

線型回帰分析における

観測誤差の影響と誤差の評価について

三枝義清

- 一 問題
 - 二 線型回帰分析における観測誤差の問題
 - (一) proxy variables を確定変数とみなした場合
 - (二) errors in variables model
 - 三 観測誤差の評価—農産生産統計を対象にして
 - (一) 農産センサスとの比較
 - (二) 比較の方法
 - (三) errors model の設定

一 問 題

線型回帰モデルを実際に推定しようとする場合、考慮すべき説明変数のなかには観測不能のものや、利用しうるデータがないものが出てくる。それらの変数を説明変数として回帰モデルの中に組み入れようとすれば、観測可能な別の変数で代理させねばなるまい。たとえば期待価格や期待収量を説明変数にした農産物の供給分析では、これらの変数を過去数年の実現値の加重平均値で代用させるだろう。農産物の価格指数や気象指数を説明変数にしようとする場合——これらの変数は観測不能ではないが——適当な統計データがないために、入手可能な別の指数で代用することがしばしば行なわれる。厳密に考えれば、トレンド項を除けば大半の説明変数は真の変数に対する代

理変数 (proxy variable) とみななければならない。多くの場合、作付面積や収量に関する統計データは標本誤差や回答誤差などの観測誤差を含むものであるから、われわれの利用するこれらの観測値も true variable の proxy variable とみなすことができる。proxy variable と true variable の関係はさまざまで、

$$\text{proxy variable} = \text{true variable} + \text{observation error} \dots\dots\dots (1.1)$$

という簡単なものもあるし、期待価格の場合のように複雑な関係をもつものもある。

回帰分析で説明変数に proxy variable を使った場合、推定されたパラメーターは真のパラメーター (true variable を説明変数にした時の) とどのような関係をもつだろうか。二の〔1〕では proxy variable を確定変数とみなした場合を取り扱い、二の〔2〕では proxy variable が一個で、特に (1.1) のような関係がある場合についてパラメーターの推定値の分布を論じた。三は観測誤差を評価する際の一つの試みを述べたものである。

二 線型回帰分析における観測誤差の問題

従属変数 Y は次のような線型回帰式にしたがうものとする。

$$Y_j = \sum_{i=1}^k \beta_i X_{ij} + u_j \dots\dots\dots (2.1)$$

ただし、 $j=1, 2, \dots, n$

説明変数はすべて確定変数で、誤差項 u_j は

$$E(u_j) = 0 \quad E(u_j^2) = \sigma^2 \quad E(u_j u_k) = 0 \quad (j \neq k)$$

とする。2.1式を行列表形式にして次のように表わすことにする。

$$y = \sum_{i=1}^k \beta_i x_i + u \dots \dots \dots (2 \cdot 2)$$

4895E

$$y = Z\beta + u \dots \dots \dots (2 \cdot 3)$$

ただし

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ni} \end{pmatrix} \quad i = 1, 2, \dots, k$$

$$u = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad Z = [z_{11}, \dots, z_{nk}]$$

true variable x_i の代わりに別の変数 x_i を proxy variable として ($i=1, 2, \dots, k$) β の最小二乗推定を行なうたしよう。最小二乗推定値を $\hat{\beta} = [\hat{\beta}_1, \dots, \hat{\beta}_k]'$ とすれば

$$\hat{\beta} = [X'X]^{-1}X'y = [X'X]^{-1}X'(Z\beta + u)$$

ただし $x_i = [x_{i1}, \dots, x_{ni}]'$ ($i=1, 2, \dots, k$)

$$X = [x_1, \dots, x_k]$$

(1) Proxy variables を確定変数とみなした場合

x_i ($i=1, 2, \dots, k$) を確定変数とみなすことにすれば

線型回帰分析における観測誤差の影響と誤差の評価について

線形回帰分析における観測誤差の分散・共分散行列の逆行列として

式

$$[X'X]^{-1}X'Z = \begin{bmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{bmatrix} = A \cdots \cdots \cdots (2.4)$$

と置く。又

$$E(\hat{\beta}) = A\beta \cdots \cdots \cdots (2.5)$$

又 $\hat{\beta}$ の 4 階の $[a_{11}, \dots, a_{kk}]'$ は

$$[a_{11}, \dots, a_{kk}]' = [X'X]^{-1}X'z_1$$

である。よって z_1, \dots, z_k は回帰分析データの回帰係数ベクトルである。

よって $z_1 = z_1, \dots, z_k = z_k$ とする。

$$a_{11} = 1, \quad a_{j1} = 0 \quad (j \neq 1)$$

又 $\hat{\beta}$ の 2 階の $a_{11} = z_1, \dots, z_k$ とする。

$$X'X = \begin{bmatrix} z_1'z_1 & z_1'z_2 \\ \tilde{X}'z_1 & \tilde{X}'z_2 \end{bmatrix} \quad \text{ただし} \quad \tilde{X} = [x_1, \dots, x_k]$$

$$\begin{bmatrix} \tilde{X}'z_1 \\ \tilde{X}'z_2 \end{bmatrix}$$

よ

$$\begin{bmatrix} a_{11} \\ \vdots \\ a_{k1} \end{bmatrix} = \begin{bmatrix} z_1'z_1 & z_1'z_2 \\ \tilde{X}'z_1 & \tilde{X}'z_2 \end{bmatrix}^{-1} \begin{bmatrix} z_1'z_1 \\ \tilde{X}'z_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 0_{k-1} \end{bmatrix}$$

ただし 0_{k-1} は $k-1$ 次のゼロベクトル

となる。

従って特に $x_i = x_i'$ ($i=1, 2, \dots, k-1$), $x_k = x_k$ とすれば $E(\hat{\beta})$ と β の関係は次のように表わされる。

$$E(\hat{\beta}) = \begin{bmatrix} I_{k-1} & \tilde{a}_k \\ 0 & a_k \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \beta_k \end{bmatrix}$$

あるいは

$$E(\hat{\beta}) = \tilde{\beta} + \beta_k \tilde{a}_k \dots \dots \dots (2.6)$$

$$E(\hat{\beta}_k) = \beta_k a_k \dots \dots \dots (2.7)$$

ただし $\tilde{\beta} = [\beta_1, \dots, \beta_{k-1}]'$ $\hat{\beta} = [\hat{\beta}_1, \dots, \hat{\beta}_{k-1}]'$

$\tilde{a}_k = [a_{1k}, \dots, a_{k-1k}]'$

従って

$$x_k = \sum_{i=1}^{k-1} x_i a_{ik} + x_k a_{kk} + \text{residual}$$

と表わける。

2・5式にみるように proxy variables により推定された係数の期待値 $E(\hat{\beta}_k)$ は、一般にオリジナルな係数 β_k 、 \dots , β_k の一次結合である。従って 2・4式の行列 A に関する知識がない限り β に関する先験的情報(符号条件とか係数間の制約条件など)が与えられたとしても、それを係数の推定に利用することはできない。

$\hat{\beta}$ と β の関係を更に追求するには proxy variable と true variable の関係を特定化せねばならないが、ここでは 1・1式のような関係がある場合について考察を進めてみよう。問題を単純化するために説明変数の中の一個に

線型回帰分析における観測誤差の影響と誤差の評価について

この proxy variable が使われているものとする。そして proxy variable x は次のようなモデル (errors in variables model) に従う確率変数であると想定する。

$$x_j = \xi_j + v_j \dots \dots \dots (2 \cdot 8)$$

$$j = 1, 2, \dots, n$$

$$\text{ただし } E(v_j) = 0 \quad E(v_j^2) = \sigma_v^2 \quad E(v_j v_{j'}) = 0 \quad (j \neq j')$$

ξ_j は確実変数だが、その値は未知な変数。

従属変数 y_j の従う線型回帰式を書き改めて次のように表わしておく。説明変数 x_1, \dots, x_k はその値が既知な確定変数であるが、 x_j はその値が未知な確定変数である。誤差項 u_j についての仮定は 2・1 式と同じ。

$$y_j = \alpha_0 + \sum_{i=1}^k \alpha_i x_{ij} + \beta \xi_j + u_j \dots \dots \dots (2 \cdot 9)$$

$$j = 1, 2, \dots, n$$

この場合には観測系列 y_j, x_j は 2・9 式、2・8 式という equation systems で規定される内生変数であるが、 x_j を v_j で代用して最小二乗法で 2・9 式の係数を推定したとすれば最小二乗推定値 \hat{y}_j はどのような分布に従うであろうか。

II errors in variables model

1 単純な場合について

まず次のような単純な場合を考察してみる。

$$y = \beta \xi + u \dots \dots \dots (2 \cdot 10)$$

$$x = \xi + v \dots \dots \dots (2 \cdot 8)$$

ただし、 u, v は正規変数ベクトルで

$$E(u) = E(v) = 0_n \quad E(u_j v_j) = 0 \quad (j \neq j')$$

u, v の共分散行列 Σ_{uu}, Σ_{vv} は

$$\Sigma_{uu} = \sigma_u^2 I_n \quad \Sigma_{vv} = \sigma_v^2 I_n$$

とする。

そこでやはり proxy variable x を使った時の β の最小二乗推定値 $\hat{\beta}$ は

$$\hat{\beta} = \frac{\sum_{j=1}^n y_j x_j / \sum_{j=1}^n x_j^2}$$

であるが、 $\hat{\beta}$ の確率密度関数やキーマントとして既に Richardson and De-Min Wu [7] Sawa [4] 及び Takouchi [6] により明らかにされている。 $\hat{\beta}$ のキーマントは [5] の結果を利用すれば直ちに導けるが、それによる $\hat{\beta}$ の期待値 $E(\hat{\beta})$ は

$$E(\hat{\beta}) = \beta 2\alpha g_n(\alpha) \quad \text{ただし、} n > 1$$

で、2・7式に対応する関係式が得られる。従って $\hat{\beta}$ の偏りは

$$E(\hat{\beta} - \beta) = -\beta [1 - 2\alpha g_n(\alpha)] \dots \dots \dots (2 \cdot 11)$$

上式の α および $g_n(\alpha)$ の定義は次の通りである。

線型回帰分析における観測誤差の影響と誤差の評価について

$$\sigma^2 \equiv \sum_{j=1}^n \xi_j^2 / 2\sigma^2$$

$$g_n(z) \equiv e^{-z^2} \sum_{k=0}^{\infty} \frac{1}{k!} \frac{z^k}{(n+2k)}$$

$\hat{\beta}$ の片冪 II 階矩 (M. S. E) は

$$E(\hat{\beta} - \beta)^2 = \beta^2 \left[1 + (2z^2 - 5z)g_n(z) + 2zg_{n-2}(z) - 2z^2g_{n+2}(z) \right] \\ + \frac{\sigma^2}{a_1^2} g_{n-2}(z) \dots \dots \dots (2 \cdot 12)$$

ただし $n > 2$

$$E(\hat{\beta} - \beta)^2 = \frac{(n-3)\beta^2}{2} + \left\{ \frac{\sigma^2}{a_1^2} - \beta^2 \left[(n-3)z + \frac{(n-2)(n-5)}{2} \right] \right\} g_{n-2}(z) \dots \dots (2 \cdot 13)$$

上記の $g_n(z)$ は confluent hypergeometric function ${}_1F_1(a, b, z)$ からの関係があるから

$$ne^z g_n(z) = {}_1F_1\left(\frac{n}{2}, \frac{n}{2} + \frac{1}{2}, z\right)$$

より $g_n(z)$ は z の偶関数である。

$$E(\hat{\beta} - \beta) = -\beta \left[1 - \frac{2z}{n} e^{-z^2} {}_1F_1\left(\frac{n}{2}, \frac{n}{2} + 1, z\right) \right] \dots \dots \dots (2 \cdot 14)$$

ただし

$$E(\hat{\beta} - \beta) = -\beta e^{-z} F_1\left(\frac{n}{2} - 1, \frac{n}{2}, z\right) \dots\dots\dots (2 \cdot 15)$$

ただし

$$E(\hat{\beta} - \beta)^2 = \frac{(n-3)}{2} \beta^2 + \frac{1}{n-2} \left\{ \frac{\sigma_1^2}{\sigma_2^2} - \beta^2 \left[(n-3)z + \frac{(n-5)(n-2)}{2} \right] \right\} \\ \times e^{-z} F_1\left(\frac{n}{2} - 1, \frac{n}{2}, z\right) \dots\dots\dots (2 \cdot 16)$$

なお(24)

$$E(\hat{\beta} - \beta)^2 = \frac{1}{n-2} \left(\frac{\sigma_1^2}{\sigma_2^2} + \beta^2 \right) e^{-z} F_1\left(\frac{n}{2} - 1, \frac{n}{2}, z\right) \\ + \beta^2 \left(\frac{n-3}{n-2} \right) e^{-z} F_1\left(\frac{n}{2} - 2, \frac{n}{2}, z\right) \dots\dots\dots (2 \cdot 17)$$

$\hat{\beta}$ の確率密度関数については [3] および [4] を参照されたい。(+) 2・15式・2・17式と合流型超幾何関数の性質を利用することにより $\hat{\beta}$ の相対的な偏りを M.S.E. をつぎの事柄を導くことができる。

- (1) σ_1 の相対的偏りは
 (2) $z > 0, n > 1$ ならば

$$F_1\left(\frac{n}{2} - 1, \frac{n}{2}, z\right) > 0$$

だから一般に

$$\frac{E(\hat{\beta} - \beta)}{\beta} < 0$$

である。

(ii) n を固定して $\hat{\beta}$ の相対的偏りを z の関数とみると、その絶対値は z が増加するにたがって減少する。何故ならば

$$r \equiv \sum_{j=1}^n \xi_j^2 / n\sigma^2 \quad \text{すなわち} \quad z = 2nr$$

$n \searrow \infty$

$$\frac{\partial}{\partial r} \left[e^{-z} F_1 \left(\frac{n}{2} - 1, \frac{n}{2}, z \right) \right] = -4e^{-z} F_1 \left(\frac{n}{2} - 1, \frac{n}{2} + 1, z \right) < 0$$

ただし、 $n > 1$

であるから

$$\frac{\partial}{\partial r} \left| \frac{E(\hat{\beta} - \beta)}{\beta} \right| < 0$$

$n \searrow \infty$ 。

よ

$$\lim_{r \rightarrow \infty} \frac{E(\hat{\beta} - \beta)}{\beta} = 0, \quad \lim_{r \rightarrow 0} \frac{E(\hat{\beta} - \beta)}{\beta} = -1$$

であるので、如何なる n に対しても $\hat{\beta}$ の相対的偏りは常に 0 と -1 の間にあることがわかる。従って $E(\hat{\beta})$ の符

号は β の符号と一致することになる。

これは合流型超幾何関数の積分表示式 2・18 式を使って 2・14 式を書き改めれば直ぐに導ける。

$$F_1\left(\frac{n}{2}, \frac{n}{2}+1, z\right) = \frac{n}{2} z^{-\frac{n}{2}} \int_0^z e^{t^{\frac{n}{2}-1}} dt \dots \dots \dots (2 \cdot 18)$$

$$\frac{E(\hat{\beta}-\beta)}{\beta} = -1 + z^{-1-\frac{n}{2}} e^{-z} \int_0^z e^{t^{\frac{n}{2}-1}} dt \dots \dots \dots (2 \cdot 19)$$

$$\lim_{r \rightarrow \infty} z^{-1-\frac{n}{2}} e^{-z} \int_0^z e^{t^{\frac{n}{2}-1}} dt = 1$$

だから

$$\lim_{r \rightarrow \infty} \frac{E(\hat{\beta}-\beta)}{\beta} = 0$$

(二) 2・19 式を使って、 n が充分大きい場合の、 $\hat{\beta}$ の相対的偏りの漸近的近似式を求めると

$$e^{-z^{-1-\frac{n}{2}}} \int_0^z e^{t^{\frac{n}{2}-1}} dt = \int_0^z e^t \left(1 - \frac{t}{z}\right)^{\frac{n}{2}-1} dt = \int_0^{2n\tau} e^{-t} \left(1 - \frac{t}{2n\tau}\right)^{\frac{n}{2}-1} dt$$

ただし、 $z = 2n\tau$

$\lim_{n \rightarrow \infty} \sum_{j=1}^n \frac{e^{t_j}}{n}$ が収束するとせば、 n が充分大きい場合には

線型回帰分析における観測誤差の影響と誤差の評価について

$$\int_0^{2n\tau} e^{-t} \left(1 - \frac{t}{2n\tau}\right)^{\frac{n}{2}-1} dt \sim \int_0^{\infty} e^{-(1+\frac{t}{2n\tau})^2} dt = \frac{\tau}{1+\tau}$$

従って

$$\frac{E(\hat{\beta} - \beta)}{\beta} \sim \frac{-1}{1+\tau}$$

$$\text{ただし, } \tau = \sum_{j=1}^n \sigma_j^2 / n.$$

また Richardson and De-Min Wu [3] は次のようにより詳しい漸近式を導いている。

$$\frac{E(\hat{\beta} - \beta)}{\beta} = \frac{-1}{1+\tau} \left[1 - \frac{2}{n} \frac{\tau^2}{(1+\tau)^2} + \dots \right]$$

(2) 上の平均二乗誤差 (M.S.E.) はつぎの通り

(1) 2・17 式より知られるように M.S.E. は n , β^2 , σ_1^2/σ^2 を含む 4 つの要因に依存している。

$$\frac{\partial E(\hat{\beta} - \beta)^2}{\partial \tau} = \frac{2n}{n-2} \left[-\frac{2}{n} \sigma^{-2} F_1 \left(\frac{n}{2} - 1, \frac{n}{2} + 1, \tau \right) - \frac{4(n-3)}{n} F_1 \left(\frac{n}{2} - 2, \frac{n}{2} + 1, \tau \right) \right] < 0$$

ただし, $n > 2$

であるから M.S.E. は τ の減少関数である。他方 β^2 , σ_1^2/σ^2 に関して増加関数であることは 2・19 式より明らかである。

(3) [3] には β の分散 $V(\hat{\beta})$ について次のような漸近式が導出されている。

$$V(\hat{\beta}) = \frac{1}{n-2} \left[\frac{\sigma^2/\sigma^2}{1+r} + \frac{\beta^2(1+r^2)}{(1+r)^2} \right] + \dots$$

↑ β の相対的偏りや M.S.E. の数値的な評価にあたっては〔3〕に掲げられている数値表が便利である。

2 外変数が追加された場合

観測値 (x_j, y_j) , $j=1, 2, \dots, n$ が 2・8 式および 2・9 式にしたがう場合を考察する。行列形式に書き改めると

$$y = Z\alpha + \beta\xi + u \dots \dots \dots (2 \cdot 20)$$

$$x = \xi + v \dots \dots \dots (2 \cdot 21)$$

$$\text{ただし, } Z = \begin{bmatrix} 1 & z_{11} & \dots & z_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nk} \end{bmatrix} \quad \alpha = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_k \end{bmatrix}$$

Z の階数は $(k+1)$ とする。

2・20 式で v の代わりに x を proxy variable に使って β の最小二乗推定値 $\hat{\beta}$ を求めたとすれば、 $\hat{\beta}$ はどんな分布に従うだろうか。以下に述べるようにこの問題は二の口の 1 で扱った問題に還元できるので、そこで導出した結果を少し修正するだけで処理できる。

↑ β を書き改めると

$$\hat{\beta} = \frac{x'Qy}{x'Qx} \dots \dots \dots (2 \cdot 22)$$

$$\text{ただし, } Q = I_n - Z(Z'Z)^{-1}Z'$$

他方を $\alpha_1, \dots, \alpha_k$ に回帰させた時の最小二乗法的な回帰式を

$$\hat{\xi} = Z\alpha + e \dots \dots \dots (2 \cdot 23)$$

とせよ。 α は次のように表わせる。

$$\hat{\alpha} = \frac{(e' + v')Q(\beta e + u)}{(e' + v')Q(e + u)} \dots \dots \dots (2 \cdot 24)$$

N の置換は $k+1$ である。 $Z'(Z'Z)^{-1}Z'$ は置換は $k+1$ のイデンプotent matrix) である。

$$Z'(Z'Z)^{-1}Z' = P' \begin{bmatrix} I_{k+1} & 0 \\ 0 & 0 \end{bmatrix} P$$

ただし, P は n 次の直交行列, I_{k+1} は $(k+1)$ 次の恒等行列.

従って $a'Qb$ なる一次形式は置換は $N \cdot N$ での N の $(k+1)$ 成分である。^(*)

$$a'Qb = c'_{(\omega)} d_{(\omega)} \dots \dots \dots (2 \cdot 25)$$

ただし, $c_{(\omega)}, d_{(\omega)}$ は $c \equiv P\alpha$, $d \equiv P\beta$ をそれぞれ次のような二つの部分に分割してえられる ($n-k-$

1) 次の部分ベクトルである。

$$\begin{aligned} c_{(\omega)} &= \begin{bmatrix} c_{(\omega)} \\ \vdots \\ c_{(\omega)} \end{bmatrix} & c_{(\omega)} &= \begin{bmatrix} c_{k+1} \\ \vdots \\ c_n \end{bmatrix} \\ d_{(\omega)} &= \begin{bmatrix} d_{(\omega)} \\ \vdots \\ d_{(\omega)} \end{bmatrix} & d_{(\omega)} &= \begin{bmatrix} d_{k+1} \\ \vdots \\ d_n \end{bmatrix} \end{aligned}$$

2・25式を使って2・24式を更に変形すると

$$\hat{\beta} = \frac{h_{(2)}k_{(2)}}{h_{(2)}h_{(2)}} \dots \dots \dots (2 \cdot 26)$$

ただし、 $h \equiv P(e+v) \equiv [h_{(1)}, h_{(2)}]^T$

$k \equiv P(\beta e + u) \equiv [k_{(1)}, k_{(2)}]^T$

ベクトル h 、 k は、お互いに独立に正規分布にしたがう確率ベクトルであって、期待値が

$$E(h) = P e \equiv [\varphi_{(1)}, \varphi_{(2)}]^T$$

$$E(k) = \beta P e \equiv [\beta \varphi_{(1)}, \beta \varphi_{(2)}]^T$$

共分散行列 Σ_{hh} , Σ_{kk} が

$$\Sigma_{hh} = \sigma_1^2 I_n \quad \Sigma_{kk} = \sigma_2^2 I_n$$

であるから

$$k_{(2)} = \beta \varphi_{(2)} + u \dots \dots \dots (2 \cdot 27)$$

$$h_{(2)} = \varphi_{(2)} + v$$

ただし、 u 、 v は(2・8)、(2・10)の時と同じ確率ベクトル。

となる。従って $\hat{\beta}$ は2・27式の β を最小二乗法で推定してゐることにほかならぬ。

$$\varphi_{(2)} \varphi_{(2)}' = e' Q e' = e^{(\prime)} e$$

であるから二の(1)で述べた結果を使えば、 $\hat{\beta}$ は次のような偏りをもつことがわかる。

線型回帰分析における観測誤差の影響と誤差の評価について

$$\frac{F(\hat{\beta} - \beta)}{\beta} = -\exp(-z) F_1 \left(\frac{n-k-1}{2} - 1, \frac{n-k-1}{2}, z \right) \dots \dots \dots (2 \cdot 28)$$

ただし、 $n > k + 2$

$$z \equiv e/2\sigma_1^2 \dots \dots \dots (2 \cdot 29)$$

他方 α の M.S.E は $\sigma^2 \cdot 19$ 式あるいは $\sigma^2 \cdot 17$ 式の n を $n-k-1$ 、 β' 、 α を $\sigma^2 \cdot 29$ 式で置換するだけで評価できる。
 e は α の回帰式 $2 \cdot 22$ 式の残差平方和だから、 k 個の説明変数による決定係数を R^2 とおけば

$$\alpha = \frac{\sum_{j=1}^n (\xi_j - \bar{\xi})^2 (1 - R^2)}{2\sigma_1^2}$$

従って、その他の条件を一定とすれば R^2 が大きい程 $\hat{\beta}$ の相対的偏りの絶対値は増大することになる。たとえ

ば

$$n = 11, k = 2, \sum_{j=1}^n (\xi_j - \bar{\xi})^2 / (n-1)\sigma_1^2 = 4$$

とする。[σ] の数値表によれば、 $R^2 = 0, 0.5, 0.8$ の時の $\hat{\beta}$ の相対的偏りは、それぞれ $-0.14, -0.28, -0.47$ となる。

$$\text{和}(\sigma^2) : F(a, b, z) \equiv \sum_{k=0}^a \frac{(a)_k}{k!} \frac{z^k}{(b)_k} = \sum_{k=0}^a \frac{a(a+1) \dots (a+k-1)}{k(b+1) \dots (b+k-1)} \frac{z^k}{k!}$$

(2)(c) 合流型超幾何関数としてこの次の関係式を利用する。

$${}_2F_1(a+1, b+1, z) = b[F_1(a+1, b, z) - F_1(a, b, z)]$$

$$\left(\frac{N-2}{N}\right)u^2 = \left(\frac{N-4}{2} + \frac{N-4}{2}\right) \left[\frac{N}{2} - 1, \frac{N}{2}, \frac{N}{2} \right] + \frac{N}{2} \left[\frac{N}{2} - 2, \frac{N}{2}, \frac{N}{2} \right]$$

(4) いずれも

$$y = \alpha + \beta \xi + u, \quad x = \xi + v$$

というモデルのもとでの u, v の分布である。

(5) $x = \beta \xi + v$ であれば $\beta \xi = \eta$ とおくとにより

$$y = \alpha_0 + \sum_{i=1}^k \alpha_i x_i + \left(\frac{\beta}{\gamma}\right) \eta + u$$

$$x = \eta + v$$

となつて 2・20 式、2・21 式と同じ形式のモデルになる。

(6) $a = P'c, b = P'd$ ならば

$$a'Qb = c'P[L_n - Z(Z'Z)^{-1}Z']P'd = c'd - c' \begin{bmatrix} I_{k+1} & O \\ O & O \end{bmatrix} d = c'(a)d(a)$$

(7) 2・25 式より

$$\varphi'(a)\varphi(a) = c'Qc = \xi'Q\xi = c'e$$

三 観測誤差の評価——農業生産統計を対象にして

proxy variable が 2・8 式のようなモデルに従うものとすれば、観測誤差に関する情報をもつことが望ましい。たとえば大雑把にでも a の大きさがわかれば、それを利用して β のより偏りの小さい最小二乗推定値を得ることができる。

最近、東南アジア諸国の農産物供給を対象にした計量経済学的な研究が発表されているが、Behrman [1] の計
 観測誤差分析における観測誤差の影響と誤差の評価について

測にみるように農業統計の不備と不正確さが一つの障害になっているように思える。生産統計の多くは表式調査や素朴な面接調査を通じて作成されているので、調査誤差の影響を無視できない。第三節では既成の統計データをもとにして調査誤差を暫定的に評価しようとする試みを述べたものである。事例としてタイ、フィリピンの米の生産統計を扱ってある。

(一) 農業センサスとの比較

Behrman〔1〕はタイの主要農産物の生産量統計の信頼性を評価するに当たって、次の三通りの方法を採用して

第1表 ベンガル州のジャート生産量に関する諸統計の比較 (1944/45 および 1945/46)

(単位: 1,000 bales)

	1944/45	1945/46
1. 消費量		
{ 工場で	6,000	6,308
{ 移出	1,050	2,213
{ 村で	600	600
2. 計	7,650	9,121
3. 前年よりの繰越し	324	697
4. 移入	598	862
5. (バランス) 消費統計による生産量	6,728	7,562
6. 農林省推計 (悉皆)	4,895	6,304
7. 標本調査	6,480	7,540
8. 6—5; 5 に対する%	-27.2	-16.6
9. 7—5; 5 に対する%	-3.6	-0.3

資料: Mahalanobis, P. C. and Lahiri, D. B., "Analysis of errors in censuses and Survey with special reference to experience in India," *Bulletin of the International Statistical Institute*, Vol. 38, Part 2, 1961.

いる。(イ)一つは農林省の生産量統計から誘導される消費量を別の消費量統計と比較する方法。(ロ)その統計が表式調査や面接調査によっている場合には、標本実測調査による推計値と比較する。たとえば単位面積当たり収量を坪刈りによる推計値と比較してチェックする方法。第一表はP. C. Mahalanobis がインド・ベンガル州のジャートの生産量について三種の統計数字を比較した結果を紹介したものであるが、上記の方法の典型的な例である。

農林省推計はジュート畑を悉皆的に見回って検見した結果にもとづいている。第一表の七の標本調査の推計量は坪刈りと面積実測をもとにしたものである。(b) 第三の方法は特定年次に限られるが、農業センサスによる結果との比較によるものである。

以下述べるのはこの方法によって調査誤差を検出しようとする試みである。事例としてタイ、フィリピンの米生産統計を扱うので、対象とする統計系列の概略を述べておきたい。フィリピンの一九五三年までの公式統計は *Philippine Agriculture Statistics, Vol. I, II* に集録されているデータが唯一のものであるが、これはそれまで各省各部署が業務統計として発表していたものを農業天然資源省 (DANR) が統一集大成したものである。 *Philippine Agriculture Statistics* の一九四八年度 (作物年度) の米の作付面積、生産量、単位面積当たり収量を一九四八年に行なわれた農業センサスの結果と比較するわけであるが、DANR が行なった生産統計の編集手続きは明らかでない。農業センサスはセンサス局によるもので、一九四八年度 (一九四七・七・一〜一九四八・六・三〇) の生産活動を対象にして農場面積が一〇ヘクタール以上の面積をもつ農家を調査している。一九五四年からは DANR によって標本面接調査による作物家畜調査が開始され、この調査から各作物に関する生産統計や家畜飼育頭数などのカルトな統計が作成されるようになっていいる。

タイの農産物に関するカルトな生産統計は *muban-tambal-amphur-changwad* という行政機構を通り、農林省の *Rice Department* や *Division of Agricultural Economics* で集計されるわけであるが、米の生産統計は *muban (village)* の長である *phuyai-bau* が彼の *village* について述べた報告数字が基礎になっている。タイについては一九六〇年度の米の生産量統計を農業センサスの結果と比較する。

(三) 比較の方法

二つの統計系列を比較するには、農家あるいは村といった最終の調査単位で行なうのが望ましいが、公表されているデータはたかだか province といった大きな行政区轄に関するものであるから、実際には province レベルの比較しかできない。フィリピンでは一九四八年度について province レベルの比較を、タイでは一九六〇年度について changwad レベルの比較をすることが出来る。

ここでは province (changwad) 間で二つの統計系列がどのような相違をしめすかに注目する。なお、フィリピンにおける province の数は五〇で、これらの province が九個の region に分類されている。タイでは七一の changwad が四個の region に分類されている。

ところで比較した結果の判定基準であるが、たとえば米の作付面積を province レベルで比較した場合、二つの統計系列の間の一致度が高ければカレントな生産統計は一致度の低いものに比べて better であると考ええる。センサスの結果が正確であるとする保証はないから、一致度の高いことはカレントな生産統計の正確さを意味するものではないが、一致度が高ければそれだけ信頼しうるものと認めようという考えである。

1 不一致度の尺度

二つの統計系列の一致度の測り方にはいろいろあるが、ここでは次のように表わす。i 番目の province (changwad) のカレントな生産統計によるデータを X_i で、農業センサスの結果を Y_i で表わす。

$$g = \frac{1}{n} \sum_{i=1}^n (X_i - Y_i)^2 \dots \dots \dots (3.1)$$

ただし、 n は province (changwad) の数。

3・1式の g は面接調査における回答誤差の変動を分析するのによく使われる統計量であるが、その場合には X_i は同じ調査条件のもとで面接を二回反覆してえられる被調査者の回答を意味する。米国のセンサス局が最近統計調査の信頼性に関する組織的研究を実施しているが、3・1式の g を gross difference rate と呼んでいる。またこの統計量は予測の正確度を測る場合にも使われている。この時には X_i が予測値を、 Y_i が実際値を表わす。統計系列の比較を年次ごと作物ごとあるいは国ごとに行なうてその結果を比べようとするには、3・1式の g を指数形式にしておくのが便利である。Theil〔6〕の手法に従って次に指数化する。Theilは“inequality coefficient”(不平等度)と呼び U という記号を使っている。

$$U = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - Y_i)^2} \sqrt{\left[\frac{\sum_{i=1}^n X_i^2}{n} + \frac{\sum_{i=1}^n Y_i^2}{n} \right]} \dots\dots\dots (3 \cdot 2)$$

U が大きい程二つの統計系列の間の不一致度は大きいと判断するわけだが、この根拠は U —指数の幾何学的性質にもとづくものである。⁽¹⁾

ところで3・1式の g には種々の要因が関与している。それぞれの統計系列を生成する調査手続きや調査環境の相違、標本抽出法が使われておれば抽出誤差、面接調査であれば被調査者の回答誤差、その他種々の要因が影響している。米国センサス局の回答誤差の研究では g の要因分析をするために、センサス局モデル⁶が設定されているが、ここでもそれに相当する errors model を前提することになる。その議論は後回しにして、まず3・1式の g をTheil〔6〕に従って次のように分解してみる。

$$g = (\bar{X} - \bar{Y})^2 + (S_x - S_y)^2 + 2(1-r)S_x S_y \dots \dots \dots (3 \cdot 3)$$

ただし、 r は相関係数、 S_x, S_y はそれぞれ標準偏差。

3.3 式の右辺の各項を $\left\{ \frac{\sum X_i^2}{n} + \frac{\sum Y_i^2}{n} \right\}$ で強じたものをそれぞれ U_M^2, U_S^2 であり、 U_C^2 とすれば、3.2

では次のように表わされる。

$$U^2 = U_M^2 + U_S^2 + U_C^2$$

Theil [6] は U_M^2, U_S^2 であり、 U_C^2 ありきの性質を論述している。

$$U_M = \sqrt{U_M^2} : \text{中心値の不等による偏不等度}$$

$$U_S = \sqrt{U_S^2} : \text{変動の不等による偏不等度}$$

$$U_C = \sqrt{U_C^2} : \text{共変関係の不等による偏不等度}$$

$$U_M = U_M^2 / U^2 : \text{bias proportion}$$

$$U_S = U_S^2 / U^2 : \text{Variance "}$$

$$U_C = U_C^2 / U^2 : \text{Covariance "}$$

U_M が調査手続きや調査環境などの相違による系統的な不整合を表わす成分であることは明らかであろう。もし二つの統計系列が、同じ調査手続きに従って類似の調査環境のもとで作成されたものとするれば、bias proportion の U_M は $U_M/0$ となろう。他方、調査手続きや調査環境が均しくても U_S や U_C が消失することはないだろう。

二つの統計系列は同じ項目を観測しているから、調査がお互いに無関係に行なわれたとしても統計系列の間にはあ

第2表 タイ、フィリピンの米生産統計のU-指数

	タイ (1960年度)			フィリピン(1948年度)			日本 (1965年度) 作付面積
	作付面積	生産量	収量	作付面積	生産量	収量	
U	0.0572	0.0685	0.0985	0.1652	0.2001	0.1160	0.0539
U_M	0.0160	0.0011	0.0228	0.0316	0.0549	0.0506	0.0471
U_S	0.0010	0.0007	0.0030	0.0317	0.0130	0.0020	0.0035
U_C	0.0549	0.0047	0.0958	0.1590	0.1920	0.1044	0.0259
U^M	0.0785	0.0003	0.0536	0.0365	0.0752	0.1904	0.7646
U^C	0.0003	0.0001	0.0009	0.0367	0.0042	0.0002	0.0041
U^S	0.9212	0.9996	0.9455	0.9268	0.9206	0.8094	0.2313
n	50			71			45

注 1. タイの収量は収穫面積当たり、フィリピンの収量は作付面積当たりの収穫量である。

2. 日本では農業センサスの収穫面積 (X) を作物調査の作付面積 (Y) と比較。

る直線的な相関関係が成立する。実際の観測過程には偶然変動を含めて種々の観測誤差が混入するから、二つの観測系列間の相関関係は誤差の程度に応じて不完全なものになる。この不完全さは U_C あるいは U^C を通じて測ることができる。 U_M あるいは U^M は観測誤差の大きさの差を表わすものであるが、次に述べる計測例ではこの成分の寄与は無視しうる程度である。

2 計測結果

米の作付面積、生産量および収量について二つの統計系列の間の U -指数を計算してみる。各成分も併せて第二表に掲げてあるが、参考のためにわが国の水稲についての U -指数を掲げてある。日本の U -指数は一九六五年度について農業センサスの収穫面積を作物調査による作付面積と比較したもので、北海道を除く四五都府県を対象にしている。ただし、農業センサスの結果を X^* で、作物調査の推計値を X で表わしている。タイやフィリピンの U -指数を検討するための対照として役立つであろう。

作付面積の U -指数は日本とタイの間ではその差は僅少であるが、指数の成分をみると著しい相違のあることがわかる。タイで

は bias proportion の U_M が 0.0785 であるのに対して、日本の U_M は 0.7646 となっている。日本の場合、作物調査による推計は実測によっているので、二つの統計系列間では調査方法の相違による影響が大きく、この相違が規則だった系統的誤差を生成する主たる源泉になっている。しかも、統計系列間の共変関係の程度が高いので、その結果として bias proportion はタイに比べて著しく増大するわけである。タイの場合、調査手続きでは、一方は面接調査、他は表式調査といった差があるわけだが、その差は β を構成する主要因にはなっていない。系統的な不整合を生成するにしてもその程度は小さいか、あるいはその現われ方が *changwad* を通じて一様ではない。

タイの場合、作付面積だけでなく生産量、収量についても共変関係の不完全さが、 U 指数の主要因になっている。この点はフィリピンでも同じだが、タイに比べて U^2 がより大きい点が特長である。

iii) errors model の設定

ii) の 2 でみたようにタイ、フィリピンの指数は主として二つの統計系列間の共変関係の乱れによっている。この傾向は米以外の作物についても同じであろう。作物間でどのような差が出てくるかを検討するのも興味ある問題であるが、ここでは共変関係の不完全さを別の観点から分析してみたい。

各 province (*changwad*) の真の米作付面積を Z_i , $i=1, 2, \dots, n$ とおくと、 $X_i - Z_i$ および $Y_i - Z_i$ は一般に Z_i の大きさに依存するであろう。この関係を次のようにモデル化する。

$$X_i = (\alpha + \Delta\alpha_i) Z_i$$

$$Y_i = (\beta + d\beta_i)Z_i$$

$$i = 1, 2, \dots, n$$

$d\alpha_i, d\beta_i$ はいずれも期待値がゼロ、ある分散をもつ確率変数。つまり X_i, Y_i の Z_i に対する偏倚率はそれぞれ期待値 α, β をもつ分布に従って変動するものと想定する。これらの期待値は地域ごとに異なるのであるが、まず各地域を通じて一定とみなして議論を進めてゆく。

(X) 系列についてみると

$$\log X_i = \log_e X_i = \log_e Z_i + \log_e \alpha + \log_e \left(1 + \frac{d\alpha_i}{\alpha}\right)$$

ゆえに

$$\log_e \left(1 + \frac{d\alpha_i}{\alpha}\right) \sim \frac{d\alpha_i}{\alpha} - \frac{1}{2} \left(\frac{d\alpha_i}{\alpha}\right)^2$$

と近似し、ゆえに $d\alpha_i/\alpha$ の分散は province (changwad) を通じて一定と仮定して C_α^2 とおけば

$$E \left\{ \log_e \left(1 + \frac{d\alpha_i}{\alpha}\right) \right\} \sim 1 - \frac{1}{2} C_\alpha^2$$

$$V \left\{ \log_e \left(1 + \frac{d\alpha_i}{\alpha}\right) \right\} \sim C_\alpha^2$$

$$\text{ただし, } C_\alpha^2 = V \left(\frac{d\alpha_i}{\alpha} \right) \quad i = 1, 2, \dots, n$$

C_α^2 は $d\alpha_i$ の relative variance にほかならない。

したがって $\{X_i\}$ 系列に対して次のような errors model が導かれる。

$$\log_e X_i = \log_e Z_i + \phi + u_i, \dots \dots \dots (3 \cdot 4)$$

$$\text{ただし, } \phi = \log_e \alpha - \frac{C_a^2}{2}$$

$$E(u_i) = 0 \quad E(u_i^2) = C_a^2$$

同じ要領で $\{Y_i\}$ 系列についても次のような errors model を設ける。

$$\log_e Y_i = \log_e Z_i + \varphi + v_i, \dots \dots \dots (3 \cdot 5)$$

$$\text{ただし, } \varphi = \log_e \beta - \frac{C_\beta^2}{2}$$

$$E(v_i) = 0 \quad E(v_i^2) = C_\beta^2$$

$$C_\beta^2 = V\left(\frac{\Delta \beta_i}{\beta}\right) \quad i=1, 2, \dots, n$$

統計系列の $\{X_i\}$, $\{Y_i\}$ が表式調査や面接調査のような観測過程を通じて生成される場合には、 α , β は一般に 1 から乖離した値をとるたさう。

3・4 式、3・5 式の ϕ , φ はそれぞれの観測過程の系統的な偏りを表わす項である。 C_a^2 , C_β^2 は真値 Z_i に対する偏倚率の province (changwad) 間変動を表わす relative Variance であるが、 α の変動は調査手続きや被調査者側の知識状態のあいまいさによって生ずるものであるから、観測過程の不規則性または不安定性をしめす尺度とみることができる。 C_a , C_β をそれぞれの統計系列の不安定度と呼ぶことにする。

わが国の米作付面積のように $\{Y_i\}$ 系列が標本実測調査で推定されている場合を考えてみよう。

推定値が不偏性をもてば $\beta=1$ とおける。 C_p^2 は標本抽出誤差による Y_i の relative Variance であるから⁽²⁾、わが国の果段階の水稲作付面積の C_p^2 は大体 0.5% とみてよい。あとでしめすように農業センサスによる水稲収穫面積の不安定度 C_p^2 は約 4% とみられる。この場合には悉皆調査であるから不安定度の主源泉は回答誤差における不規則性である。

C_p および C_p^2 を求める手続きは次の通りである。

$$\log_{10} X_i = x_i, \quad \log_{10} Y_i = y_i,$$

よって

$$y_i - x_i = d_i$$

を求めると、3・4式、3・5式から d_1, d_2, \dots, d_n は3・6式にしめすところの期待値 d 、分散 σ^2 をもつ母集団分布からの任意標本とみなすことができる。

$$\bar{d} = \bar{y} - \bar{x} = (\log_{10} \beta - \log_{10} \alpha) - \frac{M}{2} (C_p^2 - C_x^2) \dots \dots \dots (3 \cdot 6)$$

$$\sigma^2 = M^2 (C_p^2 + C_x^2 - 2C_{xp})$$

$$\text{ただし, } C_{xp} = E(\Delta \beta, \Delta \alpha) / \alpha \beta$$

$$M = 0.43429$$

\bar{d}, σ^2 は d_1, d_2, \dots, d_n より計算する標本平均 \bar{d} および標本分散 S_d^2 で推定できる。第三表の(1)、(3)欄にそれ

第3表 不安定度の計算 (米の作付面積)

	\bar{d}	β/α	S_d	不安定度	備 考
日 本	0.04960	1.120	$180,673 \times 10^{-7}$	41.3×10^{-3}	Y_j : 実測調査 X_j : センサス
タ イ	0.02013	1.0480	$7,428 \times 10^{-4}$	110.3×10^{-3}	Y_j : センサス X_j : 表式調査
フィリピン	0.01275	1.030	$16,966 \times 10^{-4}$	276.2×10^{-3}	Y_j : センサス X_j : DANRによる。

縦型回帰分析における観測誤差の影響と誤差の評価について

二八

らの計測結果が掲げられている。 \bar{d} , S_d を 3・6 式の \bar{d} , s^2 に代入することにより R^2 や不安定度を近似的に評価できる。

1 計測結果

日本の場合には既述の通り $C_2 = 5 \times 10^{-3}$, $C_{a^2} = 0$ とおけるから

$$C_a^2 = \frac{\sigma^2}{M} - C_a^2$$

よって $180,673 \times 10^{-7}$ と推定されるから、農業センサスの不安定度 (水稲面積に関する)

$C_a = 41.3 \times 10^{-3}$ と推定される。

$$\log_{10} \left(\frac{\beta}{\alpha} \right) = \bar{d} + \frac{M}{2} (C_a^2 - C_a^2)$$

よって $\log_{10} \left(\frac{\beta}{\alpha} \right) = 4924 \times 10^{-5}$ と推定される。

タイ、フィリピンの場合には $C_{a^2} = 0$ とおくことは正しくないであろう。それぞれの統計系列における観測過程がどのように関連しているかによるわけだが、 $C_{a^2} > 0$ とみなすのが普通であろう。また $C_a > C_a$ と考えられるが、ここでは

$$C_a = C_a = C, C_{a^2} = 0$$

とこの仮定を置いて、 $S_d / \sqrt{2M}$ で共通の不安定度を求めることにする。 $C_{a^2} > 0$ であれば不安定度を過小に評価している結果になる。第三表の(4)欄は以上のような手続きで

求めたものである。タイ、フィリピンの β/α は $\bar{d} = \log_{10} \left(\frac{\beta}{\alpha} \right)$ にしたがって求めた。

第二表の U_i で見ても明らかであるが、タイ、フィリピンに比べてわが国の農業センサスにおける水稻収穫面積の不安定度は小さい。過小率の全国平均は 〇・八九であるが、県間変動は relative Standard deviation でみて約四%である。タイ、フィリピンのカレントな生産統計における米作付面積の偏倚率を知る客観的なデータはないが Behrman [1] の計算によると一九六二年度のタイの米生産量の過小率は八〇〜八五%とみられている。フィリピンの過小率は村岡 [2] によると一九五一年から一九六五年までの年平均で八六%と推計されている。タイ、フィリピンのカレントな生産統計で問題になるのは全国平均の偏倚率よりも、偏倚率がしめす province (changwad) 間の変動の大きさではないかと思う。先に導いた不安定度でみるとタイが日本の三倍、フィリピンが七倍とかなりの大きさになっている。

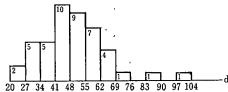
第一図は参考として各国の \bar{d}_i に関するデータをヒストグラムにまとめたものである。

2 地域差について

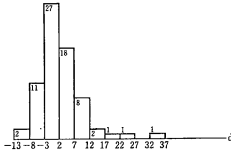
これまで述べた不安定度の計測では各 province (changwad) の X_i , Y_i の偏倚率の期待値は一定とみたが、これらの期待値は地域間で異なるだろう。その結果 3・6 式の σ_{ii} は地域間で変動する。たとえばフィリピンにおいて province をミンダナオとその他の二つの地帯に分けて、地帯別に \bar{d}_i を求めると、ミンダナオが $\bar{d}_i = 0.14965$ 、その他が $\bar{d}_i = -0.02232$ となる (第一図の (i) を参照)。これは極端なケースであるが、地域差の存在は無視できない。この影響は S_{ii} の中に混入するから不安定度を推定するためには地域差の影響をできるだけ除去しておくの

第1図 米作付面積に関する d_t のヒストグラム

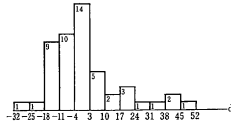
(イ) 日本(横軸の単位: $\frac{1}{100}$)



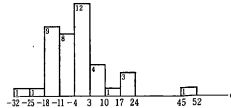
(ロ) タイ(横軸の単位: $\frac{1}{100}$)



(ハ) フィリピン(横軸の単位: $\frac{1}{100}$)



(ニ) フィリピン(ミンダナオを除く)(横軸の単位: $\frac{1}{100}$)



が望ましい。フィリピンでは九地域に、タイでは四地域に分類⁽³⁾されているのでこれらの地域を対象にしてその影響を検討してみる。

地域 k に属する province (changwad), i に関する測定値を X_{kij} , Y_{kij} Δ 区別 Δ

$$d_{kij} = \log_{10} Y_{kij} - \log_{10} X_{kij}$$

とおく。そして地域差の影響を次のようにモデル化する。

$$d_{kij} = \delta + \tau_k + \epsilon_{kij} \dots \dots \dots (3 \cdot 7)$$

ただし, $h=1, 2, \dots, m$

$$j=1, 2, \dots, n_h$$

3・7式の δ は二つの統計系列の log-difference の全国平均、 τ_h は地域差をしめす項であって $\sum_{h=1}^m \tau_h = 0$ とおく。
 ϵ_{hj} は残差項で $E(\epsilon_{hj})=0$, $E(\epsilon_{hj}^2)=\sigma^2$ と仮定する。

地域間変動および σ^2 に関する推定は次のような分散分析表に従って行なうことができる。

分散分析表

要因	平方和	自由度	平均平方	期待値
地域間	$\sum_h n_h (\bar{d}_h - \bar{d}_{..})^2$	$m-1$	S_b^2	$\sigma^2 + \sum_h n_h \tau_h^2 / (m-1)$
地域内	$\sum_h \sum_j (\bar{d}_{hj} - d_{hj})^2$	$n-m$	S_w^2	σ^2
総計	$\sum_h \sum_j (d_{hj} - \bar{d}_{..})^2$	$n-1$	S_d^2	

$$\text{ただし, } \bar{d}_h = \frac{1}{n_h} \sum_{j=1}^{n_h} d_{hj}$$

$$\bar{d}_{..} = \frac{1}{n} \sum_h \sum_j d_{hj}$$

3・7式と3・4式、3・5式および3・6式との関連は次の通りである。 X_{hj} の真値 Z_{hj} に対する偏倚率が従う分布の平均を a_h , relative Variance を C_h^2 , 一方 Y_{hj} の偏倚率の分布の平均を β_h , relative Variance

を C_{β}^2 とする。

$$\hat{\sigma}^2 = \log_{10} \left(\frac{\beta}{\alpha} \right) - \frac{M}{2} (C_{\beta}^2 - C_{\alpha}^2)$$

ただし、 β は β_A の、 α は α_A の幾何平均……………(3・8)

$$r_A = \log_{10} \left(\frac{\beta_A}{\alpha_A} / \frac{\beta}{\alpha} \right)$$

$$V(r_{Aj}) = \sigma^2 = M(C_{\beta}^2 + C_{\alpha}^2 - 2C_{r_A})$$

3・8式にみるように3・7式の地域差を表す項は β_A と α_A の間の log-difference として定義されるものである。

第四表および第五表はタイ、フィリピンに関する分散分析表を掲げたものである。有意水準1%で地域差の有意性検定をすれば、タイでは F-value が五・七一で有意差が認められる。フィリピンについてはミンダナオとその他に分けて地域差をテストした結果が第五表であって明らかに差が認められる。然しミンダナオを除いた八地域間では有意な差は検出できなかった。

まあと同じように

$$C_{\alpha} = C_{\beta} = C, \quad C_{r_A} = 0$$

の仮定を置けば分散分析表で求めた S_{α}^2 を用い $S_{\alpha} / \sqrt{2M}$ で不安定度 C を推定できる。地域差が消去されているので、三の目の1で求めた不安定度よりも小さく出てくる。地域差を除去したあとの不安定度は、タイが 110×10^{-4} 、フィリピンが 255×10^{-4} である。それにしてもわが国の農業センサスの四%と対照するとかなり大きい。日本の水

第4表 分散分析表 (タイの米作付面積, 1960年度)

要因	平方和	自由度	平均平方	
地域間	0.0786,754	3	0.0262,251	$F=5.71$
地域内	0.3075,265	67	0.0045,899	$F_{0.05}=2.75$ $F_{0.01}=4.10$
総計	0.3862.019	70	0.0055,172	

第5表 分散分析表 (フィリピンの米作付面積, 1948年度)

要因	平方和	自由度	平均平方	
地域間	0.2365,876	1	0.2365,876	$F=9.67$
地域内	1.1739,053	48	0.0244,564	$F_{0.05}=4.04$
総計	1,4104,928	49	0.0287,856	$F_{0.01}=7.19$

種にみるような例は特異かも知れない。農家自身水田面積について明確なデータをもっているので、農業センサスにおける農家の回答は実測値より偏倚しても、その偏倚の仕方は規則的である。タイの(X)系列は村長による村レベルの報告にもとづいているから、不安定度が日本に比べて大きいという事実は、偏倚の仕方がより不規則であるということにほかならないが、これは報告調査のもつ一般的な欠陥ともみられる。観測過程として恣意的な要素が多く、各種の調査誤差に曝され易いので、偏倚の仕方が単一のパターンをもたぬのであろう。三の(一)で指摘したが、フィリピンの一九四八年度の(X)系列はデータ・ソースが唯一でなくDANKによって編集されたものである。フィリピンの不安定度が二五%と大きいのはそれが関係しているものと思われる。

3 不安定度の役割

終わりに三の(一)のところで述べた統計データの評価の方法(一)との関連を明らかにしておくたい。

方法(一)によって主要作物に関するカレントな生産量統計の偏倚率がえられる。三の(二)および(三)で述べた手続きによって生産量のU-指数と近似的な不安定度がえられるが、不安定度は方法(一)で求めた偏倚率に対す

標準偏差の役割をもつものとみられる。偏倚率はいくつかの仮定のもとで導かれるものであるから、その正確度は客観的に評価し難い。客観的な評価は望めないにしても主観的な判断にもとづく評価は行なわねばならない。偏倚率の幅についての知識がないとすれば、三の(三)で求めた不安定度が手掛りになる。たとえば方法(1)によってタイのある年度の米生産量の偏倚率が八五%、すなわち $\alpha = 0.85$ と定められたとしよう。もし

$$C_a = 10\%, n = 71 \text{ (changwad 6 歳)}$$

とすれば、 α の従う分布の relative な標準偏差は ⁽⁴⁾

$$\frac{C_a}{\sqrt{n}} = \frac{0.10}{\sqrt{71}} = 0.0119$$

従って α の標準偏差は $0.85 \times 0.0119 = 0.01$.

α に対する信頼係数九五%の信頼区間は

$$0.85 - 0.02 \wedge \alpha \wedge 0.85 + 0.02$$

で、 $0.83 \wedge \alpha \wedge 0.87$ となる。この信頼区間は n についての事前的な確信度を表わしたものであるから、方法(1)によって生産量に関する客観的なデータが入手されたならばそれを使って修正すればよい。上記の信頼区間はそれまでの暫定値として意味をもつであろう。

(注一) 観測値の組 $(X_1, \dots, X_n), (Y_1, \dots, Y_n)$ を表わす n 次元空間の点を P とする(次頁の図を参照)。数分 OP, OQ は PQ の線分を n 等分する $|OP|, |OQ|, |PQ|$ となる。

$$|OP| = \sqrt{\sum X_i^2}, |OQ| = \sqrt{\sum Y_i^2}, |PQ| = \sqrt{\sum (X_i - Y_i)^2}$$

だから

$$U = \frac{|PQ|}{|OP| + |OQ|}$$

点Pと点Qの距離が二点の長さより小さければ一致性を認めるとみなすわけである。
 $\beta = 0 \leq U \leq 1$ である。

$$(2) \quad V(Y_i) = Z_i V(\Delta\beta_i).$$

$$\beta = 1 \text{ である}$$

$$\frac{V(Y_i)}{Z_i^2} = V\left(\frac{\Delta\beta_i}{\beta}\right) = C_{\beta^2}$$

(3) 地域分類 (諸県の統計は province である) 長changwud の諸級)

Ilocos (5) Cagayan (4) Central Luzon (7) Southern Tagalog (8)

Bicol (6) Eastern Visayas (4) Western Visayas (6) N. & E.

Mindanao (6) S. & W. Mindanao (4)

Central Region (35) North-eastern Region (15) Northern Re-

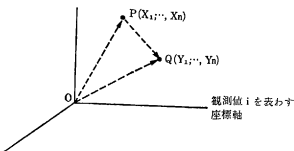
gion (7) Southern Region (14)

$$(4) \quad X_i = (\alpha + \Delta\alpha_i) Z_i, \quad i = 1, 2, \dots, n$$

$\Delta\alpha_i \sim Z_i$ である無相関であると仮定

$$\sum_{i=1}^n X_i = (\alpha + \bar{\Delta\alpha}) \sum_{i=1}^n Z_i$$

$$\text{したがって, } \bar{\Delta\alpha} = \frac{1}{n} \sum_{i=1}^n \Delta\alpha_i$$



引用文献

〔一〕 Jere R. Behrman, *Supply Response in Underdeveloped Agriculture—A Case study of four major annual*

線形回帰分析における観測誤差の影響と補正の手法について

crops in Thailand, 1957-1963, North-Holland, 1968, Chapter 7.

- [2] 村岡徳人「東南アジアの統計評価試験(Ⅳ)——フィリピンのみ——」(アジア経済研究所編『アジア研究』第一一巻第三号、昭和四五年四月)、九〇～九九頁。

[3] David H. Richardson and De-Min Wu, "Least Squares and Grouping Method Estimators in the Errors in Variables Model," *Journal of the American Statistical Association*, Vol. 65, No. 330, 1970, pp. 724~748.

[4] Takamitsu Sawa, "The Exact Sampling Distribution of Ordinary least Squares and Two-Stage least Squares Estimators," *Journal of the American Statistical Association*, Vol. 64, No. 327, 1969, pp. 923~937.

[5] Ko Takeuchi, "Exact Sampling Moments of The Ordinary least Squares, Instrumental Variable, and Two-Stage least Squares Estimators," *International Economic Review*, Vol. 11, No. 1, 1970, pp. 1~12.

[6] Henri Theil, *Economic Forecasts and Policy*, North-Holland, 1958, pp. 31~42.