

社会調査における回答者属性データの地域集計

01102840 会津大学 出水田智子 IZUMITA Tomoko

1. はじめに

世論調査などで社会的な集団の特性、意識、行動などを標本調査として実施するときの中心的な課題は、調査対象者の回答比率を通してその背後の母集団全体の比率を推定することである。このとき必要なことは、標本から得られた結果から母集団の統計値が精度よく推定されることである。

このような推定比率の精度を決定するのは、調査対象者の抽出方法によって生じる母比率と標本比率の差だけではない。調査の精度は母比率と回答比率の差で評価されるため、むしろ回収率の低さや不完全な回答結果あるいは作業ミスなどの影響による標本比率と回答比率の差の開きが影響する。さらに結果を比率で表現するために、1集計単位あたりの標本数の大きさやばらつきによる影響も無視できない時がある。

そこで本稿では、分析結果に含まれる調査誤差のうち、回答者属性データの地域集計作業が母比率推定に与える影響について計算例を基に検討する。

2. 調査誤差の定義

調査対象全体の母比率推定における調査誤差は、以下のように標本誤差  $\epsilon_S$  と非標本誤差  $\epsilon_R$  の和で定義される。

$$\begin{aligned} \epsilon_0 &= p - \hat{p} \\ &= (p - \hat{p}) + (\hat{p} - \hat{p}) \\ &= \epsilon_S + \epsilon_R \end{aligned} \tag{1}$$

ただし  $p$  は母比率、 $\hat{p}$  は標本比率、 $\hat{p}$  は回答比率である。母集団数、標本数、回答数をそれぞれ  $N, n, t$ 、その中の各個体の値を  $\theta_j$  ( $j=1, \dots, N$ )、 $\theta^*_j$  ( $j=1, \dots, n$ )、 $\theta^{\cdot}_k$  ( $k=1, \dots, t$ )、とすると、

$$p = \sum \theta_j / N$$

$$\hat{p} = \sum \theta^*_j / n$$

$$\hat{p} = \sum \theta^{\cdot}_k / t \tag{2}$$

となり、各集団におけるその値の平均に等しくなる。

集計誤差は全体を一つとして求めた(1)の結果と、異なる単位で集計した結果を比較することで求めることが出来る。そこで、母集団、標本、回答を共通な属性によって  $L$  個の層に分割し集計することをここでは考える。 $h=1, \dots, L$  のとき、層別の母集団数、標本数、回答数を  $N_h, n_h, t_h$ 、その中の各個体の値を  $\theta_{hi}, \theta^*_{hi}, \theta^{\cdot}_{hk}$  とすると、集計単位別の各集団の比率はそれぞれ

$$p_h = \sum \theta_{hi} / N_h$$

$$\hat{p}_h = \sum \theta^*_{hi} / n_h$$

$$\hat{p}_h = \sum \theta^{\cdot}_{hi} / t_h \tag{3}$$

となり、先ほどと同様に各層の値の平均に等しくなる。ここで母集団における各層の構成比率を  $w_h = N_h/N$ 、標本集団における構成比率を  $w^*_h = n_h/n$ 、回答集団における各層の構成比率を  $w^{\cdot}_h = t_h/t$  とおくと、集計後に算出される母集団、標本、回答における集計による比率の差は

$$\epsilon_a = p - \sum w_h p_h$$

$$\epsilon_a^* = \hat{p} - \sum w^*_h \hat{p}_h$$

$$\epsilon_a^{\cdot} = \hat{p} - \sum w^{\cdot}_h \hat{p}_h \tag{4}$$

となる。

そこで(1)と(4)から、集計化による誤差項  $\epsilon_a$  が考慮された調査誤差の推定モデル  $\epsilon$  を以下のように定義する。

$$\begin{aligned} \epsilon &= (p - \sum w_h p_h) - (\hat{p} - \sum w^*_h \hat{p}_h) \\ &+ ((\hat{p} - \sum w^*_h \hat{p}_h) - (\hat{p} - \sum w^{\cdot}_h \hat{p}_h)) \\ &= (p - \hat{p}) + (\hat{p} - \hat{p}) + (\sum w^{\cdot}_h \hat{p}_h - \sum w_h p_h) \\ &= \epsilon_S + \epsilon_R + (\sum w^{\cdot}_h \hat{p}_h - \sum w_h p_h) \\ &= \epsilon_S + \epsilon_R + \epsilon_a \end{aligned} \tag{5}$$

### 3. 計算例

使用する調査データの概要は以下の通りである。

- ・ 調査目的：会津若松市市民の生活意識調査
- ・ 標本抽出方法：系統抽出法(7つの初期値)
- ・ 抽出枠：住民基本台帳記載の20歳以上の男女
- ・ 母集団数：91420 (=N)
- ・ 標本数：2037 (=n\*)
- ・ 調査方法：無記名回答による郵送
- ・ 調査時期：2001年12月

回収された解答用紙(913)のうち、回答者属性データが網羅されている有効回答(n=787)を対象に、性別や年齢など回答者属性データに関する比率を異なる2つの地域スケール(郵便番号別、地区別)で集計したときの推定値と真値の差を求めた。詳細な計算結果は当日紹介する。

今回利用したデータは地域別世帯順に並んだ抽出枠から等間隔に抽出したため、母数が不明な場合は理論的な式を使った誤差推定ができない。今回は行政の統計で確認した正確な母比率を用いて

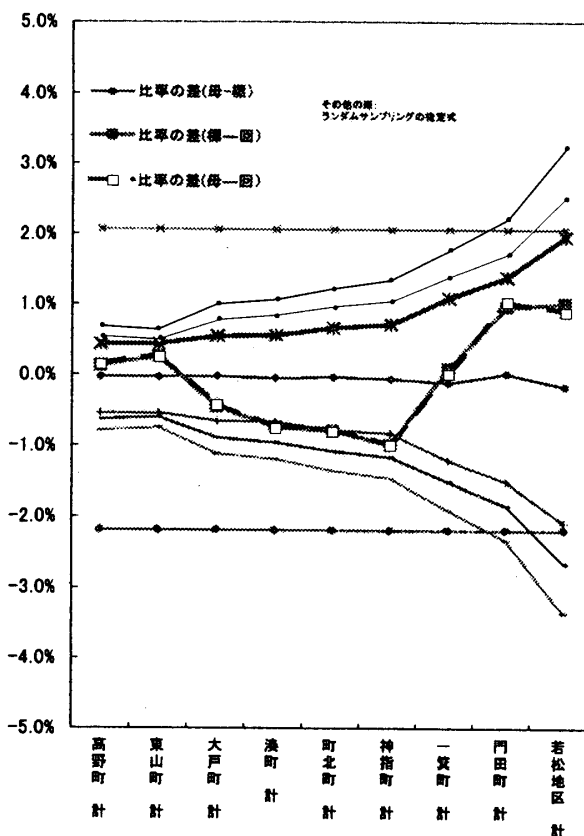


図1 地区別集計における男性比率の差と推定値

誤差を得た。念のため事後的に、ランダム・サンプリングの式に標本比率をあてはめて、近似的な区間推定も行った(図1)。

集計単位別の母比率と、回答比率がわかれば(5)式によって無回答などによる誤差と集計誤差を区別できることが確認された。

### 4. おわりに

本稿では、標本誤差  $\epsilon_s$  がある程度わかっているときに非標本誤差  $\epsilon_r$  と集計誤差を区別して求めるための評価式を導出した。今回データが限られたため、十分な検討ができなかったため、今後さらにいろいろなデータで検証作業をおこなう予定である。分散に基づく標本誤差推定では  $\epsilon_s$  の符号が不明なためにそれぞれの絶対値の和をとる場合が多い。このため調査誤差は(1)よりも大きな値で推定されているが、むしろ予測不可能な  $\epsilon_r$  の評価が問題である。本来の調査目的である母集団の真の値(例えば母比率)が既知でなければ厳密な意味で  $\epsilon_s$ 、 $\epsilon_r$ 、 $\epsilon_a$  の3種類の誤差を回答結果から分離することはできないが、今回の計算例のように補助情報などによって真の値や標本誤差が既知の回答者属性データで非標本誤差と集計誤差を分離することは可能である。少なくとも回答者属性に関するクロス分析結果を地域全体に一般化する場合や、リサンプリング標本精度の評価において利用可能であろう。今後は居住地情報の集約によって生じる集計誤差が調査結果にどのように影響するか、空間データの視点からを検討したい。

### 参考文献

- [1] NHK 放送文化研究所世論調査部編(1996)：「世論調査事典」，大空社。
- [2] 会津大学文化研究センター(2002)：会津若市における生活と文化に関するアンケート調査—基礎集計—。会津大学文化研究センター研究年報第8号，pp. 92-158。