

问题回复

这个问题可以归结为两个多项分布概率是否相等的检验问题或者列联表的齐性检验

假设

疾病组的数据 n_1, n_2, \dots, n_r 服从多项分布 $\mathcal{M}(n, p)$

对照组的数据 k_1, k_2, \dots, k_r 服从多项分布 $\mathcal{M}(k, q)$

其中:

$$n_1 + n_2 + \dots + n_r = n, \quad k_1 + k_2 + \dots + k_r = k$$

$$p = (p_1, p_2, \dots, p_r), \quad q = (q_1, q_2, \dots, q_r)$$

$$p_1 + \dots + p_r = 1, \quad q_1 + \dots + q_r = 1$$

$$p_j, q_j \geq 0, j = 1, \dots, r$$

(在您的论文 $r = 8$)

论文的问题归结为: 利用上述数据进行假设检验

$$H_0: p = q \leftrightarrow H_1: p \neq q$$

检验方法为: 令

$$v_j = \frac{n_j + k_j}{n + k}, j = 1, \dots, r$$

$$X^2 = (n + k) \left\{ \frac{(n_1 - nv_1)^2}{nv_1} + \frac{(n_2 - nv_2)^2}{nv_2} + \dots + \frac{(n_r - nv_r)^2}{nv_r} \right. \\ \left. + \frac{(k_1 - kv_1)^2}{kv_1} + \frac{(k_2 - kv_2)^2}{kv_2} + \dots + \frac{(k_r - kv_r)^2}{kv_r} \right\}$$

则当样本量比较大时 X^2 近似服从自由度为 $r - 1$ 的卡方分布。

给定检验水平 α ，查找到卡方分布(下)分位点 $X_{1-\alpha}^2(r-1)$ ，当 $X^2 > X_{1-\alpha}^2(r-1)$ 时否定原假设。

参考资料：

杨振海，《拟合优度检验》，安徽教育出版社，1993（P71，推论 3.1）

注：1. 书中问题为列联表的齐性检验

2. 估计有些英文书也有

差”三级，而年龄分为儿童、中青年和老年三组，抽查结果汇总
下面数据表：

效果 \ 年龄	儿童	中青年	老年	效果总计
显著	58	38	32	128
一般	28	44	45	117
较差	23	18	14	55
年龄总计	109	100	91	300

由上表及公式(3.14)计算得

$$\chi^2 = \sum \sum (Z_{ij} - \frac{z_{ij}z_{i.}}{300})^2 \cdot \frac{300}{z_{i.}z_{.j}} = 13.9$$

取显著水平 $\alpha=0.05$ ，查表得 $\chi^2_{4}(0.05)=9.49 < 13.9$ ，即有充分证据说明疗效与年龄有关，事实上 $13.9 > \chi^2_{4}(0.01)=12.277$ ，显著水平 < 0.01 。

(2) 二维列联表：齐性检验

另一类重要问题是齐性检验，数据形式和独立性检验相同，但统计概念是完全不同的。所谓齐性检验是检验若干总体是否有相同的分布。这一问题描述如下：设有 c 个离散总体，每个总体均取 r 个值，不失一般性这些值是 $1, 2, \dots, r$ ，第 i 个总体的概率函数是 P_i ：

$$P_i = (p_{i1}, \dots, p_{ir})^T, \quad \sum_{j=1}^r p_{ij} = 1, \quad p_{ij} > 0, \quad i=1, \dots, c$$

今将第 i 个总体记为 W_i ，则

$$P(W_i = j) = p_{ij} \quad 1 \leq j \leq r$$

对 W_i 观察 n_i 次， W_i 取 j 值的次数记为 $X_{ij}, j=1, \dots, r$ ，则 $\sum_{j=1}^r X_{ij} = n_i$ 。将 c 个总体的独立观察结果可列成表3.3：

$$n = \sum_{i=1}^c n_i, \quad X_{.j} = \sum_{i=1}^c X_{ij}$$

表3.3

总体	值	1	2	...	r	样本量
1		X_{11}	X_{12}	...	X_{1r}	n_1
2		X_{21}	X_{22}	...	X_{2r}	n_2
...		\vdots	\vdots		\vdots	
C		X_{c1}	X_{c2}	...	X_{cr}	n_r
和		$X_{\cdot 1}$	$X_{\cdot 2}$		$X_{\cdot r}$	n

上表的形式和表3.1相同，但统计概念是不相同的。两者差别在于对齐性检验数据，各行数据观察彼此是独立的，可人为控制观察第几个总体，而对独立性检验数据，没有这个性质。观察一个对象要同时记录两个属性（例如3.2例，同时观察疗效和年龄段），简言之， $n_i = \sum_{j=1}^r x_{ij}$ 不是随机的，而对独立性检验，该量是随机的。因此对齐性检验不能套用定理3.2。对独立性检验，不能去掉任何一行数据，去掉了数据就不完全，就成了基于不完全数据的独立性检验，是完全不同的统计问题。而对齐性检验，去掉任何一行均不影响问题的性质和处理方法，仅是用 $c-1$ 代替 c 罢了。

齐性检验就是根据表2的数据，检验各总体是否有相同分布。精确提法如下：令

$$X = (X_{11}, X_{12}, \dots, X_{1r}, X_{21}, \dots, X_{2r}, \dots, x_{c1}, \dots, x_{c1})^T$$

则

$$P(X=x) = P(x, \theta) = \prod_{i=1}^c (n_i! \prod_{j=1}^r \frac{1}{x_{ij}!} p_{ij}^{x_{ij}})$$

其中：

$$\theta = (p_{11}, p_{12}, \dots, p_{1r}, \dots, p_{c1}, \dots, p_{cr})^T = (P_1^T, P_2^T, \dots, P_c^T)^T$$

那么齐性检验问题是参数的检验问题：

$$H_0: P_i = P_0 \quad i=1, \dots, c; \quad H_1: P_i \quad 1 \leq i \leq c \text{ 不全相等}$$

其中: $\mathbf{P}_i = (p_{i1}, p_{i2}, \dots, p_{ir})^T, 0 \leq i \leq c.$

今用似然比检验解决这一问题。对给定 $\mathbf{X} = \mathbf{x}$ 时, 对数似然函数为

$$l(\theta) = \text{constant} + \sum_{i=1}^c \sum_{j=1}^r x_{ij} \log p_{ij} \quad (3.15)$$

且满足约条件

$$\sum_{j=1}^r p_{ij} = 1, i=1, \dots, c$$

用乘子法可求 p_{ij} 的 MLE 为

$$\hat{p}_{ij} = x_{ij}/n_i \quad i=1, \dots, c, j=1, \dots, r$$

当 $H_0: \mathbf{P}_i = \mathbf{P} \quad 1 \leq i \leq c$ 成立时, 对数似然函数为

$$l_0(\theta) = \text{constant} + \sum_{i=1}^c \sum_{j=1}^r x_{ij} \log p_i$$

且满足约束条件

$$\sum_{j=1}^r p_i = 1$$

于是求得 p_{ij} 的 MLE 为

$$\hat{p}_{ij} = x_{ij}/n, j=1, \dots, r, n = \sum_{i=1}^c n_i \quad (3.16)$$

因此, 似然比统计量为

$$\begin{aligned} LR &= 2l(\hat{\theta}_n) - 2l_0(\hat{\theta}_0) \\ &= 2 \left\{ \sum_{i=1}^c \sum_{j=1}^r x_{ij} \log \left(\frac{x_{ij}}{n_i} \right) - \sum_{i=1}^c \sum_{j=1}^r x_{ij} \log \frac{x_{\cdot j}}{n} \right\} \end{aligned} \quad (3.17)$$

其中 $\hat{\theta}_n$ 是 θ 的 MLE, $\hat{\theta}_0$ 是在 H_0 成立时 θ 的 MLE. 当 LR 过大时拒绝 H_0 . 现给出 (3.17) 式的等价形式, 并 H_0 成立的条件下, 求其极限分布.

令 $\mathbf{P} = \sum_{i=1}^c \frac{n_i}{n} \mathbf{P}_i = (p_1, \dots, p_r), H_0$ 成立意味着 $\mathbf{P} = \mathbf{P}_0$

但反之不对.

由大数定理

$$\frac{x_{ij}}{n_i} \rightarrow p_{ij} \quad a, s, n_i \rightarrow \infty \quad i=1, \dots, c, \quad j=1, \dots, r \quad (3.18)$$

$$\frac{x_{.j}}{n} \rightarrow p_j \quad a, s, n \rightarrow \infty \quad j=1, \dots, r$$

$$\begin{aligned} LR &= 2 \left\{ \sum_{j=1}^r \sum_{i=1}^c x_{ij} \log \frac{x_{ij}}{n_i} - \sum_{i=1}^r x_{.i} \log \frac{x_{.i}}{n} \right\} \\ &= 2 \sum_{j=1}^r \sum_{i=1}^c x_{ij} \log \frac{x_{ij}}{n_i p_j} + 2 \sum_{i=1}^r \sum_{j=1}^c x_{ij} \log p_j \\ &\quad - 2 \sum_{j=1}^r x_{.j} \log \frac{x_{.j}}{n} \\ &= 2 \sum_{j=1}^r \sum_{i=1}^c x_{ij} \log \frac{x_{ij}}{n_i p_j} - 2 \sum_{j=1}^r x_{.j} \log \frac{x_{.j}}{n p_j} \end{aligned}$$

由(3.18)和Taylor公式, 若 $n_i, 1 \leq i \leq c$ 同时 $\rightarrow \infty$ 时,

$$2 \sum_{j=1}^r \sum_{i=1}^c x_{ij} \log \frac{x_{ij}}{n_i p_j} = 2 \sum_{j=1}^r \sum_{i=1}^c x_{ij} \log(1 +$$

$$\frac{x_{ij} - n_i p_j}{n_i p_j})$$

$$= 2 \sum_{j=1}^r \sum_{i=1}^c x_{ij} \left\{ \left(\frac{x_{ij} - n_i p_j}{n_i p_j} \right) - \frac{1}{2} \frac{(x_{ij} - n_i p_j)^2}{(n_i p_j)^2} \right.$$

$$\left. + O_p \left(\left(\frac{x_{ij} - n_i p_j}{n_i p_j} \right)^3 \right) \right\}$$

$$= 2 \sum_{j=1}^r \sum_{i=1}^c \frac{(x_{ij} - n_i p_j)^2}{n_i p_j} - \sum_{j=1}^r \sum_{i=1}^c \frac{(x_{ij} - n_i p_j)^2}{n_i p_j}$$

$$\cdot \left(1 + \frac{x_{ij} - n_i p_j}{n_i p_j} \right)$$

$$= \sum_{j=1}^r \sum_{i=1}^c (x_{ij} - n_i p_j)^2 / n_i p_j$$

同理可证

$$2 \sum_{j=1}^r x_{.j} \log \frac{x_{.j}}{n} = \sum_{j=1}^r (x_{.j} - n p_j)^2 / n p_j$$

于是当 $n_i, 1 \leq i \leq c, n \rightarrow \infty$ 时

$$LR = \sum_{i=1}^c \sum_{j=1}^r (x_{ij} - n_i p_j)^2 / n_i p_j - \sum_{j=1}^r (x_{.j} - n p_j)^2 / n p_j \quad (3.19)$$

关于LR的极限分布有以下定理.

定理3.4 若

(1) $H_0: P_2 = P \quad 1 \leq i \leq c, P = (p_1, \dots, p_r)^T$ 成立,

(2) $\lim n_i/n = \alpha_i, 1 \leq i \leq c.$

则LR的极限分布为 $\chi^2_{(r-1)} (c=1).$

证明 令

$$Y_{ij} = \sqrt{n_i} \left(\frac{X_{ij}}{n_i} - p_i \right) / \sqrt{p_i}$$

$$Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{ir})^T \quad 1 \leq i \leq c$$

$$Y = (Y_1, Y_2, \dots, Y_c)^T$$

则(3.19)式可写作

$$LR = \sum_{i=1}^c Y_i^T Y_i - \left(\sum_{i=1}^c \sqrt{\frac{n_i}{n}} Y_i \right)^T \left(\sum_{i=1}^c \sqrt{\frac{n_i}{n}} Y_i \right) \quad (3.20)$$

由引理3.1和假设条件

$$Y_i \xrightarrow{L} N_r(\theta, I_r - \sqrt{P} \sqrt{P}^T) \quad 1 \leq i \leq c$$

LR的极限分布与变量 $R = \sum_{i=1}^c Z_i^T Z_i - \left(\sum_{i=1}^c \sqrt{\alpha_i} Z_i \right)^T$

$\cdot \left(\sum \sqrt{\alpha_i} Z_i \right)$ 的分布相同, 其中: Z_1, \dots, Z_c 相互独立且.

$$Z_i \sim N(\theta, I_r - \sqrt{P} \sqrt{P}^T)$$

取 c 阶正交阵 P^1 , 其第一列为 $(\sqrt{\alpha_1}, \sqrt{\alpha_2}, \dots, \sqrt{\alpha_c}).$

令

$$V = (V_1, V_2, V_c)^T = PZ \equiv P(Z_1, \dots, Z_c)^T$$

则 V_1, \dots, V_c 相互独立, $\text{Var } V_i = I - P^1 \sqrt{P} \sqrt{P} P \equiv \Sigma, \Sigma^2 = \Sigma.$

则

$$\begin{aligned} & \sum_{i=1}^c Z_i^T Z_i - \left(\sum_{i=1}^c \sqrt{\alpha_i} Z_i \right)^T \left(\sum_{i=1}^c \sqrt{\alpha_i} Z_i \right) \\ &= \sum_{i=1}^c V_i^T V_i \end{aligned}$$

由引理3.2

$$V_i^T V_i \sim \chi^2_{r-1}$$

(3.20)式其 V_i 的相互独立性知 $\sum_{i=2}^c V_i^T V_i \sim \chi^2_{(c-1)(r-1)}$
定理得证.

推论3.1 若定理3.4的假设条件成立, 则

$$X^2 = n \sum_{i=1}^c \sum_{j=1}^r (X_{ij} - n_i X_{.j})^2 / n_i X_{.j} \quad (3.21)$$

的极限分布是 $\chi^2_{(r-1)(c-1)}$.

证明

$$\begin{aligned} X^2 &= n \sum_{i=1}^c \sum_{j=1}^r \left\{ (X_{ij} - np_i)^2 + \left(np_i - \frac{n_i X_{.j}}{n} \right)^2 + \right. \\ &\quad \left. + 2(X_{ij} - np_i) \left(np_i - \frac{n_i X_{.j}}{n} \right) \right\} / n_i X_{.j} \\ &= n \sum_{i=1}^c \sum_{j=1}^r \left(np_i - \frac{n_i X_{.j}}{n} \right)^2 / n_i X_{.j} \\ &= n \sum_{j=1}^r \sum_{i=1}^c n_i \left(p_i - \frac{X_{.j}}{n} \right)^2 / X_{.j} = \sum_{j=1}^r \frac{1}{X_{.j}} \cdot (X_{.j} - np_i)^2 \\ &= \sum_{j=1}^r (X_{.j} - np_i)^2 / np_i \quad (\text{因为 } \frac{X_{.j}}{n} \rightarrow p_i, a, s). \\ &\quad + 2n \sum_{i=1}^c \sum_{j=1}^r (X_{ij} - np_i) \left(np_i - \frac{n_i X_{.j}}{n} \right) / n_i X_{.j} \\ &= 2n \sum_{i=1}^c \sum_{j=1}^r (X_{ij} - np_i) \left(p_i - \frac{X_{.j}}{n} \right) / X_{.j} \\ &= 2n \sum_{j=1}^r \left(p_i - \frac{X_{.j}}{n} \right) (X_{.j} - np_i) / X_{.j} \\ &= -2 \sum_{j=1}^r (X_{.j} - np_i)^2 / X_{.j} \\ &= -2 \sum_{j=1}^r (X_{.j} - np_i)^2 / np_i \end{aligned}$$

因此有

$$\begin{aligned} X^2 &= \sum_{i=1}^c \sum_{j=1}^r n (X_{ij} - np_i)^2 / n_i X_{ij} - \\ &\quad \sum_{j=1}^r (X_{.j} - np_i)^2 / np_i \end{aligned}$$