

有意抽出時の標本誤差とサンプルサイズ、母分散の関係

太郎丸博（京都大学）

1 問題

サンプリング理論は無作為抽出した標本を前提に作られており、有意抽出した標本（以下、有意標本と略称）の性質についてはざっと検索した限り研究が見当たらない。しかし、多くの学問分野で有意標本は活用されており、その性質について理解することは、それらの研究成果を理解する上で重要である。この報告では、有意標本からある変数 Y の母集団における平均値（以下、母平均と略称）を推測する際に、標本サイズと母分散が標本誤差（次節で定義）にどのように影響するのか論じる。

この研究では以下の3つを主張する。

1. 有意標本でも標本誤差は、母集団での標準偏差（以下、母標準偏差と略称）に反比例する。
2. 有意抽出された標本（以下、有意標本）でも、標本誤差は、サンプルサイズの平方根に反比例する。
3. 有意抽出した標本の標本誤差は、無作為抽出した場合の標本誤差より大きくなる。

1については次節で証明する。2と3については一般的な答えを出す用意がないので、ある特定の条件で作りに出した50個の有意標本を例に検討する。2については有意標本でも *independently and identically distributed* (iid) という条件をつければ証明できることを補論で示す。

2 標本誤差と母標準偏差の関係

Y という変数の母平均 $m(Y)$ を推測しようとしているとする。母集団から N 個の事例 i ($i = 1, \dots, N$) を有意抽出し、標本とする。これを M 回繰り返すとする。 j 番目の標本の Y の平均を $m(Y_j)$ ($j = 1, \dots, M$) とする。 j 番目の標本の i 番目の事例の Y の値を Y_{ij} とすると、

$$m(Y_j) = \frac{\sum_{i=1}^N Y_{ij}}{N} \quad (1)$$

と定義する。また、標本誤差 $s(Y)$ は、

$$s(Y) = \sqrt{\frac{\sum_{j=1}^M (m(Y_j) - m(Y))^2}{M-1}} \quad (2)$$

と定義する。ここで証明したいのは以下の定理である。

定理 1: $s(aY) = a \times s(Y)$

証明

Y を a 倍すると、母平均も a 倍になる（証明省略、統計の教科書をあたれ）。すなわち、

$$m(aY) = a \times m(Y) \quad (3)$$

である。同様に有意標本であっても Y_{ij} を a 倍すると、以下のように標本平均も a 倍になる。

$$m(aY_j) = \frac{\sum_{i=1}^N a Y_{ij}}{N} = a \times \frac{\sum_{i=1}^N Y_{ij}}{N} = a \times m(Y_j) \quad (4)$$

それゆえ、以下のように Y を a 倍すると、標本誤差も a 倍になる。

$$\begin{aligned} s(aY) &= \sqrt{\frac{\sum_{j=1}^M (m(aY_j) - m(aY))^2}{M-1}} = \sqrt{\frac{\sum_{j=1}^M (a \times m(Y_j) - a \times m(Y))^2}{M-1}} \\ &= \sqrt{\frac{a^2 \sum_{j=1}^M (m(Y_j) - m(Y))^2}{M-1}} = a \sqrt{\frac{\sum_{j=1}^M (m(Y_j) - m(Y))^2}{M-1}} = a \times s(Y) \quad (5) \end{aligned}$$

以上で証明終わり。

3 標本誤差とサンプルサイズの関係 データと方法

2021年5月23～24日に、私の講義を受講している学生にオンラインで以下のような有意抽出をするように指示した。「1～100の数字の中から10個好きな数字を自由に選んで下さい（同じ数字を2回以上選ばないでください）」。その結果50人から回答を得た。これらの数字を母集団における事例のIDとみなすと、 $M = 50$, $N = 10$ である。理解を容易にするために、50人の学生が科学社会学の研究のために、ある分野で出版された100本の論文の中から、それぞれ10本ずつ有意抽出し、被引用回数の母平均を推測しようとしていると考える。各IDの論文の被引用回数は、母平均 = 21、母標準偏差 = 22で、ID番号の二次関数に近似するように設定した。

実際には一人あたり10本抽出しただけだが、その10個の中から n 個のサンプルを無作為抽出（非復元リサンプリング）し、 n 個のサンプルを50回有意抽出したとみなす。この n 個のサンプルに関して標本誤差を計算する。 n を2から10まで変化させたときの標本誤差の大きさを計算する。以上を500回繰り返して、標本誤差と標本サイズの関係を検討する。

最後に、同じ母集団から無作為標本を有意標本と同様に $M = 50$, $N = 10$ だけ抽出し、その際の標本誤差を計算する。これを500回繰り返して、その結果を有意標本の標本誤差と比較する。

シミュレーションの結果

被引用回数の標本平均を50個の標本についてすべて計算し、それらのヒストグラムを描いたのが、図1左パネルである。母平均よりもやや高めになっている。

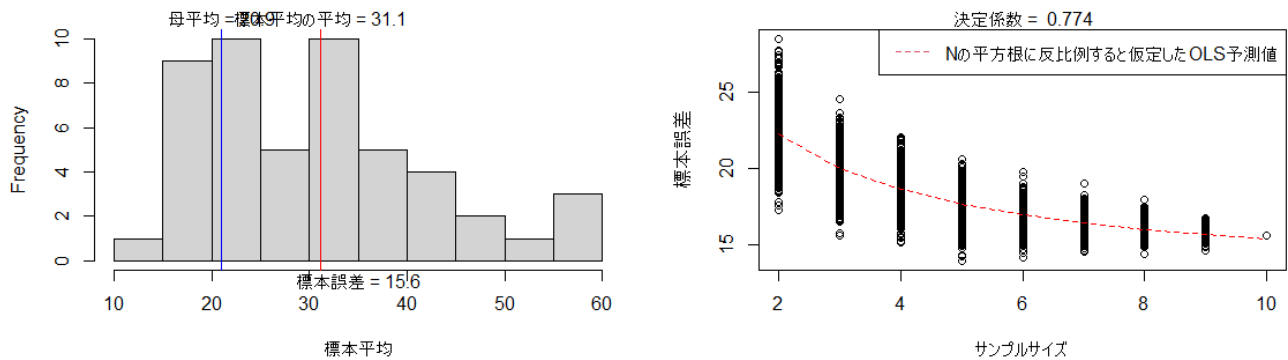


図1 標本平均のヒストグラム (左パネル) と有意標本のサイズと標本誤差の関係 (右パネル)

非復元リサンプリングによって 50 個の有意標本のサイズを 2~10 まで変化させ、それぞれの抽出誤差を計算することをさらに 500 回繰り返した結果が図 1 右パネルである。標本サイズの平方根に反比例して標本誤差が変化している。標本サイズに単純に回帰させた場合の決定係数は 0.662 である。

4 無作為抽出した場合の標本誤差

被引用回数が正規分布してれば、標準誤差は 7 である。有意抽出と同じサイズ、抽出回数で 500 回シミュレーションしてみると、標本誤差は最低で 5 最大で 8.9、無作為標本の標本誤差の平均は 7 であり、有意抽出標本よりもずっと小さくなった。

5 補論

Y_{ij} は期待値が μ で、independently and identically distributed (iid) だと仮定する。また M と N は中心極限定理による十分な近似が得られるほど大きな値をとると仮定する。サンプルサイズが N_1, N_2 ($N_1 < N_2$) のときの標本誤差を $s_1(Y)$, $s_2(Y)$ とすると、証明したい命題は、

定理 2: $s_1(Y) > s_2(Y)$

である。サンプルサイズが N_1, N_2 のときの j 番目の標本平均を $m_1(Y_j)$, $m_2(Y_j)$ とすると、

$$\begin{aligned}
 s_1(Y)^2 &= \frac{\sum_{j=1}^M (m_1(Y_j) - m(Y))^2}{M-1} = \frac{\sum_{j=1}^M (m_1(Y_j)^2 - 2m_1(Y_j) \cdot m(Y) + m(Y)^2)}{M-1} \\
 &= \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{2 \cdot m(Y)}{M-1} \sum_{j=1}^M m_1(Y_j) + \frac{M}{M-1} m(Y)^2 \quad (6)
 \end{aligned}$$

である。iid の仮定と中心極限定理より、 $\sum_{j=1}^M m_1(Y_j) = M\mu$ であるから、これを式 (6) に代入して、

$$s_1(Y)^2 = \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{2M \cdot m(Y)}{M-1} \mu + \frac{M}{M-1} m(Y)^2 \quad (7)$$

である。同様にして、

$$s_2(Y)^2 = \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 - \frac{2M \cdot m(Y)}{M-1} \mu + \frac{M}{M-1} m(Y)^2 \quad (8)$$

式(7)から式(8)をひくと、

$$s_1(Y)^2 - s_2(Y)^2 = \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 \quad (9)$$

である。

Y_{ij} の分散を v とすると、中心極限定理により、 $m_1(Y_j)$ の分散は

$$v(m_1(Y_j)) = \frac{v}{N_1} = \frac{\sum_{j=1}^M (m_1(Y_j) - \mu)^2}{M-1} = \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{M}{M-1} \mu^2 \quad (10)$$

であり、同様に $m_2(Y_j)$ の分散は

$$v(m_2(Y_j)) = \frac{v}{N_2} = \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 - \frac{M}{M-1} \mu^2 \quad (11)$$

である。 $N_1 < N_2$ だから、式(10)、式(11)より、

$$\begin{aligned} \frac{v}{N_1} &> \frac{v}{N_2} \\ v(m_1(Y_j)) &> v(m_2(Y_j)) \\ \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{M}{M-1} \mu^2 &> \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 - \frac{M}{M-1} \mu^2 \\ \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 &> \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 \\ \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 &> 0 \quad (12) \end{aligned}$$

である。式(9)と式(12)より、

$$\begin{aligned} s_1(Y)^2 - s_2(Y)^2 &= \frac{1}{M-1} \sum_{j=1}^M m_1(Y_j)^2 - \frac{1}{M-1} \sum_{j=1}^M m_2(Y_j)^2 > 0 \\ s_1(Y)^2 &> s_2(Y)^2 \\ s_1(Y) &> s_2(Y) \quad (13) \end{aligned}$$

以上で証明終わり。