

⑥ 層別とその推定値への影響

青山博次郎

サンプリングに於て、層別を行うことは通常のようになっている。この時分析の段階ではAという性質について層別した各層の統計値(例えは標本平均)を敬し、いのであるが、適当な資料が予め得られなかつたためにBという性質を層別してサンプリングを行うことがある。このようにして得られた data を発表に際してはAという性質で層別して行う場合どれ位の偏りを生ずるかを考えてみることにしよう。

簡単な実例では学級数の各段階毎の標識の標本平均を計算するのに、サンプリングでは資料の関係から生徒数を用いた場合などがこれに当る。

扱て標識をYとし、性質Bについて層別し、これより比例抽出法を用いてサンプルを抽出したとする。このサンプルを更に性質Aについて層別しなおすのであるから二回抽出法の公式が使われる。全平均 \bar{Y} の推定値 \bar{y} は

$$\bar{y} = \frac{1}{n} \sum_j \sum_k y_{jk} = \frac{1}{n} \sum_j \sum_k \sum_i y_{ijk} \quad (1)$$

但しjは性質Aによる層、iは性質Bによる層を表わし、nはサンプル数である。

このとき

$$E(\bar{y}) = E\left(\frac{1}{n} \sum_i \sum_k y_{ik}\right) = \frac{1}{n} \sum_i \frac{N_i}{N_i} \sum_{k=1}^{N_i} Y_{ik} = \frac{1}{N} \sum_i \sum_k Y_{ik} = \bar{Y} \quad (2)$$

となつて unbiased である。

また分散は

$$V(\bar{y}) = \frac{\sigma^2}{n} - \frac{1}{n} \sum_i P_i (\bar{Y}_i - \bar{Y})^2 \quad (3)$$

であつて A による層別の影響はない。(註)

次に各層毎の推定平均について考えると

$$\bar{y}_i = \frac{1}{n_i} \sum_j \sum_k y_{ijk} = \frac{1}{n_i} \sum_k y_{ik} \quad (4)$$

$$E(\bar{y}_i) = \bar{Y}_i \quad (5)$$

となるから \bar{y}_i については unbiased なることは勿論であるが

$$\bar{y}_j = \frac{1}{n_j} \sum_i y_{ij} = \frac{1}{n_j} \sum_i \sum_k y_{ijk} \quad (6)$$

は unbiased とならない。その理由に断るまでもなく分母の n_j が確率変数だからである。

さて一般に二つの確率変数 X, Y について

$$\frac{X}{Y} = \frac{E(X)}{E(Y)} + \frac{XE(Y) - YE(X)}{E^2(Y)} + R, \quad R \text{ は誤差項} \quad (7)$$

が成立つ。そこで

$$\begin{aligned} E\left(\sum_k \sum_i y_{ijk}\right) &= \sum_i n_{ij} \sum_{(n_{ij} \dots n_{ij}^*)} P_{ij}^{n_{ij}} P_{ij'}^{n_{ij'}} \dots P_{ij^*}^{n_{ij^*}} \frac{n_i!}{n_{ij}! n_{ij'}! \dots n_{ij^*}!} \cdot \frac{1}{N_j} \sum_{k=1}^{N_j} Y_{ijk} \\ &= \sum_i \frac{n_i}{N_i} Y_{ij} = \frac{n}{N} Y_j \end{aligned} \quad (8)$$

但し $P_{ij} = N_{ij}/N_i$ 等, $(n_{ij} \dots n_{ij}^*)$ は $n_{ij} + \dots + n_{ij^*} = n_i$ なる

ときの n_{ij}, \dots, n_{ij^*} のすべての組合せについての和を示す。

① 著者：二回抽出法について，統計研講究録，第6巻第6号

(昭和25年9月) 参照。

記号や条件もこの論文中的ものを互用いてある。

同様にして

$$E(n_j) = \frac{n}{N} N_j \quad (9)$$

従つて (7) により, (8), (9) を用いて

$$E(\bar{y}_j) \doteq \frac{E(\sum_i \sum_r y_{ijr})}{E(n_j)} = \frac{Y_j}{N_j} = \bar{Y}_j \quad (10)$$

となるから近似的に unbiased である。

次に (7) の $E(R)$ を評価してみると

$$E(R) \doteq \frac{E(X)}{E^3(Y)} (E(Y^2) - E^2(Y)) = \frac{E(X) \nabla(Y)}{E^3(Y)} \quad (11)$$

従つて同様の計算を行うと

$$V(n_j) = E(n_j^2) - E^2(n_j) = \frac{n}{N} N_j - \frac{n}{N} \sum_i \frac{N_{ij}^2}{N_i} \quad (12)$$

$$\therefore E(R) \doteq \frac{N}{n} \frac{\bar{Y}_j}{N_j} \left(1 - \sum_i \frac{N_{ij}^2}{N_i N_j}\right) = O\left(\frac{1}{n}\right) \quad (13)$$

即ち

$$E(\bar{y}_j) = \bar{Y}_j + O\left(\frac{1}{n}\right) \quad (14)$$

之が (10) の代りに用いるべき式である。

また \bar{y}_j の分散は

$$V\left(\frac{X}{Y}\right) \doteq \frac{E(X - \varphi Y)^2}{E^2(Y)}, \quad \text{但し } \varphi = \frac{E(X)}{E(Y)} \quad (15)$$

を用いて計算を行えば y_j と n_j の相関係数 ρ として (通常は正となることが多から大か目にはなして)

$$V(\bar{y}_j) = \frac{1}{N_j^2} \sum_i \frac{N_{ij}^2 \sigma_{ij}^2}{n_{ij}} + \frac{N}{n N_j^2} \sum_i N_{ij} \left(1 - \frac{N_{ij}}{N_i}\right) \left(\bar{Y}_{ij}^2 + \bar{Y}_j^2 + \frac{\sigma_{ij}^2}{n_{ij}}\right) \quad (16)$$

となること分かる。