

ダミー変数と数量化法への応用

青山 博次郎

(1965年4月受付)

Dummy Variable and Its Application to the Quantification Methods

Hirojiro AOYAMA

In the quantification method developed by Dr. Hayashi the classification and prediction problems are handled by making use of the effective transformation to the quantitative variable from the qualitative data.

In this paper the author treated this quantification method from the unified approach of dummy variables in the canonical correlation technique. This approach makes it easy not only to handle the given data by electronic computer and also to evaluate the sampling error of quantification results.

Institute of Statistical Mathematics

1. 緒 言

質的なものを一定の立場から数量化し、分類や予測に用いる方法が数量化法とよばれ、林知己夫氏によって発展させられている[1]。ここではダミー変数[2]を利用して、正準相関係数を用いて統一的に取扱う方法についてのべてみる。これによれば、電子計算機による処理の点で、また誤差評価の点で便利となる。

2. 外的基準が量的の場合

外的基準変数が数量であって、質的な要因を用いて量的予測を行なう場合を考えてみる。例えば性別、学歴別などで態度尺度値がどのようになるかを予測するといったような場合がこれに当る。

いまダミー変数 X_1, X_2, \dots, X_k を用いて、外的基準変数 Y を予測する式が

$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + u \quad (1)$$

で表わされるものとしよう。

ここで X_1, X_2, \dots, X_k は 1、または 0 をとる変数で、上例でいえば、 X_1 は男なら 1、女なら 0、 X_2 は女なら 1、男なら 0 なる値をとり、 X_3 は大学・高専卒なら 1、然らざるときは 0、 X_4 は高校卒なら 1、然らざるときは 0、 X_5 は小・中学校卒なら 1、然らざるときは 0 なる値をとるものと考えればよい。

また n 人のランダム・サンプルについての X_1, X_2, \dots, X_k の値を要素とする行列 \mathbf{X} を

$$\mathbf{X} = \begin{pmatrix} X_{11} & X_{21} & \cdots & X_{k1} \\ X_{12} & X_{22} & \cdots & X_{k2} \\ \vdots & & & \vdots \\ X_{1n} & X_{2n} & \cdots & X_{kn} \end{pmatrix} \quad (2)$$

また外的基準変数 \mathbf{Y} の値を要素とするベクトルを

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n) \quad (3)$$

とおく。

そうすると $\beta = (\beta_1, \beta_2, \dots, \beta_k)$ の推定値 $\hat{\beta}$ は

$$\hat{\beta} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \quad (4)$$

で表わせる。

この結果は、上例について示せば数量化法で得られる結果

$$\begin{pmatrix} f_{1.} & 0 & f_{11} & f_{12} & f_{13} \\ 0 & f_{2.} & f_{21} & f_{22} & f_{23} \\ f_{11} & f_{21} & f_{.1} & 0 & 0 \\ f_{12} & f_{22} & 0 & f_{.2} & 0 \\ f_{13} & f_{23} & 0 & 0 & f_{.3} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{pmatrix} = \begin{pmatrix} y_{1.} \\ y_{2.} \\ y_{.1} \\ y_{.2} \\ y_{.3} \end{pmatrix} \quad (5)$$

に等しい。ここで $f_{1.}$ は男の人数、 $f_{2.}$ は女の人数、 $f_{.1}$ は大学・高専卒の人数、 \dots 、 f_{11} は大学・高専卒の男子の人数、 f_{12} は高校卒の男子の人数、 f_{13} は小・中学校卒の男子の人数、 \dots 、を示し、 $y_{1.}$ は男子の態度尺度値の合計、 $y_{2.}$ は女子の態度尺度値の合計、 $y_{.1}$ は大学・高専卒のものの態度尺度値の合計、 \dots である。

このことは、(5) の左辺の行列が $n\mathbf{X}'\mathbf{X}$ に等しく、(5) の右辺のベクトルは $\mathbf{X}'\mathbf{Y}$ に等しいことから容易に分る。

また質的な変量 X_j については、測定誤差はないから、回帰式 (1) の誤差 \mathbf{u} のみを考えればよく、 \mathbf{u} の分散を σ^2 とおくと

$$D^2(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}, \quad \text{ただし} \quad \sigma^2 \mathbf{I} = E(\mathbf{u}\mathbf{u}')$$

が得られ、これを用いて数量化によって得られる $\hat{\beta}$ の誤差評価ができる事になる。

3. 外的基準変数が質的のものである場合

外的基準変数が質的のもので、予測の型が質的の要因から質的のもの（分類）を予測する場合を考えてみる。

先ず正準相関係数の理論によると、変数 $X_1, X_2, \dots, X_p, X_{p+1}, X_{p+2}, \dots, X_{p+q}$ ($p \leq q$) について

$$\left. \begin{array}{l} \zeta_i = \sum_{j=1}^p \alpha_{ij} X_j \quad (i=1, 2, \dots, p) \\ \eta_l = \sum_{m=1}^q \beta_{lm} X_{p+m} \quad (l=1, 2, \dots, q) \end{array} \right\} \quad (7)$$

なる一次変換を考え、 ζ_i, η_l の分散を 1 とし、共分散（相関係数） $R = \sum_{j,m} \alpha_{ij} \beta_{jm} \sigma_{jm}$ を最大にするような α_{ij}, β_{lm} を求めたとき、この相関係数 R が正準相関係数とよばれている。

このとき α_{ij}, β_{lm} の満足する関係式は

$$\left. \begin{array}{l} \sum_{m=1}^q \beta_{lm} \sigma_{jm} - \lambda \sum_{k=1}^p \alpha_{ik} \sigma_{jk} = 0, \quad i, j = 1, 2, \dots, p \\ \sum_{j=1}^p \alpha_{ij} \sigma_{jm} - \lambda \sum_{t=1}^q \beta_{lt} \sigma_{mt} = 0, \quad i = 1, 2, \dots, p \\ l, m = 1, 2, \dots, q \end{array} \right\} \quad (8)$$

となる [3]。

いま大きさ n のサンプルについて、 X_j ($j=1, 2, \dots, p$)、 X_{p+m} ($m=1, 2, \dots, q$) の値を X_{j1}, \dots, X_{jn} および $X_{p+m,1}, \dots, X_{p+m,n}$ 、それぞれの平均を \bar{X}_j および \bar{X}_{p+m} とし、

$$\mathbf{H}_0 = \begin{pmatrix} X_{11} & X_{21} & \cdots & X_{p1} \\ X_{12} & X_{22} & \cdots & X_{p2} \\ \vdots & \vdots & & \vdots \\ X_{1n} & X_{2n} & \cdots & X_{pn} \end{pmatrix} \quad (9)$$

$$\mathbf{Z}_0 = \begin{pmatrix} X_{p+1,1} & \cdots & X_{p+q,1} \\ X_{p+1,2} & \cdots & X_{p+q,2} \\ \vdots & & \vdots \\ X_{p+1,n} & \cdots & X_{p+q,n} \end{pmatrix} \quad (10)$$

$$\mathbf{H}_m = \begin{pmatrix} \bar{X}_1 & \bar{X}_2 & \cdots & \bar{X}_p \\ \vdots & \vdots & & \vdots \\ \bar{X}_1 & \bar{X}_2 & \cdots & \bar{X}_p \end{pmatrix} \quad (9)'$$

$$\mathbf{Z}_m = \begin{pmatrix} \bar{X}_{p+1} & \bar{X}_{p+2} & \cdots & \bar{X}_{p+q} \\ \vdots & \vdots & & \vdots \\ \bar{X}_{p+1} & \bar{X}_{p+2} & \cdots & \bar{X}_{p+q} \end{pmatrix} \quad (10)'$$

$$\mathbf{H} = \mathbf{H}_0 - \mathbf{H}_m, \quad \mathbf{Z} = \mathbf{Z}_0 - \mathbf{Z}_m \quad (11)$$

とおくと、(8) 式は

$$\left. \begin{array}{l} \mathbf{H}' \mathbf{Z} \beta - \lambda \mathbf{H}' \mathbf{H} \alpha = 0 \\ \mathbf{Z}' \mathbf{H} \alpha - \lambda \mathbf{Z}' \mathbf{Z} \beta = 0 \end{array} \right\} \quad (12)$$

と書くことができる。従って $\mathbf{H}' \mathbf{H}$, $\mathbf{Z}' \mathbf{Z}$ が非特異ならば (12) 式より

$$\left. \begin{array}{l} [\mathbf{Z}' \mathbf{H} (\mathbf{H}' \mathbf{H})^{-1} \mathbf{H}' \mathbf{Z}] \beta = \lambda^2 (\mathbf{Z}' \mathbf{Z}) \beta \\ [\mathbf{H}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{H}] \alpha = \lambda^2 (\mathbf{H}' \mathbf{H}) \alpha \end{array} \right\} \quad (13)$$

が得られる。勿論 (13) の第 1 式から β が分れば (12) 式より α を求めてよい。

さて数量化の問題に立ち戻って、外的基準が p 個のグループ別けであり、質的の要因が q 個のカテゴリー（すべての項目の各カテゴリーに一連番号を付して考える）より成立っているものとすると、 X_j を第 j 群に属すれば 1, 然らざれば 0 なる値をとるダミー変数と考え、 X_{p+m} を第 m カテゴリーに反応すれば 1, 然らざれば 0 なる値をとるダミー変数とする。また質的要因は項目数 K , それぞれのカテゴリーを m_1, m_2, \dots, m_K とし、 $q = \sum_{t=1}^K m_t$ とした。

このとき容易に分るように第 j 群の大きさを f_j とする

$$\mathbf{H}_0' \mathbf{H}_0 = \begin{pmatrix} f_{1.} & & & 0 \\ & f_{2.} & \ddots & \\ & & \ddots & \\ 0 & & & f_{p.} \end{pmatrix} \quad (14)$$

となるから $\mathbf{H}_0' \mathbf{H}_0$ は非特異行列で

$$(\mathbf{H}_0' \mathbf{H}_0)^{-1} = \begin{pmatrix} \frac{1}{f_{1.}} & & & 0 \\ & \frac{1}{f_{2.}} & \ddots & \\ & & \ddots & \\ 0 & & & \frac{1}{f_{p.}} \end{pmatrix} \quad (15)$$

また $\delta_i(p+m)$ をサンプル i がカテゴリー m に反応するとき 1, 然らざるとき 0 として

$$\mathbf{Z}_0 = \begin{pmatrix} \delta_1(p+1) & \cdots & \delta_1(p+q) \\ \vdots & & \vdots \\ \delta_n(p+1) & \cdots & \delta_n(p+q) \end{pmatrix} \quad (16)$$

とおけるから, $f_{p+m, p+m'}$ を第 m, m' カテゴリーに反応する数, f_{p+m} を第 m カテゴリーに反応する数とすると

$$\begin{aligned} \mathbf{Z}_0' \mathbf{Z}_0 = & \left(\begin{array}{c|cc|c|cc} f_{p+1.} & 0 & f_{p+1, p+m_1+1} & \cdots & f_{p+1, p+m_1+m_2} & \cdots & f_{p+1, p+q-m_K+1} & \cdots & f_{p+1, p+q} \\ \vdots & & \vdots & & \vdots & \cdots & \vdots & & \vdots \\ 0 & f_{p+m_1.} & f_{p+m_1, p+m_1+1} & \cdots & f_{p+m_1, p+m_1+m_2} & \cdots & f_{p+m_1, p+q-m_K+1} & \cdots & f_{p+m_1, p+q} \\ \hline \vdots & & f_{p+m_1+1.} & 0 & \cdots & \cdots & f_{p+m_1+1, p+q-m_K+1} & \cdots & f_{p+m_1+1, p+q} \\ \vdots & & 0 & f_{p+m_1+m_2.} & \cdots & \cdots & f_{p+m_1+m_2, p+q-m_K+1} & \cdots & f_{p+m_1+m_2, p+q} \\ \hline \vdots & & \cdots \\ \vdots & & \cdots & \cdots & \cdots & \cdots & f_{p+q-m_K+1.} & 0 & \cdots \\ \vdots & & \cdots & \cdots & \cdots & \cdots & 0 & f_{p+q.} & \cdots \end{array} \right) \\ & \equiv \mathbf{B} \quad (17) \end{aligned}$$

また $f_{j, p+m}$ は第 j 群でカテゴリー m に反応するものの数とすると

$$\mathbf{Z}_0' \mathbf{H}_0 = \begin{pmatrix} f_{1, p+1} & \cdots & f_{p, p+1} \\ \vdots & & \vdots \\ f_{1, p+q} & \cdots & f_{p, p+q} \end{pmatrix} \quad (18)$$

従って

$$\mathbf{Z}_0' \mathbf{H}_0 (\mathbf{H}_0' \mathbf{H}_0)^{-1} = \begin{pmatrix} \frac{f_{1, p+1}}{f_{1.}} & \cdots & \frac{f_{p, p+1}}{f_{p.}} \\ \vdots & & \vdots \\ \frac{f_{1, p+q}}{f_{1.}} & \cdots & \frac{f_{p, p+q}}{f_{p.}} \end{pmatrix} \quad (19)$$

$$\mathbf{Z}_0' \mathbf{H}_0 (\mathbf{H}_0' \mathbf{H}_0)^{-1} \mathbf{H}_0' \mathbf{Z}_0 = \begin{pmatrix} \sum_i \frac{f_{i, p+1}^2}{f_{i.}} & \sum_i \frac{f_{i, p+1} f_{i, p+2}}{f_{i.}} & \cdots & \sum_i \frac{f_{i, p+1} f_{i, p+q}}{f_{i.}} \\ \cdots & \cdots & \cdots & \cdots \\ \sum_i \frac{f_{i, p+q} f_{i, p+1}}{f_{i.}} & \sum_i \frac{f_{i, p+q} f_{i, p+2}}{f_{i.}} & \cdots & \sum_i \frac{f_{i, p+q}^2}{f_{i.}} \end{pmatrix} \equiv \mathbf{A} \quad (20)$$

を得られる。

ここで (12) を満足する α, β を求める代りに

$$\left. \begin{array}{l} \mathbf{H}_0' \mathbf{Z}_0 \beta - \lambda \mathbf{H}_0' \mathbf{H}_0 \alpha = \mathbf{0} \\ \mathbf{Z}_0' \mathbf{H}_0 \alpha - \lambda \mathbf{Z}_0' \mathbf{Z}_0 \beta = \mathbf{0} \end{array} \right\} \quad (12-1)$$

を満足する α, β を求めれば、同時に (12) を満足することが証明される。

それには

$$\mathbf{H}'\mathbf{H} = (\mathbf{H}_0' - \mathbf{H}_m')(\mathbf{H}_0 - \mathbf{H}_m) = \mathbf{H}_0'\mathbf{H}_0 - \mathbf{H}_m'\mathbf{H}_0 - \mathbf{H}_0'\mathbf{H}_m + \mathbf{H}_m'\mathbf{H}_m$$

より

$$\left. \begin{array}{l} \mathbf{H}_m'\mathbf{Z}_0\beta = \lambda \mathbf{H}_m'\mathbf{H}_0\alpha \\ \mathbf{H}_m'\mathbf{Z}_m\beta = \lambda \mathbf{H}_m'\mathbf{H}_m\alpha \end{array} \right\} \quad (12-2)$$

が満足されることを示せばよい。

(12-2) の第1式については

$$\mathbf{H}_m'\mathbf{Z}_0 = \begin{pmatrix} f_{p+1}, \bar{X}_1 & f_{p+2}, \bar{X}_1 & \cdots & f_{p+q}, \bar{X}_1 \\ f_{p+1}, \bar{X}_2 & f_{p+2}, \bar{X}_2 & \cdots & f_{p+q}, \bar{X}_2 \\ \cdots & \cdots & \cdots & \cdots \\ f_{p+1}, \bar{X}_p & f_{p+2}, \bar{X}_p & \cdots & f_{p+q}, \bar{X}_p \end{pmatrix} \quad (21)$$

$$\mathbf{H}_m'\mathbf{H}_0 = \begin{pmatrix} f_1, \bar{X}_1 & f_2, \bar{X}_1 & \cdots & f_p, \bar{X}_1 \\ f_1, \bar{X}_2 & f_2, \bar{X}_2 & \cdots & f_p, \bar{X}_2 \\ \cdots & \cdots & \cdots & \cdots \\ f_1, \bar{X}_p & f_2, \bar{X}_p & \cdots & f_p, \bar{X}_p \end{pmatrix} \quad (22)$$

であるから

$$(f_{p+1}, f_{p+2}, \dots, f_{p+q})\beta = \lambda(f_1, f_2, \dots, f_p)\alpha \quad (12-2')$$

となるが、これは (12-1) の第1式 (p 個あり) を加え合わせたものと同一である。

同様に

$$\mathbf{H}_m'\mathbf{Z}_m = \begin{pmatrix} \frac{f_1, f_{p+1}}{n} & \frac{f_1, f_{p+2}}{n} & \cdots & \frac{f_1, f_{p+q}}{n} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{f_p, f_{p+1}}{n} & \frac{f_p, f_{p+2}}{n} & \cdots & \frac{f_p, f_{p+q}}{n} \end{pmatrix} \quad (23)$$

$$\mathbf{H}_m'\mathbf{H}_m = \begin{pmatrix} \frac{f_1^2}{n} & \frac{f_1, f_2}{n} & \cdots & \frac{f_1, f_p}{n} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{f_p, f_1}{n} & \frac{f_p, f_2}{n} & \cdots & \frac{f_p^2}{n} \end{pmatrix} \quad (24)$$

であるから、(12-2) の第2式は上と同様に (12-2') 式と一致する。何れにせよ、(12-1) の第1式を満足する α, β は (12-2) 式を満足し、また (12-1) の第2式を満足する α, β も (12-2) 式を満足するから、結局 (12-1) を満足する α, β は (12) を満足する。かくて (12) を解く代りに (12-1)、従って (13) の代りに次の (13-1) 式を満足する α, β を求めればよいことになる。

$$\left. \begin{array}{l} [\mathbf{Z}_0'\mathbf{H}_0(\mathbf{H}_0'\mathbf{H}_0)^{-1}\mathbf{H}_0'\mathbf{Z}_0]\beta = \lambda^2(\mathbf{Z}_0'\mathbf{Z}_0)\beta \\ [\mathbf{H}_0'\mathbf{Z}_0(\mathbf{Z}_0'\mathbf{Z}_0)^{-1}\mathbf{Z}_0'\mathbf{H}_0]\alpha = \lambda^2(\mathbf{H}_0'\mathbf{H}_0)\alpha \end{array} \right\} \quad (13-1)$$

この第1式は上述の記号を用いると

$$\mathbf{A}\beta = \lambda^2 \mathbf{B}\beta \quad (25)$$

となり、これを解いて得られる β が数量化によって求められるカテゴリーに対する数値である。

以上の変形の中で、 \mathbf{Z} のみはそのままで用いるときは $\mathbf{Z} = \mathbf{Z}_0 - \mathbf{Z}_m$ を用いて

$$\mathbf{Z}'\mathbf{Z} = \mathbf{Z}_0'\mathbf{Z}_0 - \mathbf{Z}_m'\mathbf{Z}_m = \mathbf{B} - \begin{pmatrix} \frac{f_{p+1}}{n} & \frac{f_{p+1} \cdot f_{p+2}}{n} & \dots & \frac{f_{p+1} \cdot f_{p+q}}{n} \\ \dots & \dots & \dots & \dots \\ \frac{f_{p+1} \cdot f_{p+q}}{n} & \frac{f_{p+2} \cdot f_{p+q}}{n} & \dots & \frac{f_{p+q}}{n} \end{pmatrix} \equiv \mathbf{B}^* \quad (26)$$

$$\mathbf{Z}'\mathbf{H}_0 = \begin{pmatrix} f_{1,p+1} - \frac{f_1 \cdot f_{p+1}}{n} & f_{2,p+1} - \frac{f_2 \cdot f_{p+1}}{n} & \dots & f_{p,p+1} - \frac{f_p \cdot f_{p+1}}{n} \\ \dots & \dots & \dots & \dots \\ f_{1,p+q} - \frac{f_1 \cdot f_{p+q}}{n} & f_{2,p+q} - \frac{f_2 \cdot f_{p+q}}{n} & \dots & f_{p,p+q} - \frac{f_p \cdot f_{p+q}}{n} \end{pmatrix} \quad (27)$$

$$\begin{aligned} & \mathbf{Z}'\mathbf{H}_0(\mathbf{H}_0'\mathbf{H}_0)^{-1}\mathbf{H}_0'\mathbf{Z} \\ &= \begin{pmatrix} \sum_i \frac{f_{i,p+1}^2}{f_{i*}} - \frac{f_{p+1}^2}{n} & \sum_i \frac{f_{i,p+1} f_{i,p+2}}{f_{i*}} - \frac{f_{p+1} f_{p+2}}{n} & \dots & \sum_i \frac{f_{i,p+1} f_{i,p+q}}{f_{i*}} - \frac{f_{p+1} f_{p+q}}{n} \\ \dots & \dots & \dots & \dots \\ \sum_i \frac{f_{i,p+1} f_{i,p+q}}{f_{i*}} - \frac{f_{p+1} f_{p+q}}{n} & \dots & \sum_i \frac{f_{i,p+q}^2}{f_{i*}} - \frac{f_{p+q}^2}{n} & \end{pmatrix} \\ &\equiv \mathbf{A}^* \end{aligned} \quad (28)$$

となり、この場合も

$$\left. \begin{array}{l} \mathbf{H}_0'\mathbf{Z}\beta = \lambda \mathbf{H}_0'\mathbf{H}_0\alpha \\ \mathbf{Z}'\mathbf{H}_0\alpha = \lambda \mathbf{Z}'\mathbf{Z}\beta \end{array} \right\} \quad (12-3)$$

を満足する α, β を求めれば $\mathbf{Z}'\mathbf{H}_m = 0$ となり、

$$\mathbf{H}_m'\mathbf{Z}\beta = \lambda(\mathbf{H}_m'\mathbf{H}_0 + \mathbf{H}_0'\mathbf{H}_m - \mathbf{H}_m'\mathbf{H}_m)\alpha \quad (12-4)$$

を満足し、従ってまた (12) を満足することが証明できる。このときは (25) に対応する式は

$$\mathbf{A}^*\beta = \lambda \mathbf{B}^*\beta \quad (29)$$

となり、これは林氏の導出した式に当る [1]。

特にグループが 2 つの時は

$$\mathbf{A}^* = \mathbf{A} - \mathbf{Z}_m'\mathbf{Z}_m = \frac{f_1 \cdot f_2}{n} \mathbf{aa}'$$

ただし

$$\mathbf{a}' = \left(\frac{f_{1,p+1}}{f_1} - \frac{f_{2,p+1}}{f_2}, \frac{f_{1,p+2}}{f_1} - \frac{f_{2,p+2}}{f_2}, \dots, \frac{f_{1,p+q}}{f_1} - \frac{f_{2,p+q}}{f_2} \right)$$

となるから (29) の代りに

$$\mathbf{B}^*\beta = \mathbf{a} \quad (29)'$$

を解けばよいことになる。

更に全体の平均を 0 とおくときは $\mathbf{Z}_m\beta = 0$ となることが容易に証明されるから (29)' は

$$\mathbf{B}\beta = \mathbf{a} \quad (29)''$$

となる。

さてこのようにして得られた α, β の推定値を $\hat{\alpha}, \hat{\beta}$ とすると、この各要素の標本分散は

$$D^2(\hat{\alpha}_j) = \frac{\sigma_j^4}{2n}, \quad D^2(\hat{\beta}_m) = \frac{\sigma_m^4}{2n} \quad (30)$$

となる。何となれば [3] によると $\alpha_j, \beta_m, \sigma_{jk}$ の標本推定値を $\hat{\alpha}_j, \hat{\beta}_m, \hat{\sigma}_{jk}$ で表わすとき

$$2 \sum_{j,k} \hat{\sigma}_{jk} \hat{\alpha}_j \Delta \hat{\alpha}_k + \sum_{j,k} \hat{\alpha}_j \hat{\alpha}_k \Delta \hat{\sigma}_{jk} = 0$$

一般性を失なうことなく $\hat{\alpha}_1=1, \hat{\alpha}_2=\dots=\hat{\alpha}_p=0$ とおくと

$$2\Delta\hat{\alpha}_1 + \Delta\hat{\sigma}_{11} = 0$$

$$\therefore D^2(\hat{\alpha}_1) = E(\Delta\hat{\alpha}_1)^2 = \frac{1}{4} E(\Delta\hat{\sigma}_{11})^2 = \frac{\sigma_1^4}{2n}$$

同様にして

$$D^2(\hat{\beta}_1) = E(\Delta\hat{\beta}_1)^2 = \frac{1}{4} E(\Delta\hat{\sigma}_{p+1,p+1})^2 = \frac{\sigma_{p+1}^4}{2n}$$

が得られるからである。

数量化の問題においては、最大の正準相関係数に対する α_{ij}, β_{lm} の組を求めるときは、 X_j, X_{p+m} の分散の平方 σ_j^4, σ_m^4 の $1/2n$ 倍が標本分散となる。それ故 σ_j^4 は第 j 群の割合を t_j とすれば $t_j^2(1-t_j)^2$ となる。

例. 生徒の試験に合格、不合格を予測する 2 つの要因 I, II について、それぞれ 3, 2 のカテゴリに分かれ次表のような結果が得られたとする ([4] の例題より)。要因 I の各カテゴリに対応するダミー変数を X_1, X_2, X_3 、その係数を $\beta_1, \beta_2, \beta_3$ とし、要因 II のカテゴリに

要因 I		X_1		X_2		X_3		計
要因 II		X_4	X_5	X_4	X_5	X_4	X_5	
合 格		1	0	1	2	1	5	10
不 合 格		3	2	2	2	1	0	10
計		4	2	3	4	2	5	20

対応するダミー変数を X_1, X_2, X_3 、その係数を $\beta_1, \beta_2, \beta_3$ とし、上述の Z_0, H_0 を作ると次のようになる。

$$Z_0' = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdots X_1 \\ \cdots X_2 \\ \cdots X_3 \\ \cdots X_4 \\ \cdots X_5$$

$$H_0' = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdots \text{合格なら } 1 \\ \cdots \text{不合格なら } 1$$

従って

$$Z_0' Z_0 = \begin{pmatrix} 6 & 0 & 0 & 4 & 2 \\ 0 & 7 & 0 & 3 & 4 \\ 0 & 0 & 7 & 2 & 5 \\ 4 & 3 & 2 & 9 & 0 \\ 2 & 4 & 5 & 0 & 11 \end{pmatrix} \equiv B, \quad Z_0' H_0 = \begin{pmatrix} 1 & 5 \\ 3 & 4 \\ 6 & 1 \\ 3 & 6 \\ 7 & 4 \end{pmatrix}$$

$$H_0' H_0 = \begin{pmatrix} 10 & 0 \\ 0 & 10 \end{pmatrix}, \quad (H_0' H_0)^{-1} = \begin{pmatrix} 1/10 & 0 \\ 0 & 1/10 \end{pmatrix}$$

$$Z_0' H_0 (H_0' H_0)^{-1} H_0' Z_0 = \frac{1}{10} \begin{pmatrix} 26 & 23 & 11 & 33 & 27 \\ 23 & 25 & 22 & 33 & 37 \\ 11 & 22 & 37 & 24 & 46 \\ 33 & 33 & 24 & 45 & 45 \\ 27 & 37 & 46 & 45 & 45 \end{pmatrix} \equiv A$$

$$\mathbf{A}^* = \mathbf{A} - \mathbf{Z}_m' \mathbf{Z}_m = \frac{1}{20} \mathbf{aa}' , \quad \text{ただし} \quad \mathbf{a}' = (-4 \ -1 \ 5 \ -3 \ 3)$$

$$\mathbf{B}^* = \mathbf{B} - \mathbf{Z}_m' \mathbf{Z}_m = \mathbf{B} - \frac{1}{20} \begin{pmatrix} 36 & 42 & 42 & 54 & 66 \\ 42 & 49 & 49 & 63 & 77 \\ 42 & 49 & 49 & 63 & 77 \\ 54 & 63 & 63 & 81 & 99 \\ 66 & 77 & 77 & 99 & 121 \end{pmatrix} = \mathbf{B} - \frac{1}{20} \mathbf{bb}'$$

$$\text{ただし,} \quad \mathbf{b}' = (6 \ 7 \ 7 \ 9 \ 11)$$

数量化法で全体の平均を 0 とおくときは (29)'' より

$$\mathbf{B}\beta = \mathbf{a}$$

これを解いて

$$\begin{aligned} \hat{\beta}_1 &= -0.46809k, & \hat{\beta}_2 &= -0.01520k, & \hat{\beta}_3 &= 0.79939k, \\ \hat{\beta}_4 &= -0.29787k, & \hat{\beta}_5 &= 0 \quad (k \text{ は比例定数}) \end{aligned}$$

が得られる。

これはまた直接 $\mathbf{A}\beta = \lambda^2 \mathbf{B}\beta$ において、 \mathbf{B} の階数が 4 だから $\beta_5 = 0$ とおき、 \mathbf{A}, \mathbf{B} の第 5 行、第 5 列を除いた行列を $\mathbf{A}_0, \mathbf{B}_0$ とおいて、 $\mathbf{A}_0\beta = \lambda^2 \mathbf{B}_0\beta$ を解いてもよい。このときは \mathbf{B}_0^{-1} を両辺に乘じ、 $\mathbf{B}_0^{-1}\mathbf{A}_0$ の 1 でない最大の固有根 $\lambda^2 = 0.3389$ に対応する $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4$ を求めると、上述の結果が得られる。

またこのとき

$$\begin{aligned} \hat{\sigma}_1^2 &= \frac{6}{20} \left(1 - \frac{6}{20} \right) = 0.21, & \hat{\sigma}_2^2 &= \hat{\sigma}_3^2 = \frac{7}{20} \left(1 - \frac{7}{20} \right) = 0.2275, \\ \hat{\sigma}_4^2 &= \frac{9}{20} \left(1 - \frac{9}{20} \right) = 0.2475, & \hat{\sigma}_5^2 &= \frac{11}{20} \left(1 - \frac{11}{20} \right) = 0.2475 \end{aligned}$$

であるから

$$D(\hat{\beta}_1) = 0.03320, \quad D(\hat{\beta}_2) = D(\hat{\beta}_3) = 0.03597, \quad D(\hat{\beta}_4) = D(\hat{\beta}_5) = 0.03913$$

が β の推定値の標本誤差となる。

4. 外的基準のない場合

これは質的の要因を用いて適当なグループ分けを行ない、その予測が最も効率的に行われるようとする場合で、林氏のレッテルの選択におけるサンプルの分類がこれに当る [1]。

この場合は前節と同様にして q 個のレッテルの中特定の l_i 個をえらぶグループ i のサンプルの数を s_i 人とし、グループの順に整理して

$$\mathbf{H}_0' = \begin{pmatrix} \underbrace{1 \ 1 \ \cdots \ 1}_{s_1 l_1} & & & & 0 \\ & \underbrace{1 \ 1 \ \cdots \ 1}_{s_2 l_2} & & & \vdots \\ & & \ddots & & \\ 0 & & & \underbrace{1 \ 1 \ \cdots \ 1}_{s_p l_p} & \end{pmatrix} \quad (31)$$

のようになるものとする（表現の便宜上で、実際のデータでは、このようにする必要はない）。即ち異なる選択パターンを p 組とし、ダミー変数 X_1, \dots, X_p を考えるのである。対象としたサンプルの数は n であるが $(n = \sum_{i=1}^p s_i)$ 、重複してレッテルを選んだ延人数 $\sum_{i=1}^p s_i l_i = n \bar{l}$ を考え

たことになる。

また $\delta_i(p+m)$ は第 i 群が第 m 番目のレッテルをえらべば 1, 然らざれば 0 とし,

$$\mathbf{Z}_0' = \left(\begin{array}{c|c} \begin{matrix} \underbrace{\delta_1(p+1) \cdots \delta_1(p+1)}_{s_1} & 0 \\ \vdots & \vdots \\ 0 & \underbrace{\delta_1(p+q) \cdots \delta_1(p+q)}_{s_1} \end{matrix} & | \\ \hline \begin{matrix} \underbrace{\delta_2(p+1) \cdots \delta_2(p+1)}_{s_2} & 0 \\ \vdots & \vdots \\ 0 & \underbrace{\delta_2(p+q) \cdots \delta_2(p+q)}_{s_2} \end{matrix} & | \\ \hline \dots & \dots \end{array} \right) \quad (32)$$

ただし $\delta_i(p+m)=0$ となるときは、これを含む列は除くものとし、列の数は $n\bar{l}$ になっているものとする。

このとき

$$\mathbf{H}_0' \mathbf{H}_0 = \begin{pmatrix} s_1 l_1 & & & 0 \\ & s_2 l_2 & & \\ & & \ddots & \\ 0 & & & s_p l_p \end{pmatrix} \quad (33)$$

$$(\mathbf{H}_0' \mathbf{H}_0)^{-1} = \begin{pmatrix} \frac{1}{s_1 l_1} & & & 0 \\ & \frac{1}{s_2 l_2} & & \\ & & \ddots & \\ 0 & & & \frac{1}{s_p l_p} \end{pmatrix} \quad (34)$$

$$\mathbf{Z}_0' \mathbf{Z}_0 = \begin{pmatrix} \sum_i s_i \delta_i(p+1) & & & 0 \\ & \sum_i s_i \delta_i(p+2) & & \\ & & \ddots & \\ 0 & & & \sum_i s_i \delta_i(p+q) \end{pmatrix} \equiv \mathbf{B} \quad (35)$$

$$\mathbf{H}_0' \mathbf{Z}_0 = \begin{pmatrix} s_1 \delta_1(p+1) & s_1 \delta_1(p+2) & \cdots & s_1 \delta_1(p+q) \\ \dots & & & \\ s_p \delta_p(p+1) & s_p \delta_p(p+2) & \cdots & s_p \delta_p(p+q) \end{pmatrix} \quad (36)$$

$$\therefore \mathbf{Z}_0' \mathbf{H}_0 (\mathbf{H}_0' \mathbf{H}_0)^{-1} \mathbf{H}_0' \mathbf{Z}_0$$

$$= \begin{pmatrix} \sum_i \frac{s_i \delta_i^2(p+1)}{l_i} & \sum_i \frac{s_i \delta_i(p+1) \delta_i(p+2)}{l_i} & \cdots & \sum_i \frac{s_i \delta_i(p+1) \delta_i(p+q)}{l_i} \\ \dots & & & \\ \sum_i \frac{s_i \delta_i(p+1) \delta_i(p+q)}{l_i} & \sum_i \frac{s_i \delta_i(p+2) \delta_i(p+q)}{l_i} & \cdots & \sum_i \frac{s_i \delta_i^2(p+q)}{l_i} \end{pmatrix} \equiv \mathbf{A} \quad (37)$$

従って前節と同様にして

$$A\beta = \lambda^2 B\beta \quad (38)$$

を満足する β を求め,

$$\mathbf{Z}_0' \mathbf{H}_0 \boldsymbol{\alpha} = \lambda \mathbf{Z}_0' \mathbf{Z}_0 \boldsymbol{\beta} \quad (39)$$

より a を求めればよい。 a はグループに与えられる数量になる。

林氏の場合 [1] は前節と同様に Z_0 の代りに Z を用いることに当っている。このときも $Z = Z_0 - Z_m$ を用いて

$$\mathbf{Z}'\mathbf{Z} = \mathbf{Z}_0'\mathbf{Z}_0 - \mathbf{Z}_m'\mathbf{Z}_m = \mathbf{B} - \mathbf{Z}_m'\mathbf{Z}_m = (b_{jk}^*) \equiv \mathbf{B}^* \quad (40)$$

ただし

$$\left. \begin{aligned} b_{jk}^* &= -\sum_i s_i \delta_i(p+j) \delta_i(p+k) / (n\bar{l}) \quad (j \neq k) \\ b_{jj}^* &= \sum_i s_i \delta_i(p+j) - \left(\sum_i s_i \delta_i(p+j) \right)^2 / (n\bar{l}) \end{aligned} \right\} \quad (41)$$

$$\mathbf{Z}'\mathbf{H}_0 = \begin{pmatrix} s_1\delta_1(p+1) - s_1l_1\bar{X}_{p+1} & s_2\delta_2(p+1) - s_2l_2\bar{X}_{p+1} & \cdots & s_p\delta_p(p+1) - s_pl_p\bar{X}_{p+1} \\ \dots & \dots & \dots & \dots \\ s_1\delta_1(p+q) - s_1l_1\bar{X}_{p+q} & s_2\delta_2(p+1) - s_2l_2\bar{X}_{p+q} & \cdots & s_p\delta_p(p+q) - s_pl_p\bar{X}_{p+q} \end{pmatrix} \quad (42)$$

$$\mathbf{Z}' \mathbf{H}_0 (\mathbf{H}_0' \mathbf{H}_0)^{-1} \mathbf{H}_0' \mathbf{Z} = (a_{lm}) \equiv \mathbf{A}^* \quad (43)$$

ただし

$$a_{im} = \sum_i \frac{s_i \delta_i(p+t) \delta_i(p+m)}{l_i} - \frac{1}{nl} \left(\sum_i s_i \delta_i(p+t) \right) \left(\sum_i s_i \delta_i(p+m) \right) \quad (44)$$

が得られる.

このときは

$$A^* \beta = \lambda^2 B^* \beta \quad (45)$$

を満足する β を求め

$$Z' H_0 \alpha = \lambda Z' Z \beta \quad (46)$$

より α を求めることができる。結果は (38) から求めたものと同じである。

この場合の推定値 $\hat{\alpha}$, $\hat{\beta}$ の標本分散は前節と同様にして

$$D^2(\hat{\alpha}_j) = \sigma_j^4 / 2n\bar{l}, \quad D^2(\hat{\beta}_m) = \sigma_m^4 / 2n\bar{l}$$

となるが、重複してサンプルの数が数えられているので、 $n\bar{l}$ の代りに n を用いる方がよいだろう。

例. 10人の人が5種のレッテルをえらんだところ、3人のもの（第1群）は第1, 2, 4番のレッテルをえらび、4人のもの（第2群）は第1, 3番のレッテルをえらび、残りの3人（第3群）は第3, 4, 5番のレッテルをえらんだとする。このレッテルに与えられる $\beta = (\beta_1, \beta_2, \beta_3, \beta_4, \beta_5)$ および群に与えられる $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ を求めてみよう。

このとれ

$$\mathbf{Z}_0' \mathbf{Z}_0 = \begin{pmatrix} 7 & & & & 0 \\ & 3 & & & \\ & & 7 & & \\ 0 & & & 6 & \\ & & & & 3 \end{pmatrix} \equiv \mathbf{B}$$

$$\mathbf{H}_0' = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

$$\mathbf{H}_0' \mathbf{H}_0 = \begin{pmatrix} 9 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 9 \end{pmatrix}, \quad (\mathbf{H}_0' \mathbf{H}_0)^{-1} = \begin{pmatrix} \frac{1}{9} & 0 & 0 \\ 0 & \frac{1}{8} & 0 \\ 0 & 0 & \frac{1}{9} \end{pmatrix}$$

$$\mathbf{Z}_0' \mathbf{H}_0 (\mathbf{H}_0' \mathbf{H}_0)^{-1} \mathbf{H}_0' \mathbf{Z}_0 = \begin{pmatrix} 3 & 1 & 2 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 \\ 2 & 0 & 3 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 0 & 0 & 1 & 1 & 1 \end{pmatrix} \equiv \mathbf{A}$$

$$\mathbf{C} \equiv \mathbf{B}^{-1} \mathbf{A} = \begin{pmatrix} \frac{3}{7} & \frac{1}{7} & \frac{2}{7} & \frac{1}{7} & 0 \\ \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} & 0 \\ \frac{2}{7} & 0 & \frac{3}{7} & \frac{1}{7} & \frac{1}{7} \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & \frac{1}{6} \\ 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} = \frac{1}{42} \begin{pmatrix} 18 & 6 & 12 & 6 & 0 \\ 14 & 14 & 0 & 14 & 0 \\ 12 & 0 & 18 & 6 & 6 \\ 7 & 7 & 7 & 14 & 7 \\ 0 & 0 & 14 & 14 & 14 \end{pmatrix}$$

この \mathbf{C} の固有値は $1, \frac{10}{21}, \frac{8}{21}, 0, 0$ である。そこで $\mathbf{C}\beta = \lambda^2 \beta$ において、 $\lambda^2 = \frac{10}{21}$ とおき、

\mathbf{C} の第 4 行、第 4 列を除いた行列を \mathbf{C}^* とおいて ($\beta_4=0$ とおいた)

$$\mathbf{C}^*(\beta_1 \beta_2 \beta_3 \beta_5)' = \mathbf{0}$$

を解き

$$\hat{\beta}_1 = -\frac{3}{7}, \quad \hat{\beta}_2 = -1, \quad \hat{\beta}_3 = \frac{3}{7}, \quad \hat{\beta}_4 = 0, \quad \hat{\beta}_5 = 1$$

が得られる。

α については $\mathbf{Z}_0' \mathbf{H}_0 \alpha = \lambda \mathbf{Z}_0' \mathbf{Z}_0 \beta$ より上の $\hat{\beta}$ の値を用いて

$$\hat{\alpha}_1 = -\sqrt{\frac{10}{21}}, \quad \hat{\alpha}_2 = 0, \quad \hat{\alpha}_3 = \sqrt{\frac{10}{21}}$$

が得られる。

$\hat{\beta}$ についての標本誤差を求めるに、

$$\hat{\sigma}_1^2 = \hat{\sigma}_3^2 = \frac{7}{26} \left(1 - \frac{7}{26} \right) \doteq 0.197$$

$$\hat{\sigma}_2^2 = \hat{\sigma}_5^2 = \frac{3}{26} \left(1 - \frac{3}{26} \right) \doteq 0.102$$

$$\hat{\sigma}_4^2 = \frac{6}{26} \left(1 - \frac{6}{26} \right) \doteq 0.178$$

であるから、 $2n\bar{l}$ の代りに $2n=20$ を用いて

$$D(\hat{\beta}_1) = D(\hat{\beta}_3) \doteq 0.197/\sqrt{20} \doteq 0.0441$$

$$D(\hat{\beta}_2) = D(\hat{\beta}_5) \doteq 0.102/\sqrt{20} \doteq 0.0228$$

$$D(\hat{\beta}_4) \doteq 0.178/\sqrt{20} \doteq 0.0398$$

となる。

同様に

$$D(\hat{\alpha}_1) = D(\hat{\alpha}_3) \doteq 0.0505$$

$$D(\hat{\alpha}_2) = 0.0476$$

レッテルの選択において、 \mathbf{H}_0 , \mathbf{Z}_0 の作成は林氏の方法に従ったが、重複して数えないやり方で \mathbf{Z}_0 , \mathbf{H}_0 を作ることもできるし、また、交互作用の求められるようにダミー変数の数を増加して取扱かうこともできる。

統計数理研究所

参考文献

- [1] 林 知己夫: 数量化の方法と成分分析法・因子分析法, 昭和37年度工業統計講座
- [2] J. Johnston: Econometric Methods, McGraw Hill, 1960.
- [3] M. G. Kendall: The Advanced Theory of Statistics, Vol. II, 1948.
- [4] 青山博次郎: 教育統計学 (p. 199) 産業図書, 昭和32年