



**Universidade Estadual de Campinas
Instituto de Matemática, Estatística e Computação Científica
Departamento de Estatística**

Relatório - Parte II Trabalho Final de ME613

**Eliane Ramos de Siqueira RA:155233
Guilherme Pazian RA:160323
Henrique Capatto RA:146406
Murilo Salgado Razoli RA:150987**

Professor: Caio Lucidius Naberezny Azevedo

Campinas-SP, 08 de Dezembro de 2016

1. Introdução

Este presente trabalho é parte do escopo da disciplina ME613 - Regressão Linear e visa a aplicação dos conhecimentos adquiridos em sala de aula. Grande parte do trabalho foi baseado e adaptado a partir dos materiais disponibilizados na página do curso: http://www.ime.unicamp.br/~cnaber/Material_ME613_2S_2016.htm.

O conjunto de dados a ser analisado está disponível em na página do curso sob o nome “reg3.dat”. Neste são descritas as seguintes variáveis referentes a 50 estados norte-americanos: (i) **estado** (nome do estado), (ii) **pop** (população estimada em julho de 1975), (iii) **percap** (renda percapita em 1974 em USD), (iv) **analf** (proporção de analfabetos em 1970), (v) **expvida** (expectativa de vida em anos 1969-70), (vi) **crime** (taxa de criminalidade por 100000 habitantes 1976), (vii) **estud** (porcentagem de estudantes que concluem o segundo grau 1970), (viii) **ndias** (número de dias do ano com temperatura abaixo de zero grau Celsius na cidade mais importante do estado) e (ix) **area** (área do estado em milhas quadradas). O objetivo do estudo é tentar explicar a variável **expvida** usando um modelo de regressão normal linear dadas as variáveis explicativas **percap**, **analf**, **crime**, **estud**, **ndias** e **dens**, em que **dens**=pop/area (densidade da população estimada em julho de 1975 por área do estado em milhas quadradas).

Todos os modelos neste trabalho foram ajustados via metodologia de mínimos quadrados ordinários, veja Azevedo (2016), e todas suas respectivas análises residuais foram realizadas, conforme Paula (2013). A menos que seja citado o contrário, todas as variáveis que constam na base de dados serão referidas por sua descrição completa ou pelo nome que consta no banco de dados (estes foram supracitados em negrito). Denotaremos também:

Y_i - i-ésima observação da variável expvida

x_{1i} - i-esima observação da variável percap

x_{2i} - i-esima observação da variável analf

x_{3i} - i-esima observação da variável crime

x_{4i} - i-esima observação da variável estud

x_{5i} - i-esima observação da variável ndias

x_{6i} - i-esima observação da variável dens

2. Análise descritiva

Tabela 1: Estatísticas Descritivas

variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
percap	3098.0000	3993.0000	4519.000	4436.000	4814.0000	6315.000
analf	0.5000	0.6250	0.950	1.170	1.5750	2.800
crime	1.4000	4.3500	6.850	7.378	10.6800	15.100
estud	37.8000	48.0500	53.250	53.110	59.1500	67.300
ndias	0.0000	66.2500	114.500	104.500	139.8000	188.000
dens	0.0006	0.0277	0.069	0.188	0.1443	2.684
expvida	67.9600	70.1200	70.680	70.880	71.8900	73.600

Dado a natureza quantitativa dos dados, é de grande interesse que identifiquemos as relações entre as covariáveis explicativas e a variável resposta, a fim de tornar essas relações visuais, foi construído gráficos de dispersão entre a variável resposta e entre todas as covariáveis de interesse. Estes estão representados na figura XX.

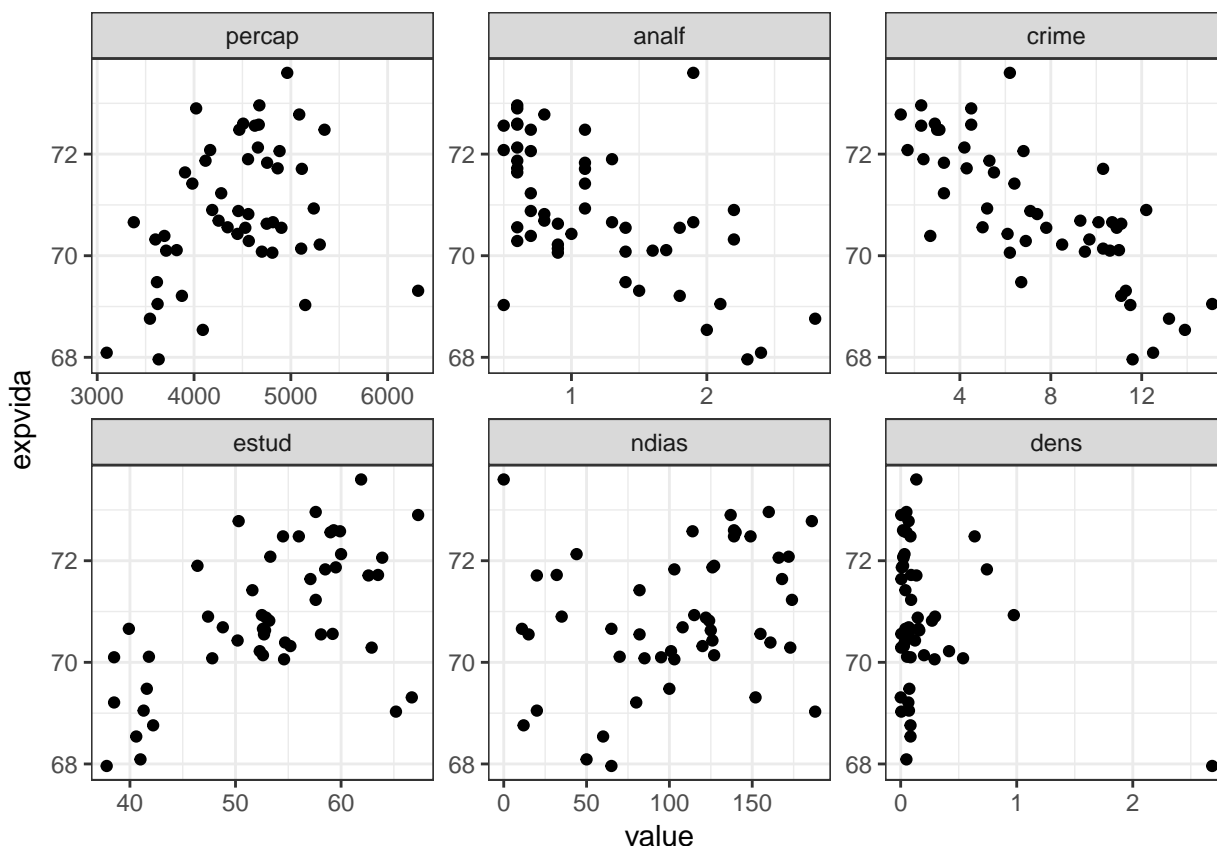


Figura 1: Gráficos de dispersão

A partir da Figura 1, podemos identificar visualmente a relação supracitada, a qual identificamos que todas as covariáveis parecem ter uma relação linear significativa, as quais existem relações lineares negativas entre a variável resposta e as covariáveis “analf” e “crime”, assim como relações lineares negativas entre a variável resposta “expvida” e as covariáveis “percap”, “estud”, “ndias” e “dens” (a menos a um ponto).

3. Análise Inferencial

Dada a constatação visual de relações lineares entre a variável resposta “expvida” e todas as covariáveis explicativas, é razoável iniciar um modelo que contemple todas estas covariáveis. Visando poder comparar diretamente os coeficientes do modelo, optamos por introduzir as covariáveis com médias e variâncias iguais, tornando-as adimensionais. Portanto, vamos iniciar a análise com o seguinte modelo:

Modelo 1:

$$Y_i = \beta_0 + \sum_{j=1}^6 \beta_j \left(\frac{x_{ji} - \bar{x}_j}{s_j} \right) + \varepsilon_i \begin{cases} i = 1, \dots, 50 \\ j = 1, \dots, 6 \end{cases}$$

Onde: $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, $x_j = \frac{1}{50} \sum_{i=1}^{50} x_{ji}$ e $s_j = \sqrt{\frac{1}{50} \sum_{i=1}^{50} (x_{ji} - \bar{x}_j)^2}$

- Y_i : Expectativa de vida em anos (1969-70).
- β_0 : Expectativa de vida esperada em anos (1969-70) para valores de covariáveis iguais às suas respectivas médias.
- $\frac{\beta_1}{s_1}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "renda percapita (em 1974 em USD)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_2}{s_2}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "proporção de analfabetos (em 1970)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_3}{s_3}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "taxa de criminalidade (por 100000 habitantes 1976)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_4}{s_4}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "porcentagem de estudantes que concluem o segundo grau (1970)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_5}{s_5}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor de "número de dias do ano com temperatura abaixo de zero grau Celsius na cidade mais importante do estado" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_6}{s_6}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "densidade da população estimada em julho de 1975 por área do estado em milhas quadradas" em uma unidade, mantendo-se as demais covariáveis fixas.

Pode-se observar a partir das figuras 2 e 3 que o modelo não teve um ajuste adequado. No gráfico A não observação tendências, o que pode indicar uma independência dos dados. Já no gráfico B, observa-se uma leve heterocedasticidade, dado que entre os valores -2 e -1 de Valores Ajustados vê-se uma discrepância entre os resíduos studentizados do que entre os resíduos compreendidos entre os valores de -1 a 2. Porém, um fator de confundimento que pode ser levado em consideração é a flutuação amostral pois o banco de dados possui apenas cinquenta observações. É possível observar uma assimetria positiva no histograma apresentado (gráfico C da figura XX), assim como uma tendência no gráfico de envelopes (figura XX) o que nos dá indícios de falta de normalidade nos resíduos. Neste caso, deveria-se procurar modelos alternativos que levem em consideração uma distribuição assimétrica. Dado o escopo do curso, procede-se com as análises posteriores.

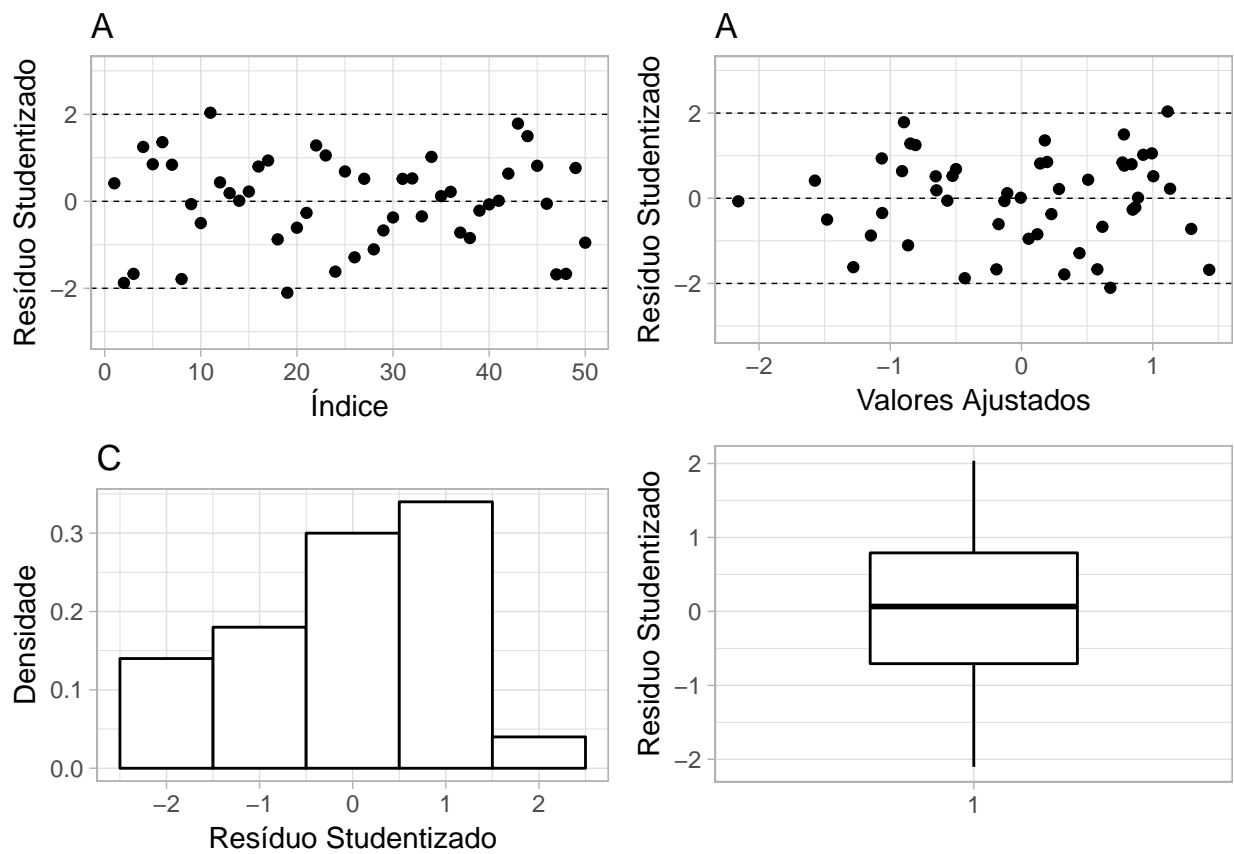


Figura 2: Análise de Resíduo para o modelo 1

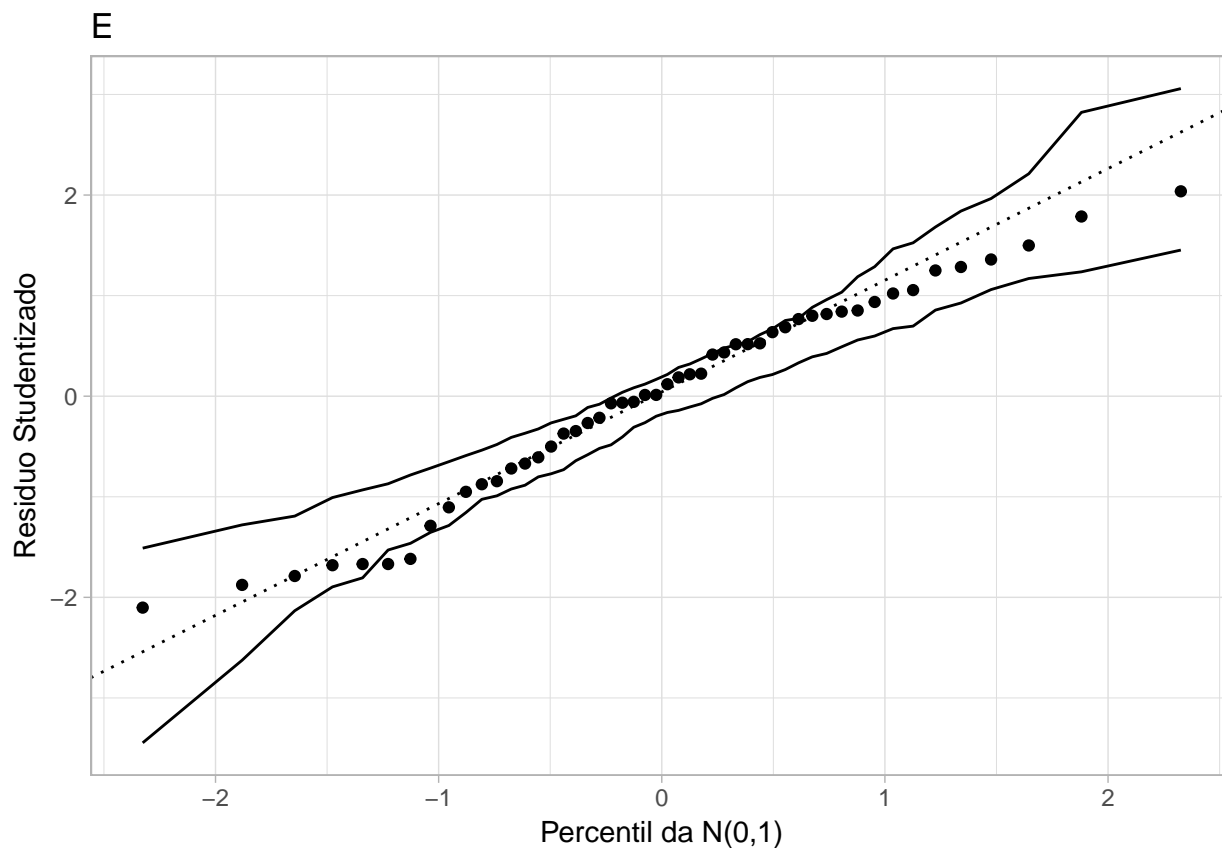


Figura 3: Gráfico de envelope para o resíduos studentizado para o modelo 1

Note, pela tabela 2 (lembrando que os valores da tabela estão arredondados em quatro casas decimais) que somente os parâmetros relacionados às covariáveis “crime”, “ndias” e “dens” significativos à um nível de significancia = 0,10. A partir desse resultado, temos a indicação de que um modelo reduzido pode ser mais apropriado. Contudo, desejamos também, obter o modelo que melhor se ajusta aos dados.

Tabela 2: Estimativas dos parâmetros, intervalo de confiança e teste de nulidade para o modelo 1

Estimativa	EP	Valor T V	alor p
0.0000	0.0774	0.0000	1
0.1007	0.1046	0.9631	0.3409
-0.0240	0.1497	-0.1601	0.8735
-0.7910	0.1137	-6.9535	<0.0001
0.1829	0.1246	1.4677	0.1495
-0.2910	0.1076	-2.7058	0.0097
-0.1571	0.0848	-1.8518	0.0709

3.1 Seleção dos modelos

A técnica escolhida para nos auxiliar a selecionar o modelo normal linear homocedástico que melhor se ajusta aos dados foi a técnica stepwise, veja Azevedo (2016).

A aplicação da metodologia stepwise, começando com o modelo só com o intercepto ou começando com o modelo completo, indicou, em ambos os casos que o modelo com que melhor se ajusta é o modelo que leva em consideração somente as covariáveis “crime”, “estud”, “ndias” e “dens”.

Temos agora o seguinte modelo:

Modelo 2:

$$Y_i = \beta_0 + \sum_{j=3}^6 \beta_j \left(\frac{x_{ji} - \bar{x}_j}{s_j} \right) + \varepsilon_i \begin{cases} i = 1, \dots, 50 \\ j = 3, \dots, 6 \end{cases}$$

Onde $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, $x_j = \frac{1}{50} \sum_{i=1}^{50} x_{ji}$ e $s_j = \sqrt{\frac{1}{50} \sum_{i=1}^{50} (x_{ji} - \bar{x}_j)^2}$

- Y_i : Expectativa de vida em anos (1969-70).
- β_0 : Expectativa de vida esperada em anos (1969-70) para valores de covariáveis iguais às suas respectivas médias.
- $\frac{\beta_3}{s_3}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "taxa de criminalidade (por 100000 habitantes 1976)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_4}{s_4}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "porcentagem de estudantes que concluem o segundo grau (1970)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_5}{s_5}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor de "número de dias do ano com temperatura abaixo de zero grau Celsius na cidade mais importante do estado" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_6}{s_6}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "densidade da população estimada em julho de 1975 por área do estado em milhas quadradas" em uma unidade, mantendo-se as demais covariáveis fixas.

Em anuência ao ajuste do modelo anterior, este modelo também não apresentou um ajuste adequado, uma vez que se observa nas figuras YY e XX indícios de mal ajuste semelhantes aos observados na análise residual para o modelo anterior. Porém uma ressalva pode ser feita de que no gráfico B pode-se observar que as discrepâncias entre os resíduos mudam conforme o valor ajustado aumenta

Neste caso, deveria-se procurar modelos alternativos que levem em consideração uma distribuição assimétrica. Novamente, dado o escopo do curso, procede-se com as análises posteriores.

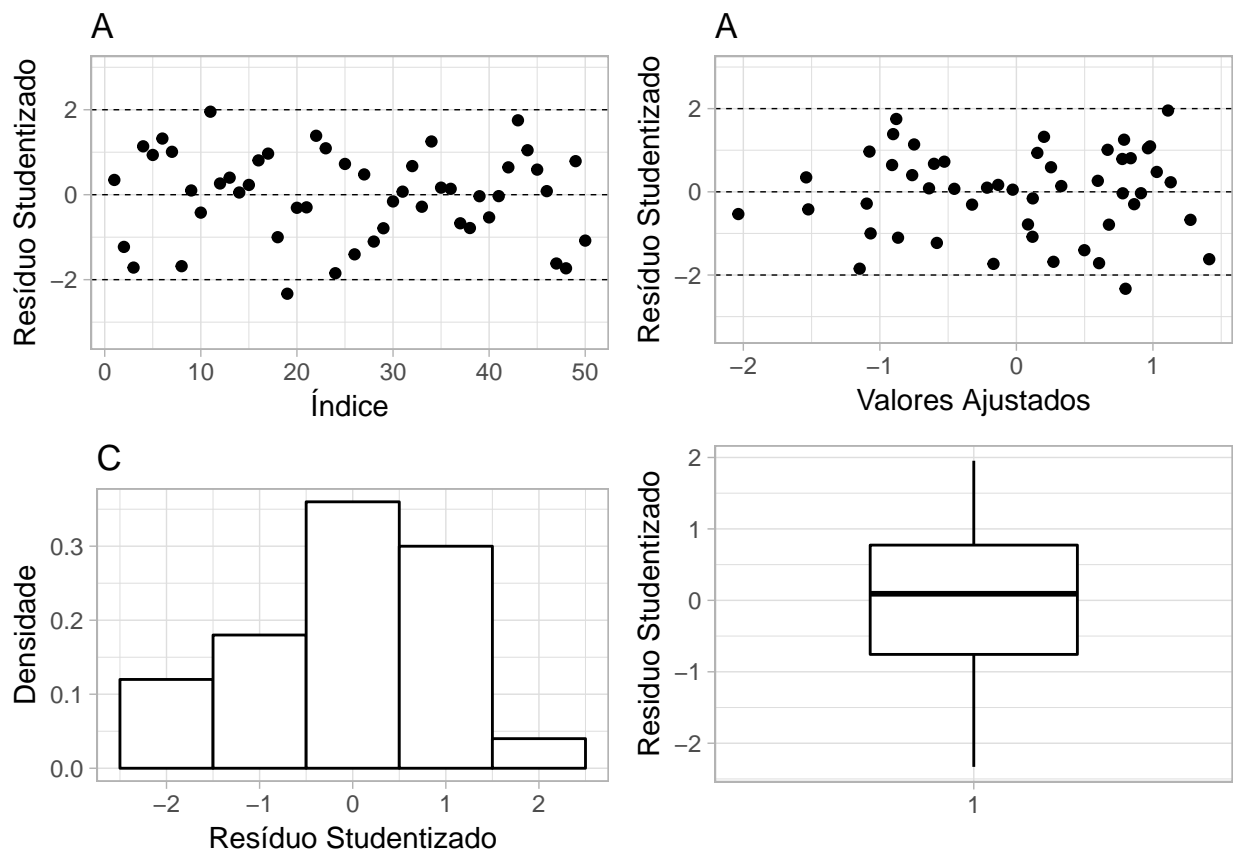


Figura 4: Análise residual para o modelo 2

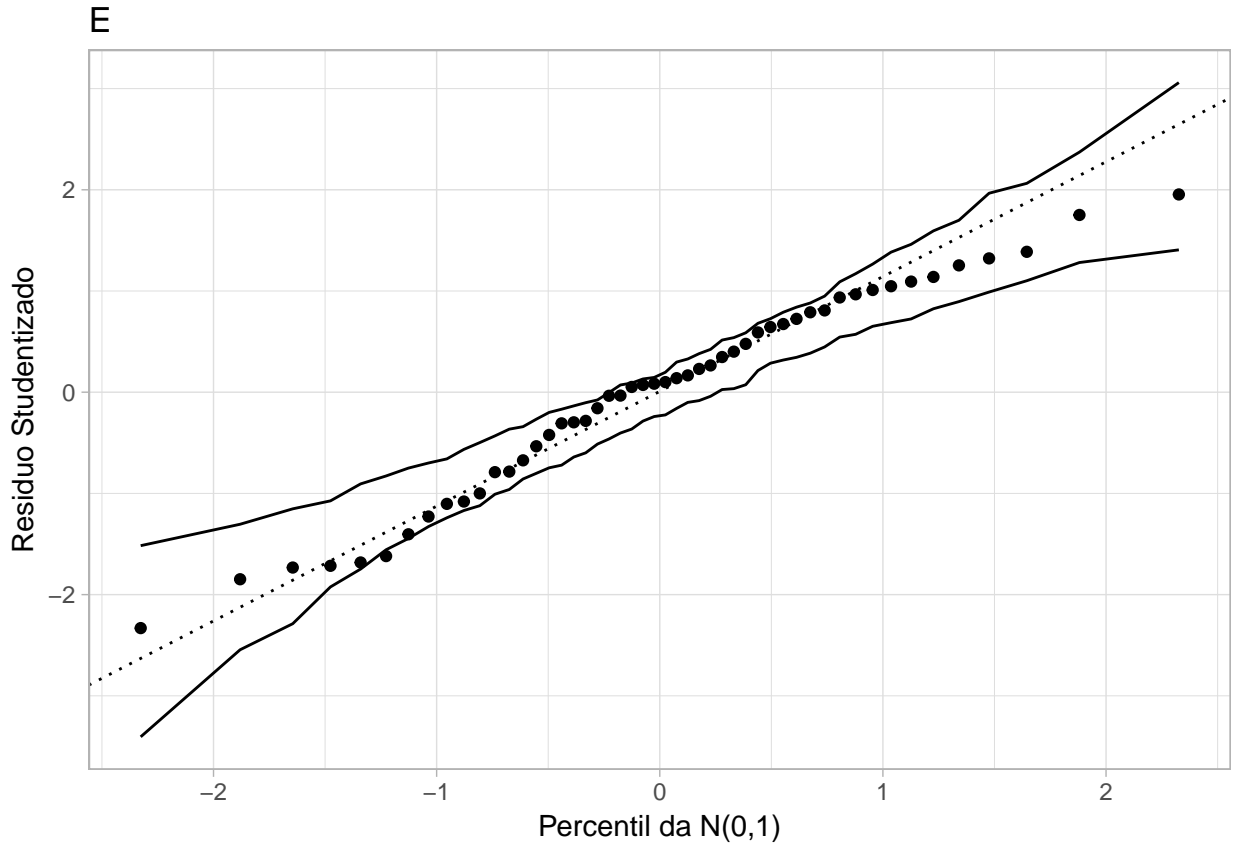


Figura 5: Gráfico de envelope para o resíduos studentizado para o modelo 2

Note que a tabela acima acusa que o parâmetro β_0 não é significativo no modelo para um nível de significância de 0,10. Sendo assim, vamos ajustar o modelo anterior sem intercepto.

Modelo 3:

$$Y_i = \sum_{j=3}^6 \beta_j \left(\frac{x_{ji} - \bar{x}_j}{s_j} \right) + \varepsilon_i \begin{cases} i = 1, \dots, 50 \\ j = 3, \dots, 6 \end{cases}$$

Onde $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, $x_j = \frac{1}{50} \sum_{i=1}^{50} x_{ji}$ e $s_j = \sqrt{\frac{1}{50} \sum_{i=1}^{50} (x_{ji} - \bar{x}_j)^2}$

- Y_i : Expectativa de vida em anos (1969-70).
- $\frac{\beta_3}{s_3}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "taxa de criminalidade (por 100000 habitantes 1976)" em uma unidade, mantendo-se as demais covariáveis fixas.
- $\frac{\beta_4}{s_4}$: Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "porcentagem de estudantes que concluem o segundo grau (1970)" em uma unidade, mantendo-se as demais covariáveis fixas.

- β_{s_5} : Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor de "número de dias do ano com temperatura abaixo de zero grau Celsius na cidade mais importante do estado" em uma unidade, mantendo-se as demais covariáveis fixas.
- β_{s_6} : Incremento (positivo ou negativo) na expectativa de vida em anos (1969-70) esperada quando se aumenta o valor da "densidade da população estimada em julho de 1975 por área do estado em milhas quadradas" em uma unidade, mantendo-se as demais covariáveis fixas.

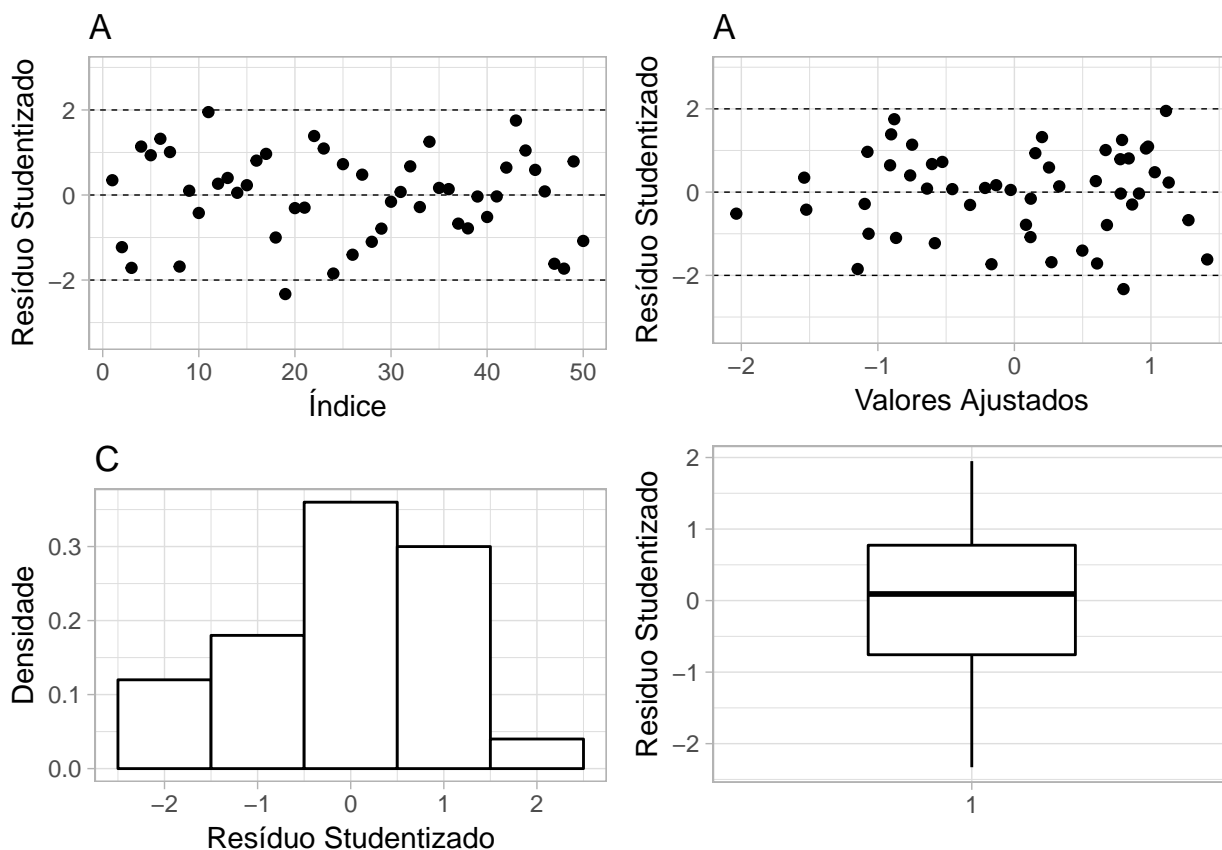


Figura 6: Análise residual para o modelo 3

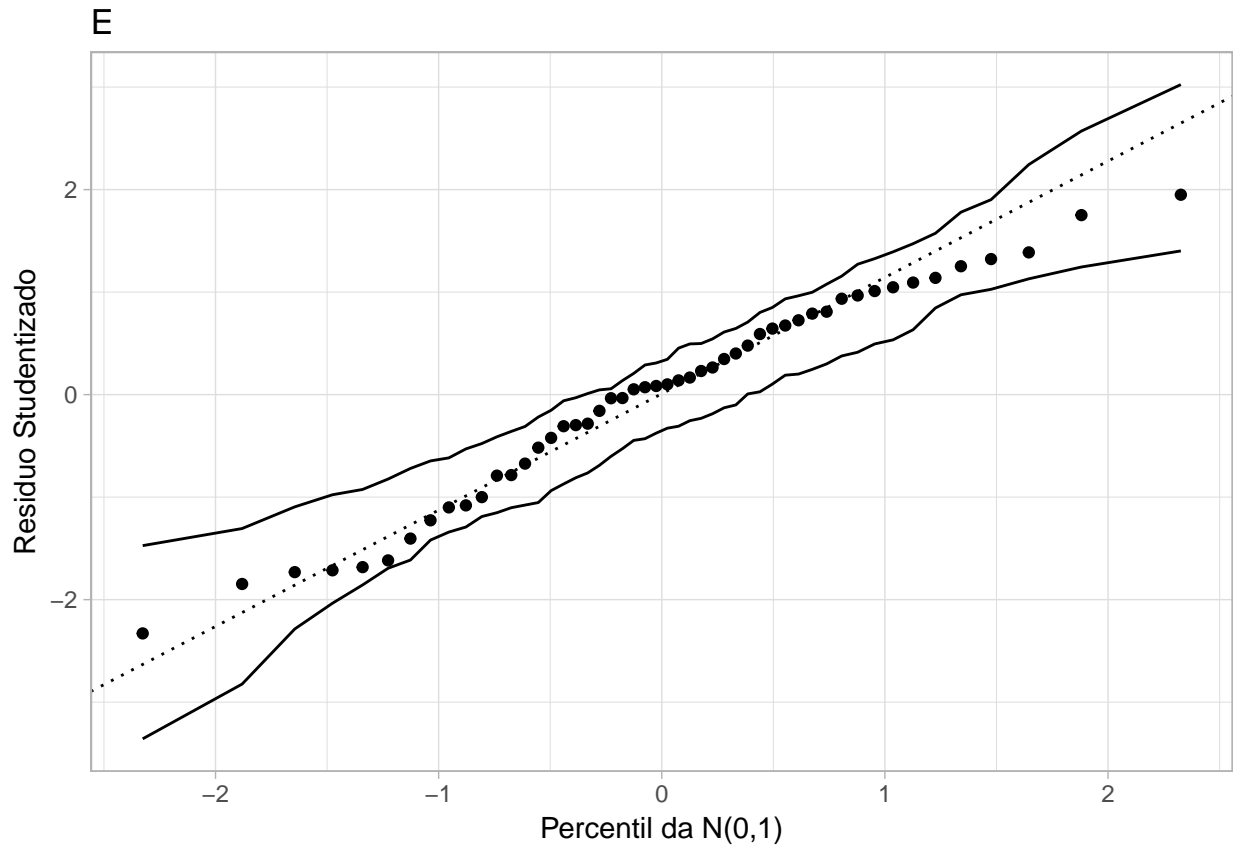


Figura 7: Gráfico de envelope para o resíduos studentizado para o modelo 3

A tabela XX abaixo mostra as principais medidas de comparação entre os dois últimos modelos (com e sem intercepto). Observamos nela, valores menores de AIC e BIC para o modelo sem intercepto, o que nos dá uma indicação de que o modelo sem intercepto é preferível.

Tabela 2: Estimativas dos parâmetros, intervalo de confiança e teste de nulidade para o modelo 1

	Modelo 2	Modelo 3
AIC	87.25904	85.25904
BIC	98.73118	94.81915
Log-Verossimilhança	-37.62952	-37.62952
R ²	0.73090	0.73090
R ² ajustado	0.70690	0.70740

Multicolinearidade

Mesmo o modelo não apresentando muitas covariáveis, é interessante, a fim de assegurar a validade dos resultados obtidos (desconsidere que o modelo não teve um bom ajuste), verificar também a possibilidade de se ter multicolinearidade presente no modelo ajustado. Com o propósito de identificar alguma indicação de multicolinearidade no modelo, podemos analisar os coeficientes de correlação linear entre as covariáveis “crime”, “estud”, “ndias” e “dens”. Porém, observamos na tabela XX (abaixo) que nenhum par de covariáveis apresenta coeficiente de correlação deveras alto, assim como a natureza das covariáveis presentes no modelo não apresentarem nenhum indício de serem fontes de multicolinearidade.

Tabela 3: Tabela dos coeficientes de correlação linear entre as covariáveis

crime	estud	ndias	dens	
crime	1.0000000	-0.4879710	-0.5388834	0.1110318
estud	-0.4879710	1.0000000	0.3667797	-0.2667561
ndias	-0.5388834	0.3667797	1.0000000	-0.1329066
dens	0.1110318	-0.2667561	-0.1329066	1.0000000

Além disso, Sabe-se que se a razão entre o maior autovalor (λ_{max}) e o menor (λ_{min}) (o chamado índice de condição) $K = \frac{\lambda_{max}}{\lambda_{min}}$ for maior que mil, geralmente, há indícios de multicolinearidade. Fazendo isso para o caso do problema em questão, tem-se que: $K = 4.8048$. Como $K < 1000$, portanto descartaremos a hipótese de multicolinearidade dos dados.

Alavancagem

Observando a figura XX podemos identificar a existência de pontos distantes dos demais, o que é uma indicação de alavancagem.

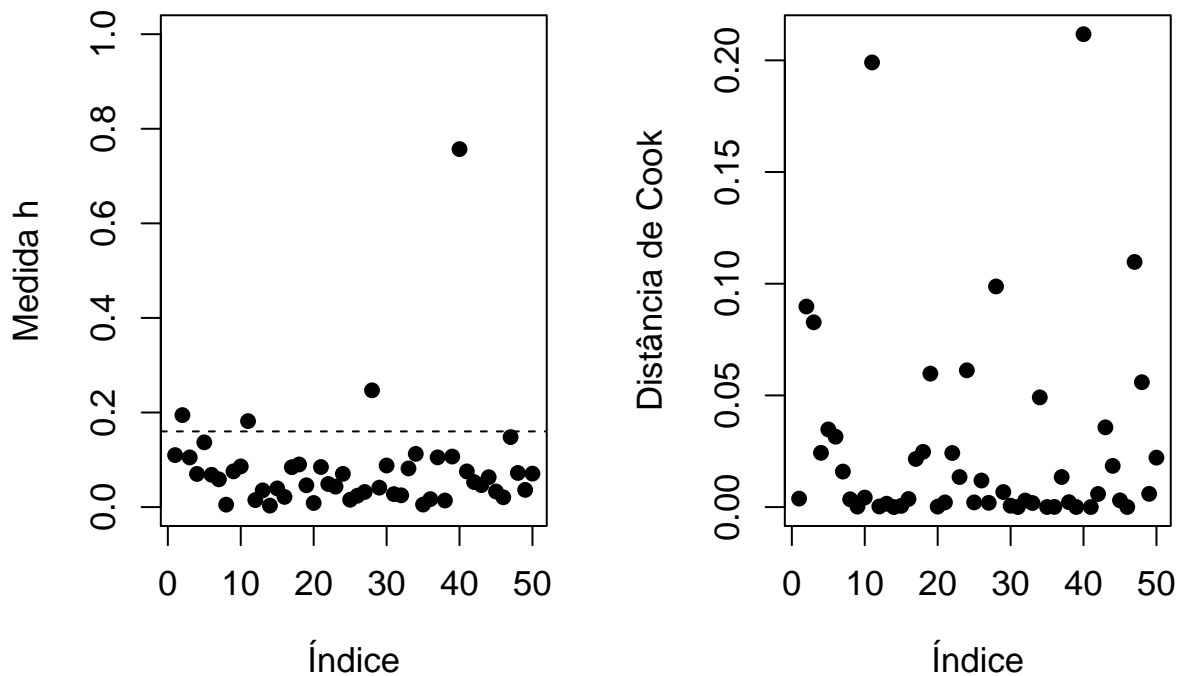


Figura 8: ESCOLHER LEGENDA

Vemos para o gráfico de Medida h , o ponto que se destaca é o de índice 40, com valor de 0.7569594. Já no de Distância de Cook, temos dois pontos em destaque: um é o do mesmo índice, com valor 0.211707, comparações entre os valores das duas medidas não consistentes e portanto inválidas, outro ponto é o de índice 11, com valor 0.1990533.

Logo, como identificamos esses pontos, vamos montar um modelo reduzido, como feito acima mas excluindo primeiramente as observações de índice 40 pois nos dois casos é o ponto mais elevado, possivelmente o que mais influência no modelo. caso ocorra uma melhora de ajuste, faremos para o outro identificado pela Distância de Cook.

Segue abaixo o ajuste do modelo sem as observações 11 e 40

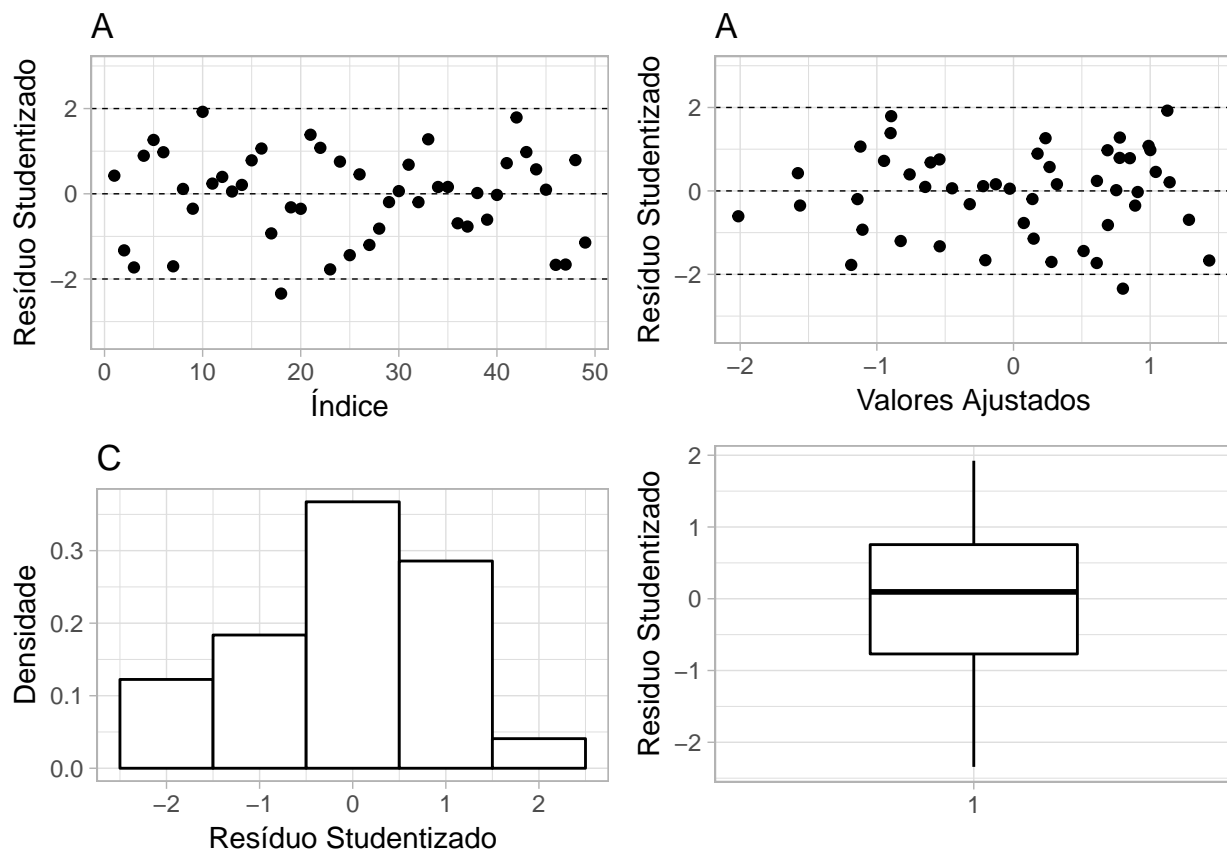


Figura 9: Gráfico de resíduos para o modelo 4 CONFERIR

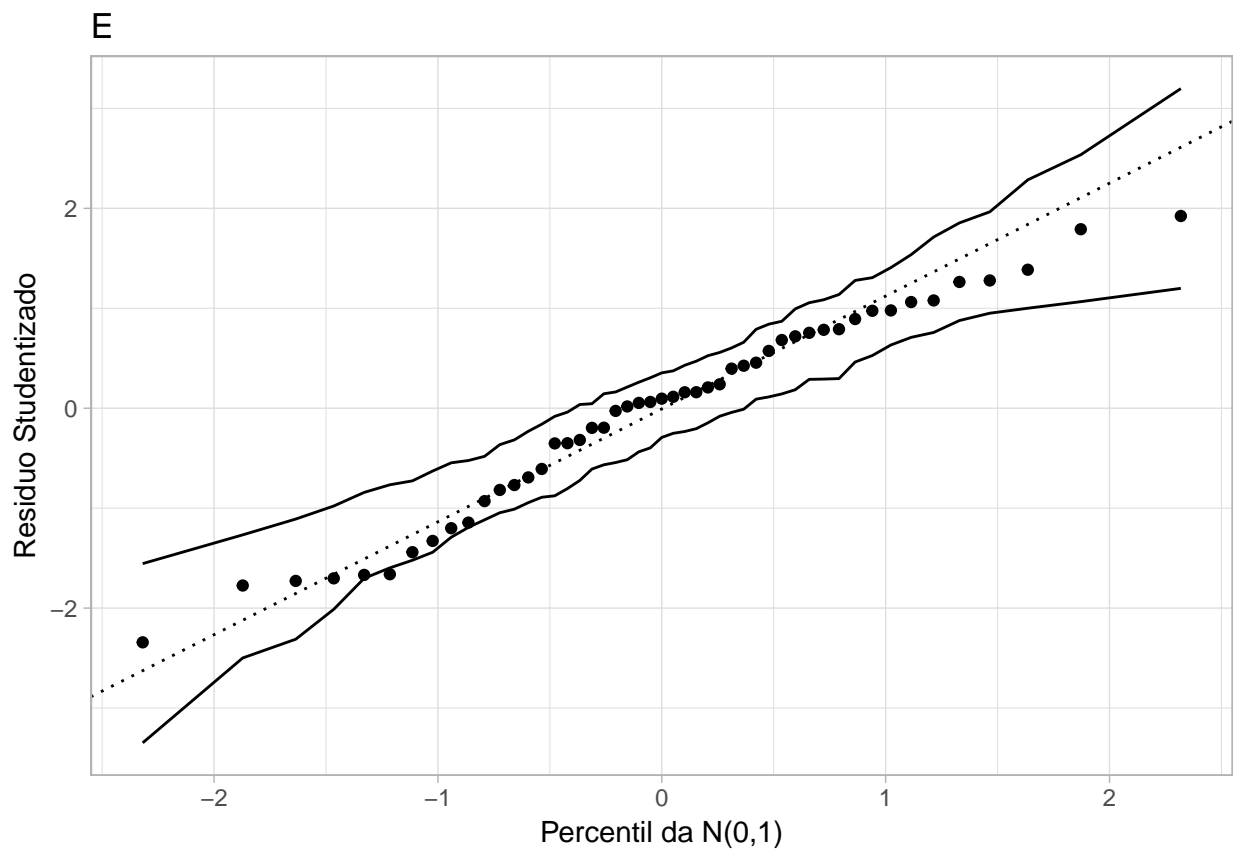


Figura 10: Gráfico de envelope para o resíduos studentizado para OQUE?

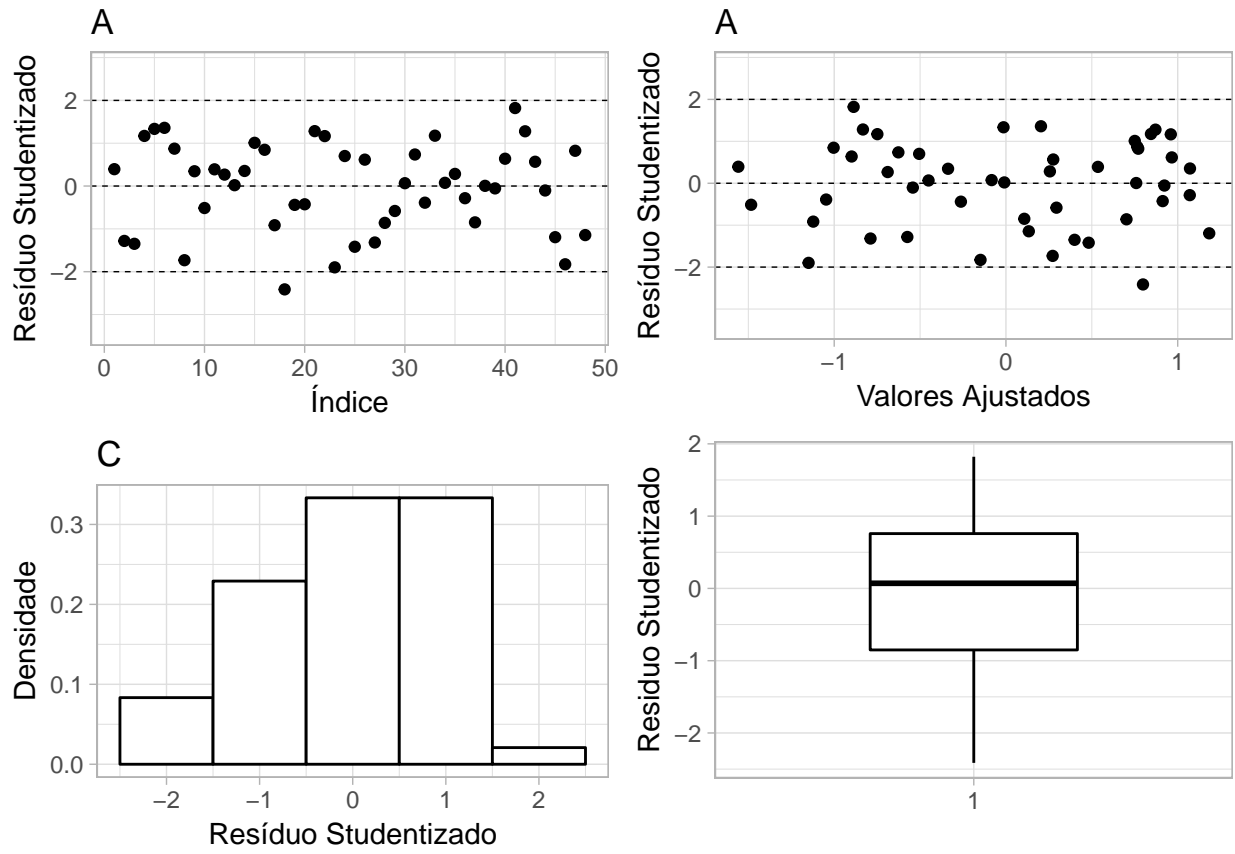


Figura 11: Gráfico de resíduos para Oque?

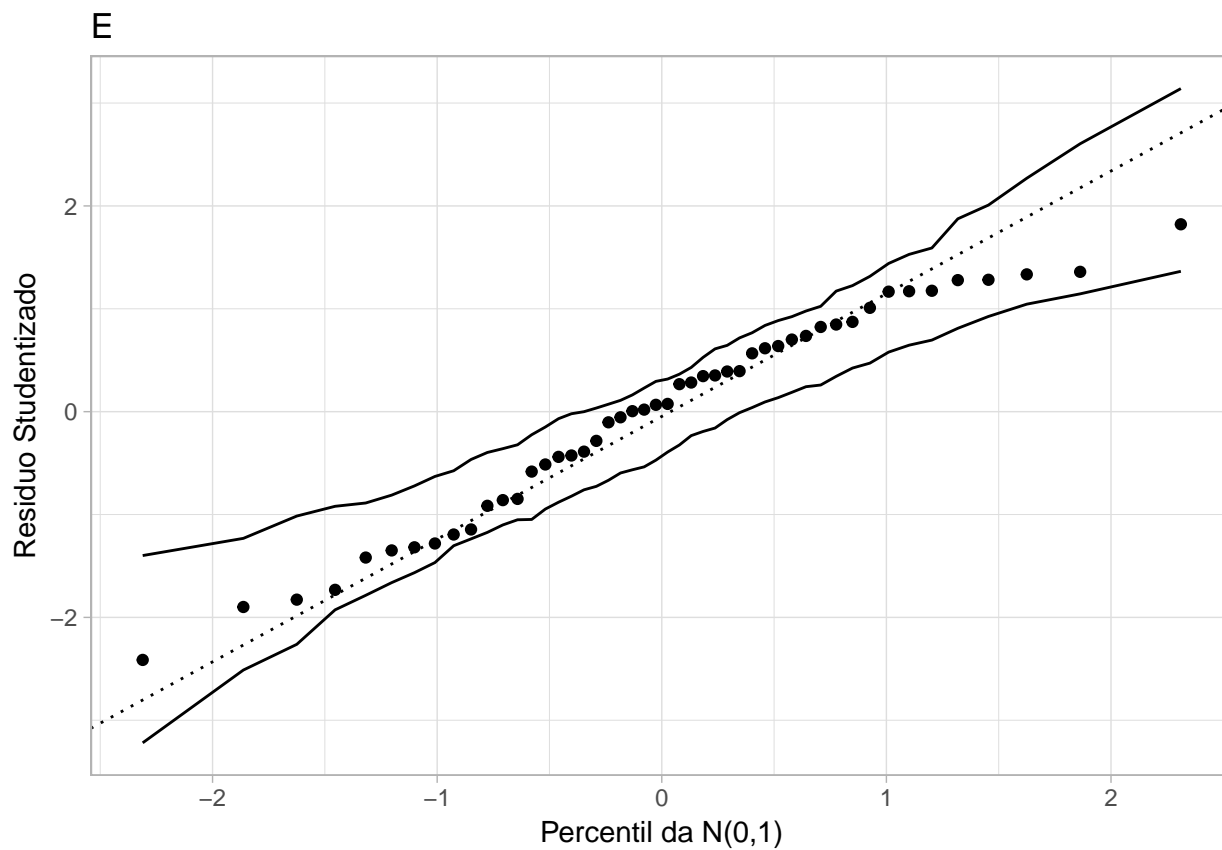


Figura 12: Gráfico de envelopes para ...

Tabela 4: Tabela de Comparação

Modelo.1	Modelo.2	Modelo.3	Modelo.4	Modelo.5	
AIC		90.07085	87.25904	85.25904	83.35105 79.76038
BIC		105.36704	98.73118	94.81915	92.81015 89.11639
Log-Verossimilhança		-37.03543	-37.62952	-37.62952	-36.67552 -34.88019
R2		0.73720	0.73090	0.73090	0.73830 0.70070
R2 ajustado		0.70050	0.70690	0.70740	0.71500 0.67350

4. Conclusões

Como visto, nenhum modelo destes propostos teve um bom ajuste ao conjunto de dados não obstante utilizaremos o modelo quatro pois é o que apresenta as melhores estimativas de comparação de modelo, ou seja baixo AIC E BIC, alta Log-Verossimilhança. Podemos observar também que há uma relação

5. Referências Bibliográficas

- Azevedo, C. L. N (2016). Notas de aula sobre planejamento e análise de experimentos, http://www.ime.unicamp.br/~cnaber/Material_ME613_2S_2016.htm
- Faraway, J. J. (2014). Linear Models with R, Second Edition, Chapman e Hall/CRC Texts in Statistical Science
- Draper, N. R. and Smith, H. (1998). Applied regression analysis, third edition. New York, NY: John Wiley e Sons.
- Paula, G. A. (2013). Modelos de regressão com apoio computacional, versão pré-eliminar https://www.ime.usp.br/~giapaula/texto_2013.pdf