

サンプルセレーションバイアスとその補正 (統計)

1. Sample Selection Bias と対処法

ただし観測されない

(1) 設定: (貸金 W_i^*) =
$$\begin{cases} \beta_0 + \beta_1 (\text{教育年数 } X_i) + \varepsilon_i & \text{if } S_i = 1 \\ (\beta_0 + \beta_1 (\text{教育年数 } X_i) + \varepsilon_i) & \text{if } S_i = 0 \end{cases} \quad (*)$$

・ (働いているか S_i) =
$$\begin{cases} 1 & \text{if } (\text{貸金 } W_i^*) - (\text{留保貸金 } W_i) = \gamma_0 + \gamma_1 (\text{貯金 } Z_i) + u_i \\ 0 & \text{if } (\text{貸金 } W_i^*) - (\text{留保貸金 } W_i) = \gamma_0 + \gamma_1 (\text{貯金 } Z_i) + u_i < 0 \end{cases}$$

ただし、 $S_i = 1$ のときのみ W_i^* を観測可能、

$S_i = 0$: W_i^* を観測測りできない。

また、 $\begin{bmatrix} \varepsilon_i \\ u_i \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_\varepsilon^2 & \rho_{\varepsilon u} \\ \rho_{\varepsilon u} & 1 \end{bmatrix}\right)$ とする。

※もちろん $X_i \equiv Z_i$ であっても問題ない。

(2) 操作: 観測測りされたデータ、今回は、 $i = 1, 3$ に限定して β を推定すると、

i	W_i^*	S	X_i	Z_i
1	9	1	16	2
2	8	0	12	1
3	10	1	20	0
4	8	0	14	1

推定値 $\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} 5 \\ 0.25 \end{pmatrix}$ を得たが、...

(3) 問題点: Selection bias (バイアス形成過程に注目)

推定量 $\hat{\beta}$ ($S_i = 1$ に制限した data) は bias か? $E[\hat{\beta}(\text{observed data})] \neq \beta$

どのようなメカニズムでどのようなバイアスが発生しているのか?
↑
大きさの

$$E[W_i | S_i = 1, X_i] = E[W_i^* | S_i = 1, X_i] = E[W_i^* | \gamma_0 + \gamma_1 Z_i + u_i > 0, X_i]$$

$S_i = 1$ に制限したデータを所有した貸金 W_i の期待値、

$$= E[\beta_0 + \beta_1 X_i + \varepsilon_i | u_i > -\gamma_0 - \gamma_1 Z_i, X_i]$$

$$= \beta_0 + \beta_1 X_i + E[\varepsilon_i | u_i > -\gamma_0 - \gamma_1 Z_i, X_i]$$

$$= \beta_0 + \beta_1 X_i + \rho_{\varepsilon u} E[u_i | u_i > -\gamma_0 - \gamma_1 Z_i, X_i]$$

(ただし設定より $u_i \sim N(0, 1)$) (*)

ここで一般に、 $Z \sim N(0, 1)$ のとき、 $E[Z | Z > c]$ の値を求めよ。

$$\begin{aligned}
 E[Z | Z > c] &= \int_{-\infty}^{\infty} z \cdot f_{Z|I}(z | I_{(Z > c)} = 1) dz \quad (\text{conditional Expectation の定義}) \\
 &= \int_{-\infty}^{\infty} z \cdot \frac{f_{Z,I}(z, I_{(Z > c)} = 1)}{P(I_{(Z > c)} = 1)} dz \quad (\text{conditional Prob の定義}) \\
 &= \int_{-\infty}^{\infty} z \cdot \frac{f_{Z,I}(z, Z > c)}{P(\{Z > c\})} dz \\
 &= \frac{1}{P(\{Z > c\})} \int_c^{\infty} z f_Z(z) dz \\
 &= \frac{1}{P(\{Z < -c\})} \int_c^{\infty} z f_Z(z) dz \quad (\text{①より対称性から } Z \sim N(0, 1) \text{ より } Z \text{ と } -Z \text{ が同じ分布を持つ})
 \end{aligned}$$

部分積分公式： $f \cdot g = f'g + f \cdot g'$ に $f(x) = -1, g(x) = \exp(-\frac{x^2}{2})$ と考えよ。

$$\begin{aligned}
 \int_c^{\infty} \underbrace{(-1)}_f \underbrace{(-z) \exp(-\frac{z^2}{2})}_{g'} dz &= \left[\underbrace{-\exp(-\frac{z^2}{2})}_f \underbrace{1}_g \right]_c^{\infty} - \int_c^{\infty} \underbrace{0}_{f'} \underbrace{\exp(-\frac{z^2}{2})}_g dz \\
 &= \left[\exp(-\frac{z^2}{2}) \right]_c^{\infty} = \exp(-\frac{c^2}{2}) \quad \text{②}
 \end{aligned}$$

①, ②より一般に $Z \sim N(0, 1)$ のとき、

$$E[Z | Z > c] = \frac{1}{\Phi(-c)} \cdot \frac{1}{\sqrt{2\pi}} \exp(-\frac{c^2}{2}) = \frac{\phi(-c)}{\Phi(-c)} \quad \text{を得る。}$$

これを(*)に適用すると、

$$\begin{aligned}
 E[W_i | S_i = 1, X_i] &= \beta_0 + \beta_1 X_i + \rho_{E,u} E[Z_i | Z_i > -\gamma_0 - \gamma_1 Z_i] \\
 &= \beta_0 + \beta_1 X_i + \rho_{E,u} \underbrace{\frac{\phi(\gamma_0 + \gamma_1 Z_i)}{\Phi(\gamma_0 + \gamma_1 Z_i)}}_{\text{③}} \quad \text{を得る。}
 \end{aligned}$$

③: 逆 Mills 比。

$$\text{よって } E[W_i | S_i = 1, X_i] = \underbrace{\beta_0 + \beta_1 X_i}_{\text{真値}} + \underbrace{\rho_{E,u} \frac{\phi(\gamma_0 + \gamma_1 Z_i)}{\Phi(\gamma_0 + \gamma_1 Z_i)}}_{\text{Bias}}$$

* $\rho_{E,u} = 0$ ならば Sample Selection Bias はない。

<解釈>

Inv Mill Ratio $\frac{\phi(\gamma_0 + \gamma_1 z_i)}{\Phi(\gamma_0 + \gamma_1 z_i)}$ の図による解釈はRコード参照IMRが大きい $\Leftrightarrow \gamma_0 + \gamma_1 z_i$ が小さい (C + Rコード)
(小さい) (大きい) \Leftrightarrow cut-offで切り捨てられることが起きやすくなり、
(にくくなる)Cond Expe $E[z_i | z_i > \gamma_0 + \gamma_1 z_i]$ も大きくなる
(小さくなる) \Leftrightarrow Bias も大きく (小さく) なる。(4) 解決策: どうやって Sample Selection Bias をとり除いて、
Unbiased Estimator $\hat{\beta}(\text{Data})$ をつくるか?

Heckman's Two-step sample selection correction:

(I) 全ての観測値を用いて、まず、プロビットモデル

$$P(\{s_i = 1 | z_i\}) = P(\gamma_0 + \gamma_1 z_i + u_i > 0 | z_i) \quad (u_i \sim N(0, 1))$$

の推定量 $\hat{\gamma}_{MLE}(\text{DATA})$ を得る。

(例1)

z	s	z
1	1	2
2	0	1
3	1	0
4	0	1

が ML は求め

(例2)

これをもとに、Inverse Mills Ratio $\hat{\lambda}_i$ の実現値を計算する。(例1) $\hat{\gamma}_{MLE}$ の実現値: $\begin{bmatrix} \hat{c}_{1MLE} \\ \hat{c}_{2MLE} \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$ とすると、例1は上の表で97%

$$\hat{\lambda}_1 = \frac{\phi(0.5 + 2 \cdot 0.5)}{\Phi(0.5 + 2 \cdot 0.5)} \approx 0.1387898$$

推定値

$$\hat{\lambda}_3 = \frac{\phi(0.5 + 0 \cdot 0.5)}{\Phi(0.5 + 0 \cdot 0.5)} \approx 0.5091604 \quad (\text{例2終})$$

- (II) Selected Sample (つまり、 $S_i = 1$ となっている国票データ) を用いて、
OLS wage on X_i and $\hat{\lambda}_i$.

(1991) λ_i	S_i	W_i^*	X_i	Z_i	$\hat{\lambda}_i$
1	1	9	16	2	0.1388
2	0	•	12	1	0.2896
3	1	10	20	0	0.5092
4	0	•	14	1	0.2896
⋮	⋮	⋮	⋮	⋮	⋮

← 4%

変数 $X, \hat{\lambda}$ を W^* に OLS. Unbiased Estimator $\hat{\beta}$ を得る.

(1991 年).

これより、 $W_{\text{miss}} = \underbrace{b_0 + b_1 X + b_2 (-\hat{\lambda})}_{\text{補完}}$ で求めると bias なくなる.

(5) 補

働いていない人が、
仮に働いたときの賃金

⇒ 一般化

(I) 全てのデータを用いて、

$$P(\{S_i = 1\} | Z_i) = P(\{S_i = 1\} | \gamma_0 + \gamma_1 Z_i + \varepsilon_i) \quad (\varepsilon_i \sim N(0, 1) \text{ 独立})$$

$\hat{\gamma}_{\text{MLE}}(\text{All Data})$ を求める。これは MLE の性質が一貫性を満たす。

$$\text{よって } \hat{\lambda}_i = \frac{\phi(Z_i' \hat{\gamma})}{\Phi(Z_i' \hat{\gamma})} \quad \text{for } S_i = 1$$

も一貫性をもち推定量として、定めることが出来る。

実際のデータを用いて、推定値 $\hat{\lambda}_i$ を計算する。

(II) ここで (3) の Selection Bias 形成過程を、まづ (1) の設定のもとでは、

$$\underbrace{E[W_{S_i=1} | S_i=1, X_i]}_{\text{今実際には働いている人が、}} = \beta_0 + \beta_1 X_i + \underbrace{\beta_2 \frac{\phi(\gamma_0 + \gamma_1 Z_i)}{\Phi(\gamma_0 + \gamma_1 Z_i)}}_{\equiv P_{\varepsilon u}}$$

もし、条件 (4) のもとで就業したら得る賃金の期待値

$$\underbrace{E[W_{S_i=1} | S_i=0, X_i]}_{\text{実際には働いていない人から、}} = E[\beta_0 + \beta_1 X_i + \varepsilon_i | S_i=0, X_i]$$

もし条件 (4) のもとで就業したら得る賃金の期待値

$$= \beta_0 + \beta_1 X_i + E[\varepsilon_i | \gamma_0 + \gamma_1 Z_i + u_i < 0] \quad (u_i \sim N(0,1))$$

$$= \beta_0 + \beta_1 X_i + P_{\varepsilon u} E[u_i | u_i < -\gamma_0 - \gamma_1 Z_i]$$

$$= \beta_0 + \beta_1 X_i + P_{\varepsilon u} \frac{\int_{-\infty}^{-\gamma_0 - \gamma_1 Z_i} u_i f_u(u_i) du_i}{P(\{u_i < -\gamma_0 - \gamma_1 Z_i\})}$$

$$= \beta_0 + \beta_1 X_i + P_{\varepsilon u} \frac{\int_{-\infty}^{-\gamma_0 - \gamma_1 Z_i} u_i f(u_i) du_i}{1 - P(\{u_i < \gamma_0 + \gamma_1 Z_i\})}$$

ここで、

$$\int_{-\infty}^c \underbrace{(-1)}_f \underbrace{(-u)}_{f'} \exp\left(-\frac{u^2}{2}\right) du = \underbrace{\left[-\exp\left(-\frac{u^2}{2}\right)\right]}_f^c = \underbrace{\left[\exp\left(-\frac{u^2}{2}\right)\right]}_g^c = -\exp\left(-\frac{c^2}{2}\right) \quad \text{よって、上式は、}$$

$$E[W_{S_i=1} | S_i=0, X_i] = \beta_0 + \beta_1 X_i + P_{\varepsilon u} \frac{-\phi(\gamma_0 + \gamma_1 Z_i)}{1 - \Phi(\gamma_0 + \gamma_1 Z_i)} \\ \equiv \beta_0 + \beta_1 X_i + \beta_2 \frac{-\phi(\gamma_0 + \gamma_1 Z_i)}{1 - \Phi(\gamma_0 + \gamma_1 Z_i)}$$

よって $S_i=1$ となるデータに限定して OLS wage on X_i , λ_i をすると、
真のモデル (4) のパラメーター β_1 を推定できる。

これをを用いて $E[W_{S_i=1} | S_i=0, X_i]$ もバイアスなく推定できる。