

Hot Deck 補完の下での 標本平均の漸近正規性について

青山学院大学経済学部准教授
元山 齊

1. はじめに

多くの社会・経済データやそれらの集計データの基礎となるデータの多くは対象となる母集団に対する全数または標本調査に基づいて作成されている。しかしながら、調査においては、調査対象者の不在ないしは調査拒否等による調査対象者からの回答を得られないこと (unit non-response) や個別の質問項目が sensitive であることや質問文の意図が分かりにくいこと等による無回答による欠損値が存在する (item non-response)。

欠損値に対処する方法として、統計作成の実務においては長い間、補完(補定、代入、imputation)と呼ばれる手法で、欠損値に対してもっともらしいと考えられる数値を代入する方法が長い間用いられてきた。補完の方法については、回答が得られたレコードの平均値を代入する方法(平均値補完)や関心のある変数と相関の高い補助変数による回帰の予測値を代入する方法(回帰補完)などが提案されているが、それらの中でも Hot Deck 補完は欠損値をその欠損値が含まれるデータセット内の、欠損値に何らかの意味で似ていると考えられる完全データの値によって補完する方法であり、マッチング代入とよばれる補完方法において欠損値と同じデータセットを用いる特別な場合と考えることもできる。

欠損値補完を行ったデータセットに基づいて統計的推論を行う際には、統計量の分布の標準偏差や分布そのものについて知ることが重要となるが、Abadie and Imbens(2012)は、統計的因果推論でしばしば用いられる Matching Estimator がマルチンゲールという性質を持った確率変数列から構成されることを示し、マルチンゲール中心極限定理に基づいて Matching Estimator の漸近正規性を証明した。また、その導出の際に用いられた考え方が Cell Hot Deck とよばれる方法で補完を行った際の標本平均にも適用できることを示し Cell Hot Deck 補完を行った際の標本平均の漸近正規性を証明した。本稿は Abadie and Imbens(2012)の当該箇所の概略をまとめたノートである。本稿は上記論文の数理的行間の一部を埋めた以外は、理論的にも視点的にも新しい結果の提示をしていくのではなく、上記論文の流れに沿った紹介にとどまっているが、今後、応用面でさらに有益な結果を導出するための足掛かりとなることを意図している。

2. Cell Hot Deck 補完とその下での標本平均の漸近正規性

いま関心のある変数を Y 、その変数と関係の強い共変量を X とし、大きさ N の標本 $(Y_i, X_i), i = 1, \dots, N$ が調査によって調べられたとする。しかしながら $Y_i, i = 1, \dots, N$ の一部には何らかの理由で欠損値が生じているとする。このようなとき、 Y の値と関係のある共変量 X の値の同じセル内の回答された完全なレコードの値で補完することがしばしば行われる。

いま共変量 X の値に基づいて、観測値の全体が重なることのない T 個のクラス(セル) $C_1, \dots, C_t, \dots, C_T$ に分割されているとする。そして Y の欠損値は同じセル内の完全なレコードによって補完されるとする。

Hot Deck 補完には上記の Cell Hot Deck 以外の方法も存在するが (Andridge and Little(2010))、応用上、最も多く用いられ一般的な枠組みと考えられる Cell Hot Deck について取り扱うことにする。

W を定義関数(指標関数)で Y が観測されたとき、すなわち完全データであったときに 1、それ以外に 0 という値をとるとする。ここで、 $X \in C_t$ を与えたときに Y は (X, W) と独立で同一の分布であると仮定する (cell mean assumption)、この仮定から同一セル内の Y の分布の平均は共通であることがわか

る。また、標本抽出は単純無作為抽出を仮定する。

いま、 $\mu = E[Y]$ 、 $\mu(x) = E[Y|X = x]$ 、 $\mu_t = E[Y|X \in C_t]$ 、 $\sigma_t^2 = V(Y|X \in C_t)$ とする。また、 $j(i)$ を*i*番目の観測値に対する Y の値を補完するのに用いられた観測値の添え字とする(*i*番目の観測値の Y の値が観測されたときは $j(i) = i$ となる)。上記の補完を行った上での Y の標本平均を \bar{Y} とすると

$$\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_{j(i)} = \frac{1}{N} \sum_{i=1}^N W_i(1 + K_{N,i})Y_i \quad (1)$$

と表すことができる。ここで $K_{N,i}$ は不完全なレコードを補完するために*i*番目の観測値が何回用いられたかを表す数である。(Note: $K_{N,i}$ の値は補完がどのように行われるか、例えば同じセルの完全なレコードから無作為に選ぶ方法(random cell hot deck)や、同じセル内で直前に観測された完全なレコードで補完する方法(sequential cell hot deck)によって変わる。)

ここで(1)式および $\sum_{i=1}^N W_i(1 + K_{N,i})\mu(X_i) = \sum_{i=1}^N \mu(X_{j(i)})$ より、

$$\bar{Y} - \mu = \frac{1}{N} \sum_{i=1}^N (Y(X_i) - \mu) + \frac{1}{N} \sum_{i=1}^N W_i(1 + K_{N,i})(Y_i - \mu(X_i)) + \frac{1}{N} \sum_{i=1}^N (\mu(X_{j(i)}) - \mu(X_i)) \quad (2)$$

が成立する。ここで cell mean assumption より、すべての*i*に対して $\mu(X_{j(i)}) - \mu(X_i) = 0$ となるので(2)式の右辺第3項は0となる。いま、セル C_t に含まれる要素の数を N_t として、 $K_{N,i}$ の2次のモーメントの存在($E[K_{N,i}^2] < \infty$)を仮定し、

$$\sigma^2 = E \left[\sum_{t=1}^T \left(\frac{N_t}{N} \right) (\mu_t - \mu)^2 \right] + E \left[\sum_{t=1}^T \left(\frac{N_t}{N} \right) \sigma_t^2 \sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i(1 + K_{N,i})^2 \right] \quad (3)$$

とおき

$$\xi_{N,k} = \begin{cases} \frac{1}{\sigma\sqrt{N}} (Y(X_k) - \mu) & 1 \leq k \leq N \\ \frac{1}{\sigma\sqrt{N}} W_{k-N}(1 + K_{N,k-N})(Y_{k-N} - \mu(X_{k-N})) & N+1 \leq k \leq 2N \end{cases} \quad (4)$$

と定義すると、(2)より

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{N}} = \sum_{k=1}^{2N} \xi_{N,k} \quad (5)$$

と表すことができる。

ここで、 $\mathcal{X}_N = \{X_1, \dots, X_N\}$ と $\omega_N = \{W_1, \dots, W_N\}$ とおき、 σ -加法族の列 $F_{N,k}$ を $F_{N,k} = \sigma\{W_1, \dots, W_k, X_1, \dots, X_k\}$ ($1 \leq k \leq N$)、 $F_{N,k} = \sigma\{\mathcal{X}_N, \omega_N, Y_1, \dots, Y_k\}$ ($N+1 \leq k \leq 2N$)と定義すると、 $F_{N,k}$ は単調非減少な σ -加法族の列(フィルトレーション)となる。ここで $\sigma\{\cdot\}$ は、 $\{\cdot\}$ 内の確率変数を可測にするような最小の σ -加法族である。

このとき、 $\sum_{j=1}^i \xi_{N,i}$ ($1 \leq i \leq 2N$)は $F_{N,i}$ ($1 \leq i \leq 2N$)についてマルチンゲールとなる。よって、マルチンゲール中心極限定理と漸近分散が1となることを示すことができることから

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{N}} \xrightarrow{d} N(0, 1) \quad (6)$$

を得る。ここで \xrightarrow{d} は分布収束を表す。Abadie and Imbens(2012)の得た結果を定理の形でまとめると以下となる(Abadie and Imbens(2012), Theorem2)。

定理 (1) $\{X_1, \dots, X_N\}$ は母集団から無作為に抽出され、(2) $P(W = 1|X \in C_t) > 0, t = 1, \dots, T$ 、(3) $X \in C_t (t = 1, \dots, T)$ を与えたときに Y は (X, W) と独立で、(4) $V(Y) > 0$ で、(5) ある $\delta > 0$ に対して $E[|Y|^{2+\delta}] < \infty$ であるとき、(6)が成立する。

最後に σ^2 の推定方法を考える。最初に Y について欠損値を補完したことを考慮に入れない通常の分散推定量

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_{j(i)} - \bar{Y})^2 \quad (7)$$

を考える。各セルごとの $Y_{j(i)}$ の平均を $\bar{Y}_t = \sum_{i=1}^N Y_{j(i)} 1_{\{X_i \in C_t\}} / N_t$ とおく。各セルを分散分析における個別母集団(級)と考えると上の分散は級間分散(セルの間の平均値 \bar{Y}_t の分散)と級内分散(各セル内の $Y_{j(i)}$ の値の分散)に分解されることから

$$\left| \hat{\sigma}^2 - \sum_{t=1}^T \left(\frac{N_t}{N}\right) (\mu_t - \mu)^2 - \sum_{t=1}^T \left(\frac{N_t}{N}\right) \sigma_t^2 \right| \xrightarrow{p} 0 \quad (8)$$

を得る。ここで \xrightarrow{p} は確率収束を表す。 σ^2 の定義式((3)式)と(8)式左辺の絶対値内の第2項、第3項を比較することより $\hat{\sigma}^2$ に $\sum_{t=1}^T \left(\frac{N_t}{N}\right) \sigma_t^2 \left(\sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (1 + K_{N,i})^2 - 1\right)$ と漸的に同等な項を加えることで σ^2 の一致推定量が構成されることが予想される。

ここで、 $\sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (1 + K_{N,i}) = N_t$ より

$$\frac{1}{N_t} \sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (1 + K_{N,i})^2 = 1 + \frac{1}{N_t} \sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (K_{N,i}^2 + K_{N,i}) \quad (9)$$

が成り立つことから σ^2 の推定量として、

$$\begin{aligned} \hat{\sigma}_{\text{adj}}^2 &= \hat{\sigma}^2 + \frac{1}{N} \sum_{t=1}^T \left(\sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (K_{N,i}^2 + K_{N,i}) \right) \hat{\sigma}_t^2 \\ &= \hat{\sigma}^2 + \sum_{t=1}^T \left(\frac{N_t}{N} \right) \left(\frac{1}{N_t} \sum_{i=1}^N 1_{\{X_i \in C_t\}} W_i (K_{N,i}^2 + K_{N,i}) \right) \hat{\sigma}_t^2 \end{aligned} \quad (10)$$

が提案される。ここで、 $\hat{\sigma}_t^2$ はセル C_t 内の観測されたレコードによる Y の分散の推定量である。

3. まとめ

本稿ではAbadie and Imbens(2012)によって示された Cell Hot Deck 補完を行った際の標本平均の漸近正規性の概略を記した。彼らの結果は類似の補完方法に容易に拡張が可能な汎用性の高いものであり、今後、彼らの結果をより一般的な標本抽出の枠組みや標本平均以外の統計量に拡張することで、応用上の可能性がさらに広がっていくものと期待される。

参考文献

- Abadie, A and G.W. Imbens(2012). A Martingale Representation for Matching Estimators. *Journal of the American Statistical Association*, **107(498)**, 833-843.
- Andridge, R.R. and R.J.A. Little(2010). A Review of Hot Deck Imputation for Survey Non-response. *International Statistical Review*, **78(1)**, 40-64.