

敏感性问题调查的一种定量研究方法——谈随机乘法模型

李红梅

(北方工业大学理学院, 北京, 100041)

摘要: 在敏感性问题的市场调查中, 人们对随机化回答技术的要求越来越高, 本文在乘法模型的基础上提出一改进模型——随机乘法模型, 不仅使估计量的精度有一定的保证, 而且操作简便, 估计量形式简单, 同时又对被调查者有一定的保护, 该模型在敏感性问题数量特征的调查中, 是一种比较理想的方法。

关键词: 敏感性问题, 随机化回答技术, 乘法模型, 随机乘法模型

1. 引言

随着社会经济的发展, 人们对信息范围与质量的要求逐步提高, 在市场调查中, 涉及到的敏感性问题也越来越多。在敏感性问题的调查中, 如果采用一般的直接询问方式进行调查, 是根本无法获得真实数据。拒绝回答和故意给出不真实信息是产生样本估计偏差的非抽样误差的两种主要来源。为了克服这种搜集数据的困难, 1965 年沃纳推出沃纳模型, 从而开创了随机化回答技术的先河。之后的 40 年, 随机化回答技术日臻完善, 已经广泛地应用于各种敏感性问题的调查研究中, 该领域的研究学者不断推陈出新, 产生了一系列的模型。然而, 在这众多的随机化回答模型中, 并非都是十分完善的, 也存在优劣之分。综合来看, 一个优良的随机化回答模型, 应满足以下几点要求:

- 1、随机化装置操作简便
- 2、估计量形式简单
- 3、估计精度有一定保证
- 4、能够最大限度地保护被调查者隐私

基于以上条件, 本文在介绍乘法模型的基础上, 提出一个改进模型, 可以称之为随机乘法模型, 是一种操作简单, 行之有效的办法, 具有一定的推广应用价值。

2. 乘法模型介绍

2.1 模型的设计及参数的估计

在乘法模型中, 要求回答者将其敏感特征数值(X)与一个来自已知分布的随机数值(Y)相乘。观测值 Z 可被表示为: $Z=XY$ (X、Y 相互独立), 故

其均值为: $\mu_z = \mu_x \mu_y$, 方差为: $\sigma_z^2 = \mu_y^2 \sigma_x^2 + \mu_x^2 \sigma_y^2 + \sigma_x^2 \sigma_y^2$ 。

于是 μ_x 的估计量为: $\hat{\mu}_{x0} = \hat{\mu}_z / \mu_y$ 。

在无放回简单随机抽样条件下, 估计量方差为:

$$Var(\hat{\mu}_{x0}) = \frac{\sigma_z^2}{n\mu_y^2} \times \frac{N-n}{N-1} = \frac{(1-f)S_z^2}{n\mu_y^2}.$$

因为 Z_1, Z_2, \dots, Z_n 为简单随机样本, 所以我们可取估计量:

$\hat{\mu}_z = \bar{z}, \hat{S}_z^2 = s_z^2$, 且容易证明 $\hat{\mu}_z, \hat{S}_z^2$ 分别为 μ_z, S_z^2 的无偏估计

以样本均值 \bar{z} 代替 $\hat{\mu}_z$ 得: $\hat{\mu}_{x0} = \bar{z} / \mu_y$

估计量方差的无偏估计为: $\hat{Var}(\hat{\mu}_{x0}) = \frac{(1-f)s_z^2}{n\mu_y^2}$

2.2 乘法模型评价

乘法模型通过引入另一无关变量 Y 使回答结果 Z 与敏感性问题答案 X 产生较大差异, 从而较大幅度地保护了被调查者的隐私。同时该模型随机化装置及估计量的形式都比较简单, 这是该模型的优点。但是, 一般情况下, 该模型精度较差, 这点通过与一般简单估计的比较可以发现。

当 Y=1 时, 乘法模型即为一般的简单估计, 其方差为: $Var(\hat{\mu}_x) = \frac{\sigma_x^2}{n} \times \frac{N-n}{N-1}$

乘法模型与一般简单估计方差之比为:

$$\frac{Var(\hat{\mu}_{x0})}{Var(\hat{\mu}_x)} = 1 + \frac{\mu_x^2 \sigma_y^2}{\mu_y^2 \sigma_x^2} + \frac{\sigma_y^2}{\mu_y^2} = 1 + \left(\frac{\mu_x^2}{\sigma_x^2} + 1\right) \left(\frac{\sigma_y}{\mu_y}\right)^2$$

显然, 乘法模型估计精度低于一般简单估计, 通常情况下, 乘法模型的抽样估计误差会因引入了无关变量而显著地增加。提高乘法模型估计精度的有效办法是降低随机变量 Y 的变异系数, 然而这又与容易引起被调查者怀疑产生矛盾。能否在不改变无关问题的前提下, 降低引入无关问题的影响程度呢? 下面介绍的随机乘法模型便很好地解决了这个问题。

3. 随机乘法模型介绍

3.1 模型的设计与参数的估计

随机乘法模型是在乘法模型的基础上提出的一个改进模型, 它的基本设计思路如下:

- 1) 通过计算机产生一均值为 μ_y 的随机数 Y, 可以适当界定 Y 值的范围。
- 2) 在一个容器内放置若干个外观完全相同的小球, 小球上面分别标有数字 1 或 0。其中标有数字 1 的小球所占的比例为 P。
- 3) 被调查者从容器内随机抓取一小球, 如果小球号码为 1, 则被调查者回答敏感性问题 X 与 μ_y 的乘积, 即 $X\mu_y$; 号码为 0 时, 要求被调查者回答敏感性问题 X 与随机产生的 Y 值的乘积, 即 XY。

抓取小球与产生随机数的过程是隐蔽的, 研究者只能看到被调查者给出的最终回答 Z。

抓取小球的过程可以用符合 0, 1 分布的随机变量 ε 来表示, 抽到 1 号小球, 表示为 $\varepsilon = 1$, 抽到 0 号小球, 表示为 $\varepsilon = 0$, $p(\varepsilon = 1) = P$, 并且 X, Y 与 ε 相互独立。

用模型表示, 则得: $Z = \varepsilon X \mu_y + (1 - \varepsilon) XY$

两边取期望, 由 X, Y 与 ε 的独立性及条件期望的概念, 得:

$$\begin{aligned}\mu_z &= P\mu_x\mu_y + (1-P)\mu_x\mu_y = \mu_x\mu_y \\ \sigma_z^2 &= E[\varepsilon X \mu_y + (1 - \varepsilon) XY]^2 - [E(\varepsilon X \mu_y + (1 - \varepsilon) XY)]^2 \\ &= P\mu_y^2 E(X^2) + (1-P)E(X^2)E(Y^2) - (\mu_x\mu_y)^2 \\ &= \mu_y^2 \sigma_x^2 + (1-P)(\sigma_x^2 \sigma_y^2 + \mu_x^2 \sigma_y^2)\end{aligned}$$

于是 μ_x 的估计量为: $\hat{\mu}_{x1} = \hat{\mu}_z / \mu_y$

以样本均值 \bar{z} 代替 $\hat{\mu}_z$ 得: $\hat{\mu}_{x1} = \bar{z} / \mu_y$,

在无放回简单随机抽样条件下, 估计量方差为:

$$Var(\hat{\mu}_{x1}) = \frac{\sigma_z^2}{n\mu_y^2} \times \frac{N-n}{N-1} = \frac{(1-f)S_z^2}{n\mu_y^2}$$

$$\text{估计量方差的无偏估计为: } \hat{Var}(\hat{\mu}_{x1}) = \frac{(1-f)s_z^2}{n\mu_y^2}$$

3.2 随机乘法模型评价

随机乘法模型吸收了乘法模型的优点, 模型随机化装置及估计量的形式都比较简单, 同时还能够较大幅度地保护被调查者的隐私。然而, 在估计精度方面, 是否有显著提高呢? 我们先进行以下比较:

同样, 当 $Y=1$ 时, 随机乘法模型即为一般的简单估计

随机乘法模型与一般简单估计方差之比为:

$$\frac{Var(\hat{\mu}_{x1})}{Var(\hat{\mu}_x)} = 1 + (1-P) \left(\frac{\mu_x^2 \sigma_y^2}{\mu_y^2 \sigma_x^2} + \frac{\sigma_y^2}{\mu_y^2} \right) = 1 + (1-P) \left(\frac{\mu_x^2}{\sigma_x^2} + 1 \right) \left(\frac{\sigma_y}{\mu_y} \right)^2$$

显然, 由于无关问题的引入, 随机乘法模型估计精度仍然会低于一般简单估计, 但通过与乘法模型方差的对比可以看出, 比值大于 1 的部分多出一个系数 $(1-P)$, 从而使随机乘法模型的估计精度比乘法模型有了显著的提高。提高随机乘法模型估计精度不仅可以通过降低随机变量 Y 的变异系数而且可以提高 P 的值来实现。可见, 随机乘法模型不仅使估计量的精度有一定的保证, 而且操作简便, 估计量形式简单, 同时又对被调查者有一定的保护, 所以该模型在敏感性问题数量特征的调查中, 是一种比较理想的方法。

4. 随机乘法模型的应用

为调查某地区高校教师隐性收入状况, 设计敏感性问题 X : “你平均每月的隐性收入数额大概是多少?” 用随机乘法模型进行估计。首先通过计算机编程, 设计 $\mu_y = 68$ 的随机数 Y , 为保证一定精度, 可以界定 Y 值的范围, 比如 Y 在 $[0, 136]$ 之间, 并确定标有数字 1 的小球所占的比例为 $P=0.7$ 。然后让被调查者从容器内随机抓取小球, 并按上述规则回答问题。

已知在该地区范围内采用无放回的简单随机抽样抽选了 1000 名高校教师, 调查结果如

$$\text{下: } \bar{z} = 53175, s_z^2 = 9.2965 \times 10^8. \text{ 所以, } \hat{\mu}_{x1} = \frac{\bar{z}}{\mu_y} = \frac{53175}{68} \approx 782 \text{ (元)},$$

$$\hat{Var}(\hat{\mu}_{x1}) = \frac{(1-f)s_z^2}{n\mu_y^2} \approx \frac{9.2965 \times 10^8}{1000 \times 68^2} \approx 201.0489, \quad \sqrt{\hat{Var}(\hat{\mu}_{x1})} \approx 14.1792 \text{ (元)}$$

取 $\alpha = 0.05$, 即 $F(t)=95\%$ 时, $t=1.96$, 可确定隐性收入的区间范围为:

$$[782 - 1.96 \times 14.1792, 782 + 1.96 \times 14.1792], \text{ 即 } [754, 810] \text{ 元}$$

即, 可以以 95% 的概率把握程度保证该地区高校教师的隐性收入平均每人每月在 $[754, 810]$ 元范围之内。