# 理论新探

2004年第7期(总第175期)

## ■蒋志华 徐琼梅

# 沃纳模型

### 一、沃纳模型的基本思想

沃纳模型是沃纳(Warner)于 1965 年首先提出的一种随机化回答模型。该 模型是为了解决社会经济现象中敏感性 问题而采用的一种随机化回答技术。模 型设计的基本思想是:为了调查某个敏 感性问题,同时列出两个存在相关关系 的问题制成卡片(一个问题是具有某种 特征为卡片 A,另一个问题是不具有某 种特征为卡片 B),被调查者随机抽取卡 片(不告诉调查人员)进行回答,回答后, 卡片不退还调查人员。这种设计,使调查 人员不知道被调查者在回答哪一个问 题,在一定程度上消除了被调查者的担 心和顾虑,使他们参与调查,从而实现了 被调查者提供真实情况的自我保护。同 时,调查人员还可以通过对所有调查结 果的汇总,利用概率原理推算得到总体 中具有该特征人数比例的估计值。其模 型公式为:

 $\pi_{\alpha}$ 的估计值为: $\hat{\pi}_{\alpha}$ =1/(2P-1)n<sub>1</sub>/n-(1-p)/(2p-1)=[n<sub>1</sub>/n-(1-p)]/(2p-1) (p  $\neq$  1/2)

 $\pi_{\alpha}$ 的方差估计量为  $:V(\hat{\pi}_{\alpha})=\hat{\pi}_{\alpha}(1-\hat{\pi}_{\alpha})$  /n+p(1-p)/n(2p-1)²

 $\pi_{\alpha}$  的置信区间为  $\hat{\pi}_{\alpha}\pm t\sqrt[4]{V(\hat{\pi}_{\alpha})}$ 上式中  $\pi_{\alpha}$  为具有卡片 A 特征的人数占卡片总数( 样本总数 )的比重;

p 为 A 卡片数所占卡片总数(样本总数)的比重;

n 为卡片总数( 样本总数 );

 $n_1$  为所有抽中卡片 A 和卡片 B 并回答" 是"的人数。

二、沃纳模型应用的典型案例

如果要调查学生在考试中是否有过 作弊行为,考虑到有这种行为的学生不 愿承认,因此可以采用沃纳模型进行实 验如下:

第一步,设计卡片 A"我在考试中有过作弊行为"和卡片 B"我在考试中没有作弊行为"。当然,在制作卡片时,调查者对两种卡片的比例是知道的,卡片 A占的比例为 p( 注意 p 不能为 1/2 ) ,卡

片 B 占的比例为(1-p)。设具有卡片 A 特征人数(或单位)的比例为  $\pi_\alpha$ ,样本容量为 n,其中回答"是"的人数为  $n_1$ ( $n_1$  既包括抽中卡片 A 回答"是",也包括抽中卡片 B 回答"是")。

第二步,由被调查者随机抽取一张 卡片上的问题回答,并统计出回答"是"和回答"否"的人数。

第三步,运用沃纳模型公式测算出 总体中具有某特征人数比例的估计值。

现以统计系 1 班和管理系 2 班和计算机系 3 班的学生进行实验,其实验结果如下:

统计系 1 班的设计和实验结果为: P=1/4 n=41  $n_1=23$ ; 管理系 2 班的设计和实验结果为 P=1/4 n=23  $n_1=18$  ,计算机系 3 班的设计和实验结果为 :P=9/20 , n=20  $n_1=8$ 。试以 95%的把握程度对上述各班学生在考试中有过作弊行为人数的比例进行点估计和置信区间估计。

由把握程度为 95% ,所以 t=1.96 ,将 各班调查结果代入沃纳模型公式可得:

统计系 1 班的 
$$\hat{\pi}_{\alpha} = \frac{1}{(2 \times 1/4)-1} \times$$

 $\frac{23}{41} - \frac{1 - (1/4)}{(2 \times 1/4) - 1} = 0.3781 = 37.81\%$ 

即该班学生在考试中有过作弊行为 人数比例的点估计值为 37.81%。

 $V (\hat{\pi}_{\alpha}) = 0.2781 (1 - 0.3781)/41 + \frac{1/4(1-1/4)}{41(2\times1/4-1)^2}0.024$ 

 $\hat{\pi}_{\alpha} \pm t \sqrt{V(\hat{\pi}_{\alpha})} = 0.3781 \pm 1.96 \sqrt{0.024}$ =0.3781±0.3036

则有:总体 $\hat{\pi}_{\alpha}$ 的置信区间为:  $0.0745 \leqslant \hat{\pi}_{\alpha} \leqslant 0.6817$ ,即有 95%的把握程度推断,该系学生在考试中有过作弊行为人数比例在 7.45%至 68.17%之间。同理可得:管理系 2 班的 $\hat{\pi}_{\alpha}$ =-0.0652=-6.52%;即该班学生在考试中有过作弊行为人数比例的点估计值为-6.52%,这一计算结果没有意义。

而  $V(\hat{\pi}_{\alpha})=-0.0356$  ,总体 $\hat{\pi}_{\alpha}$ 的置信 区间没有办法进行估计。 计算机系 3 班的 $\hat{\pi}_{\alpha}$ =1.5=150%;即该系学生在考试中有过作弊行为人数比例的点估计值为 150%。这一计算结果没有意义。

而  $V(\hat{\pi}_{\alpha})$ =1.2, $\hat{\pi}_{\alpha}$ ±t  $\sqrt{V(\hat{\pi}_{\alpha})}$ =1.5±2.1471,即总体 $\hat{\pi}_{\alpha}$ 的置信区间为:0.647 $\leq$  $\hat{\pi}_{\alpha}$  $\leq$ 13.6471,即有95%的把握程度推断,该系学生在考试中有过作弊行为人数比例在6.47%至364.71%之间。

### 三、沃纳模型的三个缺陷

从上述三个典型案例可以发现,沃 纳模型虽然解决了在一定程度上消除被 调查者的担心和顾虑,使敏感性问题有 办法通过随机抽取卡片的形式得到较为 真实的调查结果。但是,该模型存在三个 明显的缺陷:一是在设计的两个问题时, 这两个问题存在相关关系,如卡片 A 我 有过漏税行为和卡片 B 我没有漏税行为 是相互关联的两个问题, 使被调查者仍 然可能有怀疑而不予合作。因此,沃纳模 型没有完全消除被调查者的担心和顾 虑,使调查成功的可能性不大;二是在沃 纳模型公式中,要求 p 不能等于 1/2,否 则沃纳模型就无法使用。但是从消除被 调查者顾虑的角度考虑 ,应使 p 等于 1/ 2,才能保证两种卡片的抽中的机会均 等。三是即使满足沃纳模型公式中 p 不 能等于 1/2 的条件, 也会得出没有实际 意义得估计值。为了避免前两个缺陷,西 蒙斯(Simmons)在沃纳模型的基础上提 出了一种修改的方法,即西蒙斯模型。但 是,西蒙斯模型实际上仍未解决p值的 取值问题(本文不做讨论)。本文试图就 沃纳模型的基本原理和公式推导过程提 出产生沃纳模型缺陷的原因及其改进方 法。

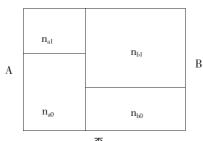
### 四、沃纳模型的改进方法

改进方法之一:根据沃纳模型的基本原理及其模型公式推导提出的限制使 用条件法。

若将沃纳模型的基本原理及其模型 公式推导过程用下图来说明,则有:

# 的缺陷及其改進

是



呇

上图中 ,n 为卡片总数 ( 样本总数 ) ,即  $n=n_{al}+n_{a0}+n_{bl}+n_{l0}$  ;

 $n_{al}$  为抽中卡片 A 并回答" 是 "的人数;

 $n_{a0}$  为抽中卡片 A 并回答 否 的人数;  $n_{b1}$  为抽中卡片 B 并回答 是 的人数;  $n_{b0}$  为抽中卡片 B 并回答 否 的人数;  $n_1$  为所有抽中卡片 A 和卡片 B 并回答 是 "的人数,即  $n_1+n_{a1}+n_{b1}$ 

p 为 A 卡片数所占卡片总数 ( 样本总数 )的比重 ,即  $p=(n_a|+n_a)/n$ 

 $(\hat{\pi}_{\alpha})$ 为具有卡片 A 特征的人数占卡片 总数( 样本总数 )的比重 ,即 $\hat{\pi}_{\alpha}$ = $(n_{ai}+n_{bo})/n$  根据概率论中的全概率定律可得到以下公式:

将公式 1 中的  $n_{al}/(n_{al}+n_{s0})$  视为 $\hat{\pi}_{\alpha}$ , 并将  $n_{bl}/(n_{bl}+n_{b0})$ 视为 $(1-\hat{\pi}_{\alpha})$ 可得  $n_{l}/n=(n_{al}+n_{bl})/n=p\hat{\pi}_{\alpha}+(1-p)(1-\hat{\pi}_{\alpha})$ ,即:

$$\begin{split} \hat{\pi}_{\alpha} = &1/(2p-1)(n_1/n) - [(1-p)/(2p-1)] = \\ &[n_1/n - (1-p)]/(2p-1) \ (p \neq 1/2) \\ &\boxed{\text{可见 ,公式 2 就是沃纳模型公式。但}} \end{split}$$

是,这是建立在 $\hat{\pi}_{\alpha}$ = $n_{al}$ / $(n_{al}+n_{al})$ , $(1-\hat{\pi}_{\alpha})$ = $n_{bl}$ / $(n_{bl}+n_{bl})$ 的基础上才能成立的。实际上就是说, 沃纳模型公式成立的条件是:一是将样本总数(包括卡片 A 和卡片 B)中回答"是"的人数占卡片总数的比重,用卡片 A 中回答"是"的人数占卡片 B 的比重来代替;二是将卡片 B 中回答"是"的人数占卡片 A 中回答"是"的人数占卡片 A 的比重来用卡片 A

代替 ,即(  $1-\hat{\pi}_{\alpha}$  )= $1-n_{al}/(n_{al}+n_{a0})=n_{a0}/(n_{al}+n_{a0})$ 

 $n_{a0}$ )= $n_{b1}$ /( $n_{b1}$ + $n_{b0}$ )。显然 这种代替只有当 n 足够大及抽样充分随机化时才是合理的 ,当 n 很小时 ,分组得到的子样本与样本总体之间可能存在很大的差异 ,极大的抽样误差将导致这一代替过程失效。在上述典型案例中,管理系 2 班和计算机系 3 班的 n 都小于 30 ,才会出现运用沃纳模型计算结果没有意义的现象。因此,在运用沃纳模型进行敏感性问题的处理时,必须在满足上述两个条件的同时,还应考虑 p (1-p)值与  $n_i$ /n 之间的数量关系 .即:

当 p > 1/2 时 (2p-1) 应为正值。而  $n_1/n > (1-p)$  则 $\hat{\pi}_{\alpha}$  的估计值为正值 ,采用 沃纳模型有意义,但如果  $n_1/n$  与(1-p) 相差很小,也会使结果大于 100%而没有 意义;如  $n_1/n < (1-p)$ ,则 $\hat{\pi}_{\alpha}$  的估计值为负值 采用沃纳模型没有意义。

当 p < 1/2 时 (2p-1) 应为负值。而  $n_1/n > (1-p)$  则 $\hat{\pi}_\alpha$  的估计值为负值,采用 沃纳模型没有意义(如管理系 2 班),如  $n_1/n < (1-p)$  则 $\hat{\pi}_\alpha$  的估计值为正值,采用 沃纳模型有意义( 如统计系 1 班),但如果  $n_1/n$  与( 1-p )相差很小,也会使结果大于 100%而没有意义( 如计算机系 3 班 )。

改进方法之二:通过访问过程的设计来得到 $\hat{\pi}_a$ 的估计值法。

上述方法是在肯定了沃纳模型的合理性的前提下提出的改进意见。但在实际应用中,利用沃纳模型很容易得到一些不合理的fn。估计值,如负值、等于可值、大于 1 的值。上述通过对沃纳模型股的改进方法,虽然可以避免不合理的出现,但即使得到了一个有效的值(在 0 和 1 之间),仍然不能确定这一值的可信性。我们不妨跳出建立沃纳模型公的思维圈,从访问过程的设计来考察其估计值的计算。

若用沃纳模型的基本思想来估计计 算机系 3 班学生作弊的比例如下:

从上述运用沃纳模型公式计算的计 算机系 3 班学生作弊比例的估计值为 150%,这显然又与事实不符。我们不妨 对这8个"是"的答案进行两个极端的假 设性分配:

第一种情况:若 n<sub>al</sub>=8 ,则 n<sub>a0</sub>=1 ,n<sub>bl</sub>= 0 ,n<sub>bo</sub>=11;

第二种情况 :若  $n_{bl}$ =8 ,则  $n_{a0}$ =9 , $n_{al}$ = 0 , $n_{bo}$ =3。

那么  $\hat{\pi}_{\alpha}$  的估计值公式为:

 $\hat{\pi}_{\alpha} = (n_{al} + n_{b0})/n \tag{3}$ 

则有:第一种情况的 $\hat{\pi}_{\alpha}=19/20=$ 95%,第二种情况的 $\hat{\pi}_{\alpha}=3/20=15\%$ 。

以上只是作了两个极端的假设 ,类似的假设可以做出很多种 , 从上述假设我们可以看出由公式 3 计算的 $\hat{\pi}_{\alpha}$  的估计值无论如何不会大于 1 , 也不会为负值和 0。但是 ,在公式 3 中  $n_{al}$  和  $n_{lo}$  得数值又怎样获得呢?

我们仍可设计两类卡片 A 和 B ,卡 片 A 写的敏感性问题,卡片 B 既可以是 与A相反相关的问题也可以是与A无 关的非敏感性问题(最好选择后者,参见 西蒙斯模型),卡片的设计者、发放者和 收回者是不同的三个人,卡片的外观设 计成完全相同的样式并密封,调查前充 分的说明本次调查的组织方式,调查后 当面统计结果。调查者回答时不署名,且 问题与答案都收回,再通过简单的对照 统计就可以得到被调查问题中抽中卡片 A 并回答"是"的人数 nal ,同样可得到抽 中卡片 B 并回答" 否"的人数 nia。可见, 这一方法仍然采用的随机化的技术,调 查者虽然从回收的卡片一眼可以看出被 调查者在调查主题上的特征,但调查者 并不知道这一卡片是谁提交的,因此还 是无从知道每个调查者的具体特征,而 只知道所有被调查者具有某一特征的比 例。当然,这种方法必须以被调查者对方 法本身的理解为前提,故在调查前必须 详细的说明整个调查的设计组织过程, 并让被调查者保证回答自己的真实情 况,若有人对调查的保密性有怀疑则不 能参加调查.

> (作者单位/成都信息工程学院) (责任编辑/亦 民)