cdot 2 \lambda e^{-0.2 \lambda} \cdot \lambda e^{-0.6 \lambda} \\ = e^{-98.8 \lambda} \lambda^2 e^{-0.8 \lambda}$$

$$\ln L = -98.8 + 2 \lambda$$

$$\text{令 } d(\ln L) / d \lambda = -98.8 + 2 \lambda = 0$$

$$\therefore \lambda^{\wedge} = 298.8$$

所求索赔频率的极大似然估计为

$$q^{\wedge} = 1 - e^{-\lambda^{\wedge}} = 0.02$$

6. 在某一特定年度中，1000 份保险单中发生赔款 140 起。假定某一保单在一定时期遇到赔款的次数服从泊松分布。估测某特定保单持有人在 9 个月中不遇到损失的概率。

解理赔次数服从泊松分布，因此，我们假定某一保单在三个月内遇到的索赔次数是一个以 q 为均值的泊松变量，并且假定每一时间区段赔款相互独立。从而一份保单在 12 个月内总索赔次数是以 $4q$ 为均值的泊松变量。对于 1000 份保单，其均值为 $4000q$ 。通过赔款次数 140，我们估计 q 为 $140/4000 = 0.035$ ，

那么一保单在 9 个月中遇到索赔次数是以 0.105 (0.035×3) 为均值的泊松变量，则由泊松分布可知，其无理赔的概率为

$$(0.105)^0 e^{-0.105} = e^{-0.105} = 0.90$$

7. 下表概括了某保险公司 100 个赔款样本的赔款额状况。若该分布适用反对数正态分布模型，求其参数 μ 和 σ 的矩估计值，并估计一笔赔款规模超过 4000 美元的概率。

| 索赔额 | 频数 | 索赔额 | 频数 |
|-----------|----|-----------|----|
| 0~400 | 22 | 400~800 | 24 |
| 400~800 | 24 | 800~1200 | 32 |
| 800~1200 | 32 | 1200~1600 | 21 |
| 1200~1600 | 21 | 1600~2000 | 10 |
| 1600~2000 | 10 | 2000~2400 | 6 |
| 2000~2400 | 6 | 2400~2800 | 3 |
| 2400~2800 | 3 | 2800~3200 | 1 |
| 2800~3200 | 1 | 3200~3600 | 1 |
| 3200~3600 | 1 | 3600 以上 | 0 |
| 总数 | | 100 | |

解平均索赔额 $= 200 \times 22 + 600 \times 24 + \dots + 3400 \times 11 = 12162$
 方差 $= 200^2 \times 22 + 600^2 \times 24 + \dots + 3400^2 \times 11 - 12162^2 = 362944$
 而反对数正态分布的均值与方差分别为 $e^{\mu + \frac{1}{2}\sigma^2}$ 和 $e^{2\mu + \sigma^2} (e^{\sigma^2} - 1)$

$$e^{\mu + \frac{1}{2}\sigma^2} = 12162$$

$$e^{2\mu + \sigma^2} (e^{\sigma^2} - 1) = 362944$$

解得 $\mu = 6.993$ $\sigma = 0.469$

某一特定索赔 X 超过 4000 美元的概率等于 $\ln X$ 大于 8.29 的概率为

$$1 - \Phi\left(\frac{8.29 - 6.993}{0.469}\right) = 1 - \Phi(2.77) = 0.00280$$

8. 下表给出了 4000 份保单的赔款经历情况，其中每一保单承担风险期限为一年。假定每份保单经历赔款的次数服从泊松分布，用最大似然估计原理估计泊松分布的参数。

| 赔款次数 | 被观察的保单数 | 赔款次数 | 被观察的保单数 |
|------|---------|------|---------|
| 0 | 3288 | 3 | 34 |
| 1 | 642 | 4 | 6 |
| 2 | 266 | 5 | 3 |
| 总和 | | 4000 | |

解未发生索赔的概率为 $e^{-q} q^0 / 0!$

发生一次概率为 $e^{-q} q^1 / 1!$

发生二次概率为 $e^{-4000q} \frac{4000^2 q^2}{2!}$!

发生三次概率为 $e^{-4000q} \frac{4000^3 q^3}{3!}$!

\therefore 似然函数

$$L = e^{-4000q} \frac{4000^0 q^0}{0!} \frac{4000^1 q^1}{1!} \frac{4000^2 q^2}{2!} \frac{4000^3 q^3}{3!} \dots$$

$$e^{-4000q} \frac{4000^0 q^0}{0!} \frac{4000^1 q^1}{1!} \frac{4000^2 q^2}{2!} \frac{4000^3 q^3}{3!} \dots$$

$$e^{-4000q} \frac{4000^0 q^0}{0!} \frac{4000^1 q^1}{1!} \frac{4000^2 q^2}{2!} \frac{4000^3 q^3}{3!} \dots$$

$$e^{-4000q} \frac{4000^0 q^0}{0!} \frac{4000^1 q^1}{1!} \frac{4000^2 q^2}{2!} \frac{4000^3 q^3}{3!} \dots$$

$$= e^{-4000q} \frac{4000^0 q^0}{0!} \frac{4000^1 q^1}{1!} \frac{4000^2 q^2}{2!} \frac{4000^3 q^3}{3!} \dots$$

$$L = -4000q + 786 \ln q - 52.91$$

$$\text{令 } d(\ln L)/dq = -4000 + 786/q$$

$$\therefore q^* = 786/4000 = 0.196$$

9. 求上题中数据所代表的风险种类的赔款频率的置信度为 95% 的置信区间。

解由中心极限定理的正态逼近结论可知，

索赔总次数是以 $4000q$ 为均值， $4000q$ 为方差的正态变量。

$$\therefore N(4000q, 4000q) \sim N(0, 1)$$

令 $N = 486$ 并使 $N - 4000q$ 等于标准正态分布上、下 2.5% 点，即

$$786 - 4000q = \pm 1.96 \sqrt{4000q}$$

$$\therefore 16000000q^2 - 6303366q + 617796 = 0$$

$$\therefore q = 0.183 \text{ 或 } 0.211$$

即 q 的 95% 置信区间为 0.183 至 0.211

10. 某保险公司在某种业务中有 100 项索赔的样本，已在题 7 中给出。其索赔频率约为 0.015，计算当对此索赔经历设定完全可信性 ($k=0.05$, $p=0.9$) 时未到期责任的最小容量。

解由第 7 题可知，索赔额均值 m 为 1216，方差 σ^2 为 362944。因此 $\sigma/m=0.50$ ，

将此数字及对应 $k=0.05$, $p=0.9$ 的值 (见下表)

代入公式 $nF = n_0 [1 + (\sigma/m)^2]$,

可得 nF 为 1353，而已知索赔频率为 0.015，所以可以推算具有完全可信的最小未到期责任业务量约 90000 张保单。

k

p

0.30.20.10.050.01

0.93068271108227060

0.954396384153738416

0.9970166663265466358

0.99912027110834331108274

注： 对应于 k 与 p 的可信性参数 n_0 。

11. 在上题中的保险公司在其未到期业务中有 19307 份保险单。其总的风险保费为 366833 元，一年内索赔总额为 340575 元。试求保险公司在下一年对每份保险合同应收取的风险保费的可信性估计值。

解由上题，完全可信性的最小期望索赔次数为 1353，因此，对于索赔频率为 0.015 而未到期责任为 19307 份保单的期望索赔次数为 290。

由 $z = n_z n F = 0.46$

对 19303 份保单应收取的总风险保费的可信性估计

按公式 $C^z = (1-z) C^C + zC$

算出 $C^z = 366833z = 0.46C = 340575$

$\therefore C^z = 354754$

\therefore 保险公司在下一年对每一张保险合同应收取的风险保费可信性估计为 $354754/19307 = 18.37$ 元，而目前每一份收取的保费为 19.00 元。

12. 在城市机动车辆保险保单持有人的报案中，90% 的事故发生在市内，10% 则发生在市外。而非城市保单持有人的报案中，15% 的事故发生在市内，85% 发生在市外。已知该保险公司的保单持有人 80% 居住在市內。现得到一市外发生交通事故的报告，问该保单持有人系市內居住者的概率为多少？

解以 D 表示居住地，

L 表示事故发生地，D 分为城市 (u) 和乡村 (r)，L 分为 u 和 r。

那么

$$P(L=u|D=u) = 0.9$$

$$P(L=r|D=u) = 0.1$$

$$P(L=u|D=r) = 0.15$$

$$P(L=r|D=r) = 0.85$$

$$P(D=u) = 0.8$$

$$P(D=r) = 0.2$$

$$\therefore P(D=u|L=r)$$

$$= P(D=u) P(L=r|D=u) + P(D=r) P(L=r|D=r)$$

$$= 0.8 \times 0.1 + 0.2 \times 0.85$$

$$= 0.32$$

13. 设下表中的理赔记录用韦伯分布来拟合，试用 0.2 和 0.7 分位点估计参数 r 。

0.1

0.2

0.3

0.4

0.5

0.6

0.8

0.95

0.98

1.0

解 0.2 分位点为 0.25

0.7 分位点为 0.875

分别令 $0.2 = 1 - e^{-cxr}$

$0.7 = 1 - e^{-cxr}$

得 $x = -\ln 0.8c1r$ 和 $x = -\ln 0.3c1r$

将 0.25 和 0.875 代入上面两式有

$$0.25 = -\ln 0.8c1r$$

$$0.875 = -\ln 0.3c1r$$

$$\therefore r = 1.35$$

14. 已知总体分布为 $\Gamma(\alpha, \lambda)$ ，其中参数 λ 的先验分布为 $\exp(\lambda')$ ，试问其后验分布是什么？

$$\begin{aligned} \text{解 } f(\lambda | x_1, x_2, \dots, x_n) &\propto f(x_1, x_2, \dots, x_n | \lambda) f(\lambda) \\ &\propto \prod_{i=1}^n (\lambda^\alpha e^{-\lambda x_i}) \cdot e^{-\lambda'} \quad \lambda = \lambda n^\alpha e^{-x_i \lambda} \cdot e^{-\lambda'} \end{aligned}$$

$$= \lambda n^\alpha e^{-(x_i + \lambda') \lambda}$$

即其分布为 $\Gamma(n\alpha, \sum x_i + \lambda' + 1)$

15. 设某保险人经营某种车辆险，对过去所发生的 1000 次理赔情况作了记录，平均理赔为 2200，又按赔付金额分为 5 档，各档中的记录次数如下。

理赔额

0~1000

1000~2000

2000~3000

3000~4000

4000~5000

5000 以上

次数

200

300

250

150

100

0

试用 χ^2 分布检验判断能否用指数分布模拟个别理赔额的分布（假设置信水平为 99.5%）。

解先假设个别理赔额 $X \sim \exp(\lambda)$

用矩方法或最大似然法都得到 $\lambda \hat{=} 1/\bar{x} = 1/2200$

为了计算 E_i ，先计算个别理赔额落入每一个档次内的概率，

比如 2000~3000 内概率为

$$\begin{aligned} \int_{2000}^{3000} \lambda e^{-\lambda x} dx &= e^{-2000\lambda} - e^{-3000\lambda} \\ &= 0.1472 \end{aligned}$$

在 2000~3000 内平均次数 $E_2 = 147.2$

类似地算出 $E_1 = 365.3$ ， $E_3 = 231.8$ ， $E_4 = 93.4$ ， $E_5 = 59.3$ ， $E_6 = 103$

χ^2 统计量的值为

$$\chi^2 = (200-365.3)^2/365.3 + (300-231.8)^2/231.8 + \dots + (0-103)^2/103 = 331.89$$

由于只有一个参数，并且记录完整。 χ^2 分布的自由度为 $6-1-1=4$ ，查 χ^2 分布表在置信水平 99.5%下值为 14.86，远低于观察值。

故应拒绝原假设，即选择指数分布不恰当。

16. 根据上题所给出的损失记录，画出频率直方图和频率折线图。问：

- (1) 这两条曲线和概率论中什么曲线相似？
- (2) 根据频率直方图和频率折线图，你准备用什么分布来模拟个别理赔额分布？

解频率直方图如下。

频率折线图如下。

由于这两条曲线具有尾部较厚的特点，可以考虑用广义帕雷托分布模拟。

17. 对于承保某一类风险，某保险公司将其保费计算基于平均赔款额为 1200 美元的假定之上。1980 年中，这类保险项目发生了 1243 次赔款，平均额为 1283.70 美元，其观察值的标准差等于 1497.31 美元。是否有证据能说明该公司在计算保险费率时所假定的平均赔款额太低了呢？

解根据中心极限定理及独立随机变量之和的性质，我们假设这 1243 起索赔支出数的分布近似于均值为 1243μ ，方差为 $1243\sigma^2$ 的正态分布

即均值为 $1243 \times 1200 = 1491600$

标准差近似为 $1243 \times (1497.31)^2 = 52789$

1980 年实际总支出额为 $1283 \times 1243 = 159639$ 。在原假设条件下，真正的平均索赔额为 1200，则 1243 起赔款中支出为 1595639 美元，或更多数目的概率为

$$1 - \Phi(1595639 - 1491600 / \sqrt{52789}) = 1 - \Phi(1.97) = 0.024$$

显然偏低，即平均索赔额假设太低了。

18. 试证明采用绝对误差函数时，参数 θ 的贝叶斯估计是后验分布的中位数。

$$\begin{aligned} \text{证： } E[\text{Loss}(\hat{\theta}, \theta)] &= \int_{-\infty}^{+\infty} \text{Loss}(\hat{\theta}, \theta) f(\theta | x) d\theta \\ &= \int_{-\infty}^{+\infty} |\hat{\theta} - \theta| f(\theta | x) d\theta \\ &= \int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) f(\theta | x) d\theta + \\ &\quad \int_{\hat{\theta}}^{+\infty} \theta - \hat{\theta} f(\theta | x) d\theta \\ &= \int_{-\infty}^{\hat{\theta}} \hat{\theta} f(\theta | x) d\theta - \int_{-\infty}^{\hat{\theta}} \theta f(\theta | x) d\theta \\ &\quad + \int_{\hat{\theta}}^{+\infty} \theta f(\theta | x) d\theta - \int_{\hat{\theta}}^{+\infty} \hat{\theta} f(\theta | x) d\theta \\ &= \hat{\theta} - 2 \int_{\hat{\theta}}^{+\infty} \theta f(\theta | x) d\theta + \int_{-\infty}^{\hat{\theta}} \theta f(\theta | x) d\theta + \int_{\hat{\theta}}^{+\infty} \theta f(\theta | x) d\theta \\ \therefore \hat{\theta} &= E[\text{Loss}(\hat{\theta}, \theta)] = 1 - 2 \int_{\hat{\theta}}^{+\infty} \theta f(\theta | x) d\theta + 2 \int_{-\infty}^{\hat{\theta}} \theta f(\theta | x) d\theta \end{aligned}$$

$$E(\hat{\theta} | X) = \int_{-\infty}^{+\infty} \theta f(\theta | X) d\theta$$

令 $E[\text{Loss}(\hat{\theta}, \theta)] = 0$ 有

$$\int_{-\infty}^{+\infty} \theta f(\theta | X) d\theta = 1/2$$

即 $\hat{\theta}$ 为后验分布的中位数。

19. 设 θ 是某险种的索赔频率，抽取 8 份保单，发现在有效期内有三三次索赔，假如先验分布为

- (1) $\theta \sim U(0, 1)$;
- (2) $\theta \sim f(\theta) = 2(1-\theta), 0 < \theta < 1$ 。

分别求 θ 的后验分布。

解 (1) $P(X=3 | \theta) = C_{38} \theta^3 (1-\theta)^5$

$$P(\theta | X=3) = \frac{C_{38} \theta^4 (1-\theta)^6}{\int_{-\infty}^{+\infty} C_{38} \theta^4 (1-\theta)^6 d\theta}$$

$\therefore \theta$ 后验分布为 Beta(4, 6)

$$(2) P(\theta | X=3) = \frac{C_{38} \theta^4 (1-\theta)^6}{\int_{-\infty}^{+\infty} C_{38} \theta^4 (1-\theta)^6 d\theta}$$

其中 C 为常数，故 θ 后验分布为 Beta(4, 7)

20. 假设某一工业新产品的寿命服从均匀分布，参数为 e ，其密度函数为 $f(t) = 1/e$ 。现对 7 件这样的产品进行观察，得到如下的产品余命为 3, 4, 5, 7, 8, 10, 12。试求 e 的最小二乘估计。

解由 $f(t) = 1/e$ 可知

$$S(t) = e^{-t/e}$$

$$SS = \sum_{i=1}^7 (t_i - S(t_i))^2 = \sum_{i=1}^7 (t_i - e^{-t_i/e})^2$$

$$= \sum_{i=1}^7 (t_i - e^{-t_i/e})^2$$

$$dSS/de = 2 \sum_{i=1}^7 (t_i - e^{-t_i/e}) \cdot (-1/e^2) e^{-t_i/e}$$

令 $dS_{de}=0$

可得到 e 的最小二乘估计表达式

$$e^{\wedge} = \frac{\sum_{i=1}^n t_i^2}{\sum_{i=1}^n t_i - \sum_{i=1}^n t_i^2} (S^{\wedge}(t_i) + S^{\wedge}(t_{i-1}))$$

具体计算如下表所示。

$$\frac{t_i^2 S^{\wedge}(t_i) + t_i^2 S^{\wedge}(t_{i-1})}{[S^{\wedge}(t_i) + S^{\wedge}(t_{i-1})]}$$

| |
|---------------------|
| 1390.92.85 |
| 24160.83.4 |
| 35250.73.75 |
| 47490.54.2 |
| 58640.43.6 |
| 6101000.12.5 |
| 7121440.00.6 |
| Σ 4940720.90 |

故有 $e^{\wedge} = 40749 - 20.9 = 14.48$