

ME720 - Modelos Lineares Generalizados

Parte 9 - Inferência em MLGs

Profa. **Larissa Avila Matos**

Estimação dos parâmetros

Estimação dos parâmetros

Uma vez definido cada componente do modelo, obteremos expressões gerais para a função de verossimilhança e para as distribuições assintóticas dos estimadores de máxima verossimilhança dos parâmetros para os MLGs.

Para n observações independentes, temos que a função log-verossimilhança do modelo é dada por $\ell(\boldsymbol{\beta}, \phi) = \sum_{i=1}^n \ell_i(\boldsymbol{\beta}, \phi)$, onde

$$\ell_i(\boldsymbol{\beta}, \phi) = \frac{[y_i \theta_i - b(\theta_i)]}{a(\phi)} + c(y_i, \phi).$$

Então,

$$\ell(\boldsymbol{\beta}, \phi) = \sum_{i=1}^n \frac{[y_i \theta_i - b(\theta_i)]}{a(\phi)} + \sum_{i=1}^n c(y_i, \phi).$$

Para um GLM $\eta_i = \sum_{j=1}^p \beta_j x_{ij} = g(\mu_i)$, com função de ligação g . Portanto, o sistema de equações de verossimilhança para β é dado por

$$\frac{\partial \ell(\beta, \phi)}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \beta_j} = 0, \quad \forall j.$$

Para diferenciar a log-verossimilhança, usamos a regra da cadeia,

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}.$$

Uma vez que,

$$\frac{\partial \ell_i}{\partial \theta_i} = \frac{[y_i - b'(\theta_i)]}{a(\phi)} \quad \text{e} \quad \mu_i = b'(\theta_i), \quad \text{Var}(Y_i) = b''(\theta_i)a(\phi),$$

temos que

$$\frac{\partial \ell_i}{\partial \theta_i} = \frac{(y_i - \mu_i)}{a(\phi)} \quad \text{e} \quad \frac{\partial \mu_i}{\partial \theta_i} = b''(\theta_i) = \frac{\text{Var}(Y_i)}{a(\phi)}.$$

Também uma vez que $\eta_i = \sum_{j=1}^p \beta_j x_{ij}$, então $\frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$.

Finalmente, $\eta_i = g(\mu_i) \Rightarrow \frac{\partial \mu_i}{\partial \eta_i}$, depende da função de ligação para o modelo escolhido.

Resumindo, temos que

$$\begin{aligned} \frac{\partial \ell_i}{\partial \beta_j} &= \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j} \\ &= \frac{(y_i - \mu_i)}{a(\phi)} \frac{a(\phi)}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ij} = \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} \end{aligned}$$

A soma das n observações produz o sistema de equações de verossimilhança para um MLG.

Equações de verossimilhança para um MLG

Equações de verossimilhança para um MLG:

$$\frac{\partial \ell(\boldsymbol{\beta}, \phi)}{\partial \beta_j} = \sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} = 0, \quad j = 1, 2, \dots, p, \quad \text{e}$$

$$\frac{\partial \ell(\boldsymbol{\beta}, \phi)}{\partial \phi} = 0,$$

onde $\eta_i = \sum_{j=1}^p \beta_j x_{ij} = g(\mu_i)$ para uma função de ligação g .

Seja V a matriz diagonal de variâncias das n observações, e seja D uma matriz diagonal com os elementos de $\frac{\partial \mu_i}{\partial \eta_i}$.

Para as expressões de um MLG $\eta = X\beta$, com matriz de planejamento X , as equações de verossimilhança têm a forma

$$XDV^{-1}(y - \mu) = 0.$$

Apesar de β não aparecer nessas equações, ele aparece implicitamente através de $\mu_i = g^{-1}(\sum_{j=1}^p \beta_j x_{ij})$.

Diferentes funções de ligação gera diferentes conjuntos de equações.

Essas equações são funções não lineares de β , e esse problema deve ser resolvido iterativamente. Há várias opções de algoritmos numéricos para resolução de sistemas não lineares: Newton-Raphson, Escore de Fisher (EF), Nelder-Mead, BFGS entre outros.

Exercício

- Encontrar as equações de máxima verossimilhança para o Modelo Poisson Log Linear.

Relação entre média e variância

Os sistema de equações de máxima verossimilhança depende da distribuição de Y_i somente pela média ($\mathbb{E}(Y_i)$) e pela variância ($\text{Var}(Y_i)$).

Além disso, a variância depende da média pela forma

$$\text{Var}(Y_i) = V(\mu_i),$$

para alguma função $V(\cdot)$.

Ou seja, a relação entre a média e a variância caracteriza a distribuição de Y_i .

Exemplo: Se Y_i tem distribuição pertencente a família exponencial e $\text{Var}(Y_i) = \mu_i$, então necessariamente Y_i tem distribuição de Poisson.

Distribuição Assintótica de $\hat{\beta}$

Através das propriedades de máxima verossimilhança, e sob condições de regularidade, para n grande o estimador de máxima verossimilhança $\hat{\beta}$ de β para um MLG é eficiente e tem distribuição Normal.

Temos que a matrix de informação \mathbf{I} , tem elementos dados por

$$\begin{aligned}\mathbb{E} \left(-\frac{\partial^2 \ell_i}{\partial \beta_j \partial \beta_k} \right) &= \mathbb{E} \left(\frac{\partial \ell_i}{\partial \beta_j} \frac{\partial \ell_i}{\partial \beta_k} \right) \\&= \mathbb{E} \left(\frac{(Y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} \frac{(Y_i - \mu_i) x_{ik}}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} \right) \\&= \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)^2} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 \underbrace{\mathbb{E}((Y_i - \mu_i)^2)}_{=\text{Var}(Y_i)} = \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2.\end{aligned}$$

A matriz \mathbf{I} é chamada de matriz de informação esperada.

Então,

$$\mathbb{E} \left(-\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k} \right) = \sum_{i=1}^n \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2$$

Seja, \mathbf{W} uma matriz diagonal com elementos dados por

$$w_i = \frac{(\partial \mu_i / \partial \eta_i)^2}{\text{Var}(Y_i)}$$

$$\Rightarrow \mathbf{I} = \mathbf{X}' \mathbf{W} \mathbf{X}$$

A forma de \mathbf{W} depende da função de ligação $g(\cdot)$, uma vez que $\partial \mu_i / \partial \eta_i = g'(\mu_i)$.

Portanto, a matriz de covariância de $\hat{\boldsymbol{\beta}}$ é dada pela inversa da matriz de informação \mathbf{I} .

Distribuição Assintótica de $\hat{\beta}$ para um MLG $\eta = X\beta$:

$\hat{\beta}$ tem distribuição aproximadamente normal $N(\beta, (X'WX)^{-1})$,

onde W é a matriz diagonal com elementos $w_i = \frac{(\partial\mu_i/\partial\eta_i)^2}{\text{Var}(Y_i)}$.

A matriz de covariância assintótica é estimada por $\widehat{\text{Var}}(\hat{\beta}) = (X'\widehat{W}X)^{-1}$, onde \widehat{W} é W avaliado em $\hat{\beta}$.

Obs:

- 1 Para a FE com parâmetro de escala, θ e ϕ são parâmetros ortogonais.
- 2 $\hat{\beta}$ e $\hat{\phi}$ são assintoticamente independentes.

Matriz de covariância para os valores ajustados

O preditor linear estimado é dado por

$$\hat{\eta} = \mathbf{X}\hat{\beta}.$$

Para n grande, temos

$$\text{Var}(\hat{\eta}) = \mathbf{X}'\text{Var}(\hat{\beta})\mathbf{X} \approx \mathbf{X}'(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}.$$

Podemos obter a variância assintótica de $\hat{\mu}$ ($\text{Var}(\hat{\mu})$) por $\text{Var}(\hat{\eta})$, através do método delta. Então,

$$\text{Var}(\hat{\mu}) \approx \mathbf{D}'\text{Var}(\hat{\eta})\mathbf{D} \approx \mathbf{D}\mathbf{X}'(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}\mathbf{D},$$

onde \mathbf{D} é uma matriz diagonal com elementos $\partial\mu_i/\partial\eta_i$.

Exercício: Pesquisar sobre o método delta.

Exercício

- Voltando ao exemplo da Poisson, encontre os elementos da matriz W .

Estimação de β

Como encontramos os estimadores de máxima verossimilhança dos parâmetros de um MLG?

Problema: O sistema de equações são geralmente não-lineares em β .

Solução: Métodos iterativos para resolver sistema de equações não lineares. Focaremos em dois métodos:

- *Newton-Raphson*

- *Escore de Fisher*

Método de Newton-Raphson

O algoritmo Newton-Raphson, método de Newton-Raphson, foi desenvolvido por Isaac Newton e Joseph Raphson e tem o objetivo estimar as raízes de uma função.

Suponha que queremos encontrar a solução da equação $g(x_0) = 0$, onde g é uma função diferenciável. Dado um número x próximo de x_0 , segue da expansão em série de Taylor em torno de x que

$$0 = g(x_0) \approx g(x) + g'(x)(x_0 - x).$$

Resolvendo para x_0 , temos

$$x_0 \approx x - \frac{g(x)}{g'(x)}.$$

Assim, dado um valor estimado x_t , então podemos ter um novo valor estimado x_{t+1} por

$$x_{t+1} \approx x_t - \frac{g(x_t)}{g'(x_t)}.$$

Este procedimento é repetido para $t = 1, 2, 3, \dots$ até $|g(x_t)/g'(x_t)|$ ser suficientemente pequeno.

Método de Newton-Raphson

Voltando ao nosso problema, queremos encontrar a solução da equação $\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0}$.

No passo t , o processo iterativo ($t = 0, 1, 2, \dots$) aproxima $\ell(\boldsymbol{\beta})$ próximo de $\boldsymbol{\beta}^{(t)}$ pela expansão em série de Taylor de segunda ordem,

$$\ell(\boldsymbol{\beta}) \approx \ell(\boldsymbol{\beta}^{(t)}) + \mathbf{u}^{(t)}(\boldsymbol{\beta} - \boldsymbol{\beta}^{(t)}) + \left(\frac{1}{2}\right) (\boldsymbol{\beta} - \boldsymbol{\beta}^{(t)})' \mathbf{H}^{(t)} (\boldsymbol{\beta} - \boldsymbol{\beta}^{(t)}),$$

onde $\mathbf{u}^{(t)}$ e $\mathbf{H}^{(t)}$ são \mathbf{u} e \mathbf{H} avaliados em $\boldsymbol{\beta}^{(t)}$ respectivamente, com

- $\mathbf{u} = \left(\frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_1}, \frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_2}, \dots, \frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_p} \right)'$, e
- \mathbf{H} a matriz *Hessiana*, onde $H_{ij} = \frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \beta_i \partial \beta_j}$.

Obs: $\boldsymbol{\beta}^{(0)}$ valor inicial (chute inicial).

Resolvendo, $\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \approx \mathbf{u}^{(t)} + \mathbf{H}^{(t)}(\boldsymbol{\beta} - \boldsymbol{\beta}^{(t)}) = \mathbf{0}$, temos a seguinte aproximação

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - (\mathbf{H}^{(t)})^{-1} \mathbf{u}^{(t)},$$

assumindo que $\mathbf{H}^{(t)}$ é não singular.

O procedimento descrito é repetido até que mudanças em $\ell(\boldsymbol{\beta}^{(t)})$ entre ciclos sucessivos são suficientemente pequenas.

Para muitos MLGs, a matriz Hessian é negativa definida, e a log verossimilhança é uma função estritamente côncava. Então, as estimativas de máxima verossimilhança dos parâmetros do modelo existem e são únicas sob condições bastante gerais. A convergência de $\boldsymbol{\beta}^{(t+1)}$ para $\hat{\boldsymbol{\beta}}$ na vizinhança de $\hat{\boldsymbol{\beta}}$ é rápida.

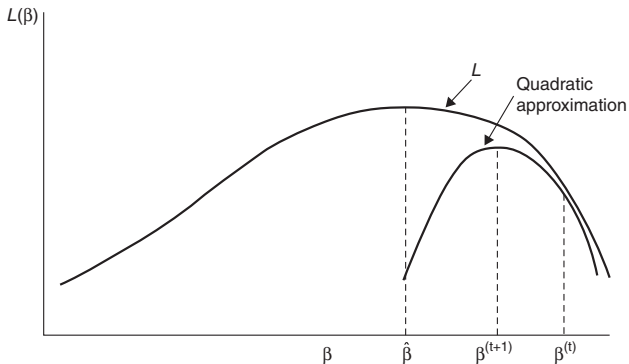


Figure 4.2 Illustration of a cycle of the Newton–Raphson method.

Figura ilustrativa do livro texto (Agresti, A. (2015)).

Método de Escore de Fisher

O método de Escore de Fisher é um método iterativo alternativo para resolver o sistema de equações de verossimilhança.

A diferença entre o método de Escore de Fisher e o método de Newton-Raphson está na maneira como escolhemos a matriz *Hessiana*.

O método de Escore de Fisher usa o valor esperado da matriz *Hessiana*, chamada de matriz de informação esperada, enquanto método de Newton-Raphson usa a própria matriz *Hessiana*, chamada de matriz de informação observada.

Portanto, temos a seguinte aproximação

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + (\mathbf{I}^{(t)})^{-1} \mathbf{u}^{(t)},$$

onde $\mathbf{I}^{(t)}$ é \mathbf{I} avaliado em $\boldsymbol{\beta}^{(t)}$, ou seja $\mathbf{I}^{(t)}$ tem elementos $-\mathbb{E}(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \beta_i \partial \beta_j})$ avaliado em $\boldsymbol{\beta}^{(t)}$.

Os algoritmos de Escore de Fisher e de Newton-Raphson são idênticos para os MLGs que usam a função de ligação canônica (Nelder e Wedderburn, 1972).

Newton-Raphson e Escore de Fisher

Newton-Raphson e Escore de Fisher para o parâmetro de uma binomial

Exemplo para o método de Newton-Raphson e Escore de Fisher de um problema mais simples para o qual sabemos a resposta.

Newton-Raphson e Escore de Fisher

Newton-Raphson e Escore de Fisher para o parâmetro de uma binomial

Exemplo para o método de Newton-Raphson e Escore de Fisher de um problema mais simples para o qual sabemos a resposta.

Seja Y uma a.a de uma distribuição $Bin(n, p)$, a função de log-verossimilhança é dada por

$$\ell(p) = \log(p^{nY}(1-p)^{n-nY}) = ny \log(p) + (n - ny) \log(1-p).$$

Sabemos que

$$u = \frac{\partial \ell}{\partial p} = \frac{(ny - np)}{p(1-p)} \quad \text{e} \quad H = \frac{\partial^2 \ell}{\partial p^2} = - \left[\frac{ny}{p^2} + \frac{n - ny}{(1-p)^2} \right].$$

Maximizando a log-verossimilhança temos que o estimador de MV para p é $\hat{p} = y$.

Cada passo do algoritmo de Newton-Raphson é dado por

$$p^{(t+1)} = p^{(t)} + \left[\frac{ny}{(p^{(t)})^2} + \frac{n - ny}{(1 - p^{(t)})^2} \right]^{-1} \frac{(ny - np^{(t)})}{p^{(t)}(1 - p^{(t)})}.$$

- Se, $p^{(0)} = \frac{1}{2} \Rightarrow p^{(1)} = y$.
- Quando $p^{(t)} = y$, temos que $p^{(t+1)} = y$, ou seja, não temos nenhuma alteração da iteração t para $t + 1$, o qual é o estimador de MV.

Calculando a esperança de H , temos que

$$\mathbb{I} = \frac{n}{p(1-p)},$$

e cada passo do algoritmo de Escore de Fisher é dado por

$$\begin{aligned} p^{(t+1)} &= p^{(t)} + \left[\frac{n}{p^{(t)}(1-p^{(t)})^2} \right]^{-1} \frac{(ny - np^{(t)})}{p^{(t)}(1-p^{(t)})} \\ &= p^{(t)} + (y - p^{(t)}) = y \end{aligned}$$

- Ou seja, $p^{(t+1)}$ é o estimador de MV após apenas uma iteração e ficando nesse valor em todas as iterações.

Estimação do parâmetro de escala

- 1 Maximizar $\ell(\beta, \phi)$ com respeito a ϕ . (Muito sensível a suposição da distribuição)
- 2 Sabemos que $\text{Var}(Y_i) = a(\phi)V(\mu_i)$, então

$$\frac{\mathbb{E}((Y_i - \mu_i)^2)}{V(\mu_i)} = a(\phi)$$

$$\Rightarrow \hat{a}(\phi) = \frac{1}{n-p} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_i)^2}{\text{Var}(\hat{\mu}_i)}.$$

Teste de hipóteses

Teste de hipóteses

A inferência para MLGs tem três abordagens usando a função de verossimilhança:

- 1 Teste da razão de verossimilhanças;
- 2 Teste de Wald;
- 3 Teste de Escore.

Focaremos em testes

$$H_0 : \beta = \beta_0 \quad \text{vs.} \quad H_1 : \beta \neq \beta_0.$$

Teste da razão de verossimilhanças

O teste da razão de verossimilhanças usa a função de verossimilhança através da razão de (1) seu valor L_0 em β_0 , e (2) seu valor máximo L_1 nos valores de β permitindo que H_0 ou H_1 sejam verdadeiras.

A razão $\Lambda = \frac{L_0}{L_1} \leq 1$, uma vez que L_0 resulta da maximização em um valor restrito de β . A estatística do teste da razão de verossimilhanças é dada por

$$-2 \log(\Lambda) = -2 \log \left(\frac{L_0}{L_1} \right) = -2(\ell_0 - \ell_1),$$

onde ℓ_0 e ℓ_1 denotam as funções de log-verossimilhança maximizadas. Sob condições de regularidade, $-2 \log(\Lambda)$ possui distribuição qui-quadrado quando $n \rightarrow \infty$, com $g.l = 1$ Então,

$$-2 \log(\Lambda) \sim \chi_1^2, \quad n \rightarrow \infty.$$

O p - valor é a probabilidade da qui-quadrado ser maior igual ao valor da estatística do teste observado, ou seja, p - valor $= \mathbb{P}(\chi_1^2 \geq \chi_{obs})$.

Este teste se estende diretamente a vários parâmetros. Por exemplo, para $\beta = (\beta_0, \beta_1)$, considere $H_0 : \beta_0 = 0$.

Então, L_1 é a função de verossimilhança calculada no valor de β pelo qual os dados seriam mais prováveis e L_0 é a função de verossimilhança calculada no valor de β para o qual os dados seriam mais prováveis quando $\beta = 0$.

O grau de liberdade da distribuição da estatística do teste é igual à diferença nas dimensões dos espaços paramétricos em $H_0 \cup H_1$ e em H_0 .

O teste também se estende à hipótese linear geral

$$H_0 : C\beta = 0,$$

uma vez que as restrições lineares implicam um novo modelo que é um caso especial do original.

Teste de Wald

Erros padrões obtidos a partir da inversa da matriz de informação dependem dos valores desconhecidos dos parâmetros.

Quando substituimos as estimativas irrestritas de MV (ou seja, não assumindo a hipótese nula), obtemos um erro padrão estimado (SE) de $\hat{\beta}$.

Para $H_0 : \beta = \beta_0$ a estatística do teste usando esse erro padrão estimado não nulo,

$$z = \frac{\hat{\beta} - \beta_0}{SE},$$

é chamada de estatística de Wald e segue aproximadamente uma distribuição normal padrão quando $\beta = \beta_0$.

Além disso, z^2 tem aproximadamente uma distribuição qui-quadrado com 1 grau de liberdade, $z^2 \sim \chi_1^2$.

Para vários parâmetros $\beta = (\beta_0, \beta_1)$, para testar $H_0 : \beta_0 = \mathbf{0}$, a estatística qui-quadrado de Wald é

$$\widehat{\beta}_0' \left[\text{Var}(\widehat{\beta}_0) \right]^{-1} \widehat{\beta}_0,$$

onde $\widehat{\beta}_0$ é a estimativa de máxima verossimilhança irrestrita de β_0 e $\text{Var}(\widehat{\beta}_0)$ é um bloco da matriz de covariância estimada irrestrita de $\widehat{\beta}$.

Teste de Escore

O teste de Escore, usa a inclinação (ou seja, a função Escore) e a curvatura esperada da função de log-verossimilhança, avaliada no valor da hipótese nula β_0 .

A estatística do teste é dada por

$$\frac{[\partial \ell(\beta)/\partial \beta_0]^2}{-\mathbb{E} [\partial^2 \ell(\beta)/\partial \beta_0^2]} \sim \chi_1^2.$$

Para vários parâmetros $\beta = (\beta_0, \beta_1)$, para testar $H_0 : \beta_0 = \mathbf{0}$, a estatística qui-quadrado do teste de Escore é uma forma quadrática baseada no vetor de derivadas parciais da função de log-verossimilhança e na inversa da matriz de informação, ambas avaliadas nas estimativas sob H_0 .

Função Desvio

Função Desvio ou Desvio

A função desvio é útil para termos uma idéia da adequabilidade do modelo (verificação da qualidade do ajuste).

Sem perda de generalidade, seja $\ell(\boldsymbol{\mu}; \mathbf{y}) \equiv \ell(\boldsymbol{\beta}, \phi; \mathbf{y})$ a log verossimilhança associada ao modelo em estudo (a log verossimilhança expressada em termos da média), então

$$\ell(\boldsymbol{\mu}; \mathbf{y}) = \sum_{i=1}^n \ell(\mu_i; y_i),$$

em que $\mu_i = g^{-1}(\eta_i)$ e $\eta_i = \mathbf{x}_i' \boldsymbol{\beta} = \sum_{j=1}^p x_{ij} \beta_j$.

Seja também $\ell(\mathbf{y}; \mathbf{y})$ a log verossimilhança do modelo saturado ($n = p$), em que cada média é representada por ela mesma.

A função $\ell(\boldsymbol{\mu}; \mathbf{y})$ é estimada por

$$\ell(\mathbf{y}; \mathbf{y}) = \sum_{i=1}^n \ell(y_i; y_i),$$

ou seja, a estimativa de máxima verossimilhança de μ_i fica nesse caso dada por $\hat{\mu}_i = y_i$.

Quando $p < n$, denotamos a estimativa de $\ell(\boldsymbol{\mu}; \mathbf{y})$ por $\ell(\hat{\boldsymbol{\mu}}; \mathbf{y})$.

O modelo saturado, explica toda variação pelo preditor linear do modelo.

Um modelo perfeito parece bom, mas o modelo saturado não é útil. Ele não suaviza os dados ou tem as vantagens de um modelo mais simples devido à sua parcimônia.

No entanto, muitas vezes serve como linha de base para comparação com outros modelos, como para verificar a bondade de ajuste.

A qualidade do ajuste de um MLG é avaliada através da função desvio

$$D^*(\mathbf{y}; \hat{\boldsymbol{\mu}}) = 2 [\ell(\mathbf{y}; \mathbf{y}) - \ell(\hat{\boldsymbol{\mu}}; \mathbf{y})],$$

que é uma distância entre o logaritmo da função de verossimilhança do modelo saturado (com n parâmetros) e do modelo sob investigação (com p parâmetros) avaliado na estimativa de máxima verossimilhança $\hat{\boldsymbol{\beta}}$.

Uma vez que $\ell(\hat{\boldsymbol{\mu}}; \mathbf{y}) \leq \ell(\mathbf{y}; \mathbf{y})$, então $D^*(\mathbf{y}; \hat{\boldsymbol{\mu}}) \geq 0$.

Um valor pequeno para a função desvio indica que, para um número menor de parâmetros, obtemos um ajuste tão bom quanto o ajuste com o modelo saturado.

Denotando por $\hat{\theta}_i = \theta_i(\hat{\mu}_i)$ e $\tilde{\theta}_i = \theta_i(\tilde{\mu}_i)$ as estimativas de máxima verossimilhança de θ_i para os modelos com p parâmetros ($p < n$) e saturado ($p = n$), respectivamente, temos que a função $D^*(\mathbf{y}; \boldsymbol{\mu})$ fica, alternativamente, dada por

$$\begin{aligned} D^*(\mathbf{y}; \hat{\boldsymbol{\mu}}) &= 2 \sum_{i=1}^n \left[\frac{y_i \tilde{\theta}_i - b(\tilde{\theta}_i)}{a(\phi)} \right] - 2 \sum_{i=1}^n \left[\frac{y_i \hat{\theta}_i - b(\hat{\theta}_i)}{a(\phi)}; \mathbf{y} \right] \\ &= 2 \sum_{i=1}^n \left[\frac{y_i (\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i)}{a(\phi)} \right]. \end{aligned}$$

Usualmente $a(\phi) = \frac{\phi}{w_i}$, então

$$\frac{D(\mathbf{y}; \hat{\boldsymbol{\mu}})}{\phi} = 2 \sum_{i=1}^n w_i \left[y_i (\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i) \right],$$

$\frac{D(\mathbf{y}; \hat{\boldsymbol{\mu}})}{\phi}$ é chamado de desvio escalonado. A estatística $D(\mathbf{y}; \hat{\boldsymbol{\mu}})$ é chamada de desvio.

- **Obs.:** O desvio no R sai com o nome de deviance após o ajuste do modelo e o número de graus de liberdade correspondente é dado por $n - p$.

Exemplo Livro

O que afeta o preço de venda de uma casa?

Os dados a seguir mostram as observações sobre as vendas de casas recentes em Gainesville, Flórida. O arquivo de dados para as 100 vendas de casas estão no site do livro texto.

As variáveis listadas são:

- o preço de venda (em milhares de dólares),
- o tamanho da casa (em pés quadrados),
- a conta anual do imposto predial (em dólares),
- o número de quartos,
- o número de banheiros, e
- se a casa é nova.

Como essas 100 observações são de uma única cidade, não podemos usá-las para fazer inferências sobre as relações em geral.

Mas, para fins ilustrativos, os tratamos como uma amostra aleatória de uma população conceitual de vendas de casas nesse mercado e analisamos como o preço de venda parece se relacionar com essas características.

A seguir são apresentados os dados de 5 casas do conjunto de dados e algumas análises iniciais.

```
Houses <- as.data.frame(read.table('Houses.txt', header = T))
Houses[1:5,]
```

	case	taxes	beds	baths	new	price	size
1	1	3104	4	2	0	279.9	2048
2	2	1173	2	1	0	146.5	912
3	3	3076	4	2	0	237.7	1654
4	4	1608	3	2	0	200.0	2068
5	5	1454	3	3	0	159.9	1477

```
str(Houses)
```

```
'data.frame':  100 obs. of  7 variables:
 $ case : int  1 2 3 4 5 6 7 8 9 10 ...
 $ taxes: int  3104 1173 3076 1608 1454 2997 4054 3002 6627 320 ...
 $ beds : int  4 2 4 3 3 3 3 3 5 3 ...
 $ baths: int  2 1 2 2 3 2 2 2 4 2 ...
 $ new  : int  0 0 0 0 0 1 0 1 0 0 ...
 $ price: num  280 146 238 200 160 ...
 $ size : int  2048 912 1654 2068 1477 3153 1355 2075 3990 1160 ...
```

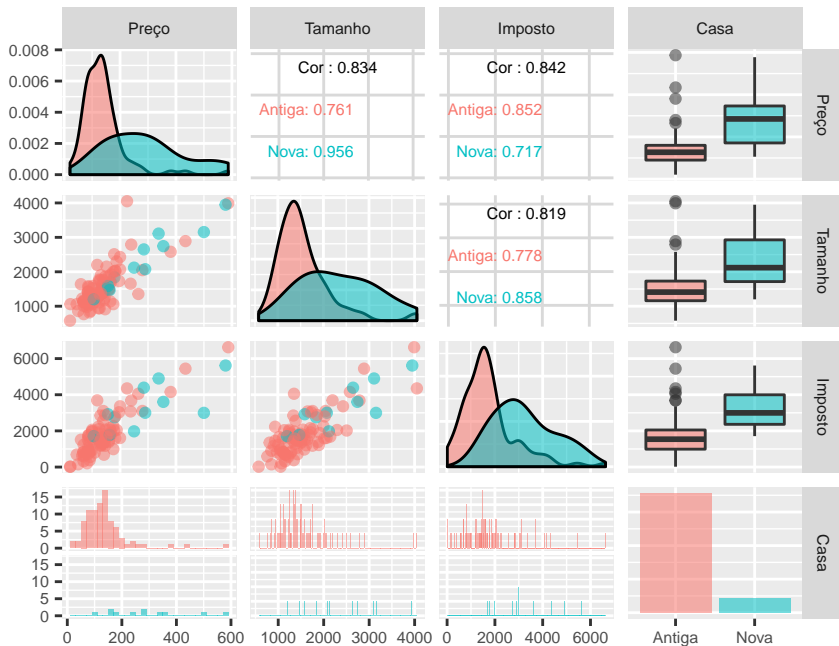
```
cbind(mean(price), sd(price), mean(size), sd(size))
```

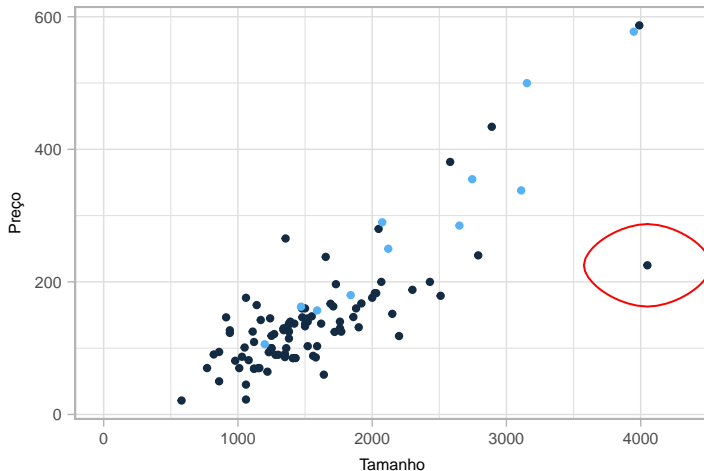
```
      [,1]      [,2]      [,3]      [,4]  
[1,] 155.331 101.2622 1629.28 666.9417
```

```
table(new)
```

```
new
```

```
 0  1  
89 11
```





Temos aproximadamente uma tendência linear crescente para o preço de venda em função do tamanho. Uma exceção é um preço de venda relativamente baixo para uma residência muito grande que não era nova (observação 64 no campo dos dados). Apenas 11 casas da amostra eram novas, portanto o impacto dessa variável é pouco claro.

Em seguida, ajustamos o modelo $E(Y_i) = \beta_0 + \beta_1 \text{ size}_{i1} + \beta_2 \text{ size}_{i2}$, tendo efeitos aditivos dessas variáveis explicativas.

```
fit <- lm(Houses$price ~ Houses$size + Houses$new)
summary(fit)
```

Call:

```
lm(formula = Houses$price ~ Houses$size + Houses$new)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-205.102	-34.374	-5.778	18.929	163.866

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-40.230867	14.696140	-2.738	0.00737 **
Houses\$size	0.116132	0.008795	13.204	< 2e-16 ***
Houses\$new	57.736283	18.653041	3.095	0.00257 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 53.88 on 97 degrees of freedom

Multiple R-squared: 0.7226, Adjusted R-squared: 0.7169

F-statistic: 126.3 on 2 and 97 DF, p-value: < 2.2e-16

Resultando num ajuste: $\hat{\mu} = -40.231 + 0.116 \text{ size}_{i1} + 57.736 \text{ size}_{i2}$.

- [Notas](#) de aula do Prof. Gilberto de Paula.
- Agresti, A. (2015). *Foundations of Linear and Generalized Linear Models*. Wiley series in probability and statistics.