

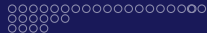
Conteúdo

- 1 Introdução
 - Modelo de Regressão Clássico
 - Modelo de Regressão Generalizado
- 2 Exemplos
- 3 Método *GEE*
 - Estruturas de Correlação
 - Métodos de Diagnóstico
- 4 Dados Ausentes
- 5 Aplicação
 - Regressão Logística
 - Aplicação no R
- 6 Referências



Introdução

- Modelagem de uma variável resposta pertencente a um grupo (*cluster*)



Introdução

- Modelagem de uma variável resposta pertencente a um grupo (*cluster*)
- Leva em consideração a correlação inerente ao grupo



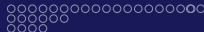
Introdução

- Modelagem de uma variável resposta pertencente a um grupo (*cluster*)
- Leva em consideração a correlação inerente ao grupo
- Modelos **Multivariados** vs Modelos **Marginais**



Introdução

- Modelagem de uma variável resposta pertencente a um grupo (*cluster*)
- Leva em consideração a correlação inerente ao grupo
- Modelos **Multivariados** vs Modelos **Marginais**
- *Generalized Estimating Equations*



Introdução

- Modelagem de uma variável resposta pertencente a um grupo (*cluster*)
- Leva em consideração a correlação inerente ao grupo
- Modelos **Multivariados** vs Modelos **Marginais**
- *Generalized Estimating Equations*

Exemplos: estudos longitudinais, estudos clínicos, pesquisas de satisfação, dentre outros.

Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

- \mathbf{X} é a uma matriz de contantes e de posto completo;

Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

- \mathbf{X} é a uma matriz de contantes e de posto completo;
- $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2)$;

Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

- \mathbf{X} é a uma matriz de constantes e de posto completo;
- $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2)$;
- Homocedasticidade;

Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

- \mathbf{X} é a uma matriz de contantes e de posto completo;
- $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2)$;
- Homocedasticidade;
- $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$

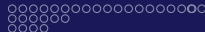
Modelo de Regressão Clássico

O modelo de regressão linear múltiplo com p variáveis explicativas é dado por

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

Suposições:

- \mathbf{X} é a uma matriz de contantes e de posto completo;
- $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2)$;
- Homocedasticidade;
- $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$
- $\text{Cov}(\epsilon_i, \epsilon_j) = 0$



Modelos Lineares Generalizados

No modelo de regressão generalizado temos um **componente aleatório** representado por um conjunto de variáveis aleatórias **independentes** Y_1, \dots, Y_n provenientes de uma mesma distribuição que faz parte de da **família exponencial** com médias μ_1, \dots, μ_n , ou seja,

$$E(Y_i) = \mu_i, i = 1, \dots, n,$$



Modelos Lineares Generalizados

No modelo de regressão generalizado temos um **componente aleatório** representado por um conjunto de variáveis aleatórias **independentes** Y_1, \dots, Y_n provenientes de uma mesma distribuição que faz parte de da **família exponencial** com médias μ_1, \dots, μ_n , ou seja,

$$E(Y_i) = \mu_i, i = 1, \dots, n,$$

um **componente sistemático** dado por

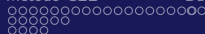
$$\eta_i = \sum_{j=1}^p x_{ij} \beta_j = \mathbf{x}_i^T \boldsymbol{\beta} \quad \text{ou} \quad \boldsymbol{\eta} = \mathbf{X} \boldsymbol{\beta},$$

Modelos Lineares Generalizados

e uma **função de ligação** que liga o componente aleatório ao componente sistemático, ou seja, relaciona a média ao preditor linear

$$\eta_i = g(\mu_i)$$

Problemática: como modelar dados correlacionados que violam a suposição de independência?



Exemplos

- 1 Em um estudo de pesquisa ocular de Baltimore, mais de 5000 pessoas com 40 anos ou mais receberam um exame visual como parte de um estudo de prevalência baseado na população de distúrbios oculares [1].
 - **Objetivo:** identificar variáveis demográficas que estão associadas à perda de visão. Os dados estão disponíveis para ambos os olhos de todos os sujeitos. Um único modelo de regressão expressando a deficiência visual em termos de variáveis demográficas aborda os objetivos científicos. No entanto, os dois olhos da mesma pessoa provavelmente não são independentes, porque muitas causas de deficiência são binoculares. Essa associação deve ser considerada.

Método *GEE*

Quais os passos para se utilizar o método?

- 1 Definir um modelo marginal para a média da resposta, $(E[Y_i])$, e escolher $g(\mu)$

Método GEE

Quais os passos para se utilizar o método?

- 1 Definir um modelo marginal para a média da resposta, $(E[Y_i])$, e escolher $g(\mu)$
- 2 Supor como $Var[Y_i]$ depende de $E[Y_i]$

Método GEE

Quais os passos para se utilizar o método?

- 1 Definir um modelo marginal para a média da resposta, $(E[Y_i])$, e escolher $g(\mu)$
- 2 Supor como $Var[Y_i]$ depende de $E[Y_i]$
- 3 Propor um comportamento da correlação no grupo

Método *GEE*

Quais os passos para se utilizar o método?

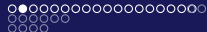
- 1 Definir um modelo marginal para a média da resposta, $(E[Y_i])$, e escolher $g(\mu)$
- 2 Supor como $Var[Y_i]$ depende de $E[Y_i]$
- 3 Propor um comportamento da correlação no grupo
- 4 Estimar através do modelo marginal e da estrutura da correlação supostos

Método *GEE*

Quais os passos para se utilizar o método?

- 1 Definir um modelo marginal para a média da resposta, $(E[Y_i])$, e escolher $g(\mu)$
- 2 Supor como $Var[Y_i]$ depende de $E[Y_i]$
- 3 Propor um comportamento da correlação no grupo
- 4 Estimar através do modelo marginal e da estrutura da correlação supostos
- 5 Ajustar os erros padrão e obter estimativas mais robustas

- Com uma única observação para cada sujeito ($r_i = 1$), um modelo linear generalizado (McCullagh & Nelder, 1983) pode ser aplicado.



- Com uma única observação para cada sujeito ($r_i = 1$), um modelo linear generalizado (McCullagh & Nelder, 1983) pode ser aplicado.
- Com observações repetidas, a **correlação** para um determinado sujeito deve ser levada em consideração.

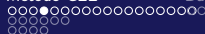
- Com uma única observação para cada sujeito ($r_i = 1$), um modelo linear generalizado (McCullagh & Nelder, 1983) pode ser aplicado.
- Com observações repetidas, a **correlação** para um determinado sujeito deve ser levada em consideração.
- Uma dificuldade na análise de dados longitudinais **não-Gaussianos** é a falta de uma classe rica de modelos, como o Normal Multivariado, para a distribuição conjunta de y_{it} , $t = 1, 2, \dots, r_i$. Assim, métodos de verossimilhança não estão disponíveis, exceto nos poucos casos mencionados acima.

- Com uma única observação para cada sujeito ($r_i = 1$), um modelo linear generalizado (McCullagh & Nelder, 1983) pode ser aplicado.
- Com observações repetidas, a **correlação** para um determinado sujeito deve ser levada em consideração.
- Uma dificuldade na análise de dados longitudinais **não-Gaussianos** é a falta de uma classe rica de modelos, como o Normal Multivariado, para a distribuição conjunta de y_{it} , $t = 1, 2, \dots, r_i$. Assim, métodos de verossimilhança não estão disponíveis, exceto nos poucos casos mencionados acima.
- Para contornar o problema é introduzindo uma estrutura de correlação na função score, produzindo um novo sistema de equações para estimar β .



Suposições para a estimação dos parâmetros de um modelo de regressão

Modelo Linear Clássico	Modelo Linear Generalizado	GEE
$i.i.d.$ $y \sim \mathcal{N}(\mu; \sigma^2)$ $\mu = \mathbf{x}^\top \boldsymbol{\beta}$	$i.i.d.$ $y \sim Fam.Expo.$ $g(\mu) = \mathbf{x}^\top \boldsymbol{\beta}$	não $y \sim V(\mu), \phi, R(\rho)$ $g(\mu) = \mathbf{x}^\top \boldsymbol{\beta}$



Definição

Seja $\mathbf{Y}_i = (y_{i1}, \dots, y_{ir_i})^\top$ o vetor resposta multivariado para a i -ésima unidade experimental, $i = 1, \dots, n$, e assumindo que apenas a distribuição marginal de Y_{it} é conhecida e expressa da seguinte maneira:

$$f(y; \theta_{it}, \phi) = \exp\{\phi[y_{it}\theta_{it} - b(\theta_{it})] + c(y_{it}; \phi)\}, \quad (2)$$

em que

$$\mathbb{E}(Y_{it}) = b'(\theta_{it}) = \mu_{it},$$

$$\mathbb{V}ar(Y_{it}) = \phi^{-1} V(\mu_{it}),$$

$$V(\mu_{it}) = b''(\theta_{it}) = \frac{d\mu_{it}}{d\theta_{it}},$$

$$\phi > 0.$$

Podemos então definir um modelo linear generalizado para cada instante t acrescentando uma função de ligação

$$g(\mu_{it}) = \eta_{it}, \quad (3)$$

em que $\eta_{it} = \mathbf{x}_{it}^\top \boldsymbol{\beta}$ é o preditor linear, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ é um vetor de parâmetros desconhecidos a serem estimados, e $\mathbf{x}_{it} = (x_{it1}, \dots, x_{itp})^\top$ representa a matriz de delineamento para a i -ésima unidade experimental no tempo.

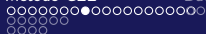
Função Escore e Matriz de Informação

A **função escore** e a **matriz de informação** para β , ignorando a estrutura de correlação intragrupo, ficam respectivamente como

$$U_{\beta} = \phi \sum_{i=1}^n D_i^{\top} V_i^{-1} (y_i - \mu_i), \quad (4)$$

$$K_{\beta\beta} = \phi \sum_{i=1}^n D_i^{\top} V_i D_i, \quad (5)$$

onde $D_i = W_i^{1/2} V_i^{1/2} X_i$.



Substituindo D_i na expressão da **função escore**, temos que

$$\begin{aligned}
 U_{\beta} &= \phi \sum_{i=1}^n (\mathbf{w}_i^{1/2} \mathbf{v}_i^{1/2} \mathbf{x}_i)^{\top} \mathbf{v}_i^{-1} (\mathbf{y}_i - \mu_i) \\
 &= \phi \sum_{i=1}^n \mathbf{x}_i^{\top} (\mathbf{v}_i^{1/2})^{\top} (\mathbf{w}_i^{1/2})^{\top} \mathbf{v}_i^{-1} (\mathbf{y}_i - \mu_i) \\
 &= \phi \sum_{i=1}^n \mathbf{x}_i^{\top} \mathbf{v}_i^{1/2} (\mathbf{w}_i^{1/2})^{\top} \mathbf{v}_i^{-1} (\mathbf{y}_i - \mu_i) \\
 &= \phi \sum_{i=1}^n \mathbf{x}_i^{\top} \mathbf{w}_i^{1/2} \mathbf{v}_i^{-1/2} (\mathbf{y}_i - \mu_i).
 \end{aligned}$$

O processo é análogo para a **matriz de informação**.



- $D_i = W_i^{1/2} V_i^{1/2} X_i$
- X_i é uma matriz $r_i \times p$ de linhas \mathbf{x}_{it}^\top
- $W_i = \text{diag}\{w_{i1}, \dots, w_{ir_i}\}$, $w_{it} = (d\mu_{it}/d\eta_{it})^2 / V_{it}$
- $V_i = \text{diag}\{V_{i1}, \dots, V_{ir_i}\}$
- $\mathbf{y}_i = (y_{i1}, \dots, y_{ir_i})^\top$
- $\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{ir_i})^\top$

Quando há **ligação canônica** a função escore e a matriz de informação de Fisher ficam dadas, respectivamente, por:

$$\mathbf{U}_\beta = \phi \sum_{i=1}^n \mathbf{x}_i^\top (\mathbf{y}_i - \boldsymbol{\mu}_i) \quad (6)$$

$$\mathbf{K}_{\beta\beta} = \phi \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{v}_i \mathbf{x}_i \quad (7)$$



Exemplo: Vamos supor inicialmente que os dados não são correlacionados e que a matriz de correlação correspondente ao i -ésimo grupo é denotado por \mathbf{R}_i . Portanto, teremos $\mathbf{R}_i = \mathbf{I}_{r_i}$.



Exemplo: Vamos supor inicialmente que os dados não são correlacionados e que a matriz de correlação correspondente ao i -ésimo grupo é denotado por \mathbf{R}_i . Portanto, teremos $\mathbf{R}_i = \mathbf{I}_{r_i}$. A matriz de variância-covariância para Y_i , por definição é dada por

$$\text{Var}(Y_i) = \phi^{-1} \mathbf{V}_i^{1/2} \mathbf{R}_i \mathbf{V}_i^{1/2}.$$

Equações de Estimação Generalizadas (EEGs)

Em 1986 Kung-Yee Liang e Scott L. Zeger publicaram o artigo *Longitudinal data analysis using generalized linear models*, onde propuseram uma matriz de correlação dada por $R_i(\rho)$, em que $\rho = (\rho_1, \dots, \rho_p)^\top$ é um **vetor de parâmetros de perturbação** que não depende de β .

GEE

Matriz de variância-covariância

Se a verdadeira correlação entre os elementos de \mathbf{Y}_i for dada por $\mathbf{R}_i(\boldsymbol{\rho})$, temos que a matriz de variância-covariância pode ser escrita como

$$\boldsymbol{\Omega}_i = \phi^{-1} \mathbf{V}_i^{1/2} \mathbf{R}_i(\boldsymbol{\rho}) \mathbf{V}_i^{1/2}. \quad (8)$$

$\mathbf{R}_i(\boldsymbol{\rho})$ é uma matriz $r_i \times r_i$ que depende do número finito de parâmetros $\boldsymbol{\rho} = (\rho_1, \dots, \rho_q)^\top$, sendo denominada matriz trabalho.

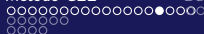
GEE

Para estimarmos β devemos resolver o seguinte sistema de equações:

$$\mathbf{S}_{\beta}(\hat{\beta}_G) = 0,$$

$$\mathbf{S}_{\beta} = \sum_{i=1}^n \mathbf{D}_i^{\top} \boldsymbol{\Omega}_i(\mathbf{y}_i - \boldsymbol{\mu}_i). \quad (9)$$

$\mathbf{S}_{\beta}(\beta)$ também depende de ϕ e $\boldsymbol{\rho} = (\rho_1, \dots, \rho_q)^{\top}$ que são estimados separadamente de β .



Estimação

O processo iterativo para a estimação de β , que é uma modificação do método score de Fisher, é dado por:

$$\beta_G^{(m+1)} = \beta_G^{(m)} + \left\{ \sum_{i=1}^n \mathbf{D}_i^{(m)\top} \boldsymbol{\Omega}_i^{(m)} \mathbf{D}_i^{(m)} \right\}^{-1} \times$$

$$\left[\sum_{i=1}^n \mathbf{D}_i^{(m)\top} \boldsymbol{\Omega}_i^{-(m)} (\mathbf{y}_i - \boldsymbol{\mu}_i^{(m)}) \right]$$

Definição

Supondo que $\hat{\rho}$ e $\hat{\phi}$ são estimadores consistentes de ρ e ϕ , respectivamente, temos que

$$\sqrt{n}(\hat{\beta}_G - \beta) \xrightarrow{d} \mathcal{N}_p(\mathbf{0}, \Sigma) \quad (10)$$

em que

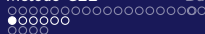
$$\Sigma = \lim_{n \rightarrow \infty} \left[n \left(\sum_{i=1}^n D_i^\top \Omega_i^{-1} D_i \right)^{-1} \times \left\{ \sum_{i=1}^n D_i^\top \Omega_i^{-1} \text{Var}(Y_i) \Omega_i^{-1} D_i \right\} \times \left(\sum_{i=1}^n D_i^\top \Omega_i^{-1} D_i \right)^{-1} \right]$$

Definição

Se a matriz de correlação $\mathbf{R}_i(\boldsymbol{\rho})$ é definida corretamente, um estimador consistente para $\mathbf{Var}(\hat{\boldsymbol{\beta}}_G)$ é dado por $\mathbf{H}^{-1}(\hat{\boldsymbol{\beta}}_G)$, em que

$$\mathbf{H}(\hat{\boldsymbol{\beta}}_G) = \sum_{i=1}^n (\hat{\mathbf{D}}_i^\top \hat{\boldsymbol{\Omega}}_i^{-1} \hat{\mathbf{D}}_i), \quad (11)$$

com $\hat{\mathbf{D}}_i$ sendo avaliado em $\hat{\boldsymbol{\beta}}_G$ e $\hat{\boldsymbol{\Omega}}_i$ avaliado em $(\hat{\phi}, \hat{\boldsymbol{\rho}}, \hat{\boldsymbol{\beta}}_G)$



Matriz de correlação não estruturada

Quando a matriz de correlação \mathbf{R}_i é não estruturada teremos $r_i(r_i - 1)/2$ parâmetros a serem estimados. Denotando $\mathbf{R}_i = \{R_{ijj'}\}$, o j, j' -ésimo elemento de \mathbf{R}_i poderá ser estimado por

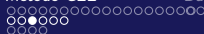
$$\hat{R}_{jj'} = \frac{\phi}{n} \sum_{i=1}^n \frac{(y_{ij} - \hat{\mu}_{ij})}{\sqrt{\hat{V}_{ij}}} \frac{(y_{ij'} - \hat{\mu}_{ij'})}{\sqrt{\hat{V}_{ij'}}}. \quad (13)$$



Matriz de correlação simétrica ou permutável

Nesse caso, assumimos $\mathbf{R}_i = \mathbf{R}_i(\rho)$, em que o (j, j') -ésimo elemento de \mathbf{R}_i é dado por $R_{ijj'} = 1$, para $j = j'$, e $R_{ijj'} = \rho$, para $j \neq j'$. Um estimador consistente para ρ é dado por:

$$\hat{\rho} = \frac{\phi}{n} \sum_{i=1}^n \frac{1}{r_i(r_i - 1)} \sum_{j=1}^{r_i} \sum_{j'=1, j' \neq j}^{r_i} \frac{(y_{ij} - \hat{\mu}_{ij})}{\sqrt{\hat{V}_{ij}}} \frac{(y_{ij'} - \hat{\mu}_{ij'})}{\sqrt{\hat{V}_{ij'}}}. \quad (14)$$



Autoregressiva

Matriz de correlação autorregressiva

Assumimos $\mathbf{R}_i = \mathbf{R}_i(\rho)$, em que o (j, j') -ésimo elemento de \mathbf{R}_i é dado por $R_{ijj'} = 1$, para $j = j'$, e $R_{ijj'} = \rho^{|j-j'|}$, para $j \neq j'$. Um estimador consistente para ρ é dado por:

$$\hat{\rho} = \frac{\phi}{n} \sum_{i=1}^n \frac{1}{(r_i - 1)} \sum_{j=1}^{r_i-1} \frac{(y_{ij} - \hat{\mu}_{ij})}{\sqrt{\hat{V}_{ij}}} \frac{(y_{i(j+1)} - \hat{\mu}_{i(j+1)})}{\sqrt{\hat{V}_{i(j+1)}}}. \quad (15)$$



1 Identidade:

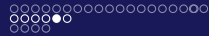
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

1 Identidade:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

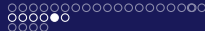
2 Simétrica ou permutável:

$$\begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix}$$



1 Autoregressiva:

$$\begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{bmatrix}$$

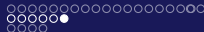


1 Autoregressiva:

$$\begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{bmatrix}$$

2 Não estruturada:

$$\begin{bmatrix} 1 & \rho_{12} & \rho_{13} \\ \rho_{12} & 1 & \rho_{23} \\ \rho_{13} & \rho_{23} & 1 \end{bmatrix}$$



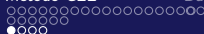
Parâmetro de dispersão

Estimação de ϕ

O parâmetro ϕ^{-1} pode ser estimado consistentemente por

$$\hat{\phi}^{-1} = \frac{1}{N - p} \sum_{i=1}^n \sum_{j=1}^{r_i} \frac{(y_{ij} - \hat{\mu}_{ij})^2}{\hat{V}_{ij}}, \quad (16)$$

Em que $N = \sum_{i=1}^n r_i$. Testes de hipóteses para β ou para subconjuntos de β podem ser desenvolvidos através de estatísticas tipo **Wald** com a matriz de variância-covariância estimada \hat{V}_G .



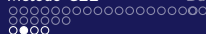
Métodos de Diagnóstico

Resíduos

Em EEGs os **resíduos de Pearson** são calculados da seguinte maneira

$$\hat{r}_{P_{ij}} = \frac{\mathbf{e}_{ij}^{\top} \hat{\mathbf{A}}^{1/2} (\hat{\mathbf{V}}_i \hat{\mathbf{W}}_i)^{-1} (\mathbf{y}_i - \hat{\boldsymbol{\mu}}_i)}{\sqrt{1 - \hat{h}_{ijj}}} \quad (17)$$

para $i = 1, \dots, n$ e $j = 1, \dots, r_i$, em que $\mathbf{A}_i^{1/2} = \phi \mathbf{W}_i^{1/2} \mathbf{R}_i^{-1} \mathbf{W}_i^{1/2}$ é uma matriz de dimensão $r_i \times r_i$, \mathbf{e}_{ij}^{\top} é um vetor de dimensão $1 \times r_i$ de zeros com 1 na j -ésima posição.

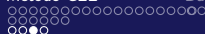


Métodos de Diagnóstico

Alavancagem

Duas medidas de **alavancagem** são usualmente aplicadas em EEGs. Medidas de alavancagem referente ao j -ésimo indivíduo do i -ésimo grupo, dada por \hat{h}_{ijj} e medida de alavancagem referente ao i -ésimo grupo, definida por

$$\hat{h}_i = \frac{1}{r_i} \sum_{j=1}^{r_i} h_{ijj}. \quad (18)$$

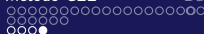


Métodos de Diagnóstico

Observações influentes

Uma versão aproximada da **distância de Cook** para avaliar o impacto da eliminação individual das observações na estimativa de $\hat{\beta}_G$ é dada por

$$LD_{ij} = \frac{\hat{h}_{ijj}}{(1 - \hat{h}_{ijj})} \hat{r}_{P_{ij}}^2. \quad (19)$$



Métodos de Diagnóstico

Critério de seleção de modelos

Uma proposta de critério de seleção de modelos em EEGs é dada por:

$$QIC = -2Q(\hat{\beta}_G) + 2\text{tr}(\hat{\mathbf{V}}_G \hat{\mathbf{H}}_{1I}), \quad (20)$$

em que $\hat{\beta}_G$ é a estimativa de quase-verossimilhança para uma matriz específica de correlação $\mathbf{R}_i(\rho)$ e \mathbf{H}_{1I} é a matriz \mathbf{H}_1 avaliada sob a estrutura de independência.



Dados Ausentes

- Em um estudo longitudinal, alguns sujeitos podem abandonar o estudo antes de seu término, talvez por se mudarem para outra cidade ou por algum motivo que os faça não querer mais participar.
- Uma análise estatística que exclui os dados dos sujeitos para os quais há dados ausentes pode resultar na perda de muitas informações e em erros padrão maiores.
- Analisar apenas os dados observados, como se não houvesse dados ausentes, pode resultar em estimativas de parâmetros viesadas.

Dados Ausentes

Missing Completely at Random (MCAR)

- Se a probabilidade de que uma resposta esteja ausente for a mesma para todos os sujeitos, independentemente dos valores das variáveis explicativas, os dados são *missing completely at random* (**MCAR**).

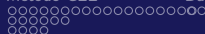
Dados Ausentes

Missing Completely at Random (MCAR)

- Se a probabilidade de que uma resposta esteja ausente for a mesma para todos os sujeitos, independentemente dos valores das variáveis explicativas, os dados são *missing completely at random* (**MCAR**).

Missing at Random (MAR)

- Se a probabilidade de que a resposta esteja ausente variar de acordo com o tempo, mas não variar de acordo com a resposta para sujeitos com os mesmos valores das variáveis explicativas, então os dados não são MCAR, mas são *missing at random* (**MAR**).



Distribuição Binomial

Seja Y_{it} uma variável aleatória com distribuição binomial. Podemos expressá-la na forma da família exponencial:

$$\begin{aligned} f(y; \pi_{it}, \phi) &= \binom{n_{it}}{y_{it}} \pi_{it}^{y_{it}} (1 - \pi_{it})^{(n_{it} - y_{it})} \\ &= \exp \left\{ \log \binom{n_{it}}{y_{it}} + y_{it} \log(\pi_{it}) + (n_{it} - y_{it}) \log(1 - \pi_{it}) \right\} \\ &= \exp \left\{ \log \binom{n_{it}}{y_{it}} + y_{it} \log(\pi_{it}) + n_{it} \log(1 - \pi_{it}) - y_{it} \log(1 - \pi_{it}) \right\} \\ &= \exp \left\{ y_{it} \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right) + \log \binom{n_{it}}{y_{it}} + n_{it} \log(1 - \pi_{it}) \right\} \\ &= \exp \left\{ y_{it} \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right) + n_{it} \log(1 - \pi_{it}) + \log \binom{n_{it}}{y_{it}} \right\}. \end{aligned}$$

Portanto,

$$\phi = 1.$$

$$\theta_{it} = \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right)$$

$$\Rightarrow \exp(\theta_{it}) = \frac{\pi_{it}}{1 - \pi_{it}}$$

$$\Rightarrow \exp(\theta_{it}) - \pi_{it} \exp(\theta_{it}) = \pi_{it}$$

$$\Rightarrow \exp(\theta_{it}) = \pi_{it} + \pi_{it} \exp(\theta_{it})$$

$$\Rightarrow \exp(\theta_{it}) = \pi_{it} (1 + \exp(\theta_{it}))$$

$$\Rightarrow \pi_{it} = \frac{\exp(\theta_{it})}{1 + \exp(\theta_{it})}.$$

$$c(y_{it}; \phi) = \log \left(\frac{n_{it}}{y_{it}} \right).$$

$$b(\theta_{it}) = -n \log(1 - \pi_{it})$$

$$\Rightarrow -n \log \left\{ 1 - \frac{\exp(\theta_{it})}{1 + \exp(\theta_{it})} \right\}$$

$$\Rightarrow -n \log \left\{ \frac{1 + \exp(\theta_{it}) - \exp(\theta_{it})}{1 + \exp(\theta_{it})} \right\}$$

$$\Rightarrow -n \log \left\{ \frac{1}{1 + \exp(\theta_{it})} \right\}$$

$$\Rightarrow -n \{ \log(1) - \log(1 + \exp(\theta_{it})) \}$$

$$\Rightarrow n \log(1 + \exp(\theta_{it})).$$



Esperança

$$\begin{aligned}\mathbb{E}(Y_{it}) &= b'(\theta_{it}) = [n_{it} \log(1 + \exp(\theta_{it}))]' \\ &= n_{it} \frac{1}{1 + \exp(\theta_{it})} \exp(\theta_{it}) \\ &= n_{it} \pi_{it}.\end{aligned}$$

$$\begin{aligned}\mu_{it} &= n_{it} \pi_{it} \\ \Rightarrow \pi_{it} &= \frac{\mu_{it}}{n_{it}}.\end{aligned}$$

Variância

$$\mathbb{V}ar(Y_{it}) = \frac{b''(\theta)}{\phi} = V(\mu_{it}) = n_{it}\pi_{it}(1 - \pi_{it}).$$

$$\Rightarrow n_{it} \frac{\mu_{it}}{n_{it}} \left(1 - \frac{\mu_{it}}{n_{it}}\right)$$

$$\Rightarrow \mu_{it} \left(\frac{n_{it}}{n_{it}} - \frac{\mu_{it}}{n_{it}}\right)$$

$$\Rightarrow \mu_{it} \left(\frac{n_{it} - \mu_{it}}{n_{it}}\right)$$

$$\Rightarrow \frac{\mu_{it}}{n_{it}} (n_{it} - \mu_{it}).$$

Ligação Canônica

Usando a função de ligação canônica, conhecida como logit, temos que

$$\theta_{it} = \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right) = \log \left(\frac{\mu_{it}/n_{it}}{1 - \mu_{it}/n_{it}} \right) = \log \left(\frac{\mu_{it}}{n_{it} - \mu_{it}} \right).$$

$$\eta_{it} = g(\pi_{it}) = \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right) = \mathbf{X}_{it}\boldsymbol{\beta}.$$

$$\eta_{it} = g(\pi_{it}) = \log \left(\frac{\pi_{it}}{1 - \pi_{it}} \right) = \mathbf{X}_{it}\beta.$$

$$\eta_{it} = g \left(\frac{\mu_{it}}{n_{it}} \right) = \log \left(\frac{\mu_{it}}{n_{it} - \mu_{it}} \right) = \mathbf{X}_{it}\beta.$$

$$v_i = \frac{d\mu_{it}}{d\theta_{it}} = \mu_{it} \left(\frac{n_{it} - \mu_{it}}{n_{it}} \right).$$

$$v_i = \frac{d\mu_{it}}{d\theta_{it}} = \mu_{it} \left(\frac{n_{it} - \mu_{it}}{n_{it}} \right).$$

$$w_i = \left(\frac{d\mu_{it}}{d\eta_{it}} \right)^2 \frac{1}{V(\mu_{it})} = \frac{1}{V(\mu_{it})[g'(\mu_{it})]^2} = \frac{\mu_{it}}{n_{it}}(n_{it} - \mu_{it}).$$

Aplicação

Um exemplo de aplicação do ajuste do método *GEE* é apresentado no seguinte link: [▶ Link](#).

Referências I



TIELSCH, J. M. Blindness and Visual Impairment in an American Urban Population: The Baltimore Eye Survey. *Archives of Ophthalmology*, v. 108, n. 2, p. 286, fev. 1990. ISSN 0003-9950. Disponível em:
<<http://archophth.jamanetwork.com/article.aspx?doi=10.1001/archophth.108.2.286>>



AGRESTI, A. *An introduction to categorical data analysis*. Third edition. Hoboken, NJ: John Wiley & Sons, 2019. (Wiley series in probability and statistics). ISBN 9781119405276 9781119405283.

Referências II



LIANG, K.-Y.; ZEGER, S. L. Longitudinal data analysis using generalized linear models. *Biometrika*, v. 73, n. 1, p. 13–22, 04 1986. ISSN 0006-3444. Disponível em:
<<https://doi.org/10.1093/biomet/73.1.13>>.



GUEORGUIEVA, R. V.; AGRESTI, A. A correlated probit model for joint modeling of clustered binary and continuous responses. *Journal of the American Statistical Association*, Taylor & Francis, v. 96, n. 455, p. 1102–1112, 2001.

Referências III



INTRODUCTION to Generalized Estimating Equations — STAT 504. <https://online.stat.psu.edu/stat504/lesson/12/12.1>.

(Acesso em 18/05/2024).



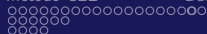
LIANG, K. Y.; ZEGER, S. L. Regression Analysis for Correlated Data. *Annual Review of Public Health*, v. 14, n. 1, p. 43–68, maio 1993. ISSN 0163-7525, 1545-2093. Disponível em:

<<https://www.annualreviews.org/doi/10.1146/annurev.pu.14.050193>>



DOWNLOAD R-4.3.1 for Windows. The R-project for statistical computing. Disponível em:

<<https://cran.r-project.org/bin/windows/base/>>.



Referências IV



HALEKOH, U.; HØJSGAARD, S.; YAN, J. The R Package **geepack** for Generalized Estimating Equations. *Journal of Statistical Software*, v. 15, n. 2, 2006. ISSN 1548-7660. Disponível em: <<http://www.jstatsoft.org/v15/i02/>>.



CAREY, V. J. et al. *gee: Generalized Estimation Equation Solver*. 2024. Disponível em: <<https://cran.r-project.org/web/packages/gee/index.html>>.

Referências V



DAVERN, M. et al. *General Social Survey 1972-2024. [Machine-readable data file]*. Chicago: NORC, 2024. Principal Investigator, Michael Davern; Co-Principal Investigators, Rene Bautista, Jeremy Freese, Pamela Herd, and Stephen L. Morgan. Sponsored by National Science Foundation. NORC ed. Chicago: NORC, 2024: NORC at the University of Chicago [producer and distributor]. Data accessed from the GSS Data Explorer website at gssdataexplorer.norc.org. Disponível em: <<https://gssdataexplorer.norc.org>>.