

MAT-466: Ecuaciones de estimación generalizadas

Felipe Osorio

fosorios.mat.utfsm.cl

Departamento de Matemática, UTFSM



Ecuaciones de estimación generalizadas

Considere $\mathbf{Y}_1, \dots, \mathbf{Y}_K$ vectores aleatorios independientes, con $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top$ y suponga

$$f(y; \theta_{ij}, \phi) = \exp[\phi\{y\theta_{ij} - b(\theta_{ij})\} + c(y; \phi)],$$

donde $E(Y_{ij}) = \mu_{ij} = b'(\theta_{ij})$, $\text{var}(Y_{ij}) = \phi^{-1}V_{ij}$, $V_{ij} = d\mu_{ij}/d\theta_{ij}$. Además,

$$g(\mu_{ij}) = \eta_{ij}, \quad \eta_{ij} = \mathbf{x}_{ij}^\top \boldsymbol{\beta}.$$

Bajo el supuesto de independencia, las funciones score son dadas por

$$U_K(\boldsymbol{\beta}) = \phi \sum_{i=1}^K \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}^\top} \right)^\top \mathbf{V}_i^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}_i) = \mathbf{0}, \quad (1)$$

estas ecuaciones pueden ser reescritas como:

$$U_K(\boldsymbol{\beta}) = \phi \sum_{i=1}^K \mathbf{X}_i^\top \boldsymbol{\Delta}_i (\mathbf{Y}_i - \boldsymbol{\mu}_i) = \mathbf{0},$$

donde $\boldsymbol{\Delta}_i = \partial \boldsymbol{\eta}_i / \partial \boldsymbol{\mu}_i^\top = \text{diag}(\partial \eta_{i1} / \partial \mu_{i1}, \dots, \partial \eta_{in_i} / \partial \mu_{in_i})$.



Sea $\hat{\beta}_{\text{Ind}}$ la solución de la ecuación (1), esto lleva al siguiente resultado.

Resultado 1:

El estimador $\hat{\beta}_{\text{Ind}}$ es consistente y $\sqrt{K}(\hat{\beta}_{\text{Ind}} - \beta)$ es asintóticamente normal con media cero y matriz de covarianza

$$\text{Cov}(\hat{\beta}_{\text{Ind}}) = \lim_{K \rightarrow \infty} K \{A(\beta)\}^{-1} B(\beta) \{A(\beta)\}^{-1},$$

donde

$$A(\beta) = \sum_{i=1}^K X_i^\top \Delta_i V_i \Delta_i X_i, \quad B(\beta) = \sum_{i=1}^K X_i^\top \Delta_i \text{Cov}(Y_i) \Delta_i X_i.$$



Observación:

Un estimador consistente para $\text{Cov}(\hat{\beta}_{\text{Ind}})$ es obtenido mediante substituir $\text{Cov}(Y_i)$ por $(Y_i - \mu_i(\hat{\beta}_{\text{Ind}}))(Y_i - \mu_i(\hat{\beta}_{\text{Ind}}))^{\top}$. Es decir,

$$\hat{\text{Cov}}(\hat{\beta}_{\text{Ind}}) = \{A(\beta)\}^{-1} \left(\sum_{i=1}^K X_i^{\top} \hat{\Delta}_i \hat{r}_i \hat{r}_i^{\top} \hat{\Delta}_i X_i \right) \{A(\beta)\}^{-1},$$

donde $r_i = Y_i - \mu_i(\beta)$.



Sea $R_i(\alpha)$ matriz simétrica tal que corresponde a una matriz de correlación y considere α un vector que caracteriza $R_i(\alpha)$ que es llamada **matriz de correlación de trabajo**.

Sea

$$\Sigma_i(\theta) = \phi V_i^{1/2} R_i(\alpha) V_i^{1/2}.$$

Esto llevó a Liang y Zeger (1986)¹ a definir:

$$\Psi_K(\beta) = \phi \sum_{i=1}^K \left(\frac{\partial \mu_i}{\partial \beta^\top} \right)^\top \{ V_i^{1/2} R_i(\alpha) V_i^{1/2} \}^{-1} (Y_i - \mu_i) = 0. \quad (2)$$

Usando que $D_i = \partial \mu_i / \partial \beta^\top$ puede ser escrito como

$$D_i = V_i \Delta_i X_i, \quad \text{o bien} \quad D_i = W_i^{1/2} V_i^{1/2} X_i,$$

lleva a diversas formas de definir $\Psi_K(\beta)$.

¹Biometrika 73, 13-22

Ecuaciones de estimación generalizadas

El algoritmo de estimación, adopta la forma:

$$\beta^{(r+1)} = \beta^{(r)} + \left(\sum_{i=1}^K D_i^{(r)\top} \Sigma_i^{-(r)} D_i^{(r)} \right)^{-1} \sum_{i=1}^K D_i^{(r)\top} \Sigma_i^{-(r)} (Y_i - \mu_i^{(r)}),$$

que puede ser escrito como un problema de IGLS, como:

$$\beta^{(r+1)} = \left(\sum_{i=1}^K D_i^{(r)\top} \Sigma_i^{-(r)} D_i^{(r)} \right)^{-1} \sum_{i=1}^K D_i^{(r)\top} \Sigma_i^{-(r)} Z_i^*,$$

con $Z_i^* = D_i \beta + Y_i - \mu_i$.

Además, podemos escribir

$$\beta^{(r+1)} = \left(\sum_{i=1}^K X_i^\top W_i^{1/2} R_i^{-1}(\alpha) W_i^{1/2} X_i \right)^{-1} \sum_{i=1}^K X_i^\top W_i^{1/2} R_i^{-1}(\alpha) W_i^{1/2} Z_i,$$

con $Z_i = \eta_i + W_i^{-1/2} V_i^{-1/2} (Y_i - \mu_i)$.



Observación:

- ▶ Note que $\hat{\beta}_{\text{GEE}}$ como solución de $\Psi_K(\beta) = \mathbf{0}$ no depende de ϕ .
- ▶ Podemos hacer que $\Psi_K(\beta)$ depende solamente de β mediante substituir α por un estimador \sqrt{K} -consistente, es decir

$$\sqrt{K}(\hat{\alpha} - \alpha) = O_p(1).$$

De ahí que

$$\Psi_K^*(\beta) = \sum_{i=1}^K \left(\frac{\partial \mu_i}{\partial \beta^\top} \right)^\top \{V_i^{1/2} R_i(\hat{\alpha}) V_i^{1/2}\}^{-1} (Y_i - \mu_i) = \mathbf{0}.$$

Recuerde que:

Para $\{X_n\}$, $\{Y_n\}$ secuencias de variables aleatorias entonces $X_n = O_p(Y_n)$ si existe $\epsilon > 0$, $M = M(\epsilon)$ y $n_0 = n_0(\epsilon)$, tal que

$$P(|X_n| \leq M|Y_n|) \geq 1 - \epsilon, \quad \text{para todo } n > n_0.$$



Resultado 2:

Suponga que

- (i) $\hat{\alpha}$ es $K^{1/2}$ -consistente dado β y ϕ .
- (ii) $\hat{\phi}$ es $K^{1/2}$ -consistente dado β .
- (iii) $|\partial \hat{\alpha}(\beta, \phi)| \leq K(\mathbf{Y}, \beta)$ que es $O_p(1)$.

Entonces $\sqrt{K}(\hat{\beta}_{\text{GEE}} - \beta)$ es asintóticamente normal con media cero y matriz de covarianza

$$\text{Cov}(\hat{\beta}_{\text{GEE}}) = \lim_{K \rightarrow \infty} K \{A(\beta)\}^{-1} B(\beta) \{A(\beta)\}^{-1},$$

donde

$$A(\beta) = \sum_{i=1}^K D_i^\top \Sigma_i^{-1} D_i, \quad B(\beta) = \sum_{i=1}^K D_i^\top \Sigma_i^{-1} \text{Cov}(\mathbf{Y}_i) \Sigma_i^{-1} D_i.$$



Ecuaciones de estimación generalizadas

Sea

$$r_{ij} = (Y_{ij} - \hat{\mu}_{ij}) / \{V(\hat{\mu}_{ij})\}^{1/2}, \quad \hat{\mu}_{ij} = \mu_{ij}(\hat{\beta}),$$

para $i = 1, \dots, K$; $j = 1, \dots, n_i$.

Entonces, podemos considerar

$$\hat{\phi}^{-1} = \frac{1}{N - p} \sum_{i=1}^K \sum_{j=1}^{n_i} r_{ij}^2,$$

con $N = \sum_{i=1}^K n_i$. Este estimador es $K^{1/2}$ -consistente siempre que los cuartos momentos de Y_{ij} sean finitos.



Suponga, sin pérdida de generalidad, que $n_i = n$. Considere los siguientes ejemplos,

Ejemplo:

Considere $\mathbf{R}_i(\boldsymbol{\alpha}) = \mathbf{I}$, entonces obtenemos las ecuaciones de estimación de independencia dadas en (1).

Ejemplo:

Sea $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top$ donde $\alpha_j = \text{corr}(Y_{ij}, Y_{i,j+1})$, $j = 1, \dots, n-1$. Un estimador natural para α_j es

$$\hat{\alpha}_j = \frac{\phi}{K-p} \sum_{i=1}^K r_{ij} r_{i,j+1}.$$



Ejemplo:

Suponga $\text{corr}(Y_{ij}, Y_{ij'}) = \alpha$ para todo $j \neq j'$. Esto corresponde a la estructura de equicorrelación. De este modo,

$$\hat{\alpha} = \frac{\phi}{\sum_{i=1}^K n_i(n_i - 1)/2 - p} \sum_{i=1}^K \sum_{j > j'} r_{ij} r_{ij'}$$

Ejemplo:

Sea $\text{corr}(Y_{ij}, Y_{ij'}) = \alpha^{|j-j'|}$, que es llamada estructura de correlación AR-1. Como $E(r_{ij} r_{ij'}) \approx \alpha^{|j-j'|}$ entonces α puede ser estimado desde la regresión de $\log(r_{ij} r_{ij'})$ sobre $\log |j - j'|$.



Algunos paquetes en R para estimación de GEE:

gee: Generalized Estimation Equation Solver.
(Carey, V.J., portado a R por Lumley, T. y Ripley, B).

geepack: Solve Generalized Estimating Equations.
McDaniel, L.S., Henderson, N.C., Rathouz, P.J. (2013).
Fast pure R implementation of GEE: application of the Matrix package.
The R Journal 5/1, 181-187.

geeM: Generalized Estimating Equation Package.
Halekoh, U., Højsgaard, S., Yan, J. (2006).
The R package geepack for generalized estimating equations.
Journal of Statistical Software 15/2, 1-11.



Datos de incontinencia urinaria (Preisser y Qaqish, 1999)

- ▶ Datos provienen de un estudio para evaluar el impacto de la **incontinencia urinaria** sobre la vida de pacientes ancianos (**GUIDE**) mayores de 76 años.
- ▶ La **respuesta es binaria**, indicando si el individuo siente que su **rutina diaria** se ve **afectada** por pérdidas accidentales de orina.
- ▶ Datos obtenidos para **137 ancianos** agrupados en **38 prácticas médicas** (cluster), los datos son **desbalanceados** (de 1 a 8 pacientes por cluster).
- ▶ Se dispone de 5 regresores, **sexo**, **edad**, **accidentes diarios**, **severidad** y número de veces que usa el **baño** diariamente.
- ▶ Se consideró un **enlace logístico** para el siguiente modelo

$$\text{logit}(\mu_{ij}) = \beta_0 + \beta_1 \text{Sexo} + \beta_2 \text{Edad} + \beta_3 \text{Accidentes} + \beta_4 \text{Severo} + \beta_5 \text{Baño}.$$

Además se asumió una estructura de **equicorrelación** para $R_i(\alpha)$.



Datos de incontinencia urinaria (Preisser y Qaqish, 1999)

Resultados del ajuste mediante GEE para los datos GUIDE.

Variable	Estimación	Err.Est.	Z
Intercepto	-3.054	0.959	-3.185
Sexo	-0.745	0.600	-1.242
Edad	-0.676	0.561	-1.205
Accidentes	0.392	0.093	4.202
Severo	0.812	0.359	2.263
Baño	0.108	0.099	1.090

Se obtuvo además $\hat{\alpha} = 0.0932$.

