

Modelos de regressão para dados de contagem inflacionados de zeros

Prof. Caio Azevedo

Exemplo 16: Número de artigos produzidos por alunos de Doutorado em bioquímica

- Corresponde aos dados relacionados à 915 alunos de pós-graduação em bioquímica.
- Disponível no pacote *pscl* do *R*, sob o nome “bioChemists”.
- Variáveis medidas:
 - art: quantidade de artigos produzidos nos últimos 3 anos de doutorado pelo aluno (PhD) (resposta).
 - fem: sexo do aluno - masculino ou feminino (explicativa).
 - mar: estado civil - solteiro ou casado (explicativa).

Exemplo 16: Número de artigos produzidos por alunos de Doutorado em bioquímica

- kid5: número de filhos com 5 ou menos anos de idade (explicativa).
- phd: prestígio do departamento onde o aluno desenvolveu seus estudos (valor observado entre 0 e 5) (explicativa).
- ment: quantidade de artigos produzidos nos últimos de 3 anos pelo orientador (explicativa)

Comentários

- Espera-se uma concentração (inflacionamento) de zeros, dado que, em geral, os alunos tendem a publicar (mais) artigos depois de ter finalizado o Doutorado.
- Neste caso, os modelos de regressão para dados de contagem **Poisson** e **Binomial-negativo**, por exemplo, podem não ser apropriados.
- Alternativa: modificar esses modelos a fim de contemplar o inflacionamento no valor zero.

Modelos probabilísticos

- Seja Y a vad (variável aleatória discreta) que representa a contagem de interesse.
- Defina

$$\begin{aligned}P(Y = y) &= g_Y(y) \\ &= [\pi + (1 - \pi)f_Z(0)] \mathbb{1}_{\{0\}}(y) + (1 - \pi)f_Z(y)\mathbb{1}_{\{1,2,\dots\}},\end{aligned}$$

em que $f_Z(\cdot)$ representa a função de probabilidade de uma vad discreta de interesse (Poisson, geométrica, binomial negativa, etc).

Modelos probabilísticos

- Notação: $Y \sim IZ\{\pi, f_Z(y)\}$ em que $f_Z(\cdot)$ representa o modelo probabilístico original.
- Por exemplo $Y \sim IZ\{\pi, \text{Poisson}(\mu)\}$, ou seja:

$$g_Y(y) = [\pi + (1 - \pi)e^{-\mu}] \mathbb{1}_{\{0\}}(y) + (1 - \pi) \frac{e^{-\mu} \mu^y}{y!} \mathbb{1}_{\{1,2,\dots\}}.$$

Modelos probabilísticos

- Neste caso π representa a probabilidade de inflacionamento.
- Note que

$$\begin{aligned}\sum_{y=0}^{\infty} g(y) &= \pi + (1 - \pi)f(0) + (1 - \pi) \sum_{y=1}^{\infty} [f(y)] \\ &= \pi + (1 - \pi)f(0) + (1 - \pi)(1 - f(0)) \\ &= \pi + (1 - \pi) = 1.\end{aligned}$$

Modelos probabilísticos

- Valor esperado

$$\mathcal{E}(Y) = \sum_{y=0}^{\infty} yg(y) = (1 - \pi) \sum_{y=1}^{\infty} yf(y) = (1 - \pi)\mathcal{E}(Z),$$

em que $\mathcal{E}(Z)$ representa a esperança da variável original (não inflacionada).

- Poisson e binomial-negativa (esta última sob a parametrização adotada em nosso curso): $\mathcal{E}(Y) = (1 - \pi)\mu$.

Modelos probabilísticos

■ Variância

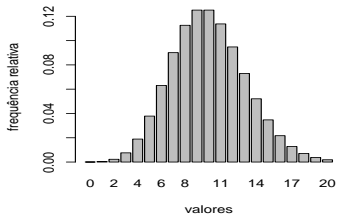
$$\begin{aligned}\mathcal{V}(Y) &= \mathcal{E}(Y^2) - \mathcal{E}^2(Y) = \sum_{y=0}^{\infty} y^2 g(y) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) \sum_{y=1}^{\infty} y^2 g(y) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) \mathcal{E}(Z^2) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) [\mathcal{E}(Z^2) - (1 - \pi) \mathcal{E}^2(Z)].\end{aligned}$$

■ Poisson: $\mathcal{V}(Y) = (1 - \pi)\mu(1 + \mu\pi)$.

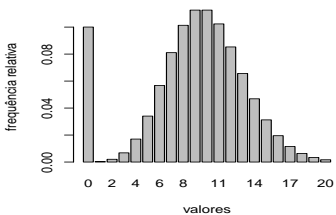
■ Binomial-negativa: $\mathcal{V}(Y) = (1 - \pi)\mu \left(1 + \frac{\mu}{\phi} + \mu\pi \right)$.

Exemplos de funções de probabilidade

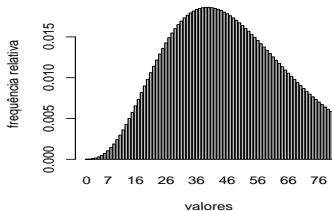
Poisson, $\mu = 10$



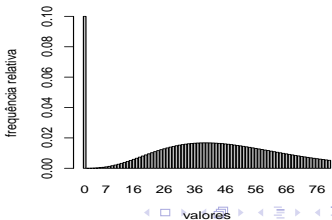
Poisson inflacionada, $\mu = 10, \pi = 0.1$



Binomial-negativa, $\mu = 10, \phi = 5$

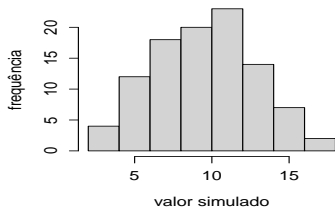


Binomial-negativa, $\mu = 10, \phi = 5, \pi = 0,1$

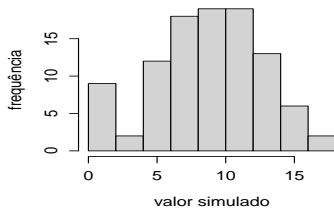


Histograma de valores simulados

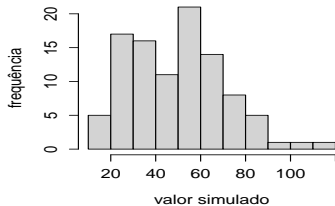
Poisson, $\mu = 10$



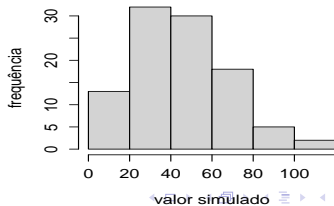
Poisson inflacionada, $\mu = 10, \pi = 0,1$



Binomial-negativa, $\mu = 10, \phi = 5$



Binomial-negativa, $\mu = 10, \phi = 5, \pi = 0.1$



Forma alternativa de representação

- Seja $U \sim \text{Bernoulli}(\pi)$, em que:
 - $Y|U = 1 \sim X_1$, em que $P(X_1 = 0) = 1$, ou seja $P(Y = 0|U = 1) = 1$.
 - e $Y|U = 0 \sim X_2$, em que X_2 segue a distribuição à contagem original (Poisson, geométrica, binomial-negativa). Assim

$$P(Y = y, U = u) = \pi^u [(1 - \pi)f_Z(y)]^{1-u}. \quad (1)$$

- Pode-se provar que $P(Y = y) = \sum_{u=0}^1 P(Y = y, U = u) = [\pi + (1 - \pi)f(0)] \mathbb{1}_{\{0\}}(y) + (1 - \pi)f(y) \mathbb{1}_{\{1,2,\dots\}}$ (exercício).

Modelo de regressão para contagens inflacionadas de zeros

- Consideramos $Y_1, \dots, Y_n \stackrel{ind.}{\sim} IZ\{\pi_i, f_{Z_i}(y_i)\}$.
- Podemos modelar tanto os parâmetros relativos à distribuição $f_Z(\cdot)$ quanto o parâmetro π .
- Por exemplo, se $Z \sim \text{Poisson}(\mu_i)$ então podemos considerar

$$\begin{aligned}\ln(\mu_i) &= \mathbf{x}_i' \boldsymbol{\beta}, \\ \text{logito}(\pi_i) &= \ln\left(\frac{\pi_i}{1 - \pi_i}\right) = \mathbf{w}_i' \boldsymbol{\gamma}.\end{aligned}$$

Cont.

- Estruturas similares podem ser utilizadas considerando-se outras funções de ligação e/ou outras distribuições para a contagem original.
- Se houver mais parâmetros na distribuição de Z (caso da distribuição binomial negativa), eles também podem ser modelados através de uma estrutura de regressão.

Inferência via MV

- Estimação paramétrica:
 - Otimização direta da log-verossimilhança.
 - Algoritmo EM, otimizando-se a log-verossimilhança aumentada definida pela distribuição conjunta de $(Y, U)'$.
- Utilizaremos a segunda opção por, dentre outros fatores, ser, em geral, mais estável.

log-verossimilhança

- Sejam θ os parâmetros da contagem original (Z) e γ os parâmetros associados à probabilidade de inflacionamento (ou seja, $\pi = h(\gamma)$).
- Seja ainda v_i , que assume o valor 1 se $y_i = 0$ e 0 caso contrário. Então, a verossimilhança (aumentada) e a log-verossimilhança (aumentada) são dadas, respectivamente, por

$$L(\theta, \gamma) = \prod_{i=1}^n \left\{ [\pi_i + (1 - \pi_i) f_{Z_i}(0; \theta)]^{v_i} [(1 - \pi_i) f_{Z_i}(y_i; \theta)]^{1-v_i} \right\}$$
$$l(\theta, \gamma) = \sum_{i=1}^n \left\{ v_i \ln [\pi_i + (1 - \pi_i) f_{Z_i}(0; \theta)] + (1 - v_i) [\ln(1 - \pi_i) + \ln f_{Z_i}(y_i; \theta)] \right\}. \quad (2)$$

log-verossimilhança

- Exemplo: Poisson com ligação log sem modelar π

$$l(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \left\{ v_i \ln \left[\pi + (1 - \pi) e^{-\mathbf{x}'_i \boldsymbol{\beta}} \right] + (1 - v_i) \left[\ln(1 - \pi) - e^{\mathbf{x}'_i \boldsymbol{\beta}} + y_i \mathbf{x}'_i \boldsymbol{\beta} - \ln(y_i!) \right] \right\}.$$

- Exemplo: Poisson com ligação log modelando π (logito)

$$l(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \left\{ v_i \ln \left[\frac{e^{\mathbf{w}'_i \boldsymbol{\gamma}} + e^{-\mathbf{x}'_i \boldsymbol{\beta}}}{1 + e^{\mathbf{w}'_i \boldsymbol{\gamma}}} \right] + (1 - v_i) \left[-\ln \left(1 + e^{\mathbf{w}'_i \boldsymbol{\gamma}} \right) - e^{\mathbf{x}'_i \boldsymbol{\beta}} + y_i \mathbf{x}'_i \boldsymbol{\beta} - \ln(y_i!) \right] \right\}$$

log-verossimilhança (completada/aumentada)

- Consiste em utilizar as observações (y_i) e as variáveis lantes (u_i) , ou seja, a distribuição conjunta de $(Y_i, U_i)'$, $i=1,2,\dots,n$.
- De (1) temos que a verossimilhança aumentada é dada por

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \prod_{i=1}^n \left\{ \pi_i^{u_i} [f_{Z_i}(y_i; \boldsymbol{\theta})(1 - \pi_i)]^{1-u_i} \right\}.$$

- Consequentemente, a log-verossimilhança aumentada é dada por:

$$l(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \left\{ u_i \ln \pi_i + (1 - u_i) [\ln(1 - \pi_i) + \ln f_{Z_i}(y_i; \boldsymbol{\theta})] \right\}. \quad (3)$$

Estimação

- Note que a expressão (3) é mais tratável do que a expressão (2).
- Contudo, não observamos as variáveis U_i . Neste caso, podemos usar, por exemplo, o algoritmo EM para obtermos as estimativas de máxima verossimilhança.
- A obtenção das emv através da maximização direta da logverossimilhança (3) é computacionalmente mais complicada.
- Adicionalmente, podemos utilizar a informação de Fisher (observada ou esperada) para calcularmos os erros-padrão

Cont.

- A obtenção das emv através via algoritmo EM é mais simples (2), mas requer o cálculo da esperança condicional:

$$\mathcal{E}(U_i|y_i, \theta, \gamma).$$

- Além disso, para calcularmos os erros-padrão devemos utilizar a abordagem de [Louis \(1982\)](#) ou a de [Meilijson \(1989\)](#).

Análise residual

- Resíduo padronizado: $R_i = \frac{Y_i - \hat{\mu}_i}{\sqrt{\mathcal{V}(Y_i)}}$.
- Resíduo quantílico aleatorizado (**aqui**).
- Construir os cinco gráficos usuais: resíduo x índice, resíduo x valores preditos, histograma, valores preditos x valores observados e gráfico de envelopes.
- Mesmo sob o bom ajuste do modelo não, necessariamente, espera-se que os resíduos apresentem distribuição normal (mesmo que aproximadamente).
- Voltaremos agora ao Exemplo 16.

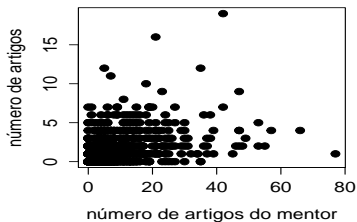
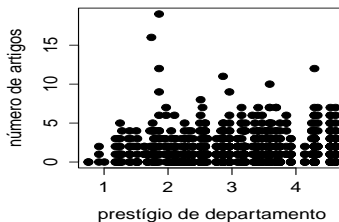
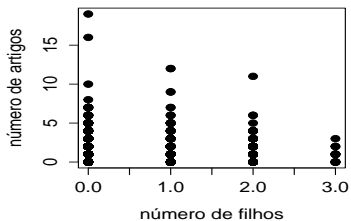
Análise de dados

- Para o ajuste dos modelos inflacionados podemos usar a função *zeroinfl* do pacote *pscl*.
- Voltaremos agora ao Exemplo 16.

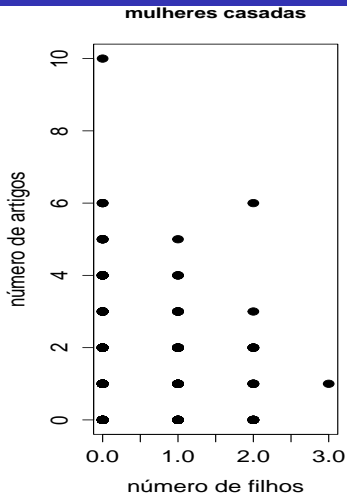
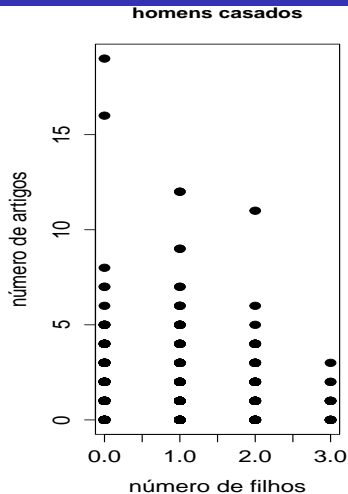
Medidas resumo

MR	masculino		feminino	
	solteiro	casado	solteiro	casado
Média	1,95	1,86	1,39	1,54
DP	2,01	2,23	1,51	1,59
Var	4,05	4,97	2,28	2,52
CV(%)	103,38	119,58	108,80	102,88
CA	1,15	3,00	1,32	1,41
Curtose	3,51	18,04	4,50	6,14
Min.	0	0	0	0
Max.	7	19	7	10

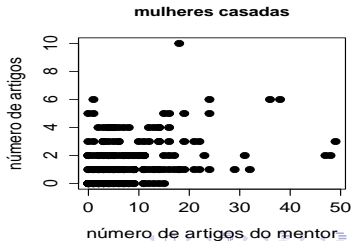
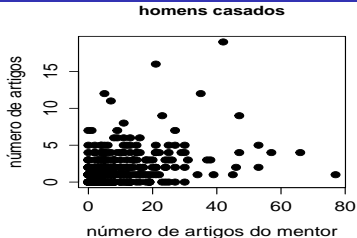
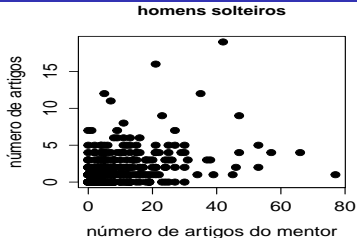
Gráficos de dispersão



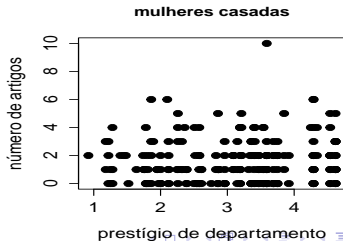
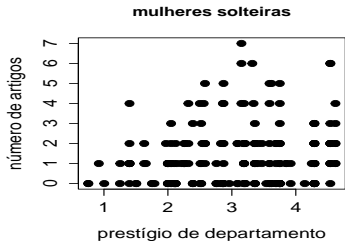
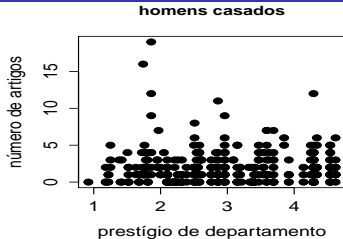
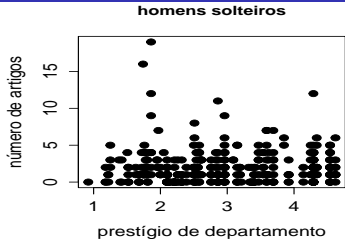
Gráficos de dispersão (número de filhos)



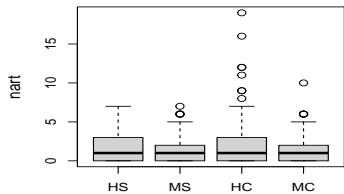
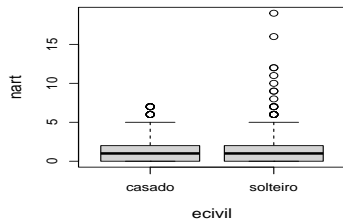
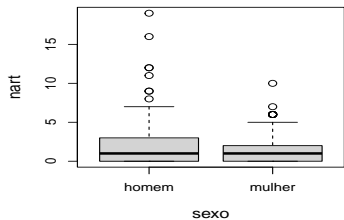
Gráficos de dispersão (número de artigos do orientador)



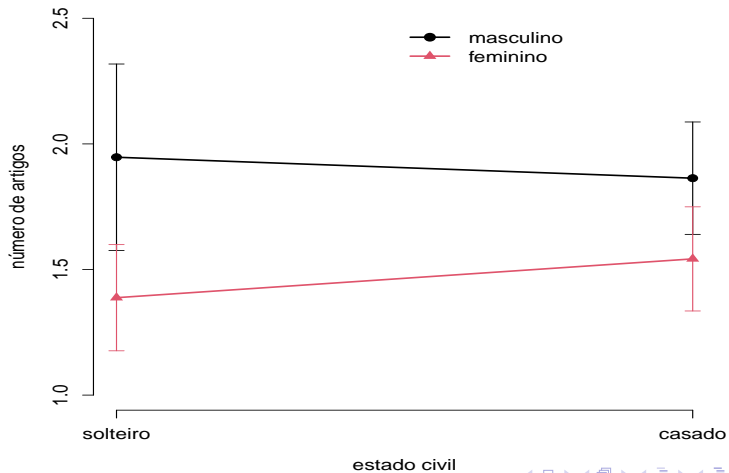
Gráficos de dispersão (prestígio do departamento)



Box-plots



Gráficos de perfis



Voltando ao Exemplo 16

Modelos

$$M1 : Y_{ijk} \stackrel{ind.}{\sim} \text{Poisson}(\mu_{ijk}),$$

$$M2 : Y_{ijk} \stackrel{ind.}{\sim} \text{BN}(\mu_{ijk}, \phi),$$

$$M3 : Y_{ijk} \stackrel{ind.}{\sim} \text{IZ}\{\pi, \text{Poisson}(\mu_{ijk})\},$$

$$M4 : Y_{ijk} \stackrel{ind.}{\sim} \text{IN}\{\pi, \text{BN}(\mu_{ijk}, \phi)\},$$

em que sexo ($i = 1$ (masculino), 2 (feminino)), estado civil ($j=1$ (solteiro), 2 (casado)), $k = 1, 2, \dots, n_{ij}$ (aluno).

Voltando ao Exemplo 16

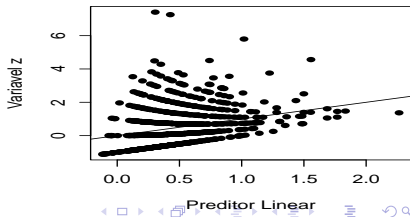
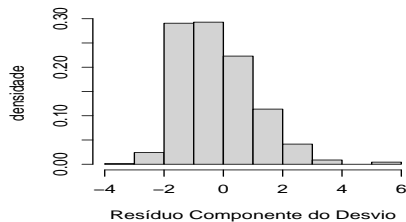
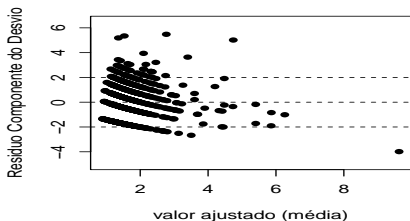
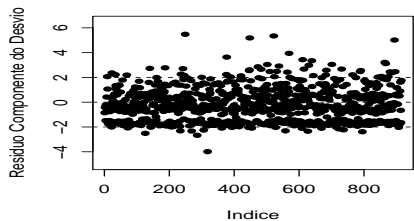
- Em todos os modelos

$$\ln(\mu_{ijk}) = \alpha + \beta_i + \gamma_j + \delta_1(x_{1ijk} - \bar{x}_1) + \delta_2(x_{2ijk} - \bar{x}_2) + \delta_3(x_{3ijk} - \bar{x}_3),$$

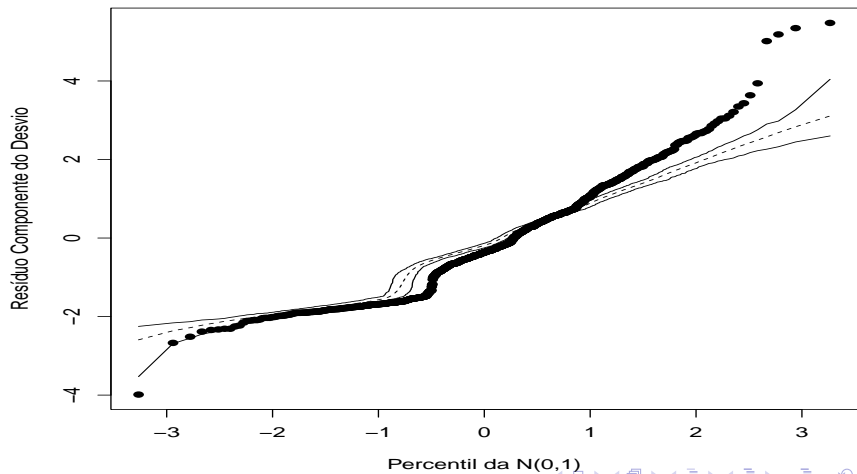
em que $\bar{x}_r = \frac{1}{n_{ij}} \sum_{j=1}^2 \sum_{i=1}^2 \sum_{k=1}^n x_{rijk}$, $r = 1, 2, 3$ e $\beta_1 = \gamma_1 = 0$, kid5
(x_{1sjk}), phd (x_{2sjk}), ment (x_{3sjk}),

- Nos modelos (3) e (4) $\ln\left(\frac{\pi}{1-\pi}\right) = \theta$.
- Exercício: interpretar os parâmetros adequadamente.

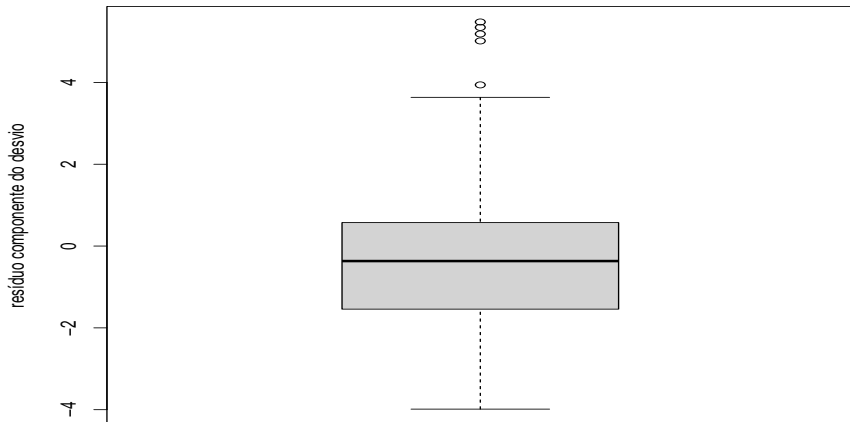
Gráficos de diagnóstico (M1)



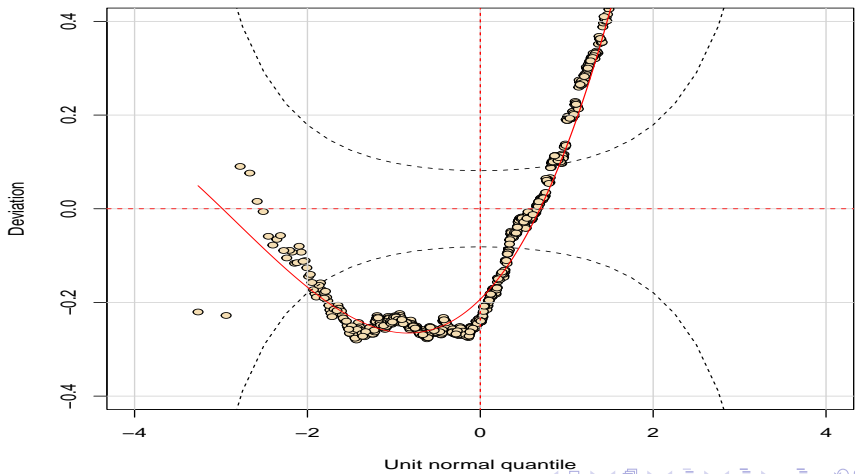
Gráficos de envelopes (M1)



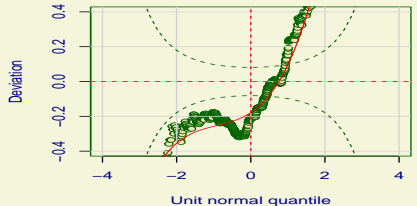
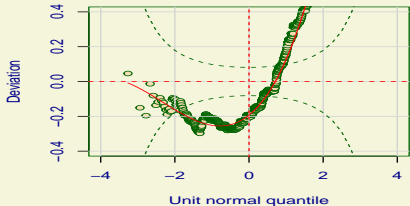
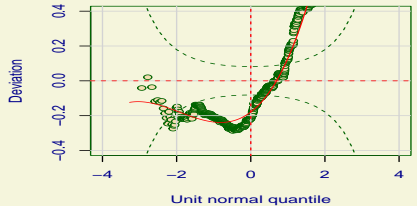
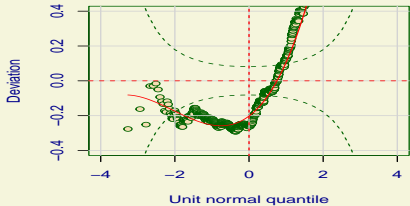
Boxplot do RCD (M1)



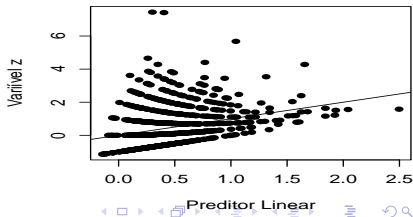
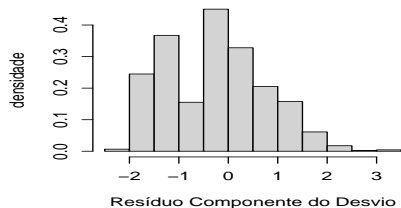
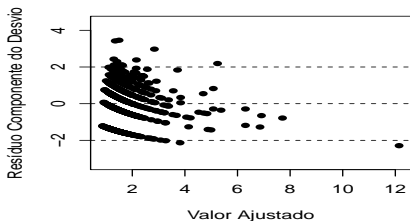
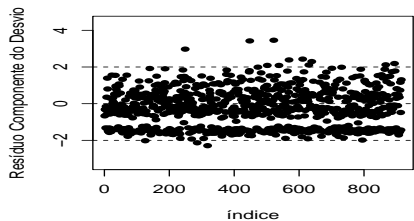
Wormplot (M1)



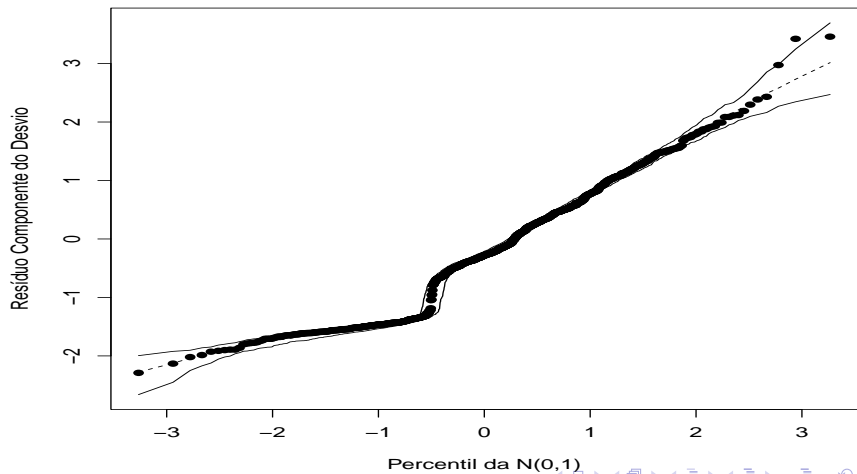
Wormplots (M1)



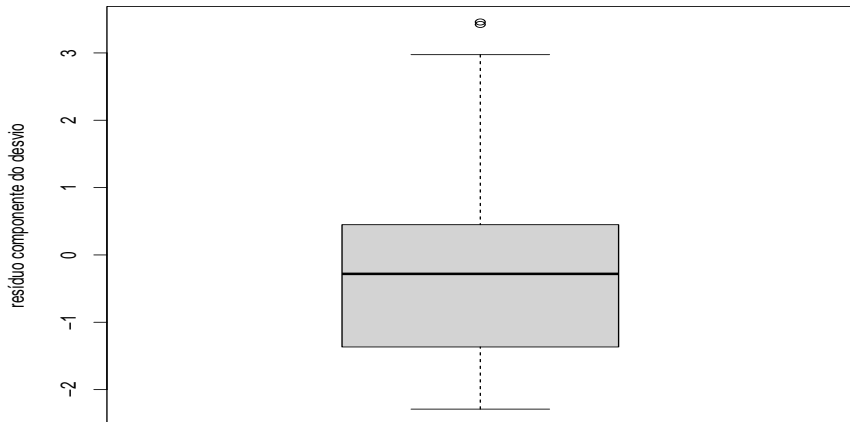
Gráficos de diagnóstico (M2)



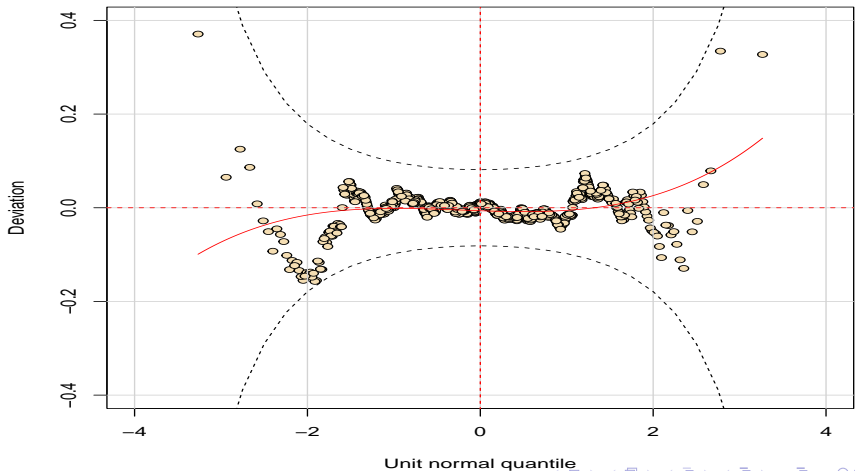
Gráficos de envelopes (M2)



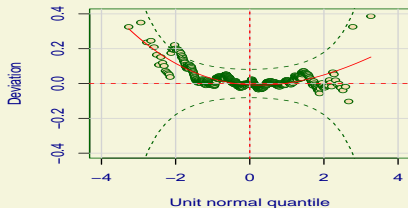
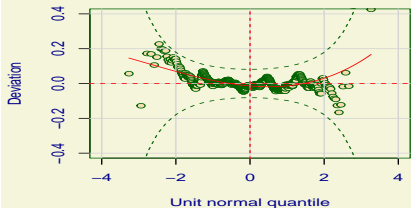
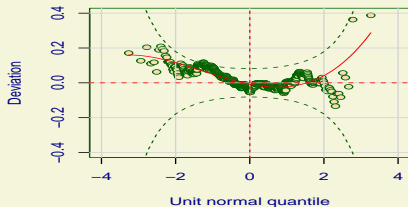
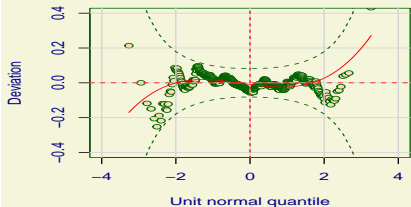
Boxplot do RCD (M2)



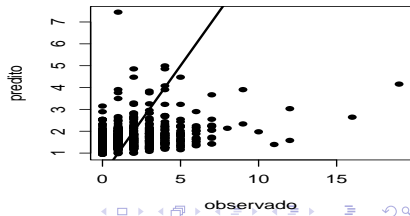
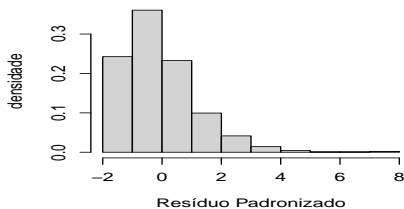
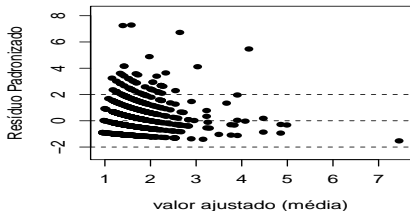
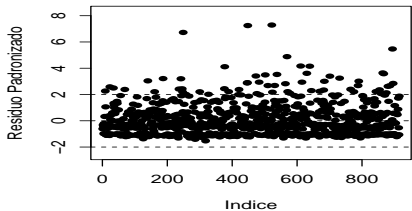
Wormplot (M2)



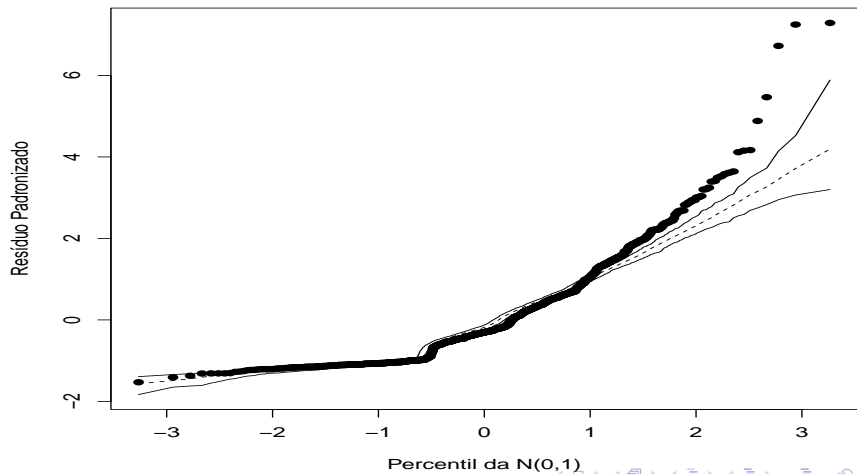
Wormplots (M2)



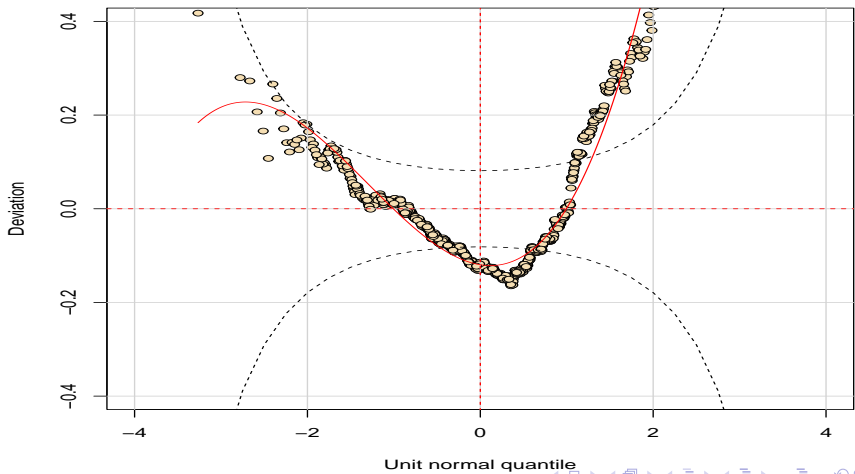
Gráficos de diagnóstico (M3)



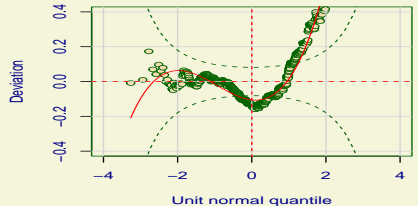
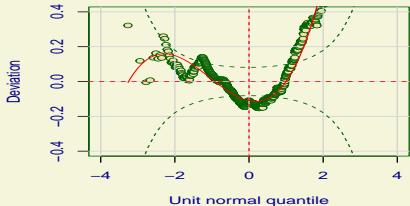
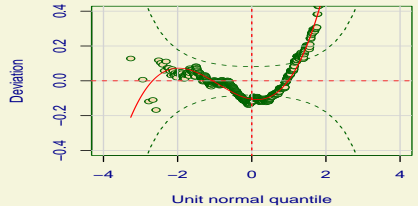
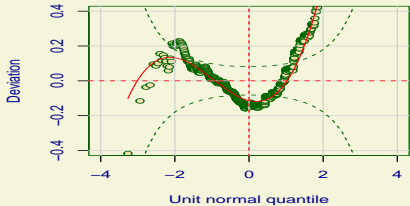
Gráficos de envelopes (M3)



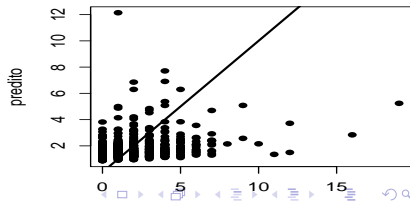
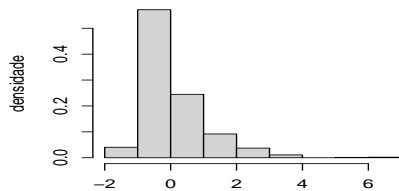
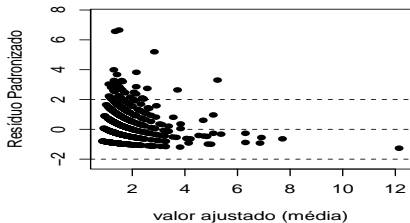
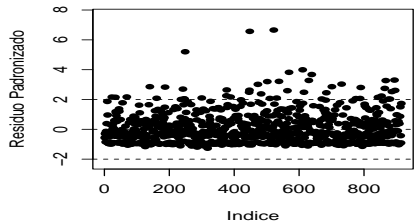
Wormplot (M3)



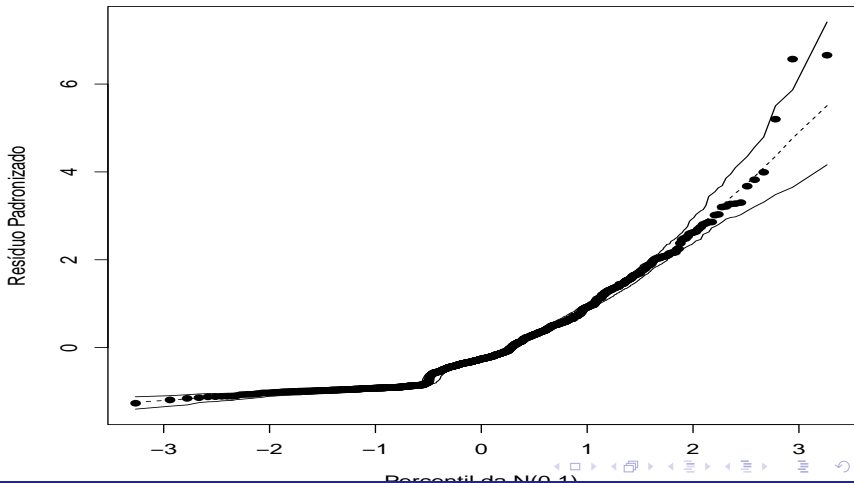
Wormplots (M3)



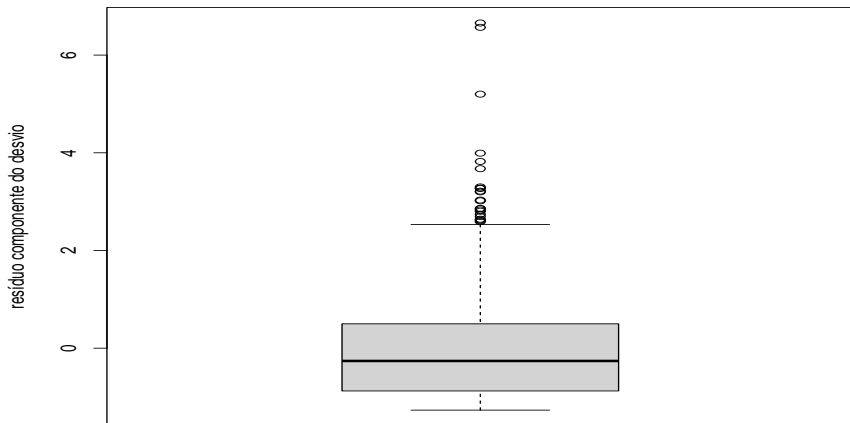
Gráficos de diagnóstico (M4)



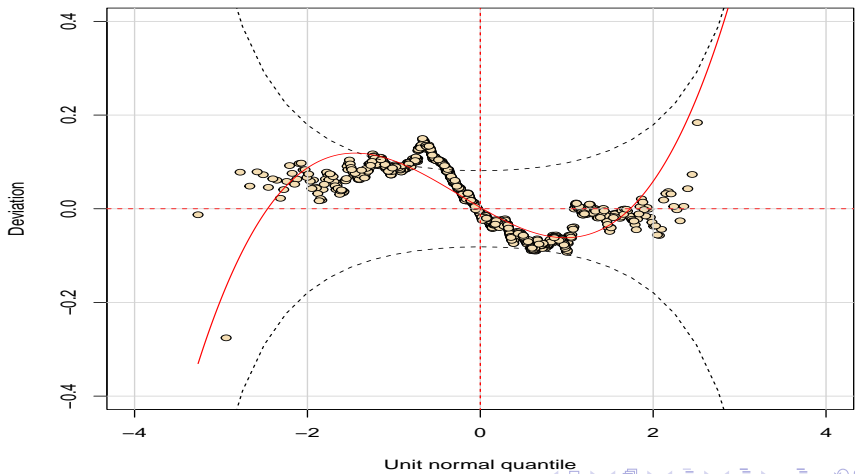
Gráficos de envelopes (M4)



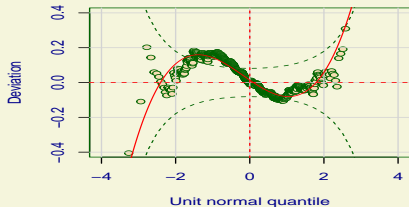
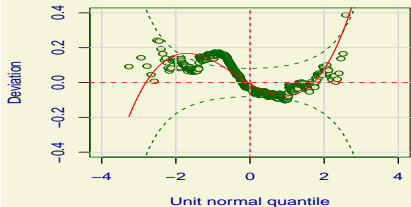
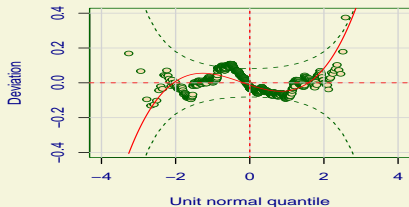
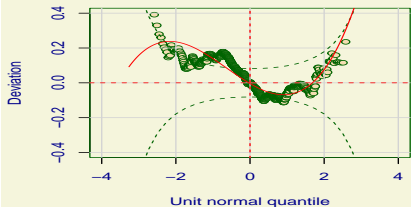
Boxplot do RCD (M4)



Wormplot (M4)



Wormplots (M4)



Estatísticas de comparação de modelos

Modelo	AIC	BIC	AIC _c	SABIC	HQCIC	CAIC
Poisson	3314,11	3343,03	3314,21	3323,97	3325,15	3349,03
BN	3135,92	3169,65	3136,04	3147,42	3148,79	3176,65
PIZ	3255,57	3289,30	3255,69	3267,07	3268,44	3296,30
BNIZ	3139,92	3183,29	3140,12	3154,70	3156,47	3192,29

Estimativas (M1)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
α	0,48	0,06	[0,37;0,59]	8,36	< 0,0001
β_2	-0,22	0,05	[-0,33;-0,12]	-4,11	< 0,0001
γ_2	0,16	0,06	[0,03;0,28]	2,53	0,0114
δ_1	-0,18	0,04	[-0,26;-0,11]	-4,61	< 0,0001
δ_2	0,03	< 0,01	[0,02;0,03]	12,73	< 0,0001
δ_3	0,01	0,03	[-0,04;0,06]	0,49	0,6271

Estimativas (M2)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
α	0,47	0,08	[0,32;0,62]	6,22	< 0,0001
β_2	-0,22	0,07	[-0,36;-0,07]	-2,98	0,0029
γ_2	0,15	0,08	[-0,01;0,31]	1,83	0,0668
δ_1	-0,18	0,05	[-0,28;-0,07]	-3,34	0,0008
δ_2	0,03	< 0,00	[0,02;0,04]	9,05	< 0,0001
δ_3	0,02	0,04	[-0,06;0,09]	0,43	0,6703
ϕ	2,26	0,27	[1,73;2,80]	-	-

Estimativas (M3)

Parâmetro	Est,	EP	IC(95%)	Estat, Z	p-valor
α	0,67	0,06	[0,54;0,79]	10,41	< 0,0001
β_2	-0,23	0,06	[-0,35;-0,12]	-3,95	0,0001
γ_2	0,13	0,07	[0,00;0,26]	2,00	0,0460
δ_1	-0,17	0,04	[-0,26;-0,09]	-3,94	0,0001
δ_2	0,02	< 0,00	[0,02;0,03]	9,97	< 0,0001
δ_3	< 0,01	0,03	[-0,05;0,06]	0,09	0,9294
π	0,16	0,06	[0,05;0,27]		

Estimativas (M4)

Parâmetro	Est,	EP	IC(95%)	Estat, Z	p-valor
α	0,47	0,08	[0,32 ; 0,62]	6,24	< 0,0001
β_2	-0,22	0,07	[-0,36 ; -0,07]	-2,98	0,0029
γ_2	0,15	0,08	[-0,01 ; 0,31]	1,83	0,0668
δ_1	-0,18	0,05	[-0,28 ; -0,07]	-3,32	0,0009
δ_1	0,03	< 0,00	[0,02 ; 0,04]	8,38	< 0,0001
δ_1	0,02	0,04	[-0,06 ; 0,09]	0,42	0,6718
ϕ	2,26	0,27	[1,73 ; 2,80]		
π	< 0,01	0,14	[-0,28 ; 0,28]		

Comentários

- As vezes superdispersão e excesso de zeros produzem comportamentos semelhantes nos dados.
- Isso se reflete, em algumas vezes, na semelhança dos ajustes de modelos que contemplam pelo menos uma dessas características.
- Os modelos M2 (superdispersão) e M4 (superdispersão + excesso de zeros) apresentaram bons e equivalentes ajustes, com uma ligeira vantagem par ao modelo M2.
- Neste caso, devido aos resultados (veja também a estimativa de baixa magnitude de π para o modelo M4), considerar a superdispersão (sem excesso de zeros) foi suficiente.

Comentários

- Exercício: Pesquisar sobre análise de influência para modelos zero-inflacionados.
- Exercício: A partir do resultados do ajuste do M3, verificar qual(is) variável(is) pode(m) ser eliminada(s) e, a partir desse modelo reduzido, considerar possíveis interações até chegar a um modelo final.