

第7回 重回帰分析 (6.1–6.3)

村澤 康友

2024年5月21日

今日のポイント

1. 結果を処置ダミーに回帰すれば平均処置効果が求まる。結果と処置の両方に影響する共変量が存在する場合は重回帰分析で共変量調整を行う。
2. 自由度修正済み決定係数は $\bar{R}^2 := 1 - [\text{RSS}/(n - k)]/[\text{TSS}/(n - 1)]$ 。2つの説明変数の相関が極めて高く、それらの回帰係数の OLS 推定値が不安定になる問題を多重共線性という。
3. 説明変数の欠落によって生じる OLS 推定量の偏りを欠落変数バイアスという。
4. ある説明変数の偏回帰係数の OLS 推定値は、その説明変数を残りの説明変数に回帰した OLS 残差に被説明変数を単回帰しても求まる。
5. 誤差項が無相関で分散が均一な線形回帰モデルを古典的線形回帰モデルという。被説明変数の線形関数で表される推定量を線形推定量という。不偏な線形推定量を線形不偏推定量という。分散が最小となる線形不偏推定量を最良線形不偏推定量 (BLUE) という。古典的線形回帰モデルの回帰係数の OLS 推定量は BLUE (ガウス=マルコフ定理)。

目次

1	重回帰分析	1
1.1	ダミー変数 (pp. 53, 166)	1

1.2	共変量調整 (p. 131)	2
1.3	MM (= OLS) 推定量 (p. 133)	2
1.4	自由度修正済み決定係数 (p. 135)	2
1.5	多重共線性 (p. 138)	3
2	欠落変数バイアス (p. 139)	3
3	偏回帰	3
3.1	重回帰モデル	3
3.2	MM (= OLS) 推定量	3
3.3	OLS 残差	3
3.4	偏回帰 (p. 158)	4
4	OLS 推定量の性質	5
4.1	古典的線形回帰モデル (p. 146)	5
4.2	MM (= OLS) 推定量	5
4.3	ガウス=マルコフ定理 (p. 146)	5
5	今日のキーワード	6
6	次回までの準備	6

1 重回帰分析

1.1 ダミー変数 (pp. 53, 166)

ある条件に該当するか否かの2値変数はベルヌーイ確率変数で表せる。すなわち

$$D := \begin{cases} 1 & \text{該当} \\ 0 & \text{非該当} \end{cases}$$

処置の有無を D 、結果を Y とする。

定義 1. ある条件に該当するなら1, しないなら0とした変数を**ダミー変数**という。

定義 2. 処置群と対照群に対する効果の差を**処置**

(介入) 効果という。

定義 3. 処置効果の平均を平均処置効果 (*Average Treatment Effect, ATE*) という。

注 1. 処置群と対照群の母平均の差に等しい。すなわち Y に対する D の ATE は

$$ATE = E(Y|D = 1) - E(Y|D = 0)$$

実験データなら母平均の差の推測のみ (2 標本問題)。

注 2. $\mu_0 := E(Y|D = 0)$, $\mu_1 := E(Y|D = 1)$ とすると

$$\begin{aligned} E(Y|D) &= D\mu_1 + (1 - D)\mu_0 \\ &= \mu_0 + (\mu_1 - \mu_0)D \\ &= \mu_0 + ATE \cdot D \end{aligned}$$

これは単回帰モデル。したがって 2 標本問題は単回帰分析で実行できる。また k 標本問題は重回帰分析で実行できる (=分散分析)。

1.2 共変量調整 (p. 131)

実験データと異なり、観察データでは D を直接コントロールできない。 Y と D の両方に影響する変数 X が存在する場合、 Y の (D, X) 上への重回帰モデルを考える。

$$E(Y|D, X) = \alpha + ATE \cdot D + \beta X$$

定義 4. 関心の対象外の説明変数を**共変量**という。

定義 5. 分析の際に共変量の影響を調整することを**共変量調整**という。

1.3 MM (= OLS) 推定量 (p. 133)

次の重回帰モデルを考える。

$$E(Y|X_1, \dots, X_k) = \alpha + \beta_1 X_1 + \dots + \beta_k X_k$$

回帰の誤差項は $U := Y - E(Y|X_1, \dots, X_k)$ 。

定理 1.

$$E(U|X_1, \dots, X_k) = 0$$

証明. 復習テスト。 □

定理 2.

$$E(U) = E(X_1 U) = \dots = E(X_k U) = 0$$

証明. 復習テスト。 □

注 3. $U = Y - \alpha - \beta_1 X_1 - \dots - \beta_k X_k$ を代入すると

$$\begin{aligned} E(Y - \alpha - \beta_1 X_1 - \dots - \beta_k X_k) &= 0 \\ E(X_1(Y - \alpha - \beta_1 X_1 - \dots - \beta_k X_k)) &= 0 \\ &\vdots \\ E(X_k(Y - \alpha - \beta_1 X_1 - \dots - \beta_k X_k)) &= 0 \end{aligned}$$

この連立方程式が解けるなら、 $(\alpha, \beta_1, \dots, \beta_k)$ は MM 法で推定できる (OLS と同値)。

1.4 自由度修正済み決定係数 (p. 135)

決定係数は

$$R^2 = 1 - \frac{RSS}{TSS}$$

ただし

$$\begin{aligned} TSS &:= \sum_{i=1}^n (y_i - \bar{y})^2 \\ RSS &:= \sum_{i=1}^n e_i^2 \end{aligned}$$

推定する係数の数 (=定数項を含む説明変数の数) を k とすると、RSS は k の減少関数。また一般に $k \geq n$ なら RSS は 0。したがって R^2 は説明変数の選択に役立たない。

定義 6. 自由度修正済み決定係数は

$$\bar{R}^2 := 1 - \frac{RSS/(n - k)}{TSS/(n - 1)}$$

注 4. 無作為標本なら

$$E\left(\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2\right) = \text{var}(y_i)$$

$\text{var}(u_i|x_i) = \text{var}(u_i)$ なら

$$E\left(\frac{1}{n-k} \sum_{i=1}^n e_i^2\right) = \text{var}(u_i)$$

したがって \bar{R}^2 は $1 - \text{var}(u_i) / \text{var}(y_i)$ の推定量 (値) となっている。ただし

$$\begin{aligned} E(\bar{R}^2) &= 1 - E\left(\frac{[1/(n-k)] \sum_{i=1}^n e_i^2}{[1/(n-1)] \sum_{i=1}^n (y_i - \bar{y})^2}\right) \\ &\neq 1 - \frac{E([1/(n-k)] \sum_{i=1}^n e_i^2)}{E([1/(n-1)] \sum_{i=1}^n (y_i - \bar{y})^2)} \\ &= 1 - \frac{\text{var}(u_i)}{\text{var}(y_i)} \end{aligned}$$

1.5 多重共線性 (p. 138)

次の重回帰モデルを考える。

$$E(Y|X, Z) = \alpha + \beta X + \gamma Z$$

ここで $X = Z$ とすると、任意の w について

$$E(Y|X, Z) = \alpha + w(\beta + \gamma)X + (1-w)(\beta + \gamma)Z$$

すなわち X, Z の係数は一意に定まらない。 $Z = a + bX$ でも同様。

より一般的に、次の重回帰モデルを考える。

$$E(Y|X_1, \dots, X_k) = \alpha + \beta_1 X_1 + \dots + \beta_k X_k$$

ここで $X_1 = a + b_2 X_2 + \dots + b_k X_k$ の場合も係数は一意に定まらない。

定義 7. 実質的に同じ説明変数が 2 つあり、それらの回帰係数が定まらない問題を **完全な多重共線性** という。

定義 8. 2 つの説明変数の相関が極めて高く、それらの回帰係数の OLS 推定値が不安定になる問題を (準) **多重共線性** という。

2 欠落変数バイアス (p. 139)

次の重回帰モデルを考える。

$$E(Y|X, Z) = \alpha + \beta X + \gamma Z$$

ここで $E(Z|X) = a + bX$ とし、 Z を説明変数に含めない、繰り返し期待値の法則より

$$\begin{aligned} E(Y|X) &= E(E(Y|X, Z)|X) \\ &= E(\alpha + \beta X + \gamma Z|X) \\ &= \alpha + \beta X + \gamma E(Z|X) \\ &= \alpha + \beta X + \gamma(a + bX) \\ &= \alpha + \gamma a + (\beta + \gamma b)X \end{aligned}$$

すなわち X の回帰係数は β でなく $\beta + \gamma b$ となる。

定義 9. 説明変数の欠落によって生じる OLS 推定量の偏りを **欠落変数バイアス** という。

3 偏回帰

3.1 重回帰モデル

$(1+k)$ 変量データを $\{(y_i, x_{i,1}, \dots, x_{i,k})\}_{i=1}^n$ とする。 y_i の $(x_{i,1}, \dots, x_{i,k})$ 上への重回帰モデルは

$$E(y_i|x_{i,1}, \dots, x_{i,k}) = \beta_1 x_{i,1} + \dots + \beta_k x_{i,k}$$

β_1 の推定を考える (β_2, \dots, β_k に関心はない)。

3.2 MM (= OLS) 推定量

繰り返し期待値の法則より

$$E(x_{i,1}(y_i - \beta_1 x_{i,1} - \dots - \beta_k x_{i,k})) = 0$$

⋮

$$E(x_{i,k}(y_i - \beta_1 x_{i,1} - \dots - \beta_k x_{i,k})) = 0$$

$(\beta_1, \dots, \beta_k)$ の MM (= OLS) 推定量を (b_1, \dots, b_k) とすると

$$\frac{1}{n} \sum_{i=1}^n x_{i,1}(y_i - b_1 x_{i,1} - \dots - b_k x_{i,k}) = 0$$

⋮

$$\frac{1}{n} \sum_{i=1}^n x_{i,k}(y_i - b_1 x_{i,1} - \dots - b_k x_{i,k}) = 0$$

3.3 OLS 残差

y_i の回帰予測は

$$\hat{y}_i := b_1 x_{i,1} + \dots + b_k x_{i,k}$$

OLS 残差は

$$\begin{aligned} e_i &:= y_i - \hat{y}_i \\ &= y_i - b_1 x_{i,1} - \dots - b_k x_{i,k} \end{aligned}$$

定理 3.

$$\sum_{i=1}^n x_{i,1} e_i = \dots = \sum_{i=1}^n x_{i,k} e_i = 0$$

証明. 復習テスト. □

系 1.

$$\sum_{i=1}^n \hat{y}_i e_i = 0$$

証明. 変形すると

$$\begin{aligned} \sum_{i=1}^n \hat{y}_i e_i &= \sum_{i=1}^n (b_1 x_{i,1} + \cdots + b_k x_{i,k}) e_i \\ &= b_1 \sum_{i=1}^n x_{i,1} e_i + \cdots + b_k \sum_{i=1}^n x_{i,k} e_i \end{aligned}$$

前定理より各項は 0. \square

3.4 偏回帰 (p. 158)

$x_{i,1}$ の $(x_{i,2}, \dots, x_{i,k})$ 上への重回帰モデルを考える. すなわち

$$E(x_{i,1} | x_{i,2}, \dots, x_{i,k}) = \gamma_2 x_{i,2} + \cdots + \gamma_k x_{i,k}$$

繰り返し期待値の法則より

$$E(x_{i,2}(x_{i,1} - \gamma_2 x_{i,2} - \cdots - \gamma_k x_{i,k})) = 0$$

\vdots

$$E(x_{i,k}(x_{i,1} - \gamma_2 x_{i,2} - \cdots - \gamma_k x_{i,k})) = 0$$

$(\gamma_2, \dots, \gamma_k)$ の MM (= OLS) 推定量を (c_2, \dots, c_k) とすると

$$\frac{1}{n} \sum_{i=1}^n x_{i,2}(x_{i,1} - c_2 x_{i,2} - \cdots - c_k x_{i,k}) = 0$$

\vdots

$$\frac{1}{n} \sum_{i=1}^n x_{i,k}(x_{i,1} - c_2 x_{i,2} - \cdots - c_k x_{i,k}) = 0$$

$x_{i,1}$ の回帰予測は

$$\hat{x}_{i,1} := c_2 x_{i,2} + \cdots + c_k x_{i,k}$$

OLS 残差は

$$\begin{aligned} x_{i,1}^* &:= x_{i,1} - \hat{x}_{i,1} \\ &= x_{i,1} - c_2 x_{i,2} - \cdots - c_k x_{i,k} \end{aligned}$$

OLS 残差の性質より

$$\sum_{i=1}^n x_{i,2} x_{i,1}^* = \cdots = \sum_{i=1}^n x_{i,k} x_{i,1}^* = 0$$

かつ

$$\sum_{i=1}^n \hat{x}_{i,1} x_{i,1}^* = 0$$

補題 1.

$$\sum_{i=1}^n \hat{x}_{i,1} e_i = 0$$

証明. 変形すると

$$\begin{aligned} \sum_{i=1}^n \hat{x}_{i,1} e_i &= \sum_{i=1}^n (c_2 x_{i,2} + \cdots + c_k x_{i,k}) e_i \\ &= c_2 \sum_{i=1}^n x_{i,2} e_i + \cdots + c_k \sum_{i=1}^n x_{i,k} e_i \end{aligned}$$

前定理より各項は 0. \square

定理 4 (偏回帰).

$$b_1 = \frac{\sum_{i=1}^n x_{i,1}^* y_i}{\sum_{i=1}^n x_{i,1}^{*2}}$$

証明. 補題より

$$\begin{aligned} \sum_{i=1}^n x_{i,1} e_i &= \sum_{i=1}^n (\hat{x}_{i,1} + x_{i,1}^*) e_i \\ &= \sum_{i=1}^n \hat{x}_{i,1} e_i + \sum_{i=1}^n x_{i,1}^* e_i \\ &= \sum_{i=1}^n x_{i,1}^* e_i \\ &= \sum_{i=1}^n x_{i,1}^* (y_i - b_1 x_{i,1} - \cdots - b_k x_{i,k}) \\ &= \sum_{i=1}^n x_{i,1}^* y_i - b_1 \sum_{i=1}^n x_{i,1}^* x_{i,1} \\ &\quad - b_2 \sum_{i=1}^n x_{i,1}^* x_{i,2} - \cdots - b_k \sum_{i=1}^n x_{i,1}^* x_{i,k} \\ &= \sum_{i=1}^n x_{i,1}^* y_i - b_1 \sum_{i=1}^n x_{i,1}^* x_{i,1} \\ &= \sum_{i=1}^n x_{i,1}^* y_i - b_1 \sum_{i=1}^n x_{i,1}^* (\hat{x}_{i,1} + x_{i,1}^*) \\ &= \sum_{i=1}^n x_{i,1}^* y_i - b_1 \sum_{i=1}^n x_{i,1}^* \hat{x}_{i,1} - b_1 \sum_{i=1}^n x_{i,1}^{*2} \\ &= \sum_{i=1}^n x_{i,1}^* y_i - b_1 \sum_{i=1}^n x_{i,1}^{*2} \end{aligned}$$

左辺 = 0 より b_1 について解けば結果が得られる. \square

注 5. 定理より β_1 の OLS 推定量 b_1 は以下の手順でも求まる.

1. $x_{i,1}$ を $(x_{i,2}, \dots, x_{i,k})$ 上へ回帰し, OLS 残差 $x_{i,1}^*$ を求める.
2. y_i を $x_{i,1}^*$ 上へ回帰.

したがって b_1 は, $(x_{i,2}, \dots, x_{i,k})$ と相関する部分を取り除いた上での y_i と $x_{i,1}$ の関係を表す.

4 OLS 推定量の性質

4.1 古典的線形回帰モデル (p. 146)

$(1+k)$ 変量データを $((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n))$ とする. ただし $\mathbf{x}_i := (x_{i,1}, \dots, x_{i,k})'$. $x_{i,1} := 1$ を定数項とすると, y_i の \mathbf{x}_i 上への重回帰モデルは

$$\begin{aligned} E(y_i | \mathbf{x}_i) &= \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} \\ &= \boldsymbol{\beta}' \mathbf{x}_i \\ &= \mathbf{x}_i' \boldsymbol{\beta} \end{aligned}$$

または

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + u_i \\ E(u_i | \mathbf{x}_i) &= 0 \end{aligned}$$

すなわち重回帰モデルをベクトルで表記すれば, 定数項なしの単回帰モデルと同様に扱える.

定義 10. $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ を所与として u_1, \dots, u_n が無相関で分散が均一な線形回帰モデルを **古典的線形回帰モデル** という.

注 6. すなわち

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + u_i \\ E(u_i | \mathbf{x}_1, \dots, \mathbf{x}_n) &= 0 \\ \text{var}(u_i | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \sigma^2 \\ \text{cov}(u_i, u_j | \mathbf{x}_1, \dots, \mathbf{x}_n) &= 0 \quad \text{for } i \neq j \end{aligned}$$

4.2 MM (= OLS) 推定量

繰り返し期待値の法則より

$$E(\mathbf{x}_i u_i) = \mathbf{0}$$

$u_i = y_i - \mathbf{x}_i' \boldsymbol{\beta}$ を代入すると

$$E(\mathbf{x}_i (y_i - \mathbf{x}_i' \boldsymbol{\beta})) = \mathbf{0}$$

$\boldsymbol{\beta}$ の MM (= OLS) 推定量を \mathbf{b} とすると

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i' \mathbf{b}) = \mathbf{0}$$

すなわち

$$\sum_{i=1}^n \mathbf{x}_i y_i = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \mathbf{b}$$

逆行列を用いて連立方程式を解くと

$$\mathbf{b} = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \sum_{i=1}^n \mathbf{x}_i y_i$$

定理 5.

$$E(\mathbf{b} | \mathbf{x}_1, \dots, \mathbf{x}_n) = \boldsymbol{\beta}$$

証明. 省略 (定数項のない単回帰モデルと同じ). □

系 2.

$$E(\mathbf{b}) = \boldsymbol{\beta}$$

証明. 省略 (繰り返し期待値の法則). □

定理 6. 古典的線形回帰モデルなら

$$\text{var}(\mathbf{b} | \mathbf{x}_1, \dots, \mathbf{x}_n) = \sigma^2 \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1}$$

証明. 省略 (定数項のない単回帰モデルと同じ). □

4.3 ガウス=マルコフ定理 (p. 146)

定義 11. 被説明変数の線形関数で表される推定量を **線形推定量** という.

注 7. \mathbf{b} は y_1, \dots, y_n の線形関数だから線形推定量.

定義 12. 不偏な線形推定量を **線形不偏推定量** という.

注 8. $E(\mathbf{b}) = \boldsymbol{\beta}$ より \mathbf{b} は線形不偏推定量.

定義 13. 分散が最小となる線形不偏推定量を **最良線形不偏推定量 (Best Linear Unbiased Estimator, BLUE)** という.

定理 7 (ガウス=マルコフ定理). 古典的線形回帰モデルの回帰係数の OLS 推定量は BLUE.

証明. 省略 (行列を使うと簡単). □

5 今日のキーワード

ダミー変数, 処置 (介入) 効果, 平均処置効果 (ATE), 共変量, 共変量調整, 自由度修正済み決定係数, 完全な多重共線性, (準) 多重共線性, 欠落変数バイアス, 偏回帰, 古典的線形回帰モデル, 線形推定量, 線形不偏推定量, 最良線形不偏推定量 (BLUE), ガウス=マルコフ定理

6 次回までの準備

復習 教科書第 6 章 1-3 節, 復習テスト 7

予習 教科書第 6 章 4-5 節