

# 整群抽样总体均值的

在经济学研究中,数据分析是很重要的。可是如何获得数据,用哪种抽样方法得到数据是需要精心选择的。本文提出了运用整群抽样的方法,来减少数据分析时的误差。

整群抽样的样本是按一定方式,由某些群中所有的单元组成的。使用整群抽样的主要原因是为了节约成本,因为群内单元都比较集中,所以比分散的单位更便于调查。举例来说,人口调查中先取户,再调查户内人口,比对人直接用简单随机抽样成本小,因此整群抽样在社会经济抽样调查中被广泛应用。由于整群抽样的精度与群的性质有很大关系,我们通常的分群原则是:尽可能使群间方差小,群内方差大。但在实际问题中常基于某种方便的原则,按自然形成的单位定义群,如按行政区、部门行业等自然形成的群划分。这就导致了群内方差过小,群间方差过大。对于这种整群抽样,我们就要利用回归的方法来提高整群的估计精度。

本文就是在群内所含单元数目已知的条件下,利用抽样中的两种信息,一为主要特征信息,另一为辅助特征信息,建立它们之间的回归方程来进行估计。

设总体的主要统计特征为 $Y$ ,与之相关的辅助统计特征为 $X$ ,总体中含有 $N$ 个群,样本中有 $n$ 个群,总体第 $i$ 个群中含有 $M_i$ 个次级单元( $i=1,2,\dots,N$ ),共有 $M_0=\sum_{i=1}^N M_i$ 个次级单元,样本中第 $i$ 个群的次级单元数为 $m_i(i=1,2,\dots,n)$ ,对群进行简单随机抽样。

## 一、平均群的不加权回归估计

根据定义,总体均值回归估计量可如下表述 $\bar{y}_q = \bar{y} + b(\bar{X} - \bar{x})$  (1)

其中 $\bar{X}$ 已知, $b$ 为参数,可以为事先设定的常数 $b_0$ ,也可以是某个特定的统计量 $b_1$ 。

$$\text{一般取 } b_1 = \frac{\sum_{i=1}^n (\bar{y}_i - \bar{y})(\bar{x}_i - \bar{x})}{\sum_{i=1}^n (\bar{x}_i - \bar{x})^2}$$

则 $E(\bar{y}_q) = E(\bar{y} + b(\bar{X} - \bar{x})) = \bar{Y} + E(b(\bar{X} - \bar{x}))$ ,所以 $\bar{y}_q$ 是有偏估计量。

若 $b=b_0$ 已知,则 $\bar{y}_q$ 作为 $\bar{Y}$ 的估计量是无偏的。

$$\text{令 } \bar{M} = M_0/N, M_0 = \sum_{i=1}^N M_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n \frac{1}{M} \sum_{j=1}^{m_i} y_{ij} = \frac{1}{nM} \sum_{i=1}^n \sum_{j=1}^{m_i} y_{ij}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n \frac{1}{M} \sum_{j=1}^{m_i} x_{ij} = \frac{1}{nM} \sum_{i=1}^n \sum_{j=1}^{m_i} x_{ij}$$

$$\bar{y}_i = \frac{1}{M} \sum_{j=1}^{m_i} y_{ij}, \bar{x}_i = \frac{1}{M} \sum_{j=1}^{m_i} x_{ij}$$

$$\bar{y}_i = \frac{1}{M} \sum_{j=1}^{m_i} y_{ij}, \bar{X}_i = \frac{1}{M} \sum_{j=1}^{m_i} x_{ij}$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n \bar{y}_i, \bar{x} = \frac{1}{n} \sum_{i=1}^n \bar{x}_i$$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N \bar{X}_i, \bar{Y} = \frac{1}{N} \sum_{i=1}^N \bar{Y}_i$$

$$\bar{y}_q = \bar{y} + b(\bar{X} - \bar{x})$$

定理1  $\bar{x}, \bar{y}$  是 $\bar{X}, \bar{Y}$ 的无偏估计量。

证 将 $\bar{Y} = \frac{1}{N} \sum_{i=1}^N \bar{Y}_i$ 看作总体均值

$\bar{y} = \frac{1}{n} \sum_{i=1}^n \bar{y}_i$ 看作样本均值

由于对群的抽取是简单随机的,所以 $\bar{y}$ 是 $\bar{Y}$ 的无偏估计量。

同理可证 $\bar{x}$ 是 $\bar{X}$ 的无偏估计量。

引理1  $\bar{x}, \bar{y}$  分别是抽自总体均值为 $\bar{X}, \bar{Y}$ 的有限总体的简单随机样本的均值,样本容量为 $n$ ,则对非负整数 $k, m$ 有:

$$E(\bar{y} - \bar{Y})^k = O(n^{-(k+1)/2})$$

$$E(\bar{y} - \bar{Y})^k (\bar{x} - \bar{X})^m = O(n^{-(k+m+1)/2})$$

引理2: 设 $B = \frac{\sum_{i=1}^n (\bar{Y}_i - \bar{Y})(\bar{X}_i - \bar{X})}{\sum_{i=1}^n (\bar{X}_i - \bar{X})^2}$ 是有限

总体回归系数, $\bar{x}$ 是简单随机样本 $\bar{x}_i$ 的

均值,记 $\varepsilon_i = (\bar{Y}_i - \bar{Y}) - B(\bar{X}_i - \bar{X})$

则 $\bar{\varepsilon} = 0$

$$\sum_{i=1}^n \varepsilon_i (\bar{X}_i - \bar{X}) = 0$$

$$E\left[\frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{X}_i - \bar{X})\right]^2 = O\left(\frac{1}{n}\right)$$

$$E\left[\frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{X}_i - \bar{X})\right]^4 = O\left(\frac{1}{n^2}\right)$$

引理3: 在简单随机抽样下,记

$$b' = \frac{\sum_{i=1}^n (\bar{y}_i - \bar{y})(\bar{x}_i - \bar{x})}{\sum_{i=1}^n (\bar{x}_i - \bar{x})^2} \text{ 则}$$

$$b' = B + \frac{\sum_{i=1}^n \varepsilon_i (\bar{x}_i - \bar{x})}{\sum_{i=1}^n (\bar{x}_i - \bar{x})^2}$$

$$E(b' - B) = O\left(\frac{1}{\sqrt{n}}\right)$$

$$E(b' - B)^2 = O\left(\frac{1}{n}\right)$$

$$E(b' - B)^4 = O\left(\frac{1}{n^2}\right)$$

定理2: 在简单随机抽样下,由以上引理

可以得到,整群回归估计量 $\bar{y}_q$ 有以下性质:

$$(1) B(\bar{y}_q) = E(\bar{y}_q) = E(\bar{y}) - \bar{Y}$$

$$= \frac{1-f}{(n-1)S_x^2} \frac{\sum_{i=1}^n \varepsilon_i (\bar{X}_i - \bar{X})}{N-1} + O\left(\frac{1}{n^{3/2}}\right) = O\left(\frac{1}{n}\right)$$

$$(2) \text{MSE}(\bar{y}_q) = \frac{1-f}{n} S_x^2 (1-p^2) + O\left(\frac{1}{n^{3/2}}\right) =$$

$$O\left(\frac{1}{n}\right)$$

$$(3) E\left(\frac{1}{n-2} \sum_{i=1}^n [(\bar{y}_i - \bar{y})(\bar{x}_i - \bar{x}) - b'(\bar{x}_i - \bar{x})]^2\right)$$

$$= S_x^2 (1-p^2) + O\left(\frac{1}{\sqrt{n}}\right)$$

$$\text{其中 } S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{Y}_i - \bar{Y})^2, S_x^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{X}_i - \bar{X})^2$$

$$S_{yx}^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{Y}_i - \bar{Y})(\bar{X}_i - \bar{X}),$$

# 回归估计

■金莹 牛美玲 汤银才

$$\rho^2 = \frac{S_{YX}^2}{S_X^2 S_Y^2}$$

由(1)知  $\bar{y}_q = \bar{y}_i + b'(\bar{X} - \bar{x})$  是  $\bar{Y}$  的有偏估计量。

推论1: 当  $n$  很大时,  $V(\bar{y}_q) \approx \text{MSE}(\bar{y}_q) \approx \frac{1-f}{n} S_Y^2 (1-\rho^2)$

证:  $\text{MSE}(\bar{y}_q) = V(\bar{y}_q) + [B(\bar{y}_q)]^2 = V(\bar{y}_q) + O(\frac{1}{n^2})$

$$V(\bar{y}_q) = \text{MSE}(\bar{y}_q) + O(\frac{1}{n^2}) = \frac{1-f}{n} S_Y^2 (1-\rho^2) + O(\frac{1}{n^{3/2}}) + O(\frac{1}{n^2})$$

$$= \frac{1-f}{n} S_Y^2 (1-\rho^2) + O(\frac{1}{n^{3/2}})$$

$$\text{故 } V(\bar{y}_q) \approx \text{MSE}(\bar{y}_q) \approx \frac{1-f}{n} S_Y^2 (1-\rho^2)$$

推论2:  $\hat{V}(\bar{y}_q) = \frac{1-f}{n} \frac{1}{n-2} \sum_{i=1}^n [(y_i - \bar{y}) - b'(\bar{x}_i - \bar{x})]^2$  是  $V(\bar{y}_q)$  的近似无偏估计量。

证:  $E(\frac{1-f}{n} \frac{1}{n-2} \sum_{i=1}^n [(y_i - \bar{y}) - b'(\bar{x}_i - \bar{x})]^2) = \frac{1-f}{n} [S_Y^2 (1-\rho^2) + O(\frac{1}{\sqrt{n}})] \approx V(\bar{y}_q)$

$$\text{因此 } \hat{V}(\bar{y}_q) = \frac{1-f}{n} \frac{1}{n-2} \sum_{i=1}^n [(y_i - \bar{y}) - b'(\bar{x}_i - \bar{x})]^2$$

$$\text{推论3: } \bar{y}_q \text{ 是可用的。}$$

$$\text{证: } S \frac{B(\bar{y}_q)}{\sqrt{\text{MSE}(\bar{y}_q)}} = \frac{O(\frac{1}{n})}{\sqrt{O(\frac{1}{n})}}$$

$$= O(\frac{1}{\sqrt{n}})$$

因此  $\bar{y}_q$  是可用的

二、平均群的加权回归估计

在群大小不等的情况下, 群的样本量大小不同, 在整个实验中所处的地位不平等, 这样用这些数据作回归时, 就不

能把他们同等看待, 用不加权的估计量精度可能会不高, 必须有所侧重, 这就需要用到加权回归来计算。

根据加权回归理论, 取

$$b^* =$$

$$\frac{\sum_{i=1}^n m_i y_i \bar{x}_i - \frac{1}{M_n} \sum_{i=1}^n m_i y_i \sum_{i=1}^n m_i \bar{x}_i}{\sum_{i=1}^n m_i \bar{x}_i^2 - \frac{1}{M_n} (\sum_{i=1}^n m_i \bar{x}_i)^2}$$

$$= \frac{\sum_{i=1}^n m_i y_i}{\sum_{i=1}^n m_i \bar{x}_i}, \quad y = \frac{\sum_{i=1}^n m_i y_i}{\sum_{i=1}^n m_i \bar{x}_i}$$

$$\text{其中 } M_n = \sum_{i=1}^n m_i$$

则  $\bar{Y}$  的加权回归估计量为

$$\bar{y}_q = \bar{y} + b^*(\bar{X} - \bar{x}), \text{ 同样它也是有偏估计。}$$

方差的估计为

$$\hat{V}(\bar{y}_q) = \frac{1-f}{n} \frac{1}{n-1} \sum_{i=1}^n \frac{m_i}{M_n} [(y_i - \bar{y}) - b^*(\bar{x}_i - \bar{x})]^2$$

群大小不等时, 整群抽样下简单估计的估计量为  $\bar{y} = \frac{\sum_{i=1}^n y_i / n}{\sum_{i=1}^n M_i / N}$ , 其抽

样方差为

$$V(\bar{y}) = \frac{1-f}{M} \frac{1-f}{n} \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 = \frac{1-f}{n} \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2 = \frac{1-f}{n} S_Y^2$$

$$\text{其中 } Y_i = \sum_{j=1}^{M_i} Y_{ij}, \bar{Y} = \sum_{i=1}^N Y_i / N$$

与不加权的估计量方差相比

$$\frac{1-f}{n} S_Y^2 (1-\rho^2) - \frac{1-f}{n} S_Y^2 = \frac{1-f}{n} S_Y^2 \cdot \rho^2 < 0$$

所以, 不加权的回归估计量方差小, 其估计精度优于整群抽样的简单估计量。

例: 以1999年全国工资与消费情况的抽样调查为例, 比较群大小不等情况下四种估计量。

把全国31个省作为不同群体, 从中

简单随机抽取8个省的资料(见下表)。

$$\text{已知 } N=31, n=8, \sum_{i=1}^N M_i=11278, \sum_{i=1}^n M_i=$$

$$4429.6, \sum_{i=1}^n X_i=11473, \bar{X}=1.0173$$

全国8省职工工资、消费调查表

| $m_i$    | $x_i$    | $y_i$    |
|----------|----------|----------|
| 职工人数(万人) | 工资总额(万元) | 消费总额(万元) |
| 403      | 566.3762 | 302.1887 |
| 379.4    | 230.1061 | 132.5237 |
| 352.6    | 252.3911 | 129.1108 |
| 327.1    | 544.3271 | 269.7819 |
| 408.9    | 266.4392 | 159.5450 |
| 320.4    | 304.0596 | 168.7447 |
| 305.9    | 206.4519 | 106.5245 |
| 335.4    | 232.4657 | 132.5920 |

(1) 不加权估计

$$\text{平均群回归估计系数 } b' = 0.4889$$

$$\text{平均群回归估计量 } \bar{y}_{lk} = 0.4212 \text{ (万元)}$$

$$\text{平均群回归估计量方差 } v_1 = 7.8309 \times 10^{-5}$$

(2) 加权估计

$$\text{平均群加权回归估计系数 } b^* = 0.4902$$

$$\text{平均群加权回归估计量 } \bar{y}_{lk} = 0.5435 \text{ (万元)}$$

$$\text{平均群加权回归估计量方差 } v_2 = 6.1718 \times 10^{-5}$$

(3) 简单估计

$$\text{估计量 } \bar{y}_{\text{简单}} = 0.4814 \text{ (万元)}$$

$$\text{估计量方差 } v_3 = 0.0036$$

从这个例子我们可以看出, 在群大小不等的情况下, 本文提出的新的回归估计量及新的加权回归估计量比简单估计量的方差小, 即精度高, 估计效果好。因此, 在群大小不等的情况下, 我们可以用本文提出的平均群回归估计量及其加权回归估计量, 估计效果会更好。

(作者单位/上海师范大学数理信息学院, 华东师范大学)  
(责任编辑/李友平)