

# 線形回帰分析における部分影響力評価

竹内 秀一

Assessment of Partial Influence in Linear Regression

Hidekazu TAKEUCHI

There are two procedures to assess the influence of observations for explanatory variables in linear regression. One is based on the case deletion procedure in variable selection problems, and the other on selecting some variables before the case deletion. The latter is called the partial influence procedure. This paper gives a new expression of the partial influence measure proposed by Cook and Weisberg [5] to assess the influence of observations for the selected variables. The new expression of the partial influence measure consists of Cook's distance and a similar influence measure. Furthermore a cut-off point for the new expression is derived by using that for Cook's distance. For the single case with one observation and one variable deleted, the cut-off point for the new expression is also compared with the size-adjusted cut-off point proposed by Belsley, Kuh and Welsch [2] and Chatterjee and Hadi [4].

## 1 はじめに

線形回帰分析においては、説明変数に対する観測値の影響力評価 (assessment of influence) を検討する場合に大きく分けて二つの方法がある。一つは説明変数の選択問題 (以下では「変数選択問題 (variable selection problems)」とする) として観測値の影響力評価を考える場合であり、影響力を調べたい観測値集合を除去してから変数選択された説明変数集合に関する影響力を評価する方法である。もう一つは、手順としてはこの逆になるが、説明変数集合を選択 (選定) してから影響力を調べたい観測値集合を除去して影響力の評価をする方法 (以下では「部分影響力 (partial influence) 評価」とする) である。この二つの方法は、除去される対象の順序が異なるだけであるが、影響力評価の立場がまったく異質で

あるので、一般に影響力評価の結果も異なる。

前者の影響力評価方法については、Léger and Altman [7] や Takeuchi [10] あるいは竹内 [11] などの研究例があるが、後者についてはあまり研究されていない。もちろん、Belsley, Kuh and Welsch [2] や Cook and Weisberg [6] などの影響力評価（または回帰診断）に関する代表的な著書においては紹介されているが、研究論文としては Cook and Weisberg [5] や Chatterjee and Hadi [4] くらいしか見受けられない。最近では、Castillo, Hadi, Conejo and Fernández-Canteli [3] の論文の中で提案されている新規の方法の比較対象として取り上げられているが、影響力評価方法として踏み込んだ議論はされていない。

そこで本論文では、従来から提案されている部分影響力を評価するための診断統計量 (influence measure) を再検討し、新たな表現を導入する。また、この新表現に対してデータ数に基づいて調整された (size-adjusted) 打ち切り点 (たとえば、竹内 [12] を参照) を提案する。さらに、通常、一つの説明変数を除去し、そのつぎに一つの観測値を除去することにより部分影響力を評価するための診断統計量を算出しているが、これを一般化して複数の説明変数集合および複数の観測値集合を除去する場合へと拡張することについても検討をする。

本論文の構成は以下のとおりである。2 節では、線形回帰モデルおよび観測値の影響力評価で用いられる各種の定義を与える。3 節において、従来から提案されている部分影響力を評価するための診断統計量を与える。4 節では、その診断統計量の新表現を導入する。5 節では、提案する新表現に基づく打ち切り点を導出し、従来の打ち切り点との比較をする。最後の 6 節は全体のまとめである。

## 2 定義

ここでは、線形回帰モデルとして、

$$y = X\beta + \varepsilon$$

を考える。このとき、 $y$  は  $n \times 1$  の目的変数ベクトル、 $X$  は  $n \times q$  のフルランクの説明変数行列、 $\beta$  は  $q \times 1$  の回帰係数ベクトル、そして  $\varepsilon$  は  $n \times 1$  の誤差ベクトルであり、正規分布  $N(0, \sigma^2 I_n)$  に従うものとする。ただし、 $I_n$  は  $n$  次の単位行列を表す。また、 $\beta$  の最小 2 乗推定量は  $\hat{\beta} = (X'X)^{-1}X'y$  として得られ、 $\sigma^2$  の不偏推定量は  $\hat{\sigma}^2 = e'e/(n-q)$  となる。ただし、「 $'$ 」は行列あるいはベクトルの転置を表し、 $e$  は残差ベクトルであり、 $e = y - X\hat{\beta} = (I_n - H)y$  である。このとき、 $H$  は説明変数行列から構成されるハット行列 (hat matrix)  $H = X(X'X)^{-1}X'$  であり、その第  $i$  対角成分  $h_{ii}$  がてこ比である。ただし、 $1/n \leq h_{ii} < 1$  とする。

さらに、残差ベクトルの第  $i$  成分  $e_i$  を標準化した  $t_i = e_i/(\hat{\sigma}\sqrt{1-h_{ii}})$  を標準化残差 (内的スチューデント化残差) と呼び、 $t_i$  の定義式において、 $\hat{\sigma}$  の代わりに  $\hat{\sigma}_{(i)}$  を用いた  $t_i^* = e_i/\{\hat{\sigma}_{(i)}\sqrt{1-h_{ii}}\}$  をスチューデント化残差 (外的スチューデント化残差) と呼ぶ。ここで、添

字の  $(\cdot)$  は  $n$  個の観測値の中から除去される観測値番号または観測値集合を表し、 $\hat{\sigma}^2$  および  $\hat{\sigma}_{(i)}^2$  の関係式は、

$$\hat{\sigma}_{(i)}^2 = \frac{n-q-t_i^2}{n-q-1} \hat{\sigma}^2$$

である。よって、 $t_i$  および  $t_i^*$  の関係式は、

$$t_i = t_i^* \sqrt{\frac{n-q}{n-q-1+t_i^{*2}}} \quad (2.1)$$

となる。

つぎに、観測値の影響力評価を行うときに必要な各種の定義を与える。Cook and Weisberg [6] から、観測値の影響力評価を行うための典型的な診断統計量である Cook の距離は、除去される  $m$  個の観測値集合  $I = \{i_1, i_2, \dots, i_m\}$  に対して、 $y' = (y'_{(I)} \ y'_I)$ 、 $X' = (X'_{(I)} \ X'_I)$ 、それに  $e' = (e'_{(I)} \ e'_I)$  と分割することにより

$$CD_I = \frac{(\hat{\beta} - \hat{\beta}_{(I)})' X' X (\hat{\beta} - \hat{\beta}_{(I)})}{q \hat{\sigma}^2} = \frac{(\hat{y} - \hat{y}_{(I)})' (\hat{y} - \hat{y}_{(I)})}{q \hat{\sigma}^2} = \frac{e'_I (I_m - H_I)^{-1} H_I (I_m - H_I)^{-1} e_I}{q \hat{\sigma}^2} \quad (2.2)$$

と定義される。ただし、

$$H = \begin{pmatrix} X_{(I)} (X' X)^{-1} X'_{(I)} & X_{(I)} (X' X)^{-1} X'_I \\ X_I (X' X)^{-1} X'_{(I)} & X_I (X' X)^{-1} X'_I \end{pmatrix}$$

であり、 $H_I = X_I (X' X)^{-1} X'_I$  とする。

さらに、(2.2) 式の代替表現として、Takeuchi [8] は以下のような表現を提案した。

$$CD_I = c'_I \Pi_I^{-\frac{1}{2}} R_I^{-1} (\Pi_I + I_m - R_I) R_I^{-1} \Pi_I^{-\frac{1}{2}} c_I \quad (2.3)$$

ただし、

$$\begin{aligned} \Pi_I &= [\text{diag}(H_I)] [\text{diag}(I_m - H_I)]^{-1} \\ R_I &= [\text{diag}(I_m - H_I)]^{-\frac{1}{2}} (I_m - H_I) [\text{diag}(I_m - H_I)]^{-\frac{1}{2}} \\ c_I &= \frac{1}{\sqrt{q}} \Pi_I^{-\frac{1}{2}} t_I \end{aligned}$$

であり、

$$t_I = \frac{1}{\hat{\sigma}} [\text{diag}(I_m - H_I)]^{-\frac{1}{2}} e_I$$

である。ここで、 $\text{diag}(\cdot)$  は  $(\cdot)$  内の正方行列の対角成分を取り出し、非対角成分をすべて 0 とする対角行列を表す。

特に、 $I = \{i\}$  の場合、標準化残差および  $t$  を利用して (2.2) 式は

$$CD_i = \frac{t_i^2}{q} \cdot \frac{h_{ii}}{1 - h_{ii}} \quad (2.4)$$

と簡略化して表現される。また、(2.3) 式の代替表現については、

$$CD_i = c_i^2$$

と表現される。さらに、 $m$  個の観測値に関する影響力の単純な和については、

$$\sum_{i \in I} CD_i = c'c_I \quad (2.5)$$

と表記することもある。

### 3 部分影響力を評価するための診断統計量

部分影響力を評価するための診断統計量として、代表的な二つの指標の定義式を与える。どちらの定義も、説明変数集合  $J$  (または第  $j$  番目の説明変数) を除去した後に、観測値集合  $I$  (または第  $i$  番目の観測値) を除去することにより部分影響力を評価するという手法を採用している。これは結果的に、説明変数集合  $J$  に対する観測値集合  $I$  の部分影響力の評価指標に相当する。

#### 3.1 Belsley, Kuh and Welsch の定義

Belsley, Kuh and Welsch [2] は第  $j$  番目の説明変数に対する第  $i$  番目の観測値の影響力を評価するために、以下のような診断統計量を提案した。それは、第  $i$  番目の観測値を除去したときに、第  $j$  番目の説明変数に対する回帰係数に生じる変化量を、その標準偏差 (分散の平方根) で規準化した指標として

$$DFBETAS_{ij} = \frac{\hat{\beta}_j - \hat{\beta}_{(i)j}}{\sqrt{\text{var}(\hat{\beta}_j)}} = \begin{cases} \frac{t_i}{\sqrt{1-h_{ii}}} \cdot \frac{w_{ij}}{\sqrt{W_j'W_j}} \\ \frac{t_i^*}{\sqrt{1-h_{ii}}} \cdot \frac{w_{ij}}{\sqrt{W_j'W_j}} \end{cases} \quad (3.1)$$

と与えられる。ただし、 $\text{var}(\cdot)$  は  $(\cdot)$  内の分散を表し、 $W_j = (I - H_{[j]})X_j$  である。ここで、

$$H_{[j]} = X_{[j]}(X_{[j]}'X_{[j]})^{-1}X_{[j]} \quad \text{および} \quad w_{ij} = x_{ij} - x_{i[j]}(X_{[j]}'X_{[j]})^{-1}X_{[j]}'X_j$$

である。このとき、添字の  $[\cdot]$  は  $q$  個の説明変数の中から除去される説明変数番号または説明変数集合を表す。

(3.1)式の定義において、第一表現が Belsley, Kuh and Welsch [2] の定義であり、簡略化された第二表現は Chatterjee and Hadi [4] の定義である。正確には、第二表現の第二式(下)は、平方した場合が定義されているが、定義式を統一的に扱うためにこのように再定義する。また、Belsley, Kuh and Welsch [2] では(3.1)式の定義において、回帰係数(ベクトルの第  $j$  成分)の添字は  $\hat{\beta}_{(i)j}$  ではなく  $\hat{\beta}_{j(i)}$  と表記している。けれども、本論文では、説明変数行列  $X$  やハット行列  $H$  などの添字と同様に、観測値番号(または観測値集合)のつぎに説明変数番号(または説明変数集合)の順に統一した。

(3.1)式の第二表現において、第一式(上)および第二式(下)の相違は、 $\sigma^2$  の推定量とし

て、 $\hat{\sigma}^2$ を適用するのか $\hat{\sigma}_{(i)}^2$ を適用するのかという点である。この相違が $t_i$ と $t_i^*$ の差異として現れている。 $t_i$ および $t_i^*$ の関係式は(2.1)式のとおりであり、 $|t_i| \leq 1$ 程度であれば両者に大きな違いはない(逆に、 $|t_i| > 1$ の場合は変動が大きくなりやすくなる)。一般に、 $t_i$ を利用した場合は保守的な結果になり、 $t_i^*$ を利用した場合は劇的に大きな変化をもたらす結果になることがある。このため、通常の部分影響力評価においては、解析結果が安定している第一式を利用することが多い。しかしながら、Belsley, Kuh and Welsch [2] は第二式を部分影響力の評価式として定義している。

### 3.2 Cook and Weisberg の定義

Cook and Weisberg [5] [6] は、Belsley, Kuh and Welsch [2] とは異なる観点から、一般化された部分影響力評価のための診断統計量を提案した。説明変数集合を特定する行列  $L$  に基づき、観測値集合  $I$  の影響力を調べるために、以下のような診断統計量(距離規準)を導入したのである。二つのベクトル  $\hat{\Psi} = L\hat{\beta}$  および  $\hat{\Psi}_{(I)} = L\hat{\beta}_{(I)}$ 、ただし、行列  $L$  は大きさが  $\ell \times q$  でランクは  $\ell$  であるものと定義し、これらの間の距離を  $D_I(\hat{\Psi})$  とする。つまり、

$$D_I(\hat{\Psi}) = \frac{(\hat{\Psi} - \hat{\Psi}_{(I)})' [L(X'X)^{-1}L']^{-1} (\hat{\Psi} - \hat{\Psi}_{(I)})}{\ell \hat{\sigma}^2} \quad (3.2)$$

とする。特に、 $I = \{i\}$  および  $L = [0 \cdots 01]$ 、つまり  $\ell = 1$  の場合に、たとえば、第  $j$  列(ここでは、最後の  $q$  列目とみなす) についての部分影響力の評価式は

$$D_i(\hat{\Psi}) = \frac{(\hat{\beta} - \hat{\beta}_{(i)})' L' [L(X'X)^{-1}L']^{-1} L (\hat{\beta} - \hat{\beta}_{(i)})}{\hat{\sigma}^2} = \frac{t_i^2}{1 - h_{ii}} (h_{ii} - h_{ii(j)}) \quad (3.3)$$

となる。(3.3)式の簡単な式変形により第  $j$  番目の説明変数に対する  $D_i(\hat{\Psi})$  が  $DFBETAS_{ij}$  (の2乗、正確には第二表現の第一式の2乗) に一致することがわかる(付録Aを参照)。

### 4 診断統計量の新表現

前節の部分影響力を評価するための診断統計量との比較をするために、(3.2)式の  $D_I(\hat{\Psi})$  の表記を変えて  $D_{IJ}$  とし、全説明変数の数  $q$  と関連付けて以下のように定義し直す。すなわち、 $\ell$  個の説明変数集合  $J = \{j_1, j_2, \dots, j_\ell\}$  に対する  $m$  個の観測値集合  $I$  の部分影響力を測定するための診断統計量は、 $X = (X_{[I]} X_J)$  とし、

$$\hat{\beta}_{(I)} = \begin{pmatrix} \hat{\beta}_{(I)[I]} \\ \hat{\beta}_{(I)J} \end{pmatrix}$$

と分割するとき

$$\begin{aligned}
 D_{ij} &\equiv \frac{(\hat{\beta}_j - \hat{\beta}_{(i)})' [\text{Var}(\hat{\beta}_j)]^{-1} (\hat{\beta}_j - \hat{\beta}_{(i)})}{\ell} \\
 &= \frac{q}{\ell} \cdot \frac{(\hat{\beta} - \hat{\beta}_{(i)})' X' (H - H_{[i]}) X (\hat{\beta} - \hat{\beta}_{(i)})}{q\hat{\sigma}^2} = \frac{q}{\ell} \cdot \frac{(\hat{y} - \hat{y}_{(i)})' (H - H_{[i]}) (\hat{y} - \hat{y}_{(i)})}{q\hat{\sigma}^2}
 \end{aligned} \tag{4.1}$$

と定義される（式変形は付録 B を参照）。ただし、 $\text{Var}(\cdot)$  は  $(\cdot)$  内のベクトルの分散共分散行列を表し、 $H_{[i]} = X_{[i]}(X_{[i]}'X_{[i]})^{-1}X_{[i]}$  である。

(4.1) 式は (2.2) 式における Cook の距離を部分影響力評価のために拡張しているとも考えられる。ベクトルに挟まれた重み行列の形式は異なるが、よく似た表現になっている。この点を明確にするために、さらに、(4.1) 式に対して別の式変形を試みる。ハット行列は  $H = H_{[i]} + W_j(W_j'W_j)^{-1}W_j'$  と分解できる。ただし、 $W_j = (I_n - H_{[i]})X_j$  である。これを利用すると (4.1) 式の第二表現から、

$$\begin{aligned}
 D_{ij} &= \frac{q}{\ell} \cdot \frac{(\hat{\beta} - \hat{\beta}_{(i)})' X' W_j (W_j'W_j)^{-1} W_j' X (\hat{\beta} - \hat{\beta}_{(i)})}{q\hat{\sigma}^2} \\
 &= \frac{q}{\ell} \cdot \frac{(\hat{\beta} - \hat{\beta}_{(i)})' \begin{bmatrix} 0 & 0 \\ 0 & X_j'(I_n - H_{[i]})X_j \end{bmatrix} (\hat{\beta} - \hat{\beta}_{(i)})}{q\hat{\sigma}^2} \\
 &= \frac{q}{\ell} (CD_I - CD_{ij}^*)
 \end{aligned} \tag{4.2}$$

と変形することができる。ただし、

$$CD_{ij}^* \equiv \frac{e_i'(I_m - H_i)^{-1}H_{[i]}(I_m - H_i)^{-1}e_i}{q\hat{\sigma}^2} \tag{4.3}$$

であり、

$$H_{[i]} = X_{[i]}(X_{[i]}'X_{[i]})^{-1}X_{[i]}$$

とする。このとき、

$$X_{[i]} = \begin{pmatrix} X_{(i)[i]} \\ X_{I[i]} \end{pmatrix}$$

と分割している。

この結果、 $D_{ij}$  は典型的な 2 次形式であるから、 $D_{ij} \geq 0$  となる。つまり、(4.2) 式の第三表現から

$$CD_I \geq CD_{ij}^* (\geq 0)$$

となることがわかり、(2.2) 式の Cook の距離  $CD_I$  の最小値は (4.3) 式の  $CD_{ij}^*$  のとる値によって決まることになる。従来の研究においては、Barret and Gray [1] が  $CD_I$  の上限について議論をしているが、この点を含めた研究については別の機会に検討する。

## 5 打ち切り点の比較

部分影響力を評価するための診断統計量である(3.2)式についての一般的な打ち切り点は提案されていない。けれども、特定の観測値一つを除去した場合の(3.1)式については、Belsey, Kuh and Welsch [2] や Chatterjee and Hadi [4] により打ち切り点が導出されている。(3.1)式の打ち切り点として、第  $j$  番目の説明変数に対する第  $i$  番目の観測値が

$$\text{DFBETAS}_{ij}^2 = D_i(\hat{\Psi}) > \frac{4}{n} \quad (5.1)$$

となる場合を影響力の大きい観測値と判定するものとして導かれている。(5.1)式の導出は、(3.1)式と本質的に同じ(3.3)式において、標準化残差が  $|t_i| > 2$  (スチューデント化残差の場合も  $|t_i^*| > 2$ )、てこ比については  $h_{ii} = 1/n$  および  $h_{ii} - h_{ii(j)} = 1/n$ 、その上で、 $n-1 \cong n$  の近似が適用されている。てこ比への代入については、最初の「 $h_{ii} = 1/n$ 」はてこ比  $h_{ii}$  の平均  $q/n$  について  $q=1$ 、つまり一つの説明変数を想定した場合であり、つぎの「 $h_{ii} - h_{ii(j)} = 1/n$ 」は  $h_{ii} - h_{ii(j)}$  の  $n$  個の観測値についての合計が1になることから、その平均を適用した場合であるので、導出経過はまったく異なるのである。このため、てこ比に関するこれら二つの打ち切り点の導出自体に問題点(不自然さ)があることは否定できない。

そこで、こうした問題点を改善するために、提案する新表現に基づく診断統計量について、まずこの新表現の特徴を利用した打ち切り点を一般的な場合について導出する。つぎに観測値と説明変数がそれぞれ一つずつ除去される場合の部分影響力評価において、従来の打ち切り点および新表現に基づく打ち切り点の具体的な比較検討をする。

(4.2)式の第三表現に着目し、Cook の距離  $CD_I$  および  $CD_{I_j}^*$  の打ち切り点をそれぞれ導出する。Cook の距離  $CD_I$  については、いくつかの打ち切り点が提案されているが、ここでは、Takeuchi [9] における Welsch-Kuh の距離の打ち切り点導出と同様の方法を与える。その導出方法に従えば、(2.3)式において、 $CD_I \cong c'c_I$  と近似することになる。つまり、 $m$  個の観測値集合  $I$  を構成する観測値それぞれの単独の Cook の距離を単純に合計したものとして近似をするのである。すると、(2.4)式および(2.5)式から

$$CD_I \cong c'c_I = \sum_{i \in I} \frac{t_i^2}{q} \cdot \frac{h_{ii}}{1-h_{ii}}$$

となる。ここで、標準化残差 ( $t$  分布に従う統計量) についてはその2乗(つまり  $F$  分布に従う統計量) の期待値  $t_i^2 = (n-q)/(n-q-2)$  を、てこ比についてはその平均  $h_{ii} = q/n$  を代入して打ち切り点を導く。よって、

$$\sum_{i \in I} \frac{t_i^2}{q} \cdot \frac{h_{ii}}{1-h_{ii}} = \frac{m}{q} \cdot \frac{n-q}{n-q-2} \cdot \frac{q}{1-\frac{q}{n}} = \frac{m}{n-q-2} \quad (5.2)$$

となる。

(5.2)式の導出方法と同様に、(4.3)式の  $CD_{ij}^*$  についても打ち切り点を類推する。Cook の距離の定義式である (2.2) 式の第三表現における  $H_I$  が (4.3) 式における  $H_{I|J}$  に置き換わっただけであり、 $H_{I|J}$  の対角成分の和が  $q-l$  であるので、この平均  $(q-l)/n$  を利用すればよいのである。よって、 $CD_{ij}^*$  の近似式から

$$CD_{ij}^* \cong \sum_{i \in I} CD_{ij}^* = \sum_{i \in I} \frac{t_i^2}{q} \cdot \frac{h_{ii|J}}{1-h_{ii}} = \frac{m}{q} \cdot \frac{n-q}{n-q-2} \cdot \frac{\frac{q-l}{n}}{1-\frac{q-l}{n}} = \frac{m}{q} \cdot \frac{q-l}{n-q-2} \quad (5.3)$$

となる。

したがって、上記のことから、(4.2)式の  $D_{ij}$  の打ち切り点は、係数  $q/l$  を考慮して(5.2)式および(5.3)式から

$$D_{ij} > \frac{m}{n-q-2} \quad (5.4)$$

となる。(5.4)式は一般的な場合であるので、第  $j$  番目の説明変数に対する第  $i$  番目の観測値に限定して比較を行うことにする。つまり、 $I=\{i\}$  および  $J=\{j\}$  の場合に対して比較をするので、(5.4)式は、

$$D_{ij} > \frac{1}{n-q-2} \quad (5.5)$$

となる。この(5.5)式と従来の打ち切り点である(5.1)式のそれぞれの右辺についての差を計算すると、

$$(5.1)式 - (5.5)式 = \frac{3n-4(q-2)}{n(n-q-2)} \quad (5.6)$$

となる。(5.6)式の分母は明らかに正の値であるので、分子の大小関係だけが問題になる。分子についても、 $q \geq 2$  ( $q=2$  は「定数項+説明変数」の単回帰) であるので、

$$n > \frac{4}{3}(q-2)$$

の場合に、提案する新しい打ち切り点(5.5)式が従来の(5.1)式よりも小さくなり、より厳密な打ち切り点になるといえる。

一般的な場合については比較対象となる打ち切り点が存在しないので、明確なことは言えないが、特定の一つの説明変数に対する複数個の観測値の影響力評価については、一つの観測値の場合と同様の結果になることが予想される。だが、複数の説明変数に対する複数個の観測値の影響力評価については、 $CD_{ij}^*$  の打ち切り点をより精密に検討する必要があると思われるので、簡単に判断をすることはできないであろう。



## 6 まとめ

本論文では、部分影響力を評価するための診断統計量について、従来の指標に関する新表現を提案し、その新表現に基づく打ち切り点を導出した。新表現は観測値の影響力評価（回帰診断）においてよく利用される Cook の距離とよく似た表現形式になっている。このため、通常の回帰診断を実施すればそれと連動させて部分影響力を測定し評価することも可能となる。また、この新表現に対して新たに導出された打ち切り点が、従来の打ち切り点よりも、ある条件下で常に小さくなることが示された。従来の打ち切り点は、大き目に設定されており、実際のデータ解析での影響力評価において、あまり有効ではないという欠点が指摘されている。けれども、この新しい打ち切り点により部分影響力評価のための診断統計量が、実用上、より利用しやすいものになったと考えられる。

今後の課題としては、変数選択問題と部分影響力の評価方法における説明変数集合と観測値集合の除去手順の違いによる立場の違いをより厳密に検討し、それぞれの相違点あるいは類似点を明確にする必要がある。また、二つの評価方法の相補性についても、一般的な診断統計量との関連性から、実用上の観点も視野に入れて検討することが必要であると考えられる。加えて、この二つの課題とは異なるが、Cook の距離の上限と下限の議論についても、部分影響力評価と関連付けて研究を進めていきたいと考えている。

## 付録 A : (3.1) 式および (3.3) 式の同一性

(3.1)式の第二表現の第一式を2乗すると

$$DFBETAS_{ij}^2 = \frac{t_i^2}{1-h_{ii}} \cdot \frac{w_{ij}^2}{W_j W_j}$$

となる。(3.3)式への式変形を考える上で、 $w_{ij}^2/W_j W_j$ をハット行列 H の成分により表現することが問題になる。そこで、 $X = (X_{[1]} X_j)$  とし、

$$X = \begin{bmatrix} X_{(i)[1]} & X_{(i)j} \\ x_{i[1]} & x_{ij} \end{bmatrix}$$

と分割すると

$$H = H_{[1]} + \frac{(I_n - H_{[1]}) X_j X_j' (I_n - H_{[1]})}{X_j' (I_n - H_{[1]}) X_j} = H_{[1]} + \frac{W_j W_j'}{W_j' W_j}$$

であるので、

$$H - H_{[1]} = \frac{W_j W_j'}{W_j' W_j}$$

となる。この両辺の行列成分は当然等しいので、第  $i$  対角成分についても

$$h_{ii} - h_{ii[j]} = \frac{w_{ij}^2}{W_j W_j}$$

となる。ただし、

$$W_j = \begin{bmatrix} X_{(i)j} - X_{(i)[j]}(X'_{[j]}X_{[j]})^{-1}X'_{[j]}X_j \\ x_{ij} - x_{i[j]}(X'_{[j]}X_{[j]})^{-1}X'_{[j]}X_j \end{bmatrix} = \begin{pmatrix} W_{(i)j} \\ w_{ij} \end{pmatrix}$$

である。したがって、

$$\text{DFBETAS}_{ij}^2 = \frac{t_i^2}{1 - h_{ii}}(h_{ii} - h_{ii[j]}) = D_i(\hat{\Psi})$$

となり、(3.1)式の第二表現の第一式を2乗したものと(3.3)式は一致することがわかる。

### 付録 B: (4.1) 式の導出

(3.2)式を  $\hat{\beta}$  および  $\hat{\beta}_{(i)}$  を使って表すと

$$D_I(\hat{\Psi}) = \frac{(\hat{\beta} - \hat{\beta}_{(i)})'L'[L(X'X)^{-1}L']^{-1}L(\hat{\beta} - \hat{\beta}_{(i)})}{\ell\sigma^2} \quad (\text{B.1})$$

となる。ここで、 $L = (O \quad I_\ell)$ 、 $A = (X'_{[j]}X_{[j]})^{-1}$  それに  $B = [X'_j(I_n - H_{[j]})X_j]^{-1}$  とすると、

$$\begin{aligned} L(X'X)^{-1}L' &= (O \quad I_\ell) \begin{bmatrix} A + AX'_{[j]}X_jBX'_jX_{[j]}A & -AX'_{[j]}X_jB \\ -BX'_jX_{[j]}A & B \end{bmatrix} \begin{pmatrix} O \\ I_\ell \end{pmatrix} \\ &= B = \frac{1}{\sigma^2} \cdot \sigma^2(X'X)^{-1} = \frac{1}{\sigma^2} \text{Var}(\hat{\beta}_j) \end{aligned}$$

であり、

$$L(\hat{\beta} - \hat{\beta}_{(i)}) = (O \quad I_\ell) \begin{pmatrix} \hat{\beta}_{[j]} - \hat{\beta}_{(i)[j]} \\ \hat{\beta}_j - \hat{\beta}_{(i)j} \end{pmatrix} = \hat{\beta}_j - \hat{\beta}_{(i)j}$$

である。よって、

$$D_I(\hat{\Psi}) = \frac{(\hat{\beta}_j - \hat{\beta}_{(i)j})'[\text{Var}(\hat{\beta}_j)]^{-1}(\hat{\beta}_j - \hat{\beta}_{(i)j})}{\ell} = D_{Ij}$$

である。ただし、 $\sigma^2$  をその不偏推定量  $\hat{\sigma}^2$  で置き換えている。これで(3.2)式と(4.1)式の第一表現は一致することが証明された。

つぎに、(3.2)式から(4.1)式の第二表現および第三表現への式変形を示す。(3.2)式の分子に着目すると(B.1)式の分子から

$$\begin{aligned} &(\hat{\beta} - \hat{\beta}_{(i)})'L'[L(X'X)^{-1}L']^{-1}L(\hat{\beta} - \hat{\beta}_{(i)}) \\ &= [X(\hat{\beta} - \hat{\beta}_{(i)})]'X(X'X)^{-1}L'[L(X'X)^{-1}X'X(X'X)^{-1}L']^{-1}L(X'X)^{-1}X'[\hat{\beta} - \hat{\beta}_{(i)}] \\ &= (\hat{y} - \hat{y}_{(i)})'C(C'C)^{-1}C'(\hat{y} - \hat{y}_{(i)}) \end{aligned} \quad (\text{B.2})$$

となる。ただし、 $\hat{y} = X\hat{\beta}$ 、 $\hat{y}_{(i)} = X\hat{\beta}_{(i)}$  それに  $C = X(X'X)^{-1}L'$  である。ここで、上記の式変形過程の結果から

$$(C'C)^{-1} = X'_j(I_n - H_{[j]})X_j = B^{-1}$$

であり、また

$$C = X(X'X)^{-1}L' = (I_n - H_{[j]})X_j[X'_j(I_n - H_{[j]})X_j]^{-1} = (I_n - H_{[j]})X_jB$$

であるので、

$$\begin{aligned} C(C'C)^{-1}C' &= (I_n - H_{[j]})X_j[X'_j(I_n - H_{[j]})X_j]^{-1}X'_j(I_n - H_{[j]}) \\ &= W_j(W'_jW_j)^{-1}W'_j = H - H_{[j]} \end{aligned}$$

となる。よって、(B.2)式は

$$\begin{aligned} (\hat{\beta} - \hat{\beta}_{(j)})'L'[L(X'X)^{-1}L']^{-1}L(\hat{\beta} - \hat{\beta}_{(j)}) &= (\hat{y} - \hat{y}_{(j)})'(H - H_{[j]})(\hat{y} - \hat{y}_{(j)}) \\ &= (\hat{\beta} - \hat{\beta}_{(j)})'X(H - H_{[j]})X(\hat{\beta} - \hat{\beta}_{(j)}) \end{aligned}$$

と式変形でき、(3.2)式の分子は(4.1)式の定数項  $q/\ell$  部分を除き、それぞれ第三表現および第二表現の分子と一致する。

以上のことより、(3.2)式から(4.1)式を導出できることが証明された。

#### 参考文献

- [1] Barrett, B. E. and Gray, J. B. (1992), Efficient computation of subset influence in regression, *Journal of Computational and Graphical Statistics*, 1, 271-286.
- [2] Belsley, D. A., Kuh, E. and Welsch, R. E. (1980), *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*, Wiley: New York.
- [3] Castillo, E., Hadi, A. S., Conejo, A. and Fernández-Canteli, A. (2004), A general method for local sensitivity analysis with application to regression models and other optimization problems, *Technometrics*, 46, 430-444.
- [4] Chatterjee, S. and Hadi, A. S. (1986), Influential observations, high leverage points and outliers in linear regression, *Statistical Science*, 1, 379-416.
- [5] Cook, R. D. and Weisberg, S. (1980), Characterizations of an empirical influence function for detecting influential cases in regression, *Technometrics*, 22, 495-508.
- [6] Cook, R. D. and Weisberg, S. (1982), *Residuals and Influence in Regression*, Chapman and Hall: New York.
- [7] Léger, C. and Altman, N. (1993), Assessing influence in variable selection problems, *Journal of the American Statistical Association*, 88, 547-556.
- [8] Takeuchi, H. (1991), Detecting influential observations by using a new expression of Cook's distance, *Communications in Statistics—Theory and Methods*, 20, 261-274.
- [9] Takeuchi, H. (1992), Regression diagnostics using a new expression of Welsch-Kuh distance, 静岡県立大学経営情報学部報「経営と情報」, 4, 17-26.
- [10] Takeuchi, H. (2002), Assessment of influence of individual observations on prediction mean square errors in variable selection problems, *Journal of the Japan Statistical Society*,

## 線形回帰分析における部分影響力評価

32, 43-55.

- [11] 竹内秀一 (2003), 変数選択問題における観測値除去法に基づく診断統計量, 人文自然科学論集, 116号, 23-36.
- [12] 竹内秀一 (2005), 線形回帰における尤度距離による影響力評価, 人文自然科学論集, 119号, 19-30.