

杂 感 篇

如何让人讲真话

——谈敏感问题的抽样调查

边 重

民意调查即调查老百姓对某事物的态度。若从著名的盖洛甫民意测验所成立算起,它已有50年的历史了。这几年国内也有不少单位以及报刊杂志搞了民意调查,其内容五花八门:婚姻恋爱态度、人生观、消费倾向、交通习惯甚至对国家大政方针的态度等等。但是总的来讲,这些调查的设计者都把注意力集中在如何设计问卷上,而很少或根本没有考虑如何正确地选择有代表性的调查对象,如何更精确地计算调查结果,如何降低各种干扰所产生的误差。在有些人看来这些技术无非是些“锦上添花”,可有可无的事。这里,我们向读者介绍一个“高招”,它可以让被调查者毫无顾虑地回答那些他不想回答的问题。它也表明,这类纯技术性的方法可以做成别的手段不易做成的事。

假设你是一个学校的教师。你发现最近一个时期学生考试作弊现象增加。为了

说明这一发现,以便提醒校方采取措施。你搞了一项调查。问卷是:

“你考试作过弊吗?”

然后,按简单随机抽样的方式选取了 $n = 100$ 名学生,让他们回答上述问题。用 n_A 记回答“是”的人数。于是全体学生的作弊率 π_A 可用 f_A 来估计。

$$f_A = n_A / n \quad (1)$$

实施上述调查方案的一个前提是:每个被调查对象都如实回答问题。但是所调查的问题有点敏感性,也就是说被调查者对如实回答有些顾虑。因此,我们实际计算得到的 f_A 比真实的 π_A 要小。按习惯作法,可以做点思想工作以减少顾虑。例如,讲明:“这次调查结果不涉及对个人评价。”这时,有人会认为:“谁知道以后你会不会用这材料整我呀!”另一种处理办法是不要求被调查者署名。但这也不能排除人家怕“查笔迹”的耽心。“作弊调

也要用计算器,那么得到 S_3 只需在最后按一次 (σ_{x-1}) 键,并不增加工作量。估计生产过程的标准差是保证产品质量的重要环节,应当力求估计精确,一百几十个数据来之不易,要化相当的时间和成本,而最后只为了图简便不充分利用这些数据可提供的信息,这未免是舍本求末。而实

际上也没有真正得到简便。

综上所述,我们建议用样本标准差 S_3 作为生产过程标准差 σ 的估计式,同时在编制 X 图、 \bar{X} 图和 \bar{X} 图时也利用 S_3 , 来确定控制界限。这个建议不仅具有理论上的优越性,也具有实践上的可行性。

查”事情还不严重.要是更敏感的问题就
更难让人讲实话了.而这在许多调查中,
特别是在社会调查中是经常遇到的.

好了,你不妨试一下下面的方法.

第一步,先做一些纸条.其中一部份写上

问题A:我考试作过弊.

另一些写

问题B:我考试没有作过弊.

这两类纸条数量之比为 $A:B=P:(1-P)$.

一般 $P \approx 1/2$.

第二步,把这两类纸条混在一起,放在一个口袋里.

第三步,让 n 个学生一个一个过来.每个人随机摸一个纸条.只让他自己看一下问题A还是B.(其他的人,包括调查者都不知道是哪个问题)然后,针对纸条上的问题,给出“是”或者“不是”的回答.最后放回纸条,混合后让下一个被调查者重复上面的程序.

第四步,你只要统计一下回答“是”的人数,记为 m_A .

第五步,利用已知的 P , n 和 m_A ,按下式计算

$$\hat{\pi}_A = \frac{m_A/n - (1-P)}{2P-1} \quad (2)$$

例如,在你抽查的 $n=100$ 人之中有 $m_A=66$ 个人回答“是”,而30%的纸条上有问题A,即 $P=0.3$.于是

$$\begin{aligned} \hat{\pi}_A &= \frac{m_A/n - (1-P)}{2P-1} \\ &= \frac{66/100 - (1-0.3)}{2 \times 0.3 - 1} = 0.1 \end{aligned}$$

即你的学生中有十分之一的人作过弊.

这个方法的优点在于,每一个被调查的人都明白:其他人不知道他回答的是问题A还是问题B.因此从他的回答中并不能推断这个人是否作过弊.这样,他的个人

秘密仍属于他自己,我们却可以利用 m_A ,即回答“是”的次数这样一个信息对总体作弊率进行估计,这不是有点“妙”吗?当然,这种方法的成功很大程度上要依赖被调查者对该方法的理解,这一定要预先讲明白.一个糊里糊涂的被调查者固守着准备说慌,那再好的方法也无济于事.

公式(2)是怎么来的?其推导十分初等:

$$\begin{aligned} \text{定义: } Y &= \begin{cases} 1 & \text{抽到问题A} \\ 0 & \text{抽到问题B} \end{cases} \text{ 以及} \\ X &= \begin{cases} 1 & \text{回答“是”} \\ 0 & \text{回答“不是”} \end{cases} \end{aligned}$$

于是利用全概率公式可得

$$\begin{aligned} P(\text{回答“是”}) &= P(X=1) \\ &= P(X=1|Y=1) \cdot P(Y=1) \\ &\quad + P(X=1|Y=0) \cdot P(Y=0) \\ &= \pi_A \cdot P + (1-\pi_A)(1-P) \end{aligned} \quad (3)$$

用 m_A/n 来估计 $P(X=1)$.于是

$$m_A/n = \pi_A \cdot P + (1-\pi_A)(1-P) \quad (4)$$

在上式中解出 π_A ,即得(2).

有得必有失.这个方法虽然好,但我们每问一个人所获得的“信息量”确实也少了.不妨算一算来证实直观的猜测.

什么叫一个事件的信息量?它的定义是

信息量 = 该事件概率的负对数(以2为底) (5) 信息量的单位称为比特.若A是“掷一枚硬币正面向上.”这样一个事件, $P(A)=1/2$.则A事件发生的信息量为

$$I(A) = -\log_2 1/2 = 1 \text{ 比特}$$

若记B为事件“某架飞机坠毁,” $P(B)=1/10000$.于是

$$I(B) = -\log_2 1/10000 = 13.3 \text{ 比特}$$

这也可以说明为什么报纸把客机坠毁看成是新闻而掷一个硬币是正面不当做新闻.

充其量不过是B的信息量大而已.

若记事件E: “某学生作过弊”. 我们知道 $P(E) = 0.1$. 这就是说当某人讲他作过弊时, 他提供给我们的信息量为

$$I(E) = -\log_2 0.1 = 3.3 \text{ 比特}$$

而事件F “在学生抽一张纸条后, 回答“是” 则

$$\begin{aligned} P(F) &= P(X = 1) \\ &= \pi_A P + (1 - \pi_A)(1 - P) \\ &= 0.1 \times 0.3 + 0.9 \times 0.7 = 0.66 \end{aligned}$$

于是

$$I(F) = -\log_2 0.66 = 0.6 \text{ 比特}$$

信息量这么小! 怪不得人家回答时不觉得

有什么损失呢.

用信息量来分析一种调查方法并不常见. 更多的是分析估计量的方差. 如果读者有兴趣不妨去看W.Cochran《抽样技术》一书第13章第17节. 当然, 更深刻的问题是信息量与方差之间的联系, 这就有点理论味了.

本文题目叫做“如何让人讲真话”. 初看起来有点像“如何做政治思想工作”一类文章. 不过现在不是提倡政工也要现代化吗?

问题: 请就一般情况证明

$$I(E) \geq I(F).$$

(上接第48页)

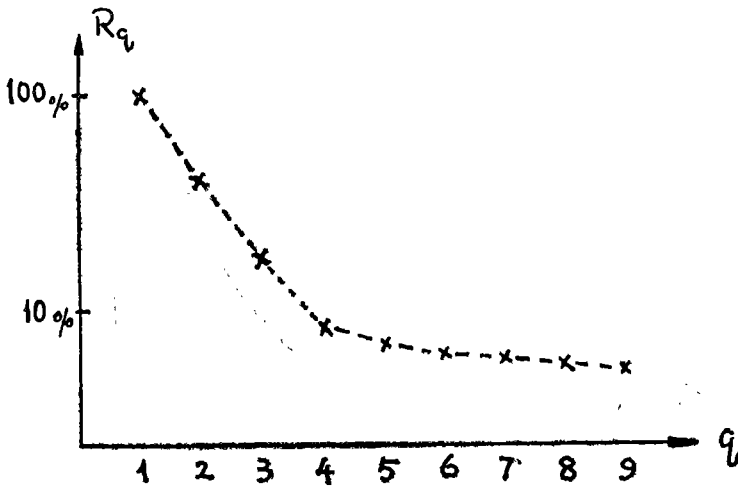


图 5.3 聚类特性曲线 (纵坐标是对数)

意慎用. 不过因为确实没有别的更好的办法, 所以上述做法还是值得我们注意的. 看来最好的办法还是把分类结果从物理背景的实际含义出发作一番检查, 看看是否可以接受.