

# 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 (Ph.D.) (resposta).
  - fem: gênero 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 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 de zeros.

# Modelos probabilísticos

- Seja  $Y$  a vad (variável aleatória discreta) que representa a contagem de interesse.
- $P(Y = y) = g_Y(y) = [\pi + (1 - \pi)f_Z(0)] \mathbb{1}_{\{0\}}(y) + (1 - \pi)f_Z(y)\mathbb{1}_{\{1,2,\dots\}}$ , em que  $f_Z(\cdot)$  representa a função de probabilidade de uma vad discreta de interesse (Poisson, geométrica, binomial negativa, etc).
- Notação:  $Y \sim IZ\{\pi, f_Z(y)\}$  em que  $f_Z(\cdot)$  representa o modelo probabilístico original. Por exemplo  $Y \sim IZ\{\pi, \text{Poisson}(\mu)\}$ .
- Exemplo:

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

# Modelos probabilísticos

- Neste caso  $\pi$  representa a probabilidade de inflacionamento.
- Note que

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

# 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 & binomial-negativa:  $\mathcal{E}(Y) = (1 - \pi)\mu$ .

# Modelos probabilísticos

## ■ Variância

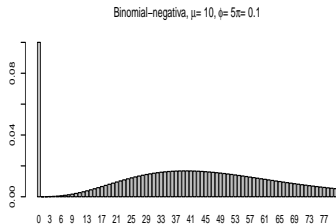
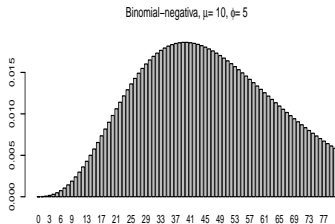
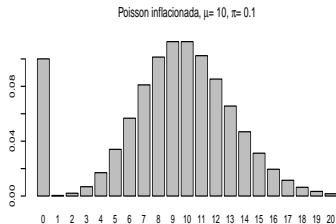
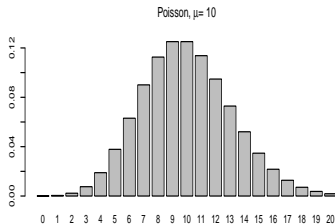
$$\begin{aligned}\mathcal{V}(Y) &= \mathcal{E}(Y^2) - \mathcal{E}^2(Y) = \sum_{y=0}^{\infty} y^2 g(y) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) \sum_{y=1}^{\infty} y^2 g(y) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) \mathcal{E}(Z^2) - (1 - \pi)^2 \mathcal{E}^2(Z) \\ &= (1 - \pi) [\mathcal{E}(Z^2)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

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

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

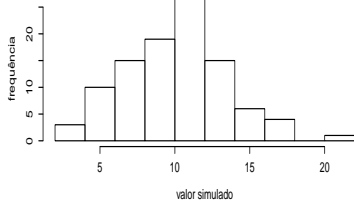


# Exemplos de funções de probabilidade

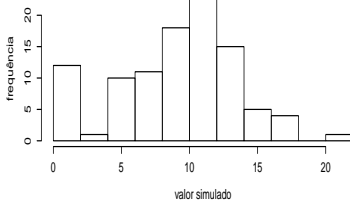


# Histograma de valores simulados

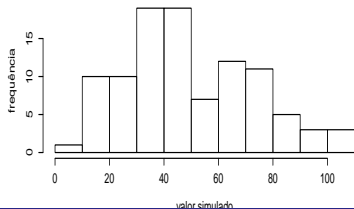
Poisson,  $\mu = 10$



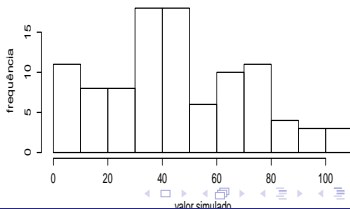
Poisson inflacionada,  $\mu = 10, \pi = 0.1$



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  $\pi$ .
- Por exemplo, se  $Z \sim \text{Poisson}(\mu_i)$  então podemos considerar  $\ln(\mu_i) = \mathbf{x}'_i \boldsymbol{\beta}$  e  $\text{logito}(\pi_i) = \ln\left(\frac{\pi_i}{1 - \pi_i}\right) = \mathbf{w}'_i \boldsymbol{\gamma}$ .

# Inferência via MV

- 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.
- Estimação paramétrica:
  - Otimização direta da log-verossimilhança.
  - Algoritmo EM, otimizando-se a log-verossimilhança aumentada definida pela distribuição conjunta de  $(Y, U)$ .

## log-verossimilhança

- Sejam  $\theta$  os parâmetros da contagem original ( $Z$ ) e  $\gamma$  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

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

# log-verossimilhança

- Exemplo: Poisson com ligação log sem modelar  $\pi$

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

- Exemplo: Poisson com ligação log modelando  $\pi$  (logito)

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

## 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

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

- Log-verossimilhança aumentada

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



- Note que a expressão (3) é mais tratável do que a expressão (2). Contudo, não observamos as variáveis  $u_j$ . Neste caso, devemos 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. Podemos utilizar a informação de Fisher (observada ou esperada) para calcularmos os erros-padrão
- 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).

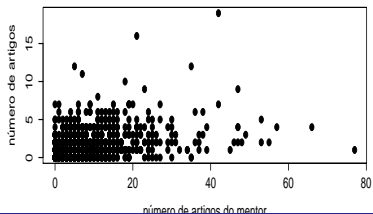
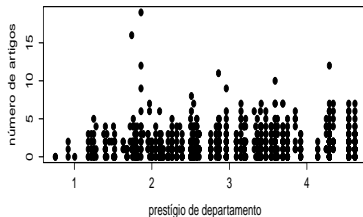
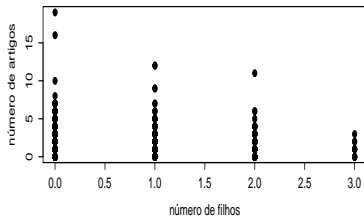
# Análise residual

- Resíduo padronizado:  $R_i = \frac{Y_i - \hat{\mu}_i}{\sqrt{V(Y_i)}}$ .
- 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.

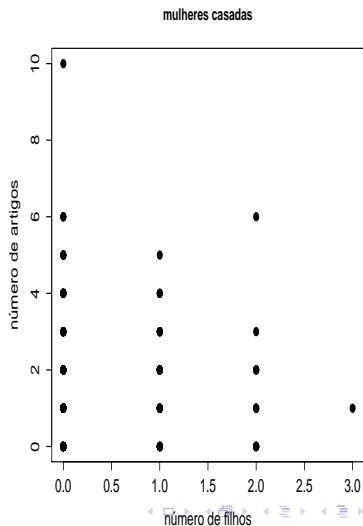
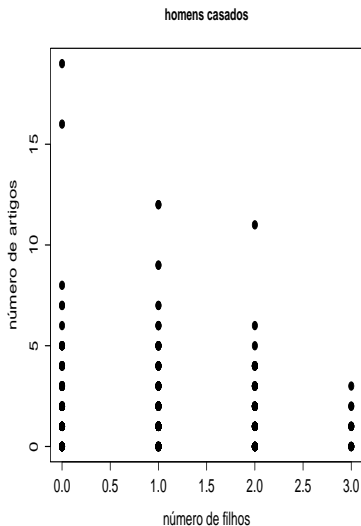
# Medidas resumo

Gênero	Estado civil	Média	DP	Var.	CV(%)	Min.	Max.	n
masculino	solteiro	1,95	2,01	4,05	103,38	0	7	113
	casado	1,86	2,23	4,97	119,58	0	19	381
feminino	solteiro	1,39	1,51	2,28	108,80	0	7	196
	casado	1,54	1,59	2,52	102,88	0	10	225

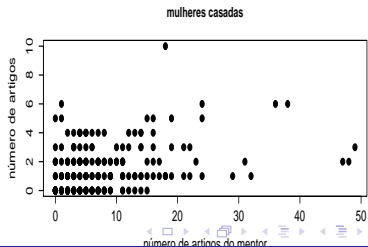
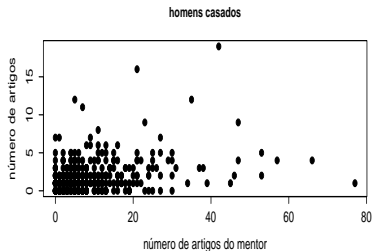
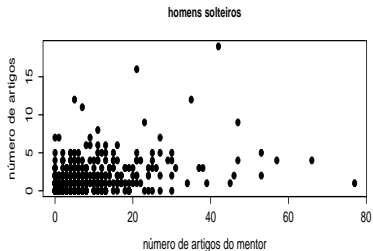
# Gráficos de dispersão



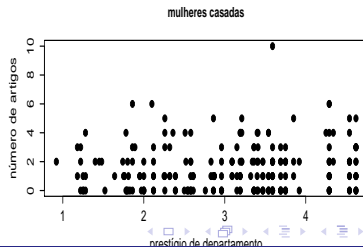
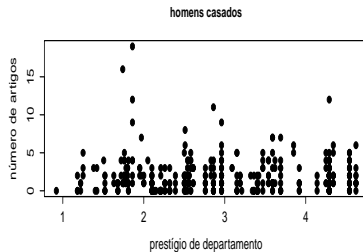
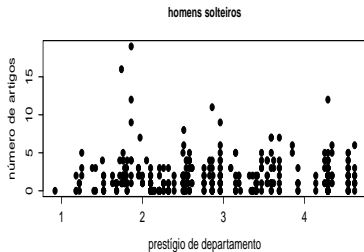
# Gráficos de dispersão (número de filhos)



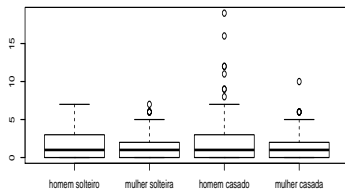
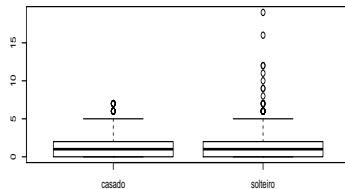
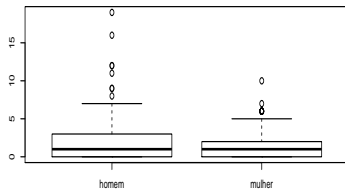
# Gráficos de dispersão (número de artigos do orientador)



# Gráficos de dispersão (prestígio do departamento)

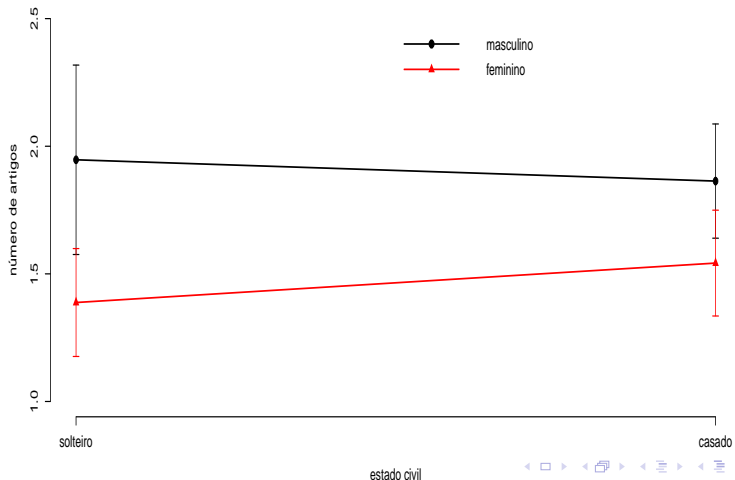


# Box-plots





# Gráficos de perfis



## Voltando ao Exemplo 16

### Modelos

$Y_{ijk} \stackrel{ind.}{\sim} \text{Poisson}(\mu_{ijk})(1); Y_{ijk} \stackrel{ind.}{\sim} \text{BN}(\mu_{ijk}, \phi)(2);$

$Y_{ijk} \stackrel{ind.}{\sim} \text{IZ}\{\pi, \text{Poisson}(\mu_{ijk})\}(3); Y_{ijk} \stackrel{ind.}{\sim} \text{IN}\{\pi, \text{BN}(\mu_{ijk}, \phi)\}(4)$

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

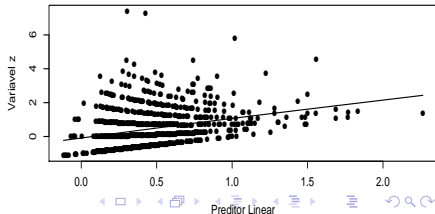
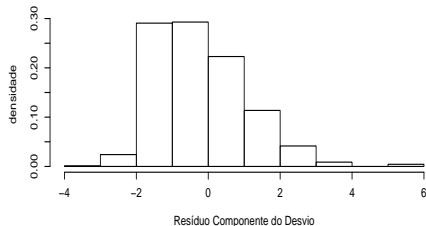
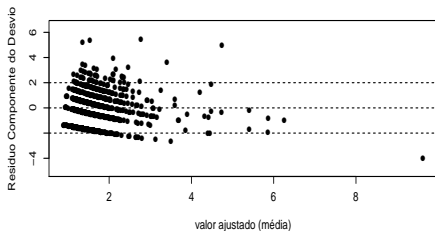
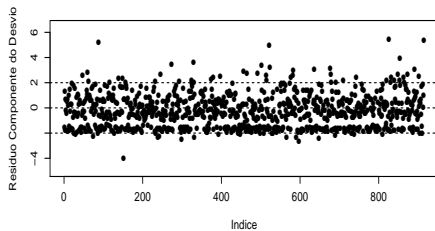
- Em todos os modelos

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

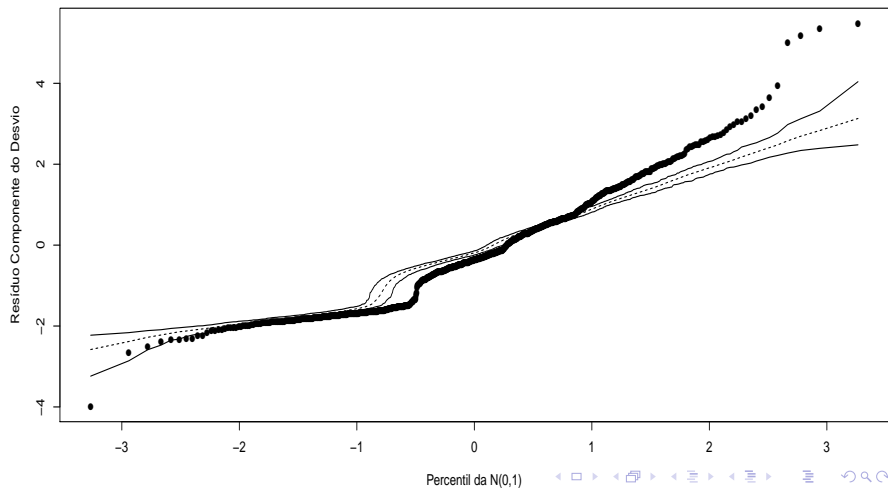
em que  $x_r = \frac{1}{n_{ij}} \sum_{k=1}^{n_{ij}} x_{rijk}$ ,  $r = 1, 2, 3$  e  $\beta_1 = \gamma_1 = 0$ .

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

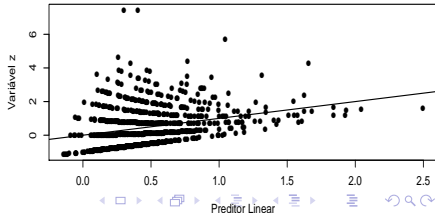
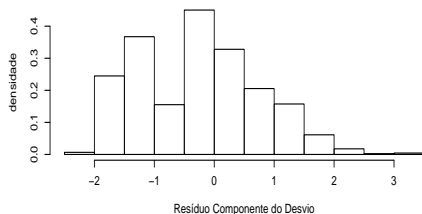
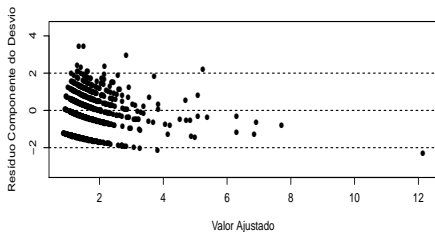
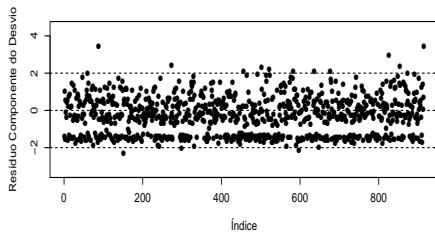
# Gráficos de diagnóstico (M1)



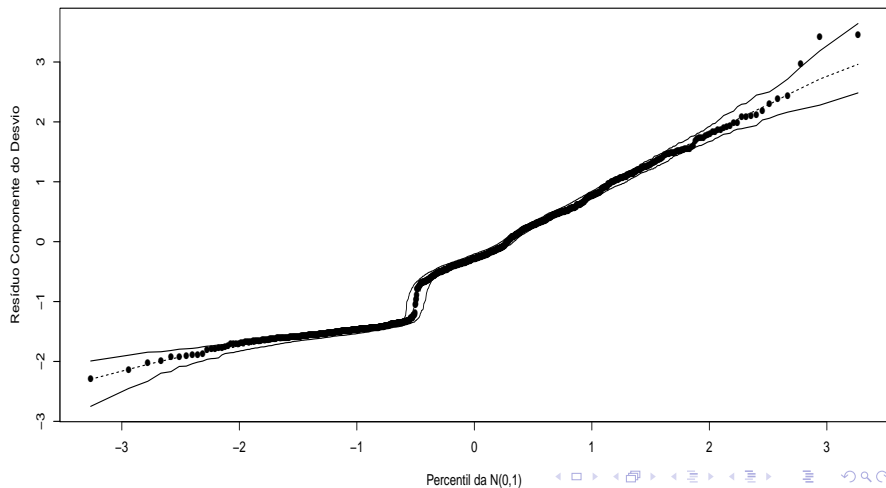
# Gráficos de envelopes (M1)



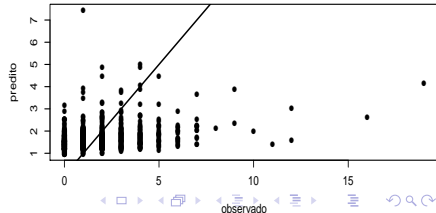
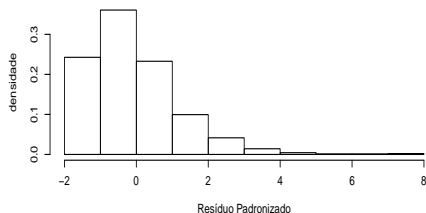
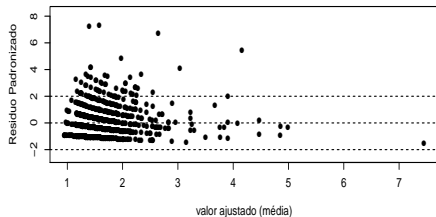
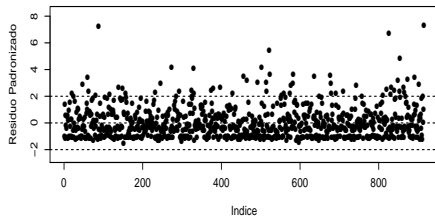
# Gráficos de diagnóstico (M2)



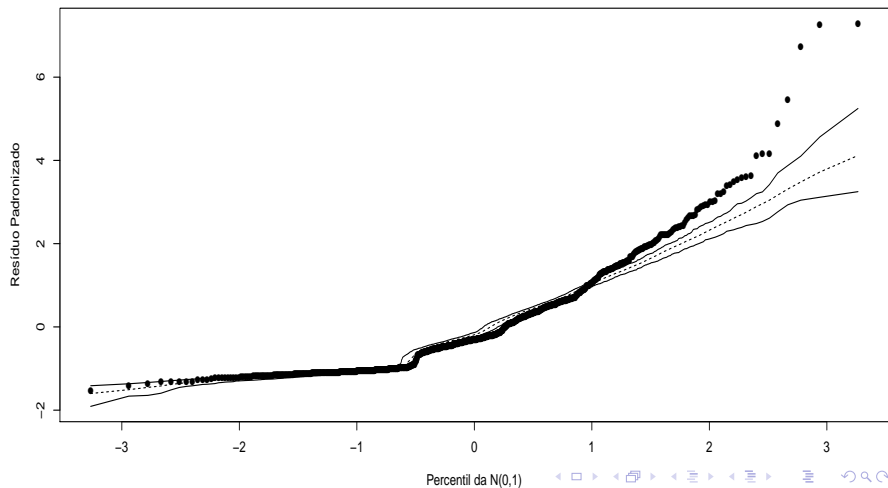
# Gráficos de envelopes (M2)



# Gráficos de diagnóstico (M3)

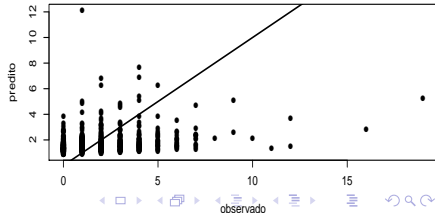
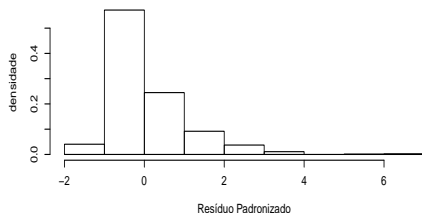
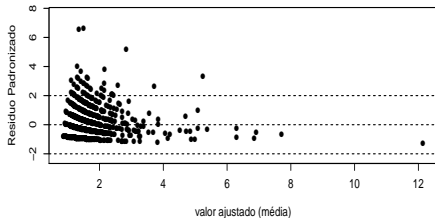
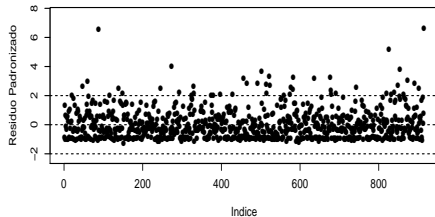


# Gráficos de envelopes (M3)

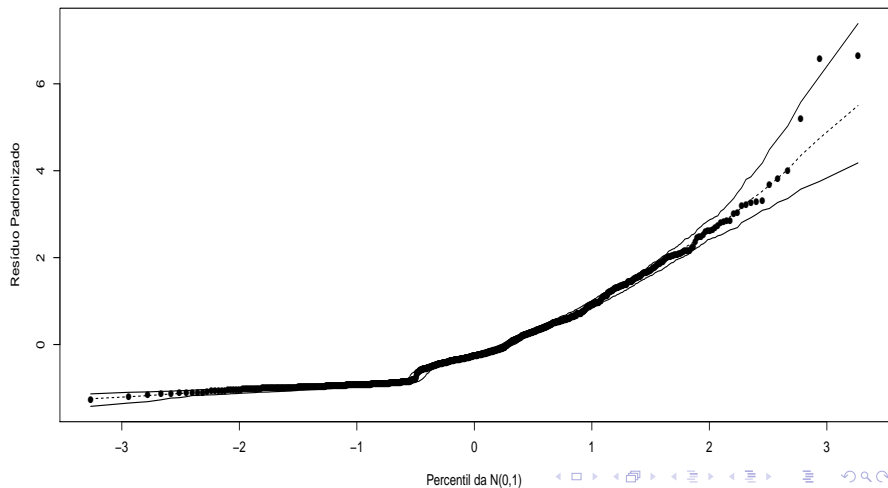




# Gráficos de diagnóstico (M4)



# Gráficos de envelopes (M4)



# Estatísticas de comparação de modelos

Modelo	AIC	BIC	desvio	p-valor (desvio)
1	3314,11	3343,03	1634,37	< 0,0001
2	3135,92	3169,65	1004,28	0,0148
3	3255,57	3289,30	-	-
4	3137,92	3176,47	-	-

# Estimativas (M1)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
$\alpha$	0,48	0,06	[0,37 ; 0,59]	8,36	< 0,0001
$\beta_2$	-0,18	0,04	[-0,26 ; -0,11]	-4,61	< 0,0001
$\gamma_2$	0,026	0,002	[0,022 ; 0,029]	12,73	< 0,0001
$\delta_1$	0,01	0,03	[-0,04 ; 0,06]	0,49	0,6271
$\delta_2$	-0,22	0,05	[-0,33 ; -0,12]	-4,11	< 0,0001
$\delta_3$	0,16	0,06	[0,03 ; 0,28]	2,53	0,0114

## Estimativas (M2)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
$\alpha$	0,47	0,08	[0,32 ; 0,62]	6,22	< 0.0001
$\beta_2$	-0,18	0,05	[-0,28 ; -0,07]	-3,34	0.0008
$\gamma_2$	0,03	< 0,01	[0,02 ; 0,04]	9,05	< 0.0001
$\delta_1$	0,02	0,04	[-0,06 ; 0,09]	0,43	0,6703
$\delta_2$	-0,22	0,07	[-0,36 ; -0,07]	-2,98	0,0029
$\delta_3$	0,15	0,08	[-0,01 ; 0,31]	1,83	0,0668
$\phi$	2,26	0,27	[1,73 ; 2,80]	-	-

## Estimativas (M3)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
$\alpha$	0,67	0,06	[ 0,54 ; 0,79]	10,41	< 0,0001
$\beta_2$	-0,17	0,04	[-0,26 ; -0,09]	-3,94	0,0001
$\gamma_2$	0,02	< 0,01	[0,02 ; 0,03]	9,97	0,0001
$\delta_1$	0,003	0,029	[ -0,053 ; 0,058]	0,09	0,9294
$\delta_2$	-0,23	0,06	[-0,35 ; -0,12]	-3,95	0,0001
$\delta_3$	0,13	0,07	[ <0,01 ; 0,26]	2,00	0,0460
$\pi$	0,16	0,07	[0,02 ; 0,30]	0,00	< 0,0001

# Estimativas (M4)

Parâmetro	Est.	EP	IC(95%)	Estat. Z	p-valor
$\alpha$	0,47	0,08	[0,32 ; 0,62]	6,24	<0,0001
$\beta_2$	-0,18	0,05	[-0,28 ; -0,07]	-3,33	0,0009
$\gamma_2$	0,03	< 0,01	[0,02 ; 0,04]	8,38	<0,0001
$\delta_1$	0,02	0,04	[-0,06 ; 0,09]	0,42	0,6716
$\delta_2$	-0,22	0,07	[-0,36 ; -0,07]	-2,98	0,0029
$\delta_3$	0,15	0,08	[-0,01 ; 0,31]	1,83	0,0677
$\phi$	2,26	0,37	[1,54 ; 2,99]	0,00	<0,0001
$\pi$	$4,11 \times 10^{-6}$	35,17	-	-	-

# 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  $\pi$  para o modelo M4), considerar a superdispersão (sem excesso de zeros) foi suficiente.