

二、三维列联表确切概率检验的 有关算法*

李 建 立

(上海第二医科大学生物统计教研室)

AN ALGORITHM IN THE EXACT PROBABILITY TEST IN TWO-AND THREE-DIMENSIONAL CONTINGENCY TABLES

Li Jian-li

(Section of Biostatistics Shanghai Second Medical University)

Abstract

In this paper an algorithm for finding out all possible two-and three-dimensional contingency tables with the fixed marginals is given. In this algorithm, a network consisting of nodes and arcs is constructed and several paths are set up. Corresponding to each node and each path a unique nonnegative value is defined and the conditions for these values are given. From these values all possible contingency tables can be found. Furthermore, the formulas of calculating P values in exact probability tests are given.

对于一个列联表当样本较小时,就会产生用 χ^2 统计量或 G^2 统计量来作关于交互作用的显著性检验不够精确的问题,于是就需作确切概率检验。本文就二、三维列联表在某些条件下的确切概率计算的有关算法加以讨论。

§ 1. 列联表的样本分布

若对所抽的样本大小没有任何约束,且每个格子的观察频数 $\{X_{\theta}\}$ 可看作服从独立的 Poisson 分布,各自的均值为 $\{m_{\theta}\}$,其中 θ 表示下标集: $\{i_1, i_2, \dots, i_n\}$ 。对于一个 n 维列联表 $\{X_{\theta}\}$,其密度函数(联合分布密度函数)为

$$f(\{X_{\theta}\}) = \prod_{\theta} \frac{m_{\theta}^{X_{\theta}} e^{-m_{\theta}}}{X_{\theta}!} \quad (1)$$

有时往往要考虑样本大小固定的独立 Poisson 抽样, 假设 $\sum_{\theta} X_{\theta} = N$. 可知 $\{X_{\theta}\}$ 的分布为多项分布, 其密度函数为

$$f(\{X_{\theta}\}) = \frac{N!}{\prod_{\theta} X_{\theta}!} \prod_{\theta} \left(\frac{m_{\theta}}{N}\right)^{X_{\theta}}, \quad (2)$$

其中 m_{θ} 为 X_{θ} 的期望值, $m_{\theta}/N = p_{\theta}$ ($\sum_{\theta} p_{\theta} = 1$).

在实验设计中, 会确定各组的总数, 亦即固定组态(边际和)大小 $c_{\theta_i} (\theta_i \subseteq \theta)$, c_{θ_i} 即 $\left\{\sum_{\theta-\theta_i} X_{\theta}\right\}$ (其中 $\theta - \theta_i$ 表示下标集 θ 与下标集 θ_i 之差. $\sum_{\theta-\theta_i}$ 表示对 $\theta - \theta_i$ 中所有下标求和), 其结果得到的分布是多个多项分布之积. 推导如下:

因为 c_{θ_i} 的边际分布为

$$f\left(\left\{\sum_{\theta-\theta_i} X_{\theta}\right\}\right) = \frac{N!}{\prod_{\theta-\theta_i} \left(\sum_{\theta-\theta_i} X_{\theta}\right)!} \prod_{\theta-\theta_i} \left(\frac{\sum_{\theta-\theta_i} m_{\theta}}{N}\right)^{\sum_{\theta-\theta_i} X_{\theta}},$$

所以条件分布为

$$\begin{aligned} f(\{X_{\theta}\} \mid \left\{\sum_{\theta-\theta_i} m_{\theta} = \sum_{\theta-\theta_i} X_{\theta}\right\}) &= f(\{X_{\theta}\}) / f\left(\left\{\sum_{\theta-\theta_i} X_{\theta}\right\}\right) \\ &= \prod_{\theta_i} \left[\frac{\left(\sum_{\theta-\theta_i} X_{\theta}\right)!}{\prod_{\theta-\theta_i} X_{\theta}!} \prod_{\theta-\theta_i} \left(\frac{m_{\theta}}{\sum_{\theta-\theta_i} X_{\theta}}\right)^{X_{\theta}} \right]. \end{aligned} \quad (3)$$

§ 2. 二、三维表的确切概率计算

二维列联表 $\{X_{ij}\}$ 有以下结果^[4]:

结果 1. 若原假设 $H_0: \mu_{12} = 0$ 成立, 即二变量相互独立, 且有条件: 边际和 $X_{i.}$ 与 $X_{.j}$ 固定, 则分布密度函数为

$$f^*(\{X_{ij}\}) = \frac{\prod_i X_{i.}! \prod_j X_{.j}!}{\prod_{ij} X_{ij}! N!}, \quad (4)$$

其中 $X_{i.} = \sum_j X_{ij}$, $X_{.j} = \sum_i X_{ij}$, $N = X_{..} = \sum_{ij} X_{ij}$.

为考察观察到的列联表 $\{X_{ij}^*\}$ 的条件显著性, 可计算下述概率值:

$$P = S(\{X_{ij}^*\}) = \sum f^*(\{X'_{ij}\}), \quad (5)$$

其中求和是对所有 $\{X'_{ij}\}$ 进行的, 而 $f^*(\{X'_{ij}\}) \leq f^*(\{X_{ij}^*\})$.

根据(4)式,因为 $X_{i.}$ 与 $X_{.j}$ 固定,所以上述选择 $\{X'_{ij}\}$ 的条件可改写为

$$\frac{1}{\prod_{ij} X'_{ij}!} \leq \frac{1}{\prod_{ij} X''_{ij}!}, \quad (6)$$

即

$$\prod_{ij} \frac{X'_{ij}!}{X''_{ij}!} \geq 1. \quad (7)$$

对于事先给定检验水平 α , 作双侧检验,若 $P \leq \alpha$, 则拒绝原假设.

三维列联表 $\{X_{ijk}\}$ 的情况较复杂,本文只就二种情况作讨论:有以下结果:

结果 2. 若原假设 $H_0: u_{11} = u_{12} = u_{21} = u_{22} = 0$ 成立,即三变量相互独立,且有条件: 边际和 $X_{i..}, X_{.j.}$ 及 $X_{...k}$ 固定,则分布密度函数为

$$f^*(\{X_{ijk}\}) = \frac{\prod_i X_{i..}! \prod_j X_{.j.}! \prod_k X_{...k}!}{\prod_{ijk} X_{ijk}! (N!)^3}, \quad (8)$$

其中 $N = X_{...} = \sum_{ijk} X_{ijk}$, $X_{i..} = \sum_k X_{ijk}$, $X_{.j.} = \sum_i X_{ijk}$, $X_{...k} = \sum_{ij} X_{ijk}$.

推导如下:

根据边际和固定条件,若原假设成立,即三变量相互独立,可推得

$$\begin{aligned} f(\{X_{ijk}\} | \{m_{i..} = X_{i..}, m_{.j.} = X_{.j.}, m_{...k} = X_{...k}\}) \\ = f(\{X_{ijk}\} | \{m_{i..} = X_{i..}\}) / f(\{X_{.j.}\}) / f(\{X_{...k}\}) \\ = \frac{\prod_i X_{i..}! \prod_{ijk} (m_{ijk})^{X_{ijk}} \prod_j X_{.j.}! N^N \prod_k X_{...k}! N^N}{\prod_{ijk} X_{ijk}! \prod_i (X_{i..})^{X_{i..}} N! \prod_j (m_{.j.})^{X_{.j.}} N! \prod_k (m_{...k})^{X_{...k}}} \end{aligned} \quad (9)$$

当 H_0 成立时,应有 $m_{ijk} = \frac{X_{i..} X_{.j.} X_{...k}}{N^2}$, 即可由(9)推得

$$\begin{aligned} f^*(\{X_{ijk}\}) \\ = \frac{\prod_i X_{i..}! \prod_j X_{.j.}! \prod_k X_{...k}! \prod_i (X_{i..})^{X_{i..}} \prod_j (X_{.j.})^{X_{.j.}} \prod_k (X_{...k})^{X_{...k}} N^{2N}}{\prod_{ijk} X_{ijk}! (N!)^3 N^{2N} \prod_i (X_{i..})^{X_{i..}} \prod_j (X_{.j.})^{X_{.j.}} \prod_k (X_{...k})^{X_{...k}}} \\ = \frac{\prod_i X_{i..}! \prod_j X_{.j.}! \prod_k X_{...k}!}{\prod_{ijk} X_{ijk}! (N!)^3}. \end{aligned} \quad (10)$$

对于一个观察到的列联表 $\{X''_{ijk}\}$, 在原假设(三变量相互独立)成立及边际和 $X_{i..}$, $X_{.j.}$, $X_{...k}$ 固定的条件下的条件概率值记为 $f^*(\{X''_{ijk}\})$.

为考察观察到的列联表 $\{X_{ijk}^*\}$ 的条件显著性,可计算下述概率值:

$$P = S(\{X_{ijk}^*\}) = \sum f^*(\{X'_{ijk}\}), \quad (11)$$

其中求和是对所有 $\{X'_{ijk}\}$ 进行的,而 $f^*(\{X'_{ijk}\}) \leq f^*(\{X_{ijk}^*\})$.

根据(11)式,因为 $X_{i..}, X_{.i.}, X_{..k}$ 固定,所以上述选择 $\{X'_{ijk}\}$ 的条件可改写为

$$\prod_{ijk} \frac{X'_{ijk}}{X_{ijk}^*} \geq 1. \quad (12)$$

对于事先给定检验水平 α ,作双侧检验,若 $P \leq \alpha$, 则拒绝原假设.

结果 3. 若原假设 $H_0: u_{12} = u_{13} = u_{23} = 0$ 成立,即二变量(变量 1 与变量 2 及变量 1 与变量 3)完全独立,且有条件: 边际和 $X_{i..}$ 与 $X_{.ik}$ 固定,则分布密度函数为

$$f^*(\{X_{ijk}\}) = \frac{\prod_i X_{i..}! \prod_{jk} X_{.jk}!}{\prod_{ijk} X_{ijk}! N!}. \quad (13)$$

推导如下:

根据边际和固定条件,若原假设成立,则可推得

$$\begin{aligned} f(\{X_{ijk}\} | \{m_{i..} = X_{i..}, m_{.ik} = X_{.ik}\}) \\ = f(\{X_{ijk}\} | \{m_{i..} = X_{i..}\}) / f(\{X_{.ik}\}) \\ = \frac{\prod_i X_{i..}! \prod_{ijk} (m_{ijk})^{X_{ijk}} \prod_{ik} X_{.ik}! N^N}{\prod_{ijk} X_{ijk}! \prod_i (X_{i..})^{X_{i..}} \prod_{ik} (m_{.ik})^{X_{.ik}} N!}. \end{aligned} \quad (14)$$

当 H_0 成立时,应有 $m_{ijk} = \frac{X_{i..} X_{.ik}}{N}$, 即可由(14)推得

$$\begin{aligned} f^*(\{X_{ijk}\}) &= \frac{\prod_i X_{i..}! \prod_i (X_{i..})^{X_{i..}} \prod_{jk} (X_{.jk})^{X_{.jk}} \prod_{ik} X_{.ik}! N^N}{\prod_{ijk} X_{ijk}! N^N N! \prod_{jk} (X_{.jk})^{X_{.jk}} \prod_i (X_{i..})^{X_{i..}}} \\ &= \frac{\prod_i X_{i..}! \prod_{jk} X_{.jk}!}{\prod_{ijk} X_{ijk}! N!}. \end{aligned} \quad (15)$$

同样,若要考察观察到的 $\{X_{ijk}^*\}$ 的条件显著性,只需计算概率值

$$P = S(\{X_{ijk}^*\}) = \sum_{f^*(\{X'_{ijk}\}) < f^*(\{X_{ijk}^*\})} f^*(\{X'_{ijk}\}).$$

易见,选择 $\{X'_{ijk}\}$ 之条件为

$$\prod_{ijk} \frac{X'_{ijk}}{X_{ijk}^*} \geq 1. \quad (16)$$

§ 3. 产生各种可能的列联表的算法

在列联表某些边际和固定的情况下,可用网络来求得各种可能的列联表,此网络由节点和连线组成.

对于二维列联表 $\{X_{ij}\} (1 \leq i \leq I, 1 \leq j \leq J)$, 每一节点记为 $ND_{ij}^{(\alpha)}$ ($1 \leq i \leq I-1, 0 \leq j \leq J$), 其中 $[\theta_{ij}] = n_{ij}n_{i,j-1} \cdots n_{i0} \cdots n_{i0}$ 为节点序列号. 且可有以下分解: $[\theta_{ij}] = n_{ij}[\theta_{i,j-1}]$. 对于节点 $ND_{ij}^{(\alpha)}$ 和 $ND_{i,j-1}^{(\alpha)}$, 若有 $[\theta_{i,j-1}]^* = [\theta_{i,j-1}]$, 则这两节点间可有连线相连, 这些连线构成了通路. 对应于每一个节点和每条通路有唯一确定的一个非负整数值, 记为 $V_{ij}^{(\alpha)}$, 其中 α 表示通路号, 且令 $V_{i0}^{(\alpha)} = 0$.

定义 1. 二维列联表 $\{X_{ij}\}$ 中第 (i, j) 格第 α 种可能的频数

$$X_{ij}^{(\alpha)} = V_{ij}^{(\alpha)} - V_{ij}^{(\alpha)}{}_{-1}.$$

引理 1. 二维列联表 $\{X_{ij}\} (1 \leq i \leq 2, 1 \leq j \leq J)$ 中, 若 $X_{i.}$ 与 $X_{.j}$ 固定, 则 $V_{ij}^{(\alpha)}$ 满足以下条件:

$$\max \left(V_{ij}^{(\alpha)}{}_{-1}, \sum_{l=1}^j X_{.l} - X_{2.} \right) \leq V_{ij}^{(\alpha)} \leq \min (V_{ij}^{(\alpha)}{}_{-1} + X_{.j}, X_{1.}). \quad (17)$$

证. 在 $\{X_{ij}\}$ 中, 若 $X_{il}^{(\alpha)} (1 \leq l \leq j-1)$ 已确定, 则显然有

$$X_{ij}^{(\alpha)} \leq X_{1.} - \sum_{l=1}^{j-1} X_{il}^{(\alpha)}.$$

根据定义 1 有

$$V_{ij}^{(\alpha)} - V_{ij}^{(\alpha)}{}_{-1} \leq X_{1.} - \sum_{l=1}^{j-1} (V_{il}^{(\alpha)} - V_{il}^{(\alpha)}{}_{-1}),$$

$$V_{ij}^{(\alpha)} \leq X_{1.} - \sum_{l=1}^{j-1} (V_{il}^{(\alpha)} - V_{il}^{(\alpha)}{}_{-1}) + V_{ij}^{(\alpha)}{}_{-1},$$

所以

$$V_{ij}^{(\alpha)} \leq X_{1.}. \quad (18)$$

又因为 $X_{ij}^{(\alpha)} \leq X_{.j}$, 所以 $V_{ij}^{(\alpha)} - V_{ij}^{(\alpha)}{}_{-1} \leq X_{.j}$, 即有

$$V_{ij}^{(\alpha)} \leq X_{.j} + V_{ij}^{(\alpha)}{}_{-1}. \quad (19)$$

根据(18)和(19)有

$$V_{ij}^{(\alpha)} \leq \min (V_{ij}^{(\alpha)}{}_{-1} + X_{.j}, X_{1.}). \quad (20)$$

另一方面, 由定义 1 可知

$$V_{ij}^{(\alpha)} - V_{ij}^{(\alpha)}{}_{-1} = X_{ij}^{(\alpha)} \geq 0,$$

即

$$V_{ij}^{(\alpha)} \geq V_{ij}^{(\alpha)}{}_{-1}. \quad (21)$$

又因为 $\sum_{l=1}^j (X_{.l} - X_{1l}) \leq X_{2.}$, 即 $\sum_{l=1}^j X_{.l} - \sum_{l=1}^j (V_{il}^{(\alpha)} - V_{il}^{(\alpha)}{}_{-1}) \leq X_{2.}$, 所以

$$V_{ij}^{(\alpha)} \geq \sum_{l=1}^i X_{.l} - X_{i..} \quad (22)$$

根据(21)和(22)有

$$V_{ij}^{(\alpha)} \geq \max \left(V_{i,j-1}^{(\alpha)}, \sum_{l=1}^i X_{.l} - X_{i..} \right). \quad (23)$$

由(20)和(23)可得结论.

定理 1. 二维列联表 $\{X_{ij}\}$ ($1 \leq i \leq I, 1 \leq j \leq J$) 中, 若 $X_{i.}$ 与 $X_{.j}$ 固定, 则 $V_{ij}^{(\alpha)}$ ($1 \leq i \leq I-1, 1 \leq j \leq J$) 应满足以下条件:

$$\begin{aligned} & \max \left(V_{i,j-1}^{(\alpha)}, \sum_{l=1}^i X_{.l} - \sum_{k=1}^{i-1} V_{ki}^{(\alpha)} - \sum_{k=i+1}^I X_{k.} \right) \\ & \leq V_{ij}^{(\alpha)} \leq \min \left(V_{i,j-1}^{(\alpha)} + X_{.j} - \sum_{k=1}^{i-1} (V_{ki}^{(\alpha)} - V_{k,i-1}^{(\alpha)}), X_{i.} \right). \end{aligned} \quad (24)$$

证. 不失一般性, 假定 $X_{ij}^{(\alpha)}$ 已确定 ($1 \leq k \leq i-1, 1 \leq j \leq J$). 作一个二维列联表 $\{X'_{ij}^{(\alpha)}\}$ ($1 \leq i \leq 2$), 其中 $X'_{1i}^{(\alpha)} = X_{ij}^{(\alpha)}, X'_{2i}^{(\alpha)} = \sum_{k=i+1}^I X_{ki}^{(\alpha)}$. 于是有

$$X'_{1.} = X_{i.}, X'_{2.} = \sum_{k=i+1}^I X_{k.}, X'_{.j} = X_{.j} - \sum_{k=1}^{i-1} X_{kj}^{(\alpha)} \quad (1 \leq j \leq J). \quad (25)$$

据引理 1 有

$$\max \left(V'_{1,j-1}^{(\alpha)}, \sum_{i=1}^I X'_{.i} - X'_{2.} \right) \leq V'_{ij}^{(\alpha)} \leq \min(V'_{1,j-1}^{(\alpha)} + X'_{.j}, X'_{1i}^{(\alpha)}). \quad (26)$$

事实上应有 $V'_{ij}^{(\alpha)} = V_{ij}^{(\alpha)}, V'_{.j-1} = V_{.j-1}^{(\alpha)}$. 再将(25)代入(26), 即有

$$\begin{aligned} & \max \left(V_{i,j-1}^{(\alpha)}, \sum_{l=1}^i (X_{.l} - \sum_{k=1}^{i-1} X_{kl}^{(\alpha)}) - \sum_{k=i+1}^I X_{k.} \right) \\ & \leq V_{ij}^{(\alpha)} \leq \min \left(V_{i,j-1}^{(\alpha)} + X_{.j} - \sum_{k=1}^{i-1} X_{ki}^{(\alpha)}, X_{i.} \right). \end{aligned} \quad (27)$$

于是根据定义 1 中 $X_{ij}^{(\alpha)}$ 与 $V_{ij}^{(\alpha)}, V_{i,j-1}^{(\alpha)}$ 的关系即可得结论.

这样, 一个二维列联表可用网络的一条通路来表示.

对于三维列联表 $\{X_{ijk}\}$ ($1 \leq i \leq I, 1 \leq j \leq J, 1 \leq k \leq K$), 每一节点记为

$$ND_{ijk}^{(\theta_{ijk})} \quad (1 \leq i \leq I, 0 \leq j \leq J, 1 \leq k \leq K-1),$$

其中 $[\theta_{ijk}] = n_{ijk} n_{i,i-1,k} \cdots n_{i0k} \cdots n_{i01}$ 为节点序列号, 且可有以下分解: $[\theta_{ijk}] = n_{ijk} [\theta_{i,i-1,k}]$. 对于节点 $ND_{ijk}^{(\theta_{ijk})}$ 和 $ND_{i,j-1,k}^{(\theta_{i,j-1,k})}^*$ 若有 $[\theta_{i,i-1,k}]^* = [\theta_{i,i-1,k}]$, 则这两节点间可有连线相联, 这些连线就构成了通路. 对应于每一个节点和每条道路有唯一确定的一个非负整数值, 记为 $V_{ijk}^{(\alpha)}$, 且令 $V_{i0k}^{(\alpha)} = 0$.

定义 2. 三维列联表 $\{X_{ijk}\}$ 中第 (i, j, k) 格第 α 种可能的频数

$$X_{ijk}^{(\alpha)} = V_{ijk}^{(\alpha)} - V_{i,j-1,k}^{(\alpha)}.$$

引理 2. 三维列联表 $\{X_{ijk}\}$ ($1 \leq i \leq I, 1 \leq j \leq J, 1 \leq k \leq 2$) 中, 若 $X_{i1.}, X_{i..k}$ 及

$X_{i,k}$ 固定, 则 $V_{ij}^{(s)}$ 满足以下条件

$$\begin{aligned} \max \left(V_{ij}^{(s)} - 1, \sum_{r=1}^i X_{r,i} - \sum_{r=1}^{i-1} V_{ij}^{(s)} - \sum_{r=i+1}^j X_{r,i}, \sum_{r=1}^j X_{i,r} - X_{i,i} \right) &\leq V_{ij}^{(s)} \\ &\leq \min \left(V_{ij}^{(s)} - 1, \sum_{r=1}^{i-1} (V_{ij}^{(s)} - V_{ij}^{(s-1)}), X_{i,i}, V_{ij}^{(s)} + X_{i,i} \right). \end{aligned} \quad (28)$$

证. 在 $\{X_{ij}\}$ 中, 若 $X_{ij}^{(s)} (1 \leq s \leq j-1)$ 已确定, 则根据定理 1 应有

$$\max \left(V_{ij}^{(s)} - 1, \sum_{r=1}^i X_{r,i} - \sum_{r=1}^{i-1} V_{ij}^{(s)} - \sum_{r=i+1}^j X_{r,i} \right) \leq V_{ij}^{(s)}. \quad (29)$$

又因为

$$\begin{aligned} \sum_{r=1}^j (X_{i,r} - X_{ij}^{(s)}) &\leq X_{i,i}, \sum_{r=1}^j X_{i,r} - \sum_{r=1}^{i-1} (V_{ij}^{(s)} - V_{ij}^{(s-1)}) \leq X_{i,i}, \\ \sum_{r=1}^j X_{i,r} - V_{ij}^{(s)} &\leq X_{i,i}, \end{aligned}$$

所以

$$V_{ij}^{(s)} \geq \sum_{r=1}^j X_{i,r} - X_{i,i}. \quad (30)$$

据(29)和(30)有

$$\max \left(V_{ij}^{(s)} - 1, \sum_{r=1}^i X_{r,i} - \sum_{r=1}^{i-1} V_{ij}^{(s)} - \sum_{r=i+1}^j X_{r,i}, \sum_{r=1}^j X_{i,r} - X_{i,i} \right) \leq V_{ij}^{(s)}. \quad (31)$$

另一方面, 据定理 1 应有

$$V_{ij}^{(s)} \leq \min \left(V_{ij}^{(s)} - 1, \sum_{r=1}^{i-1} (V_{ij}^{(s)} - V_{ij}^{(s-1)}), X_{i,i} \right). \quad (32)$$

又因为 $X_{ij}^{(s)} \leq X_{ij}$, $V_{ij}^{(s)} - V_{ij}^{(s-1)} \leq X_{ij}$, 所以

$$V_{ij}^{(s)} \leq V_{ij}^{(s-1)} + X_{ij}. \quad (33)$$

据(32)和(33)有

$$V_{ij}^{(s)} \leq \min \left(V_{ij}^{(s-1)} + X_{ij} - \sum_{r=1}^{i-1} (V_{ij}^{(s)} - V_{ij}^{(s-1)}), X_{i,i}, V_{ij}^{(s-1)} + X_{ij} \right). \quad (34)$$

由(31)和(34)即可得结论.

定理 2. 三维列联表 $\{X_{ijk}\} (1 \leq i \leq I, 1 \leq j \leq J, 1 \leq k \leq K)$ 中, 若 $X_{i,i}$, $X_{i,k}$ 及 $X_{i,k}$ 固定, 则 $V_{ijk}^{(s)}$ 满足以下条件

$$\begin{aligned} \max \left(V_{ijk}^{(s)} - 1, \sum_{r=1}^i X_{r,i} - \sum_{r=1}^{i-1} V_{ijk}^{(s)} - \sum_{r=i+1}^j X_{r,i}, \right. \\ \left. \sum_{r=1}^j X_{i,r} - \sum_{r=1}^{k-1} V_{ijk}^{(s)} - \sum_{r=k+1}^k X_{i,r} \right) &\leq V_{ijk}^{(s)} \\ &\leq \min \left(V_{ijk}^{(s)} - 1, \sum_{r=1}^{i-1} (V_{ijk}^{(s)} - V_{ijk}^{(s-1)}), X_{i,i}, \right. \\ &\quad \left. \sum_{r=1}^{k-1} (V_{ijk}^{(s)} - V_{ijk}^{(s-1)}), X_{i,k} \right) \end{aligned}$$

$$V_{i,j-1,k}^{(\alpha)} + X_{ij} - \sum_{t=1}^{k-1} (V_{ij}^{(\alpha)} - V_{i,j-1,t}^{(\alpha)}). \quad (35)$$

证. 不失一般性, 假定 $X_{ij}^{(\alpha)}$ 已确定 ($1 \leq i \leq I, 1 \leq j \leq J, 1 \leq t \leq k-1$). 作一个三维列联表 $\{X_{ijk}^{(\alpha)}\} (1 \leq k \leq 2)$, 其中 $X_{ij1}^{(\alpha)} = X_{ij}^{(\alpha)}, X_{ij2}^{(\alpha)} = \sum_{t=k+1}^k X_{ij}^{(\alpha)}$. 于是有

$$\begin{aligned} X_{i1}^{(\alpha)} &= X_{i,k}, X_{i2}^{(\alpha)} = X_{i,k}, X_{i3}^{(\alpha)} = \sum_{t=k+1}^k X_{i,t}, X_{i4}^{(\alpha)} = \sum_{t=k+1}^k X_{i,t}, \\ X_{ij}^{(\alpha)} &= X_{ij} - \sum_{t=1}^{k-1} X_{ij}^{(\alpha)} (1 \leq i \leq I, 1 \leq j \leq J). \end{aligned} \quad (36)$$

据引理 2 有

$$\begin{aligned} \max \left(V_{i,j-1,1}^{(\alpha)}, \sum_{i=1}^I X_{i1}^{(\alpha)} - \sum_{t=1}^{i-1} V_{t1}^{(\alpha)} - \sum_{t=i+1}^I X_{t1}^{(\alpha)}, \sum_{i=1}^I X_{i2}^{(\alpha)} - X_{i2}^{(\alpha)} \right) &\leq V_{ij}^{(\alpha)} \\ &\leq \min \left(V_{i,j-1,1}^{(\alpha)} + X_{i1}^{(\alpha)} - \sum_{t=1}^{i-1} (V_{t1}^{(\alpha)} - V_{t,j-1,1}^{(\alpha)}), X_{i1}^{(\alpha)}, V_{i,j-1,1}^{(\alpha)} + X_{i1}^{(\alpha)} \right). \end{aligned} \quad (37)$$

事实上应有 $V_{ij}^{(\alpha)} = V_{ij}^{(\alpha)}, V_{i,j-1,1}^{(\alpha)} = V_{i,j-1,k}^{(\alpha)}$. 再将(36)代入(37)即得结论.

这样, 在边际和 $X_{i.}, X_{i,k}$ 及 $X_{.k}$ 固定的条件下, 一个三维列联表可用网络的一条通路来表示.

对于结果 2, 其中 $X_{i.}, X_{.j}, X_{.k}$ 是固定的, 那么根据定理 1 分别可由 $X_{i.}$ 和 $X_{.j}$ 得到 $X_{ij}^{(\alpha)} (1 \leq a \leq A)$, 由 $X_{i.}$ 和 $X_{.k}$ 得到 $X_{ik}^{(\alpha)} (1 \leq b \leq B)$ 以及由 $X_{.j}$ 和 $X_{.k}$ 得到 $X_{jk}^{(\alpha)} (1 \leq c \leq C)$, 其中 a, b, c 均为正整数. 适当组合 $X_{ij}^{(\alpha)}, X_{ik}^{(\alpha)}, X_{jk}^{(\alpha)}$ 就可以得到满足 $X_{i.}, X_{.j}, X_{.k}$ 固定条件的所有可能的列联表. 每一种组合对应有一个网络, 而每个网络中每条通路即表示一个可能的列联表.

对于结果 3, 类似地也可得到若干个网络, 每个网络中每条通路表示一个可能的列联表.

§ 4. 假设检验中概率值的计算

根据 §2 所述结果 1, 在二维列联表中, 在对二变量独立的原假设作检验时, 为计算概率值 P , 选择 $\{X_{ij}^{(\alpha)}\}$ 的条件为

$$\prod_{ij} \frac{X_{ij}^{(\alpha)}}{X_{ij}^{(\alpha)}} \geq 1. \quad (38)$$

于是当

$$\prod_i \prod_j X_{ij}^{(\alpha)} \geq \prod_{ij} X_{ij}^{(\alpha)} (1 \leq i^* < I, 1 \leq j^* \leq J)$$

成立时, 显然有 $\prod_{ij} \frac{X_{ij}^{(\alpha)}}{X_{ij}^{(\alpha)}} \geq 1$. 可停止向前搜索, 与相应的节点 $ND_{i^*j^*}^{(\alpha)}$ 连接的道

路都不必再考虑。这样由最后余下的各道路相应的表可算得 $P = \sum f^*(\{X'_{ij}\})$, 其中

$$\prod_{ij} (X'_{ij1}/X^*_{ij1}) < 1. \text{ 由此可得}$$

$$P = 1 - P'. \tag{39}$$

同样,对于 §2 所述结果 2 和 3,在三维列联表中,对原假设作检验时,为计算概率值 P ,选择 $\{X'_{ijk}\}$ 的条件为

$$\prod_{ijk} \frac{X'_{ijk1}}{X^*_{ijk1}} \geq 1. \tag{40}$$

于是也只要计算 $P = \sum f^*(\{X'_{ijk}\})$, 其中 $\prod_{ijk} (X'_{ijk1}/X^*_{ijk1}) < 1$. 由此即得

$$P = 1 - P'. \tag{41}$$

§ 5. 讨 论

为便于讨论,先作以下定义。

定义 3. 在产生二维列联表的网络中,对于节点 $ND_{ij}^{(0,ij)}$ 有

$$TN_{ij}^{(0,ij)} = \begin{cases} \sum TN_{i,j+1}^{(0,i,j+1)*} & (1 \leq i \leq I-1, 0 \leq j \leq J-1), \\ \sum TN_{i+1,i}^{(0,i+1,i)*} & (1 \leq i \leq I-2, j = J), \end{cases}$$

其中求和是对所有符合 $[\theta_{i,j+1}]^* = n_{i,j+1}$ 的 $ND_{i,j+1}^{(0,i,j+1)*}$ 进行的。

可见 $TN_{ij}^{(0,ij)}$ 即为网络中与节点 $ND_{ij}^{(0,ij)}$ 相连接的各道路上 $ND_{ij}^{(0,ij)}$ 以后的所有节点数。

根据 §4 中停止向前搜索的原则,若对于某节点 $ND_{i^*j^*}^{(0,i^*j^*)}$ 有

$$\prod_i \prod_j X'_{ij1} \geq \prod_{ij} X^*_{ij1},$$

则不必求出由 $ND_{i^*j^*}^{(0,i^*j^*)}$ 出发的新节点及相应的各种数值。于是可少做 $TN_{i^*j^*}^{(0,i^*j^*)}$ 轮运算(求出节点,相应的 $V_{ij}^{(0)}$ 值和 $X_{ij}^{(0)}$ 值以及 ΠX_{ij1} 值)。

现以一计算实例加以说明。现观察到一、二维列联表 $\{X_{ij}^{(0)}\}$

2	2	0
1	1	1
1	0	1

, 其边际和为

$X_{1.} = 4, X_{2.} = 3, X_{3.} = 2, X_{.1} = 4, X_{.2} = 3, X_{.3} = 2$. 在行、列边际和固定条件下,欲检验原假设:二变量相互独立。根据前述方法可得网络,列于图 1。

图中 * 表示搜索于该节点停止,---表示不再搜索的道路,连线上数字表示列联表相应格子中的数值。

根据此网络,所有可能的列联表为 39 个。而因为 $\prod_i \prod_j X^*_{ij1} = 4$, 所以根据停止

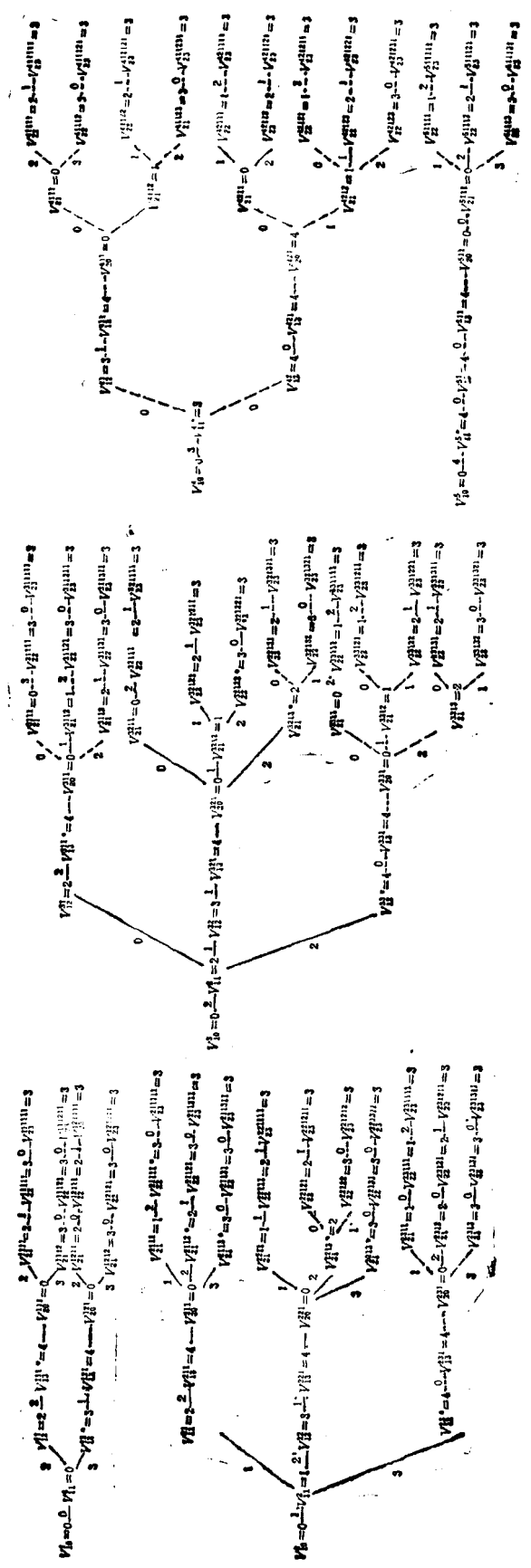


图 1 产生边和为 $X_1 = 4, X_2 = 3, X_3 = 2, X_4 = 4, X_5 = 3, X_6 = 2$ 的所有二维列联表 (3×3) 的网络

搜索原则,为作确切概率检验,实际上只需求出一张表的确切概率,求其它各表的过程都先后在中途停止。该表的确切概率值为 $\frac{4! 3! 2! 4! 3! 2!}{2! 9!} = 0.114285714$ 。于是,对原假设作检验,算得的概率值为 $P = 1 - 0.114285714 = 0.885714286$ 。据此,接受原假设,二变量相互独立。

至于三维列联表的情况可作类似讨论,在此不作赘述。

本文承史秉璋教授审阅指正,谨表谢意。

参 考 文 献

- [1] C. R. Rao, *Linear Statistical Inferences and Its Applications*, John Wiley and Sons, New York, 1973.
- [2] V. M. M. Bishop, S. E. Fienberg, P. M. Holland, *Discrete Multivariate Analysis, Theory and Practice*. MIT Press, Cambridge, Mass. 1975.
- [3] S. E. Fienberg, *The Analysis of Cross-Classified Categorical Data*. 2nd ed. MIT Press, Cambridge Mass, 1980.