



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

Relatório - Parte I 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, 06 de Dezembro de 2016

1. Introdução

O conjunto de dados analisado corresponde a informações de homens e mulheres envolvidos em exercícios regulares e apresenta para cada indivíduo, o peso(em kg) e altura(em cm) medidos e informados pelo mesmo. Além disso, o sexo de cada indivíduo também foi coletado, sendo que 112 são do sexo feminino e 88 são do sexo masculino, totalizando 200 pessoas. Os dados podem ser encontrados no R no pacote car, sob o nome “Davis”.

O objetivo é estudar o impacto da altura no peso, levando em consideração o sexo.

Utilizamos a metodologia dos modelos normais lineares homocedásticos, metodologias da qualidade do ajuste e comparação de modelos apropriados com o suporte computacional do R.

2. Análise descritiva

Observando a tabela 1, podemos ver que em média, o peso e a altura dos homens são maiores que os das mulheres. Além disso vemos valores superiores em todas as estatísticas, para os homens, incluindo uma maior variação nos dados, variação essa, mostrada pelos valores de erro padrão.

Tabela 1: Estatísticas Descritivas

Sexo	Variável	Minímo	1ºQuartil	Mediana	Média	3ºQuartil	Máximo	Erro Padrão
F	Peso	39	52,50	56	56,89	62	78	6,8905
F	Altura	148	161,50	165	164,70	169	178	5,6835
M	Peso	54	67,75	75	75,90	83	119	11,8903
M	Altura	163	173,00	178	178,00	183	197	6,4407

A figura 1 mostra um boxplot dos valores de peso por gênero e um boxplot dos valores de altura também por gênero. Observando-os podemos confirmar os padrões já identificados pelas estatísticas descritivas.

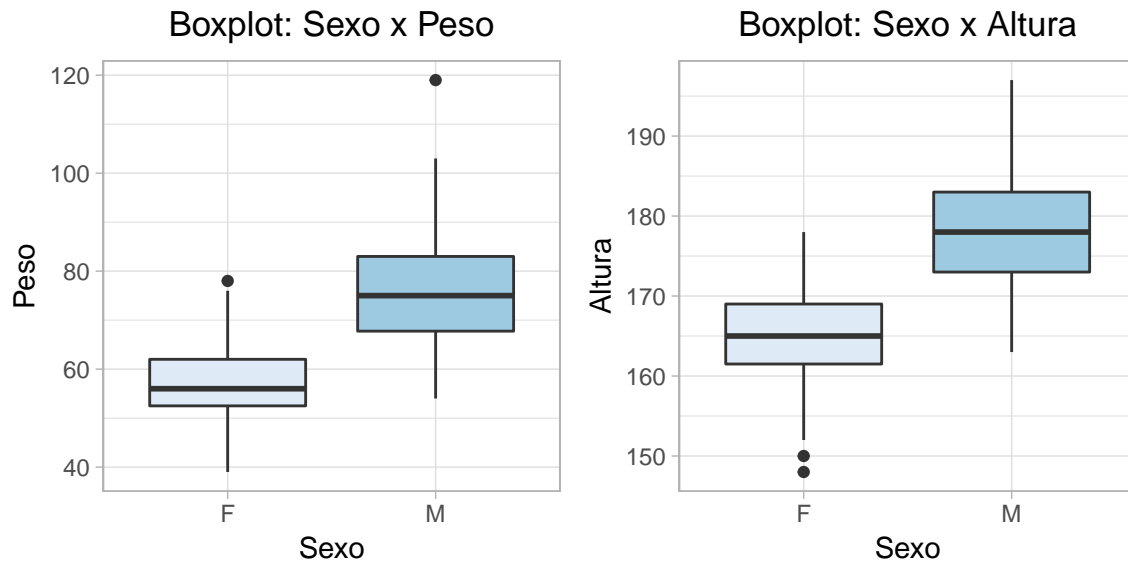


Figura 1: Boxplot Comparativo

Na figura 2 abaixo, temos um gráfico de dispersão do peso em relação a altura dos indivíduos, desconsiderando o gênero. Podemos perceber que há uma relação positiva entre a variável resposta(peso) e a variável explicativa(altura), isto é, quanto maior a altura, maior o peso do indivíduo. Tal relação pode ser razoavelmente representada por uma reta ou uma curva quadrática.

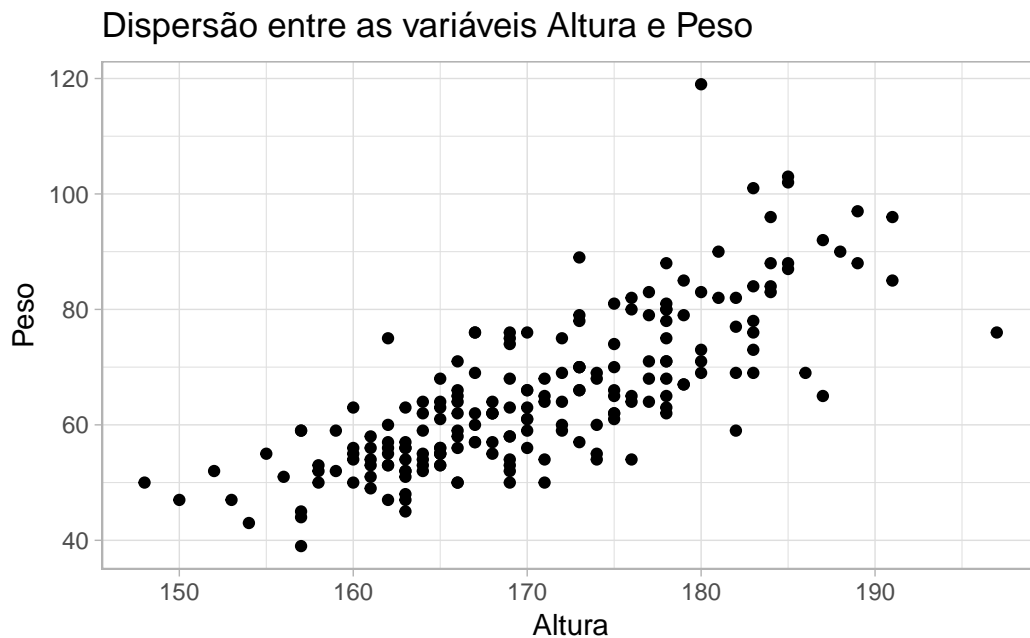


Figura 2: Dispersão entre as variáveis Altura e Peso

A figura 3 também apresenta os gráficos de dispersão entre a variável resposta e explicativa, desta vez considerando o sexo de cada indivíduo. Podemos ver que para ambos os sexos, há uma relação positiva, o peso cresce a medida que a altura cresce. Além disso é possível notar a presença de um outlier para o sexo masculino, enquanto que para o sexo feminino podemos ver em torno de três outliers. Tanto para o sexo feminino quanto o masculino podemos ainda razoavelmente representar esta relação entre peso e altura por uma reta ou uma curva quadrática.

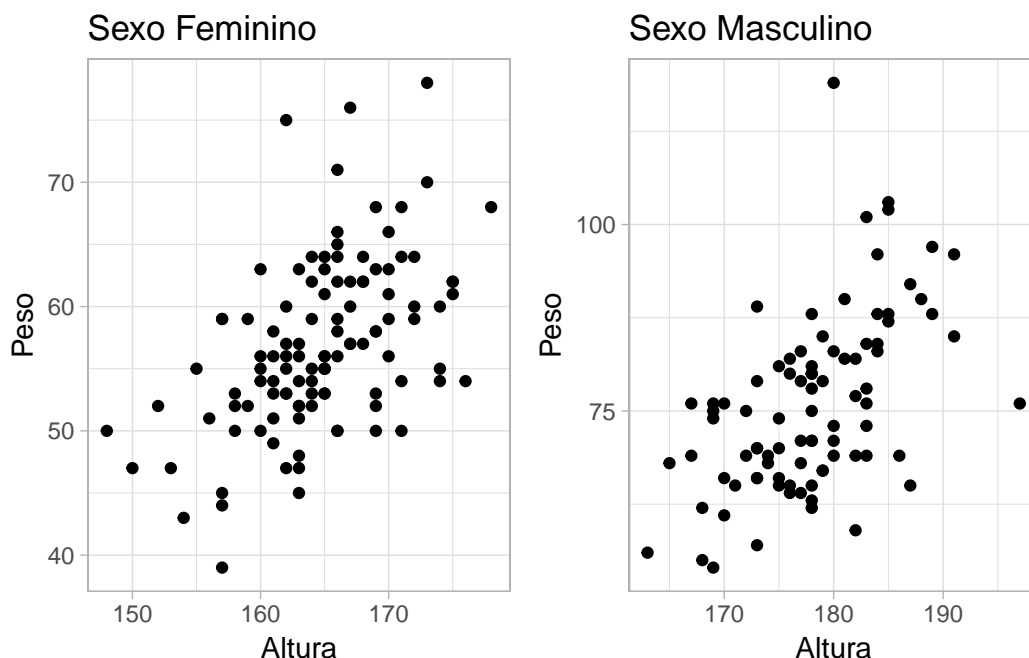


Figura 3: Dispersão entre as variáveis Altura e Peso considerando o sexo.

3. Análise Inferencial

Devido ao objetivo em questão e aos resultados da análise descritiva, vamos considerar quatro modelos e verificar qual deles é o mais reduzido e melhor se ajusta aos dados.

Para os dois primeiros, consideramos modelos com interceptos (um para cada sexo), o primeiro deles, é um modelo linear centrado na média, enquanto que o segundo é um modelo quadrático também centrado na média. Esses dois modelos nos permitem identificar se os interceptos modificam as análises sobre a significância da covariável sexo.

Para os outros dois, consideramos o modelo sem intercepto, onde o primeiro deles é modelo linear e o segundo, um modelo quadrático. Tal sugestão foi feita, pois além dos objetivos descritos acima, desejamos identificar se esses modelos fornecem informação relevante, mesmo contendo menos parâmetros.

OBS: Dado a natureza dos dados, não faz sentido observarmos algum valor diferente de zero para a variável resposta quando a variável explicativa for igual a zero, por isso, os modelos com intercepto sugeridos são todos centrados na média.

Modelo 1

$$Y_{ij} = \beta_{0i} + \beta_{1i}(x_{ij} - \bar{x}_i) + \varepsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, 199 \end{cases}$$

Onde: $\varepsilon_{ij} \sim N(0, \sigma^2)$.

- Y_{ij} : Peso do j-ésimo indivíduo do i-ésimo sexo.
- x_{ij} : Altura do j-ésimo indivíduo do i-ésimo sexo.
- β_{0i} : Peso esperado para o indivíduo do i-ésimo sexo, quando sua altura é igual ao valor da média de cada grupo.
- β_{1i} : Incremento (positivo ou negativo) no peso esperado do j-ésimo indivíduo do i-ésimo sexo, para o aumento em uma unidade na altura.

Modelo 2

$$Y_{ij} = \beta_{0i} + \beta_{1i}(x_{ij} - \bar{x}) + \beta_{2i}(x_{ij} - \bar{x})^2 + \varepsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, 199 \end{cases}$$

Onde: $\varepsilon_{ij} \sim N(0, \sigma^2)$

- Y_{ij} : Peso do j-ésimo indivíduo do i-ésimo sexo.
- x_{ij} : Altura do j-ésimo indivíduo do i-ésimo sexo.
- β_{0i} : Peso esperado indivíduo do i-ésimo sexo, quando sua altura é igual a média.
- β_{1i} : Incremento(positivo ou negativo) no peso quando se aumenta em uma unidade o valor da altura.
- $\frac{-\beta_{1i}}{2\beta_{2i}}$: Valor da altura para o qual o peso esperado é máximo ou mínimo.

Modelo 3

$$Y_{ij} = \beta_{1i}(x_{ij}) + \varepsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, 199 \end{cases}$$

Com: $\varepsilon_{ij} \sim N(0, \sigma^2)$

- Y_{ij} : Peso do j-ésimo indivíduo do i-ésimo sexo.

- x_{ij} : Altura do j-ésimo indivíduo do i-ésimo sexo.
- β_{1i} : Incremento(positivo ou negativo) no peso quando se aumenta em uma unidade o valor da altura.

Modelo 4

Considere o seguinte modelo quadrático, ainda com o fator Sexo:

$$Y_{ij} = \beta_{1i}(x_{ij}) + \beta_{2i}(x_{ij})^2 + \varepsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, 199 \end{cases}$$

Com: $\varepsilon_{ij} \sim N(0, \sigma^2)$

- Y_{ij} : Peso do j-ésimo indivíduo do i-ésimo sexo.
- x_{ij} : Altura do j-ésimo indivíduo do i-ésimo sexo.
- β_{1i} : Incremento(positivo ou negativo) no peso quando se aumenta em uma unidade o valor da altura.
- $-\frac{\beta_{1i}}{2\beta_{2i}}$: Valor da altura para o qual o peso esperado é máximo ou mínimo.

Os quatro modelos foram ajustados usando a metodologia de mínimos quadrados ordinários, mais detalhes na referência 1(Azevedo (2016)) e análises residuais foram realizadas, conforme pode ser visto na referência 4(Paula (2013)), veja as Figuras de 4 a 11. Podemos ver que o modelo 4 sugerido, apresentou um melhor ajuste, em comparação com os outros três, embora tanto este quanto os outros ajustes não sejam satisfatórios.

Para o modelo 1, na figura 4, podemos notar pelo gráfico A que uma observação se destaca das demais. Já pelo gráfico B(resíduos x valores ajustados), nota-se indícios de heterocedasticidade uma vez que, a variabilidade dos resíduos parece aumentar com o aumento dos valores ajustados. Tanto pelo gráfico C, quanto pelo gráfico D, nota-se uma leve assimetria positiva nos dados. Além disso, no gráfico D, temos a presença de outliers, o que possivelmente pode indicar a não normalidade dos dados. Além disso, o gráfico de envelopes, na figura 5, acusa um mal ajuste, devido ao comportamento levemente sistemático. Dadas as observações feitas, concluímos que o modelo não teve um ajuste adequado.

Para o modelo 2, os aspectos mencionados acima para o modelo 1 continuam presentes, com exceção do boxplot que acusou um número maior de outliers. Como podemos ver na figura 6.

Para o modelo 3, observamos na figura 8 que há uma tendência de aumento na variabilidade dos dados com o aumento dos valores ajustados no gráfico B, além disso, podemos observar dois grupos distintos, mas isto é devido à natureza dos dados, uma vez que eles estão separados por sexo. Tal tendência, indica um possível heterocedasticidade. Nota-se pelos gráficos C e D, uma leve assimetria positiva, com presença de outliers no gráfico D. O gráfico de envelopes, na figura 9, apresenta uma tendência e novamente, vários pontos fora dos envelopes, sugerindo assim, uma falta de normalidade.

Para o modelo 4, pela figura 10 podemos ver um ponto discrepante tanto no gráfico A como no gráfico B. No gráfico A, não se observa nenhuma tendência que indique dependência dos dados. No gráfico B, notamos um pequeno agrupamento dos dados e logo depois uma pequena dispersão dos mesmos, tal tendência, sugere uma possível heterocedasticidade dos dados. Ambos os gráficos C e D, mostram uma leve assimetria positiva, onde no gráfico C e no gráfico D, temos a presença de outliers, tanto abaixo como acima da média. O gráfico de envelopes para o modelo 4, na figura 11, apresenta uma tendência e vários pontos fora dos envelopes, sugerindo assim, uma falta de normalidade.

Concluimos, então, que nenhum dos quatro modelos se ajustou bem. Contudo, devido ao fato de não podermos escolher modelos para além da classe dos modelos lineares normais homocedásticos, iremos continuar comparando os modelos propostos acima.

As estatísticas de comparação dos quatro modelos podem serem vistas na tabela 2. Seguindo o critério de escolha usando as estatísticas AIC e BIC, temos novamente que o modelo 1 é o modelo mais bem ajustado. Na tabela 2 também podemos ver os valores para R^2 e R^2 ajustado para cada modelo, para todos os modelos os valores são relativamente altos, indicando que os respectivos modelos têm alto poder explicativo quanto à variabilidade dos dados. Sabemos porém que, a simples visualização dessas quantidades, por si só, não é suficiente para se medir a adequabilidade do modelo.

Na Tabela 3 apresentamos os principais resultados relativos à estimação pontual e intervalar dos quatro modelos. Podemos ver que para os modelos 1 e 3, todos os parâmetros são significativos para qualquer nível de significância usual (0,01 à 0,10). Já para o modelo 2, os parâmetros β_{01} , β_{02} , β_{11} e β_{12} são significativos para os mesmos níveis de significância citados acima, enquanto que os parâmetros β_{21} e β_{22} não são significativos (à níveis de significância usuais). Para o modelo 4, apenas o parâmetro β_{22} é significativo à qualquer nível de significância usual, enquanto que β_{21} é significativo para níveis a partir de 0,05, já os demais parâmetros não são significativos (à níveis de significância usuais).

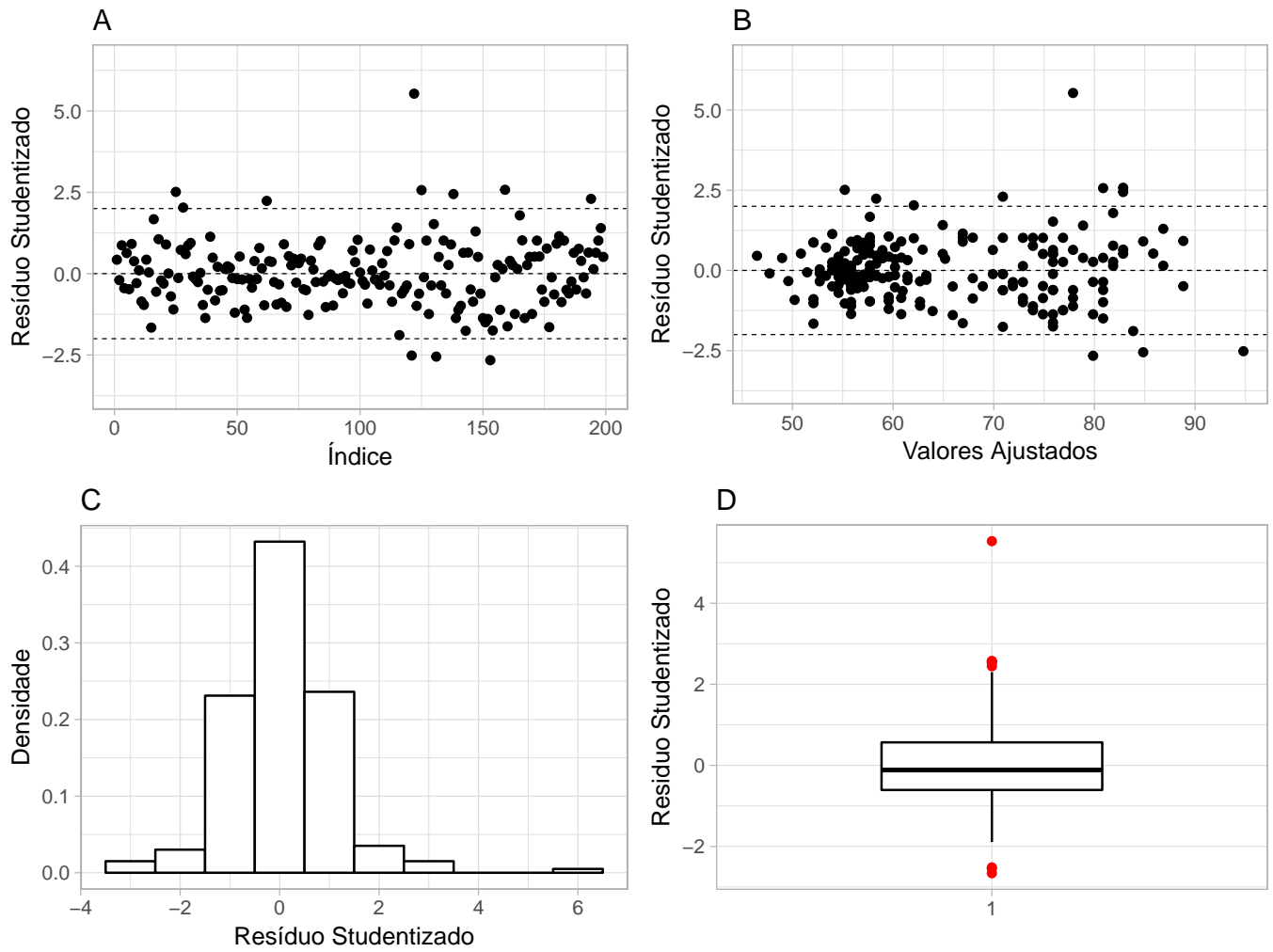


Figura 4: Análise residual para o modelo 1

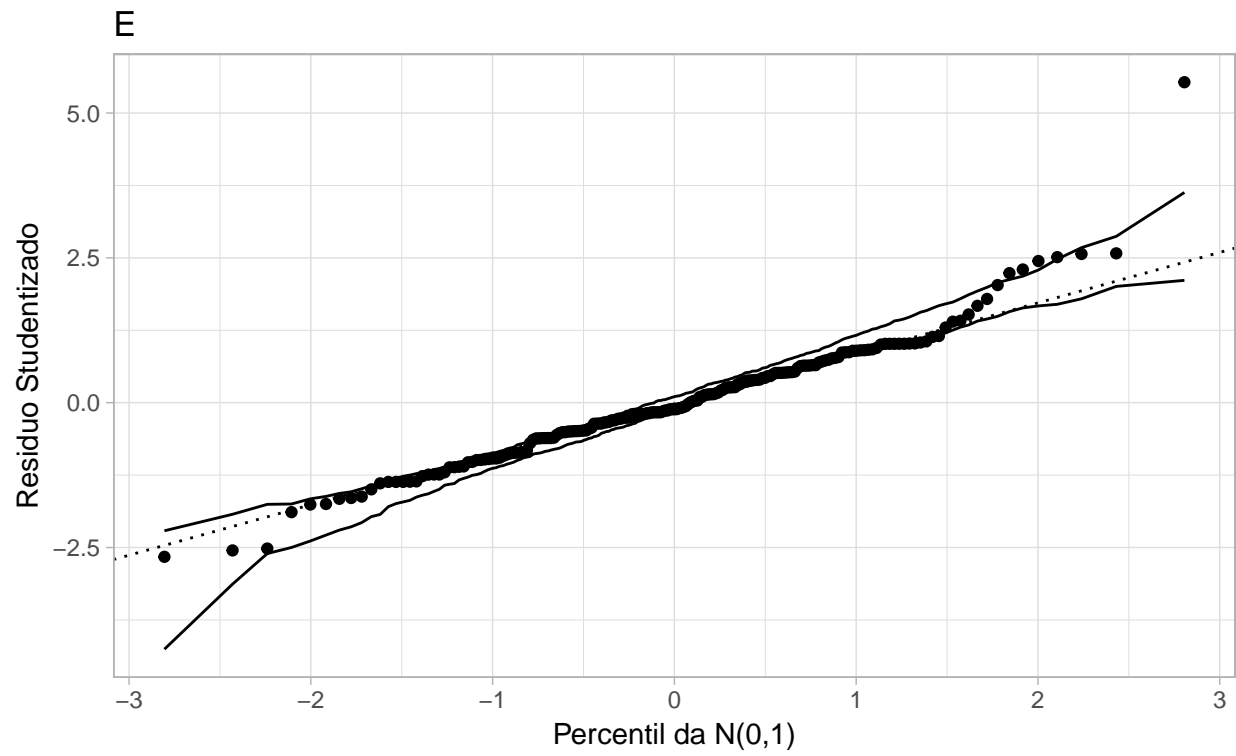


Figura 5: Gráfico de envelope para o resíduos studentizado para o modelo 1

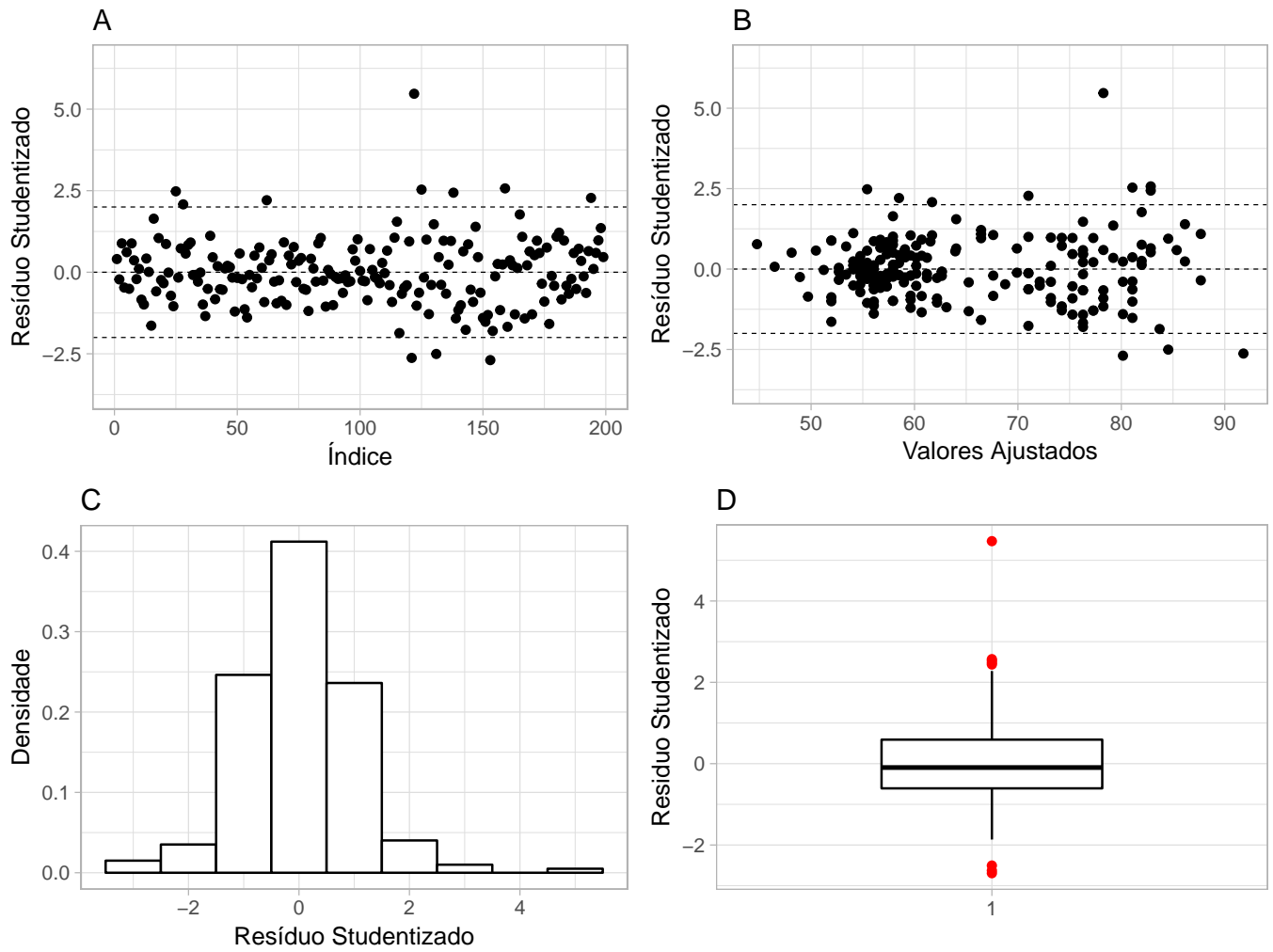


Figura 6: Análise residual para o modelo 2

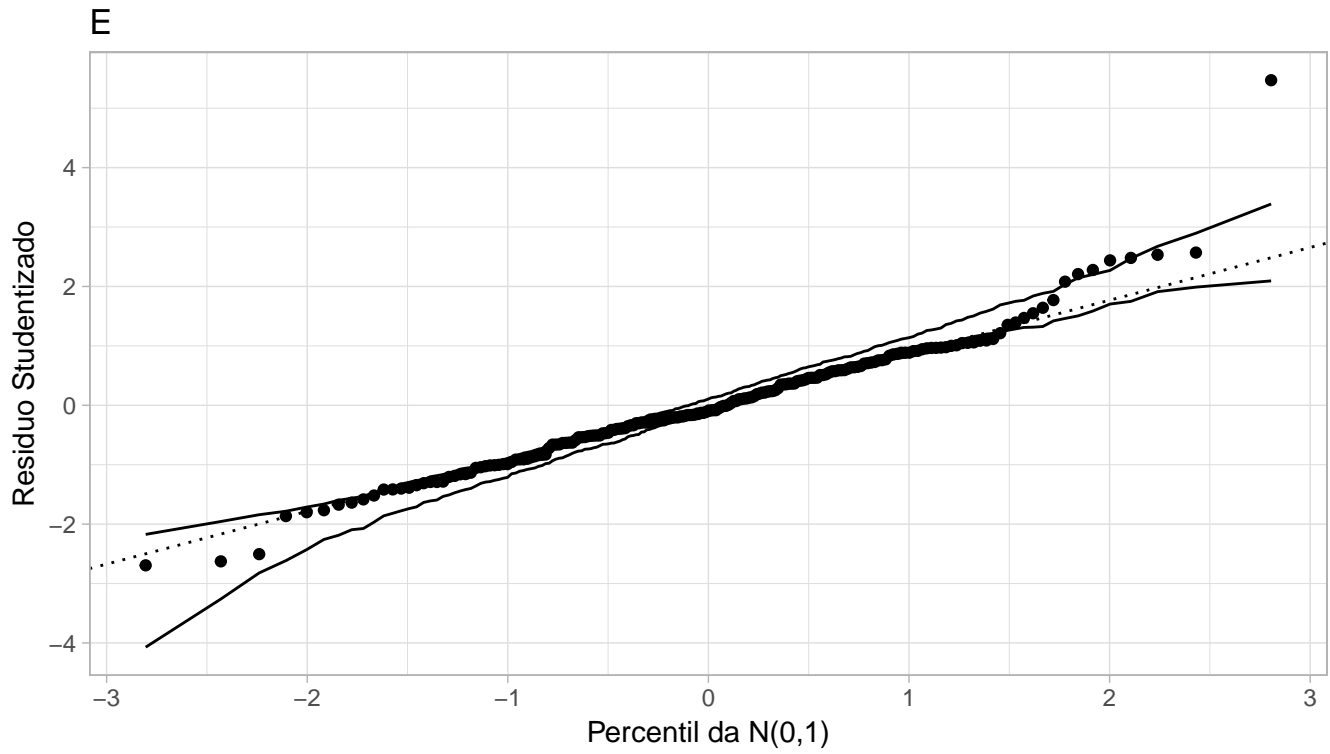


Figura 7: Gráfico de envelope para o resíduos studentizado para o modelo 2

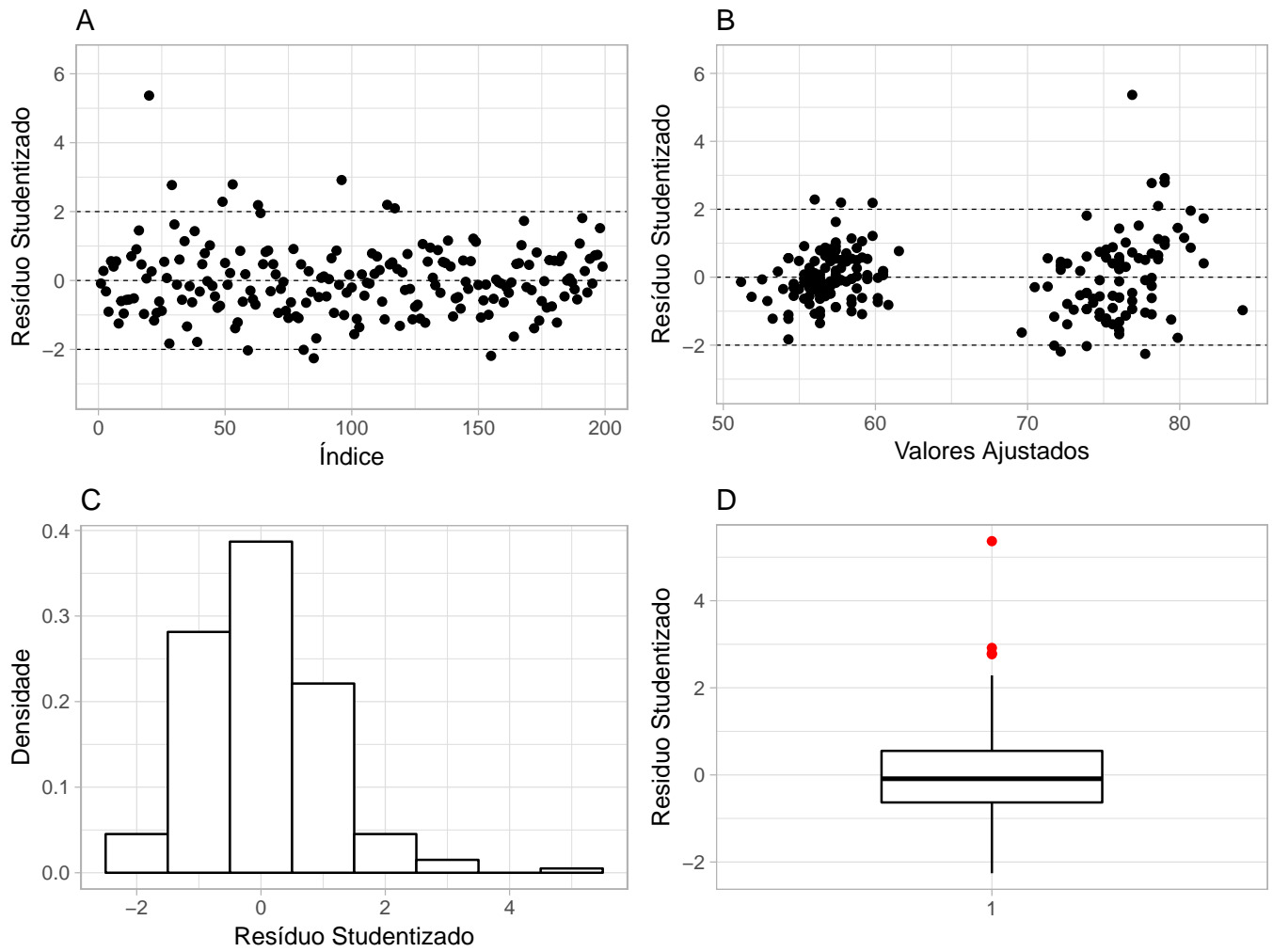


Figura 8: Análise residual para o modelo 3

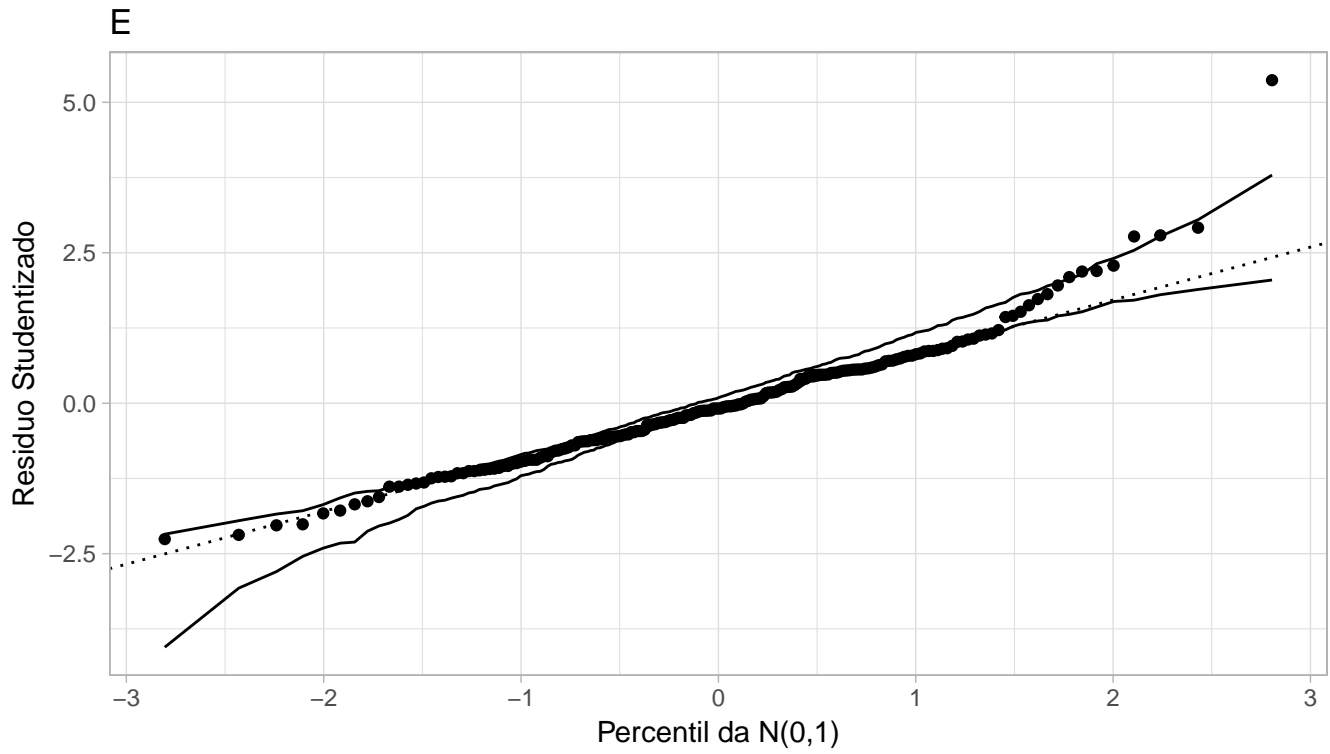


Figura 9: Gráfico de envelope para o resíduos studentizado para o modelo 3

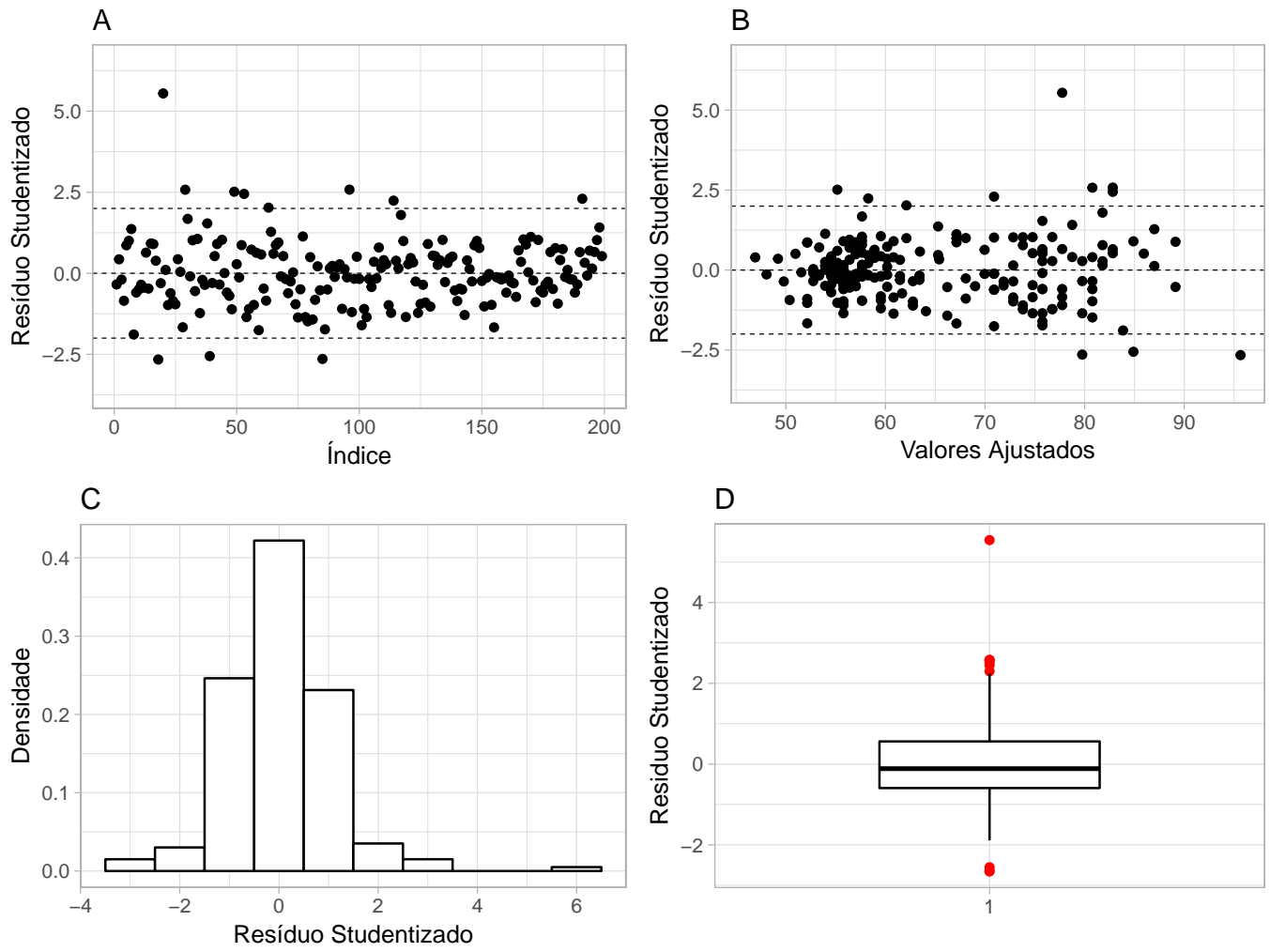


Figura 10: Análise residual para o modelo 4

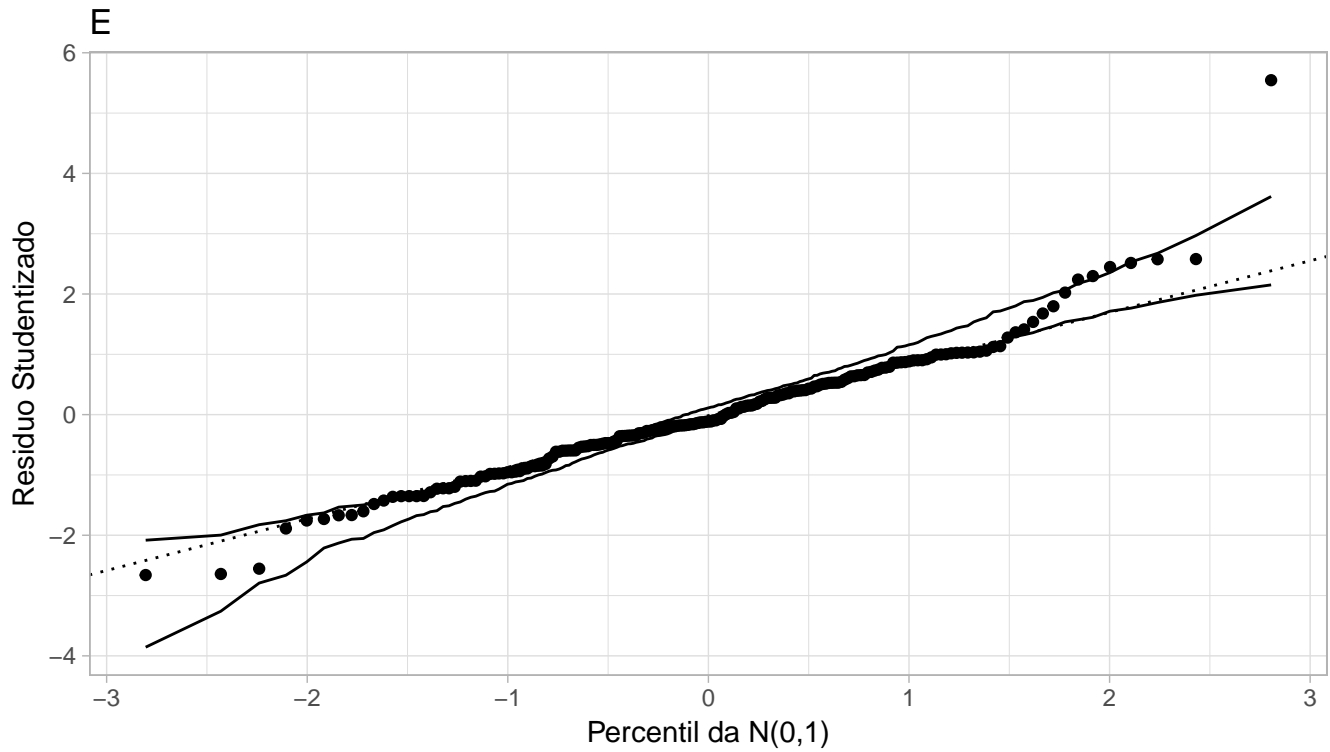


Figura 11: Gráfico de envelope para o resíduos studentizado para o modelo 4

Tabela 2: Comparação dos modelos

	Modelo 1	Modelo 2	Modelo 3	Modelo 4
AIC	1399,6915	1403,01552	1417,2977	1400,1353
BIC	1416,158,	1426,0687	1427,1776	1416,6018
Log-Verossimilhança	-694,8458	-694,5078	-705,6488	-695,0677
R^2	0,9858	0,9858	0,9841	0,9857
R^2 Ajustado	0,9855	0,9854	0,984	0,9855

Tabela 3: Estimativas dos parâmetros, intervalo de confiança e teste de nulidade

Modelo	Parâmetro	Estimativa	EP	Valor T	Valor p
Modelo 1	β_{01}	56,8919	0,7620	74,6654	<0,0001
	β_{02}	75,8977	0,8558	88,6906	<0,0001
	β_{11}	0,6229	0,1347	4,6256	<0,0001
	β_{12}	0,9956	0,1336	7,4505	<0,0001
Modelo 2	β_{01}	57,1353	0,9269	61,6432	<0,0001
	β_{02}	76,3053	1,0556	72,2870	<0,0001
	β_{11}	0,6132	0,1368	4,4831	<0,0001
	β_{12}	1,0051	0,1349	7,4532	<0,0001
	β_{21}	-0,0076	0,0164	-0,4646	0,6427
	β_{22}	-0,0099	0,0150	-0,6639	0,5075
Modelo 3	β_{11}	0,3457	0,0049	71,1911	<0,0001
	β_{12}	0,4271	0,0050	84,6343	<0,0001
Modelo 4	β_{11}	0,0690	0,1364	0,5059	0,6135
	β_{12}	-0,1354	0,1336	-1,0136	0,3120
	β_{21}	0,0017	0,0008	2,0310	0,0436
	β_{22}	0,0032	0,0007	4,2136	<0,0001

4. Conclusões

Nenhum dos quatro modelos propostos teve bom ajuste ao conjunto de dados. Porém, dado o escopo do curso (classe de modelos lineares normais homocedásticos), continuou-se com as análises optando-se por selecionar o modelo que melhor se adequa aos dados. Conclui-se então que o melhor modelo, segundo os critérios de seleção e comparação de modelos AIC, BIC e Log-Verossimilhança, é o modelo 1 pois teve os menores valores entre os valores observados de AIC's e BIC's dos modelos, e o segundo maior das Log-Verossimilhanças comparadas dos modelos. Nota-se também, pela tabela 3, que todos os coeficientes do modelo 1 são diferentes de zero, para qualquer nível de significância usual (0,01 à 0,10). Os modelos apontaram uma tendência crescentes nos dados, vistas na análise descritiva para ambos os sexos. Vê-se que o aumento em uma unidade de altura acarreta em um aumento no peso do indivíduo e este aumento, segundo os resultados obtidos, não muda sua magnitude conforme os valores de alturas observadas.

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