



Universidade Federal de Pernambuco  
Centro de Ciências Exatas e da Natureza  
Programa de Pós-Graduação em Estatística

**RODOLPHO JÓRDAN DOMINGOS QUINTELA**

**INFLUÊNCIA LOCAL EM MODELOS PARCIALMENTE LINEARES ADITIVOS  
GENERALIZADOS**

**Recife**

**2020**

**RODOLPHO JÓRDAN DOMINGOS QUINTELA**

**INFLUÊNCIA LOCAL EM MODELOS PARCIALMENTE LINEARES  
ADITIVOS GENERALIZADOS**

Dissertação apresentada ao Programa de Pós-Graduação em Estatística do Centro de Ciências Exatas e da Natureza da Universidade Federal de Pernambuco, como requisito parcial à obtenção do título de mestre em Estatística.

**Área de Concentração:** Estatística Matemática

Orientador: Dr. Roberto Ferreira Manghi

Coorientador: Dr. Francisco José de Azevedo  
Cysneiros

**Recife**

**2020**

Catálogo na fonte  
Bibliotecária Monick Raquel Silvestre da S. Portes, CRB4-1217

Q7i Quintela, Rodolpho Jórdan Domingos  
Influência local em modelos parcialmente lineares aditivos generalizados /  
Rodolpho Jórdan Domingos Quintela. – 2020.  
125 f.: il., fig., tab.

Orientador: Roberto Ferreira Manghi.  
Dissertação (Mestrado) – Universidade Federal de Pernambuco. CCEN,  
Estatística, Recife, 2020.  
Inclui referências, apêndice e anexo.

1. Estatística matemática. 2. Diagnóstico. I. Manghi, Roberto Ferreira  
(orientador). II. Título.

519.5                      CDD (23. ed.)                      UFPE- CCEN 2020 - 128

RODOLPHO JÓRDAN DOMINGOS QUINTELA

INFLUÊNCIA LOCAL EM MODELOS PARCIALMENTE LINEARES ADITIVOS  
GENERALIZADOS

Dissertação apresentada ao Programa de Pós-Graduação em Estatística da Universidade Federal de Pernambuco, como requisito parcial para a obtenção do título de Mestre em Estatística.

Aprovada em: 28 DE FEVEREIRO DE 2020

**BANCA EXAMINADORA**

---

Prof.(º) Roberto Ferreira Manghi  
UFPE

---

Prof.(a) Audrey Helen Mariz de Aquino Cysneiros  
UFPE

---

Prof.(º) Gilberto Alvarenga Paula  
USP

À minha família, por me apoiar em todo momento e por ser um porto seguro. Mãe, seu amor e cuidado me deram forças pra continuar. Pai, és meu maior exemplo.

## AGRADECIMENTOS

Acredito que a vida é fundamentada em escolhas. É claro que algumas delas não cabem a nós, ou seja, não temos controle sobre tais decisões. Notadamente, nos deparamos com questões do tipo: O que é preciso fazer para ter uma vida bem sucedida? Aristóteles dirá que é preciso fazer desabrochar as potências que são as tuas. E até chegar a este estado, inevitavelmente, nos esbarramos à fatores que não estão sob o nosso controle. Nesse momento, mediante esses choques singulares com o mundo, acabamos tomando decisões ao sabor das circunstâncias da vida.

Assim, na tentativa de encontrar a minha alegria, isto é, o desabrochar das minhas potências, entendi que não existe uma fórmula para vida. Para mim, existe uma quantidade finita de alternativas que nos são apresentadas. Obviamente optamos por aquelas que julgamos serem as melhores baseado em tudo aquilo que achamos certo. Como todo e qualquer ser humano, às vezes erramos em ter escolhido por um caminho ao invés de outro, acertamos em alguns momentos e em outras situações não resolvemos arriscar. Omitir-se é sem dúvidas a pior de todas as escolhas.

Portanto, não é fácil decidir-se sobre alguma coisa. Quase sempre estamos diante de decisões conflitantes as quais não sabemos exatamente o que fazer. Tais escolhas são complexas e ter a responsabilidade sobre elas é angustiante. Mesmo assim, diante disso tudo, temos que fazê-las. Porque também são elas que vão proporcionar instantes de vida alegres os quais não gostaríamos que acabassem tão logo, fazendo da nossa vida sempre alguma coisa digníssima de ser buscada e fantástica de ser vivida. Dessa forma, cada escolha, cada momento, cada pessoa foram importantes e, definitivamente, contribuíram para a minha formação enquanto pessoa.

Logo, em minha trajetória como aluno da pós-graduação, contei com a colaboração de inúmeras pessoas. Portanto, gostaria muito de agradecer pelo menos minimamente cada um deles como eles merecem.

- Ao professor Luis Gonzaga Pinheiro Felix. Meu primeiro amigo do mestrado. Hoje, mais que antes, tenho a plena convicção que de um jeito ou de outro, seríamos grandes amigos.
- Ranah Duarte, Larissa Lima e Adenice Ferreira, minhas conselheiras. Agradeço pelos momentos de conversas sobre relacionamentos, sobre amizade, sobre a vida. Pelas incontáveis vezes que me ajudaram nas disciplinas tornando assim, minha vida nesse mestrado menos penosa. Aprendi muito com vocês.
- Cristine Oliveira, de todos os colegas de pós, era talvez a que eu menos imaginava que nos

tornaríamos amigos. Nunca me enganei tanto e nunca fiquei tão feliz por isso. Você é uma pessoa maravilhosa.

- Calinhos, um grande amigo a quem devo muito estima.
- Fernando Luiz Maia Gomes, sempre de bem com a vida, muito obrigado pela força de sempre.
- José Jairo, o alívio cômico dessa série. Meu camarada, você é um cara que com essa alegria, bom humor e humildade é o tipo de pessoa que todo mundo quer por perto. Mas o que seria do herói sem seu fiel companheiro? João Eudes Miquéias é outro que tem alegria no coração. Continue assim que vocês vão longe.
- César Diogo foi de longe o cara mais legal que conheci na vida. Seu senso crítico e de justiça me mostrou que não é preciso ser algo que não somos só pra agradar terceiros. Meu amigo, você faz isso de uma maneira espetacularmente educada, concisa e coerente. É pra poucos!
- Nayara Luíza é a pessoa com a qual eu mais me identifiquei. Costumo dizer que ela sou eu na versão feminina, só que bem melhor.
- Anabeth Radünz e Eduardo Ensslin. Torço muito por vocês, pois os tenho como inspiração. Obrigado por tudo.
- Cris Guedes e Lucas Araújo, meus conterrâneos, obrigado por, além de tudo, trazerem um pouco da terrinha pra perto de mim. Graças a vocês, foi mais fácil lidar com a saudade de casa.
- Vinícios Scher, não tenho palavras pra descrever tanto respeito e admiração que tenho pela sua pessoa. Obrigado por todas as vezes que você quebrou meu galho. Sendo sempre muito compressível e disposto a ajudar.
- Lucas de Miranda, aqui deixo também registrado meus mais sinceros agradecimentos. Anny Kerollayne, Thalytta Evilly, Lucas David, Bruna, Tatiane Fontora, Saul De Azevêdo Souza (Shineray), Alê, Pedro Almeida, Charles, Luana. Antigos e novos amigos, serei eternamente grato à vocês. Muito Obrigado!
- Gostaria também de agradecer ao Prof. Dr. Roberto Ferreira Manghi e ao Prof. Dr. Francisco José de Azevêdo Cysneiros por me orientarem nesta dissertação.
- Um agradecimento muito especial a Profa. Dra. Audrey Helen Cysneiros. Não tenha dúvidas de que a senhora será uma ótima lembrança em minha vida. Muito obrigado por tudo.
- A todos o professores do Departamento de Estatística da UFPE e, de modo geral, a todos

os funcionários que fazem parte deste departamento. Muito Obrigado!

- Agradeço à Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) pelo apoio financeiro a este projeto.
- Não poderia deixar de agradecer a mãe do meu filho, Amélia Coelho Maciel, que me ajudou de diferentes formas para que eu pudesse terminar essa dissertação. Inclusive, nas correções, dando ideias, sugestões dentre outras coisas. A ela também devo alguns dos meus dias mais felizes, em especial, por ter me proporcionado a maior alegria de todas: ser pai. Por isso, serei eternamente grato.
- Dedico também este trabalho ao meu avô, João Oliveira (*in memoriam*).
- E finalmente, gostaria de agradecer, de forma mais que especial, à minha família. À minha mãe, Maria Regina, que sempre acreditou em mim ainda que eu mesmo não acreditasse. Ao meu pai, Francisco Wilson, e minhas irmãs, Maria Jordayane e Jordânia Horrana que me deram todo apoio às minhas decisões mesmo nos momentos mais difíceis e em escolhas que possivelmente pudessem desagradá-los. Espero um dia poder retribuir tudo que vocês fizeram por mim com todo o amor que lhes é inerente.

## RESUMO

Esta dissertação tem como objetivo propor resíduos e técnicas para a análise de diagnóstico nos Modelos Parcialmente Lineares Aditivos Generalizados (MPLAGs), tais como: alavancagem generalizada, análise de resíduos, dois quais propomos utilizar os resíduos de Pearson e resíduos aleatorizados, bem como medidas para análise de influência local sob os seguintes esquemas de perturbação: perturbação de de casos, perturbação na variável resposta e perturbação em uma das variáveis explicativas. Para isto, derivamos tais medidas fundamentados em uma vasta pesquisa bibliográfica e conceitual sobre tais métodos no contexto dos MPLAGs. Essas técnicas foram utilizadas em exemplos de aplicação a dados reais e os resultados foram discutidos a fim de avaliar o nosso estudo teórico. Para tanto, apresentamos as equações de estimação para os parâmetros do modelo através da função de verossimilhança penalizada, considerando como estrutura não paramétrica o uso de *P-splines*. Assim, definimos tal modelo, buscando apresentar algumas propriedades e vantagens que motivam o uso de *P-splines* no contexto de regressão não paramétrica. Por fim, o método iterativo *backfitting* (Gauss-Seidel) é utilizado para a obtenção das estimativas.

**Palavras-chave:** Diagnóstico. Influência Local. MPLAGs. *P-splines*.

## ABSTRACT

This dissertation aims at the proportion of residues and techniques for diagnostic analysis in Generalized Additive Partial Linear Models (GAPLMs), such as: generalized leverage, analysis of residues, two of which we propose to use Pearson's residues and randomized residues, as well as measures for analysis of local influence under the following perturbation schemes: perturbation of cases, perturbation in response variable and perturbation in one of the variables explanatory. For this, we derive such measures based on a vast bibliographic and conceptual research on such methods in the context of GAPLMs. These techniques were used in real data application examples and the results were discussed in order to evaluate our theoretical study. For that, the estimation equations for the model parameters are provided through the penalized likelihood function, considering the use of *P-splines* as a non-parametric structure. Thus, we define such a model, searching to present some properties and advantages that motivate the use of *P-splines* in the context of non-parametric regression. Finally, the iterative backfitting method (Gauss-Seidel) is used to obtain the estimates.

**Keywords:** Diagnostics. Local Influence. GAPLMs. *P-splines*.

## LISTA DE FIGURAS

<b>Figura 1</b> – Gráfico de dispersão da aceleração contra o tempo. . . . .	22
<b>Figura 2</b> – Contorno suave obtido com auxílio de um <i>splines</i> . . . . .	23
<b>Figura 3</b> – Representação de um ajuste por meio de uma <i>spline</i> linear com nós em $\xi_1 = 15$ e $\xi_2 = 30$ . . . . .	26
<b>Figura 4</b> – Aproximação de curvas por meio de <i>splines</i> . . . . .	27
<b>Figura 5</b> – Representação de um ajuste por meio de <i>spline</i> linear e cúbica ambas com nós em $\xi_1 = 1.5$ e $\xi_2 = 2.5$ . . . . .	28
<b>Figura 6</b> – Ilustração de um B- <i>spline</i> de grau 1 isolado com três nós posicionados em "x". . . . .	29
<b>Figura 7</b> – Ilustração de um B- <i>splines</i> de grau 2 isolado com cinco nós posicionados em "x". . . . .	30
<b>Figura 8</b> – Velocidade do impactador em função do tempo de contato entre o projétil e o alvo para uma energia de impacto constante. . . . .	54
<b>Figura 9</b> – Gráficos da distribuição de frequência, frequência acumulada e boxplot, respectivamente, para a variável tempo. . . . .	54
<b>Figura 10</b> – Gráficos da distribuição de frequência, frequência acumulada e boxplot, respectivamente, para a variável aceleração. . . . .	55
<b>Figura 11</b> – O valor do critério VCG, curva aproximada e intervalos de confiança pontuais de 95% para a aceleração contra o tempo. . . . .	56
<b>Figura 12</b> – Gráficos da distribuição de frequência, frequência acumulada e boxplot respectivos da a variável idade para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .	70
<b>Figura 13</b> – Gráficos com as porcentagens para cada faixa de idade e relação da idade com a quantidade de votos recebidas pelos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .	71
<b>Figura 14</b> – Estes gráficos mostram a relação da idade com os gastos e diagrama de dispersão: votos segundo os gastos na campanha dos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .	71
<b>Figura 15</b> – Gráficos com as porcentagens dos sexo na amostra e relação da sexo dos candidatos com as quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .	72

<b>Figura 16 – Gráficos com as porcentagens dos sexo na amostra e relação da sexo dos candidatos com as quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .</b>	<b>73</b>
<b>Figura 17 – Gráfico com a relação do sexo com o estado civil dos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .</b>	<b>74</b>
<b>Figura 18 – Resultado das eleições de 2006 ao Senado brasileiro segundo o estado civil. . . . .</b>	<b>74</b>
<b>Figura 19 – Gráficos com as porcentagens dos níveis de instrução na amostra e relação das quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro . . . . .</b>	<b>75</b>
<b>Figura 20 – Gráfico da curva ajustada sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.</b>	<b>77</b>
<b>Figura 21 – Gráficos de resíduos referente ao modelo ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .</b>	<b>78</b>
<b>Figura 22 – Gráficos de alavancagem e influência local sob os esquemas de perturbação de caso, perturbação na variável resposta e perturbação em umas das variáveis explicativas, respectivamente, referentes ao ajuste do modelo sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro. . . . .</b>	<b>80</b>
<b>Figura 23 – Gráficos da distribuição de frequência, frequência acumulada e boxplot respectivos da a variável preço do aluguel para o exemplo do aluguel em Munique. . . . .</b>	<b>83</b>
<b>Figura 24 – Gráficos com as proporções de cada variável categórica na amostra para o exemplo do aluguel em Munique. . . . .</b>	<b>84</b>
<b>Figura 25 – Boxplots do valor do aluguel segundo as variáveis categóricas para o exemplo do aluguel em Munique. . . . .</b>	<b>85</b>
<b>Figura 26 – Diagramas de dispersão do valor do aluguel contra as variáveis numéricas para o exemplo do aluguel em Munique. . . . .</b>	<b>86</b>
<b>Figura 27 – Gráfico das correlações entre as variáveis quantitativas para o exemplo do aluguel em Munique. . . . .</b>	<b>86</b>
<b>Figura 28 – Gráficos das curvas ajustadas sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>88</b>

<b>Figura 29 – Gráfico das curvas ajustadas para o segundo modelo sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>89</b>
<b>Figura 30 – Gráficos de resíduos referente ao modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>90</b>
<b>Figura 31 – Gráficos de alavancagem e influência local sob os esquemas de perturbação de caso, perturbação na variável resposta e perturbação em umas das variáveis explicativas, respectivamente, referentes ao ajuste do modelo sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>92</b>
<b>Figura 32 – Boxplots das variáveis presentes no baco de dados considerando o exemplo do ar poluído em Chicago. . . . .</b>	<b>94</b>
<b>Figura 33 – Diagrama de dispersão do logaritmo do número de mortes contra as demais variáveis numéricas considerando o exemplo do ar poluído em Chicago. . . . .</b>	<b>96</b>
<b>Figura 34 – Gráficos das curvas ajustadas sob distribuição Poisson para o exemplo do ar poluído em Chicago. . . . .</b>	<b>98</b>
<b>Figura 35 – Gráficos de resíduos referentes ao modelo ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago. . . . .</b>	<b>100</b>
<b>Figura 36 – Gráfico de influência local sob os esquemas de perturbação de caso e perturbação na variável resposta respectivamente para o exemplo do ar poluído em Chicago. . . . .</b>	<b>101</b>
<b>Figura 37 – Gráfico que mostra como está distribuída a reprodução das aves Cotovia-de-crista em Portugal. . . . .</b>	<b>102</b>
<b>Figura 38 – Curva aproximada e intervalos de confiança pontuais de 95% referente ao modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves. . . . .</b>	<b>104</b>
<b>Figura 39 – Gráfico de resíduos referente o modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves. . . . .</b>	<b>105</b>
<b>Figura 40 – Gráficos de de alavancagem e influência local sob os esquemas de perturbação de caso e perturbação em uma das variáveis explicativas, respectivamente, referente ao ajuste do modelo sob distribuição binomial para o exemplo da reprodução das aves. . . . .</b>	<b>106</b>

## LISTA DE TABELAS

<b>Tabela 0</b>	<b>– Resumo das estimativas referentes ao modelo ajustado para o exemplo do impacto do capacete. . . . .</b>	<b>56</b>
<b>Tabela 1</b>	<b>– Resultado das eleições de 2006 ao Senado brasileiro. . . . .</b>	<b>72</b>
<b>Tabela 2</b>	<b>– Resultado das eleições de 2006 ao Senado brasileiro segundo o estado civil. . . . .</b>	<b>73</b>
<b>Tabela 3</b>	<b>– Resultado das eleições de 2006 ao Senado brasileiro segundo o nível de instrução. . . . .</b>	<b>75</b>
<b>Tabela 4</b>	<b>– Resumo das estimativas referentes ao modelo ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro. . . .</b>	<b>76</b>
<b>Tabela 5</b>	<b>– Resumo das estimativas referentes ao modelo final ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro. . . .</b>	<b>77</b>
<b>Tabela 6</b>	<b>– Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição normal e mudança relativa em porcentagem entre parênteses para os parâmetros <math>\beta_0</math>, <math>\beta_1</math> e <math>\beta_2</math> após a retirada dos pontos para o exemplo das eleições de 2006 ao senado brasileiro. . . . .</b>	<b>81</b>
<b>Tabela 7</b>	<b>– Resumo das estimativas de máxima verossimilhança referentes ao modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>87</b>
<b>Tabela 8</b>	<b>– Resumo das estimativas de máxima verossimilhança referentes ao segundo modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique. . . . .</b>	<b>89</b>
<b>Tabela 9</b>	<b>– Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição gama e mudança relativa em porcentagem entre parênteses para os parâmetros <math>\beta_0</math>, <math>\beta_1</math>, <math>\beta_2</math>, <math>\beta_3</math>, <math>\beta_4</math>, <math>\beta_5</math> e <math>\beta_6</math> após a retirada dos pontos para o exemplo do aluguel em Munique. . . . .</b>	<b>93</b>
<b>Tabela 10</b>	<b>– Resumo das estimativas de máxima verossimilhança referentes ao modelo ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago. . . . .</b>	<b>97</b>
<b>Tabela 11</b>	<b>– Resumo das estimativas de máxima verossimilhança referentes ao modelo final ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago. . . . .</b>	<b>98</b>

<b>Tabela 12 – Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição poisson e mudança relativa em porcentagem entre parênteses para os parâmetros <math>\beta_0</math> e <math>\beta_1</math> após a retirada dos pontos para o exemplo do ar poluído em Chicago. . . . .</b>	<b>101</b>
<b>Tabela 13 – Resumo das estimativas de máxima verossimilhança referentes ao modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves . . . . .</b>	<b>103</b>
<b>Tabela 14 – Estimativas de máxima verossimilhança e mudança relativa em porcentagem entre parênteses para os parâmetros <math>\beta_0</math> e <math>\beta_1</math> do modelo binomial após a retirada dos pontos para o exemplo da reprodução das aves em Portugal.</b>	<b>106</b>

## LISTA DE CÓDIGOS-FONTE

<b>Código-fonte 1</b>	<b>– Função para ajuste de uma splines linear</b>	<b>111</b>
<b>Código-fonte 2</b>	<b>– Função para ajuste de uma splines cúbica</b>	<b>111</b>
<b>Código-fonte 3</b>	<b>– Função para ajuste do modelo</b>	<b>114</b>
<b>Código-fonte 4</b>	<b>– Obtenção de alpha</b>	<b>115</b>
<b>Código-fonte 5</b>	<b>– Plot - Pontuação VCG</b>	<b>117</b>
<b>Código-fonte 6</b>	<b>– Estimativas do modelo ótimo</b>	<b>117</b>
<b>Código-fonte 7</b>	<b>– Função para análise de influência local</b>	<b>118</b>

## SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO</b> . . . . .	<b>18</b>
1.1	ORGANIZAÇÃO DA DISSERTAÇÃO . . . . .	20
1.2	SUPORTE COMPUTACIONAL . . . . .	21
1.3	CONCEITOS PRELIMINARES . . . . .	21
<b>1.3.1</b>	<b>Splines</b> . . . . .	<b>22</b>
<b>1.3.2</b>	<b><i>Splines</i> lineares</b> . . . . .	<b>25</b>
<b>1.3.3</b>	<b><i>Splines</i> cúbicas</b> . . . . .	<b>26</b>
<b>1.3.4</b>	<b><i>B-splines</i></b> . . . . .	<b>29</b>
<b>1.3.5</b>	<b><i>P-splines</i></b> . . . . .	<b>30</b>
<b>2</b>	<b>MODELOS PARCIALMENTE ADITIVOS GENERALIZADOS</b> . . . . .	<b>33</b>
2.1	INTRODUÇÃO . . . . .	33
2.2	ESPECIFICAÇÃO DO MODELO . . . . .	34
<b>2.2.1</b>	<b>Função de verossimilhança penalizada</b> . . . . .	<b>35</b>
<b>2.2.2</b>	<b>Restrição de identificabilidade</b> . . . . .	<b>40</b>
2.3	ESTIMAÇÃO DOS PARÂMETROS . . . . .	46
2.4	INFERÊNCIA ESTATÍSTICA . . . . .	48
<b>2.4.1</b>	<b>Graus de liberdade efetivos</b> . . . . .	<b>48</b>
<b>2.4.2</b>	<b>Seleção do parâmetro de suavização</b> . . . . .	<b>49</b>
<b>2.4.3</b>	<b>Escolha do número de nós</b> . . . . .	<b>50</b>
<b>2.4.4</b>	<b>Intervalos de confiança</b> . . . . .	<b>51</b>
<b>2.4.5</b>	<b>Teste de hipóteses</b> . . . . .	<b>52</b>
<b>2.4.6</b>	<b>Impacto do capacete da motocicleta</b> . . . . .	<b>53</b>
<b>3</b>	<b>TÉCNICAS DE DIAGNÓSTICO NOS MPLAGS</b> . . . . .	<b>58</b>
3.1	ANÁLISE DE RESÍDUOS . . . . .	59
<b>3.1.1</b>	<b>Resíduos de Pearson</b> . . . . .	<b>59</b>
<b>3.1.2</b>	<b>Resíduos quantílicos aleatorizados</b> . . . . .	<b>60</b>
3.2	MEDIDAS DE ALAVANCAGEM . . . . .	60
3.3	INFLUÊNCIA LOCAL . . . . .	62
<b>3.3.1</b>	<b>Perturbação de casos</b> . . . . .	<b>65</b>
<b>3.3.2</b>	<b>Perturbação da variável resposta</b> . . . . .	<b>65</b>
<b>3.3.3</b>	<b>Perturbação em uma das variáveis explicativas</b> . . . . .	<b>66</b>

3.3.4	<b>Perturbação em vetores particionados</b>	67
4	<b>APLICAÇÕES</b>	69
4.1	ELEIÇÕES DE 2006 AO SENADO BRASILEIRO	69
4.1.1	<b>Análise exploratória dos dados</b>	69
4.1.2	<b>Ajuste do modelo</b>	75
4.1.3	<b>Diagnóstico</b>	78
4.2	ALUGUEL EM MUNIQUE	82
4.2.1	<b>Análise exploratória dos dados</b>	82
4.2.2	<b>Ajuste do modelo</b>	87
4.2.3	<b>Diagnóstico</b>	90
4.3	AR POLUÍDO EM CHICAGO	93
4.3.1	<b>Análise exploratória dos dados</b>	94
4.3.2	<b>Ajuste do modelo</b>	96
4.3.3	<b>Diagnóstico</b>	99
4.4	REPRODUÇÃO DE AVES EM PORTUGAL	102
4.4.1	<b>Ajuste do modelo</b>	103
4.4.2	<b>Diagnóstico</b>	104
5	<b>CONCLUSÕES E TRABALHOS FUTUROS</b>	107
	<b>REFERÊNCIAS</b>	108
	<b>APÊNDICE A – CÓDIGOS COMPUTACIONAIS</b>	111
	<b>ANEXO A – DERIVAÇÃO DOS ESQUEMAS DE PERTURBAÇÃO</b>	122

## 1 INTRODUÇÃO

Em geral, para descrever algum fenômeno da natureza e/ou social, considera-se algum aspecto de interesse da população em estudo como a média, variância ou algum outro parâmetro. Neste caso, busca-se alguma lei matemática capaz de modelar a influência sistemática de variáveis que se supõem ter algum efeito sobre tais parâmetros. Desta forma, assumindo como parâmetro de interesse a média das variáveis resposta, então para relacionar o efeito que determinadas variáveis explicativas tem sobre este parâmetro, pode-se considerar metodologias, tais como: regressão paramétrica, utilizada quando se sabe, mesmo que aproximadamente, algum argumento teórico que indique qual deve ser a verdadeira relação entre a média da variável resposta e as variáveis explicativas; a regressão não paramétrica, utilizada quando, por outro lado, esta relação se dá mediante a funções que se adequem aos dados e não parâmetros; ou, como uma combinação destas metodologias, a regressão semi-paramétrica.

Neste contexto, os modelos parcialmente lineares generalizados (MPLAGs) representam uma classe de modelos de regressão semi-paramétricas que tentam estimar relações complexas entre a variável resposta e as covariáveis. A singularidade desses modelos é a sua flexibilidade, visto que não há necessidade de estabelecer antecipadamente o relacionamento entre as variáveis. Nesta classe de modelos, os dados são os que determinam a forma do relacionamento. Isso tornou os MPLAGs uma ferramenta fundamental em qualquer área de pesquisa, pois possibilitou análises estatísticas mais precisas, proporcionando vislumbrar caminhos com melhores resultados.

Dentre os principais trabalhos relacionados aos modelos aditivos generalizados, convém destacar as referências de Green e Silverman (1994). Estes autores assumem como preditor, uma soma de funções suaves de covariáveis quando a distribuição da variável resposta pertence a família exponencial. Green e Yandell (1985) apresentam o conceito de estimativas de máxima verossimilhança penalizada. Aqui, os autores também discutem critérios de seleção dos parâmetros de suavização para os modelos de regressão semi-paramétricos utilizando *splines*. Eilers e Marx (1996) estudam a estimação não paramétrica por meio de *B-splines* e *P-splines*. Além destes, Fahrmeir *et al.* (2013) e Hodges (2016) também trazem um desenvolvimento completo de modelos aditivos, principalmente no que tange a modelos mistos. Ademais, recomenda-se fortemente o livro de Wood (2017). Nele, pode-se encontrar, além de uma descrição detalhada da teoria que envolve esses modelos, uma excelente parte com aplicações no R.

Porém, conforme Cordeiro e Demétrio (2008), mesmo tendo passado por uma criteriosa escolha do modelo, pode ocorrer um resultado insatisfatório para o ajuste a um conjunto de observação. Assim sendo, uma etapa seguinte muito importante da análise de um ajuste consiste em seu diagnóstico. Segundo Paula (2013), esta etapa é responsável por verificar possíveis desvios em relação as suposições feitas para o modelo, em especial, no que diz respeito ao componente aleatório e, notadamente, para a parte sistemática, bem como a presença de observações destoantes com alguma interferência desproporcional ou inferencial nos resultados do ajuste. Em outras palavras, busca-se analisar a precisão dos resultados obtidos, tentando identificar possíveis desvios de suposições decorrentes de erros de especificação do modelo — como a escolha errônea para distribuição da variável resposta, da função de variância, função de ligação, ausência ou não do parâmetro de dispersão —, ou discrepâncias isoladas — oriundas de pontos presentes nas extremidades da amplitude de validade da variável explanatória, erros de digitação, algum fator não controlado (mas relevante) que influenciou a sua obtenção ou até mesmo multicolinearidade entre variáveis (correlação alta entre variáveis) — que podem exercer influência acentuada nas estimativas dos parâmetros do modelo (CORDEIRO; DEMÉTRIO, 2008, p.135).

Neste contexto, se faz necessário uso de técnicas para o diagnóstico do modelo, por meio da análise de resíduo, alavancagem e, sobretudo, por meio da análise de influência local, proposta por Cook (1978). Desta forma, Pulgar (2009), por exemplo, discute medidas de diagnósticos para modelos mistos aditivos semiparamétricos com erros de contornos elípticos. Emami (2016) estuda a influência local de pequenas perturbações nas estimativas obtidas pelo procedimento dos mínimos quadrados penalizados de Ridge. Manghi (2016) desenvolve técnicas de diagnóstico para os Modelos Parcialmente Lineares Aditivos Generalizados para dados correlacionados.

Dentre os trabalhos mais recentes no estudo dos MPLAGs, destaca-se o trabalho de Holanda (2018), com o título: modelos lineares parciais aditivos generalizados com suavização por meio de *P-splines*. Neste trabalho é apresentado um estudo sobre os modelos parcialmente lineares aditivos generalizados utilizando *P-splines* para descrever a relação da variável resposta com as variáveis explicativas contínuas, sendo que os *P-splines* são uma combinação entre o método de suavização *B-splines* e uma penalização de diferenças. Contudo, não foi desenvolvida a análise de diagnóstico para este modelo, bem como não há nada na literatura a respeito desse tema, sendo portanto uma lacuna a ser preenchida. Sendo assim, é oportuno desenvolver um estudo a fim de preencher tal lacuna existente. Neste caso, esta pesquisa tem a pretensão de

debruçar nas técnicas de diagnóstico por meio do estudo de resíduos, alavancagem e influência local baseada na proposta de Cook (1978), para esta classe de modelos.

Portanto, justifica-se esta pesquisa tendo em vista a necessidade de ampliar e aprofundar os estudos sobre procedimentos para o diagnóstico dos modelos parcialmente lineares aditivos generalizados. De forma mais precisa, será dedicado consideráveis esforços no que tange a análise de influência local, bem como apresentar essa técnica, mostrando seu uso em um contexto de pesquisa teórico e prático. Dessa forma, os objetivos específicos são, além do desenvolvimento teórico para análise de diagnóstico por meio da técnica de influência local, fornecer um suporte computacional no R. Além disso, buscar exemplos práticos para ilustrar a teoria desenvolvida.

## 1.1 ORGANIZAÇÃO DA DISSERTAÇÃO

Esta dissertação possui 4 capítulos e estão organizados da seguinte forma. No Capítulo 1 é apresentada uma introdução que visa descrever o objetivo deste trabalho. Nela, também buscamos lembrar alguns resultados importantes na literatura que utilizaram a metodologia de análise de diagnóstico por meio da técnica de influência local desenvolvida por Cook (1978) em diferentes contextos de aplicação.

No Capítulo 2 os MPLAGs são definidos e o processo de estimação dos parâmetros é apresentado. Desta forma, enfatizamos as ideias por trás da especificação desses modelos, bem como apresentamos discussões sobre processos de estimação por meio do método de máxima verossimilhança penalizada, bem como discorremos sobre o método iterativo *backfitting*.

No Capítulo 3, visando à compreensão prévia das técnicas que envolvem a análise de diagnóstico para um modelo de regressão, apresentamos e desenvolvemos uma explanação das principais estruturas para realizar a análise de diagnóstico dos MPLAGs. Para isso, sugerimos os resíduos de Pearson e os resíduos quantílicos aleatorizados com o intuito de identificar possíveis *outliers*, bem como discutimos o uso da matriz de projeção em modelos semi-paramétricos para investigar pontos de alavanca, além de desenvolver medidas de diagnóstico sob o enfoque da técnica de influência local.

Por fim, no Capítulo 4 são apresentadas algumas aplicações a banco de dados reais, seguido de comentários e conclusões do estudo.

## 1.2 SUPORTE COMPUTACIONAL

Destaca-se que as análises dos dados apresentados neste trabalho foram realizadas por meio do ambiente de programação R, em sua versão 3.6.1 para a plataforma **Linux**, por meio de um Ambiente de Desenvolvimento Integrado (Integrated Development Environment- IDE) chamado RStudio. A linguagem e seus pacotes podem ser obtidos no endereço <http://www.r-project.org> de forma gratuita. Para mais detalhes, consultar Ihaka e Gentleman (1996) e Team (2019). Para redigir a presente dissertação, foi utilizado o sistema tipográfico  $\text{\LaTeX}$  através da IDE  $\text{\TeX}$ Studio, obtida gratuitamente pelo endereço <https://www.texstudio.org>.

## 1.3 CONCEITOS PRELIMINARES

Em muito casos, a relação funcional entre a média da variável resposta e as variáveis explicativas se dá por meio de formas mais complexas que uma simples estrutura linear. Considere por exemplo o modelo de Michaelis-Menten para cinética enzimática. Tal modelo tem dois parâmetros e uma variável independente, relacionados por  $f$ , de tal forma que:

$$f(x, \boldsymbol{\beta}) = \frac{\beta_1 x}{\beta_2 + x}.$$

Esta função é não linear, pois não pode ser escrita como uma combinação linear dos  $\beta$ 's. Para estimar os parâmetros dessa função poderia-se utilizar o método dos mínimos quadrados.

A abordagem acima é conhecida na literatura como regressão paramétrica não linear (ver BEALE 1960, PAYANDEH 1983 e BATES; WATTS 1988, por exemplo). A escolha dessas funções é motivada, em sua grande parte, pelo conhecimento prévio do tipo de relação entre as variáveis. Contudo, nem sempre existe esse conhecimento. Neste caso, a alternativa é deixar que a amostra, por si mesma, apresente possíveis relações entre as variáveis. A essa metodologia dá-se o nome de regressão não paramétrica.

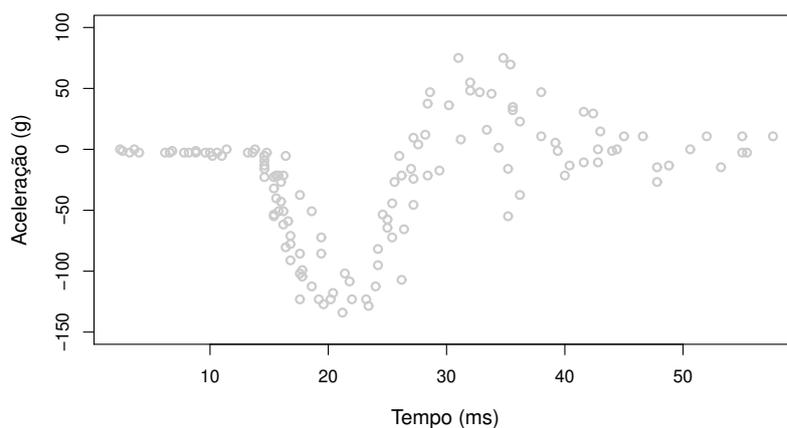
Neste contexto, para estimar  $f(\cdot)$  de forma não paramétrica, é necessário o uso de alguns métodos especiais. Dentre estes, por exemplo, estão os métodos baseados em séries ortogonais, método dos  $k$  vizinhos mais próximos (em inglês, *k-nearest neighbours*, *KNN*), regressão polinomial local, *splines*, *smoothing splines*, apenas para citar alguns. Estes métodos, assim como na regressão paramétrica não linear, são mais flexíveis em comparação com a regressão linear. Por captar melhor a relação funcional de  $f(\cdot)$ , oferecem um forte poder preditivo, bem como não há a necessidade de fazer hipóteses sobre esta relação, diferentemente da regressão paramétrica não linear. Por outro lado, são métodos cuja maior flexibilidade vem

atrelada a uma interpretabilidade mais complicada. Além disso, é necessário uma quantidade razoavelmente grande de observações a fim de obter uma aproximação para  $f(\cdot)$  suficientemente precisa. Mais detalhes sobre estes procedimentos podem ser encontrados em Izbicki e Santos (2019).

Contudo, este trabalho focará em apenas um desses métodos. Para estimar a relação funcional de  $f(\cdot)$ , será utilizado o método de suavização por meio de *P-splines*. Este procedimento é uma combinação entre o método de *B-splines* e uma penalização de diferenças aplicadas diretamente aos coeficientes adjacentes dos *B-splines*. Assim, é necessário inicialmente apresentar uma ideia introdutória sobre *splines*, bem como os conceitos do método de *B-splines* para posteriormente definirmos os *P-splines*.

### 1.3.1 Splines

Considere um problema em que o interesse é investigar o comportamento de capacetes durante o impacto na cabeça em acidentes simulados. Este exemplo será retomado com mais detalhes na **Seção 2.4.6**. Por hora, é suficiente saber que neste experimento foi considerado a aceleração da cabeça em unidades de  $g$ , para determinar força sentida pelo cérebro imediatamente após o impacto. A **Figura 1** mostra o diagrama de dispersão da aceleração contra o tempo dado em milissegundo (ms) para este experimento.



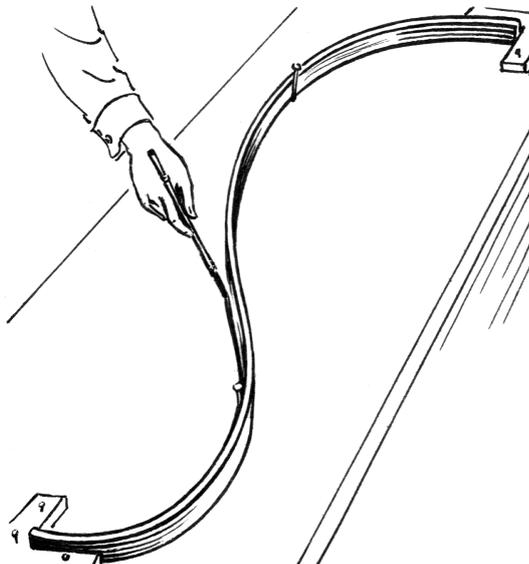
**Figura 1 – Gráfico de dispersão da aceleração contra o tempo.**

Como é possível observar, o comportamento da aceleração em função do tempo claramente não possui uma forma linear, tão pouco existe algum conhecimento *a priori* de como se dá essa relação. Neste caso, como brevemente exposto anteriormente, costuma-se utilizar

algum procedimento capaz de captar essa relação de modo que não seja feito qualquer suposição a seu respeito.

Contudo, talvez a forma mais usual seria tentar aproximar esse comportamento por meio de algum polinômio. Estes são bastante flexíveis e são bem simples de serem estimados. No entanto, como será visto mais adiante, a tentativa de ajustar curvas por meio de um único polinômio, principalmente os que possuem graus elevados, pode apresentar algumas dificuldades. Alternativamente pode-se usar *splines*, que são formados por pedaços de polinômios unidos em pontos chamados de nós e que são mais estáveis, sendo capazes de controlar melhor aspectos locais da relação funcional.

A palavra *spline* veio emprestada do desenho técnico e designa uma régua flexível que, antes dos computadores, era usada para desenhar curvas que passam por pontos pré-determinados com o intuito de auxiliar na etapa de delineamento de objetos, principalmente em projetos de engenharia, tais como cascos de navio, peças de avião, etc. É uma ferramenta que tenta, por meio de desenhos, comunicar visualmente como algo funciona ou é construído. Na **Figura 2** é apresentado um contorno obtido por meio de *splines*.



**Figura 2 – Contorno suave obtido com auxílio de um *splines*.**

Fonte: <https://fr.wikipedia.org/wiki/Spline>

Na matemática, o termo foi usado pela primeira vez em um artigo de 1946, pelo matemático romeno Isaac Jacob Schoenberg. A ideia é construir uma curva que passe por alguns pontos da amostra, chamados de nós, com o objetivo de aproximar tal curva do que se acredita ser aquela que represente adequadamente o comportamento desses dados.

De modo bastante sucinto, define-se *splines* como sendo curvas formadas por pedaços

de polinômios. Neste caso, ao invés de modelar um conjunto de observações por um único polinômio, escolhe-se pontos distintos no intervalo das observações (nós) e é definido um polinômio para cada intervalo, de forma a modelar curvas mais complexas por polinômios mais simples.

Existem algumas vantagens em se usar *splines*, dentre elas pode-se destacar a sua maior flexibilidade para o ajuste dos modelos se comparado ao modelo de regressão linear ou polinomial, a capacidade de modelar comportamentos atípicos dos dados, bem como a facilidade do ajuste quando determinados a quantidade e a posição dos nós.

Basicamente, o método consiste em substituir o vetor de entrada  $\mathbf{x}$ , por variáveis adicionais, que serão combinações de  $\mathbf{x}$  e então utilizar a aproximação linear nesse novo espaço. Assim, seja  $b_\ell(\mathbf{x})$  a  $\ell$ -ésima transformação de  $\mathbf{x}$ ,  $\ell = 1, 2, \dots, q$ , então  $f(\mathbf{x})$  assume a seguinte representação:

$$f(\mathbf{x}) = \sum_{\ell=1}^q b_\ell(\mathbf{x})\gamma_\ell, \quad (1.1)$$

para alguns valores dos parâmetros desconhecidos  $\gamma_\ell$ . Neste caso, obtém-se uma expansão linear de bases em  $\mathbf{x}$ . Portanto, após a definição de  $b_\ell(\mathbf{x})$ , o modelo passa a ser linear nessas novas variáveis gerados pelo vetor  $\mathbf{x}$  e conseqüentemente, os métodos de estimação dos parâmetros  $\gamma_\ell$  podem ser aplicados nesse novo espaço. Note que, ao contrário do que se imagina, técnicas não paramétricas evolvem a estimação de muitos parâmetros. Por essa razão, Eilers e Marx (1996) preferem "*overparametric*" *techniques* (técnicas de muitos parâmetros), embora tais parâmetros não tenham interpretação científica.

De acordo com Wood (2017, p. 120), costuma-se escolher como representação de  $b_\ell(\mathbf{x})$  funções capazes de aumentar a flexibilidade de  $f(\mathbf{x})$  e, conseqüentemente, do modelo. Um exemplo disso são os polinômios. Os polinômios são muito úteis quando o interesse reside nas propriedades de  $f$  na vizinhança de um ponto específico. Contudo, devido a sua natureza, tendem a distorcer a realidade em regiões remotas. Considerando os interesses estatísticos, o ajuste de curvas por meio de um único polinômio geralmente acarreta em problemas de multicolinearidade, isto é, associação de alta correlação entre as variáveis explicativas. Além disso, considerando polinômios com ordens elevadas, o risco de haver *overfitting* aumenta bastante. O que queremos é uma curva que se adéque a diferentes amostras vindas de uma mesma população e que seja capaz de captar adequadamente a variabilidade dos dados. Uma ideia para lidar com essas dificuldades é ajustar polinômios de menor grau e por partes.

Assim, supondo que  $\mathbf{x}$  seja unidimensional, então um polinômio por partes é obtido

separando o domínio de  $\mathbf{x}$  em intervalos conectados e definir  $f$  como um polinômio por intervalo. Neste caso, pode-se definir um polinômio por partes da seguinte maneira:

$$b_\ell(\mathbf{x}) = h_1(\mathbf{x})\mathbb{I}(\mathbf{x} < \xi_1) + h_i(\mathbf{x})\mathbb{I}(\xi_i < \mathbf{x} < \xi_{i+1}) + h_k(\mathbf{x})\mathbb{I}(\xi_k < \mathbf{x}), \quad i = 1, 2, \dots, k \quad (1.2)$$

em que  $h_i$  são polinômios conhecidos,  $\mathbb{I}(\cdot)$  uma função indicadora e  $k$  é a quantidade de nós. É importante destacar que não está sendo feito qualquer afirmação sobre a continuidade destas funções. Para garantir a continuidade, deve-se estabelecer restrições, de modo que  $f(\xi_i^-) = f(\xi_i^+)$ . Ou seja, a união dos polinômios se dá de forma suave, sem mudanças abruptas, e assim os valores da função nos nós serão os mesmo.

### 1.3.2 Splines lineares

Uma *spline* linear, também chamada de função linear segmentada, é definida por um conjunto de polinômios de grau um, unidos continuamente entre si através dos nós. Neste caso, para cada intervalo representado pelos nós, existe um segmento de reta.

Considerando  $\mathbf{x}$  um vetor ordenado de tamanho  $n$ . Assim, definindo  $n - 1$  nós,  $\xi_1, \xi_2, \dots, \xi_{n-1}$ , tem-se que o modelo mais simples de *spline* possui a seguinte forma:

$$f(x) = \beta_0 + \beta_1 x + \beta_2(x - \xi_1)_+ + \dots + \beta_n(x - \xi_{n-1})_+, \quad (1.3)$$

em que

$$(x - \xi_j)_+ = \begin{cases} (x - \xi_j), & (x - \xi_j) \geq 0 \\ 0, & (x - \xi_j) < 0, \end{cases} \quad (1.4)$$

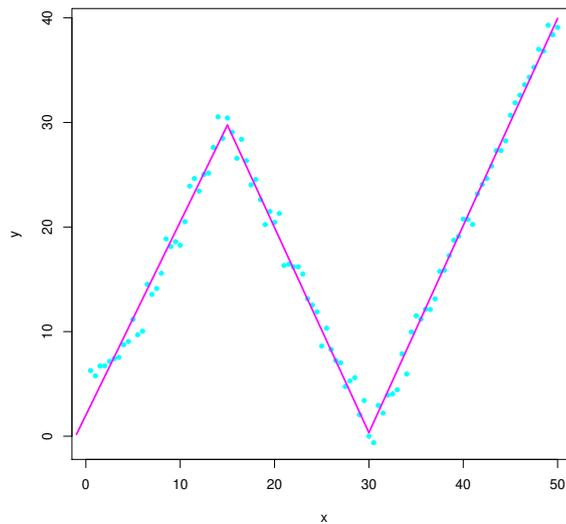
para  $j = 1, 2, \dots, n - 1$ . Pode-se notar que, conforme o valor no eixo  $x$  aumenta, novos termos são adicionados à função. Note que, se  $x < \xi_1$ , a equação irá configurar esse espaço. Isto é, entre o menor valor de  $x$  e  $\xi_1$  o valor da função será  $f(x) = \beta_0 + \beta_1 x$ . No entanto, quando  $x$  estiver entre os dois nós seguintes, a equação será  $f(x) = \beta_0 + \beta_1 x + \beta_2(x - \xi_1)$  e assim sucessivamente. Além disso, também é possível notar que nos intervalos entre os nós a *spline* linear é contínua. Porém, sua primeira derivada é descontínua nos pontos de união representado pelos nós.

Na **Figura 3** é apresentado um ajuste por meio de uma *spline* linear. Com base no gráfico de dispersão, foram escolhidos como ponto de quebra os nós  $\xi_1 = 15$  e  $\xi_2 = 30$ . Neste

modelo, tem-se a seguinte formulação:

$$f(x) = \begin{cases} \hat{\beta}_0 + \hat{\beta}_1 x, & \text{se } 0 < x \leq 15 \\ \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2(x - 15), & \text{se } 15 < x \leq 30 \\ \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2(x - 15) + \hat{\beta}_3(x - 30), & \text{se } 30 < x, \end{cases} \quad (1.5)$$

em que  $\hat{\beta}_0 = 2.02, \hat{\beta}_1 = 1.85, \hat{\beta}_2 = -3.81, \hat{\beta}_3 = 3.94$ .



**Figura 3 – Representação de um ajuste por meio de uma *spline* linear com nós em  $\xi_1 = 15$  e  $\xi_2 = 30$ .**

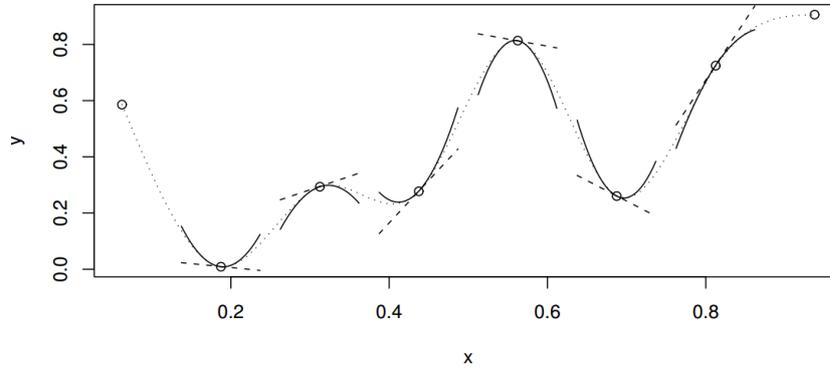
Note que uma regressão linear simples não seria capaz de estimar de forma adequada a relação funcional entre as variáveis. Isto é, observando o ajuste, a regressão por partes se sai muito melhor, dado sua proximidade aos pontos observados. O banco de dados para este exemplo de aplicação considerando *splines* lineares, bem como o código para a construção do gráfico são apresentados em **Código-fonte 1**.

### 1.3.3 *Splines* cúbicas

As *splines* cúbicas são as de menor ordem nas quais a descontinuidade nos nós são menos evidentes. Portanto, não há muitas razões capazes de motivar o uso de *splines* com graus maiores, a não ser que por algum motivo específico seja necessário mais derivadas suavizadas.

Essas *splines* podem ser dividida em dois tipos: as *splines* cúbicas restritas e *splines* cúbicas irrestritas. No caso em que a *spline* cúbica é do tipo restrita, então suas primeiras derivadas nos nós serão zero e suas caudas, isto é, a parte do polinômio antes do primeiro nó e após o último nó, serão modeladas por meio de funções lineares. A *spline* cúbica restrita é

também chamadas de *spline* natural. Caso contrario, será uma *spline* cúbica irrestrita. Na **Figura 4**, é dado uma ilustração de uma *spline* cúbica natural.



**Figura 4 – Aproximação de curvas por meio de splines.**

Fonte: Wood (2017, p. 122).

Assim, seja  $f$  uma função no intervalo  $[a, b]$ . Supondo que esse intervalo seja dividido em intervalos menores representado pelos nós, então define-se uma *spline* cúbica como sendo uma curva construída a partir de polinômios cúbicos, os quais são unidos de forma contínua nos nós. Ou seja, nos nós a primeira e a segunda derivadas da função serão contínuas. Além disso, cada seção do polinômio cúbico tem diferentes coeficientes, mas nos nós os valores das funções serão os mesmo.

Resumidamente, seja  $[\xi_i, \xi_{i+1}]$  o intervalo delimitado pelos nós e  $S$  uma *spline* cúbica utilizada para aproximar a curva  $f$ , então  $S$  possui as seguintes propriedades:

- i) para  $i = 1, 2, \dots, N$ , seja  $S(x) = s_i(x)$ , onde cada  $s_i$  é um polinômio cúbico;
- ii)  $s_i(\xi_{i-1}) = f(\xi_{i-1})$ , para  $i = 1, 2, \dots, N$ ;
- iii)  $s_i(\xi_i) = f(\xi_i)$ , para  $i = 1, 2, \dots, N$ ;
- iv)  $s'_i(\xi_i) = s'_{i+1}(\xi_i)$ , para  $i = 1, 2, \dots, N - 1$ ;
- v)  $s''_i(\xi_i) = s''_{i+1}(\xi_i)$ , para  $i = 1, 2, \dots, N$ .

Dessa forma, analogamente ao que foi feito para construir os *splines* lineares, pode-se construir as *splines* cúbicas. De modo geral, a função *spline* cúbica com  $k$  nós pode ser escrita como:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 (x - \xi_1)_+^3 + \dots + \beta_{k+1} (x - \xi_k)_+^3, \quad (1.6)$$

em que

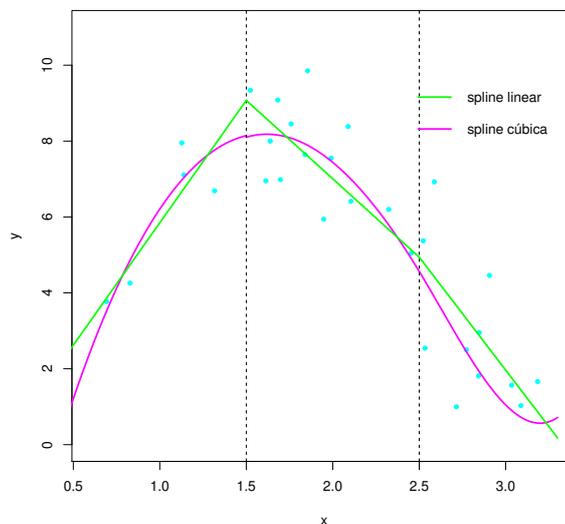
$$(x - \xi_k)_+^3 = \begin{cases} (x - \xi_k)^3, & x - \xi_k \geq 0 \\ 0, & x - \xi_k < 0. \end{cases} \quad (1.7)$$

A principal vantagem do uso de *splines* cúbicas em relação as *splines* lineares é o fato desta ser contínua até pelo menos a segunda derivada, inclusive nos nós. Isso possibilita captar pequenas ondulações no comportamento dos dados. A desvantagem é que os valores da função original não coincidem com os valores da *spline* **nas primeiras derivadas**, mesmo nos nós.

Como exemplo de aplicação, considere que os dados sejam gerados a partir da função:

$$y_i = 3 * \cos(3 + 2x_i) + 5 + \epsilon_i, \text{ para } i = 1, 2, \dots, 30,$$

em que  $\epsilon_i \sim \mathcal{N}(\mu = 0, \sigma^2 = 1)$  e  $x_i \sim \mathcal{U}(a = 0.5, b = 3.3)$ . A **Figura 5** mostra o resultado do ajuste. Os dados gerados por meio da função anterior para este exemplo de aplicação considerando *splines* cúbica, bem como o código para a construção do gráfico é apresentado no **Código-fonte 2**.



**Figura 5 – Representação de um ajuste por meio de *spline* linear e cúbica ambas com nós em  $\tilde{\zeta}_1 = 1.5$  e  $\tilde{\zeta}_2 = 2.5$ .**

Contudo, Pulgar (2009) mostra que as *splines* cúbicas não restritas costumam apresentar comportamento inadequado nas caudas, sendo estas, muito suscetíveis a mudanças. Portanto, por convenção, as funções que vem antes e depois do primeiro nó são forçadas a terem formato linear.

### 1.3.4 B-splines

Dentre todas as funções que são contínuas em  $[a, b]$  e interpola  $\{x_i, y_i\}$ , a *spline* cúbico é o interpolante mais suave possível através de qualquer conjunto de dados. Este resultado poder ser consultado em Wood (2017, p. 142-143). Existem várias maneiras de se representar uma *spline* cúbico. Uma delas é a representação B-*spline* proposta por De Boor (1978). Assim, neste caso, tomando  $b_\ell(\cdot) = B_\ell^m(\cdot)$ , define-se (1.1) como:

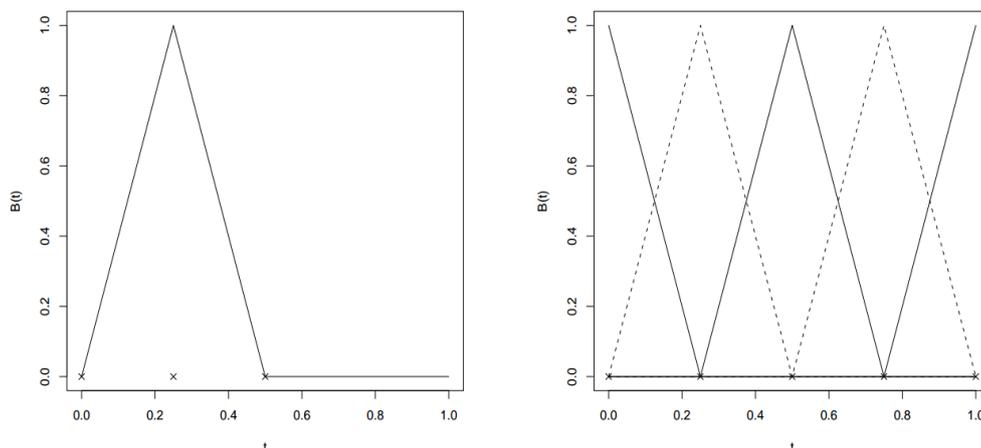
$$f(t_i) = \sum_{\ell=1}^q B_\ell^{(m)}(t_i) \gamma_\ell, \quad (1.8)$$

sendo os B-*splines*,  $B_\ell^{(m)}(t_i)$ , definidos como:

$$B_\ell^{(m)}(t_i) = \frac{t_i - \xi_\ell}{\xi_{\ell+m+1} - \xi_\ell} B_\ell^{(m-1)}(t_i) + \frac{\xi_{\ell+m+2} - t_i}{\xi_{\ell+m+2} - \xi_{\ell+1}} B_\ell^{(m-1)}(t_i) \text{ e}$$

$$B_\ell^{(0)}(t_i) = \begin{cases} 1 & \text{se } \xi_\ell \leq t_i < \xi_{\ell+1} \\ 0 & \text{c.c.} \end{cases} \quad (1.9)$$

em que  $t_i$  é uma variável que contribui de forma não linear no modelo, para  $i = 1, 2, \dots, n$ ,  $\xi_1, \xi_2, \dots, \xi_z$  são  $z$  nós equidistantes, sendo que  $m = 1, 2, \dots$  é o grau dos B-*splines* e  $z = q + m + 1$ . De acordo com Holanda (2018), a base de De Boor é capaz de controlar aspectos locais da relação funcional através do número e da posição dos nós, bem como a quantidade de parâmetros a serem estimados.

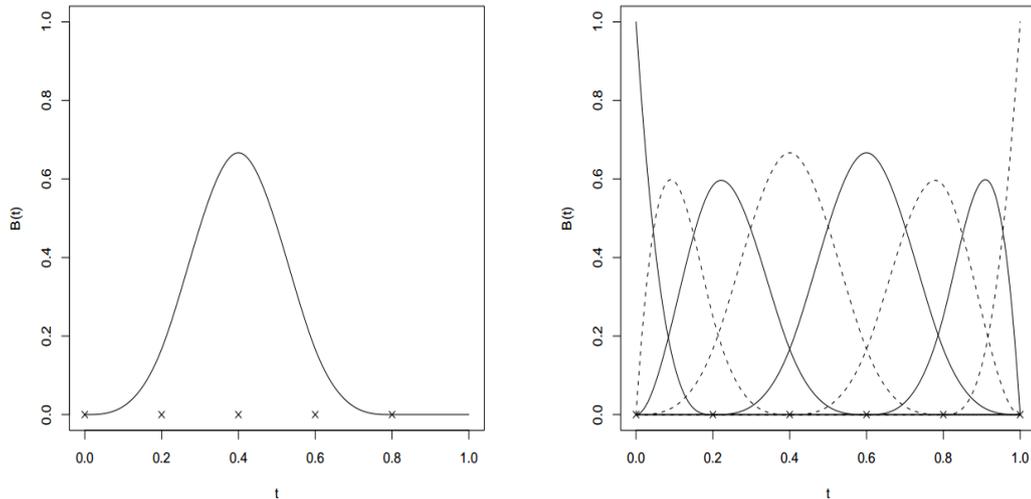


**Figura 6 – Ilustração de um B-*spline* de grau 1 isolado com três nós posicionados em "x".**

Fonte: Souza (2008, p. 35).

A **Figura 6** mostra dois gráficos com B-*splines* de grau 1, ambos baseados em três nós. Na parte direita da supra citada figura, encontramos todos os B-*splines* de grau 1 no intervalo

$[0, 1]$  com nós em  $\{0; 0.25; 0.5; 0.75; 1\}$ . Evidentemente, um número maior de *B-splines* podem ser construídos. Para isto, basta que sejam inseridos mais nós. Por sua vez, à direita da **Figura (7)** tem-se um *B-spline* de grau 2. Neste caso, a curva é formada por três pedaços de polinômios quadráticos, unidos por dois nós internos. À esquerda da **Figura 7** o leitor pode observar todos os possíveis *B-splines* de grau 2 levando em conta os mesmo nós.



**Figura 7 – Ilustração de um *B-splines* de grau 2 isolado com cinco nós posicionados em "x".**

Fonte: Souza (2008, p. 35).

Novamente é interessante observar como os polinômios se encaixam, tanto na ordem como também em suas primeiras e segundas derivadas (mas não as terceiras derivadas), o que evidencia o caráter contínuo, bem como o aspecto suave apresentado pela curva. Além disso, segundo Eilers e Marx (1996, p. 90), um *B-spline* de ordem  $m$  consiste em  $m$  pedaços de polinômios, cada um de grau  $m - 1$ . Estes pedaços de polinômios se juntam em  $m - 1$  nós internos. Nos pontos de união, as derivadas de ordem até  $m - 2$  são contínuas. Além disso, o *B-spline* é positivo no domínio abrangido por  $m + 1$  nós, mas fora disto ele é zero.

### 1.3.5 *P-splines*

De acordo com Eilers e Marx (1996), os *B-splines* possuem suavização influenciada unicamente pelo número de nós que será utilizado. Muitos nós podem levar a um *overfitting*, isto é, um super ajuste dos dados. Assim, a curva ajustada se adapta muito bem a uma amostra proveniente de alguma população, mas será pouco representativa da variabilidade dos dados. Por outro lado, uma quantidade muito pequena de nós pode acarretar problemas de *underfitting*, que

ocorre quando a curva é muito suave e não se adapta bem se quer aos dados que se dispõem. Eilers e Marx (1996) afirmam que a há intensas discussões entorno da escolha dos nós (ver FRIEDMAN; SILVERMAN, 1989 e KOOPERBERG; STONE, 1992). Contudo, as pesquisas ainda não levaram a nenhum esquema que fosse considerado pelos autores como uma forma atraente.

Para contornar essa dificuldade, O’Sullivan (1986), conforme citado por Eilers e Marx (1996), propõe usar um número relativamente grande de nós. Para controlar o ajuste da curva, é utilizada uma penalidade semelhante à pioneira para *splines* de suavização proposta por Reinsch (1967). À vista disso, ao invés de definir  $f(t_i) = y_i$ , é melhor tratar  $f(t_i)$  como  $n$  parâmetros livres da *spline* cúbica e estimá-la a fim de minimizar

$$\sum_{i=1}^n \{y_i - f(t_i)\}^2 + \alpha \int f''(t)^2 dt, \quad (1.10)$$

em que  $\alpha$  é um parâmetro de suavização usado para controlar o peso relativo a ser dado aos objetivos conflitantes de combinar os dados e produzir uma função  $f$  suave. Nota-se que, portanto, a função objetivo (1.10) trata-se, na verdade, de uma penalização na função objetivo do método dos mínimos quadrados.

A ideia de utilizar o método de estimação por mínimos quadrados penalizado pela integral do quadrado da segunda derivada de  $f$  tornou-se o padrão em grande parte da literatura (ver, por exemplo, EUBANK, 1988, WAHBA, 1990 e GREEN; SILVERMAN, 1994). Contudo, Eilers e Marx (1996) sugerem usar como penalização diferenças simples aplicadas diretamente aos parâmetros adjacentes dos *B-splines*. Os autores argumentam que essa combinação, a qual chamaram de *P-splines*, apresenta propriedades interessantes, tais como: não apresentam efeitos de contorno (ou fronteira); são uma extensão direta de modelos de regressão linear (generalizados), conservam momentos (média, variância) dos dados e têm curvas polinomiais ajustadas como limites, além de serem computacionalmente mais baratos. Assim, considerando (1.10), o interesse passa a ser minimizar:

$$\sum_{i=1}^n \{y_i - f(t_i)\}^2 + \alpha \boldsymbol{\gamma}^\top \mathbf{P}^d \boldsymbol{\gamma}, \quad (1.11)$$

em que  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_q)^\top$  são os parâmetros do *B-spline* e  $\mathbf{P}^d$  é uma matriz obtida a partir de um operador de diferenças entre os parâmetros adjacentes em  $\boldsymbol{\gamma}$ . Por meio de  $\mathbf{P}^d$  é possível controlar a dependência entre intervalos adjacentes (no domínio de  $t$ ) na estimação de  $f(t)$ . Para

mais detalhes sobre o uso dos *P-splines*, deve-se consultar Eilers e Marx (1996) e Eilers, Marx e Maria (2015).

Nos próximos capítulos será apresentada uma classe de modelos semi-paramétricos, definidos por um preditor que é modelado tanto por parâmetros como por funções suaves, tais como (1.1). Com isso, serão discutidas algumas especificações relacionadas as condições de estimação dos parâmetros, bem como a obtenção dos graus de liberdade. Além disso, será discutido procedimentos para realização de testes de hipóteses e intervalos de confiança.

## 2 MODELOS PARCIALMENTE ADITIVOS GENERALIZADOS

### 2.1 INTRODUÇÃO

Os modelos lineares constituem a forma mais simples de se estabelecer uma relação funcional entre a variável resposta e as variáveis explicativas. Esse modelo é especialmente aplicado quando a variável resposta é contínua e apresenta distribuição próxima a normalidade. Em razão dessa simplicidade, Paula (2013) lembra que durante muitos anos essa foi a principal abordagem para descrever a maioria dos fenômenos. Porém, como exemplificado por Fahrmeir *et al.* (2013), em várias situações práticas as suposições para garantir um bom uso do modelo de regressão linear clássico não são justificadas para estabelecer uma análise adequada, como por exemplo, nos casos em que a variável resposta era discreta ou, de modo mais geral, pertencente a um sub conjunto dos números reais, ou quando é necessário considerar a heteroscedasticidade.

Assim, na tentativa de contornar essas dificuldades, muitos pesquisadores costumavam aplicar alguma transformação nos dados com o objetivo de obter a normalidade, linearidade e constância da variância. Porém, essa abordagem quase sempre apresentava algumas limitações, além de tornar difícil a interpretação dos resultados. Neste sentido, surgiram na literatura alternativas mais interessantes que visavam superar essas dificuldades específicas. A regressão logística, por exemplo, foi desenvolvida no século XIX para lidar com problemas em que a variável resposta era dicotômica. Contudo, a proposta mais inovadora é atribuída a Nelder e Wedderburn (1972). Eles mostraram que muitas técnicas estatísticas, tais como regressão logística, regressão para dados de contagem, dentre outras, que antes eram estudadas separadamente, fazem parte de uma única classe de modelos de regressão. Os autores denominaram essa classe de **Modelos Lineares Generalizados** (MLG). Esse conjunto de métodos e procedimentos possibilitou um leque bem mais diversificado para a distribuição da variável resposta  $y$ . No entanto, essa abordagem, por si só, ainda não é suficiente para tratar casos em que o efeito das covariáveis é não-linear.

Para tanto, como discutido em seções anteriores, existem algumas metodologias que podem ser levadas em consideração: a primeira seria considerar uma regressão paramétrica não linear. Para isso, seria necessário saber, pelo menos de forma aproximada, qual deve ser a verdadeira relação não linear entre a média da variável resposta e as variáveis explicativas. A segunda seria deixar que os próprios dados indiquem qual deve ser essa relação. Nessa perspectiva, considera-se o uso de funções especiais que aproximem o comportamento não

linear das variáveis sem que, para isso, seja feito qualquer suposição a respeito da verdadeira relação entre a média da variável resposta e as variáveis explicativas. Neste caso, há um maior poder preditivo do modelo. Porém, há também uma perda na interpretabilidade, bem como fazer inferência a respeito dos parâmetros torna-se impossível.

Contudo, existe um interesse na comunidade estatística em entender como a resposta é afetada pela variação dos preditores ou identificar os mais importantes na relação entre a resposta e cada um deles. Foi baseado nessa perspectiva que surgiram os Modelos Parcialmente Lineares Aditivos Generalizados (MPLAG). Esta abordagem é uma forma de encontrar um equilíbrio entre interpretabilidade e flexibilidade. Na sequência, esta classe de modelos será apresentada com mais detalhes.

## 2.2 ESPECIFICAÇÃO DO MODELO

Os MPLAGs fazem parte de uma família de modelos que considera, além das funções paramétricas adotadas nos modelos lineares generalizados, uma soma de funções suaves de covariáveis a fim de relaxar a suposição de linearidade. Portanto, para definir os MPLAGs, considera-se a distribuição de  $y_i$ , para  $i = 1, 2, \dots, n$ , dada como:

$$f(y_i; \theta_i, \phi) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) \right\}, \quad (2.1)$$

de modo que, sob certas condições de regularidade, o valor esperado e a variância de  $y_i$  são dados, respectivamente, por  $\mathbb{E}(y_i) = \mu_i$  e  $\text{Var}(y_i) = \phi V_i$ , em que  $\mu_i = \partial b(\theta_i) / \partial \theta_i$  e  $V_i = \partial^2 b(\theta_i) / \partial \theta_i^2$ . O parâmetro natural (ou canônico)  $\theta_i$  pode ser expresso como  $\theta_i = \int V_i^{-1} d\mu_i = q(\mu_i)$ . Além disso, considere  $g(\cdot)$  uma função de ligação monótona e duas vezes diferenciável, tal que:

$$g(\mu_i) = \eta_i = \sum_{j=1}^p x_{ji} \beta_j + \sum_{k=1}^r f_k(t_{ik}), \quad i = 1, 2, \dots, n, \quad (2.2)$$

sendo  $f_k(\cdot)$  uma função (suave) desconhecida da variável  $t_k$  que contribui de forma não linear no modelo. Então tem-se que um MPLAG nada mais é do que uma extensão dos MLGs, quando existe apenas a parte paramétrica no preditor.

Supõe-se que  $f_k(\cdot)$  pertença a um espaço infinito-dimensional de funções duplamente contínuas e diferenciáveis e duplamente integráveis. Existem vários procedimentos para estimar esta função. Contudo, usar-se-á o método de suavização por meio de *B-splines* com função de penalização dos *P-splines*. Sendo assim, a função  $f_k(t_{ik})$  é expressa da seguinte forma:

$$f_k(t_{ik}) = \sum_{\ell=1}^{q_k} B_{k\ell}^{(m_k)}(t_{ik}) \gamma_{k\ell}, \quad i = 1, 2, \dots, n, \quad \text{e } k = 1, 2, \dots, r, \quad (2.3)$$

no qual  $q_k$  é a dimensão da base dos  $k$ -ésimo B-spline (ver De BOOR, 1978),  $z_k = q_k + m_k + 1$  é o número de nós,  $B_{k\ell}^{(m_k)}(t_{ik})$  denota  $\ell$ -ésimo componente da base do  $k$ -ésimo B-spline de ordem  $m_k + 1$  no ponto  $t_{ik}$  e  $\gamma_{\ell k}$  são os coeficientes a serem estimados, para  $i = 1, \dots, n$  e  $k = 1, \dots, r$ .

Com relação as estimativas dos parâmetros do modelo, toma-se como principal procedimento para tal finalidade uma adaptação do *Método da Máxima Verossimilhança*. Este método consiste em, baseado nos dados da amostra, encontrar a distribuição mais plausível, dentre todas aquelas definidas pelos possíveis valores de seus parâmetros, capaz de ter gerado a amostra disponível (COLOSIMO; GIOLO, 2006, p. 85). Isto é, as estimativas dos parâmetros do modelo serão obtidos pela maximização da função de verossimilhança (ou do logaritmo da função de verossimilhança). Contudo, é importante observar que a maximização direta, sem que tenha estabelecido qualquer restrição sobre a função não paramétrica  $f_k(t_{ik})$ , pode gerar um superajuste aos dados. Desta forma, é necessário incorporar uma penalização no logaritmo da função de verossimilhanças para controlar o ajuste.

### 2.2.1 Função de verossimilhança penalizada

Conforme observado anteriormente, quando se trabalha com uma matriz de base, a princípio não é necessária uma penalização no método da máxima verossimilhança (ou mínimos quadrados). Entretanto, neste caso, a suavidade da curva estimada depende do número de nós e da dimensão da base. Para não ter que se preocupar com esta escolha, deixa-se a base e o número de nós fixados e aplica-se a penalização. Esta penalização é função tanto do parâmetro de suavização como dos nós, sendo ela responsável por evitar o super ou sob ajuste dos dados.

Na literatura, é usual a penalização por meio da integral da segunda derivada da função  $f$ . Porém, uma alternativa mais atraente é usar como critério de penalização a função de penalização dos P-splines proposta por Eilers e Marx (1996), definida como:

$$\boldsymbol{\gamma}^\top \mathbf{P}^d \boldsymbol{\gamma}, \quad (2.4)$$

em que  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_q)^\top$ ,  $\mathbf{P}^d = (\Delta^d)^\top (\Delta^d)$ , de modo que  $\Delta^d$  é uma matriz com ordem  $(q - d) \times q$  obtida por meio do operador de diferenças  $\Delta^d$ , tal que  $\Delta^d \gamma_\ell = \Delta (\Delta^{d-1} \gamma_\ell)$  e  $\Delta \gamma_\ell = (\gamma_\ell - \gamma_{\ell-1})$ , sendo  $d$  a ordem dessas diferenças.

Portanto, seja  $d = 1$ , por exemplo, então tem-se que  $\Delta \gamma_\ell = (\gamma_\ell - \gamma_{\ell-1})$ . Assim, podemos obter um vetor coluna  $\Delta \boldsymbol{\gamma}$ , dado por  $\Delta \boldsymbol{\gamma} = (\Delta \gamma_2, \Delta \gamma_3, \dots, \Delta \gamma_q)^\top$ , de onde conseguimos  $\Delta \boldsymbol{\gamma} = (\gamma_2 - \gamma_1, \gamma_3 - \gamma_2, \dots, \gamma_q - \gamma_{q-1})^\top$ . Observe que  $\Delta \boldsymbol{\gamma}$  é um vetor de tamanho  $q - 1$ .

Desse modo, tem-se que a função de penalização será definida por:

$$\begin{aligned}
 (\Delta\boldsymbol{\gamma})^\top \Delta\boldsymbol{\gamma} &= \begin{bmatrix} \gamma_2 - \gamma_1 & \gamma_3 - \gamma_2 & \cdots & \gamma_q - \gamma_{q-1} \end{bmatrix} \begin{bmatrix} \gamma_2 - \gamma_1 \\ \gamma_3 - \gamma_2 \\ \vdots \\ \gamma_q - \gamma_{q-1} \end{bmatrix} \\
 &= (\gamma_2 - \gamma_1)^2 + (\gamma_3 - \gamma_2)^2 + \cdots + (\gamma_q - \gamma_{q-1})^2.
 \end{aligned}$$

Porém, é possível notar ainda que o vetor  $\Delta\boldsymbol{\gamma}$  pode ser reescrito como o produto de uma matriz por um vetor de parâmetros, de tal maneira que

$$\Delta\boldsymbol{\gamma} = \begin{bmatrix} \gamma_2 - \gamma_1 \\ \gamma_3 - \gamma_2 \\ \vdots \\ \gamma_q - \gamma_{q-1} \end{bmatrix} = \begin{bmatrix} -1 & 1 & 0 & 0 & \cdots \\ 0 & -1 & 1 & 0 & \cdots \\ 0 & 0 & 0 & -1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_q \end{bmatrix}.$$

Desta forma, a matriz  $\Delta$ , neste caso, tem dimensão  $(q - 1) \times q$  e será dada por:

$$\Delta = \begin{bmatrix} -1 & 1 & 0 & 0 & \cdots \\ 0 & -1 & 1 & 0 & \cdots \\ 0 & 0 & 0 & -1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

de modo que também vale

$$\boldsymbol{\gamma}^\top \Delta^\top \Delta\boldsymbol{\gamma} = (\gamma_2 - \gamma_1)^2 + (\gamma_3 - \gamma_2)^2 + \cdots + (\gamma_q - \gamma_{q-1})^2.$$

Por exemplo, tomando  $q = 4$ , então teremos  $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \gamma_3, \gamma_4)^\top$  e, conseqüentemente  $\Delta\boldsymbol{\gamma} = (\gamma_2 - \gamma_1, \gamma_3 - \gamma_2, \gamma_4 - \gamma_3)^\top$ . Ou seja:

$$\boldsymbol{\gamma}^\top \Delta^\top \Delta\boldsymbol{\gamma} = (\gamma_2 - \gamma_1)^2 + (\gamma_3 - \gamma_2)^2 + (\gamma_4 - \gamma_3)^2.$$

Contudo, também temos que:

$$\Delta\boldsymbol{\gamma} = \begin{bmatrix} \gamma_2 - \gamma_1 \\ \gamma_3 - \gamma_2 \\ \gamma_4 - \gamma_3 \end{bmatrix} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \\ \gamma_4 \end{bmatrix},$$

de onde obtemos

$$\Delta = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix},$$

resultando em uma matriz  $\mathbf{P} = \Delta^\top \Delta$  dada por:

$$\mathbf{P} = \begin{bmatrix} -1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

$$\mathbf{P} = \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix},$$

o que também nos leva à função de penalização a seguir

$$\boldsymbol{\gamma}^\top \Delta^\top \Delta \boldsymbol{\gamma} = (\gamma_2 - \gamma_1)^2 + (\gamma_3 - \gamma_2)^2 + (\gamma_4 - \gamma_3)^2,$$

já definida anteriormente.

Adotando a ordem  $d = 2$ , tem-se que  $\Delta^2 \boldsymbol{\gamma} = (\Delta^2 \gamma_3, \Delta^2 \gamma_4, \dots, \Delta^2 \gamma_q)^\top$ . Neste caso,  $\Delta^2 \gamma_\ell = \Delta(\Delta \gamma_\ell)$  em que

$$\begin{aligned} \Delta^2 \gamma_\ell &= \Delta(\gamma_\ell - \gamma_{\ell-1}) \\ &= \Delta \gamma_\ell - \Delta \gamma_{\ell-1} \\ &= (\gamma_\ell - \gamma_{\ell-1}) - (\gamma_{\ell-1} - \gamma_{\ell-2}) \\ &= \gamma_\ell - 2\gamma_{\ell-1} + \gamma_{\ell-2}, \end{aligned}$$

de onde conseguimos  $\Delta^2 \boldsymbol{\gamma} = (\gamma_3 - 2\gamma_2 + \gamma_1, \gamma_4 - 2\gamma_3 + \gamma_2, \dots, \gamma_q - 2\gamma_{q-1} + \gamma_{q-2})^\top$ . Assim,  $\Delta^2 \boldsymbol{\gamma}$  é um vetor de tamanho  $q - 2$ . Além disso, de forma análoga ao que foi desenvolvido anteriormente para  $d = 1$ , é possível reescrever este vetor como

$$\Delta^2 \boldsymbol{\gamma} = \begin{bmatrix} \gamma_3 - 2\gamma_2 + \gamma_1 \\ \gamma_4 - 2\gamma_3 + \gamma_2 \\ \vdots \\ \gamma_q - 2\gamma_{q-1} + \gamma_{q-2} \end{bmatrix} = \begin{bmatrix} 1 & -2 & 1 & 0 & \cdots \\ 0 & 1 & -2 & 1 & \cdots \\ 0 & 0 & 1 & -2 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_q \end{bmatrix}.$$

Logo, a matriz  $\Delta^2$  obtida com operador de diferenças de segunda ordem assume a seguinte forma:

$$\Delta^2 = \begin{bmatrix} 1 & -2 & 1 & 0 & \cdots \\ 0 & 1 & -2 & 1 & \cdots \\ 0 & 0 & 1 & -2 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

que neste caso  $\Delta^2$  é uma matriz de dimensão  $(q-2) \times q$ . Consequentemente, também leva à função de penalização será dada por:

$$\boldsymbol{\gamma}^\top (\Delta^2)^\top (\Delta^2) \boldsymbol{\gamma} = (\gamma_1 - 2\gamma_2 + \gamma_3)^2 + (\gamma_2 - 2\gamma_3 + \gamma_4)^2 + \cdots + (\gamma_{q-2} - 2\gamma_{q-1} + \gamma_q)^2.$$

Por exemplo, tomando novamente  $q = 4$ , então teremos  $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \gamma_3, \gamma_4)^\top$  e, portanto  $\Delta^2 \boldsymbol{\gamma} = (\gamma_3 - 2\gamma_2 + \gamma_1, \gamma_4 - 2\gamma_3 + \gamma_2)^\top$ . Sendo assim, a função de penalização será dada por:

$$\boldsymbol{\gamma}^\top \Delta^\top \Delta \boldsymbol{\gamma} = (\gamma_2 - \gamma_1)^2 + (\gamma_3 - \gamma_2)^2 + (\gamma_4 - \gamma_3)^2.$$

Contudo, note que:

$$\Delta^2 \boldsymbol{\gamma} = \begin{bmatrix} \gamma_3 - 2\gamma_2 + \gamma_1 \\ \gamma_4 - 2\gamma_3 + \gamma_2 \end{bmatrix} = \begin{bmatrix} 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \\ \gamma_4 \end{bmatrix},$$

de onde obtemos

$$\Delta^2 = \begin{bmatrix} 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix},$$

resultando em uma matriz  $\mathbf{P}^2 = (\Delta^2)^\top (\Delta^2)$  dada por:

$$\begin{aligned} \mathbf{P}^2 &= \begin{bmatrix} 1 & 0 \\ -2 & 1 \\ 1 & -2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & -2 & 1 & 0 \\ -2 & 5 & -4 & 1 \\ 1 & -4 & 5 & -2 \\ 0 & 1 & -2 & 1 \end{bmatrix}, \end{aligned}$$

É claro que, como observado, pode-se calcular a função de penalização diretamente sem o intermédio de  $\mathbf{P}^d$ . Porém, a matriz  $\mathbf{P}^d$  é necessária para o cálculo das estimativas dos parâmetros do modelo, como será desenvolvido a seguir. Além disso, os cálculos das medidas de diagnóstico que serão apresentadas no capítulo 3 também necessitam desta matriz.

Segundo Holanda (2018), a penalização *P-spline* é mais flexível, visto que independe do grau adotado para os *B-splines*. Desta maneira, é possível combinar qualquer ordem da penalidade com qualquer ordem das bases de De Boor. Além disso, ainda segundo a autora, pode-se utilizar outras penalidades a fim de obter algumas propriedades desejáveis, como, por exemplo: garantir a monotonicidade da função, para exigir uma forma convexa ou côncava ou para forçar os sinais das restrições.

Portanto, considere que  $\boldsymbol{\xi}^\top = (\boldsymbol{\beta}^\top, \boldsymbol{\gamma}_1^\top, \boldsymbol{\gamma}_2^\top, \dots, \boldsymbol{\gamma}_r^\top) \in \Xi \subseteq \mathbb{R}^s$ , em que  $s = p + \sum_{k=1}^r q_k$ . Isto posto, sob a suposição de independência, o estimador de  $\boldsymbol{\xi}$  será obtido como sendo o valor que maximiza a função de log-verossimilhança penalizada dada a seguir:

$$L_p(\boldsymbol{\xi}) = L(\boldsymbol{\xi}) + \frac{1}{2} \sum_{k=1}^r \alpha_k \boldsymbol{\gamma}_k^\top \mathbf{P}_k^d \boldsymbol{\gamma}_k, \quad (2.5)$$

em que  $L(\boldsymbol{\xi}) = \sum_{i=1}^n \{ \phi^{-1}[y_i \theta_i - b(\theta_i)] + c(y_i; \phi) \}$ ,  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_r)^\top$  é um vetor de parâmetros de suavização, de forma que  $\alpha_k > 0$ ,  $\boldsymbol{\gamma}_k$  são os coeficientes dos *B-splines* da  $k$ -ésima função não paramétrica,  $\mathbf{P}_k^d = (\Delta_k^d)^\top (\Delta_k^d)$  e  $\Delta_k^d$  é um operador de diferenças, tal que  $\Delta_k^d \gamma_{k\ell} = \Delta_k (\Delta_k^{d-1} \gamma_{k\ell})$  e  $\Delta_k \gamma_{k\ell} = (\gamma_{k\ell} - \gamma_{k(\ell-1)})$ , sendo  $d$  a ordem da diferença. Assim sendo, deve-se maximizar (2.5), resolvendo o seguinte sistema de equações:

$$U_p(\boldsymbol{\xi}) = \frac{\partial L_p(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = \mathbf{0}, \quad (2.6)$$

obtendo, assim, as estimativas dos parâmetros e sendo  $U_p(\boldsymbol{\xi})$  o *vetor score penalizado*.

À vista disto, tomando  $\mathbf{X}^* = (\mathbf{X}^\top, \mathbf{B}_1^\top, \dots, \mathbf{B}_r^\top)$ , em que  $\mathbf{X}$  é uma matriz  $n \times p$  formada pelos elementos  $x_{ji}$ ,  $\mathbf{B}_k$  é uma matriz de dimensão  $n \times q_k$  formada pelos componentes  $B_k^{(m_k)}(t_{ik})$  da base de De Boor, as equações de estimação de  $\boldsymbol{\xi}$  são dadas por:

$$U_p(\boldsymbol{\xi}) = \mathbf{X}^{*\top} \mathbf{W} \mathbf{D}^{-1} (\mathbf{y} - \boldsymbol{\mu}) - \mathbf{P}(\boldsymbol{\alpha}) = \mathbf{0}, \quad (2.7)$$

sendo  $\mathbf{W} = \mathbf{D} \mathbf{V}^{-1} \mathbf{D}$ , com  $\mathbf{D} = \text{diag} \{ \partial \mu_1 / \partial \eta_1, \dots, \partial \mu_n / \partial \eta_n \}$ ,  $\mathbf{V} = \text{diag} \{ V_1, \dots, V_n \}$  e  $\mathbf{P}(\boldsymbol{\alpha}) = \text{blocodiag} \{ \mathbf{0}_{pp}, \mathbf{P}^d(\boldsymbol{\alpha}) \}$ , no qual  $\mathbf{0}_{pp}$  é uma matriz de zeros de dimensão  $p \times p$ ,  $\mathbf{P}^d(\boldsymbol{\alpha}) = \mathbf{K}(\boldsymbol{\alpha}) \boldsymbol{\gamma}$  com  $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_1^\top, \dots, \boldsymbol{\gamma}_r^\top)^\top$  e  $\mathbf{K}(\boldsymbol{\alpha}) = \text{diag} \{ \alpha_1 \mathbf{P}_1^d, \dots, \alpha_r \mathbf{P}_r^d \}$ .

Conforme observado por Vanegas e Paula (2016), a matriz  $\mathbf{X}^*$  pode não ser singular. Dessa forma, Wood (2017) propõe que cada suavidade esteja sujeita a uma restrição de

centralização, aplicando uma reparametrização em termos de  $q_k - 1$  novos parâmetros  $\boldsymbol{\tau}_k$ , de tal forma que  $\boldsymbol{\gamma}_k = \mathbf{Z}_k \boldsymbol{\tau}_k$ , em que  $\mathbf{Z}_k$  é uma matriz  $q_k \times (q_k - 1)$  obtida por meio da decomposição **QR**. A seguir, o procedimento será detalhado.

### 2.2.2 Restrição de identificabilidade

A incorporação de termos suaves no ajuste de um modelo aditivo generalizado geralmente acarreta em problemas de identificabilidade (mesmo para o modelo sem intercepto). De acordo com Wood (2017), isso ocorre porque os termos suaves não são ortogonais em relação a interceptação, bem como é impossível estimar um intercepto para cada suavizador. Desta forma, o autor propõe que cada função não paramétrica esteja sujeita a uma restrição de identificabilidade. Para isso, dentre as alternativas existentes, uma restrição adequada seria impor que a soma (ou a média) dos elementos de  $\mathbf{f}_k$  seja zero, i.e.,  $\sum_{i=1}^n f_k(t_{ik}) = 0$ . Deste modo, a ortogonalidade do termo suave para o intercepto estará garantida.

Com isso, note que  $\mathbf{f}_k$  pode ser reescrito como  $\mathbf{f}_k = \mathbf{B}_k \boldsymbol{\gamma}_k$ . Assim sendo, pode-se redefinir  $\sum_{i=1}^n f_k(t_{ik}) = 0$  matricialmente como:

$$\mathbf{1}^\top \mathbf{B}_k \boldsymbol{\gamma}_k = 0,$$

em que  $\mathbf{1}^\top$  é um vetor de 1's. Especificamente, o autor sugere encontrar uma matriz  $\mathbf{Z}_k$  com  $q_k$  linhas por  $q_k - 1$  colunas ortogonais que satisfaz a seguinte condição:

$$\mathbf{1}^\top \mathbf{B}_k \mathbf{Z}_k = 0,$$

de modo que  $\boldsymbol{\gamma}_k = \mathbf{Z}_k \boldsymbol{\tau}_k$ . Diante disso, obtém-se uma nova matriz modelo para o  $k$ -ésimo termo,  $\tilde{\mathbf{B}}_k = \mathbf{B}_k \mathbf{Z}_k$ , tal que  $\mathbf{f}_k = \tilde{\mathbf{B}}_k \boldsymbol{\tau}_k$  automaticamente satisfaz a condição de centralização.

A matriz  $\mathbf{Z}_k$  não é formada explicitamente. Contudo, uma abordagem geral e elegante, porém bem simples, para obtê-la é por meio da decomposição QR (ou fatoração QR) de  $\mathbf{C}_k^\top$ , em que  $\mathbf{C}_k = \mathbf{1}^\top \mathbf{B}_k$  (ver Wood (2017, p. 47), seção 1.8.1). A decomposição QR permite fatorar uma matriz  $\mathbf{A}$  em um produto  $\mathbf{A} = \mathbf{QR}$  de uma matriz ortogonal  $\mathbf{Q}$  e uma matriz triangular superior  $\mathbf{R}$ . Uma matriz  $\mathbf{Q}$  é dita ortogonal se suas colunas são vetores unitários ortogonais, isto é, se  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{Q} \mathbf{Q}^\top = \mathbf{I}$ . De fato, encontra-se a decomposição QR de  $\mathbf{C}_k^\top$  e as últimas  $q_k - 1$  colunas da matriz ortogonal definem  $\mathbf{Z}_k$ .

É importante observar que toda matriz  $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ) possui uma fatoração QR. Se for de posto completo, então ela é única. Esta decomposição pode ser obtida por meio do algoritmo de ortogonalidade de Gram-Schmidt. Contudo, tal procedimento não é

recomendado, pois problemas computacionais podem resultar em uma matriz  $\mathbf{Q}$  não ortogonal. Assim, costuma-se usar transformações adequadas de matrizes, dentre as quais encontra-se o método da triangulação de Householder (ou reflexão de Householder) e que possui a seguinte forma:

$$\mathbf{H} = \mathbf{I} - 2 \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top \mathbf{u}}, \quad (2.8)$$

em que  $\mathbf{I}$  é uma matriz identidade de dimensão  $m \times m$ ,  $\mathbf{u} = \mathbf{x} + \text{sign}(x_1)\alpha\mathbf{e}$ , sendo que  $\mathbf{x}$  é um vetor coluna de  $\mathbf{A}$ , tal que  $\|\mathbf{x}\| = \alpha$  com  $\|\cdot\|$  representando a norma euclidiana e  $\mathbf{e} = (1, 0, 0, \dots, 0)^\top$ .

Na triangulação de Householder, a decomposição QR funciona encontrando-se matrizes  $\mathbf{H}$  apropriadas de tal maneira que ao multiplicar essas matrizes pela matriz original  $\mathbf{A}$ , obtém-se uma matriz triangular superior. No entanto,  $\mathbf{Q}$  não é obtida explicitamente, sendo necessário repetir esse procedimento na matriz  $\mathbf{R}$  resultante um certo número de vezes até obter uma matriz triangular superior. Assim, a matriz  $\mathbf{Q}$  será um produto escalar de cada matriz  $\mathbf{H}$  formada sucessivamente. Isto é:

$$\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_{m-2} \mathbf{H}_{m-1}. \quad (2.9)$$

Seja, por exemplo, encontrar a de composição QR por meio da triangulação de Householder da matriz  $\mathbf{A}$  dada por:

$$\mathbf{A} = \begin{bmatrix} 12 & -51 & 4 \\ 6 & -16 & -68 \\ -4 & 24 & -41 \end{bmatrix}. \quad (2.10)$$

Primeiramente, é necessário obter uma reflexão que transforme a primeira coluna da matriz  $\mathbf{A}$ . Neste caso, temos que  $\mathbf{x}_1 = (12, 6, -4)^\top$ . Assim,  $\alpha_1 = \|\mathbf{x}_1\| = \sqrt{12^2 + 6^2 + (-4)^2} = 14$  e  $\mathbf{e}_1 = (1, 0, 0)$ . Consequentemente, obtemos  $\mathbf{u}_1 = (12, 6, -4)^\top + (14, 0, 0)^\top = (26, 6, -4)^\top$ . Portanto, aplicando a equação (2.8), tem-se que:

$$\mathbf{H}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 2 \frac{\begin{bmatrix} 26 \\ 6 \\ -4 \end{bmatrix} \begin{bmatrix} 26 & 6 & -4 \end{bmatrix}}{\begin{bmatrix} 26 \\ 6 \\ -4 \end{bmatrix} \begin{bmatrix} 26 & 6 & -4 \end{bmatrix}}$$

$$\begin{aligned}
&= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \frac{2}{728} \begin{bmatrix} 676 & 156 & -104 \\ 156 & 36 & -24 \\ -104 & -24 & 16 \end{bmatrix} \\
&= \begin{bmatrix} -6/7 & -3/7 & 2/7 \\ 3/7 & 82/91 & 6/91 \\ 2/7 & 6/91 & 87/91 \end{bmatrix}.
\end{aligned}$$

Portanto, multiplicando a matriz  $\mathbf{H}_1$  obtida anteriormente por  $\mathbf{A}$ , consegue-se o seguinte resultado:

$$\begin{aligned}
\mathbf{H}_1\mathbf{A} &= \begin{bmatrix} -6/7 & -3/7 & 2/7 \\ 3/7 & 82/91 & 6/91 \\ 2/7 & 6/91 & 87/91 \end{bmatrix} \begin{bmatrix} 12 & -51 & 4 \\ 6 & -16 & -68 \\ -4 & 24 & -41 \end{bmatrix} \\
&= \begin{bmatrix} 14 & -30 & 14 \\ 0 & -451/13 & -854/13 \\ 0 & 474/13 & -553/13 \end{bmatrix}.
\end{aligned}$$

Note, porém, que a matriz obtida ainda não é triangular superior. Para isso, é necessário que o elemento (3,2) da matriz  $\mathbf{H}_1\mathbf{A}$  seja zero. Assim, considere  $\mathbf{A}^{(1)}$  formada por (1,1) menor da matriz  $\mathbf{H}_1\mathbf{A}$ , então temos que:

$$\mathbf{A}^{(1)} = \begin{bmatrix} -451/13 & -854/13 \\ 474/13 & -553/13 \end{bmatrix}.$$

Agora, basta aplicar o método novamente. Neste caso, temos que  $\mathbf{x}_2 = (-451/13, 474/13)^\top$ .

Logo,  $\alpha_2 \approx 50.4$  e  $\mathbf{e}_2 = (1,0)$ . Note que o sinal que acompanha o primeiro elemento de  $\mathbf{x}_2$  é negativo. Consequentemente, obtemos  $\mathbf{u}_2 = (-451/13, 474/13)^\top - (50.4, 0)^\top = (-85.0, 474/13)^\top$ .

Com isso,

$$\begin{aligned}
\mathbf{H}_2^{(1)} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} 0.84 & -0.36 \\ -0.36 & 0.15 \end{bmatrix} \\
&= \begin{bmatrix} -0.68 & 0.72 \\ 0.72 & 0.68 \end{bmatrix} \\
\mathbf{H}_2 &= \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_2^{(1)} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.68 & 0.72 \\ 0 & 0.72 & 0.68 \end{bmatrix}.
\end{aligned}$$

A primeira coluna  $(1,0,0)^\top$  e a primeira linha  $(1,0,0)$  são adicionadas para resultar em uma matriz  $\mathbf{H}_2$  com dimensões  $3 \times 3$ , a fim de manter as mesmas dimensões da matriz  $\mathbf{A}$ . Então, da mesma forma que anteriormente, deve-se multiplicar  $\mathbf{H}_2$  por  $\mathbf{H}_1\mathbf{A}$ , para obter uma nova matriz  $\mathbf{R}$  que seja triangular superior. Desse modo, temos que:

$$\begin{aligned} \mathbf{H}_2\mathbf{H}_1\mathbf{A} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.68 & 0.72 \\ 0 & 0.72 & 0.68 \end{bmatrix} \begin{bmatrix} 14 & -30 & 14 \\ 0 & -451/13 & -854/13 \\ 0 & 474/13 & -553/13 \end{bmatrix} \\ &= \begin{bmatrix} 14 & -30 & 14 \\ 0 & 50.33 & 14.46 \\ 0 & 0 & -76.91 \end{bmatrix}. \end{aligned}$$

Finalmente a matriz resultante é triangular superior. Ou seja,  $\mathbf{H}_2\mathbf{H}_1\mathbf{A} = \mathbf{R}$ , de modo que  $\mathbf{Q} = \mathbf{H}_1\mathbf{H}_2$ . Assim, considerando três casas decimais de aproximação, temos que:

$$\begin{aligned} \mathbf{Q} = \mathbf{H}_1\mathbf{H}_2 &= \begin{bmatrix} 14 & -30 & 14 \\ 0 & -451/13 & -854/13 \\ 0 & 474/13 & -553/13 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & -0.68 & 0.72 \\ 0 & 0.72 & 0.68 \end{bmatrix} \\ &= \begin{bmatrix} -0.857 & 0.502 & -0.113 \\ -0.428 & -0.573 & 0.698 \\ 0.286 & 0.647 & 0.707 \end{bmatrix}. \end{aligned}$$

Com isso, observa-se que o uso de transformações de Household é inerentemente muito simples. Além disso, é numericamente estável, sendo portanto a melhor alternativa para a obtenção da decomposição QR.

Seja agora  $\mathbf{C} = \mathbf{1}^\top \mathbf{A}$ . Então isso implica que  $\mathbf{C}^\top = \mathbf{A}^\top \mathbf{1}$ . Nosso objetivo será encontrar a decomposição QR de  $\mathbf{C}^\top$  por meio da triangulação de Householder. Para tanto, note que:

$$\begin{aligned} \mathbf{C}^\top &= \begin{bmatrix} 12 & 6 & -4 \\ 51 & -16 & 24 \\ 4 & -68 & -41 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 14 \\ 59 \\ -105 \end{bmatrix}. \end{aligned}$$

Ou seja,  $\mathbf{C}^\top$  é um vetor coluna dado pela soma dos elementos contidos nas linhas da matriz  $\mathbf{A}$ . Desse modo,  $\mathbf{x} = (14, 59, -105)^\top$ . Logo,  $\alpha \approx 121.252$  e  $\mathbf{e}_1 = (1, 0, 0)$ . Assim sendo, obtemos  $\mathbf{u} = (14, 59, -105)^\top - (121.252, 0, 0)^\top = (135.252, 59, -105)^\top$  e, utilizando três casas decimais de aproximação, tem-se que:

$$\mathbf{H} = \begin{bmatrix} -0.115 & -0.487 & 0.866 \\ -0.487 & 0.788 & 0.378 \\ 0.866 & 0.378 & 0.328 \end{bmatrix}.$$

Com isso, pré multiplicando  $\mathbf{C}^\top$  por  $\mathbf{H}$ , obtemos:

$$\mathbf{HC}^\top = \begin{bmatrix} -0.115 & -0.487 & 0.866 \\ -0.487 & 0.788 & 0.378 \\ 0.866 & 0.378 & 0.328 \end{bmatrix} \begin{bmatrix} 14 \\ 59 \\ -105 \end{bmatrix} = \begin{bmatrix} -121.252 \\ 0 \\ 0 \end{bmatrix},$$

finalizando o procedimento ai mesmo. Logo,  $\mathbf{Q} = \mathbf{H}^\top$ . Porém, note que a matriz  $\mathbf{H}$  é simétrica. Neste caso,  $\mathbf{H} = \mathbf{H}^\top$  e, portanto,  $\mathbf{Q} = \mathbf{H}$ .

Como discutido anteriormente, a matriz  $\mathbf{Z}$  não é obtida explicitamente. Mas, realizando a decomposição QR de  $\mathbf{C}^\top$ , obtemos  $\mathbf{Z}$  como sendo a matriz formada pelas últimas colunas da matriz  $\mathbf{Q}$ , excluindo-se a primeira. Assim, temos que:

$$\mathbf{Z} = \begin{bmatrix} -0.487 & 0.866 \\ 0.788 & 0.378 \\ 0.378 & 0.328 \end{bmatrix}.$$

Desse modo, tomando o produto  $\mathbf{CZ}$ , obtém-se:

$$\mathbf{CZ} = \begin{bmatrix} 14 & 59 & -105 \end{bmatrix} \begin{bmatrix} -0.487 & 0.866 \\ 0.788 & 0.378 \\ 0.378 & 0.328 \end{bmatrix} = \begin{bmatrix} 0 & 0 \end{bmatrix},$$

como era de se esperar. Além disso, embora seja redundante, considere tomar a multiplicação da matriz  $\mathbf{Z}$  pela matriz  $\mathbf{A}$ , então obtém-se o seguinte resultado:

$$\tilde{\mathbf{A}} = \mathbf{AZ} = \begin{bmatrix} 12 & -51 & 4 \\ 6 & -16 & -68 \\ -4 & 24 & -41 \end{bmatrix} \begin{bmatrix} -0.487 & 0.866 \\ 0.788 & 0.378 \\ 0.378 & 0.328 \end{bmatrix} = \begin{bmatrix} 35.847 & 30.968 \\ -41.211 & -23.134 \\ 5.364 & -7.834 \end{bmatrix},$$

em que  $\tilde{\mathbf{A}}$  será a nova matriz modelo. Observe que a soma das colunas dessa matriz é zero. Consequentemente,  $\tilde{\mathbf{A}}$  satisfaz a condição de centralização.

Por fim, podemos definir um modelo centrado matricialmente. Assim, para cada termo de suavização, a equação (2.2) pode ser escrita como:

$$g(\mu_i) = \tilde{\mathbf{X}}_i^* \tilde{\boldsymbol{\xi}}, \quad (2.11)$$

em que  $\tilde{\mathbf{X}}^* = (\mathbf{X} : \tilde{\mathbf{B}}_1 : \dots : \tilde{\mathbf{B}}_r)$  e  $\tilde{\boldsymbol{\xi}}^\top = (\boldsymbol{\beta}^\top, \boldsymbol{\tau}_1^\top, \dots, \boldsymbol{\tau}_r^\top)$ . Neste caso, a primeira coluna da matriz  $\mathbf{X}$  será formada por 1's, correspondente ao intercepto do modelo. Assim, (2.11) recai em um MLG, sendo possível escrever sua função de verossimilhança. No entanto, convém lembrar que, tomando uma quantidade  $q_k$  razoável de parâmetros com o propósito de se obter uma representação adequada das funções  $f_k$ 's desconhecidas, bem como as estimativas de  $\tilde{\boldsymbol{\xi}}$  pela maximização da verossimilhança ordinária, provavelmente provocará um sobre ajuste, sendo necessário, portanto, a incorporação de uma penalização à função de verossimilhança.

Convenientemente, utiliza-se penalidades que assumem uma forma quadrática nos coeficientes. É fácil mostrar que, devido ao fato de  $f(\cdot)$  ser linear nos parâmetros  $\boldsymbol{\gamma}_k$ , a penalidade baseada na integral do quadrado da segunda derivada de  $f(\cdot)$  pode sempre ser escrita como uma forma quadrática em  $\boldsymbol{\gamma}_k$ , como pode ser visualizado a seguir

$$\int [f_k''(t)]^2 dt = \boldsymbol{\gamma}_k^\top \mathbf{S}_k \boldsymbol{\gamma}_k, \quad (2.12)$$

em que  $\mathbf{S}_k$  é uma matriz de coeficientes conhecida. Porém, devido a reparametrização, essa penalidade se converte em  $\boldsymbol{\tau}_k^\top \bar{\mathbf{S}}_k \boldsymbol{\tau}_k$ , em que  $\bar{\mathbf{S}}_k = \mathbf{Z}_k^\top \mathbf{S}_k \mathbf{Z}_k$ . Finalmente, no caso específico dos *P-splines*, basta fazer  $\mathbf{S}_k = \mathbf{P}_k^d$ .

Diante do exposto, tomando a matriz  $\tilde{\mathbf{X}}^*$  e o vetor de parâmetros  $\tilde{\boldsymbol{\xi}}$  como definidos anteriormente, tem-se que as equações de estimação à luz do modelo reparametrizado serão, portanto, dadas por:

$$\mathbf{U}_p(\tilde{\boldsymbol{\xi}}) = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} (\mathbf{y} - \boldsymbol{\mu}) - \tilde{\mathbf{P}}(\boldsymbol{\alpha}) = \mathbf{0}, \quad (2.13)$$

sendo que  $\tilde{\mathbf{P}}(\boldsymbol{\alpha}) = \text{blocodiag} \{ \mathbf{0}_{pp}, \tilde{\mathbf{P}}^d(\boldsymbol{\alpha}) \}$ , no qual  $\mathbf{0}_{pp}$  é uma matriz de zeros de dimensão  $p \times p$ ,  $\tilde{\mathbf{P}}^d(\boldsymbol{\alpha}) = \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \boldsymbol{\tau}$  com  $\boldsymbol{\tau} = (\boldsymbol{\tau}_1^\top, \dots, \boldsymbol{\tau}_r^\top)^\top$  e  $\tilde{\mathbf{K}}(\boldsymbol{\alpha}) = \text{diag} \{ \alpha_1 \mathbf{Z}_1^\top \mathbf{P}_1^d \mathbf{Z}_1, \dots, \alpha_r \mathbf{Z}_r^\top \mathbf{P}_r^d \mathbf{Z}_r \}$ .

A título de ilustração, assumindo um modelo parcialmente aditivo generalizado com distribuição normal para a variável resposta, então segue que:

$$\mathbf{y} = \tilde{\mathbf{X}}^* \tilde{\boldsymbol{\xi}} + \boldsymbol{\epsilon}, \quad (2.14)$$

e conseqüentemente, deve-se minimizar:

$$\| (\mathbf{y} - \tilde{\mathbf{X}}^* \tilde{\boldsymbol{\xi}}) \|^2 + \tilde{\boldsymbol{\xi}}^\top \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \tilde{\boldsymbol{\xi}}. \quad (2.15)$$

O que resulta precisamente na estimativa de  $\tilde{\boldsymbol{\xi}}$ , dada da seguinte forma:

$$\hat{\boldsymbol{\xi}} = \left[ \tilde{\mathbf{X}}^{*\top} \tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right]^{-1} \tilde{\mathbf{X}}^{*\top} \mathbf{y}, \quad (2.16)$$

em que  $\tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) = \text{blocodiag} \left\{ \mathbf{0}_{pp}, \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \right\}$ . Este é o único caso em que é possível obter de forma fechada as estimativas de  $\tilde{\boldsymbol{\xi}}$ . Para outros casos é necessário um processo iterativo capaz de obtê-las numericamente, conforme será descrito a seguir.

### 2.3 ESTIMAÇÃO DOS PARÂMETROS

Para estimar os parâmetros dos MPLAGs, com suavização por meio de *P-splines*, é necessário encontrar os valores de  $\tilde{\boldsymbol{\xi}}$  que maximizam a função de log-verossimilhança (2.5). Tendo em vista o aspecto discutido anteriormente, é necessário que o modelo seja reparametrizado a fim de que este seja identificável. A partir desse modelo, obtêm-se as equações de estimação dada por (2.13). Além disso, Cordeiro e Demétrio (2008) lembram que, na grande maioria das vezes,  $\mathbf{U}_p(\tilde{\boldsymbol{\xi}})$  é um sistema de equações não-lineares que não possuem uma solução analítica. Isto é, não podem ser resolvidas diretamente para  $\tilde{\boldsymbol{\xi}}$ . Sendo, portanto, necessário algum procedimento numérico para resolvê-las. Deste modo, Green e Silverman (1994) mostram que, para  $\alpha_k$  e  $\phi$  fixos, o estimador de  $\tilde{\boldsymbol{\xi}}$  pode ser obtido iterativamente através do método dos mínimos quadrados ponderados penalizados da seguinte maneira:

$$\tilde{\boldsymbol{\xi}}^{(b+1)} = \left\{ \tilde{\mathbf{X}}^{*\top} \mathbf{W}^{(b)} \tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}^{-1} \left\{ \tilde{\mathbf{X}}^* \mathbf{W}^{(b)} \tilde{\mathbf{z}}^{(b)} \right\}, \quad b = 1, 2, \dots \quad (2.17)$$

em que  $\tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) = \text{blocodiag} \left\{ \mathbf{0}_{pp}, \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \right\}$ ,  $\tilde{\mathbf{K}}(\boldsymbol{\alpha}) = \text{diag} \left\{ \alpha_1 \mathbf{Z}_1^\top \mathbf{P}_1^d \mathbf{Z}_1, \dots, \alpha_r \mathbf{Z}_r^\top \mathbf{P}_r^d \mathbf{Z}_r \right\}$  e  $\mathbf{P}_k^d$  é uma matriz  $q_k \times q_k$  que depende apenas das diferenças entre os coeficientes adjacentes da base de De Boor.

De acordo com Manghi (2016), o estimador de  $\tilde{\boldsymbol{\xi}}$  pode ser obtido por meio do algoritmo *backfitting* (Gauss-Siedel). O processo iterativo, conforme apresentado pelo supracitado autor, é dado por:

$$\tilde{\boldsymbol{\xi}}_k^{(b+1)} = \mathbf{C}_k^{(b)} \mathbf{z}_k^{(b+1)} \quad k = 1, 2, \dots, r, \quad \text{e} \quad b = 1, 2, \dots,$$

de tal maneira que

$$\tilde{\boldsymbol{\xi}}_k^{(b+1)} = \begin{cases} \boldsymbol{\beta}^{(b+1)}, & \text{para } k = 0; \\ \tilde{\mathbf{B}}_k \boldsymbol{\tau}_k^{(b+1)}, & \text{para } k = 1, 2, \dots, r; \end{cases}$$

$$\mathbf{C}_k^{(b)} = \begin{cases} \left\{ \mathbf{X}^\top \mathbf{W}^{(b)} \mathbf{X} \right\}^{-1} \mathbf{X}^\top \mathbf{W}^{(b)}, & \text{para } k = 0; \\ \left\{ \tilde{\mathbf{B}}_k^\top \mathbf{W}^{(b)} \tilde{\mathbf{B}}_k + \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \right\}^{-1} \tilde{\mathbf{B}}_k^\top \mathbf{W}^{(b)}, & \text{para } k = 1, 2, \dots, r; \end{cases}$$

$$\mathbf{z}_k^{(b)} = \begin{cases} \tilde{\mathbf{z}}^{(b)} - \sum_{k=1}^r \tilde{\mathbf{B}}_k \boldsymbol{\tau}_k^{(b)}, & \text{para } k = 0; \\ \tilde{\mathbf{z}}^{(b)} - \mathbf{X}\boldsymbol{\beta}^{(b)} - \sum_{k < u}^r \tilde{\mathbf{B}}_k \boldsymbol{\tau}_k^{(b)} - \sum_{k > u}^r \tilde{\mathbf{B}}_u \boldsymbol{\tau}_u^{(b)}, & \text{para } k, u = 1, 2, \dots, r; \end{cases}$$

em que  $\tilde{\mathbf{z}} = \mathbf{D}^{-1}(\mathbf{y} - \boldsymbol{\mu}) + \mathbf{X}\boldsymbol{\beta} - \sum_{k=1}^r \tilde{\mathbf{B}}_k \boldsymbol{\tau}_k$ . Neste caso, Manghi (2016) afirma que  $b$  deve ser incrementado após  $k$  ser incrementado de 0 a  $r$ . O referido autor também mostra que, para  $r = 1$ , o processo iterativo para obter as estimativas de  $\tilde{\boldsymbol{\xi}}$  se resumem a seguinte forma:

$$\begin{aligned} \boldsymbol{\beta}^{(b+1)} &= \left\{ \mathbf{X}^\top \mathbf{W}^{(b)} \left[ \mathbf{I} - \mathbf{S}^{(b)} \right] \mathbf{X} \right\}^{-1} \mathbf{X}^\top \mathbf{W}^{(b)} \left[ \mathbf{I} - \mathbf{S}^{(b)} \right] \mathbf{z}^{(b)}; \\ \boldsymbol{\tau}^{(b+1)} &= \left\{ \tilde{\mathbf{B}}^\top \mathbf{W}^{(b)} \tilde{\mathbf{B}} + \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \right\}^{-1} \mathbf{X}^\top \mathbf{W}^{(b)} \left\{ \tilde{\mathbf{z}}^{(b)} - \mathbf{X}\boldsymbol{\beta}^{(b)} \right\}, \end{aligned}$$

em que  $\mathbf{S} = \tilde{\mathbf{B}} \left\{ \tilde{\mathbf{B}}^\top \mathbf{W}^{(b)} \tilde{\mathbf{B}} + \tilde{\mathbf{K}}(\boldsymbol{\alpha}) \right\}^{-1} \tilde{\mathbf{B}}\mathbf{W}$  e  $\mathbf{I}$  é a matriz identidade. Este método foi proposto por Breiman e Friedman (1985). Além deles, para mais detalhes sobre este procedimento, pode-se consultar também Buja, Hastie e Tibshirani (1989).

Por outro lado, de acordo com Wood (2017), a estimativa de  $\tilde{\boldsymbol{\xi}}$  também pode ser obtida minimizando a função

$$\left\| \sqrt{\mathbf{W}} \left( \tilde{\mathbf{z}} - \tilde{\mathbf{X}}^* \tilde{\boldsymbol{\xi}} \right) \right\|^2 + \tilde{\boldsymbol{\xi}}^\top \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \tilde{\boldsymbol{\xi}}, \quad (2.18)$$

cujo processo iterativo também recai em (2.17). Além disso, para computação prática, note que

$$\left\| \sqrt{\mathbf{W}} \left( \tilde{\mathbf{z}} - \tilde{\mathbf{X}}^* \tilde{\boldsymbol{\xi}} \right) \right\|^2 + \tilde{\boldsymbol{\xi}}^\top \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \tilde{\boldsymbol{\xi}} = \left\| \begin{bmatrix} \sqrt{\mathbf{W}} & 0 \\ 0 & \mathbf{I} \end{bmatrix} \left( \begin{bmatrix} \tilde{\mathbf{z}} \\ 0 \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{X}}^* \\ \mathbf{G} \end{bmatrix} \tilde{\boldsymbol{\xi}} \right) \right\|^2, \quad (2.19)$$

onde  $\mathbf{G}$  é uma matriz quadrada tal que  $\mathbf{G}^\top \mathbf{G} = \tilde{\mathbf{K}}^*(\boldsymbol{\alpha})$ . Esta forma é mais conveniente, tendo em vista que o autor usa a rotina da função glm do R para a obtenção das estimativas do modelo.

Ademais, para estimar o parâmetro  $\phi$ , seja  $r_i = (y_i - \mu_i) / V_i^{1/2}$ , então pode-se usar a seguinte relação:

$$\text{Var}(r_i) = \mathbb{E} \left( \left[ \frac{y_i - \mu_i}{V_i^{1/2}} \right]^2 \right) = \phi.$$

Dessa forma, segundo Manghi (2016), tem-se que uma estimativa de  $\phi$  pode ser dada tomando

$$\hat{\phi} = \frac{\sum_{i=1}^n \hat{r}_i^2}{\text{tr}(\mathbf{I} - \widehat{\mathbf{H}}^*)}, \quad (2.20)$$

em que  $\text{tr}(\mathbf{A})$  denota o traço da matriz  $\mathbf{A}$  e  $\mathbf{H}^* = \tilde{\mathbf{X}}^* \left\{ \tilde{\mathbf{X}}^{*\top} \tilde{\mathbf{W}} \tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}^{-1} \tilde{\mathbf{X}}^{*\top} \tilde{\mathbf{W}}$ .

A seguir, serão derivados os procedimentos para obtenção dos graus de liberdade e, posteriormente, procedimentos para a seleção dos parâmetros de suavização. Outrossim, será discutido a obtenção das estimativas aproximadas para os erros padrão de  $\hat{\xi}$ , bem como derivar medidas que evidencie as incertezas associadas às estimativas. Isto é, serão derivados os testes de hipóteses e intervalos de confiança.

## 2.4 INFERÊNCIA ESTATÍSTICA

### 2.4.1 Graus de liberdade efetivos

Em princípio, seria necessário estimar  $s_v = p + \sum_{k=1}^r v_k$  parâmetros, em que  $v_k = q_k - 1$ . Este corresponderia, em tese, aos graus de liberdade referente a estimação desses parâmetros. Contudo, em razão da penalização, há uma redução do espaço paramétrico. Isso implica que os graus de liberdade também sejam reduzidos.

Para a obtenção dos graus de liberdade efetivos, Hastie e Tibshirani (1990) usam a ideia de projetores lineares ou suavizadores, obtido a partir do ajuste do MPLAG. Pra isso, observa-se que na convergência do processo iterativo  $\hat{\eta} = \tilde{\mathbf{X}}^* \hat{\xi} = \hat{\mathbf{H}}^* \hat{\mathbf{z}}$ , em que

$$\mathbf{H}^* = \tilde{\mathbf{X}}^* \mathbf{A}^{*-1} \tilde{\mathbf{X}}^{*\top} \mathbf{W},$$

com  $\mathbf{A}^* = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha})$ . Com isto, Green e Silverman (1994) definem como os graus de liberdade efetivos do modelo a soma dos autovalores do suavizador. Isto é,

$$df(\boldsymbol{\alpha}) = \text{tr}(\hat{\mathbf{H}}^*).$$

Ou seja, os graus de liberdade efetivos são dados pelo traço da matriz  $\hat{\mathbf{H}}^*$ .

Logo, tomando  $\hat{\eta} = \sum_{k=0}^r \hat{\eta}_k$ , em que  $\hat{\eta}_k = \tilde{\mathbf{X}}_k \hat{\xi}_k$ , com  $\tilde{\mathbf{X}}_k = \mathbf{X}$ , se  $k = 0$  e  $\tilde{\mathbf{X}}_k = \tilde{\mathbf{B}}_k$ , para  $k = 1, 2, \dots, r$ . Além disso, seja  $\hat{\mathbf{A}}_k^{*-1}$  particionada da seguinte forma  $\hat{\mathbf{A}}_k^{*-1} = \hat{\mathbf{A}}_k^{*k,u}$ , para  $k, u = 0, 1, 2, \dots, r$ . Definido então  $\hat{\eta}_k = \hat{\mathbf{H}}_k^* \hat{\mathbf{z}}$ , Manghi (2016) mostra que pode-se decompor o suavizador  $\hat{\mathbf{H}}_k^*$  como:

$$\hat{\mathbf{H}}_k^* = \tilde{\mathbf{X}}_k \left[ \sum_{u=0}^r \hat{\mathbf{A}}_k^{*k,u} \tilde{\mathbf{X}}_u^\top \right] \hat{\mathbf{W}}, \quad (2.21)$$

para  $k = 0, 1, 2, \dots, r$ . Desse modo, para determinar os graus de liberdade efetivos associado a  $k$ -ésima componente não-paramétrica de  $\hat{\eta}$ , tem-se que:

$$df_k(\alpha_k) = \text{tr}(\hat{\mathbf{H}}_k^*).$$

Observa-se que, por meio de (2.21),  $\widehat{\mathbf{H}}_k^* = \widetilde{\mathbf{B}}_k \left\{ \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \widetilde{\mathbf{B}}_k + \widetilde{\mathbf{K}}(\alpha_k) \right\}^{-1} \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}}$ , para  $k = 1, 2, \dots, r$ . Logo, os graus de liberdade correspondente a cada função não paramétrica serão obtidos por:

$$\begin{aligned} df_k(\alpha_k) &= \text{tr} \left( \widehat{\mathbf{H}}_k^* \right) \\ &= \text{tr} \left( \widetilde{\mathbf{B}}_k \left\{ \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \widetilde{\mathbf{B}}_k + \widetilde{\mathbf{K}}(\alpha_k) \right\}^{-1} \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \right). \end{aligned}$$

Em conformidade com Eilers e Marx (1996), essa quantidade ainda pode ser reescrita da seguinte maneira:

$$\begin{aligned} df_k(\alpha_k) &= \text{tr} \left( \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \widetilde{\mathbf{B}}_k \left\{ \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \widetilde{\mathbf{B}}_k + \widetilde{\mathbf{K}}(\alpha_k) \right\}^{-1} \right) \\ &= \text{tr} \left( [\mathbf{I}_k + \alpha_k \mathbf{L}_k] \right) \\ &= \sum_{j=1}^{v_k} \frac{1}{1 + \alpha_k \lambda_j}, \end{aligned}$$

para  $k = 1, 2, \dots, r$ , em que  $\lambda_j$  são os autovalores associados à matriz não negativa definida  $\mathbf{L}_k = \mathbf{Q}_k^{-\frac{1}{2}} \mathbf{P}_k^d \mathbf{Q}_k^{-\frac{1}{2}}$  e  $\mathbf{Q}_k^{-\frac{1}{2}}$  é uma matriz positiva definida, de forma que  $\mathbf{Q}_k = \widetilde{\mathbf{B}}_k^\top \widehat{\mathbf{W}} \widetilde{\mathbf{B}}_k$ . Sendo assim, tem-se uma relação inversa entre  $df_k(\alpha_k)$  e  $\alpha_k$ . Logo, sabendo um desses parâmetros, é possível obter o outro. Conseqüentemente, pode-se derivar um procedimento capaz de encontrar o valor ótimo para os parâmetros de suavização.

#### 2.4.2 Seleção do parâmetro de suavização

Conforme mencionando anteriormente, o parâmetro de suavização é usado para ponderar a quantidade ideal de suavização a fim de produzir uma função  $f$  que não interpole os dados, mas que também não seja tão suave a ponto de não representá-los. Nesta perspectiva, é necessário algum procedimento capaz de encontrar o valor ótimo desse parâmetro.

Como critério para escolha do parâmetro de suavização, Wood (2017) utiliza o *método da validação cruzada generalizada*. Este consiste em usar parte dos dados para ajustar o modelo e o restante para avaliar a adequabilidade do ajuste a fim de minimizar:

$$\text{VCG}(\boldsymbol{\alpha}) = \frac{(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\text{tr} \left( \mathbf{I} - \widehat{\mathbf{H}}^* \right)},$$

em que  $\mathbf{I}$  é uma matriz identidade  $\sum_{k=1}^r v_k \times \sum_{k=1}^r v_k$ , com  $v_k = q_k - 1$ .

Ademais, conforme foi visto anteriormente, existe uma relação inversa entre os graus de liberdade efetivos e o parâmetro de suavização. Além disso, sabendo um deles, pode-se obter o outro. Na prática, é mais fácil estabelecer uma faixa de valores para os graus de

liberdade. Segundo Manghi (2016), para valores pequenos de  $\alpha_k$ , tem-se que  $df_k(\alpha_k) \approx v_k$ , enquanto que para altos valores de  $\alpha_k$ , tem-se que  $df_k(\alpha_k) \approx 2$ . Uma vez que  $df_k(\alpha_k) \in [2, v_k]$  e  $\alpha_k \in [0, \infty]$ , o autor sugere considerar para a seleção do vetor de parâmetros de suavização os valores  $\alpha_k \in [\alpha_{k(min)}, \alpha_{k(max)}]$ , tais que:

$$df_k(\alpha_{k(min)}) - df_k(\alpha_{k(max)}) = \gamma_k(v_k - 2),$$

tomando  $\gamma_k \in (0, 1)$ . Assim, de acordo com Manghi (2016), pode-se selecionar o vetor de suavização como sendo o vetor  $\boldsymbol{\alpha}_k^* \in [\alpha_{1(min)}, \alpha_{1(max)}] \times [\alpha_{2(min)}, \alpha_{2(max)}] \times \cdots \times [\alpha_{r(min)}, \alpha_{r(max)}]$  e usa-los para minimizar (2.4.2) e, assim, selecionar o vetor de parâmetros de suavização ótimo.

Além da validação cruzada generalizada, de acordo com Wood (2017), pode-se usar também medidas de informação, como o Critério de Informação de Akaike (AIC) e Critério de Informação Bayesiano (BIC). Para mais detalhes sobre esses critérios, ver Akaike (1974) e Schwarz *et al.* (1978), respectivamente.

### 2.4.3 Escolha do número de nós

A suavidade da curva estimada é controlada por meio de uma penalização, que é função do parâmetro de suavização  $\alpha$ . Portanto, a escolha da posição e do número de nós não é, necessariamente, crucial. De acordo com Ruppert (2002), deve haver uma quantidade suficiente de nós para o ajuste da curva, mas qualquer incremento a partir desse número mínimo necessário de nós tem pouco efeito no ajuste. Porém, existem exemplos em que uma quantidade de nós acima do mínimo necessário aumenta moderadamente o erro quadrático médio. Assim, o autor recomenda escolher uma quantidade de nós de modo a garantir que este seja suficientemente grande para o ajuste dos dados, mas que não seja tão grande que possa resultar em um erro quadrático médio maior que o necessário.

Um dos algoritmos proposto pelo autor considera utilizar 5, 10, 20, 40, 80 e 120 como valores de avaliação para a quantidade de nós necessária. Inicialmente, é ajustado um *P-spline* com quantidade de nós iguais a 5 e 10. Em cada caso, um  $\alpha$  é selecionado para minimizar o critério VCG de acordo com o número de nós. Então, compara-se os valores obtidos pelo critério VCG e escolhe-se a quantidade de nós que resultou no menor valor do critério. Assim, se para a quantidade de nós igual a 10 o critério VCG for maior que 0.98 vezes o valor do critério VCG obtido com 5 nós, então conclui-se que é improvável que novos aumentos no número de nós diminuam o VCG e usa-se 5 nós ou 10 (o que tiver menor VCG). Caso contrário, calcula-se um *P-spline* com 20 nós e compara-se o valor do critério VCG obtido com o valor do critério VCG

obtido com 10 nós. Se para a quantidade de nós igual a 20 o critério VCG for maior que 0.98 vezes o valor do critério VCG obtido com 10 nós, então usa-se 10 nós ou 20 (o que tiver menor VCG). Caso contrário, calcula-se um *P-spline* com 40 nós e assim por diante.

#### 2.4.4 Intervalos de confiança

Nas seções anteriores foram apresentados procedimentos para obter as estimativas pontuais para os parâmetros  $\tilde{\boldsymbol{\zeta}}$  e  $\boldsymbol{\alpha}$ . Entretanto, também é interessante uma evidência quantitativa a respeito da incerteza que envolve essas estimativas. Em particular, seria útil realizar testes de hipóteses a cerca dos parâmetros, bem como poder encontrar seus respectivos intervalos de confiança. Para isso, primeiramente é necessário obter os erros-padrão aproximados associados às estimativas.

Desta forma, Manghi (2016) observa que, na convergência do processo iterativo (dado  $\boldsymbol{\alpha}$ ),

$$\widehat{\tilde{\boldsymbol{\zeta}}} = \widehat{\mathbf{A}}^{*-1} \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}^* \widehat{\mathbf{z}}. \quad (2.22)$$

Assumindo que  $\boldsymbol{\mu} = \boldsymbol{\mu}_0$  é o verdadeiro vetor de médias de  $\mathbf{y}$  em que  $\widehat{\boldsymbol{\mu}} \approx \boldsymbol{\mu}_0$ , tem-se que  $\widehat{\tilde{\boldsymbol{\zeta}}} = \mathbf{A}_0^{*-1} \widetilde{\mathbf{X}}^{*\top} \mathbf{W}_0 \mathbf{z}_0 \approx \mathbf{A}^{*-1} \widetilde{\mathbf{X}}^{*\top} \mathbf{W}^* \mathbf{z}$ . Assim, sabendo que  $\mathbf{y}$  possui matriz de variância e covariância dada por  $\mathbf{V}^{-1}\phi$ , então:

$$\begin{aligned} \text{Var}(\widehat{\tilde{\boldsymbol{\zeta}}}) &\approx \mathbf{A}^{*-1} \widetilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} \text{Var}(\mathbf{y}) \mathbf{D}^{-1} \mathbf{W} \widetilde{\mathbf{X}}^* \mathbf{A}^{*-1} \\ &\approx \mathbf{A}^{*-1} \widetilde{\mathbf{X}}^{*\top} \mathbf{W} \widetilde{\mathbf{X}}^* \mathbf{A}^{*-1} \phi \end{aligned} \quad (2.23)$$

será a matriz de variância e covariância do estimador  $\widehat{\tilde{\boldsymbol{\zeta}}}$ . Com efeito, a estimativa de  $\text{Var}(\widehat{\tilde{\boldsymbol{\zeta}}})$  será dada por:

$$\widehat{\text{Var}}(\widehat{\tilde{\boldsymbol{\zeta}}}) = \widehat{\mathbf{A}}^{*-1} \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}} \widetilde{\mathbf{X}}^* \widehat{\mathbf{A}}^{*-1} \widehat{\phi}. \quad (2.24)$$

Para a normalidade de  $\mathbf{y}$  ou normalidade multivariada assintótica de  $\widetilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{z}$ , segue-se que, aproximadamente:

$$\widehat{\tilde{\boldsymbol{\zeta}}} \sim \text{N}\left(\text{IE}(\widehat{\tilde{\boldsymbol{\zeta}}}), \widehat{\text{Var}}(\widehat{\tilde{\boldsymbol{\zeta}}})\right).$$

No entanto, Wood (2017) aponta para o fato de que geralmente,  $\text{IE}(\widehat{\tilde{\boldsymbol{\zeta}}}) \neq \tilde{\boldsymbol{\zeta}}$  e, portanto, é inapropriado usar esse resultado para realizar testes de hipóteses e calcular os intervalos de confiança. Contudo, se  $\tilde{\boldsymbol{\zeta}} = 0$ , então  $\text{IE}(\widehat{\tilde{\boldsymbol{\zeta}}}) = 0$  e o resultado anterior pode ser útil para testar a inclusão de termos no modelo.

Alternativamente, pode-se usar a abordagem bayesiana para estimar a incerteza. Neste caso, de acordo com Wood (2017), a matriz de variâncias e covariâncias é baseada na distribuição a posteriori dos parâmetros que, ainda segundo o autor, é dada da seguinte forma:

$$\widehat{\text{Var}}\left(\widehat{\boldsymbol{\xi}}\right) = \widehat{\mathbf{A}}^*{}^{-1} \widehat{\boldsymbol{\phi}},$$

em que  $\widehat{\mathbf{A}}^* = \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}} \widetilde{\mathbf{X}}^* + \widetilde{\mathbf{K}}^*(\boldsymbol{\alpha})$ , o que resulta em uma distribuição posteriori dada por:

$$\widehat{\boldsymbol{\xi}} | \mathbf{Y} = \mathbf{y} \sim \text{N}\left(\widehat{\boldsymbol{\xi}}, \widehat{\text{Var}}\left(\widehat{\boldsymbol{\xi}}\right)\right).$$

Para dados não normais a aproximação tem justificativa baseada no teorema central do limite. Logo, este resultado pode ser usado diretamente para calcular intervalos de credibilidade para os parâmetros.

Além disso, a partir dos intervalos de credibilidade para o vetor  $\widehat{\boldsymbol{\xi}}$ , também é possível obter os intervalos para quantidades dele derivadas, tais como seus próprios termos suaves. De fato, observe que:

$$\widehat{\text{Var}}\left(\widetilde{\mathbf{X}}^* \widehat{\boldsymbol{\xi}}\right) = \widetilde{\mathbf{X}}^* \widehat{\mathbf{A}}^*{}^{-1} \widetilde{\mathbf{X}}^{*\top} \widehat{\boldsymbol{\phi}}.$$

Com isso, um intervalo de credibilidade assintótico pontual de 95% para  $g(\mu_i)$  pode ser dado por:  $\text{IC}(g(\mu_i); 0, 95) = \left[\widetilde{\mathbf{x}}_i^* \widehat{\boldsymbol{\xi}} \pm 2\sqrt{c_{ii}}\right]$ , em que  $c_{ii}$  é o  $i$ -ésimo elemento da diagonal de  $\widehat{\text{Var}}\left(\widetilde{\mathbf{X}}^* \widehat{\boldsymbol{\xi}}\right)$ . De igual forma, pode-se conseguir um intervalo de credibilidade para  $\mu_i$ , bastando apenas aplicar a inversa da função  $g(\cdot)$  de tal forma que  $\text{IC}(\mu_i; 0, 95) = \left[g^{-1}\left(\widetilde{\mathbf{x}}_i^* \widehat{\boldsymbol{\xi}} \pm 2\sqrt{c_{ii}}\right)\right]$ .

### 2.4.5 Teste de hipóteses

Os testes de hipóteses podem ser usados para a seleção de modelos, especialmente quando existem razões para favorecer modelos mais simples. No contexto dos MPLAGs, seja  $\boldsymbol{\tau}_k$  o vetor que contém os coeficientes para um único termo suave, então  $\mathbb{E}(\widehat{\boldsymbol{\tau}}_k) \approx 0$  se  $\boldsymbol{\tau}_k = 0$ . Desta forma, considerando testar a hipótese  $H_0 : \boldsymbol{\tau}_k = \mathbf{0}$ , então pode-se usar o resultado obtido em (2.4.4) para realizar esse teste. Assim, seja  $\text{Var}(\widehat{\boldsymbol{\tau}}_k)$  denotando o bloco de  $\text{Var}\left(\widehat{\boldsymbol{\xi}}\right)$ , definida em (2.23), correspondente à  $\boldsymbol{\tau}_k$ , tem-se que:

$$\mathbf{W}_{\boldsymbol{\tau}_k} = \widehat{\boldsymbol{\tau}}_k^\top \text{Var}(\widehat{\boldsymbol{\tau}}_k)^{-1} \widehat{\boldsymbol{\tau}}_k. \quad (2.25)$$

Que, sob  $H_0$ , segue assintoticamente uma distribuição  $\chi_{v_k}^2$ . Logo,

$$\widehat{\boldsymbol{\tau}}_k^\top \text{Var}(\widehat{\boldsymbol{\tau}}_k)^{-1} \widehat{\boldsymbol{\tau}}_k / v_k \sim F_{v_k, n-d}, \quad (2.26)$$

em que  $d$  e  $v_k$  são, respectivamente, as dimensões de  $\tilde{\boldsymbol{\xi}}$  (desconsiderando os termos paramétricos) e  $\boldsymbol{\tau}_k$ .

Além disso, substituindo  $(\mathbf{C}\hat{\boldsymbol{\tau}}_k - \mathbf{u})^\top (\mathbf{C}\text{Var}(\hat{\boldsymbol{\tau}}_k)\mathbf{C}^\top)^{-1} (\mathbf{C}\hat{\boldsymbol{\tau}}_k - \mathbf{u})$  no lugar de  $\hat{\boldsymbol{\tau}}_k^\top \text{Var}(\hat{\boldsymbol{\tau}}_k)^{-1} \hat{\boldsymbol{\tau}}_k$ , pode-se também testar hipóteses mais gerais, tais como:  $H_0 : \mathbf{C}\boldsymbol{\tau}_k = \mathbf{u}$ . Esse resultado é igualmente útil, quando se deseja testar casos que envolvem um único parâmetro, podendo usar equivalentemente os testes  $t$  ou  $N(0,1)$  como distribuições de referência.

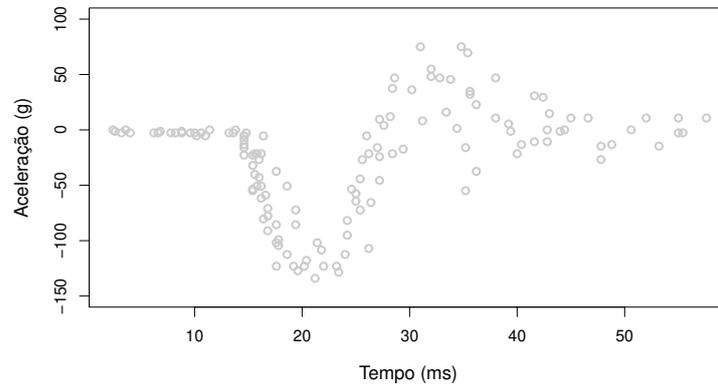
Porém, como penalizar a função de verossimilhança leva a uma redução do espaço paramétrico, Wood (2017) destaca que o cálculo do valor  $p$  é menos direto, sendo possível levar o teste a uma perda de poder. Portanto, o supracitado autor recomenda, entre outras propostas, utilizar  $df_k(\boldsymbol{\alpha})$  no lugar de  $v_k$ .

A seguir, será apresentado um exemplo prático. O intuito será apenas destacar os pormenores no que tange os resultados obtidos até aqui. Em particular, o objetivo é descrever com certo detalhe como o processo de suavização é realizado de forma prática, com indicação dos passos a serem tomados, bem como apresentar os códigos desenvolvidos em R para a obtenção dos resultados.

## 2.4.6 Impacto do capacete da motocicleta

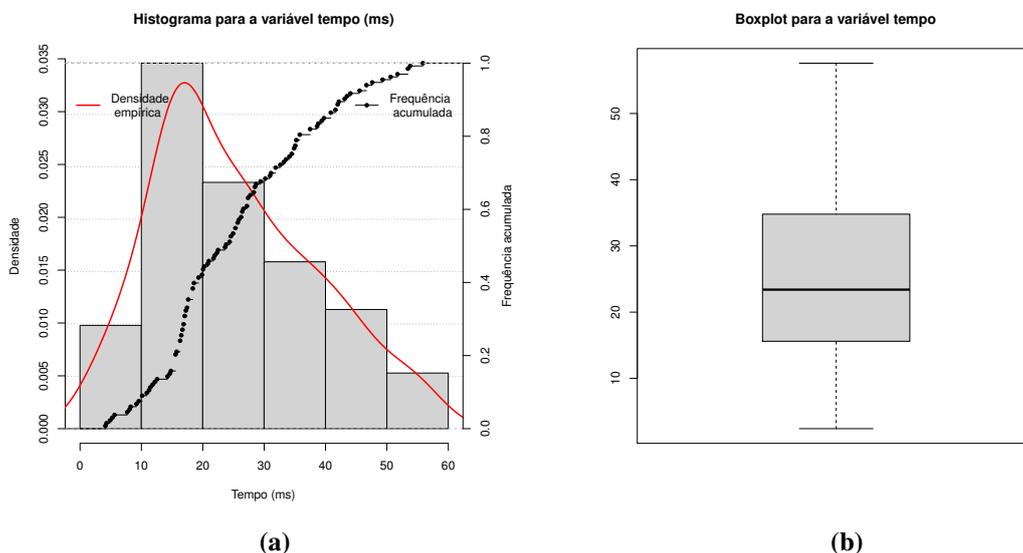
A título de exemplificação, considere novamente o exemplo exposto na **Seção 1.3.1**. Na supracitada seção, foi apresentado um experimento cujo objetivo era investigar comportamento de capacetes durante o impacto na cabeça em acidentes simulados. Para tanto, é considerada aceleração da cabeça em unidades de  $g$ , imediatamente após o impacto. A **Figura 8** mostra o diagrama de dispersão da aceleração contra o tempo dado em milissegundo (ms) para este experimento.

Estes dados podem ser encontrados em Härdle (1990), cujo estudo é baseado em uma amostra de tamanho  $n = 133$  referente a um experimento que consiste em medir a quantidade de energia absorvida por capacetes imediatamente após o impacto. Tal experimento, em geral, considera um impactador com uma determinada massa, sujeito à queda livre, ao qual é acoplada uma célula de carga que proverá pontos de dados de carga contra o tempo de contato entre o impactador e o alvo. Este experimento é descrito em detalhes por Schmidt, Mattern e Schüller (1981). Contudo, segue uma breve análise descritiva do problema.



**Figura 8 – Velocidade do impactador em função do tempo de contato entre o projétil e o alvo para uma energia de impacto constante.**

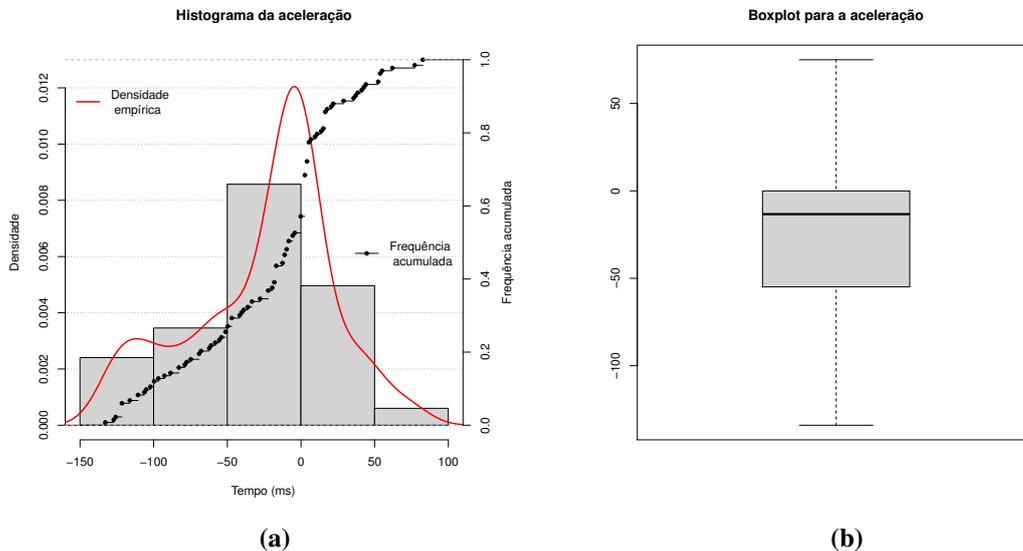
Para este conjunto de dados, observou-se que o tempo médio foi de aproximadamente  $\bar{x} = 25.18$  ms com desvio padrão de  $s = 13.13$  ms, bem como coeficiente de variação dado por  $cv = 52.15\%$ . Além disso, o tempo mínimo observado foi de  $x_{(1)} = 2.40$  ms com tempo máximo de  $x_{(n)} = 57.60$  ms, com o primeiro quartil, mediana e terceiro quartil dados respectivamente por  $Q_1 = 15.60$  ms,  $Md = 23.40$  ms e  $Q_3 = 34.80$  ms.



**Figura 9 – Gráficos da distribuição de frequência, frequência acumulada e boxplot, respectivamente, para a variável tempo.**

De igual forma, destaca-se que a média da *aceleração* foi de aproximadamente  $\bar{x} = -25.54$  g, com desvio padrão de  $s = 48.32$  g, bem como a menor e a maior acelerações observadas foram de  $x_{(1)} = -134$  g e  $x_{(n)} = 75$  g, respectivamente. Ademais, o primeiro quartil, mediana e terceiro quartil foram dados respectivamente por  $Q_1 = -54.90$  g,  $Md = -13.30$  g e

$$Q_3 = 0.00 \text{ g.}$$



**Figura 10 – Gráficos da distribuição de frequência, frequência acumulada e boxplot, respectivamente, para a variável aceleração.**

Como pode-se observar, o comportamento da aceleração em função do tempo claramente dar-se de forma não linear. Com isso, assumir o modelo normal clássico, i.e., com o predictor dado de forma linear, para analisar tal comportamento é inviável. Portanto, consideraram erros aditivos e distribuição normal para ajuste do modelo cuja variável tempo é ajustada de forma não paramétrica e assumindo  $y_i$  como sendo a aceleração, o modelo adotado será dado por

$$y_i | \text{Tempo}_i \sim N(\mu_i, \sigma^2),$$

$$g(\mu_i) = \mu_i = \eta_i = f(\text{Tempo}_i),$$

em que  $f(\cdot)$  é uma função suave definida como em (1.1).

Para o ajuste desse modelo, usou-se uma função em R dado pelo **Código-fonte 3**, baseado na teoria desenvolvida nas seções anteriores. Neste código, especificando os nós e  $\alpha$ , a função `psfit` retorna o ajuste do modelo baseado na função `glm`, o valor do critério VCG, sigma e o traço da matriz H (correspondente ao grau de liberdade do modelo). Para a escolha do  $\alpha$  ótimo pode-se usar, por exemplo, o **Código-fonte 4**. Tanto o **Código-fonte 3** quanto o **Código-fonte 4**, bem como outros códigos usados para manipulação da metodologia apresentada nas seções subsequentes, são disponibilizados no **Apêndice A**.

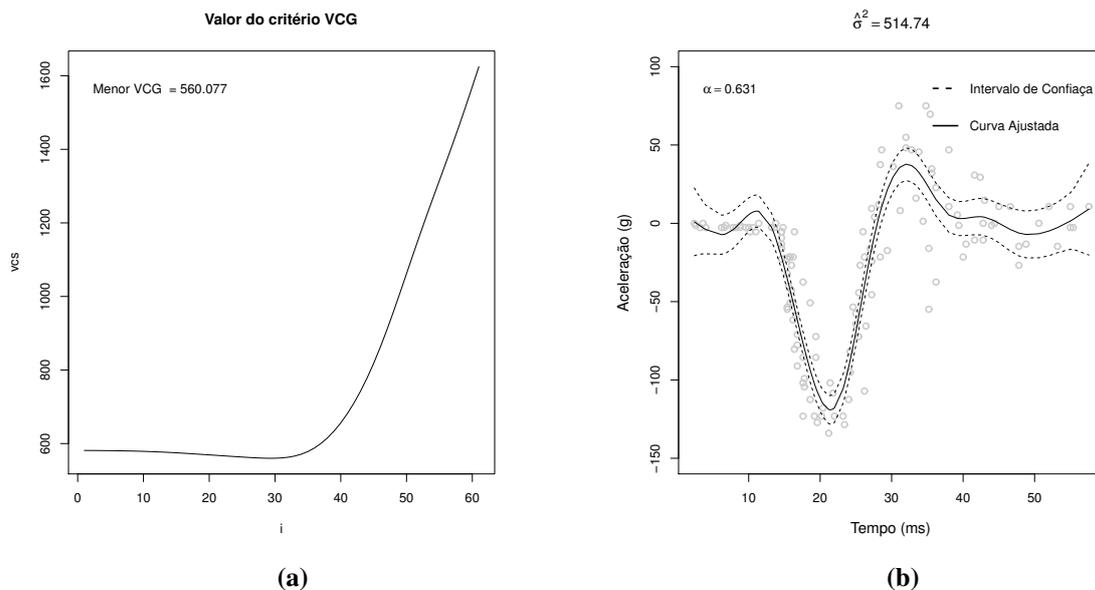
Assim, ajustando o modelo proposto, obtém-se as estimativas de máxima verossimilhança penalizadas apresentadas na **Tabela 0**. Destaca-se que para o ajuste desse modelo foi

utilizado o procedimento de *P-spline* para estimar  $f(\cdot)$  por meio de um *B-spline* cúbico com 18 nós equidistante e penalização quadrática. Além disso, foi empregado o método da validação cruzada generalizada (VCG) para estimar o parâmetro de suavização  $\alpha$ . Sobre os graus de liberdade, por definição, tem-se que  $df(\alpha) = 1 + df_1(\alpha)$ , em que 1 refere-se ao intercepto do modelo e  $df_1(\alpha)$  corresponde ao grau de liberdade da variável suavizada.

**Tabela 0 – Resumo das estimativas referentes ao modelo ajustado para o exemplo do impacto do capacete.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	-25.546	1.967	< 0.0001
Outras medidas do ajuste			
$df(\alpha)$	10.767		
$\alpha$	0.631		
$\sigma^2$	514.738		
VCG	560.077		

Portanto, por meio da **Figura 11**, verifica-se o bom ajuste da curva, bem como é possível notar uma maior variabilidade para acelerações com tempo acima de 40 e abaixo de 10 milissegundos, destacado pelo intervalo de confiança pontual de 95%, com base nos resultados apresentados na **Seção 2.4.4**.



**Figura 11 – O valor do critério VCG, curva aproximada e intervalos de confiança pontuais de 95% para a aceleração contra o tempo.**

No entanto, para verificar a viabilidade do modelo, é necessário ainda uma análise

subsequente ao ajuste. Esta é responsável por verificar afastamento das suposições adotadas, bem como corrigi-las com objetivo de obter um modelo mais fiel à realidade do fenômeno em questão. No próximo capítulo, serão detalhados os procedimentos pelos quais se faz essa análise seguida de exemplificações.

### 3 TÉCNICAS DE DIAGNÓSTICO NOS MPLAGS

Modelos estatísticos são descrições aproximadas de processos bastante complexos, conseqüentemente podem levar a resultados imprecisos. Isto posto, surge uma importante motivação para o estudo de técnicas que avaliem essa inexatidão. Para tanto, foram desenvolvidos métodos para avaliarem a qualidade de um ajuste de regressão. Esses métodos se inserem no ramo da estatística conhecida como *análise de diagnóstico* que, de acordo com Billor e Loynes (1993), se iniciou com a análise de resíduos para detectar a presença de pontos extremos e avaliar a adequação da distribuição proposta para a variável resposta.

A ideia básica da análise de diagnóstico é verificar se, após o ajuste do modelo, o resultado obtido é satisfatório. Isso é necessário, tendo em vista que erros sistemáticos, ou até mesmo algumas observações discrepantes, podem contribuir de forma inadequada na estimação dos parâmetros do modelo. Os erros sistemáticos são identificáveis e, em geral, podem ser corrigidos. Estão associados à escolha inadequada da função de ligação ou pela escolha equivocada para a distribuição da variável resposta. Com relação a ocorrência das observações discrepantes, pode ser que, de fato, estejam erradas como resultado de uma leitura incorreta, ou porque algum fator não controlado acarretou na sua obtenção, ou porque são observações que ocorreriam com pouca frequência caso fosse realizado outro processo de amostragem para obtenção de dados provenientes de uma mesma população. Em qualquer um desses casos é necessário uma análise mais cuidadosa em torno dessas observações, que exercem um peso desproporcional na obtenção das estimativas, a fim de identificar quais foram as possíveis razões ou fatores que levaram a sua ocorrência.

Na prática, pode haver uma combinação dos diferentes tipos de erros. Esse fato, conforme observado por Cordeiro e Demétrio (2008), faz com que a verificação da adequação de um modelo para um determinado conjunto de observações seja um processo realmente difícil. Contudo, foram desenvolvidas várias técnicas com o intuito de superar essa dificuldade. A análise de resíduo, por exemplo, é uma das técnicas mais antigas e mais utilizadas. Esta técnica possibilita identificar comportamentos atípicos para o componente aleatório esperado para uma determinada especificação da distribuição da variável resposta. Uma forma bastante geral para definir resíduos pode ser vista em Cox e Snell (1968).

Para lidar com observações que influenciam nas estimativas dos parâmetros do modelo, Hoaglin e Welsh (1978) propõem estudar a matriz de projeção do modelo linear padrão  $H = X(X^T X)^{-1} X$ . Observações influentes tendem a elevar os valores da diagonal desta

matriz. Com isso, pode-se identificar pontos que destoam da grande maioria das observações em  $\mathbf{X}$  e causam um peso desproporcional no valor ajustado. Por essa razão, esses pontos recebem o nome de pontos de alavanca. Existem generalizações dessa medida para a classe dos MLGs.

Segundo Paula (2013), uma das técnicas mais utilizada é a deleção de pontos que visa avaliar o impacto da retirada de uma observação particular nas estimativas da regressão. Porém, ele destaca que uma das propostas mais inovadoras na área de diagnóstico é a análise de **influência local**, desenvolvida por Cook (1978). Em seus estudos, "Cook propõe avaliar a influência conjunta das observações sob pequenas mudanças (perturbações) no modelo ou nos dados, ao invés da avaliação pela retirada individual ou conjunta de pontos" (PAULA, 2013, p. 45). Para mais detalhes sobre esta técnica e as demais citadas anteriormente, sugere-se ver Paula (2013) e Cordeiro e Demétrio (2008).

### 3.1 ANÁLISE DE RESÍDUOS

A análise de resíduos é uma maneira eficiente de verificar a qualidade do ajuste do modelo. Nela busca-se identificar valores da variável resposta que apresentem algum comportamento atípico. Isto é, que possuem comportamento fora do padrão esperado em relação ao assumido para a variável  $y_i$ . Neste caso, os resíduos são uma grande fonte de informação, visto que por meio de seu comportamento é possível detectar a presença de pontos aberrantes, identificar a relevância de um fator omitido no modelo, assim como verificar se há indícios de afastamento sérios da distribuição assumida para o erro. Em resumo, por meio da análise de resíduos, é possível investigar possíveis falhas de suposições feitas a respeito do modelo, tais como: a escolha da função de ligação ou da distribuição da variável resposta.

#### 3.1.1 Resíduos de Pearson

No modelo de regressão linear, os resíduos mais utilizados são os *studentizados*. Segundo Cordeiro e Demétrio (2008), estes resíduos, sob normalidade, possuem uma distribuição  $t$  de Student com  $(n - p - 1)$  graus de liberdade. No caso dos MLG's, os resíduos de Pearson são os mais simples e os mais comuns. Em conformidade com Green e Silverman (1994), estes resíduos também podem ser utilizados para detectar a presença de *outliers* nos MPLAG's. Com isso, sob a suposição de independência e considerando  $\phi$  fixo, estes resíduos são dados por:

$$\hat{r}_i^P = \phi^{-1/2} \left( \frac{y_i - \hat{\mu}_i}{\hat{V}_i^{1/2}} \right),$$

com  $i = 1, 2, \dots, n$ . Como  $\mathbb{E}(y_i) = \mu_i$  e  $\text{Var}(y_i) = \phi V_i$ , se  $\hat{\mu}_i = \mu_i$ , então  $\mathbb{E}(\hat{r}_i^P) \approx 0$  e  $\text{Var}(\hat{r}_i^P) \approx 1$ . Portanto, de acordo com Manghi (2016), pode-se considerar como observações extremas aquelas que possuírem altos valores dos resíduos de Pearson.

### 3.1.2 Resíduos quantílicos aleatorizados

Alternativamente aos resíduos de Pearson, pode-se utilizar os resíduos quantílicos aleatorizados, que foram introduzidos por Dunn e Smyth (1996). Estes resíduos apresentam distribuição normal, independente da distribuição da variável resposta e de sua dispersão. Assumindo que  $y_i$  é uma variável contínua, os resíduos quantílicos aleatorizados são definidos da seguinte forma:

$$\hat{r}_i^q = \Phi^{-1} \{F(y_i; \hat{\mu}_i, \phi)\}, \text{ para } i = 1, 2, \dots, n,$$

em que  $\Phi$  é a função da distribuição acumulada da normal padrão e  $F(y_i; \hat{\mu}_i, \phi)$  é a função de distribuição acumulada da distribuição da variável resposta. Como se sabe, pelo Teorema da inversão,  $F(y_i; \hat{\mu}_i, \phi) \sim U(0, 1)$ . Assim, se os parâmetros do modelo são consistentemente estimados,  $\hat{r}_i^q$  convergem para uma distribuição normal.

Para o caso em que a distribuição da variável resposta é discreta, o recurso da aleatorização é aplicado de tal forma que:

$$\hat{r}_i^q = \Phi^{-1} \{u_i\},$$

em que  $u_i$  é uma variável aleatória uniformemente distribuída no intervalo  $(a_i, b_i]$ , com  $a_i = F(y_i - 1; \hat{\mu}_i)$  e  $b_i = F(y_i; \hat{\mu}_i)$ , para  $i = 1, 2, \dots, n$ . Da mesma forma que para o caso em que a variável resposta é contínua, se os parâmetros do modelo são consistentemente estimados, então a distribuição  $\hat{r}_i^q$  converge para uma normal padrão.

## 3.2 MEDIDAS DE ALAVANCAGEM

As medidas de alavancagem são usadas para identificar pontos extremos nos regressores. Nessas circunstâncias, de acordo com Paula (2013), essas observações são aquelas que contribuem de maneira excessiva para as estimativas dos parâmetros, associados a erros grandes (em valor absoluto) ou em razão de características dos regressores. As consequências disso são modelos que não representam adequadamente a realidade. Por isso, são pouco precisos em suas conclusões.

A principal ideia que está por trás do conceito de alavancagem consiste em conhecer a influência que cada observação exerce no próprio valor ajustado  $\hat{y}_i$ . Esta influência, de acordo com Wei, Hu e Fung (1998), pode ser bem representada pela derivada  $\partial \hat{y}_i / \partial y_i$ . Assim, no caso dos MPLAGs, tem-se que:

$$\widehat{\text{GL}} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{y}^\top} = \left\{ \mathbf{D}_{\tilde{\boldsymbol{\xi}}} \left( -\ddot{\mathcal{J}}_{\tilde{\boldsymbol{\xi}}\tilde{\boldsymbol{\xi}}} \right)^{-1} \ddot{\mathbf{L}}_{\tilde{\boldsymbol{\xi}}\mathbf{y}} \right\} \Big|_{\tilde{\boldsymbol{\xi}} = \hat{\tilde{\boldsymbol{\xi}}}}, \quad (3.1)$$

em que  $\mathbf{D}_{\tilde{\boldsymbol{\xi}}} = \partial \boldsymbol{\mu} / \partial \tilde{\boldsymbol{\xi}}^\top$ ,  $\ddot{\mathcal{J}}_{\tilde{\boldsymbol{\xi}}\tilde{\boldsymbol{\xi}}} = \partial^2 L_p(\tilde{\boldsymbol{\xi}}) / \partial \tilde{\boldsymbol{\xi}}^\top \partial \tilde{\boldsymbol{\xi}}^\top$  e  $\ddot{\mathbf{L}}_{\tilde{\boldsymbol{\xi}}\mathbf{y}} = \partial^2 L_p(\tilde{\boldsymbol{\xi}}) / \partial \tilde{\boldsymbol{\xi}}^\top \partial \mathbf{y}^\top$ . Observe que

$$\mathbf{D}_{\tilde{\boldsymbol{\xi}}} = \mathbf{D}\mathbf{X}^* \text{ e } \ddot{\mathbf{L}}_{\tilde{\boldsymbol{\xi}}\mathbf{y}} = \mathbf{X}^{*\top} \mathbf{W}\mathbf{D}^{-1}.$$

Além disso, tem-se que:

$$\ddot{\mathcal{J}}_{\tilde{\boldsymbol{\xi}}\tilde{\boldsymbol{\xi}}} = \tilde{\mathbf{X}}^{*\top} [\mathbf{M}_1 + \mathbf{M}_2] \tilde{\mathbf{X}}^* - \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}),$$

que é a informação de Fisher observada desde que tenhamos  $\mathbf{M}_1 = [\mathbf{D}\dot{\mathbf{V}}^{-1}\mathbf{D} + \dot{\mathbf{D}}\mathbf{V}^{-1}]\mathbf{C}$ ,  $\dot{\mathbf{D}} = \text{diag} \{ \partial^2 \mu_1 / \partial \eta_1^2, \dots, \partial^2 \mu_n / \partial \eta_n^2 \}$ , bem como  $\dot{\mathbf{V}} = \text{diag} \{ \partial V_1 / \partial \mu_1, \dots, \partial V_n / \partial \mu_n \}$ , com  $\mathbf{C} = \text{diag} \{ C_1, \dots, C_n \}$ , no qual  $C_i = y_i - \mu_i$  e  $\mathbf{M}_2 = -\mathbf{W}$ .

Note que, tomando o valor esperado de  $-\ddot{\mathcal{J}}_{\tilde{\boldsymbol{\xi}}\tilde{\boldsymbol{\xi}}}$ , isto é, o valor esperado do negativo da informação de Fisher, então temos que  $-\mathbb{E}(\ddot{\mathcal{J}}_{\tilde{\boldsymbol{\xi}}\tilde{\boldsymbol{\xi}}}) = \tilde{\mathbf{X}}^{*\top} \mathbf{W}\tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha})$ . Portanto, obtemos de forma aproximada

$$\widehat{\text{GL}} = \widehat{\mathbf{D}}\mathbf{X}^* \left\{ \tilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}\tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}^{-1} \mathbf{X}^{*\top} \widehat{\mathbf{W}}\widehat{\mathbf{D}}^{-1}.$$

Assim, deve-se estudar os elementos da diagonal de  $\widehat{\text{GL}}$  a fim de identificar os pontos de alavancagem.

Contudo, a medida mais comum usada para detectar pontos extremos nos MPLAGs é baseada no estudo da diagonal de uma matriz que é análoga a matriz de projeção obtida como solução de mínimos quadrados de uma regressão normal linear ponderada. Observe que na convergência do processo iterativo 2.17 apresentado na **Seção 2.3** referente a um MPLAG, a estimativa dos parâmetros  $\tilde{\boldsymbol{\xi}}$  possui a seguinte forma:

$$\hat{\tilde{\boldsymbol{\xi}}} = \left\{ \tilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}\tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}^{-1} \tilde{\mathbf{X}}^* \widehat{\mathbf{W}}\hat{\tilde{\mathbf{z}}},$$

em que  $\hat{\tilde{\mathbf{z}}} = \widehat{\mathbf{D}}^{-1}(\mathbf{y} - \hat{\boldsymbol{\mu}}) + \hat{\boldsymbol{\eta}}$ , com  $\hat{\boldsymbol{\eta}} = \tilde{\mathbf{X}}^* \hat{\tilde{\boldsymbol{\xi}}}$ . Assim, este resultado pode ser interpretado como uma regressão linear ponderada de uma variável dependente  $\tilde{\mathbf{z}}$  sobre a matriz modelo  $\tilde{\mathbf{X}}^*$  usando uma matriz de pesos  $\mathbf{W}$ . Portanto, temos que:

$$\tilde{\mathbf{X}}^* \hat{\tilde{\boldsymbol{\xi}}} = \mathbf{X}^* \left\{ \tilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}\tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}^{-1} \tilde{\mathbf{X}}^* \widehat{\mathbf{W}}\hat{\tilde{\mathbf{z}}},$$

então obtemos:

$$\hat{\boldsymbol{\eta}} = \widehat{\mathbf{H}}^* \widehat{\mathbf{z}},$$

chamando  $\widehat{\mathbf{A}}^* = \left\{ \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}} \widetilde{\mathbf{X}}^* + \widetilde{\mathbf{K}}^*(\boldsymbol{\alpha}) \right\}$  e  $\widehat{\mathbf{H}}^* = \widetilde{\mathbf{X}}^* \widehat{\mathbf{A}}^{*-1} \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}$  que, neste caso, representa a matriz de projeção.

Eubank (1984) discute algumas propriedades desta matriz para modelos de regressão não-paramétricos e mostra que suas propriedades são conservadas para o caso semi-paramétrico. Com isso, seja  $\hat{h}_{ii}^*$  os elementos diagonais da matriz  $\widehat{\mathbf{H}}^* = \widetilde{\mathbf{X}}^* \widehat{\mathbf{A}}^{*-1} \widetilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}}$ . Assim sendo, a distância de uma observação em relação as demais pode ser observada por  $\hat{h}_{ii}^*$ . Segundo Manghi (2016), esses elementos determinam quanto o valor predito é influenciado pelo valor observado da variável resposta  $y_i$ , considerando o respectivo vetor observado das variáveis independentes.

Ainda segundo o referido autor, supondo que todos os pontos exercem a mesma influência nas estimativas dos parâmetros, espera-se que  $\hat{h}_{ii}^*$  apresente valor próximo à sua média  $\bar{h}^* = \text{tr}(\widehat{\mathbf{H}}^*)/n$ . Diante disso, convém examinar os pontos tais que  $\hat{h}_{ii}^* \geq 3\bar{h}^*$ . Pode-se também tirar proveito do recurso gráfico, plotando os pontos  $\hat{h}_{ii}^*$  versus  $i$ , para  $i = 1, 2, \dots, n$ , a fim de identificar os pontos de alta alavancagem.

Sejam  $\hat{w}_i$  e  $\hat{v}_i$  os elementos diagonais das matrizes  $\widehat{\mathbf{W}}$  e  $\widehat{\mathbf{V}}$ , respectivamente. Como  $\hat{w}_i$  é função de  $\hat{v}_i$  e que por sua vez é função de  $\hat{\mu}_i$ , então  $\hat{h}_{ii}^*$  também pode ser função de  $\hat{\mu}_i$ . Neste caso, Paula (2013) recomenda usar os gráficos de  $\hat{h}_{ii}^*$  contra o valores ajustados para detectar pontos de alavanca. Por fim, pode-se verificar que as matrizes  $\widehat{\mathbf{H}}^*$  e  $\widehat{\mathbf{GL}}$  coincidem, de modo que  $\hat{h}_{ii}^* = \widehat{\mathbf{GL}}_{ii}$ , à medida que o tamanho da mostra cresce. No caso de ligação canônica essa igualdade é sempre válida para qualquer tamanho amostral.

### 3.3 INFLUÊNCIA LOCAL

Ajustando um modelo a um conjunto de dados, deseja-se que as estimativas obtidas não sofram mudanças substanciais quando submetidas a pequenas perturbações no modelo proposto ou nas observações. Portanto, tomando o modelo postulado como correto e comparando as estimativas obtidas através desse modelo com as estimativas decorrente de pequenas perturbações, considere o vetor de parâmetros de suavização  $\boldsymbol{\alpha}$  um vetor fixo, bem como  $\mathcal{F}_{\boldsymbol{\alpha}}(\tilde{\boldsymbol{\xi}})$  uma função penalizada e pelo menos duas vezes diferenciável em  $\tilde{\boldsymbol{\xi}}$  referente ao MPLAG, e, além disso, a estimativa de  $\tilde{\boldsymbol{\xi}}$  seja solução da equação

$$\frac{\partial \mathcal{F}_{\boldsymbol{\alpha}}(\tilde{\boldsymbol{\xi}})}{\partial \tilde{\boldsymbol{\xi}}} = \Psi_{\boldsymbol{\alpha}}(\tilde{\boldsymbol{\xi}}) = \mathbf{0}. \quad (3.2)$$

Tem-se que, portanto,  $\Psi_\alpha(\tilde{\boldsymbol{\xi}})$  é a função escore penalizada referente ao MPLAG. Assim, para avaliar a influência de  $\boldsymbol{\omega} \in \Omega \subseteq \mathbb{R}^l$ , em que  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_l)$  é um vetor de perturbação, Cook (1978) propõe a medida de influência denominada *afastamento pela verossimilhança* definida por:

$$LD(\boldsymbol{\omega}) = 2 \left\{ \mathcal{F}_\alpha(\hat{\tilde{\boldsymbol{\xi}}}) - \mathcal{F}_\alpha(\hat{\tilde{\boldsymbol{\xi}}}_\omega) \right\}, \quad (3.3)$$

quando  $\boldsymbol{\omega}$  varia numa vizinhança de  $\boldsymbol{\omega}_0$ , em que  $\boldsymbol{\omega}_0$  é tal que  $\mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}_{\boldsymbol{\omega}_0}) = \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}})$ . Aqui,  $\hat{\tilde{\boldsymbol{\xi}}}$  e  $\hat{\tilde{\boldsymbol{\xi}}}_\omega$  constituem as estimativas de  $\tilde{\boldsymbol{\xi}}$  sob os modelos postulado e perturbado respectivamente.

O gráfico de  $LD(\boldsymbol{\omega})$  versus  $\boldsymbol{\omega}$  contém informações essenciais da influência do esquema de perturbação utilizado. Esse gráfico pode ser interpretado como a superfície geométrica formado pelos valores do vetor

$$\alpha(\boldsymbol{\omega}) = \begin{pmatrix} \boldsymbol{\omega} \\ LD(\boldsymbol{\omega}) \end{pmatrix}, \quad (3.4)$$

quando  $\boldsymbol{\omega}$  varia em  $\Omega$ . Essa superfície é chamada de gráfico de influência. Aqui, o objetivo é avaliar como a superfície de  $\alpha(\boldsymbol{\omega})$  desvia-se de seu plano tangente em torno do ponto  $\boldsymbol{\omega}_0$ , em que  $\boldsymbol{\omega}_0$  é vetor de não perturbação.

Para isso, Cook (1978) propõe avaliar o comportamento de  $LD(\boldsymbol{\omega}_0 + a\mathbf{d})$  versus  $a$ , com  $a \in \mathbb{R}$  e  $\mathbf{d}$  uma direção com norma unitária. Desde que  $LD(\boldsymbol{\omega}_0 + a\mathbf{d}) = 0$  então,  $LD(\boldsymbol{\omega})$  tem mínimo local em  $a = 0$ . Esse gráfico é chamado de linha projetada. Cada linha projetada pode ser caracterizada por uma curvatura normal  $\mathcal{C}_\mathbf{d}$  em torno de  $a = 0$ . Nesse caso, Cook (1978) sugere usar aquela que possui a maior curvatura.

Cook (1978) mostra que a curvatura normal correspondente à maior curvatura  $\mathcal{C}_{\mathbf{d}_{max}}$ , é obtida considerando a direção  $\mathbf{d}_{max}$ , dada pelo maior autovetor associada ao maior autovalor da matriz definida por:

$$\mathcal{C}_\mathbf{d} = 2 \left| \mathbf{d}^\top \mathbf{F} \mathbf{d} \right|, \quad (3.5)$$

em que  $\mathbf{F} = \Delta^\top \ddot{\mathcal{F}}_\alpha^{-1} \Delta$ , desde que as matrizes  $\Delta$  e  $\ddot{\mathcal{F}}_\alpha$  sejam definidas respectivamente como:

$$\Delta = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tilde{\boldsymbol{\xi}} \partial \boldsymbol{\omega}^\top} = \frac{\partial \Psi_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \boldsymbol{\omega}^\top} \quad (3.6)$$

e

$$\ddot{\mathcal{F}}_\alpha = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tilde{\boldsymbol{\xi}} \partial \tilde{\boldsymbol{\xi}}^\top} = \frac{\partial \Psi_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tilde{\boldsymbol{\xi}}^\top} = \tilde{\mathbf{X}}^{*\top} [\mathbf{M}_1 + \mathbf{M}_2] \tilde{\mathbf{X}}^* - \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}), \quad (3.7)$$

sob a condição de  $\mathbf{M}_1 = [\mathbf{D}\dot{\mathbf{V}}^{-1}\mathbf{D} + \dot{\mathbf{D}}\mathbf{V}^{-1}]\mathbf{C}$ ,  $\dot{\mathbf{D}} = \text{diag}\{\partial^2\mu_1/\partial\eta_1^2, \dots, \partial^2\mu_n/\partial\eta_n^2\}$ , bem como  $\dot{\mathbf{V}} = \text{diag}\{\partial V_1/\partial\mu_1, \dots, \partial V_n/\partial\mu_n\}$ ,  $\mathbf{C} = \text{diag}\{C_1, \dots, C_n\}$ , no qual  $C_i = y_i - \mu_i$  e  $\mathbf{M}_2 = -\mathbf{W}$  avaliados em  $\tilde{\boldsymbol{\xi}} = \hat{\boldsymbol{\xi}}$  e  $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ . Note que a equação 3.7 corresponde a informação de Fisher observada para os MPLAGs. Este resultado pode ser encontrado no **Anexo A**.

Tomando  $f_{ii}$  como os elementos da diagonal principal da matriz  $\mathbf{F}$  e, além disso, considerando a direção da  $i$ -ésima observação definida como:

$$C_i = 2|f_{ii}|, \quad (3.8)$$

então pode-se utilizar o gráfico de  $C_i$  versus  $i$  (ordem dos dados) como técnica de influência local para avaliar a existência de observações influentes. Sugere-se considerar a  $i$ -ésima observação como influente se seu valor for maior que duas vezes a média das medidas  $C_i$ .

Contudo, Lee, Lu e Song (2006) sugerem utilizar o desvio padrão das medidas de curvatura para definir o ponto de corte na consideração de pontos potencialmente influentes. Assim, para o ponto de corte, é sugerido utilizar  $\mathcal{PC} = c_1\bar{C} + c_2\text{sd}(C)$  em que  $\bar{C}$  e  $\text{sd}(C)$  são a média e o desvio padrão das medidas  $C_i$ , respectivamente, e  $c_1$  e  $c_2$  são constantes pré-fixadas. Uma possibilidade seria  $c_1 = 1$  e  $c_2 = 2$  ou  $c_2 = 3$ . Neste caso, considera-se como observações potencialmente influentes aquelas tais que  $C_i > \mathcal{PC}$ .

Outra alternativa para definir a curvatura normal seria utilizar a negativa da informação de Fisher esperada no lugar da informação observada. Ou seja,

$$\mathcal{I}_\alpha = -\mathbb{E}\{\ddot{\mathcal{F}}_\alpha\} = -\mathbb{E}\left\{\frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tilde{\boldsymbol{\xi}} \partial \tilde{\boldsymbol{\xi}}^\top}\right\} = \tilde{\mathbf{X}}^{*\top} \widehat{\mathbf{W}} \tilde{\mathbf{X}}^* + \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}). \quad (3.9)$$

Com isso, pode-se utilizar para análise de influência local baseada na curvatura normal a quantidade

$$C_d(\tilde{\boldsymbol{\xi}}) = \Delta^\top \mathcal{I}_\alpha^{-1} \Delta. \quad (3.10)$$

Adicionalmente, a fim de obter uma medida que seja invariante sob mudança de escala, Poon e Poon (1999) propõem a curvatura normal conformal definida como

$$\mathcal{B}_d(\tilde{\boldsymbol{\xi}}) = -\frac{\mathbf{d}^\top \Delta^\top \ddot{\mathcal{F}}_\alpha^{-1} \Delta \mathbf{d}}{\sqrt{\text{tr}(\Delta^\top \ddot{\mathcal{F}}_\alpha^{-1} \Delta)}}. \quad (3.11)$$

Assim, para toda direção unitária  $\mathbf{d}$ ,  $0 \leq \mathcal{B}_d(\tilde{\boldsymbol{\xi}}) \leq 1$ . Dessa forma, tomando a direção  $\mathbf{d} = \mathbf{e}_i$ , em que  $\mathbf{e}_i$  é um vetor de dimensão  $n \times 1$  com o valor 1 na  $i$ -ésima posição e 0 nas demais posições, pode-se investigar o gráfico de  $\mathcal{B}_{\mathbf{e}_i}$  contra seus respectivos índices e considera-se como

observações potencialmente influentes aquelas tais que  $\mathcal{B}_i > \mathcal{PC}$ , com  $\mathcal{PC} = c_1\bar{\mathcal{B}} + c_2\text{sd}(\mathcal{B})$ , em que  $\bar{\mathcal{B}}$  e  $\text{sd}(\mathcal{B})$  são, respectivamente, a média e o desvio padrão de  $\mathcal{B}_i$ .

Para avaliar a qualidade do ajuste do modelo por meio da técnica de influência local, considera-se três tipos de perturbação dados pelos casos ponderados, perturbação aditiva na variável resposta e perturbação aditiva em uma das covariáveis. Billor e Loynes (1993) resumem os vários esquemas de perturbação que podem ser introduzidos através de  $\boldsymbol{\omega}$ , os quais podem ser divididos em dois grupos:

- i) Perturbação no modelo: Neste caso, a perturbação visa modificar as suposições propostas para o modelo. Um exemplo disso, seria a violação da suposição de homoscedasticidade que, para erros normalmente distribuídos, pode ser substituída por uma suposição de heteroscedasticidade.
- ii) Perturbação nos dados: Perturbar a variável resposta ou as variáveis explanatórias são exemplos de perturbação nos dados. Neste caso, as razões pelas quais se considera a perturbação nos dados residem nos possíveis erros de medida e a existência de observações aberrantes (*outliers*), em uma proporção relativamente pequena das observações, as quais podem exercer forte influência nas estimativas dos parâmetros do modelo.

### 3.3.1 Perturbação de casos

Neste caso, o objetivo é identificar observações influentes entre todas as observações individualmente. Isto é, buscar entender a contribuição individual de cada observação dentre todas as observações no processo de estimação do modelo proposto. Neste caso, o vetor de não perturbação será  $\boldsymbol{\omega}_0$ , tal que  $\omega_i = 1$  para todo  $i, i = 1, 2, \dots, n$ . Desta forma, tem-se que as equações de estimação do modelo perturbado são dadas por:

$$\Psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}|\boldsymbol{\omega}) = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} \text{diag}\{\boldsymbol{\omega}\} (\mathbf{y} - \boldsymbol{\mu}) - \mathbf{P}(\boldsymbol{\alpha}). \quad (3.12)$$

Assim, a equação (3.6) equivale a:

$$\Delta = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} \mathbf{C}.$$

A obtenção destes resultados podem ser conferidos no **Anexo A**.

### 3.3.2 Perturbação da variável resposta

Seja  $y_i$  perturbado de tal maneira que  $y_{\omega_i} = y_i + \omega_i V_i^{1/2}$ . Aqui, o vetor de não perturbação é tal que  $\omega_i = 0$  para  $i = 1, 2, \dots, n$ . Assim, tem-se que as equações de estimação

para este caso são dadas por:

$$\Psi_{\alpha}(\xi|\omega) = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\omega}) - \mathbf{P}(\alpha), \quad (3.13)$$

em que  $\boldsymbol{\mu}_{\omega} = \boldsymbol{\mu} - \mathbf{V}^{1/2} \boldsymbol{\omega}$ . Logo,

$$\Delta = \tilde{\mathbf{X}}^{*\top} \mathbf{W} \mathbf{D}^{-1} \mathbf{V}^{1/2}.$$

A obtenção destes resultados podem ser conferidos no **Anexo A**.

### 3.3.3 Perturbação em uma das variáveis explicativas

Se, no entanto, o objetivo for perturbar o vetor  $\mathbf{x}_j$  de covariáveis da matriz de planejamento  $\mathbf{X}$ , então considere a perturbação na  $j$ -ésima coluna da matriz  $\mathbf{X}$ , tal que  $x_{j\omega_i} = x_{ji} + \sigma_j \omega_j$ , no qual  $\sigma_j$  é o desvio padrão de  $\mathbf{x}_j$ , então  $\mu_{\omega_i} = g^{-1}(\eta_{\omega_i})$ , sendo  $\eta_{\omega_i} = \beta_1 x_{1i} + \dots + \beta_j (x_{ji} + \sigma_j \omega_j) + \sum_{k=1}^r f_k(t_{ik})$ , para  $i = 1, \dots, n$ . Assim, as equações de estimação sob o modelo perturbado serão dadas por:

$$\Psi_{\alpha}(\xi|\omega) = \mathbf{X}_{\omega}^{*\top} \mathbf{W}_{\omega} \mathbf{D}_{\omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\omega}) - \mathbf{P}(\alpha) = \mathbf{X}_{\omega}^{*\top} \mathbf{D}_{\omega} \mathbf{V}_{\omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\omega}) - \mathbf{P}(\alpha), \quad (3.14)$$

em que  $\mathbf{X}_{\omega}^* = \mathbf{X}^* + \mathbf{0}_{j\omega}$ , sendo  $\mathbf{0}_{j\omega}$  uma matriz de dimensão  $n \times s_v$  com  $\sigma_j \omega$  na  $j$ -ésima coluna e 0's nas demais componentes se  $j = \ell$  ou  $\mathbf{X}_{\omega}^* = \mathbf{X}^*$  se  $j \neq \ell$ , para  $j, \ell = 1, 2, \dots, p$ .

Portanto, temos que a matriz  $\Delta$  será dada por  $\Delta = (\Delta_1, \Delta_2, \dots, \Delta_n)$ , com  $\Delta_i$  sendo obtida da seguinte forma:

$$\Delta_i = \mathbf{T}_{1\omega_i} + \mathbf{T}_{2\omega_i},$$

em que

$$\begin{aligned} \mathbf{T}_{1\omega_i} &= \mathbf{X}_{\omega}^{*\top} \mathbf{D}_{\omega} \left[ \frac{\partial \mathbf{V}_{\omega}^{-1}}{\partial \omega_i} (\mathbf{y} - \boldsymbol{\mu}_{\omega}) + \mathbf{V}_{\omega}^{-1} \frac{\partial (\mathbf{y} - \boldsymbol{\mu}_{\omega})}{\partial \omega_i} \right], \\ \mathbf{T}_{2\omega_i} &= \left[ \frac{\partial \mathbf{X}_{\omega}^{*\top}}{\partial \omega_i} \mathbf{D}_{\omega} + \mathbf{X}_{\omega}^{*\top} \frac{\partial \mathbf{D}_{\omega}}{\partial \omega_i} \right] \mathbf{V}_{\omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\omega}). \end{aligned}$$

Neste caso, temos que:

$$\begin{aligned} \frac{\partial \mathbf{V}_{\omega}^{-1}}{\partial \omega_i} &= -\mathbf{V}_{\omega}^{-1} \frac{\partial \mathbf{V}_{\omega}}{\partial \omega_i} \mathbf{V}_{\omega}^{-1}, \\ \frac{\partial \mathbf{V}_{\omega}}{\partial \omega_i} &= \frac{\partial \mathbf{V}_{\omega}}{\partial \mu_{\omega_i}} \mathbf{D}_{\omega}, \\ \frac{\partial \mathbf{V}_{\omega}}{\partial \mu_{\omega_i}} &= \beta_j \sigma_j \text{diag}\{0, \dots, \partial^2 V_{\omega_i} / \partial \mu_{\omega_i}^2, \dots, 0\}. \end{aligned}$$

Além disso,

$$\begin{aligned}\frac{\partial(\mathbf{y} - \boldsymbol{\mu}_\omega)}{\partial\omega_i} &= -\beta_j\sigma_j(0, \dots, \partial\mu_{\omega_i}/\partial\eta_{\omega_i}, \dots, 0)^\top, \\ \frac{\partial\mathbf{D}_\omega^\top}{\partial\omega_i} &= \beta_j\sigma_j\text{diag}\{0, \dots, \partial^2\mu_{\omega_i}/\partial\eta_{\omega_i}^2, \dots, 0\},\end{aligned}$$

bem como  $\frac{\partial\mathbf{X}_\omega^{*\top}}{\partial\omega_i} = \frac{\partial\mathbf{X}^{*\top}}{\partial\omega_i} + \frac{\partial\mathbf{0}_{j\omega}^\top}{\partial\omega_i}$ , resultando em

$$\frac{\partial\mathbf{X}_\omega^{*\top}}{\partial\omega_i} = \begin{cases} \mathbf{0}, & \text{para } j \neq \ell, \\ \sigma_j\dot{\mathbf{X}}_{\omega_i}^{*\top}, & \text{para } j = \ell, \end{cases} \quad (3.15)$$

em que  $\dot{\mathbf{X}}_{\omega_i}^*$  é uma matriz com dimensão  $n \times s_v$  com 1 sendo o elemento  $(i, j)$  e 0 nas demais componentes.

### 3.3.4 Perturbação em vetores particionados

Quando há interesse em uma parte específica do conjunto de parâmetros, Cook (1978) mostra que, considerando a partição  $\tilde{\boldsymbol{\xi}} = (\tilde{\boldsymbol{\xi}}_1, \tilde{\boldsymbol{\xi}}_2)$ , a função de afastamento da verossimilhança será dada por:

$$LD(\boldsymbol{\omega})_s = 2 \left\{ \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}) - \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}_{1\omega}, g(\tilde{\boldsymbol{\xi}}_{1\omega})) \right\},$$

sendo  $\tilde{\boldsymbol{\xi}}_{1\omega}$  o vetor obtido a partir de  $(\tilde{\boldsymbol{\xi}}_{1\omega}, \tilde{\boldsymbol{\xi}}_{2\omega})$  e  $g(\tilde{\boldsymbol{\xi}}_{1\omega})$  a função que, para  $\tilde{\boldsymbol{\xi}}_1$  fixado, maximiza  $\mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}_1, \tilde{\boldsymbol{\xi}}_2)$ .

Consequentemente, a superfície será dada por:

$$\alpha_s(\boldsymbol{\omega}) = \begin{pmatrix} \boldsymbol{\omega} \\ LD(\boldsymbol{\omega})_s \end{pmatrix},$$

resultando na curvatura normal segundo uma direção unitária  $\mathbf{d}$  dada por:

$$C(\tilde{\boldsymbol{\xi}}_1)_\mathbf{d} = 2 \left| \mathbf{d}^\top \mathbf{F}(\tilde{\boldsymbol{\xi}}_1) \mathbf{d} \right|,$$

para  $\mathbf{F}(\tilde{\boldsymbol{\xi}}_1) = \Delta^\top (\ddot{\mathcal{F}}_\alpha^{-1} - \mathbf{B}_{22}) \Delta$ , em que

$$\mathbf{B}_{22} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddot{\mathcal{F}}_{\alpha_{22}}^{-1} \end{pmatrix},$$

com  $\ddot{\mathcal{F}}_{\alpha_{22}}^{-1}$  sendo a inversa da matriz  $\ddot{\mathcal{F}}_{\alpha_{22}}$  obtida a partir da matriz particionada

$$\ddot{\mathcal{F}}_{\alpha} = \begin{pmatrix} \ddot{\mathcal{F}}_{\alpha_{11}} & \ddot{\mathcal{F}}_{\alpha_{12}} \\ \ddot{\mathcal{F}}_{\alpha_{21}} & \ddot{\mathcal{F}}_{\alpha_{22}} \end{pmatrix}.$$

Assim sendo, considerando a máxima curvatura  $\mathbf{d}_{max}$  obtida pelo maior autovetor associado ao maior autovalor da matriz  $\mathbf{F}(\tilde{\boldsymbol{\xi}}_1)$ , pode-se utilizar a direção do  $i$ -ésimo indivíduo tal como definido em (3.8) para revelar observações que exercem forte influência nas estimativas nos valores ajustados de  $\tilde{\boldsymbol{\xi}}_1$ .

## 4 APLICAÇÕES

Neste capítulo, por intermédio de estudos a alguns bancos de dados reais, busca-se colocar em prática a teoria desenvolvida ao longo desta pesquisa. Para tanto, foi usado o pacote MGCV implementado no R por Simon Wood. Sua documentação pode ser encontrada em <https://cran.r-project.org/web/packages/mgcv/mgcv.pdf>. Para mais detalhes sobre a teoria que envolve os modelos parcialmente lineares aditivos generalizados e sua base computacional, sugere-se consultar o livro de Wood (2017).

### 4.1 ELEIÇÕES DE 2006 AO SENADO BRASILEIRO

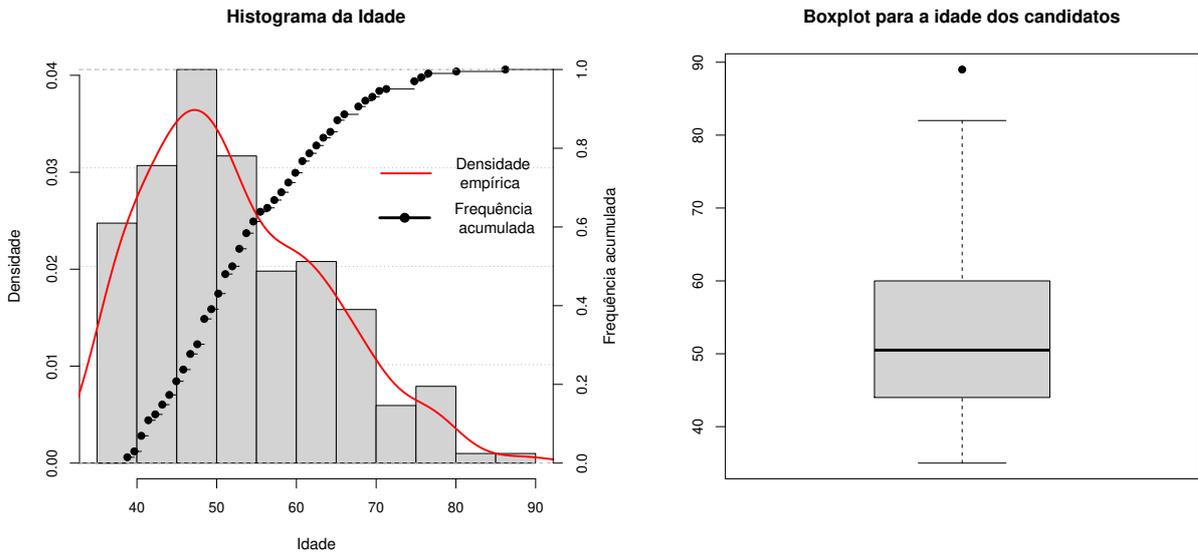
Como primeira aplicação, apresentamos uma base de dados contendo informações de candidatos ao Senado do Brasil nas eleições de 2006. No referido banco de dados, encontram-se informações tais como: O Estado do candidato (UF), o partido, quantidade de votos recebido, o sexo, o resultado da eleição (eleito ou não eleito), dentre outras. No entanto, foram consideradas para análise apenas as seguintes variáveis:

1. `log.votos`: O logaritmo dos votos. Variável resposta;
2. `log.vgastos`: O logaritmo dos valores gastos por cada candidato;
3. `escola`: Esta variável representa o nível de escolaridade e foi recodificada de modo a reduzir o número de categorias. Sendo assim, foi criada a variável `instr` (*nível de instrução*), a qual encontra-se dividida em duas categorias: 0 = Baixo e 1 = Alto. Na categoria "Baixo" estão desde pessoas alfabetizadas até pessoas com ensino médio completo. Na categoria "Alto" estão pessoas com ensino superior incompleto até ensino superior completo;
4. `idade`: Corresponde a idade do candidato;
5. `est.civil`: Estado civil: 0 = casado, 1 = divorciado, 2 = separado, 3 = solteiro e 4=viúvo;
6. `sexo`: 0 = masculino e 1 = feminino.

#### 4.1.1 Análise exploratória dos dados

Este conjunto de dados reúne 202 observações referentes aos candidatos ao Senado brasileiro no ano de 2006. Dentre os partidos como o maior número de candidatos destaca-se o PSOL com 17, PFL com 15, PDT com 14, PSDB e PSTU ambos com 13, PCB e PMDB 12 e PT com 10. Os candidatos dessa edição possuíam uma média de idade de 52.25 anos, sendo a menor delas de 35 e a maior de 89 anos. Destaca-se também que, do total dos candidatos, 30

eram do sexo feminino e 117 do sexo masculino. A maioria casados ou, pelo menos, já foram casados, bem como a maioria possuindo no mínimo ensino superior incompleto. Este é perfil geral dos candidatos. Contudo, uma análise mais aprofundada dessas variáveis será apresentada na sequência.

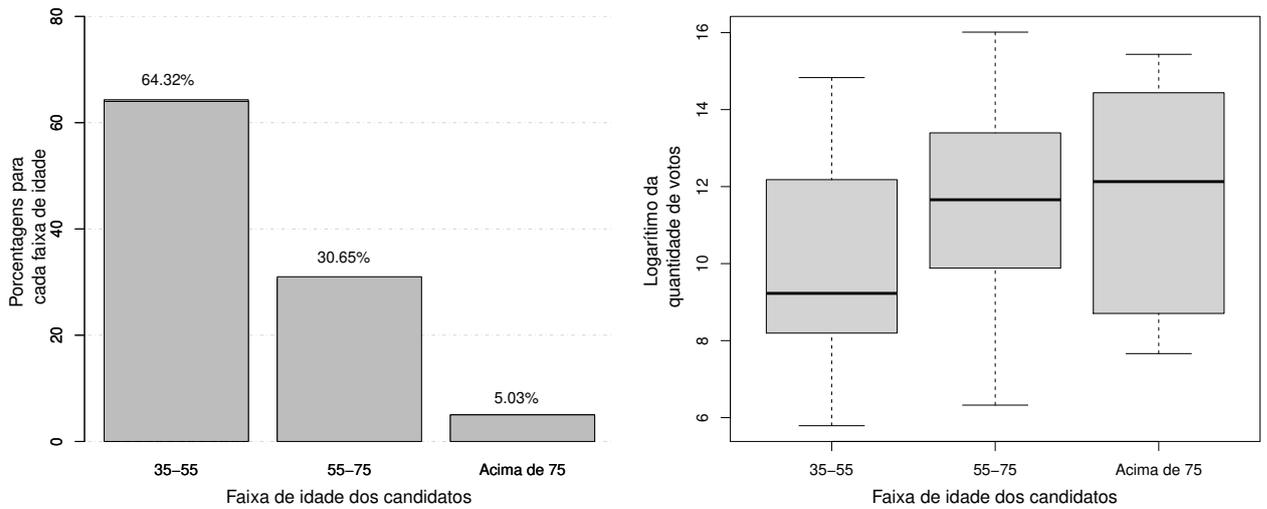


**Figura 12 – Gráficos da distribuição de frequência, frequência acumulada e boxplot respectivos da a variável idade para o exemplo das eleições de 2006 ao Senado brasileiro.**

Por meio do gráfico da distribuição, bem como pelo boxplot da idade dos candidatos apresentados na **Figura 12**, pode-se observar que esta variável possui comportamento assimétrico na direção das maiores idades. Além disso, o gráfico da frequência acumulada mostra outros quantis para a distribuição desta variável. Por meio dele pode-se observar que aproximadamente 93% dos candidatos possuem idades menores que 70 anos.

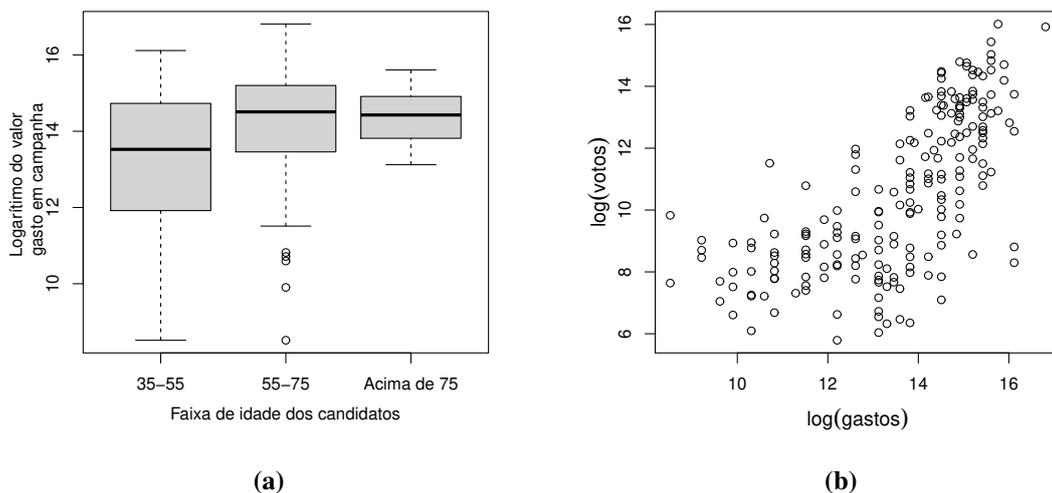
O gráfico das porcentagens para cada faixa de idade dos candidatos apresentado na **Figura 13**, em essência, possui o mesmo significado dos gráficos anteriores. Porém, em conjunto com o gráfico que relaciona os boxplot de cada faixa de idade com suas respectivas quantidades de votos recebidas, também apresentado na **Figura 13**, traz uma informação importante: mesmo possuindo a menor proporção, os candidatos acima de 75 anos conseguiram obter uma expressiva quantidade de votos, com a mediana maior que a dos candidatos mais jovens.

Por outro lado, como informação complementar, o gráfico da **Figura 14a** mostra que estes candidatos também foram os que mais gastaram em suas campanhas. Além disso, como os candidatos mais jovens gastaram menos em suas campanhas, há indícios de que possivelmente os



**Figura 13 – Gráficos com as porcentagens para cada faixa de idade e relação da idade com a quantidade de votos recebidas pelos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro.**

gastos em campanha possa ter refletido na quantidade de votos recebidas. Neste caso, o gráfico de dispersão do logaritmo dos votos contra o logaritmo dos valores gastos apresentado na **Figura 14b** pode dar uma noção intuitiva dessa relação.

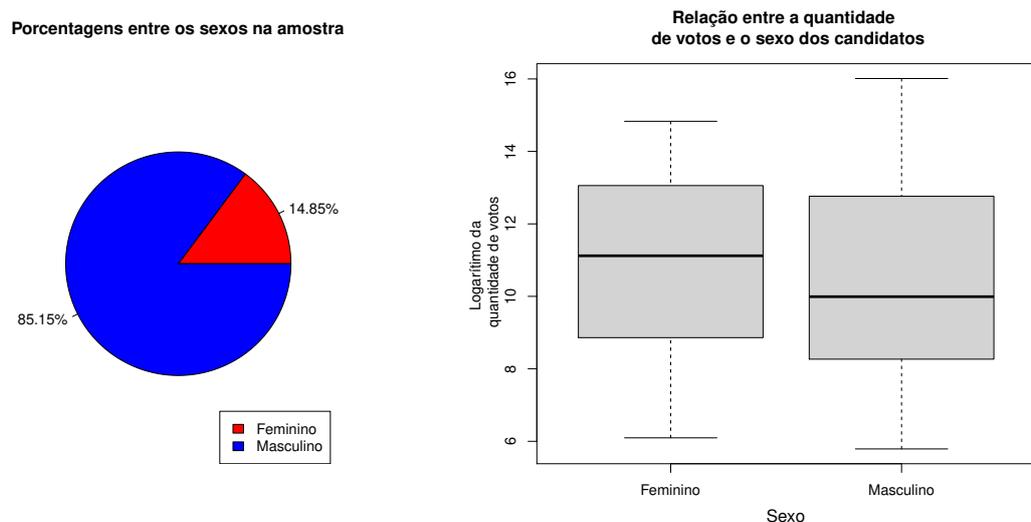


**Figura 14 – Estes gráficos mostram a relação da idade com os gastos e diagrama de dispersão: votos segundo os gastos na campanha dos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro.**

Portanto, é possível perceber uma maior concentração dos pontos nas interseções entre baixo gasto e baixo número de votos e entre alto gasto e alta votação. Ou seja, parece haver uma relação de causa e efeito entre essas variáveis, indicando que gastos maiores nas eleições

levam quantidades maiores de votos. Porém, é necessário uma análise de regressão a fim de identificar a correlação entre essas variáveis e sua significância estatística.

Tomando agora a variável sexo dos candidatos, nota-se que, embora representem apenas 14.85% da amostra, o sexo feminino consegue obter mediana de votos superior aos do sexo masculino. Porém, como evidenciado na **Tabela 1**, a maioria dos candidatos eleitos são do sexo masculino. Neste caso, é natural o questionamento sobre o fator sexo como potencialmente influente nos resultados.



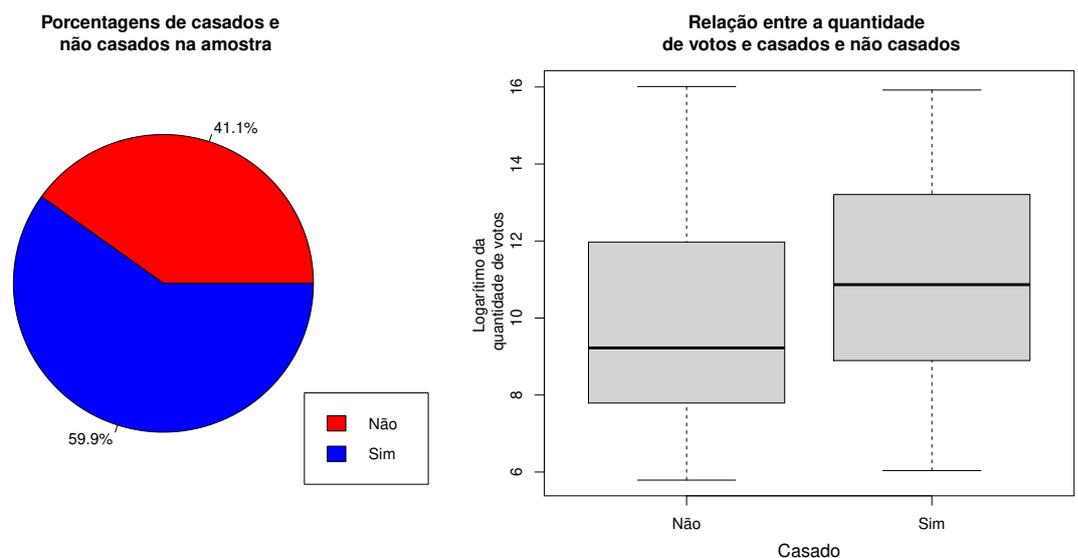
**Figura 15 – Gráficos com as porcentagens dos sexo na amostra e relação da sexo dos candidatos com as quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro.**

Sexo	Resultado		Total
	Eleitos	Não eleitos	
Feminino	4	26	30
Masculino	23	149	172
Total	27	175	202

**Tabela 1 – Resultado das eleições de 2006 ao Senado brasileiro.**

Estas constatações, no entanto, carecem de mais evidências. Considerando a normalidade dos dados, por exemplo, poderia ser aplicado um teste de hipóteses para confirmar se as médias dos votos podem ser consideradas estatisticamente diferentes em relação ao sexo ou se essa diferença é apenas resultado do processo aleatório que envolveu a obtenção dessa amostra em particular. Além disso, um teste qui-quadrado também poderia ser útil para verificar se o resultado das eleições depende do sexo do candidato.

Como o número de categorias do estado civil é relativamente grande, optou-se por reduzi-la, a fim de obter uma ideia mais clara da interação dessa variável com as demais. Neste caso, dividiu-se essa variável em duas categorias: quem era casado e que não era. Neste caso, quando a variável estado civil é relacionada com o logaritmo dos votos, percebe-se que: a proporção de casados e solteiros na amostra não difere muito, bem como os casados possuem mediana de votos superior em relação aos não casados.



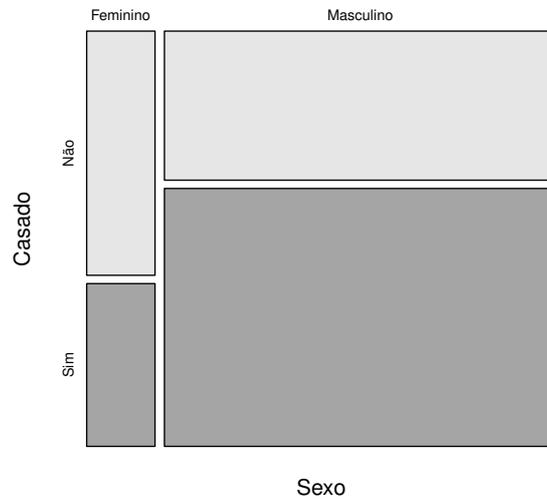
**Figura 16 – Gráficos com as porcentagens dos sexo na amostra e relação da sexo dos candidatos com as quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro.**

Além disso, apresentamos a **Tabela 2** de dupla entrada com as frequências absolutas e o gráfico da **Figura 17** que mostram quantos pontos do conjunto de dados se encaixam em cada categoria. Por meio da **Figura 17** é possível visualizar o relacionamento entre essas variáveis graficamente. Em particular, verifica-se que a maioria dos candidatos do sexo masculino são casados e o oposto ocorre com sexo feminino.

Casado	Sexo		Total
	Feminino	Masculino	
Não	18	63	81
Sim	12	109	121
Total	30	172	202

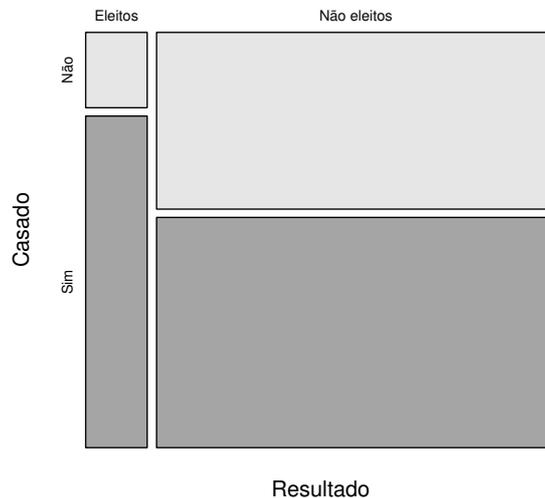
**Tabela 2 – Resultado das eleições de 2006 ao Senado brasileiro segundo o estado civil.**

Assim como realizado para a variável sexo, é interessante investigar se o fato de um candidato ser casado ou não pode ter influência nos resultados. Para isso, usamos o recurso



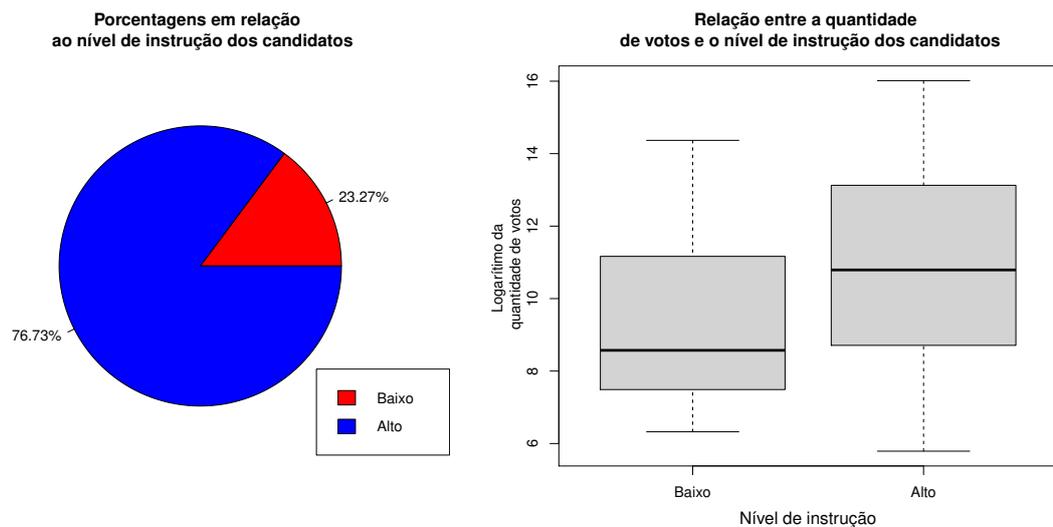
**Figura 17 – Gráfico com a relação do sexo com o estado civil dos candidatos para o exemplo das eleições de 2006 ao Senado brasileiro.**

gráfico apresentado na **Figura 18**. Portanto, pode-se identificar que os candidatos casados foram maioria dos eleitos. Neste caso também é interessante testes para diferenças entre médias e teste qui-quadrado a fim de obter resultados mais precisos para estas considerações.



**Figura 18 – Resultado das eleições de 2006 ao Senado brasileiro segundo o estado civil.**

Por fim, analisando o nível de instrução e a relação com a quantidade de votos, nota-se que há uma maior concentração de votos para os candidatos que possuem ensino superior incompleto ou completo.



**Figura 19 – Gráficos com as porcentagens dos níveis de instrução na amostra e relação das quantidades de votos por eles recebidas para o exemplo das eleições de 2006 ao Senado brasileiro**

Por meio da **Tabela 3**, nota-se que, corroborando com o gráfico da **Figura 19**, há forte indício de que um nível de educação alto pode ter efeito sobre a quantidade de votos, bem como nos resultados das eleições.

Nível de Instrução	Resultado		Total
	Eleitos	Não eleitos	
Baixo	4	43	47
Alto	23	132	155
Total	27	175	202

**Tabela 3 – Resultado das eleições de 2006 ao Senado brasileiro segundo o nível de instrução.**

Note que as análises foram feitas de forma individual, sem que fosse considerado o efeito conjunto destas variáveis no valor médio da quantidade dos votos. O que queremos é uma forma de acomodar esta informação em um modelo capaz explicar a relação presente em cada uma das covariáveis, mostrando o impacto conjuntamente dessas variáveis na quantidade de votos. Para isso, iremos propor um modelo de regressão apropriado.

#### 4.1.2 Ajuste do modelo

Para a análise desses dados, considerou-se um modelo de regressão com distribuição normal para a variável resposta. Sendo assim, definindo  $Y_i, i = 1, 2, \dots, n$ , como sendo a variável *logaritmo dos votos* para a  $i$ -ésima observação, então temos que  $Y_i \sim \mathcal{N}(\mu_i, \phi = \sigma^2)$ , de modo

que estabelecemos um MPLAG com preditor linear dado por:

$$\begin{aligned}
 g(\mu_i) &= \mu_i = \eta_i \\
 &= \beta_0 + \beta_1 \text{est.civil}_1 + \beta_2 \text{est.civil}_2 + \beta_3 \text{est.civil}_3 + \beta_4 \text{est.civil}_4 \\
 &+ \beta_5 \text{instr} + \beta_6 \text{sexo} + \beta_7 \text{idade} + f(t_i), \text{ para } i = 1, 2, \dots, n,
 \end{aligned} \tag{4.1}$$

com função de ligação *identidade*,  $n = 202$ ,  $f(\cdot)$  uma função suave para a variável  $\log.vgastos$  e variável dependente dada pelo logaritmo do número de votos recebidos pelos candidatos ao Senado em 2006.

Assim, para o ajuste deste modelo, foram empregados um *B-spline* cúbico com 10 nós equidistantes e penalização quadrática. Os graus de liberdade efetivos são obtidos pela soma dos graus de liberdade referente a função não paramétrica mais oito, referindo-se ao intercepto e os  $\beta$ 's de cada variável ajustada parametricamente. Os resultados desse ajuste estão resumidos na **Tabela 4**.

**Tabela 4 – Resumo das estimativas referentes ao modelo ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	8.738	0.655	< 0.0001
Idade	0.032	0.011	< 0.005
instrAlto	1.139	0.271	< 0.0001
est.civil <sub>1</sub>	-0.297	0.335	0.376
est.civil <sub>2</sub>	-0.353	0.421	0.403
est.civil <sub>3</sub>	-0.336	0.360	0.352
est.civil <sub>4</sub>	-0.875	0.536	0.104
Masculino	-0.726	0.325	< 0.01
Outras medidas do ajuste			
$df(\alpha)$	12.66		
$\alpha$	3.497		
$\phi$	2.399		
VCG	2.571		

Nesse modelo, observa-se que todas as categorias da variável estado civil têm uma alta probabilidade de estar apenas por acaso correlacionadas com a variável dependente, pois apresentam valor  $p > 0.05$ . Isto é, na possibilidade de poder repetir o processo de amostragem indefinidas vezes, supondo que a variável estado civil não tenha efeito sobre o resultado das eleições, a probabilidade de se obter uma amostra que indique que há uma correlação entre essas variáveis é elevada.

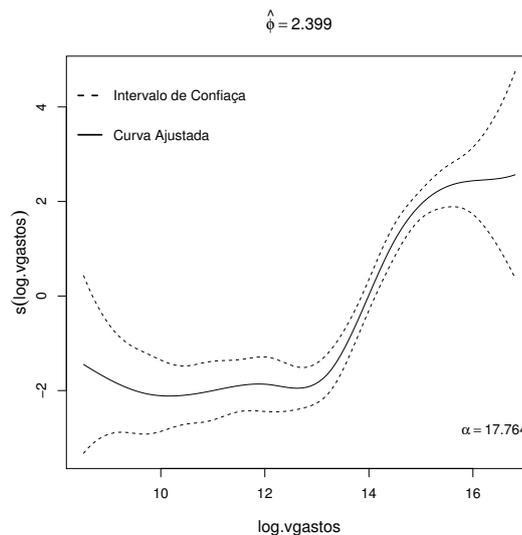
Tendo em vista que em um modelo de regressão somente devem ser adicionadas variáveis cuja correlação com a variável dependente é elevada e, principalmente, teoricamente

explicável, então esta variável deve ser retirada da análise. Assim, mantendo apenas aquelas com maior contribuição. Portanto, realizou-se um ajuste sem essa variável e o resultado foi que a variável sexo também deixou de ser significativa. Com isso, foi obtido o modelo final, o qual pode ser observado na **Tabela 5**. Graficamente, o subsequente ajuste da variável suavizada pode ser visualizado pela **Figura 20**.

**Tabela 5 – Resumo das estimativas referentes ao modelo final ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	8.148	0.574	< 0.0001
Idade	0.028	0.010	< 0.001
instrAlto	1.175	0.271	< 0.0001
Outras medidas do ajuste			
$df(\alpha)$	8.69		
$\alpha$	3.082		
$\phi$	2.425		
VCG	2.534		

Nota-se que a curva se adéqua muito bem a tendência dos dados. Assim, podemos identificar melhor a relação entre as variáveis quantidade de votos e valor gasto em campanha, bem como podemos observar que há uma maior variabilidade do número de votos em regiões de baixo e elevados gastos, destacado pelos intervalos de confiança pontuais.



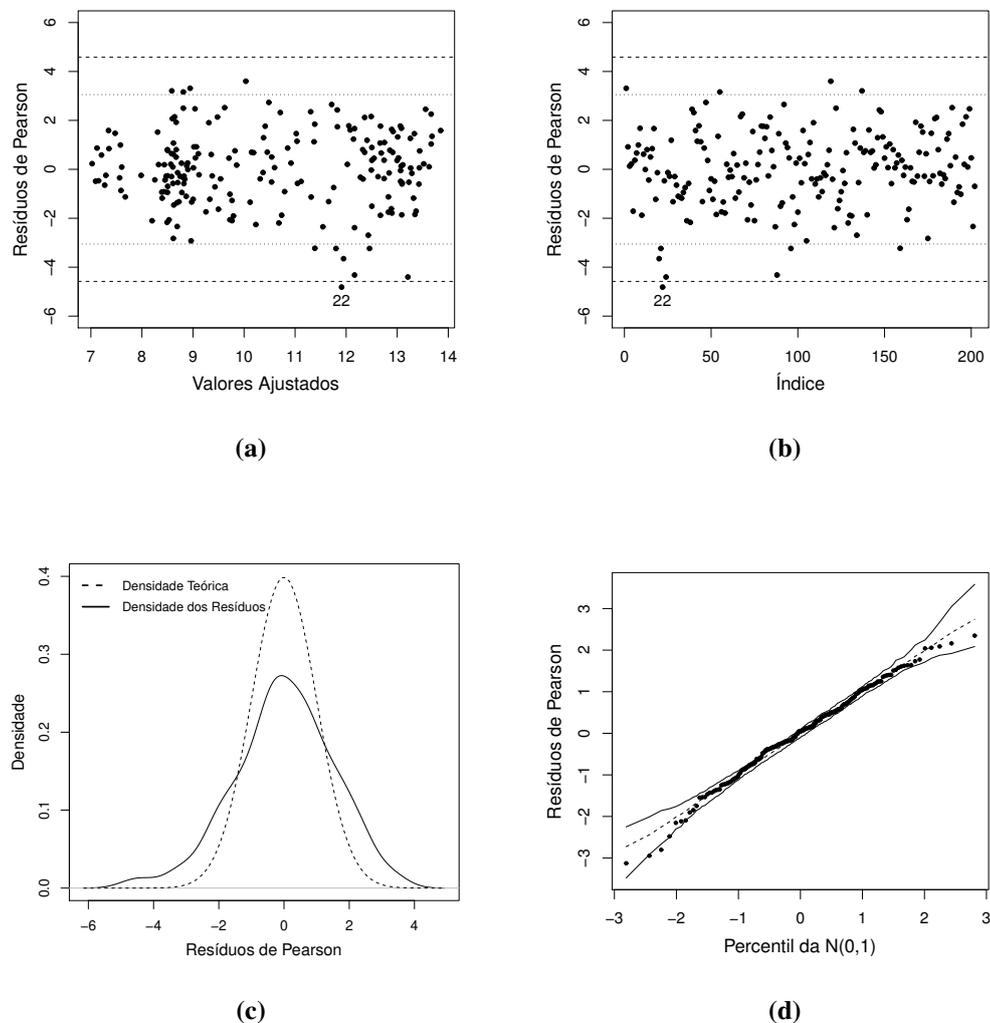
**Figura 20 – Gráfico da curva ajustada sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.**

Na sequência, dar-se-á continuidade ao estudo desse modelo mediante o seu diagnóstico, colocando em prática a teoria desenvolvida ao longo do texto para esta finalidade, começando

pela análise de resíduos, seguida pela alavancagem e finalizando pela análise de influência local.

### 4.1.3 Diagnóstico

Para avaliar a adequacidade do ajuste para o modelo proposto, inicialmente deve-se observar os gráficos da **Figura 21**, em que podemos analisar os resíduos. Por meio destes gráficos, pode-se tirar conclusões a respeito do componente estocástico do modelo. Assim, qualquer comportamento atípico referente aos resíduos pode ser um indicativo de que o modelo proposto não é adequado para representar o fenômeno em estudo.



**Figura 21 – Gráficos de resíduos referente ao modelo ajustado sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.**

Desta forma, o gráfico dos resíduos versus valores ajustados (**Figura 21a**) permite verificar a homoscedasticidade do modelo, isto é,  $\sigma^2$  constante, além de possibilitar identificar

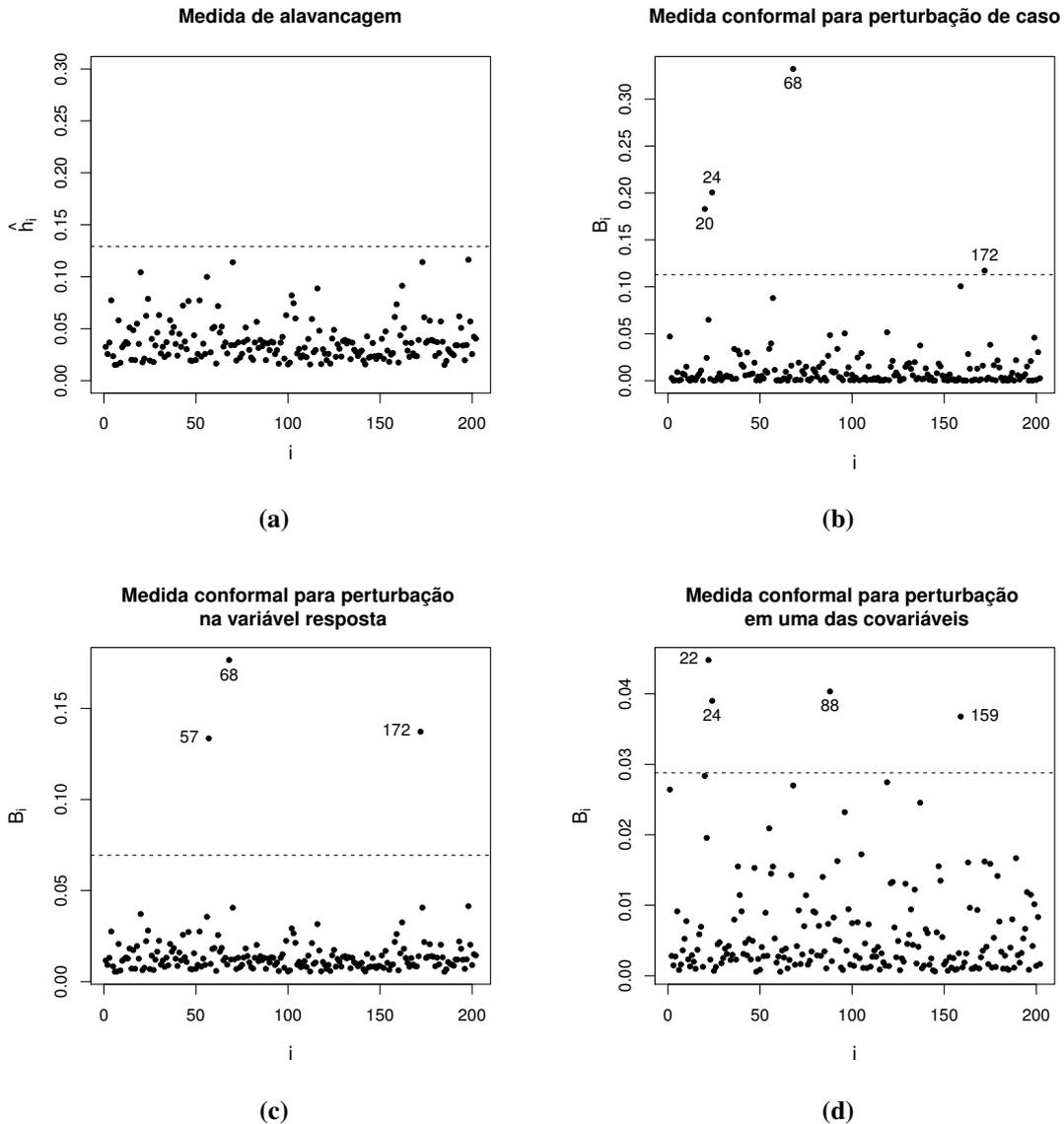
outliers em  $y$ . Assim, espera-se que este gráfico apresente os pontos distribuídos aleatoriamente em torno de zero. Se os pontos possuem alguma tendência, há indícios de que a variância dos dados não é constante, ou seja, existe indicativo de heteroscedasticidade. Já o gráfico dos resíduos versus a ordem de coleta dos dados (**Figura 21b**) avalia a hipótese de independência dos dados. Portanto, se os pontos tiverem um comportamento que se repete em determinado ponto do gráfico, temos indícios de dependência dos resíduos. Por fim, os gráficos 21c e 21d fornecem informações sobre a distribuição dos resíduos. Neste caso, espera-se um comportamento próximo a de uma distribuição normal com média zero e desvio padrão igual a um.

Portanto, é possível observar que o gráfico de resíduos quantílicos contra os valores ajustados de fato apresentam comportamento aleatório em torno de zero, indicando uma variabilidade constante dos resíduos, sendo que alguns dos pontos estão dois desvios padrão distantes da média, mas não muito distantes de três desvios padrão (linhas tracejadas), podendo representar outliers. Algumas explicações possíveis para a obtenção desses pontos, seriam erros sistemáticos, ou erro de especificação para o modelo, ou que estes pontos estão no extremo de validade para a distribuição da variável resposta. O mesmo ocorre para os resíduos versus índices, os quais não apresentam comportamento serial aparente. O único resíduo destacado foi o correspondente a observação #22 como um possível valor aberrante, porém apresenta uma leve discrepância. Além disso, os gráficos inferiores da **Figura 21** corroboram com a hipótese de normalidade dos dados, embora a densidade dos resíduos possua uma variância um pouco maior em relação a densidade teórica esperada. Com base nessa análise de resíduos, verifica-se que as suposições para o modelo estão razoavelmente satisfeitas.

Segundo o gráfico de alavancagem, destacado pela **Figura 22a**, não há pontos considerados com alta alavancagem. Ainda tomando o modelo definido anteriormente, considerando a análise de influência local sob os esquemas de perturbação de casos, perturbação na variável resposta e perturbação em uma das variáveis explicativas, tem-se que para os pontos destacados pelas **Figura 22b** e **22c** são considerado pontos potencialmente influentes, que podem estar exercendo um peso desproporcional nas estimativas do modelo.

Desse modo, é necessário avaliar como as estimativas dos parâmetros do modelo estão sendo afetadas por essas observações, tentando verificar se existe interferências substanciais, bem como mudanças inferenciais. Neste caso, deve-se ter uma atenção especial sobre estas observações, investigando quais são suas características e quais as possíveis respostas para entender o porquê de estarem influenciando de forma tão crítica nas estimativas. Assim, deve-se retirar individualmente cada uma dessas observações, realizar um novo ajuste e saber qual o

impacto da retirada dessa observação potencialmente influente nas novas estimativas.



**Figura 22 – Gráficos de alavancagem e influência local sob os esquemas de perturbação de caso, perturbação na variável resposta e perturbação em umas das variáveis explicativas, respectivamente, referentes ao ajuste do modelo sob distribuição normal para o exemplo das eleições de 2006 ao Senado brasileiro.**

Nota-se, de acordo com os gráficos, que as observações potencialmente influentes foram: #20, #22, #24, #57, #68, #88, #159 e #172. A fim de analisar o impacto das observações destacadas sob as estimativas dos parâmetros, realizamos uma análise de sensibilidade destas observações. Como já destacado, esta análise confirmatória consiste em reajustar o modelo, eliminando individualmente cada observação considerada possivelmente influente e, em seguida,

calcular as taxas de variação das estimativas dos parâmetros segundo a expressão:

$$\text{Taxa}(\theta) = \left| \frac{\tilde{\theta}^{(i)} - \hat{\theta}}{\hat{\theta}} \right| \times 100, \quad (4.2)$$

em que  $\hat{\theta}$  e  $\tilde{\theta}^{(i)}$  são as EMV no modelo completo e no modelo sem a  $i$ -ésima observação respectivamente. Com isso, obteve-se a tabela a **Tabela 6**, em que temos as observações retiradas seguidas dos valores das estimativas e as respectivas mudanças relativas em porcentagem.

**Tabela 6 – Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição normal e mudança relativa em porcentagem entre parênteses para os parâmetros  $\beta_0$ ,  $\beta_1$  e  $\beta_2$  após a retirada dos pontos para o exemplo das eleições de 2006 ao senado brasileiro.**

Observação retirada	Estimativas dos parâmetros		
	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$
#20	8.30 (1.92%)	0.03 (4.55%)	1.07 (9.20%)
#22	8.19 (0.56%)	0.03 (4.16%)	1.22 (3.70%)
#24	8.23 (1.00%)	0.03 (5.74%)	1.19 (1.40%)
#57	8.19 (0.48%)	0.03 (3.03%)	1.20 (2.19%)
#68	8.18 (0.38%)	0.03 (4.01%)	1.18 (0.46%)
#88	8.10 (0.62%)	0.03 (2.83%)	1.21 (2.75%)
#159	7.87 (3.41%)	0.03 (19.48%)	1.19 (1.17%)
#172	8.18 (0.40%)	0.03 (4.92%)	1.23 (4.32%)
#68 e #172	8.21 (0.80%)	0.03 (9.04%)	1.23 (4.74%)

O valor de referência para comparar a mudança relativa das estimativas após a retirada de um ponto potencialmente influente é  $(p/n) \times 100$ , em que  $p$  é a quantidade de parâmetros e  $n$  o tamanho da amostra. Sendo assim, espera-se que na retirada de uma observação a mudança relativa seja de aproximadamente 0.01%. Então, baseado na **Tabela 6**, temos que as observações causam impactos desproporcionais principalmente nas estimativas de  $\beta_1$  e  $\beta_2$ . Porém, não houve mudança inferencial com a retirada destas observações.

Por fim, deve-se investigar os motivos pelos quais as observações apresentadas na **Tabela 6** se destacam das demais. Assim, por exemplo, para a observação #20, tem-se que é um candidato do sexo masculino, possui 48 anos e ensino médio completo (nível de instrução

baixa). Além disso, teve o quarto maior gasto na campanha, porém tendo obtido apenas 4013 (ou  $\log(4013) = 8.3$ ) votos, ficando entre os 50 candidatos menos votados. Neste caso, há um contraste em relação a tendência apresentada no modelo.

Já para a observação #22, também acontece algo semelhante, pois tem-se que é um candidato do sexo masculino, com 48 anos e nível superior completo (nível de instrução alto). Além disso, teve gastos elevado na campanha (logaritmo igual 14.50), porém ficando também entre os 15 candidatos com menores votos recebidos.

Tomando agora a observação #159, tem-se que é um candidato do sexo feminino, possui 89 anos (a maior do banco de dados), com nível superior completo (nível de instrução alto). Além disso, teve gastos elevado na campanha (logaritmo igual 13.82), porém ficando também entre os 50 candidatos com menores votos recebidos.

## 4.2 ALUGUEL EM MUNIQUE

Neste segundo exemplo será utilizado o banco de dados disponível no R denominado *rent*. Tal banco de dados é sobre uma pesquisa realizada em abril de 1993, com o objetivo de estudar o valor do aluguel em Munique. Para este exemplo, serão consideradas as seguintes variáveis:

1. R: Valor mensal do Aluguel (em euros). Variável resposta;
2. Fl: Tamanho do apartamento (em metros quadrados);
3. A: Ano de construção (período);
4. loc: Qualidade da Localização (1: Média, 2: Boa e 3: Excelente);
5. B: Qualidade do banheiro (0: Padrão e 1: Premium) e
6. L: Qualidade da Cozinha (0: Padrão e 1: Premium).

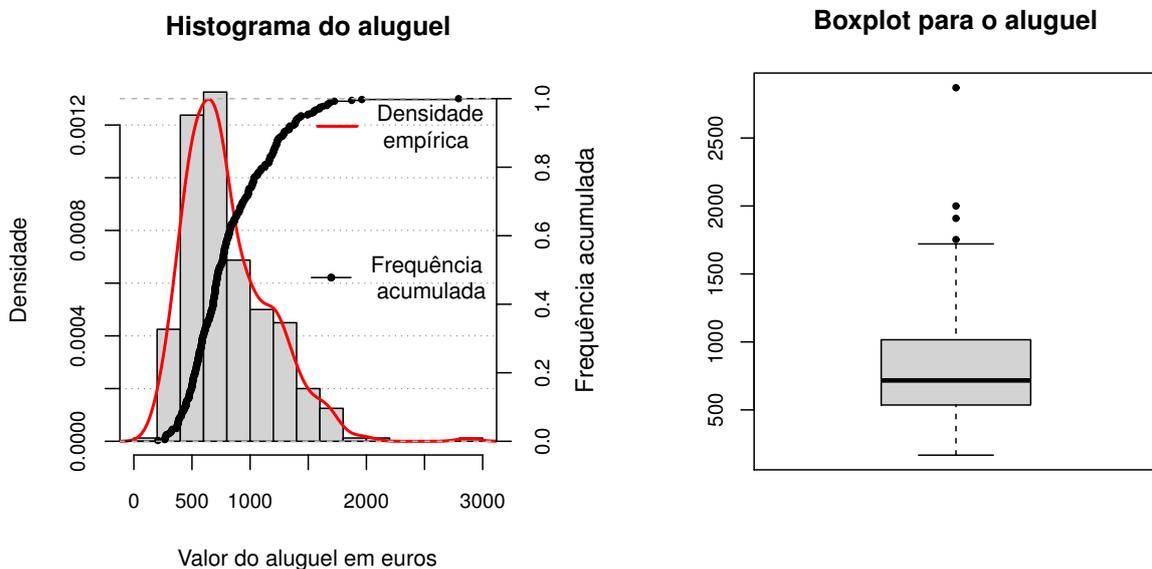
Antes de propor um modelo para explicar a relação das variáveis explicativas com a média do aluguel mensal em Munique, será apresentada uma análise descritivas dos dados. O objetivo é verificar e resumir as principais características dos dados para obter um entendimento básico sobre as relações existentes entre as variáveis analisadas.

### 4.2.1 Análise exploratória dos dados

Inicialmente, será descrito o comportamento de cada variável segundo medidas resumo, tais como: média, desvio padrão, mediana, etc., afim de obter uma ideia básica sobre o comportamento dos dados. Além disso, para melhor entendimento, serão apresentados métodos

visuais, tais como gráficos, bem como comentários sobre aspectos particulares das análises.

Para este conjunto de dados, a variável *valor mensal do aluguel* possui as seguintes características: o valor médio do aluguel mensal foi de  $\bar{x} = 804.4 \text{ €}$ , com desvio padrão igual a  $s = 364.708 \text{ €}$  e, conseqüentemente, coeficiente de variação dado por  $cv = 45.34\%$ . Além disso, os valores do menor e maior aluguel foram de  $x_{(1)} = 167.2 \text{ €}$  e  $x_{(n)} = 2869.9 \text{ €}$ , respectivamente, bem como os primeiros e terceiros quartis dados por  $Q_1 = 536.5 \text{ €}$  e  $Q_3 = 1015.8 \text{ €}$  com mediana igual a  $Md = 716.5 \text{ €}$ . A **Figura 23** apresentam um resumo gráfico dessas características.

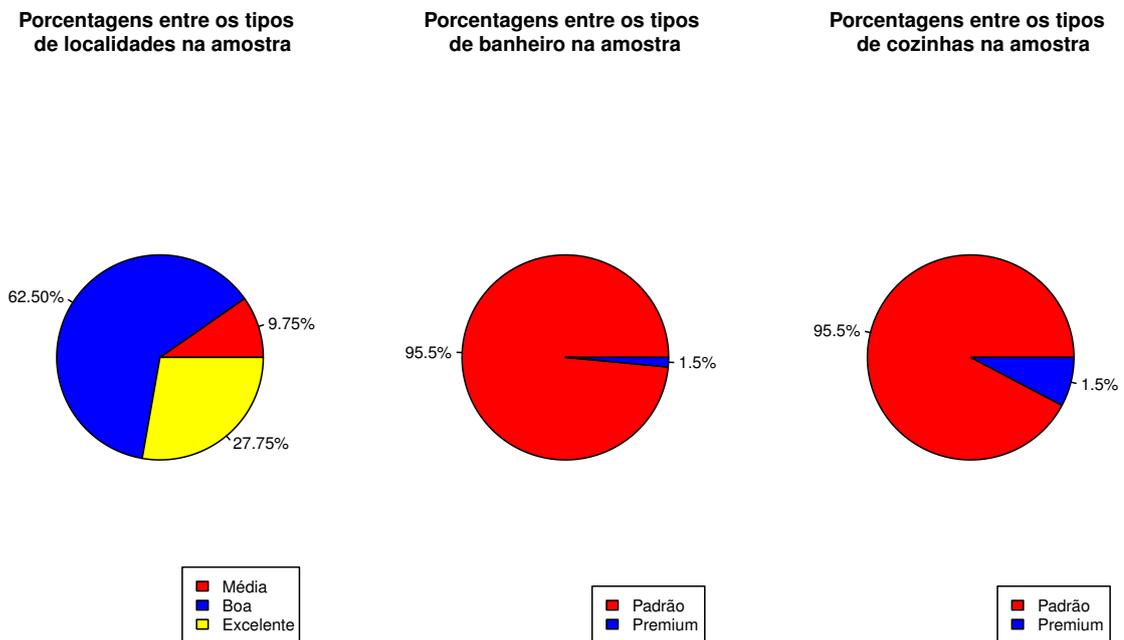


**Figura 23 – Gráficos da distribuição de frequência, frequência acumulada e boxplot respectivos da a variável preço do aluguel para o exemplo do aluguel em Munique.**

Pode-se observar que a distribuição do valor mensal do aluguel apresenta uma assimetria à direita, bem como quatro pontos que se destacam das demais observações. Os pontos destacados são: #299, #236, #312, e #26, que correspondem respectivamente aos seguintes valores dos aluguéis: 1753.2;1909.0;2000.0 e 2869.9. A esses aluguéis estão associadas as seguintes variáveis: para a observação #299, correspondente ao aluguel 1753.2 – apartamento com 83 metros quadrados, para o ano de 1957, com excelente localização, além de banheiro padrão e cozinha premium; para a observação #236, correspondente ao aluguel 1909.2 – apartamento com 104 metros quadrados, para o ano de 1981, com boa localização, além de banheiro padrão e cozinha premium; para a observação #312, correspondente ao aluguel 2000.0 – apartamento com 114 metros quadrados, para o ano de 1985, com excelente localização,

além de banheiro e cozinha padrão; para a observação #26, correspondente ao aluguel 2869.9 — apartamento com 90 metros quadrados, para o ano de 1985, com excelente localização, além de banheiro e cozinha padrão.

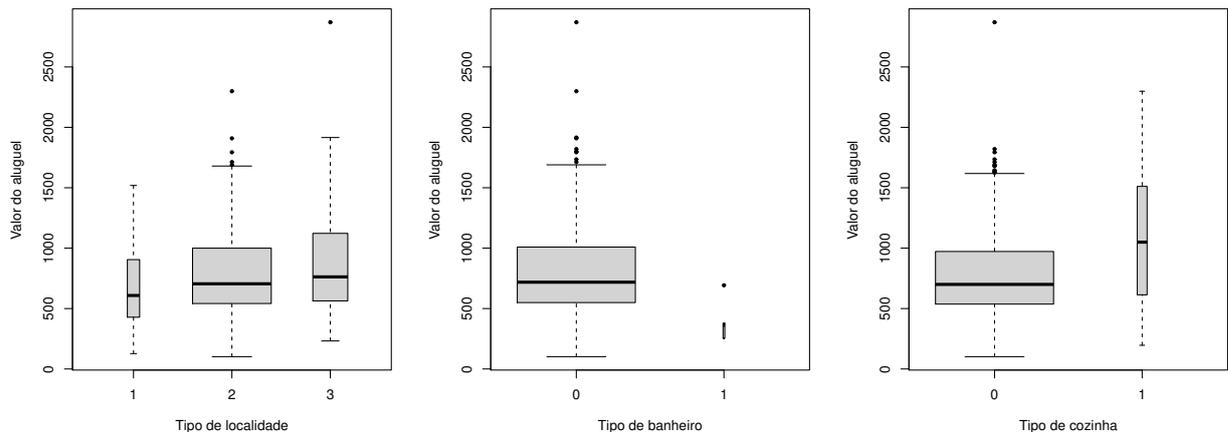
Também foi observado que as variáveis categóricas possuem suas proporções na amostra como apresentadas nos gráficos da **Figura 24**. Desta forma, pode-se concluir que na amostra os tipos mais comuns de localidade, banheiro e cozinha são, respectivamente, boa e padrão. Além disso, segundo os boxplots apresentados na **Figura 25**, nota-se que os preços dos alugueis aumentam quando a qualidade da localidade e da cozinha melhoram.



**Figura 24 – Gráficos com as proporções de cada variável categórica na amostra para o exemplo do aluguel em Munique.**

Os alugueis médios, segundo os tipos de localidade, foram de 595.40 € com desvio padrão igual a 169.97 € para o tipo de localidade média; 792.23 € com desvio padrão igual a 348.00 € para o tipo de localidade boa; 884.32 € com desvio padrão igual a 411.25 € para o tipo de localidade excelente. Conseqüentemente, os coeficientes de variação são: 28.54%, 43.92% e 46.5% respectivamente. Conclui-se, portanto, que a localidade com denominação média corresponde a categoria mais homogênea em relação as demais. De igual forma, temos que o aluguel médio segundo o tipo de banheiro foram de 811.13 € com desvio padrão igual a

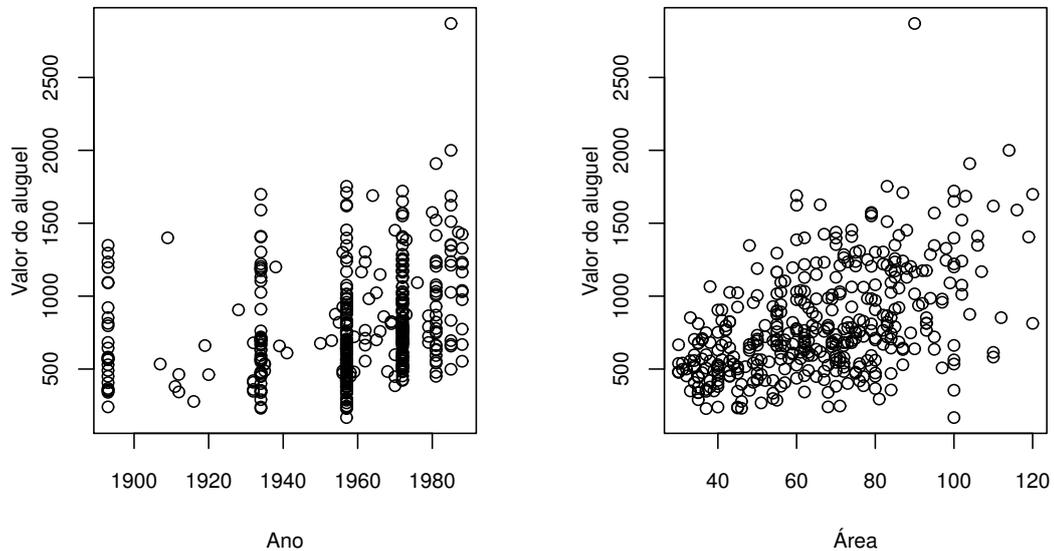
363.25 € para o tipo padrão; 362.76 € com desvio padrão igual a 73.40 € para o tipo premium. Consequentemente, os coeficientes de variação são: 44.78% e 20.23%, respectivamente. Finalmente, temos que o aluguel médio segundo o tipo de cozinha foram de 778.29 € com desvio padrão igual a 349.65 € para o tipo padrão; 1085.47 € com desvio padrão igual a 408.66 € para o tipo premium. Consequentemente, os coeficientes de variação são: 44.92% e 50.80%, respectivamente.



**Figura 25 – Boxplots do valor do aluguel segundo as variáveis categóricas para o exemplo do aluguel em Munique.**

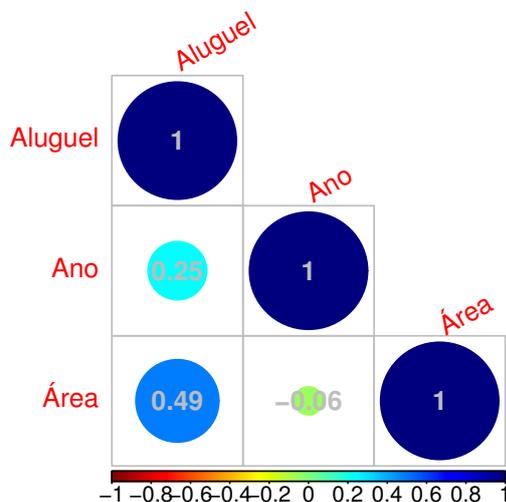
Além disso, a variável *tamanho do apartamento* possui as seguintes características: tamanho médio de  $\bar{x} = 65.98\text{m}^2$ , com desvio padrão igual a  $s = 20.50\text{m}^2$  e, conseqüentemente, coeficiente de variação dado por  $cv = 31.08\%$ . Os valores da menor e maior área foram de  $x_{(1)} = 30.00\text{m}^2$  e  $x_{(n)} = 120.00\text{m}^2$ , respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = 50.0\text{m}^2$  e  $Q_3 = 80.0\text{m}^2$  com mediana igual a  $\text{Md} = 65.50\text{m}^2$ . Portanto, os tamanhos dos apartamentos são aproximadamente homogêneos. Por fim, apresentamos os gráficos de dispersão, **Figura 26** e medidas da correlação entre as variáveis numéricas, **Figura 27**. Esses gráficos buscam dar uma ideia visual da relação presente entre as variáveis.

Assim, nota-se que, por meio do gráfico de dispersão, parece que o valor do aluguel aumenta com o passar dos anos e com o tamanho do apartamento. Além disso, de acordo com o gráfico da **Figura 27**, há uma correlação moderada entre o valor do aluguel e as variáveis ano e área do apartamento. Logo, há indícios de possível relação de causalidade entre essas variáveis. Contudo, apenas análises estatísticas mais avançadas, tais como testes de hipóteses ou modelos de regressão, poderá dizer a natureza aproximada dessa relação.



**Figura 26 – Diagramas de dispersão do valor do aluguel contra as variáveis numéricas para o exemplo do aluguel em Munique.**

Note que as análises foram feitas de forma individual. Isto é, observa-se cada variável de forma separada, desconsiderando o efeito que a presença de outras variáveis desempenham no valor médio do aluguel. Em outras palavras, essas análises não levam em conta a relação da média do aluguel segundo a presença conjunta dessas variáveis. Isso só será possível por meio de um modelo de regressão apropriado, que mostre o impacto conjuntamente dessas variáveis no preço do aluguel em Munique.



**Figura 27 – Gráfico das correlações entre as variáveis quantitativas para o exemplo do aluguel em Munique.**

Diante de todo o exposto, o objetivo daqui em diante passa a ser encontrar um modelo de regressão apropriado de modo a acomodar o efeito simultâneo dessas variáveis. Por fim, realizar a análise de diagnóstico, a fim de obter indícios de afastamento das suposições para o modelo, bem como fazer inferências a respeito dos parâmetros sob o modelo postulado.

#### 4.2.2 Ajuste do modelo

Visto que a variável resposta, *valor do aluguel* ( $R$ ), é contínua e assume valores nos reais positivos, então faz sentido assumir uma distribuição gama como parte estocástica do modelo. Sendo assim, definindo  $Y_i$ ,  $i = 1, 2, \dots, n$  como sendo o *valor do aluguel* para a  $i$ -ésima observação, então temos que  $Y_i \sim G(\mu_i, \phi)$ , de modo que a parte sistemática será dada por:

$$\begin{aligned} g(\mu_i) &= \mu_i^{-1} = \eta_i \\ &= \beta_1 + \beta_2 loc_{2i} + \beta_3 loc_{3i} + \beta_4 L_i + \beta_5 B_i + f_1(A) + f_2(Fl), \end{aligned}$$

em que  $\mu_i = \mathbb{E}(R_i)$ ,  $f_k$ , com  $k = 1, 2$  são funções suaves das covariáveis *ano de construção* ( $A$ ) e *tamanho do apartamento* ( $Fl$ ), além de uma amostra de tamanho  $n = 400$ . Tem-se também que a função de ligação adotada será a *inversa* e os resultados do ajuste podem ser conferidos na **Tabela 7** com seus respectivos erros padrão e valor  $p$ . Assim, destaca-se que as variáveis

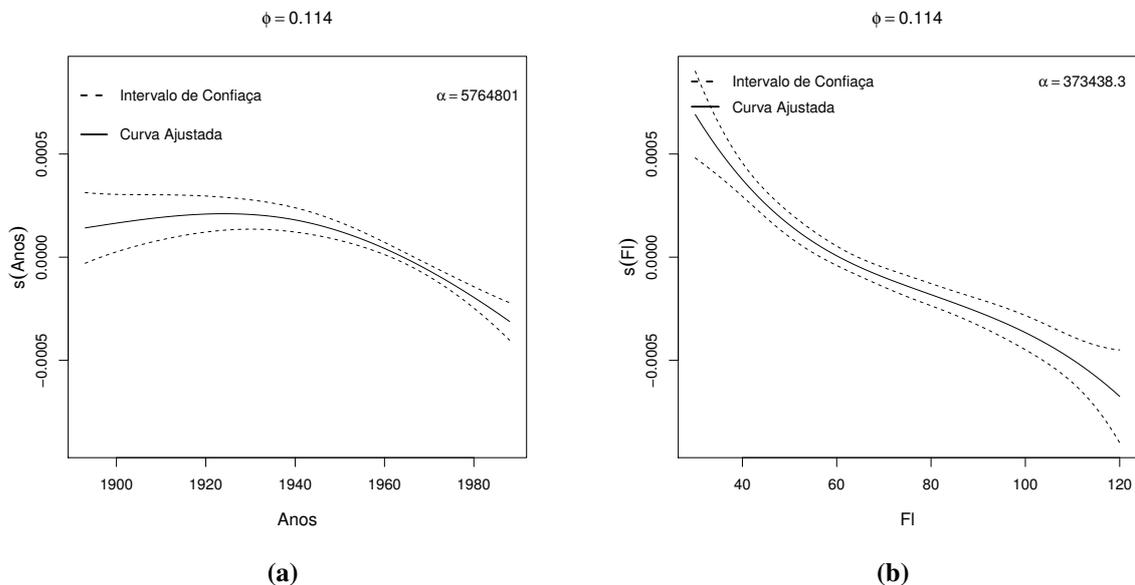
**Tabela 7 – Resumo das estimativas de máxima verosimilhança referentes ao modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	0.0017	0.00010	< 0.0001
loc2	-0.00034	0.00011	< 0.005
loc3	-0.00048	0.00011	< 0.001
B1	0.0011	0.00038	< 0.001
L1	-0.00013	0.00005	< 0.05
Outras medidas do ajuste			
$df_1(\alpha)$	2.10		
$df_2(\alpha)$	3.08		
$\alpha$	(5764801.0, 373438.3)		
$\phi$	0.114		
VCG	0.120		

"A" e "Fl" foram ambas ajustadas segundo um *P-spline* considerando um *B-spline* cúbico com 5 nós equidistantes para cada suavizador e diferenças de segunda ordem. Neste caso, observou-se que uma quantidade de nós acima de 5 aumentava o valor do critério VCG. Além disso, os

respectivos parâmetros de suavização foram obtidos por meio do método VCG. Observa-se que todas as variáveis ajustadas parametricamente foram significativas para o modelo.

Com relação aos graus de liberdade efetivos, nota-se que  $df(\alpha) = 5 + df_1(\alpha) + df_2(\alpha)$ , em que 5, por definição, corresponde aos graus de liberdade referente a parte paramétrica e  $df_k(\alpha)$  correspondendo a parte não paramétrica do modelo. Neste caso, tem-se que  $df_1(\alpha) = 2.10$  para a variável "A" e  $df_2(\alpha) = 3.08$  para a variável "FI", em que, respectivamente,  $\alpha_1 = 5,764,801$  e  $\alpha_2 = 373,438.3$ . Além disso, foi obtida uma pontuação VCG de 0.1203. As estimativas estão resumidas na **Tabela 7**. Na sequência, os gráficos do ajuste para as variáveis suavizadas são apresentados na **Figura 28**. Como a função de ligação utilizada foi a inversa,



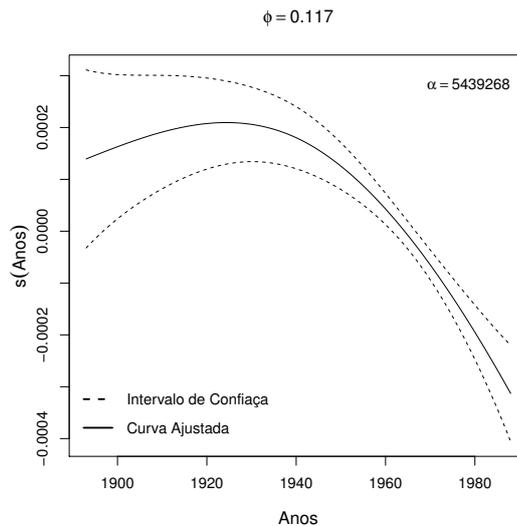
**Figura 28 – Gráficos das curvas ajustadas sob distribuição gama para o exemplo do aluguel em Munique.**

então significa que a relação entre as variáveis se dá de maneira inversa. Desta forma, como era de se esperar, a tendência é que os preços médios dos aluguéis aumentem de acordo com o tamanho do imóvel, o que é evidenciado pelo gráfico da **Figura 28b**. Note que, no entanto, esta variável se relaciona com o valor médio do aluguel de maneira aproximadamente linear, sendo necessário reavaliar a possibilidade de realizar um novo ajuste, porém considerando esta variável de forma paramétrica. Sendo assim, foi realizado outro ajuste, agora seguindo esta mudança. Isto é, considerar que a variável *Fl* seja ajustada de forma paramétrica. Os resultados podem ser consultado na **Tabela 8**.

**Tabela 8 – Resumo das estimativas de máxima verossimilhança referentes ao segundo modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	0.0024	0.000130	< 0.0001
F1	-0.00001	0.000001	< 0.0001
loc2	-0.00033	0.000108	< 0.001
loc3	-0.00046	0.000100	< 0.001
B1	0.00116	0.000400	< 0.01
L1	-0.00011	0.000059	< 0.05
Outras medidas do ajuste			
$df(\alpha)$	2.12		
$\alpha$	5439268		
$\phi$	0.117		
VCG	0.122		

Desta forma, ao observar o gráfico da **Figura 29**, pode-se notar que, de 1900 à 1940 os preços dos alugueis apresentaram queda. Esse período foi marcado por crises, recessões e guerras na Europa, prováveis fatores que levaram a esta redução dos alugueis em Munique. Após as destruições provocadas neste conturbado período, foi implantado pelos Estados Unidos o plano de recuperação da Europa, também conhecido como Plano Marshall. Este plano consistia em uma política de ajuda econômica para a reconstrução dos países aliados, tais como a Alemanha Ocidental. Com isto, a economia de Munique pôde se recuperar, retomando, dentre outros aspectos da economia, o aquecimento do mercado imobiliário.



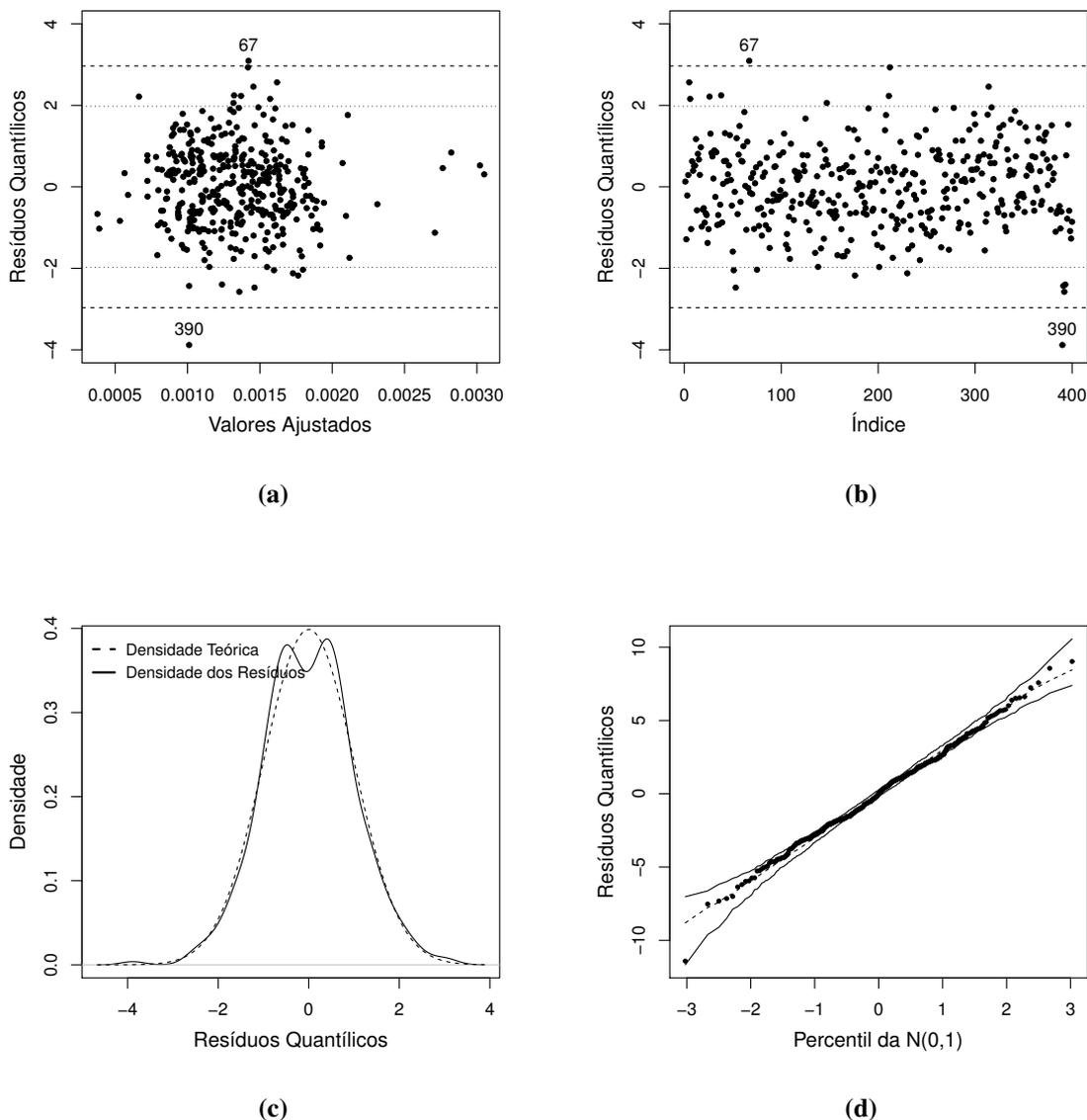
**Figura 29 – Gráfico das curvas ajustadas para o segundo modelo sob distribuição gama para o exemplo do aluguel em Munique.**

Com base no exposto, por meio da análise gráfica, nota-se que o modelo é um forte candidato a uma boa representação para estudarmos a relação do aluguel médio em função das variáveis explicativas. A seguir, será dado o seu diagnóstico a fim de confirmar a viabilidade do

modelo proposto.

### 4.2.3 Diagnóstico

A análise de resíduos para este modelo encontra-se na **Figura 30**. Verifica-se que o gráfico dos resíduos contra os valores ajustados apresentam comportamento aleatório e nenhum indício de tendência. De modo semelhante, o gráfico dos resíduos contra os índices também não apresenta comportamento atípico ao esperado. Além disso, apenas alguns poucos pontos encontram-se distantes dois desvios padrão da média, mas os resíduos associados as observações #67 e #390 ultrapassam 3 desvios padrão.



**Figura 30 – Gráficos de resíduos referente ao modelo ajustado sob distribuição gama para o exemplo do aluguel em Munique.**

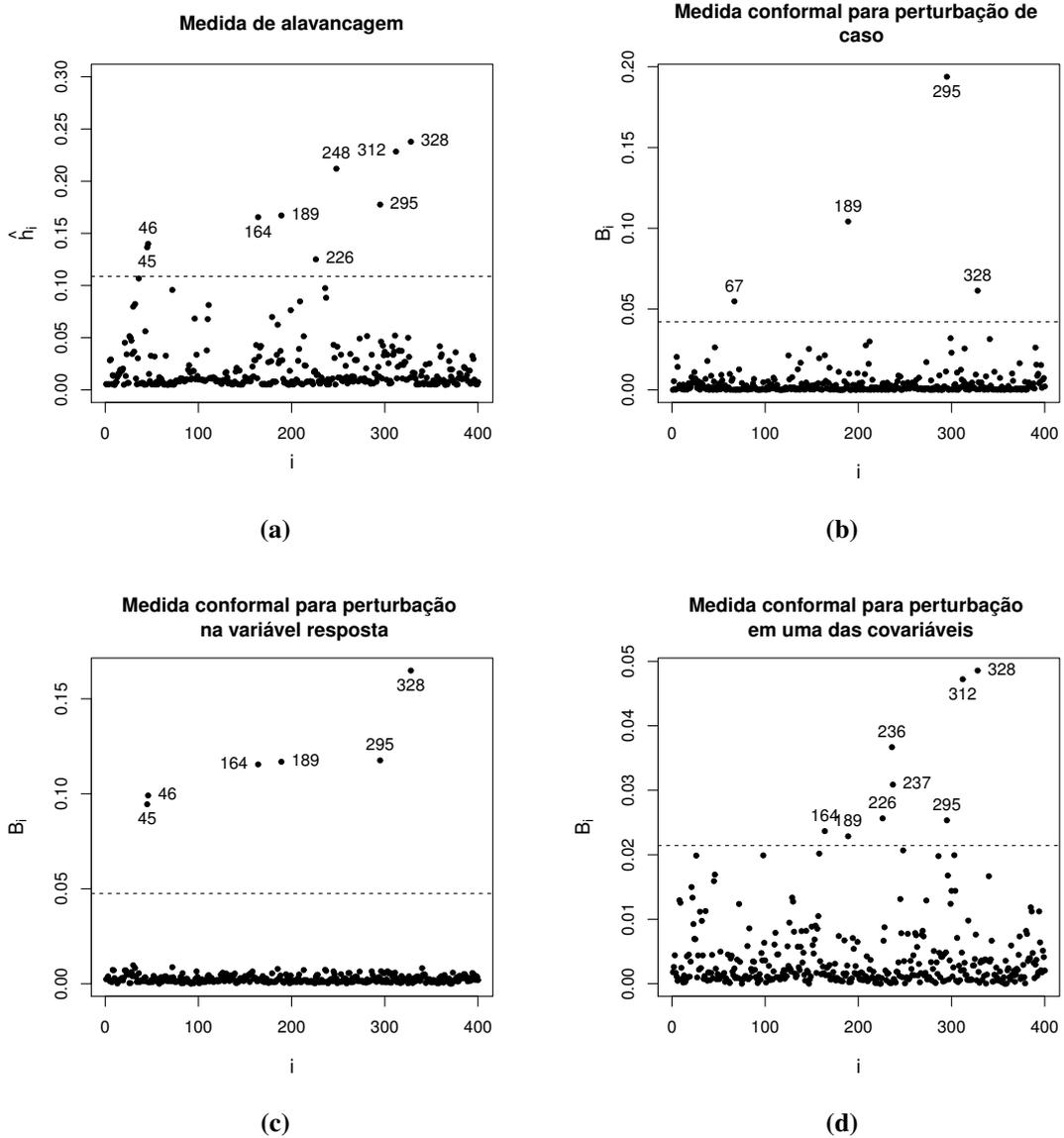
A observação #67 corresponde ao aluguel é de 1689.7 (um dos alugueis mais caros presentes no banco de dados), do ano de 1964 cuja área é de 60 metros quadrados. Dos que possuem esse tamanho, essa é a observação com o aluguel mais caro. Além disso, possui boa localidade, banheiro e cozinha padrão. Por outro lado, a observação #390 corresponde ao aluguel de 167.2 (o menor da mostra), com área de 100 metros quadrados, também uma das maiores áreas presentes na amostra. Além disso, possui boa localidade, cozinha e banheiro padrão, para o ano de 1957.

Como visto anteriormente, se o modelo for adequando, os resíduos quantílicos aleatorizados possuem uma distribuição próxima a de uma normal com média zero e desvio padrão igual a um. Logo, os gráficos inferiores da **Figura 30** dizem respeito a essa característica. Portanto, espera-se que, para um ajuste adequado, os resíduos observados contra os teóricos devem ficar concentrados dentro do envelope simulado e a densidade dos resíduos quantílicos devem se aproximar de uma normal padrão a medida que o tamanho da amostra cresce.

Dessa maneira, pode-se observar que a densidade dos resíduos apresentada está bem próxima a densidade teórica esperada, além de nenhum dos pontos estarem fora do envelope. Ou seja, não há indicações de que a função de ligação é inadequada. Portanto, segundo a análise de resíduo, o modelo proposto não apresenta desvios graves de suposição, sendo apto ao uso. Porém, seguiremos a análise por meio das medidas de alavancagem e influência local para finalizar o diagnóstico.

Segundo o gráfico de alavancagem, destacam-se na **Figura (31a)** as observações que são candidatas a pontos de alavancagem, ou aberrantes. Portanto, é necessário uma análise mais cuidadosa destas observações. Considerando a análise de influência sob os esquemas de perturbação de caso (**Figura 31b**), perturbação na variável resposta (**Figura 31c**) e perturbação em uma das variáveis explicativas (**Figura 31d**), considera-se como observações potencialmente influentes aquelas distantes da linha tracejada, que corresponde a média mais três vezes o desvio padrão de  $B_i$ . Desta forma, destacam-se no gráfico os pontos com suposta influência nas estimativas dos parâmetros, sendo necessário realizar uma análise mais cuidadosa dessas observações.

Portanto, assim como desenvolvido no exemplo anterior, devemos remover essas observações e realizar novo ajuste. Em seguida, calcular a taxa de variação das estimativas deste ajuste em relação ao ajuste com todas as observações, segundo a expressão 4.2 e observar se as estimativas mudaram de maneira muito drástica, bem como se houve mudança inferencial. Com isso, obteve-se a **Tabela 9**, com as estimativas do modelo sem as observações destacadas e



**Figura 31 – Gráficos de alavancagem e influência local sob os esquemas de perturbação de caso, perturbação na variável resposta e perturbação em umas das variáveis explicativas, respectivamente, referentes ao ajuste do modelo sob distribuição gama para o exemplo do aluguel em Munique.**

mudança relativa das estimativas.

O valor de referência para comparar a mudança relativa das estimativas após a retirada de um ponto potencialmente influente é de 1.5%. Sendo assim, espera-se que na retirada de uma observação a mudança relativa seja de aproximadamente essa porcentagem, se esta não causa influência nas estimativas. Então, na **Tabela 9** temos que as observações causam impactos desproporcionais, principalmente nas estimativas de  $\beta_4$ , correspondente ao tipo de banheiro, especialmente quando foram retiradas as observações #189, #295 e #328, bem como nas estimativas de  $\beta_5$ , correspondente a qualidade da cozinha. Em particular, quando é retirada a observação #328, a estimativa desse parâmetro aumentar 10.6%. Além disso, quando as

observações #236 e #312 foram retiradas do ajuste, houve mudança inferencial ao nível de 5% para este mesmo parâmetro.

**Tabela 9 – Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição gama e mudança relativa em porcentagem entre parênteses para os parâmetros  $\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5$  e  $\beta_6$  após a retirada dos pontos para o exemplo do aluguel em Munique.**

Observação retirada	Estimativas dos parâmetros					
	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$
#164	$2.46 \times 10^{-3}$ (0.04%)	$1.19 \times 10^{-5}$ (0.02%)	$3.33 \times 10^{-4}$ (0.28%)	$4.62 \times 10^{-4}$ (0.13%)	$1.2 \times 10^{-3}$ (5.87%)	$1.19 \times 10^{-4}$ (0.13%)
#189	$2.46 \times 10^{-3}$ (0.3%)	$1.19 \times 10^{-5}$ (0.15%)	$3.43 \times 10^{-4}$ (2.84%)	$4.72 \times 10^{-4}$ (2.05%)	$1.32 \times 10^{-3}$ (12.91%)	$1.19 \times 10^{-4}$ (0.49%)
#236	$2.46 \times 10^{-3}$ (0.07%)	$1.19 \times 10^{-5}$ (0.18%)	$3.32 \times 10^{-4}$ (0.37%)	$4.62 \times 10^{-4}$ (0.05%)	$1.17 \times 10^{-3}$ (0.12%)	$1.15 \times 10^{-4}$ (2.71%)
#295	$2.47 \times 10^{-3}$ (0.09%)	$1.19 \times 10^{-5}$ (0.03%)	$3.36 \times 10^{-4}$ (0.75%)	$4.65 \times 10^{-4}$ (0.68%)	$9.69 \times 10^{-4}$ (17.07%)	$1.19 \times 10^{-4}$ (0.01%)
#312	$2.48 \times 10^{-3}$ (0.56%)	$1.12 \times 10^{-5}$ (2.07%)	$3.27 \times 10^{-4}$ (1.94%)	$4.66 \times 10^{-4}$ (0.82%)	$1.17 \times 10^{-3}$ (0.17%)	$1.06 \times 10^{-4}$ (10.60%)
#328	$2.47 \times 10^{-3}$ (0.09%)	$1.19 \times 10^{-5}$ (0.14%)	$3.34 \times 10^{-4}$ (0.32%)	$4.64 \times 10^{-4}$ (0.43%)	$1.05 \times 10^{-3}$ (10.16%)	$1.18 \times 10^{-4}$ (0.31%)

Sendo assim, foi destacado que a observação #189 corresponde ao aluguel de 450.4 € com cozinha padrão e banheiro premium, dentre os quais o aluguel é o segundo maior. Por outro lado, a observação #236 corresponde a um apartamento com banheiro padrão e cozinha premium cujo aluguel é de 1909 €, sendo este o maior valor dentre os apartamentos com cozinha premium. Para a observação #295, foi constatado que corresponde a um apartamento com banheiro e cozinha padrão, porém em uma ótima localidade, bem como um aluguel de 233.0 €. O quarto aluguel mais baixo da amostra. Com relação a observação #312, temos que corresponde a um apartamento com banheiro e cozinhas padrão, porém em uma ótima localidade cujo aluguel, avaliado no ano de 1985, foi de 2000.0 €. O segundo maior valor da amostra. Por fim, tomando a observação #328, constatou-se que se trata de um apartamento com cozinha padrão, ótima localidade e banheiro premium, dentre os quais possui o maior aluguel.

### 4.3 AR POLUÍDO EM CHICAGO

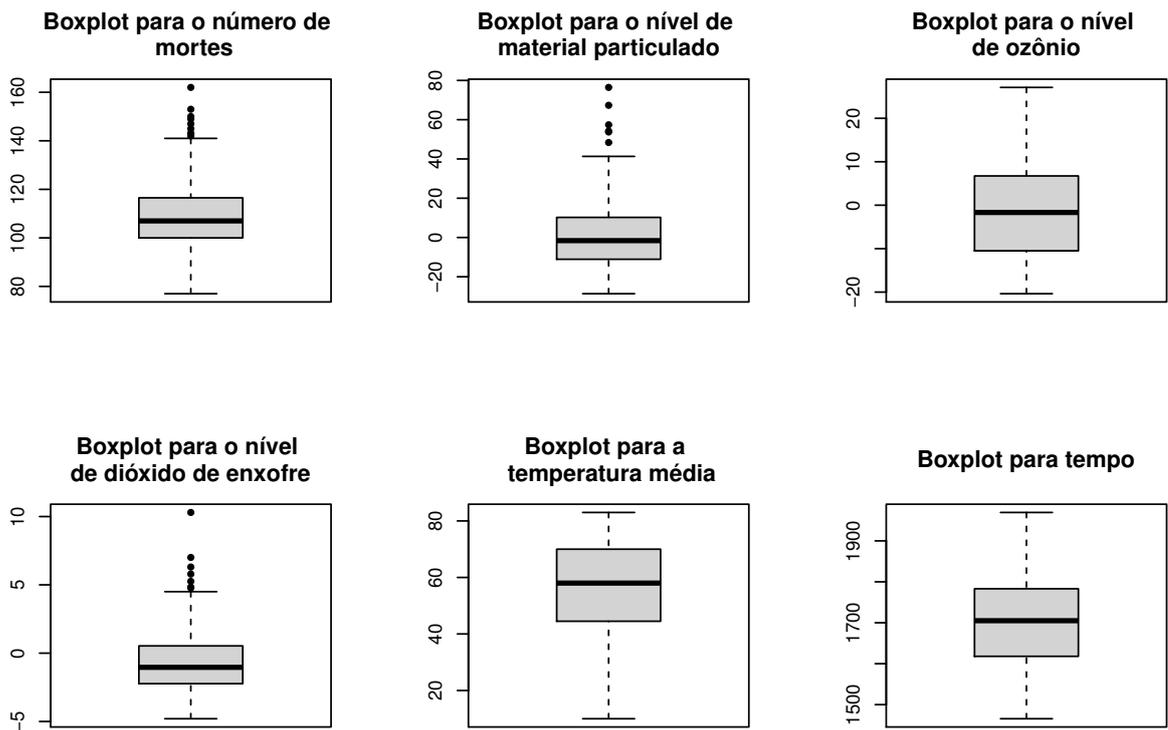
Nesta aplicação, considera-se um conjunto de dados apresentado em Wood (2017, p. 243) disponível no pacote *gamair* do R e que relaciona a taxa de mortes diárias, ao longo de vários anos, com os níveis de ozônio, níveis de dióxido de enxofre, temperatura média diária e níveis de material particulado<sup>1</sup>. Além dessas variáveis, a taxa de mortalidade tende a variar

<sup>1</sup> São partículas muito finas de sólidos ou líquidos suspensos no ar, como as que são liberados pelo escapamento de motores a diesel, por exemplo.

conforme o passar do tempo, em especial ao longo do ano, por razões que têm pouco ou nada a ver com a qualidade do ar conforme observa Wood (2017, p. 243). Assim, o objetivo desta aplicação é apresentar uma análise estatística, por meio de um modelo parcialmente linear generalizado para dados de contagem, bem como apresentar as medidas de diagnósticos a fim de verificar a qualidade do ajuste, pondo em prática as técnicas desenvolvidas ao longo das seções anteriores.

#### 4.3.1 Análise exploratória dos dados

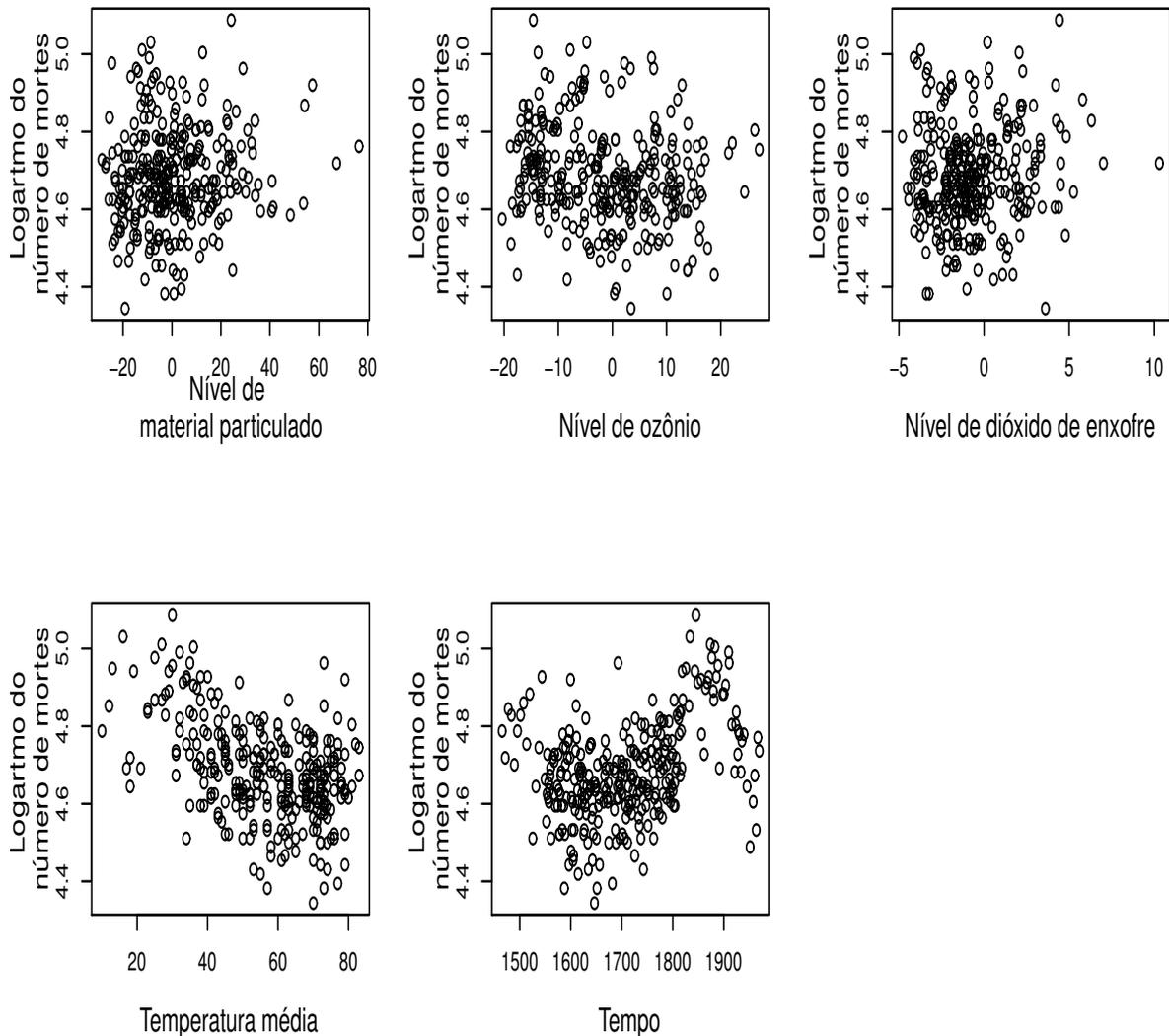
Este conjunto de dados reúne uma amostra de 300 observações. Observamos alguns boxplots que facilitam a visualização dos dados. Nota-se que algumas variáveis têm um ou mais pontos discrepantes, sendo necessário investigar o porquê do comportamento atípico.



**Figura 32 – Boxplots das variáveis presentes no baco de dados considerando o exemplo do ar poluído em Chicago.**

Cada uma das variáveis possuem as seguintes características: O *número de mortes diária* possui valor médio dado por  $\bar{x} = 109.1$  (por dia), com desvio padrão igual a  $s = 14.17$  (por dia). Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = 77$  (por dia) e  $x_{(n)} = 162$  (por dia), respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = 100$  (por dia) e  $Q_3 = 116.2$  (por dia) com mediana igual a  $Md = 107.0$  (por dia); O *material particulado* possui valor médio dado por  $\bar{x} = 1.03 \text{ mg/m}^3$ , com desvio padrão igual a  $s = 17.17 \text{ mg/m}^3$ . Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = -28.63 \text{ mg/m}^3$  e  $x_{(n)} = 76.43 \text{ mg/m}^3$ , respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = -11.082 \text{ mg/m}^3$  e  $Q_3 = 10.179 \text{ mg/m}^3$  com mediana igual a  $Md = -1.60 \text{ mg/m}^3$ ; o *nível de ozônio* possui valor médio dado por  $\bar{x} = -1.41$  ppb, com desvio padrão igual a  $s = 10.402$  ppb. Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = -20.36$  ppb e  $x_{(n)} = 27.11$  ppb, respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = -10.49$  ppb e  $Q_3 = 6.690$  ppb com mediana igual a  $Md = -1.68$  ppb; a *temperatura média diária* possui média  $\bar{x} = 56.18^\circ\text{F}$ , com desvio padrão igual a  $s = 16.28^\circ\text{F}$ . Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = 10.00^\circ\text{F}$  e  $x_{(n)} = 83.00^\circ\text{F}$ , respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = 44.75^\circ\text{F}$  e  $Q_3 = 70.00^\circ\text{F}$  com mediana igual a  $Md = 58.00^\circ\text{F}$ ; o nível de *dióxido de enxofre* possui valor médio dado por  $\bar{x} = -0.6634$  ppb, com desvio padrão igual a  $s = 2.35$  ppb. Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = -4.80$  ppb e  $x_{(n)} = 10.29$  ppb, respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = -2.22$  ppb e  $Q_3 = 0.51$  ppb com mediana igual a  $Md = -1.04$  ppb e, finalmente, o tempo em dias possui o valor médio dado por  $\bar{x} = 1708$  (por dia), com desvio padrão igual a  $s = 112.39$  (por dia). Além disso, os valores correspondentes ao menor e maior valor para esta variável foram de  $x_{(1)} = 1466$  (por dia) e  $x_{(n)} = 1970$  (por dia), respectivamente, bem como os primeiro e terceiros quartis dados por  $Q_1 = 1618$  (por dia) e  $Q_3 = 1783$  (por dia) com mediana igual a  $Md = 1705$  (por dia).

Em alguns casos pode ser conveniente trabalhar com o log da variável para obter uma maior simetria. Assim, plotamos os gráficos de dispersão para cada variável segundo o logaritmo da variável número de mortes e o resultado pode ser conferido na **Figura 33**. Assim, podemos notar que parece existir uma relação entre essas variáveis. Nota-se também que a relação entre a variável número de mortes e temperatura, bem como número de mortes e tempo possuem uma relação aparentemente não linear.



**Figura 33 – Diagrama de dispersão do logaritmo do número de mortes contra as demais variáveis numéricas considerando o exemplo do ar poluído em Chicago.**

Assim como realizado em exemplos anteriores, neste também será proposto um modelo de regressão. O objetivo é o de acomodar conjuntamente a informação contida nas variáveis explicativas e como elas impactam a média da variável resposta. Posteriormente a escolha do modelo de regressão adequando, será realizado uma análise de diagnóstico, visando identificar possíveis afastamentos de suposição, bem como prováveis erros sistemáticos.

#### 4.3.2 Ajuste do modelo

Neste exemplo, desejamos modelar uma variável de contagem, ou seja, uma variável discreta com suporte no conjunto dos inteiros não negativos. Para problemas deste tipo,

comumente a primeira alternativa de modelagem via modelo parcialmente lineares aditivos generalizados faz uso da distribuição Poisson com função de ligação logarítmica. Portanto, para ajuste desse modelo, considera-se que o número de mortes segue uma distribuição de Poisson com valor médio da taxa de mortalidade produto dessas variáveis que medem a qualidade do ar. Dessa forma, o modelo será dado por:

$$g(\mu_i) = \log(\mu_i) = \eta_i \\ = \beta_0 + \beta_1 MP_i + \beta_2 SO_{2i} + f_1(O_{3i}) + f_2(\text{tmpd}_i) + f_3(\text{Tempo}_i)$$

em que  $\mu_i = \mathbb{E}(\text{Morte}_i)$ , para  $i = 1, 2, 3, \dots, 300$ , onde  $f_k(\cdot)$ , com  $k = 1, 2$ , são funções suaves das covariáveis ozônio, temperatura média diária e tempo. As quais, acredita-se que possuem relação não linear e, portanto, devem ser estimadas de forma não paramétrica. Além disso, foi utilizada a função de ligação logarítmica.

Assim, para o ajuste deste modelo, destaca-se que foram empregados o método de *P-spline* para o ajuste da curva  $f_k(\cdot)$ , baseado em um *B-spline* cúbico com 10 nós equidistantes em todas as curvas e penalização quadrática. Os graus de liberdade efetivos são obtidos pela soma dos graus de liberdade referente a função não paramétrica mais os graus de liberdade referente a parte não paramétrica do modelo. Além disso, foi utilizado o método VCG para estimar o vetor de parâmetros de suavização  $\alpha$ . A **Tabela 10** apresenta as estimativas de máxima verossimilhança penalizada para ajuste do modelo.

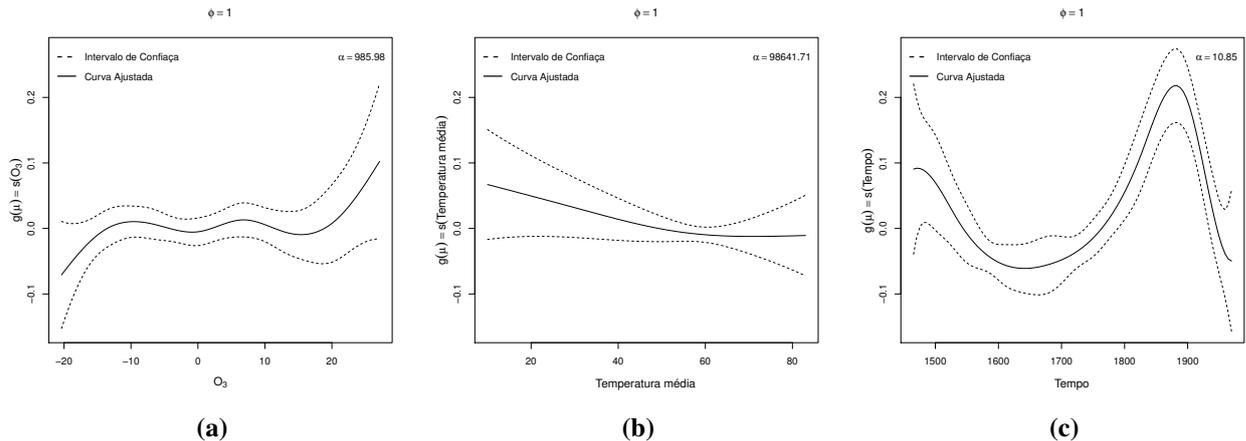
**Tabela 10 – Resumo das estimativas de máxima verossimilhança referentes ao modelo ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	4.68	0.006	< 0.05
MP	0.001	0.0004	< 0.05
SO <sub>2</sub>	-0.0007	0.003	0.8
Outras medidas do ajuste			
$df_1(\alpha)$	5.03		
$df_2(\alpha)$	1.73		
$df_3(\alpha)$	8.10		
$\alpha$	(985.98, 98641.71, 10.85)		
$\phi$	1		
VCG	0.017		

Note que a variável SO<sub>2</sub> não foi significativa para o modelo. Além disso, o custo total para a estimação dos parâmetros dos componentes paramétricos e não paramétricos foi

$df(\boldsymbol{\alpha}) = 3 + df_1(\boldsymbol{\alpha}) + df_2(\boldsymbol{\alpha}) + df_3(\boldsymbol{\alpha}) = 17.85$ , bem como critério VCG igual a 0.0222. Por meio da **Figura 34**, observamos as estimativas das curvas ajustada de forma não paramétrica.

Por fim, realizou-se mais um ajuste. No entanto, considerando apenas as variáveis material particulado, níveis de ozônio e tempo. Assim, a **Tabela final 11** resultante apresenta um resumo das estimativas e outras quantidades do referido ajuste.



**Figura 34 – Gráficos das curvas ajustadas sob distribuição Poisson para o exemplo do ar poluído em Chicago.**

Como a variável temperatura média possui grau de liberdade pequeno, bem como a curva ajustada é aproximadamente linear, então resolveu-se fazer um novo ajuste retirando a variável  $SO_2$ , que não foi significativa, assim como tomou-se a variável temperatura média sem suavizá-la. Entretanto, a avariável temperatura média também não foi significativa ao nível de significância de 5% para este segundo modelo.

**Tabela 11 – Resumo das estimativas de máxima verossimilhança referentes ao modelo final ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago.**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	4.68	0.005	< 0.0001
MP	0.001	0.0003	< 0.0001
<b>Outras medidas do ajuste</b>			
$df_1(\boldsymbol{\alpha})$	5.05		
$df_2(\boldsymbol{\alpha})$	8.28		
$\boldsymbol{\alpha}$	(965.7, 7.85)		
$\phi$	1		
VCG	0.018		

Desta forma, o modelo final será explicado de forma paramétrica pela variável material particulado, sendo possível realizar inferências sobre a relação ou saber qual o efeito da mudança de valor dessa variável em relação ao número de mortos. Por outro lado, também será explicado pelas componentes não paramétricas, perdendo um pouco do aspecto inferencial, mas ganhando muito no que tange a predição de novas observações.

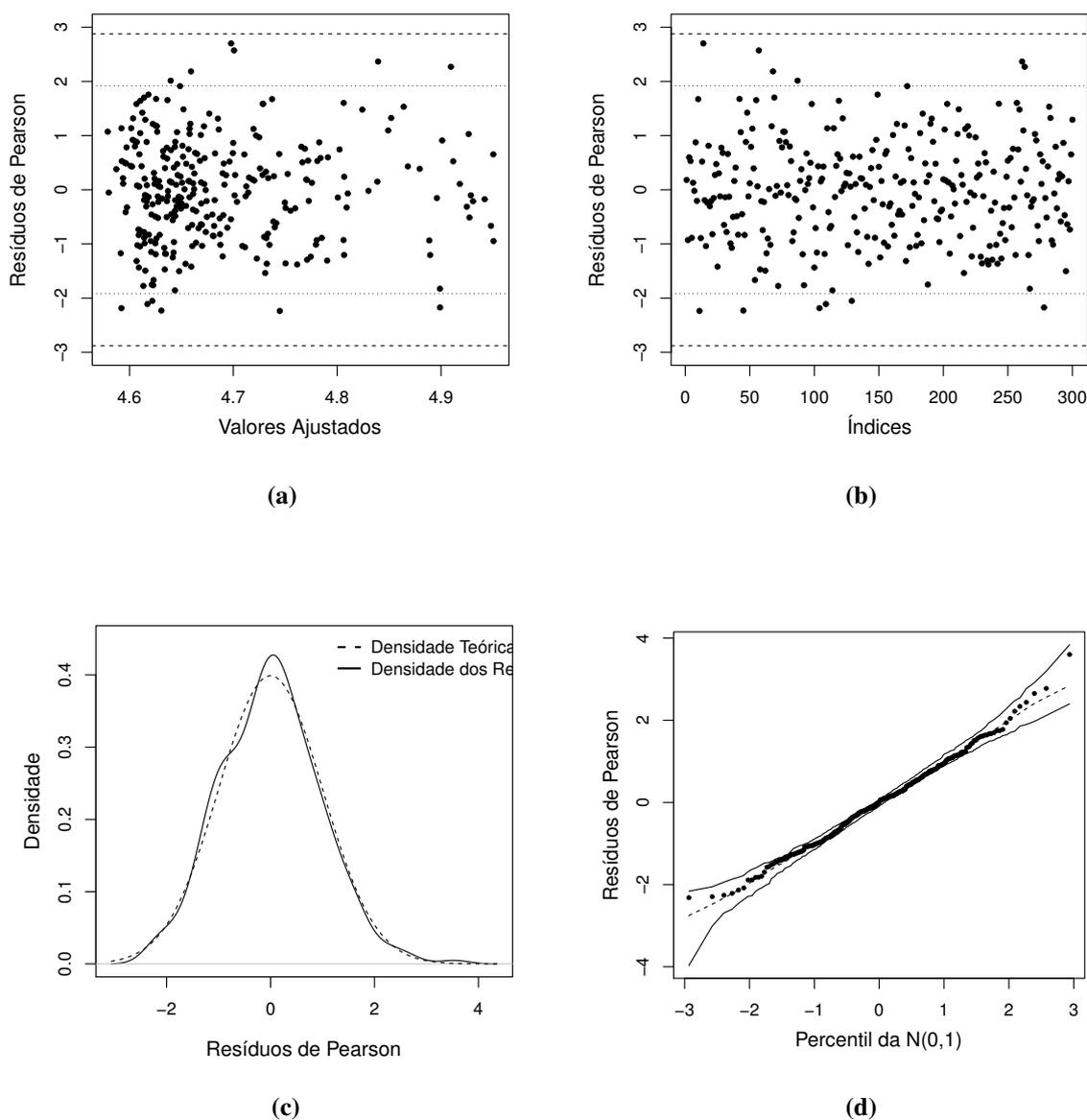
### 4.3.3 Diagnóstico

Pela análise de resíduo tem-se que, se o modelo for adequado, os resíduos quantílicos aleatorizados possuem uma distribuição próxima a de uma normal com média zero e desvio padrão igual a um. Dessa forma, espera-se que os resíduos quantílicos contra os valores ajustados demonstrados no gráfico apresentem comportamento aleatório em torno de zero e que os resíduos versus índices demonstrados no outro gráfico não apresentem comportamento serial aparente. Além disso, espera-se também que tanto a densidade quanto o envelope simulado dos resíduos apresentem um comportamento característico para essa distribuição. Assim, é possível observar que, por meio da **Figura (35)**, os gráficos dos resíduos contra os valores ajustados e o dos resíduos contra os índices, assim como os gráficos da densidade e envelope simulado apresentam comportamentos adequados.

Ademais, considerando pontos extremos aqueles que apresentam resíduos com 3 desvios padrão da média, então estes gráficos não apresentam pontos extremos. No entanto, apresentam alguns pontos de alarmes, que são merecedores de análise mais cuidadosa. Ou seja, pontos que estão acima de dois desvios padrão da média, porém abaixo de três desvios padrão. Portanto, pode-se tomar o modelo, tal como foi proposto, como um forte candidato para representar adequadamente uma aproximação da verdadeira natureza geradora dos dados.

Seguindo a análise de diagnóstico, considerando as medidas de alavancagem para o modelo proposto, observa-se que os pontos em destaque são os que possuem as alavancagens mais altas ilustradas na **Figura (36a)**. Portanto, sugere-se uma análise mais cautelosa destes pontos. Considerando a análise de influência sob os esquemas de perturbação de caso e em uma das variáveis explicativas, foram obtidos, respectivamente para cada tipo de observação, os gráficos da curvatura conformal contra a ordem das observações apresentado nas **Figuras. 36b e 36c**.

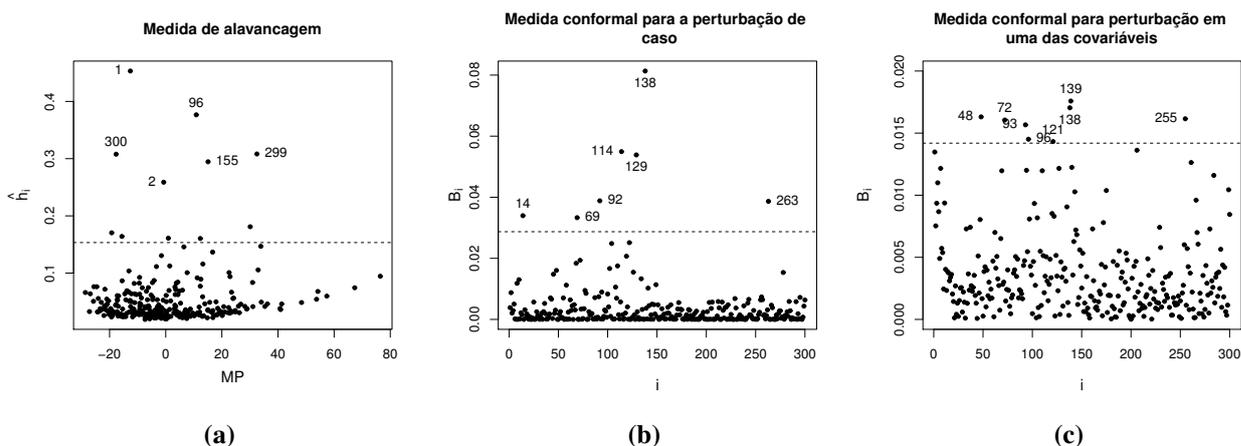
Neste caso, deve-se retirar esses pontos, realizar um novo ajuste e comparar as estimativas do novo modelo com as do modelo anterior. Assim, pela retirada desses pontos e



**Figura 35 – Gráficos de resíduos referentes ao modelo ajustado sob distribuição Poisson para o exemplo do ar poluído em Chicago.**

novo ajuste, obteve-se os resultados presentes na **Tabela 12**, onde podemos observar a mudança relativa entre o modelo completo e o modelo sem as observações potencialmente influentes. Após a confirmação de que as observações destacadas são influentes, devemos saber o porquê destas impactarem de forma desproporcional nas estimativas dos parâmetros.

O valor de referência para comparar as mudanças relativas das estimativas para este exemplo é de 0.66%. Assim, espera-se que na retirada de uma observação, as mudanças relativas sejam próximas a essa magnitude se a observação não for influente. Desta forma, nota-se que as estimativas do intercepto pouco mudam com a retirada individual destas observações, bem



**Figura 36 – Gráfico de influência local sob os esquemas de perturbação de caso e perturbação na variável resposta respectivamente para o exemplo do ar poluído em Chicago.**

como não mudou a inferência. Isto é, os parâmetros permaneceram significativos. Porém, as estimativas do parâmetro associado variável material particulado,  $\beta_1$ , aumenta quase 10% com a retirada da observação #138. Ao se investigar esta observação, nota-se que corresponde a maior quantidade de material particulado associado a um total de 143 mortes, um dos sete dias que mais houveram mortes.

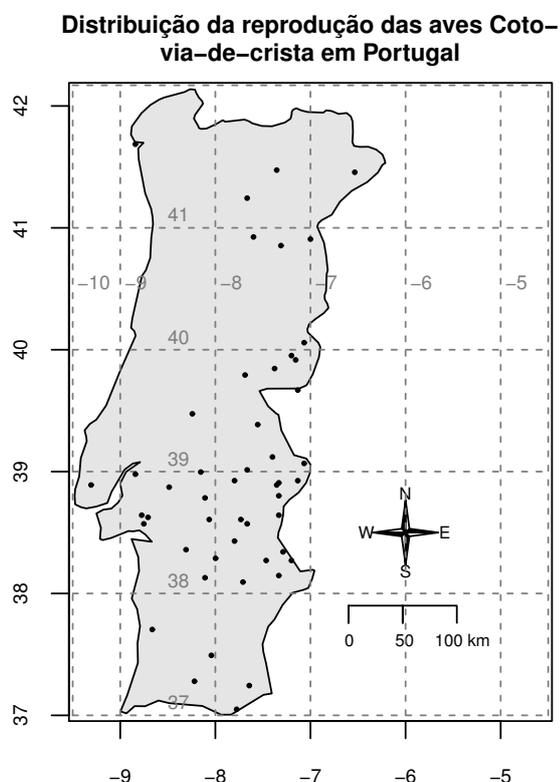
**Tabela 12 – Estimativas de máxima verossimilhança referente ao ajuste do modelo sob distribuição poisson e mudança relativa em porcentagem entre parênteses para os parâmetros  $\beta_0$  e  $\beta_1$  após a retirada dos pontos para o exemplo do ar poluído em Chicago.**

Observação retirada	Estimativas dos parâmetros	
	$\hat{\beta}_0$	$\hat{\beta}_1$
#1	4.68 (0.01%)	0.0013 (0.15%)
#96	4.69 (0.0%)	0.001 (0.06%)
#114	4.69 (0.02%)	0.001 (1.54%)
#129	4.69 (0.02%)	0.001 (0.81%)
#138	4.69 (0.02%)	0.001 (10.37%)
#139	4.69 (0.01%)	0.001 (0.20%)
#255	4.69 (0.01%)	0.001 (0.00%)

#### 4.4 REPRODUÇÃO DE AVES EM PORTUGAL

Nesta aplicação, considera-se parte de um conjunto de dados apresentado em Wood (2017). Este conjunto de dados descreve como é distribuída a reprodução de aves da espécie Cotovia-de-poupa em Portugal. O objetivo será, neste exemplo, estudar como a probabilidade de reprodução dessa espécie de ave varia no país, considerando, para isso, a localidade: "x"(em km considerando a direção leste de uma origem) e "y"(em km considerando a direção norte de uma origem). O objetivo é identificar onde as aves Cotovia-de-crista costumam se reproduzir em Portugal, considerando subdivisões do país em quadras de 2km por 2km. Para este exemplo, foi considerada uma mostra de 250 localidades. O resultado pode ser observado no gráfico da **Figura 37**.

Neste gráfico, as coordenadas  $x$  e  $y$  são usadas para indicar a presença ou ausência de aves da espécies Cotovia-de-crista se reproduzindo em Portugal. Portanto, os pontos destacados em preto indicam a presença de aves se reproduzindo naquela região (ou para aquelas coordenadas). Nota-se que para o lado leste e em regiões próximas a fronteira, bem como, principalmente, na região sul do país há uma maior concentração dessas aves se reproduzindo.



**Figura 37 – Gráfico que mostra como está distribuída a reprodução das aves Cotovia-de-crista em Portugal.**

A seguir, será proposto um modelo de regressão para tentar estimar a relação da reprodução dessas aves segundo a região definida no mapa. Em seguida, serão apresentadas as análises de diagnóstico para este modelo por meio da análise de resíduos, alavancagem e análise de influência local sob os esquemas de perturbação definidos e exemplificados em seções anteriores.

#### 4.4.1 Ajuste do modelo

Para o ajuste desse modelo, considera-se que a presença de aves nessas subdivisões segue uma distribuição binomial com média subjacente produto dessas variáveis. Ou seja, considera-se que

$$g(\mu_i) = \eta_i \quad (4.3)$$

$$= \beta_0 + \beta_1 y_i + f(x_i), \quad i = 1, 2, \dots, n. \quad (4.4)$$

Destaca-se que foram empregados o método de *P-spline* para o ajuste da curva  $f(\cdot)$  baseado em um *B-spline* cúbico com 10 nós equidistantes e penalização quadrática. Os graus de liberdade efetivos são obtidos pela soma dos graus de liberdade referente a função não paramétrica mais os graus de liberdade referente a parte paramétrica do modelo. Além disso, foi utilizado o método VCG para estimar o vetor de parâmetros de suavização  $\alpha$ , bem como a função de ligação utilizada logit. A **Tabela** (13) apresenta as estimativas de máxima verossimilhanças penalizada para ajuste do modelo.

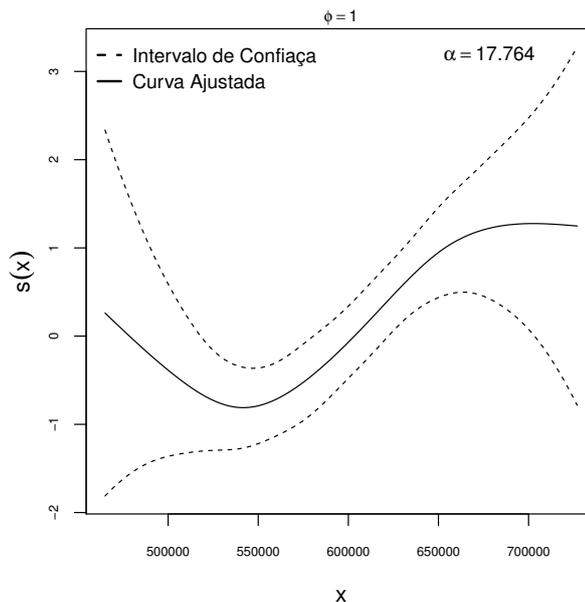
**Tabela 13 – Resumo das estimativas de máxima verossimilhança referentes ao modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves**

Efeito	Estimativa	Erro-Padrão	Valor-p
Intercepto	22.69	5.954	< 0.0001
MP	$-5.546 \times 10^{-6}$	$1.38 \times 10^{-6}$	< 0.0001
Outras medidas do ajuste			
$df(\alpha)$	5.01		
$\alpha$	17.764		
$\phi$	1		
VCG	-0.094		

Note que tanto o intercepto quanto a variável  $y$  foram significativas para o modelo. Além disso, os graus de liberdade efetivos foram estimados em  $df(\alpha) = 2 + df_1(\alpha) = 5.01$  e

que o grau de liberdade referente a variável suavizada  $x$  foi estimado em 3.01. Por fim, destaca-se que o critério VCG obtido foi de  $-0.094$ .

Por meio da **Figura 38**, podemos observar a representação gráfica do supracitado ajuste. Neste gráfico, nota-se que, de fato, a relação entre as variáveis se dá de forma mais complexa que uma simples estrutura linear. Assim, o modelo será bastante preciso em análises de predição, embora perca um pouco da interpretabilidade.

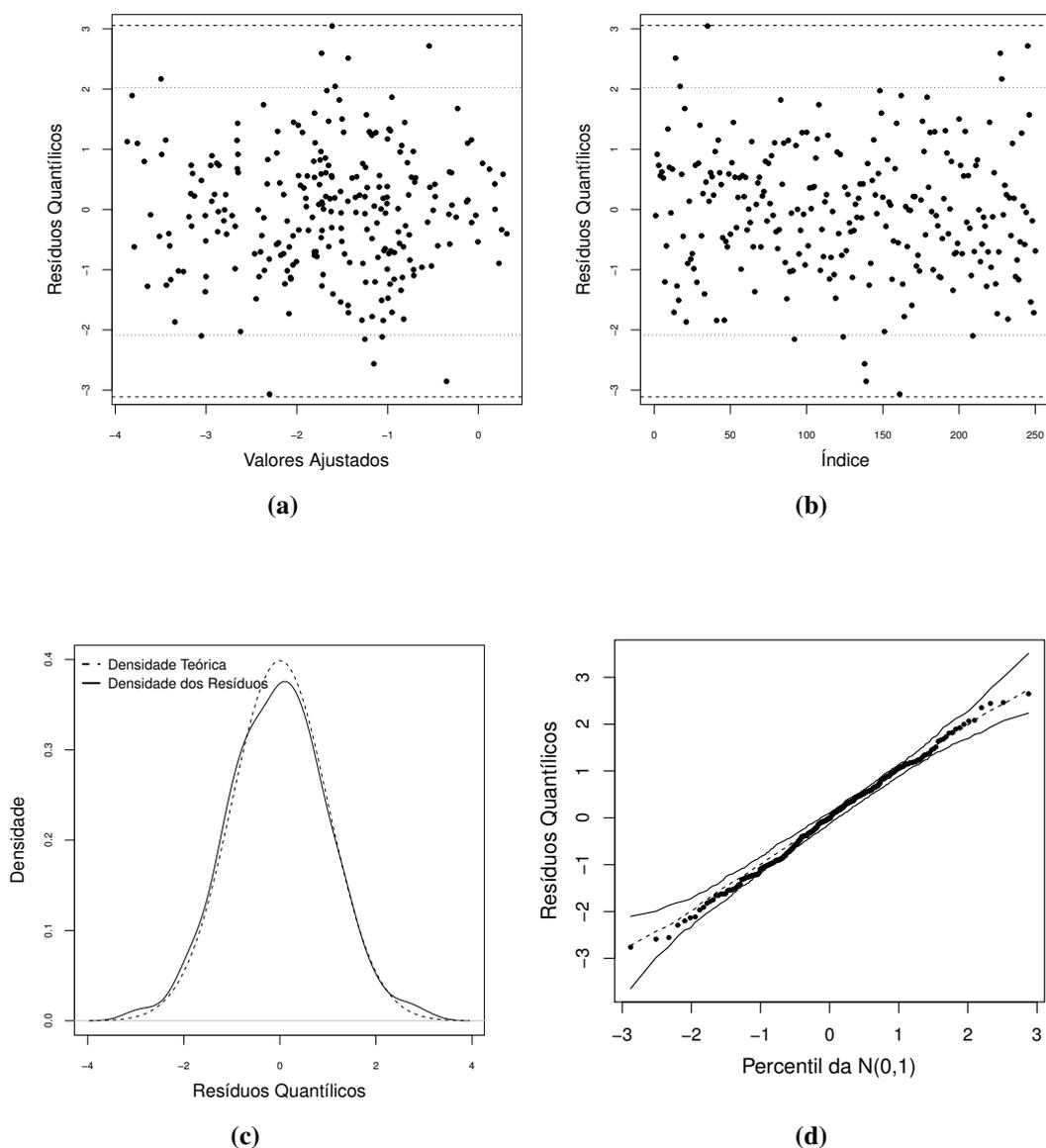


**Figura 38 – Curva aproximada e intervalos de confiança pontuais de 95% referente ao modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves.**

Nosso estudo agora se volta para verificar a validade das suposições assumidas para o modelo. Neste caso, será empregada técnicas de diagnósticos, por meio das quais será possível identificar desvios de suposições. Para tanto, faremos uso da análise de resíduos, alavancagem e análise de influência local. Superada essa fase, nosso modelo estará apto para uso.

#### 4.4.2 Diagnóstico

Para análise de resíduo, considera-se a **Figura 39** apresentada a seguir. Portanto, observa-se que, como visto anteriormente, o modelo será adequado se os resíduos quantílicos aleatorizados possuírem uma distribuição próxima a de uma normal com média zero e desvio padrão igual a um. Dessa maneira, pode-se observar que a densidade dos resíduos apresentada está bem próxima a densidade teórica esperada, além de nenhum dos pontos estarem fora do envelope.

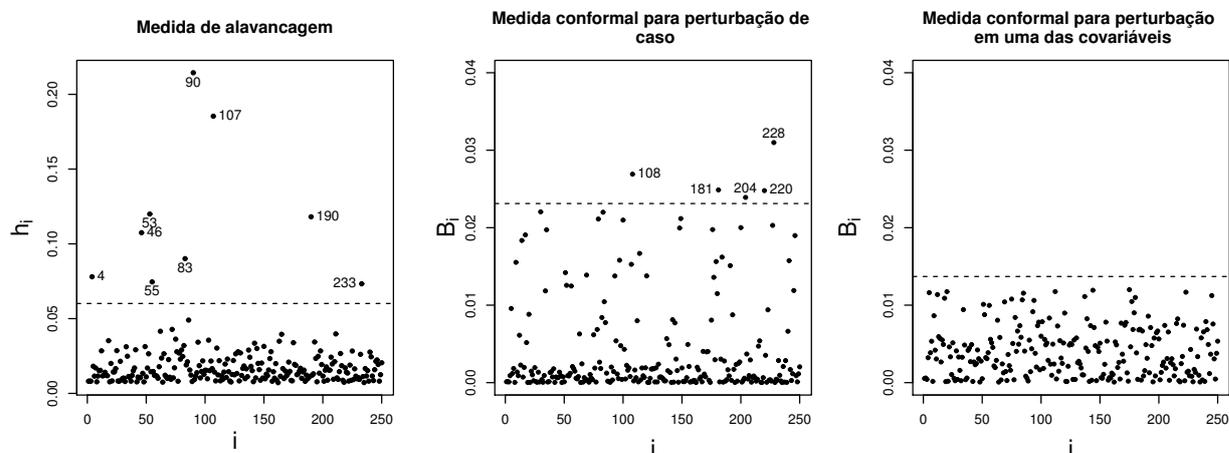


**Figura 39 – Gráfico de resíduos referente o modelo ajustado sob distribuição binomial para o exemplo da reprodução das aves.**

Portanto, segundo a análise de resíduo, o modelo proposto está apto ao uso. Porém, seguiremos a análise por meio das medidas de alavancagem e influência local para finalizar o diagnóstico. Dessa forma, pode-se tomar o modelo tal qual foi proposto como um forte candidato para representar adequadamente uma aproximação da verdadeira natureza geradora dos dados.

Seguindo a análise de diagnóstico, considerando as medidas de alavancagem para o modelo proposto, observa-se que os pontos em destaque são os que possuem as alavancagens mais altas ilustrada na **Figura 40**. Portanto, sugere-se uma análise mais cautelosa destes. Considerando a análise de influência sob os esquemas de perturbação de caso e na variável explicativa, observa-se que existem observações como possíveis observações influentes para

o gráfico correspondente a perturbação de caso. Por outro lado, não houveram observações influentes para o caso em que perturbamos a variável explicativa.



**Figura 40 – Gráficos de de alavancagem e influência local sob os esquemas de perturbação de caso e perturbação em uma das variáveis explicativas, respectivamente, referente ao ajuste do modelo sob distribuição binomial para o exemplo da reprodução das aves.**

Assim, novamente, analisamos o impacto das observações destacadas na **Figura 39** sob as estimativas dos parâmetros, por meio de uma análise de sensibilidade. Desta forma, consideramos a **Tabela 14** em que temos as observações retiradas seguidas dos valores das estimativas e as respectivas mudanças relativas em porcentagem para a análise de influência local. Para este exemplo, o valor de referência para comparar a mudança relativa das estimativas é de 0.8%. Portanto, mudanças relativas distantes desse valor confirma que a observação destacada é, de fato, influente nas estimativas dos parâmetros do modelo. Sendo assim, temos que as observações destacadas causam mudanças relativas substanciais nas estimativas dos parâmetros. Porém, não houveram mudanças inferenciais.

**Tabela 14 – Estimativas de máxima verossimilhança e mudança relativa em porcentagem entre parênteses para os parâmetros  $\beta_0$  e  $\beta_1$  do modelo binomial após a retirada dos pontos para o exemplo da reprodução das aves em Portugal.**

Estimativas dos parâmetros	Observação retirada				
	#108	#181	#204	#220	#228
$\hat{\beta}_0$	24.667 (8.705%)	24.270 (6.957%)	23.439 (3.294%)	24.228 (6.770%)	25.527 (12.495%)
$\hat{\beta}_1$	$-60.09 \times 10^{-7}$ (8.339%)	$-59.16 \times 10^{-7}$ (6.663%)	$-57.27 \times 10^{-7}$ (3.257%)	$-59.07 \times 10^{-7}$ (6.499%)	$-62.10 \times 10^{-7}$ (11.965%)

## 5 CONCLUSÕES E TRABALHOS FUTUROS

Na presente dissertação, discutimos o desenvolvimento da teoria que envolve os modelos parcialmente lineares aditivos generalizados com suavização por meio de *P-slines*. Neste contexto, desenvolvemos medidas para análise de diagnóstico por meio das principais estruturas conhecidas na literatura para este fim. Assim, sugerimos os resíduos de Pearson e os resíduos quantílicos aleatorizados com o intuito de identificar possíveis *outliers*, bem como discutimos o uso da matriz de projeção em modelos semi-paramétricos para investigar pontos de alavanca. Além disso, derivamos as expressões para realizarmos a análise de influência local, nossa principal contribuição teórica, baseando-se na proposta de Cook. Com isso, destaca-se que foram elaborados três esquemas de perturbação, a saber: perturbação no modelo (caso ponderado), perturbação na variável resposta e perturbação em uma das covariáveis, com o objetivo de estudar as estimativas do modelo sob tais esquemas de perturbação. Neste sentido, derivamos os resultados teóricos, além de elaborarmos uma implementação por intermédio do *software R* a fim de viabilizar o uso da técnica. Por fim, apresentamos quatro exemplos para que fosse possível vislumbrar o uso do método em diferentes distribuições.

Por meio dos exemplos apresentados, foi possível pôr em prática as técnicas desenvolvidas neste trabalho, bem como observar o desempenho destas medidas. Assim sendo, vimos que as técnicas de diagnóstico são extremamente úteis para ajudar a entender melhor as suposições inerentes ao modelo e são imprescindíveis para quaisquer análises estatísticas mais avançadas. Desta forma, as medidas desenvolvidas nesta dissertação são contribuições muito importantes para a comunidade estatística.

Em termos computacionais, a implementação dessas medidas a princípio parecem complicadas. Porém, o *software R* oferece suporte computacional bem sofisticado e de compreensão razoável, o que possibilitou desenvolver uma função para realizar os diagnósticos, baseada nas saídas do pacote **mgcv**. Estas funções e outros códigos usados nesta dissertação encontram-se no apêndice A, que podem motivar a criação de um pacote. Por fim, seria interessante investigar o desempenho dessas técnicas em diferentes cenários, o que só seria possível por meio de um estudo de simulação.

## REFERÊNCIAS

- AKAIKE, H. A new look at the statistical model identification. In: **Selected Papers of Hirotugu Akaike**. [S.l.]: Springer, 1974. p. 215–222.
- BATES, D. M.; WATTS, D. G. Nonlinear regression: iterative estimation and linear approximations. **Nonlinear Regression Analysis and Its Applications, John Wiley & Sons, Inc., Hoboken, NJ, USA, doi**, v. 10, n. 9780470316757, p. 9393–9441, 1988.
- BEALE, E. Confidence regions in non-linear estimation. **Journal of the Royal Statistical Society: Series B (Methodological)**, Wiley Online Library, v. 22, n. 1, p. 41–76, 1960.
- BILLOR, N.; LOYNES, R. Local influence: a new approach. **Communications in Statistics-Theory and Methods**, Taylor & Francis, v. 22, n. 6, p. 1595–1611, 1993.
- BOOR, C. D. A practical guide to splines. **Springer**, Berlin, 1978.
- BREIMANE, L.; FRIEDMAN, J. H. Estimating optimal transformations for multiple regression and correlation. **Journal of the American Statistical Association**, v. 80, p. 580–298, 1985.
- BUJA, A.; HASTIE, T.; TIBSHIRANI, R. Linear smoother and additive models. **The Annals of Statistics**, v. 17, p. 453–555, 1989.
- COLOSIMO, A. E.; GIOLO, S. R. **Análise de Sobrevivência Aplicada**. [S.l.]: Associação Brasileira de Estatística, 2006. v. 1.
- COOK, R. D. Assessment of local influence. **Journal of the Royal Statistical Society - Series B (Methodological)**, 1978.
- CORDEIRO, G. M.; DEMÉTRIO, C. G. B. **Modelos Lineares Generalizados e Extensões**. Recife - PE: Departamento de Estatística e Informática - UFRPE, 2008.
- COX, D. R.; SNELL, E. J. A general definition of residuals (with discussion). **Journal of the Royal Statistical Society B**, v. 30, p. 248–275, 1968.
- DUNN, P. K.; SMYTH, G. K. Randomized quantile residuals. **Journal of Computational and Graphical Statistics**, v. 5, p. 236–244, 1996.
- EILERS, P. H. C.; MARX, B. D. Flexible smoothing with b-splines and penalties. **Statistical Science**, 1996.
- EILERS, P. H. C.; MARX, B. D.; MARIA, D. Twenty years of p-splines. *Statistics and Operations Research Transactions*, v. 39(2), p. 149–186, 2015.
- EMAMI, H. Local influence in ridge semiparametric models. **Journal of Statistical Computation and Simulation**, Taylor & Francis, v. 86, n. 17, p. 3357–3370, 2016.
- EUBANK, R. L. The hat matrix for smoothing splines. **Statistics and Probability Letters**, v. 2, p. 9–14, 1984.
- EUBANK, R. L. **Spline Smoothing and Nonparametric Regression**. [S.l.]: Marcel Dekker, 1988.
- FAHRMEIR, L.; KNEIB, T.; LANG, S.; MARX, B. **Regression: Models, Methods and Applications**. [S.l.]: Springer, 2013.

FRIEDMAN, J.; SILVERMAN, B. W. Flexible parsimonious smoothing and additive modeling (with discussion). **Technometrics**, v. 31, p. 2–39, 1989.

GREEN, P. J.; SILVERMAN, B. W. Nonparametric regression and generalized linear models. Chapman and Hall, 1994.

GREEN, P. J.; YANDELL, B. S. Semi-parametric generalized linear models. **Lecture Notes in Statistics**, v. 32, p. 44–55, 1985.

HÄRDLE, W. Applied nonparametric regression. Cambridge Univ. Press, 1990.

HASTIE, T.; TIBSHIRANI, R. Generalized additive models. **Chapman and Hall**, 1990.

HOAGLIN, D. C.; WELSCH, R. E. The hat matrix in regression and anova. **The American Statistician**, v. 32, p. 17–22, 1978.

HODGES, J. S. **Richly parameterized linear models: additive, time series, and spatial models using random effects**. [S.l.]: Chapman and Hall/CRC, 2016.

HOLANDA, A. A. Modelos lineares parciais aditivos generalizados com suavização por meio de p-splines. Departamento de Estatística - USP, São Paulo - SP, Maio 2018.

IHAKA, R.; GENTLEMAN, R. R. A language for data analysis and graphics. **Journal of Computational and Graphical Statistics**, v. 5, p. 299–314, 1996.

IZBICKI, R.; SANTOS, T. M. **Machine Learning sob a ótica estatística: Uma abordagem preditiva para a Estatística com Exemplos em R**. [S.l.]: Departamento de Estatística - UFSCAR, 2019.

KOOPERBERG, C.; STONE, C. J. Log-spline density estimation for censored data. **J. Comput. Graph. Statist**, v. 1, p. 301–328, 1992.

LEE, S.-Y.; LU, B.; SONG, X.-Y. Assessing local influence for nonlinear structural equation models with ignorable missing data. **Computational statistics & data analysis**, Elsevier, v. 50, n. 5, p. 1356–1377, 2006.

MANGHI, R. F. Técnicas de diagnóstico em modelos parcialmente lineares aditivos generalizados para dados correlacionados. Universidade Federal de Pernambuco - Departamento de Estatística, 2016.

NELDER, J. A.; WEDDERBURN, R. W. M. Generalized linear models. **Journal of the Royal Statistical Society A**, v. 135, p. 370–384, 1972.

O’SULLIVAN, F. A statistical perspective on ill-posed inverse problems (with discussion). **Statist. Sci.**, v. 1, p. 505–527, 1986.

PAULA, G. A. **Modelos de Regressão com apoio computacional**. Universidade de São Paulo: Instituto de Matemática e Estatística, 2013. 2 p.

PAYANDEH, B. Some applications of nonlinear regression models in forestry research. **The Forestry Chronicle**, NRC Research Press Ottawa, Canada, v. 59, n. 5, p. 244–248, 1983.

PULGAR, G. M. I. **Modelos mistos aditivos semiparamétricos de contornos elípticos**. Tese (Doutorado) — Universidade de São Paulo, 2009.

REINSCH, C. Smoothing by spline functions. **Numer. Math**, v. 10, p. 177–183, 1967.

RUPPERT, D. Selecting the number of knots for penalized splines. **Journal of computational and graphical statistics**, Taylor & Francis, v. 11, n. 4, p. 735–757, 2002.

SCHMIDT, G.; MATTERN, R.; SCHÜLER, F. Biomechanical investigation to determine physical and traumatological differentiation criteria for the maximum load capacity of head and vertebral column with and without protective helmet under effects of impact. **Final report phase III, Project**, v. 65, 1981.

SCHWARZ, G. *et al.* Estimating the dimension of a model. **The annals of statistics**, Institute of Mathematical Statistics, v. 6, n. 2, p. 461–464, 1978.

SOUZA, C. P. E. Testes de hipóteses para dados funcionais baseados em distâncias: um estudo usando splines. 2008.

TEAM, R. D. C. **An Introduction to R**. [s.n.], 2019. Disponível em: <<https://cran.r-project.org/doc/manuals/r-release/R-intro.pdf>>.

VANEGAS, L. H.; PAULA, G. A. An extension of log-symmetric regression models: R codes and application. **Journal of Statistical Simulation and Computation**, v. 86, p. 1709–1735, 2016.

WAHBA, G. Spline models for observational data. SIAM, 1990.

WEI, B.-C.; HU, Y.-Q.; FUNG, W.-K. Generalized leverage and its applications. **Scandinavian Journal of statistics**, Wiley Online Library, v. 25, n. 1, p. 25–37, 1998.

WOOD, S. N. **Generalized Additive Models: an Introduction with R**. [S.l.]: CRC Press, 2017. v. 2.

## APÊNDICE A – CÓDIGOS COMPUTACIONAIS

Neste apêndice, são apresentados alguns códigos computacionais que foram implementados, utilizando a sintaxe do R, para obtenção dos resultados expostos nessa dissertação.

### CÓDIGO EM R PARA CONSTRUIR UMA SPLINE LINEAR E CÚBICA

```

1
2  pkg <- c("segmented",
3          "latticeExtra",
4          "RSADBE")
5  sapply(pkg, require,
6         character.only = TRUE)
7
8
9  data("PW_Illus")
10 data <- PW_Illus
11 names(data) <- c("x", "y")
12 data$x2 <- with(data,
13                ifelse(x <= 15, 0, x - 15))
14 data$x3 <- with(data,
15                ifelse(x <= 30, 0, x - 30))
16 m0 <- lm(y ~ x+x2+x3, data)
17 coef.m0 <- coefficients(m0)
18
19 plot(data$x, data$y, pch = 20, col = "cyan", ylab = "y", xlab = "x")
20 curve(coef.m0[1]+{coef.m0[2]*x}, add=T,
21       col = 6, lwd = 2, from = -1, to = 15)
22
23 curve(coef.m0[1]+{coef.m0[2]*x+coef.m0[3]*(x-15)},
24       col = 6, lwd = 2, add = T, from = 15, to = 30)
25 curve(coef.m0[1]+{coef.m0[2]*x+coef.m0[3]*(x-15)+coef.m0[4]*(x-30)},
26       col = 6, lwd = 2, add = T, from = 30)

```

**Código-fonte 1 – Função para ajuste de uma splines linear**

```

2  set.seed(19)
3  x = sort(round(runif(30, 0.5,3.3),3))
4  y = 3*cos(3+2*x) + 5 + rnorm(30, 0, 1)
5
6
7  data <- data.frame(y,x)
8  data$x2 <- with(data,
9  ifelse(x <= 1.5, 0, x - 1.5))
10
11 data$x3 <- with(data,
12 ifelse(x <= 2.5, 0, x - 2.5))
13
14 m0 <- lm(y ~ x+I(x**2)+I(x**3)+I(x2**3)+I(x3**3),
15 data)
16 coef.m0 <- coefficients(m0)
17
18 plot(x,y, pch = 20, col = "cyan", ylim = c(0,11), xlim = c(0.6,3.3))
19 curve(coef.m0[1]+coef.m0[2]*x+coef.m0[3]*x**2+{coef.m0[4]*x**3},
20 col = 6,
21 lwd = 2, add = T, from = 0, to = 1.5)
22
23 curve(coef.m0[1]+coef.m0[2]*x+coef.m0[3]*x**2+{coef.m0[4]*x**3+coef.m0
    [5]*(x-2)**3},
24 col = 6,
25 lwd = 2, add = T, from = 1.5, to = 2.5)
26
27 curve(coef.m0[1]+coef.m0[2]*x+coef.m0[3]*x**2+{coef.m0[4]*x**3+
28 coef.m0[5]*(x-2)**3+coef.m0[6]*(x-2.5)**3},
29 col = 6,
30 lwd = 2, add = T, from = 2.5)
31 abline(v=1.5, lty = 2)
32 abline(v=2.5, lty = 2)
33
34 m1 <- lm(y ~ x+x2+x3,
35 data)
36 coef.m1 <- coefficients(m1)
37
38 curve(coef.m1[1]+{coef.m1[2]*x}, add=T,
39 col = "green",

```

```

40 lwd = 2, from = 0, to = 1.5)
41
42 curve(coef.m1[1]+{coef.m1[2]*x+coef.m1[3]*(x-1.5)},
43 col = "green",
44 lwd = 2, add = T, from = 1.5, to = 2.5)
45 curve(coef.m1[1]+{coef.m1[2]*x+coef.m1[3]*(x-1.5)+coef.m1[4]*(x-2.5)},
46 col = "green",
47 lwd = 2, add = T, from = 2.5)
48
49 legend(2.4,10,legend = c("spline linear","spline cubica"),
50 bty = "n",lwd = 2, lty=c(1,1), col = c("green",6))

```

## Código-fonte 2 – Função para ajuste de uma splines cúbica

### FUNÇÃO PARA AJUSTE DE UMA FUNÇÃO NÃO PARAMÉTRICA

#### Exemplo 1: Impacto do capacete da motocicleta

```

##### Função para ajuste do modelo #####
#
# ----- Função para estimar as quantidades de um modelo aditivo -----#
#      generalizado para uma única função não paramétrica
#
#####
# =====
#
#              DETALHES DA FUNÇÃO
#
# A função recebe 4 objetos de entrada:
#
# 1 - A variável explicativa

```

```

# 2 - A variável resposta #
# 3 - Um vetor contendo os nós #
# 4 - Um parâmetro desuavização #
# #
# ----- (NESSA ORDEM) ----- #
# Obs.: Note que esta função é facilmente estendida para ajustes #
# considerando mais de uma função não paramétrica #
# #
# ===== #

```

```

1 library(MASS)
2 library(mosaic)
3 data(mcycle)
4 attach(mcycle)
5
6 # Funcao para o ajuste do modelo
7 psfit <- function(x,y,pord=2,k,alpha=1){
8 X <- splineDesign(knots = k, x) # Matriz de base B-spline
9 nb <- ncol(X)
10 P <- diff(diag(nb),differences=pord)
11 MatC <- rep(1, nrow(X)) % * % X
12 qrc <- qr(t(MatC))
13 Z <- qr.Q(qrc,complete=TRUE)[, (nrow(MatC)+1):ncol(MatC)]
14 B <- cbind(1,X % * % Z) # Nova matriz modelo
15
16 # Matriz de penalidades
17 Matsqrt <- function(S)
18 {
19 d <- eigen(S,symmetric=TRUE)
20 rS <- d$vector % * % abs(diag(d$values))^0.5 % * % t(d$vector)
21 return(rS)
22 }
23
24 S <- t(Z) % * % t(P) % * % P % * % Z
25 rS <- Matsqrt(S) * sqrt(alpha/16)
26 O <- matrix(0,1,1)
27 rS <- bdiag(O,rS)
28 # Ajuste

```

```

29 | X1 <- as.matrix(rbind(B,rS))
30 | k <- ncol(B)
31 | n <- nrow(B)
32 | y1 <- y;y1[(n+1):(n+k)]<-0
33 | f <- glm(y1~X1-1,family = gaussian(link = "identity"))
34 | h <- hat(f$qr)[1:n]
35 | mu <- fitted.values(f)[1:n]
36 | # Validacao Cruzada e dispersao
37 |
38 | trH <- sum(h)
39 | rss <- sum((y-mu)^2/f$family$variance(mu)[1:n])
40 | gcv <- n*rss/(n-trH)^2
41 | sigma <- sqrt(rss / (n - trH))
42 |
43 | output <- list(gcv=gcv , sigma=sigma , f=f , trH=trH , S=S, B=B)
44 | return(output)
45 | }
46 | knots <- c(-12.739855,-7.711636,-2.683418,2.344800,7.373018,
47 |           12.401236,17.429455,22.457673,27.485891,32.514109,
48 |           37.542327,42.570545,47.598764,52.626982,57.655200,
49 |           62.683418,67.711636,72.739855)
50 | fit <- psfit(times , accel ,k=knots , alpha = 1)

```

### Código-fonte 3 – Função para ajuste do modelo

```

##### Obtendo Alpha #####
#
#
# ----- Esta parte do código é utilizada para encontrar o -----#
#           valor do parâmetro de suavização ótimo
#
#####

```

```

1 | lla <- seq(-3, 3, by = 0.10)

```

```
2 cvs <- 0 * lla
3 alpha<-1
4 for (k in 1:length(lla)) {
5   alpha = 10^lla[k]
6   pn = psfit(times, accel, k=knots, alpha = alpha)
7   cvs[k] = pn$gcv
8 }
9 # alpha escolhido por cv generalizado
10 alpha.cv <- 10^(lla[which.min(cvs)])
```

#### **Código-fonte 4 – Obtenção de alpha**

```

1 #Menor VCG
2 cvs[which.min(cvs)]
3 plot(1:length(l1a),cvs,type="l",main="Pontuacao VCG",
4      xlab="i",ylab="vcs")
5 legend("topleft",legend="Menor VCG = 560.077",bty="n")

```

### Código-fonte 5 – Plot - Pontuação VCG

```

##### Estimativas do modelo ótimo #####
#
# Esta parte do código é usada para obtenção do ajuste baseado no #
# valor ótimo de alpha, bem como obter outras medidas importantes, #
# tais como: Matriz de variância e covariância, intervalos de con- #
# fiança, dentre outras. #
# #
#####

```

```

1 # Ajuste baseado no alpha otimo
2 fit.cv <- psfit(times, accel, k=knots, lambda=lam.cv)
3 # Graus de liberdade efetivos
4 df <- fit.cv$trH
5 # medidas para obter os erros padrao e intervalos de confianca
6 n <- length(accel)
7 mu <- fitted.values(fit.cv$f)[1:n]
8 eta <- fit.cv$f$family$linkinv(mu)[1:n]
9 D <- diag(fit.cv$f[]$family$mu.eta(eta)[1:n], n)
10 O <- matrix(0,1,1)
11 S <- fit.cv$S
12 K <- bdiag(lam.cv*S); K <- bdiag(O,K)
13 W <- diag(fit.cv$f$weights[1:n], n)
14 A <- solve(t(fit.cv$B) % * % W % * % fit.cv$B + K) * fit.cv$sigma^2
15 Variace <- fit.cv$B % * % A % * % t(fit.cv$B)
16 pred.int <- fit.cv$B % * % coef(fit.cv$f)
17 upr <- pred.int + (2 * sqrt(diag(Variace)))
18 lwr <- pred.int - (2 * sqrt(diag(Variace)))

```

```

19 plot(times , accel , lwd=2, lty = 1,
20 col = "gray80", ann = F, ylim = c(-150,100))
21 lines( fitted( fit.cv$f)[ order(times)]~times[ order(times)], col= 1, lwd=1)
22 lines( upr~times , lty = 2, col = 1)
23 lines( lwr~times , lty = 2, col = 1)
24 mtext(side = 1, text = "Tempo (ms)", line=3, cex = 1.2)
25 mtext(side = 2, text = "Aceleracao (g)", line=3, cex = 1.2)
26 fit.cv$sigma^2
27 title( expression( hat(sigma)^2==514.74))
28 cor = c(1,1)
29 legend("topright", legend = c("Intervalo de Confiaca", "Curva Ajustada"),
30       bty = "n", lwd = 2, lty=c(2,1), col = cor)
31 legend("topleft", legend = expression(alpha==0.631), bty = "n")
32
33 erros_padrao<-sqrt(diag(A))
34
35 t_value<-coef( fit.cv$f)[1]/erros_padrao[1]

```

### Código-fonte 6 – Estimativas do modelo ótimo

```

1 library(ISLR)
2 library(splines)
3 library(mgcv)
4 library(splines)
5 library(mosaicCalc)
6 library(mosaic)
7 library(MASS)
8 library(matlib)
9
10
11 gaminfluence.diag <- function( fit ){
12 y <- fit$y
13 w <- fit$weights
14 mu <- fitted.values( fit )
15 eta <- fit$family$linkfun(mu)
16 lambda <- fit$sp
17 family <- fit$family
18 #Matriz modelo

```

```

19 X<-model.matrix( fit )
20 # quantidade de termos suavizados
21 d <- length( fit$sp )
22 # quantidades de parametros na parte parametrica do modelo
23 p <-length( summary( fit )$p.coeff )
24 n<-length( y )
25 #derivada da funcao de variacia
26
27 deriv.variencie <- D( family$variance( mu )~mu )
28
29 # Derivada de g(eta)
30 # detalhes - mu.eta e a derivada da inversa da funcao de ligacao:
31 #           mu = ginv(eta) e mu.eta = d(ginv(eta))/d(eta) = d(mu)/d(
           eta )
32
33 deriv.mu.eta <- D( family$mu.eta( eta )~eta )
34
35 S<-D<-list ( )
36 for ( i in 1:d ){
37 #----- Matriz S para a penalizacao de cada termo suavizado
           -----
38 D[i]<-fit$smooth[[1]]$S
39 S[[i]] <- lambda[i]*D[[i]]
40 }
41
42
43 #----- Matriz O de zeros
           -----
44 O <- matrix( 0 , p , p )
45
46
47 #----- Matriz K de penalizacao
           -----
48
49 K <- bdiag( S ) ; K<-bdiag( O , K )
50
51
52 #----- Matriz L
           -----

```

```

53 D <- diag(family$mu.eta(eta),n)
54 W <- diag(w,n)
55 V <- diag(family$variance(mu),n)
56 dV <- diag(deriv.variance(mu),n)
57 dD <- diag(deriv.mu.eta(eta),n)
58 C <- diag(y - mu,n)
59 M1 <- (D % * % ginv(dV) % * % D + dD % * % ginv(V)) % * % C
60 M2 <- -W
61 L <- t(X) % * % (M1+M2) % * % X - K;L<- as.matrix(L)
62 # ----- Perturbacao de caso
    -----
63 Deltac <- t(X) % * % W % * % ginv(D) % * % C
64 Fc<- -t(Deltac) % * % ginv(L) % * % Deltac # Matriz F
65 Cic<-2*abs(diag(Fc))
66 deFc<-eigen(Fc)
67 Cdmaxc<-deFc$val[1] #maior autovalor
68 dmaxc<-deFc$vec[,1] #maior autovetor associado
69 dmaxc<-dmaxc/sqrt(Cdmaxc)
70 dmaxc<-abs(dmaxc)
71 #curvatura normal conformal
72 Bc<- abs(diag(Fc)/sqrt(sum(deFc$values^2)))
73
74
75 #----- Perturbacao na variavel Resposta
    -----
76 Deltar <- t(X) % * % W % * % ginv(D) % * % sqrt(V)
77 Fr<- -t(Deltar) % * % ginv(L) % * % Deltar # Matriz F
78 Cir<-2*abs(diag(Fr))
79 deFr<-eigen(Fr)
80 Cdmaxr<-deFr$val[1] #maior autovalor
81 dmaxr<-deFr$vec[,1] #maior autovetor associado
82 dmaxr<-dmaxr/sqrt(Cdmaxr)
83 dmaxr<-abs(dmaxr)
84 #curvatura normal conformal
85 Br<- abs(diag(Fr)/sqrt(sum(deFr$values^2)))
86
87 list(Fc=Fc, Cic = Cic , Cdmaxc = Cdmaxc, dmaxc=dmaxc, Bc=Bc, Fr=Fr, Cir = Cir ,
      Cdmaxr = Cdmaxr, dmaxr=dmaxr, Br=Br)
88 }

```

---

**Código-fonte 7 – Função para análise de influência local**

## ANEXO A – DERIVAÇÃO DOS ESQUEMAS DE PERTURBAÇÃO

Os métodos de diagnósticos por meio da análise de influência local consistem nos seguintes esquemas de perturbação: casos ponderados, perturbação na variável resposta e perturbação de uma covariável. Portanto, considere que  $\xi = (\beta, \tau_1, \tau_2, \dots, \tau_r)^\top \in \Xi \subseteq \mathbb{R}_{v_s}^s$ , em que  $s_v = p + \sum_{k=1}^r v_k$  e  $v_k - q_k - 1$ . Isto posto, sob a suposição de independência, a função de log-verossimilhança penalizada será dada por:

$$\mathcal{F}_\alpha(\xi) = L(\xi) + \frac{1}{2} \sum_{k=1}^r \alpha_k \tau_k^\top \bar{\mathbf{S}}_k \tau_k, \quad (\text{A.1})$$

em que  $L(\xi) = \sum_{i=1}^n \{ \phi^{-1}[y_i \theta_i - b(\theta_i)] + c(y_i; \phi) \}$ ,  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_r)^\top$  é um vetor de parâmetros de suavização, de forma que  $\alpha_k > 0$ ,  $\tau_k$  são os coeficientes dos *B-splines* da  $k$ -ésima função não paramétrica centralizada,  $\bar{\mathbf{S}}_k = \mathbf{Z}_k^\top \mathbf{P}_k^d \mathbf{Z}_k$ ,  $\mathbf{P}_k^d = (\Delta_k^d)^\top (\Delta_k^d)$  e  $\Delta_k^d$  é um operador de diferenças, tal que  $\Delta_k^d \gamma_{kl} = \Delta_k (\Delta_k^{d-1} \gamma_{kl})$  e  $\Delta_k \gamma_{kl} = (\gamma_{kl} - \gamma_{k(l-1)})$ , sendo  $d$  a ordem da diferença.

Para fazer uso da teoria de influência local, devemos calcular a matriz de informação observada,  $\ddot{\mathcal{F}}_\alpha$ . Para tanto, considerando  $\alpha$  e  $\phi$  fixados, tem-se que  $\ddot{\mathcal{F}}_\alpha$  é definida como:

$$\ddot{\mathcal{F}}_\alpha = \frac{\partial^2 \mathcal{F}_\alpha(\xi | \omega_0)}{\partial \xi \partial \xi^\top} = \begin{pmatrix} \ddot{\mathcal{F}}_{\alpha\beta\beta} & \ddot{\mathcal{F}}_{\alpha\beta\tau_k} \\ \ddot{\mathcal{F}}_{\alpha\tau_k\beta} & \ddot{\mathcal{F}}_{\alpha\tau_k\tau_k} \end{pmatrix}, \quad (\text{A.2})$$

em que

$$\begin{aligned} \ddot{\mathcal{F}}_{\alpha\beta\beta} &= \frac{\partial^2 \mathcal{F}_\alpha(\xi | \omega_0)}{\partial \beta_j \partial \beta_l} = \phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d^2 \theta_i}{d\mu_i^2} \left( \frac{d\mu_i}{d\eta_i} \right)^2 x_{ij} x_{il} \\ &\quad + \phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d\theta_i}{d\mu_i} \frac{d^2 \mu_i}{d\eta_i^2} x_{ij} x_{il} \\ &\quad - \phi^{-1} \sum_{i=1}^n \frac{d\theta_i}{d\mu_i} \left( \frac{d\mu_i}{d\eta_i} \right)^2 x_{ij} x_{il}, \\ \ddot{\mathcal{F}}_{\alpha\tau_k\tau_k} &= \frac{\partial^2 \mathcal{F}_\alpha(\xi | \omega_0)}{\partial \tau_{jk} \partial \tau_{kl}} = \phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d^2 \theta_i}{d\mu_i^2} \left( \frac{d\mu_i}{d\eta_i} \right)^2 \tilde{b}_{ijk} \tilde{b}_{ilk} \\ &\quad + \phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d\theta_i}{d\mu_i} \frac{d^2 \mu_i}{d\eta_i^2} \tilde{b}_{ijk} \tilde{b}_{ilk} \\ &\quad - \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \left( \frac{d\mu_i}{d\eta_i} \right)^2 \tilde{b}_{ijk} \tilde{b}_{ilk} - \alpha_k \bar{s}_{ijk} \right\}, \quad \text{para } j, l = 1, 2, \dots, v_k. \\ \ddot{\mathcal{F}}_{\alpha\beta\tau_k} &= \frac{\partial^2 \mathcal{F}_\alpha(\xi | \omega_0)}{\partial \beta_j \partial \tau_{kl}} = -\phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d^2 \theta_i}{d\mu_i^2} \left( \frac{d\mu_i}{d\eta_i} \right)^2 x_{ij} \tilde{b}_{ilk} \\ &\quad - \phi^{-1} \sum_{i=1}^n (y_i - \mu_i) \frac{d\theta_i}{d\mu_i} \frac{d^2 \mu_i}{d\eta_i^2} x_{ij} \tilde{b}_{ilk} \end{aligned}$$

$$+ \phi^{-1} \sum_{i=1}^n \frac{d\theta_i}{d\mu_i} \left( \frac{d\mu_i}{d\eta_i} \right)^2 x_{ij} \tilde{b}_{ilk},$$

$\ddot{\mathcal{F}}_{\alpha\beta\tau_k} = \ddot{\mathcal{F}}_{\alpha\tau_k\beta}^\top$  no qual  $\mu_i = g^{-1}(\eta_i)$ ,  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \tilde{\mathbf{b}}_{1i}^\top \boldsymbol{\tau}_1 + \tilde{\mathbf{b}}_{2i}^\top \boldsymbol{\tau}_2 + \dots, \tilde{\mathbf{b}}_{ri}^\top \boldsymbol{\tau}_r$ . Escrevendo em forma de matriz, temos:

$$\ddot{\mathcal{F}}_{\alpha} = \frac{\partial \Psi_{\alpha}(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tilde{\boldsymbol{\xi}}^\top} = \tilde{\mathbf{X}}^{*\top} [\mathbf{M}_1 + \mathbf{M}_2] \tilde{\mathbf{X}}^* - \tilde{\mathbf{K}}^*(\boldsymbol{\alpha}), \quad (\text{A.3})$$

sob a condição de  $\mathbf{M}_1 = [\mathbf{D}\dot{\mathbf{V}}^{-1}\mathbf{D} + \dot{\mathbf{D}}\mathbf{V}^{-1}]\mathbf{C}$ ,  $\dot{\mathbf{D}} = \text{diag}\{\partial^2\mu_1/\partial\eta_1^2, \dots, \partial^2\mu_n/\partial\eta_n^2\}$ , bem como  $\dot{\mathbf{V}} = \text{diag}\{\partial V_1/\partial\mu_1, \dots, \partial V_n/\partial\mu_n\}$ ,  $\mathbf{C} = \text{diag}\{C_1, \dots, C_n\}$ , no qual  $C_i = y_i - \mu_i$ ,  $\mathbf{M}_2 = -\mathbf{W}$  e  $\tilde{\mathbf{K}}^*(\boldsymbol{\alpha}) = \text{blocodiag}\{\mathbf{O}_{pp}, \bar{\mathbf{S}}(\boldsymbol{\alpha})\}$  avaliados em  $\tilde{\boldsymbol{\xi}} = \hat{\boldsymbol{\xi}}$  e  $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ .

### CASO PONDERADO

Para avaliar a influência das perturbações de casos, o logaritmo da função de verossimilhança perturbada é definido por:

$$\mathcal{F}_{\alpha}(\tilde{\boldsymbol{\xi}}) = L(\tilde{\boldsymbol{\xi}}) + \frac{1}{2} \sum_{k=1}^r \alpha_k \boldsymbol{\tau}_k^\top \bar{\mathbf{S}}_k \boldsymbol{\tau}_k, \quad (\text{A.4})$$

Contudo, perturbamos  $L(\tilde{\boldsymbol{\xi}})$  de tal forma que  $L(\tilde{\boldsymbol{\xi}}) = \sum_{i=1}^n \omega_i \{\phi^{-1}[y_i\theta_i - b(\theta_i)] + c(y_i;\phi)\}$ . Aqui, objetivo é identificar pontos influentes entre todas as observações. Neste caso, o vetor correspondente à não perturbação é o vetor  $\boldsymbol{\omega}_0 = (1, 1, \dots, 1)^\top$   $n$ -dimensional. Note que  $\boldsymbol{\Delta} = (\boldsymbol{\Delta}_{\boldsymbol{\beta}}^\top, \boldsymbol{\Delta}_{\boldsymbol{\tau}_1}^\top, \boldsymbol{\Delta}_{\boldsymbol{\tau}_2}^\top, \dots, \boldsymbol{\Delta}_{\boldsymbol{\tau}_k}^\top)^\top$  é uma matriz  $(p + s_v) \times n$ , avaliada em  $\tilde{\boldsymbol{\xi}} = \hat{\boldsymbol{\xi}}$  e em  $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ , onde:

$$\begin{aligned} \frac{\partial \mathcal{F}_{\alpha}(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \beta_j} &= \phi^{-1} \sum_{i=1}^n \omega_i \left\{ y_i \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_j} - \frac{db(\theta_i)}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_j} \right\} \\ &= \phi^{-1} \sum_{i=1}^n \omega_i \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) x_{ij} \right\} \quad \text{e, portanto,} \\ \frac{\partial^2 \mathcal{F}_{\alpha}(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \beta_j \partial \omega_i} &= \frac{\partial}{\partial \omega_i} \left[ \phi^{-1} \sum_{i=1}^n \omega_i \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) x_{ij} \right\} \right] \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) x_{ij} \right\}. \end{aligned}$$

Da mesma forma, tem-se que:

$$\begin{aligned} \frac{\partial \mathcal{F}_{\alpha}(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tau_{kj}} &= \phi^{-1} \sum_{i=1}^n \left\{ y_i \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \tau_{kj}} - \frac{db(\theta_i)}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \tau_{kj}} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \\ &= \phi^{-1} \sum_{i=1}^n \omega_i \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) \tilde{b}_{kij} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \end{aligned}$$

e, portanto,

$$\begin{aligned} \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \tau_{kj} \partial \omega_i} &= \frac{\partial}{\partial \omega_i} \left[ \phi^{-1} \sum_{i=1}^n \omega_i \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) \tilde{b}_{kij} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \right] \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_i - \mu_i) \tilde{b}_{kij} \right\}. \end{aligned}$$

Colocando na forma matricial, obtemos

$$\Delta_\beta = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \beta \partial \omega^\top} \phi^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{D}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \quad (\text{A.5})$$

$$\Delta_{\tau_k} = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \tau_k \partial \omega^\top} = \phi^{-1} \tilde{\mathbf{B}}_k^\top \mathbf{W} \mathbf{D}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \quad (\text{A.6})$$

em que  $\tilde{\mathbf{B}}_k$  é uma matriz de dimensão  $(n \times v_k)$  formada pelos componentes  $B_k^{(m_k)}(t_{ik})$  da base de De Boor (1978) após ser centralizada,  $i = 1, 2, \dots, n$  e  $j = 1, 2, \dots, v_k$ .

## ESQUEMA DE PERTURBAÇÃO NA VARIÁVEL RESPOSTA

Para este caso, considera-se que a variável  $y_i$  é perturbada de tal maneira que  $y_{\omega_i} = y_i + V_i^{1/2}$  em que o vetor de não perturbação  $\boldsymbol{\omega}_0$  é tal que  $\omega_i = 0 \forall i = 1, 2, \dots, n$ . Assim como no caso anterior,  $\Delta = (\Delta_\beta^\top, \Delta_{\tau_1}^\top, \Delta_{\tau_2}^\top, \dots, \Delta_{\tau_k}^\top)^\top$ , porém

$$\begin{aligned} \frac{\partial \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \beta_j} &= \phi^{-1} \sum_{i=1}^n \left\{ y_{\omega_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_j} - \frac{db(\theta_i)}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_j} \right\} \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_{\omega_i} - \mu_i) x_{ij} \right\} \quad \text{e, portanto,} \\ \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \beta_j \partial \omega_i} &= \frac{\partial}{\partial \omega_i} \left[ \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_{\omega_i} - \mu_i) x_{ij} \right\} \right] \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{\partial}{\partial \omega_i} \left[ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \right] (y_{\omega_i} - \mu_i) x_{ij} \right. \\ &\quad \left. + \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial}{\partial \omega_i} [(y_{\omega_i} - \mu_i) x_{ij}] \right\} \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} V_i^{1/2} x_{ij} \right\}. \end{aligned}$$

Da mesma forma, tem-se que:

$$\begin{aligned} \frac{\partial \mathcal{F}_\alpha(\tilde{\xi}|\omega)}{\partial \tau_{kj}} &= \phi^{-1} \sum_{i=1}^n \left\{ y_{\omega_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \tau_{kj}} - \frac{db(\theta_i)}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \tau_{kj}} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_{\omega_i} - \mu_i) \tilde{b}_{kij} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \end{aligned}$$

e, portanto,

$$\begin{aligned} \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tau_{kj} \partial \omega_i} &= \frac{\partial}{\partial \omega_i} \left[ \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} (y_{\omega_i} - \mu_i) \tilde{b}_{kij} \right\} + \frac{1}{2} \sum_{k=1}^r \alpha_k \bar{\mathbf{S}}_k \boldsymbol{\tau}_k \right] \\ &= \phi^{-1} \sum_{i=1}^n \left\{ \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} V_i^{1/2} \tilde{b}_{kij} \right\}. \end{aligned}$$

Considerando que  $\tilde{\boldsymbol{\xi}} = \hat{\tilde{\boldsymbol{\xi}}}$  e  $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ , temos que, em forma matricial

$$\Delta_{\boldsymbol{\beta}} = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \boldsymbol{\beta} \partial \omega_i} = \phi^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{D}^{-1} \mathbf{V}^{1/2} \quad (\text{A.7})$$

$$\Delta_{\boldsymbol{\tau}_k} = \frac{\partial^2 \mathcal{F}_\alpha(\tilde{\boldsymbol{\xi}}|\boldsymbol{\omega})}{\partial \tau_{kj} \partial \omega_i} = \phi^{-1} \tilde{\mathbf{B}}_k^\top \mathbf{W} \mathbf{D}^{-1} \mathbf{V}^{1/2}. \quad (\text{A.8})$$