

建設省 正員。桐越信
 北海道大学 正員 五十嵐日出夫
 北海道大学 正員 山形耕一

1.はじめに

重回帰モデルは工事計画の種々の分野で予測モデルとして数多く用いられている。重回帰モデルの構築にあたって直面する最も重要な問題は、モデル構築時にモデルに取り入れる説明変数をどのように選択するべきかという問題である。モデル構築以前に、説明変数がまえもって決定されていることはほとんどない。一般には、モデル構築時に説明変数のいくつが候補群が、モデル構築者の経験や勘などによって主観的に用意され、そのなかから試行錯誤の後に、なんらかの客観的基準により最も望ましいと思われる説明変数の組合せが選択される。

マローズのC_p-統計量は説明変数選択の客観的基準のひとつであり、モデル構築時の偏りを明示的、定式化のなかにとり入れているところにその特徴がある。本研究は、このマローズのC_p-統計量を基礎にして、さらにモデルに取り入れられた説明変数の予測時における予測誤差を考慮した予測誤差基準について考察するものである。

2.マローズのC_p-統計量¹⁾

マローズのC_p-統計量では、モデル構築時の基本構造式を

$$Y_i = \beta_0 + \sum_{k=1}^p \beta_k X_{ik} + \varepsilon_i + U_i \quad (1)$$

と設定する。(1)式で ε_i は想定されたモデルの未知の偏りを表わし、 U_i は正規分布N(0, η^2)に従うと仮定されている。ここで、 ε_i について、

$$\sum_{i=1}^n \varepsilon_i = 0, \quad k = 1, 2, \dots, p \quad (2)$$

$$\sum_{i=1}^n \varepsilon_i^2 = 0 \quad (3)$$

を仮定する。(2), (3)式を仮定すると偏回帰係数 β_k 、 β_0 の最小2乗推定量 b_k 、 b_0 はより不偏推定量となる。説明変数が($X_{i1}, X_{i2}, \dots, X_{ip}$)をとるとときの予測対象 \hat{Y}_i を(1)式に従って(4)式で定義する。

$$\hat{Y}_i = \beta_0 + \sum_{k=1}^p \beta_k X_{ik} + \varepsilon_i + U_i \quad (4)$$

(4)式において、 U_i は(1)式の U_i と独立に、ガフ

すべて互いに独立に正規分布N(0, η^2)に従うものとする。このとき予測値 \hat{Y}_i は、

$$\hat{Y}_i = b_0 + \sum_{k=1}^p b_k X_{ik} \quad (5)$$

と与えられるので、

$$C_p = \sum_{i=1}^n (\hat{Y}_i - Y_i)^2 \quad (6)$$

とおくと、

$$E[C_p] = (n+p+1)\eta^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad (7)$$

となる。

一方、この場合、残差平方和の期待値E[RSS]は、

$$E[RSS] = (n-p-1)\eta^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad (8)$$

となるので、(7), (8)式より、

$$E[C_p] = E[RSS] + 2(p+1)\eta^2 \quad (9)$$

となる。ここで、 η^2 について不偏推定量 $\hat{\eta}^2$ が求められるなら、

$$C_p = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + 2(p+1)\hat{\eta}^2 \quad (10)$$

と書きかえよ。(10)式を(9)式の推定量とすることができる。 $\hat{\eta}^2$ については、一般には、

$$\min \hat{\eta}_{\text{p}}^2 = \min \left[\frac{1}{n-p-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \right] \quad (11)$$

が採用されている。

(10)式はモデルが偏りをもつことを前提にしたうえで、そのモデルの望ましさを評価する基準であると考えることができる。(10)式で、第1項は説明変数の数の増加に伴って小さくなるが、第2項は大きくなるので、マローズのC_p-統計量はQSS, PSS, AICなどと同じようにモデルの簡略さと適合性とのトレード・オフの関係を表わしていると考えることができる。QSS²⁾, PSS³⁾, AIC⁴⁾はそれぞれ(12), (13), (14)式で表現される。

$$\begin{aligned} QSS &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + 2(p+1)\hat{\eta}^2 \\ &= (n+p+1)\hat{\eta}^2 \\ &= \frac{(n+p+1)}{(n-p-1)} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \end{aligned} \quad (12)$$

$$PSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (13)$$

ここで、 $\hat{y}_i = b_0 + \sum_{k=1}^p b_k x_{ik}$

$$AIC = n \ln \hat{\eta}^2 + 2(p+1) \quad (14)$$

(12), (13), (14)式では、モデル構築時の基本構造式を、

$$y_i = \beta_0 + \sum_{k=1}^p \beta_k x_{ik} + \varepsilon_i \quad (15)$$

と設定し、 ε_i は正規分布 $N(0, \sigma^2)$ に従うと仮定されている。(12)式のQSSにおける $\hat{\eta}^2$ は $\hat{\eta}^2$ の不偏推定量である。(13)式のPSSにおける b_0, b_k はし組目のデータ ($y_i, x_{i1}, x_{i2}, \dots, x_{ip}$) を除いた $(n-1)$ 組のデータより求めた偏回帰係数である。(14)式の AIC における $\hat{\eta}^2$ は $\hat{\eta}^2$ の最尤推定量である。したがって $\hat{\eta}^2$ は不偏推定量ではなく偏推定量であり、漸近的不偏推定量である。

3. C_p -統計量を基礎とした予測誤差基準

(4)式の予測対象年の予測にあたって、モデルに取り入れられた説明変数の予測時における予測誤差も考慮して、それを ΔX_{ik}^* とすると、予測値は

$$\hat{y}_i^* = b_0 + \sum_{k=1}^p b_k (x_{ik} + \Delta X_{ik}) \quad (16)$$

となる。したがって、予測誤差 \hat{e}_i^* は

$$\begin{aligned} \hat{e}_i^* &= y_i - \hat{y}_i^* \\ &= \beta_0 + \sum_{k=1}^p \beta_k x_{ik} + \varepsilon_i + U_i^* \\ &\quad - \{ b_0 + \sum_{k=1}^p b_k (x_{ik} + \Delta X_{ik}) \} \\ &= (\beta_0 - b_0) + \sum_{k=1}^p (\beta_k - b_k) x_{ik} \\ &\quad + \varepsilon_i + U_i^* - \sum_{k=1}^p b_k \Delta X_{ik} \end{aligned} \quad (17)$$

となる。ここで、 ΔX_{ik} の期待値 $E[\Delta X_{ik}]$ について

$$\begin{aligned} E[\Delta X_{ik}] &= 0, \quad i=1, 2, \dots, n \\ k &= 1, 2, \dots, p \end{aligned} \quad (18)$$

と仮定すると、

$$\begin{aligned} E\left[\sum_i \hat{e}_i^{*2}\right] &= E\left[\sum_{i=1}^n (y_i - \hat{y}_i^*)^2\right] \\ &= (n+p+1) \eta^2 + \sum_{i=1}^n \varepsilon_i^2 \\ &\quad + \sum_{i=1}^n V\left[\sum_{k=1}^p b_k \Delta X_{ik}\right] \\ &= E[RSS] + 2(p+1) \eta^2 \\ &\quad + \sum_{i=1}^n V\left[\sum_{k=1}^p b_k \Delta X_{ik}\right] \end{aligned} \quad (19)$$

となる。ここで、 ΔX_{ik} と ΔX_{jk} とが互いに独立であるとすると、(19)式は

$$\begin{aligned} E\left[\sum_i \hat{e}_i^{*2}\right] &= E[RSS] + 2(p+1) \eta^2 \\ &\quad + \sum_{i=1}^n V[\Delta X_{ik}] \{ (E[b_k])^2 + V[b_k] \} \end{aligned} \quad (20)$$

となる。

実際には、(20)式を

$$K = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2(p+1) \hat{\eta}^2 + \sum_{i=1}^n \sum_{k=1}^p V[\Delta X_{ik}] \{ b_k^2 + V[b_k] \} \quad (21)$$

として推定する。ここで (21)式の K を予測誤差基準とよぶこととする。

(21)式において、説明変数の数の増加に伴って、第1項は小さくなるが、第2項は大きくなるので、第1項+第2項はモデルの簡略さと適合性とのトレードオフの関係を表わしている。このことは、第1項+第2項が (10)式の C_p と同じものであることが明らかである。ところが (21)式では、説明変数の予測時における予測誤差のために、説明変数の数の増加に伴って第3項も大きくなる。このことは、(21)式の K が、説明変数の数の増加に対して (10)式の C_p よりもさらに大きなペナルティをつけることを示している。

4. おわりに

本研究では、重回帰モデル構築における説明変数選択のひとつの客観的基準であるマローズの C_p -統計量をとりあげ、これにさらにモデルに取り入れられた説明変数の予測時における予測誤差をも明示的に考慮した予測誤差基準について考察を行なった。説明変数の予測誤差をも考慮している点で予測誤差基準は予測モデルを評価するのにふさわしい評価方法といえるが、予測誤差基準によれば C_p -統計量によるよりもさらに一層予測モデルとしての重回帰モデルの簡略化の方向が示唆される。

＜参考文献＞

- 1) 竹内啓；現象と行動のながの統計数理，新星社，pp. 74~85, 1980年
- 2) 奥野忠一, 芳賀敏郎；統計量解説法, 日科技連出版社, p. 73, 1976年
- 3) 芳賀敏郎, 竹内啓, 奥野忠一；重回帰分析における変数選択の新しい基準, 品質, Vol. 6, No. 2, pp. 35~40, 1976年
- 4) 佐和隆光；回帰分析, 朝倉書店, pp. 150~153, 1979年
- 5) 山形耕一, 棚越信；交通需要予測の予測精度について, 第1回土木計画学研究発表会講演集, pp. 52~57, 1979年
- 6) 棚越信, 他；予測誤差基準による発生重回帰モデルの評価方法に関する研究, 第35回年講集, pp. 79~80