

線形回帰モデル (多重共線性, Ridge 回帰)

1 多重共線性

次の線形回帰モデルを考える。

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (1)$$

ここで, $p < n$ として,

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

とする。また, $\boldsymbol{\varepsilon}$ に関して,

$$\text{条件 1: } E[\boldsymbol{\varepsilon}] = \mathbf{0}$$

$$\text{条件 2: } V[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}_n$$

を仮定する。ただし, \mathbf{I}_n は $n \times n$ 単位行列である。このとき, \mathbf{Y} を目的変数, \mathbf{X} を説明変数とし, $\boldsymbol{\beta}$ の推測問題を考える。 $\boldsymbol{\beta}$ の最小二乗推定量 $\hat{\boldsymbol{\beta}}$ は, 次式で与えられる。

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

(1) 式の線形回帰モデルの $n \times (p+1)$ 行列 \mathbf{X} について, $(\mathbf{X}^\top \mathbf{X})^{-1}$ が正則であるために, $\text{rank}(\mathbf{X}) = \text{rank}((\mathbf{X}^\top \mathbf{X})^{-1}) = p+1$ であることが仮定されていた。この仮定を満たすためには, 観測値の個数 n が変数の数 p よりも大きいことに加え, \mathbf{X} の $p+1$ 個の列が 1 次独立であることが必要になる。この仮定が満たされない場合, $\text{rank}((\mathbf{X}^\top \mathbf{X})^{-1}) < p+1$ となり, その逆行列は存在しないことになる。つまり, 最小二乗推定量は定義されない。一般化逆行列を用いて, $\text{rank}((\mathbf{X}^\top \mathbf{X})^{-1}) < p+1$ の場合でも最小二乗推定量の定義を拡張することは可能であるが, 便宜的に拡張された最小二乗推定量の意味づけは, 必ずしも明確ではない。

行列 $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$ の階数が $p+1$ 未満であるということは,

$$a_0 \mathbf{1} + a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \cdots + a_p \mathbf{x}_p = \mathbf{0}$$

を満たす, 0 ではない定数 $a_0, a_1, a_2, \dots, a_p$ が存在することと同値である。つまり, $\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ は一次従属であることと同値である。したがって, p 個の変数に共線関係が存在することを意味する。

簡単のために, $\mathbf{x}_1 = a\mathbf{x}_2$ という関係が成り立っていると仮定して, 次の線形回帰モデルを考える。

$$\mathbf{Y} = \beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2 + \boldsymbol{\varepsilon} \quad (2)$$

(2) 式の線形回帰モデルは, $\mathbf{x}_1 = a\mathbf{x}_2$ の仮定のもとで, 次の線形回帰モデルと同等であり, 与えられた観測値に基づいて, 一方を他方から識別することは不可能である。

$$\mathbf{Y} = (\beta_1 a + \beta_2) \mathbf{x}_2 + \boldsymbol{\varepsilon}$$

\mathbf{x}_1 と \mathbf{x}_2 の各々が、 \mathbf{Y} の変動をどの程度説明できるかに関心があったとしても、 $\mathbf{x}_1 = a\mathbf{x}_2$ の共線関係のために、どうしようもない。

一般には、説明変数間に厳密な共線関係が存在することはほぼないと考えられる。実際には、共線関係が近似的に生じる場合である。つまり、すべてが 0 ではない定数 $a_0, a_1, a_2, \dots, a_p$ に対して、次式となる場合である。

$$a_0 \mathbf{1} + a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_p \mathbf{x}_p \doteq \mathbf{0}$$

このような場合、 $(\mathbf{X}^\top \mathbf{X})^{-1}$ が正則であるために、 β を最小二乗推定することは可能である。しかし、推定上の困難が生じることになる。(2) 式の線形回帰モデルのもとで、この問題について以下に述べていく。

簡単のために、 $\mathbf{X}^\top \mathbf{X}$ と $\mathbf{X}^\top \mathbf{Y}$ を次のようにする。

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix}, \quad \mathbf{X}^\top \mathbf{Y} = \begin{pmatrix} m_{1y} \\ m_{2y} \end{pmatrix}$$

(2) 式の線形回帰モデルの β_1 と β_2 の最小二乗推定量は、次の正規方程式の解によって与えられる。

$$m_{11}\hat{\beta}_1 + m_{12}\hat{\beta}_2 = m_{1y}$$

$$m_{12}\hat{\beta}_1 + m_{22}\hat{\beta}_2 = m_{2y}$$

この連立方程式の解を求めると、次式を得る。

$$\hat{\beta}_1 = \frac{m_{22}m_{1y} - m_{12}m_{2y}}{m_{11}m_{22} - m_{12}^2}, \quad \hat{\beta}_2 = \frac{m_{11}m_{2y} - m_{12}m_{1y}}{m_{11}m_{22} - m_{12}^2}$$

\mathbf{x}_1 と \mathbf{x}_2 の間に近似的な共線関係 $\mathbf{x}_1 \doteq a\mathbf{x}_2$ が成り立っている場合、 $\mathbf{X}^\top \mathbf{X}$ の行列式は近似的に 0 になる。つまり、

$$m_{11}m_{22} - m_{12}m_{21} = m_{11}m_{22} - m_{12}^2 \doteq 0$$

となるため、最小二乗推定値は求まるが、分母が 0 に近いため、その値は不安定になることが考えられる。実際、分母が 0 に近いということは、観測誤差や丸めの誤差のために、分子の値がごくわずかに変化したとき、推定値が大幅に変動する可能性を示唆している。また、 $\hat{\beta}_1$ と $\hat{\beta}_2$ の分散共分散行列は次式となる。

$$V[\hat{\beta}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{\sigma^2}{m_{11}m_{22} - m_{12}^2} \begin{pmatrix} m_{22} & -m_{12} \\ -m_{12} & m_{11} \end{pmatrix}$$

したがって、 σ^2 が小さいとしても、 \mathbf{x}_1 と \mathbf{x}_2 の間に近似的な共線関係 $\mathbf{x}_1 \doteq a\mathbf{x}_2$ が成り立っているれば、推定量の分散・共分散は著しく大きくなることがわかる。 \mathbf{x}_1 と \mathbf{x}_2 が強い正の相関をもっている ($a > 0$) ならば、 $\hat{\beta}_1$ と $\hat{\beta}_2$ の間に強い負の相関が生じる。つまり、 \mathbf{x}_1 の係数 $\hat{\beta}_1$ を過大（または過少）推定したとすれば、 \mathbf{x}_2 の係数 $\hat{\beta}_2$ を過少（または過大）推定する傾向があることを示している。

2 多重共線への対処

多重共線に対処するための一つの可能な推定上の工夫として、後に述べる Ridge 回帰がある。通常よく用いられるのは、多重共線関係にある変数群の一部を除去する方法である。仮に、 \mathbf{x}_1 と \mathbf{x}_2

が強度の多重共線関係にあるとする。つまり、 x_1 と x_2 の間に近似的な共線関係 $x_1 \doteq ax_2$ が成り立っているとする。 x_1 が既に線形回帰モデルに含まれているとき、 x_2 を追加することによりもたらされる Y への説明力の増加は大きくないと考えられる。回帰分析の目的が予測であり、 x_1 と x_2 の多重共線が構造的なものであれば、 x_2 を除去することにより、精度高く予測できることが見込まれる。多重共線の問題は、**共線関係にある変数の一部を除去することにより処理されている**、というのが現状である。

3 多重共線の数値例

多重共線性が線形回帰モデルの推定にどのような影響をもたらすのかを例示するために、表 1 のデータ (E. Hilleboe. *New York State Journal of Medicine*, **57**, 1957, 2243–2254) を用いる。 x_1 は、総カロリーに占める脂肪のカロリーの比率、 x_2 は総カロリーに占める動物性蛋白のカロリーの比率、 Y は $100 \times [\log(55 \text{ 歳から } 59 \text{ 歳までの男子 } 10 \text{ 万人当たりの心臓疾患による死亡者数}) - 2]$ である。

表 1 心臓疾患による死亡率とカロリー摂取

国名	x_1	x_2	Y
オーストラリア	33	8	81
オーストリア	31	6	55
セイロン	17	2	24
デンマーク	39	6	52
フィンランド	30	7	88
フランス	29	7	45
ドイツ	35	6	50
アイルランド	31	5	69
イスラエル	23	4	66
イタリア	21	3	45
日本	8	3	24
メキシコ	23	3	43
オランダ	37	6	38
ニュージーランド	40	8	72
ノルウェー	38	6	41
ポルトガル	25	4	38
スウェーデン	39	7	52
スイス	33	7	52
イギリス	38	6	66
アメリカ	39	8	89
カナダ	38	8	80

心臓病の発症率が動物性蛋白や脂肪の摂取と関係あることは、医学的に良く知られた事実であ

る。つまり、動物性蛋白や脂肪の摂取が多い人ほど、心臓病の発症率が高くなることが予想される。この関係を検証するために、(3) 式の線形回帰モデルにより、表 1 のデータを回帰分析してみる。

$$\mathbf{Y} = \beta_0 + \beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2 + \varepsilon \quad (3)$$

(3) 式の線形回帰モデルにより、線形回帰分析をすると表 2 の結果が得られた。

表 2 表 1 に対する線形回帰分析の結果

変数	推定値	標準誤差	t 値
\mathbf{x}_1	-0.2230	0.6563	-0.340
\mathbf{x}_2	8.0437	3.0046	2.677
定数項	16.6217	11.9115	1.395

動物性蛋白や脂肪の摂取が多い人ほど心臓病の発症率が高くなる、つまり、 \mathbf{x}_1 と \mathbf{x}_2 の回帰係数 β_1 と β_2 は正となることが予想される。しかし、予想に反して、表 2 の結果では、 β_1 の推定値は負となった。この推定結果は、脂肪の摂取が多いほど心臓病による死亡率が低い、あるいは脂肪の摂取量と心臓病による死亡率との間には有意な関係が認められない、という医学的に良く知られた事実とは異なる結論を示唆している。表 2 の結果の妥当性を検討するために、より詳しくデータを調べてみる。

3 変数の相関係数を表 3 にまとめた。表 3 からわかるように、 \mathbf{x}_1 と \mathbf{Y} には、正の相関があることがわかる。つまり、脂肪の摂取が多い人ほど、心臓病の発症率が高くなることが予想されるという、医学的に良く知られた事実と一致する結果となっている。 \mathbf{x}_1 と \mathbf{x}_2 の間に強い正の相関があることがわかる。一般に、動物性蛋白の含有量の高い食品は、脂肪も多く含む傾向がある。したがって、 \mathbf{x}_1 と \mathbf{x}_2 の間に強い正の相関があることは妥当な結果である。表 2 の結果のように、脂肪の摂取が多いほど心臓病による死亡率が低い、あるいは脂肪の摂取量と心臓病による死亡率との間には有意な関係が認められない、という医学的に良く知られた事実とは異なる結果が得られたのは、 \mathbf{x}_1 と \mathbf{x}_2 の共線関係によるものである可能性が高いと考えられる。もっとはっきりとした結論を得るためには、観測値の個数を増やす必要がある。

表 3 相関行列

	\mathbf{x}_1	\mathbf{x}_2	\mathbf{Y}
\mathbf{x}_1	1.0	0.823	0.547
\mathbf{x}_2		1.0	0.704
\mathbf{Y}			1.0

4 Ridge 回帰

説明変数間に多重共線が存在するとき、一つの有用な対応策として、最小二乗推定法に次のような改良をする方法がある。近似的な共線関係が存在するとき、推定が不安定になるそもその原因は、 $\mathbf{X}^\top \mathbf{X}$ の行列式が近似的に 0 になることであった。そこで、 $\mathbf{X}^\top \mathbf{X}$ の対角成分に正定数 k を

加えて次のような推定量を考える。

$$\hat{\beta}_k = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{Y}$$

この推定量を **Ridge 推定量**という。Ridge 推定量は、 β の不偏推定量ではない。

$$\begin{aligned} E[\hat{\beta}_k] &= E[(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{Y}] \\ &= E[(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \varepsilon)] \\ &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X}\beta \\ &\neq \beta \end{aligned}$$

しかし、 $(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1}$ は $(\mathbf{X}^\top \mathbf{X})^{-1}$ よりも、推定量の分散共分散行列が小さくなるという意味で、推定が安定すると予想される。このことは、次式から容易に確かめることが可能である。

$$\begin{aligned} V[\hat{\beta}_k] &= V[(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{Y}] \\ &= V[(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \varepsilon)] \\ &= V[(\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \varepsilon] \\ &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top V[\varepsilon] (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X} \\ &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I}_n) \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \\ &\leq \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \quad (\because \text{補足 3 を参照}) \\ &= V[\hat{\beta}] \end{aligned}$$

補足 3

任意の列ベクトル \mathbf{a} に対して次式が成り立つ。

$$\begin{aligned} &\mathbf{a}^\top [(\mathbf{X}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1}] \mathbf{a} \\ &= [\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} - \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1}] \mathbf{a} \\ &= \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} - \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{a} \\ &\geq 0 \end{aligned}$$

したがって、対称行列 $(\mathbf{X}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_{p+1})^{-1}$ は非負定値行列である。

以上の結果から、任意の列ベクトル \mathbf{a} に対して次式が成り立つ。

$$\begin{aligned} \mathbf{a}^\top (V[\hat{\beta}] - V[\hat{\beta}_k]) \mathbf{a} &\geq 0 \Leftrightarrow \mathbf{a}^\top V[\hat{\beta}] \mathbf{a} \geq \mathbf{a}^\top V[\hat{\beta}_k] \mathbf{a} \\ &\Leftrightarrow V[\mathbf{a}^\top \hat{\beta}] \geq V[\mathbf{a}^\top \hat{\beta}_k] \end{aligned}$$

列ベクトル \mathbf{a} を単位ベクトルとすれば、次式を得る。

$$V[\hat{\beta}] \geq V[\hat{\beta}_k]$$

したがって、 k をうまく選択することで、Ridge 推定量を用いることで、推定値や予測値の平均二乗誤差を小さくすることができる。 k の選択方法は、Mallows の C_p 規準に代表される情報量規

準の最小化により求められることが多い。それらの情報量規準は、予測の平均二乗誤差に基づくリスク関数の推定量であり、それを最小にする k を選ぶことで、あてはめ値の予測平均二乗誤差が小さくなることが期待できる。しかし、情報量規準を最小にする解は陽な形で求めることができないため、実際の最適化には計算機による繰り返し計算が必要になる。

参考文献

- [1] 佐和隆光. (2020). 回帰分析 (新装版). 朝倉書店.
- [2] 柳原宏和, 永井勇, 佐藤健一. (2009). 多変量一般化リッジ回帰におけるリッジパラメータ最適化のためのバイアス補正 C_p 規準. 応用統計学, **38**, 151–172.