

Data Science and Big Data

Taconeli, C.A.

09 de outubro, 2018

Por que GAMLSS?

Introdução

- GAMLSS (*Generalized additive models for location, scale and shape*) é uma recente metodologia de regressão (semi) paramétrica, que contempla uma grande variedade de distribuições, e em que qualquer um de seus parâmetros pode ser modelado em função de covariáveis.
- Permite modelar dados com dispersão não constante, diferentes níveis de assimetria e curtose, relações não lineares entre as variáveis. . .
- O pacote `gamlss` e pacotes complementares permitem ajustar modelos da classe GAMLSS para diferentes tipos e estruturas de dados.

Introdução

- Nesta aula vamos discutir alguns pontos da metodologia por meio de um exemplo referente aos preços de aluguel de imóveis em Munique, 1980 (data set `rent`, pacote `gamlss`).
- Adicionalmente, vamos explorar, de maneira preliminar, recursos computacionais implementados na biblioteca `gamlss` do R.
- Vamos ajustar uma sequência de modelos com nível crescente de complexidade, partindo de uma regressão linear e chegando a um GAMLSS.
- Os slides apresentados na sequência são complementados com scripts em R, disponíveis na página da disciplina.

Dados sobre preços de aluguel em Munique, 1980

- A base de dados dispõe de informações de 1969 imóveis disponíveis para aluguel em nove variáveis, das quais cinco serão usadas na análise:
 - R: variável resposta, é o aluguel mensal em Marcos alemães menos o custo utilitário, calculado ou estimado;
 - Fl: Área construída, em metros quadrados;
 - A: Ano de construção;
 - H: Fator com dois níveis, se o imóvel tem (0) ou não (1) aquecimento central;
 - loc: um fator que classifica a locação do imóvel como abaixo da média (1), na média (2) ou acima da média (3).

Análise exploratória

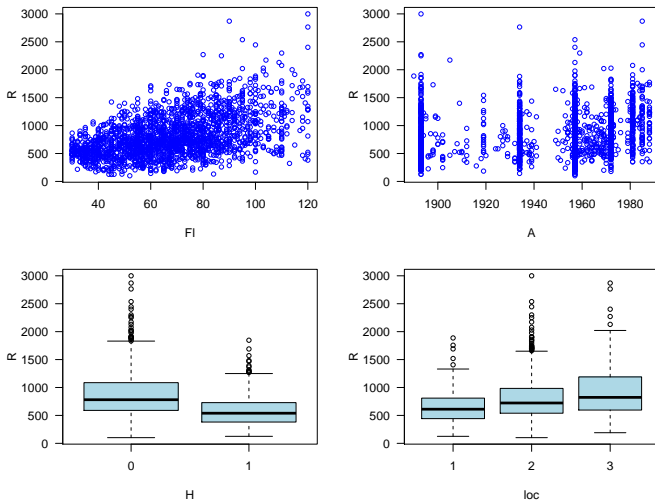


Figura 1: Gráficos para o valor de aluguel (R) vs variáveis explicativas.

Por que GAMLSS?

- Complexidade da relação entre a variável resposta e as variáveis explicativas;
- A variância dos valores de aluguel não é constante para diferentes valores das variáveis explicativas;
- Distribuição dos valores de aluguel é assimétrica, e o nível de assimetria parece variar conforme os valores das variáveis explicativas.

Modelo de regressão linear

- O primeiro modelo a ser ajustado é o de regressão linear, considerando resposta com distribuição Normal.
- Considere y a variável resposta e x_1, x_2, \dots, x_r um conjunto de r covariáveis avaliados em uma amostra de tamanho n .
- O modelo de regressão linear fica definido, para uma amostra de tamanho n , por:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_r x_{ir} + \epsilon_i,$$

com $\epsilon_i \stackrel{ind}{\sim} N(0, \sigma^2)$, para $i = 1, 2, \dots, n$.

Modelo de regressão linear

- O modelo de regressão linear pode ser representado na forma matricial por:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

onde $\mathbf{y} = (y_1, \dots, y_n)'$ é o vetor de respostas, \mathbf{X} é a matriz do modelo $n \times p$ ($p = r + 1$), $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_r)'$ é o vetor de parâmetros de regressão e $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)'$ o vetor de erros.

- Uma forma equivalente (e mais flexível) de representar o modelo de regressão linear é a seguinte:

$$y|\mathbf{x} \stackrel{ind}{\sim} N(\mu_{\mathbf{x}}, \sigma^2),$$

onde $\mu_{\mathbf{x}} = \beta_0 + \beta_1 x_1 + \dots + \beta_r x_r$.

Modelo de regressão linear

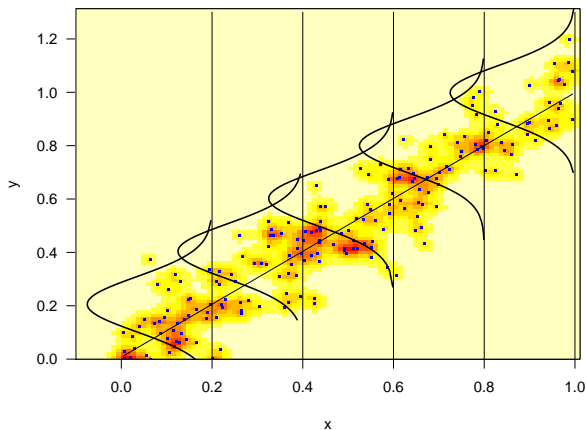


Figura 2: Ilustração de modelo de regressão linear com uma covariável.

Modelo de regressão linear

- O método de mínimos quadrados é usualmente aplicado na estimação dos parâmetros do modelo (β 's);
- O estimador de mínimos quadrados de β é o vetor $\hat{\beta}$ que minimiza a soma de quadrados dos erros, dada por:

$$SQE(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta).$$

- O vetor $\hat{\beta}$ pode ser obtido de forma analítica, resultando em:

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.$$

Modelo de regressão linear

- Sob a especificação do modelo de regressão linear, o estimador de mínimos quadrados é também o estimador de máxima verossimilhança de β .
- O estimador de máxima verossimilhança de σ^2 é dado por:

$$\hat{\sigma}^2 = \frac{SQ_{Res}}{n} = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n},$$

sendo viciado para σ^2 . Um estimador não viciado de σ^2 é dado por:

$$s^2 = \frac{SQ_{Res}}{n - p} = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n - p}.$$

Modelo de regressão linear

- Voltando à análise dos dados de aluguéis de imóveis, o seguinte modelo de regressão linear é proposto:

$$y|\mathbf{x} \sim \text{Normal}(\mu_{\mathbf{x}}, \sigma^2),$$

em que

$$\begin{aligned} \mu_{\mathbf{x}} = & \beta_0 + \beta_1 \times FI + \beta_2 \times A + \beta_3 \times I(H = 1) \\ & + \beta_4 \times I(loc = 2) + \beta_5 \times I(loc = 3), \end{aligned}$$

sendo $I(\cdot)$ é a função indicadora, tal que $I(H = 1) = 1$ para os imóveis sem aquecimento central ($H = 1$) e $I(H = 1) = 0$ para os imóveis com aquecimento central ($H=0$).

Modelo de regressão linear

- A equação do modelo de regressão linear ajustado é a seguinte:

$$\hat{\mu}_x = -2775.04 + 8.84 \times FI + 1.48 \times A - 204.76 \times I(H = 1) \\ + 134.05 \times I(loc = 2) + 209.58 \times I(loc = 3).$$

- Além disso:

$$\log(\hat{\sigma}) = 5.73165,$$

tal que $\hat{\sigma} = 308.48$.

Modelo linear generalizado

- Modelos lineares generalizados configuram extensões dos modelos de regressão linear, apresentando como diferenciais:
 - Permitem modelar respostas com distribuição pertencente à família exponencial;
 - A relação entre a média de y e as covariáveis é determinada por uma função de ligação monotônica $g(\cdot)$;
 - A estimação dos parâmetros do modelo se dá por um algoritmo de mínimos quadrados ponderados iterativamente.

Modelo linear generalizado

- Uma variável aleatória y tem distribuição pertencente à família exponencial se sua função (densidade) de probabilidade tiver a seguinte forma:

$$f(y; \mu, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right\},$$

em que θ e ϕ são os parâmetros canônico e de dispersão, respectivamente.

- A média e a variância de y são dadas, respectivamente, por $E(y) = b'(\theta)$ e $Var(y) = \phi V(\mu)$, onde $V(\mu) = b''[\theta(\mu)]$ é a chamada *função de variância*.

Modelo linear generalizado

- Dentre as principais distribuições contempladas pela teoria de MLG estão a normal ($V(\mu) = 1$), binomial ($V(\mu) = \mu(1 - \mu)$), Poisson ($V(\mu) = \mu$), Gamma ($V(\mu) = \mu^2$) e normal inversa ($V(\mu) = \mu^3$).
- Um modelo linear generalizado pode ser representado, genericamente, da seguinte forma:

$$y|\mathbf{x} \stackrel{ind}{\sim} \epsilon(\mu_{\mathbf{x}}, \phi),$$

em que ϵ denota uma particular distribuição da família exponencial, ϕ é um parâmetro de dispersão e

$$g(\mu_{\mathbf{x}}) = \beta_0 + \beta_1 x_1 + \dots + \beta_r x_r.$$

Modelo linear generalizado

- Na aplicação dos dados sobre os valores de aluguel de imóveis, vamos considerar a distribuição Gamma, que é uma alternativa para a modelagem de dados contínuos assimétricos.
- Uma variável aleatória com distribuição Gamma de média μ e parâmetro de dispersão ϕ tem a função densidade de probabilidade dada por:

$$f(y; \mu, \phi) = \frac{y^{\frac{1}{\phi}-1} \exp\left(-\frac{y}{\phi\mu}\right)}{(\phi\mu)^{1/\phi} \Gamma(1/\phi)}, \quad y > 0, \mu > 0, \phi > 0.$$

- No pacote `gamlss` a distribuição Gamma é parametrizada pelo parâmetro de escala $\sigma = \sqrt{\phi}$.

Modelo linear generalizado

- Como $\mu > 0$, uma escolha adequada para o modelo é a função de ligação logarítmica, de forma que o modelo a ser ajustado fica especificado por:

$$y|\mathbf{x} \stackrel{ind}{\sim} \text{Gamma}(\mu_{\mathbf{x}}, \sigma), \quad y > 0, \mu > 0, \sigma > 0,$$

com

$$\begin{aligned} \log(\mu_{\mathbf{x}}) = & \beta_0 + \beta_1 \times FI + \beta_2 \times A + \beta_3 \times I(H = 1) \\ & + \beta_4 \times I(loc = 2) + \beta_5 \times I(loc = 3). \end{aligned}$$

Modelo linear generalizado

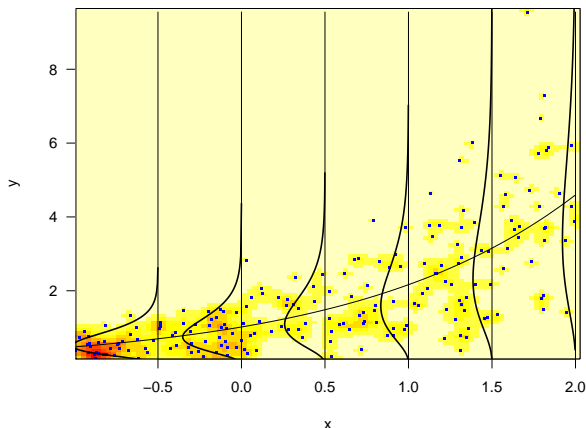


Figura 3: Ilustração de modelo de regressão Gamma com função de ligação exponencial.

Modelo linear generalizado

- O ajuste do modelo linear generalizado com resposta Gamma resulta em $y|\mathbf{x} \sim \text{Gamma}(\hat{\mu}_{\mathbf{x}}, \hat{\sigma})$, tal que:

$$\begin{aligned} \log(\hat{\mu}_{\mathbf{x}}) = & 2.8649 + 0.0106 \times FI + 0.0015 \times A - 0.3001 \times I(H = 1) \\ & + 0.1907 \times I(loc = 2) + 0.2641 \times I(loc = 3), \end{aligned}$$

com

$$\log(\hat{\sigma}) = -0.9822,$$

tal que $\hat{\sigma} = 0.3745$.

- Informações mais detalhadas sobre o ajuste encontram-se no script disponível na página da disciplina.

Modelo generalizado aditivo

- Modelos generalizados aditivos configuram extensões dos MLGs em que o efeito de ao menos uma das covariáveis é incorporado ao preditor linear através de uma função suave, sem parâmetros associados.
- Um modelo generalizado aditivo pode ser representado, de forma geral, por:

$$y|\mathbf{x} \stackrel{ind}{\sim} \epsilon(\mu_{\mathbf{x}}, \phi),$$

onde

$$g(\mu_{\mathbf{x}}) = \beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \dots + s_{j+1}(x_{j+1}) + \dots + s_r(x_r),$$

em que s_k é uma função suave não paramétrica aplicada à covariável x_k , $k = j + 1, \dots, r$.

Modelo generalizado aditivo

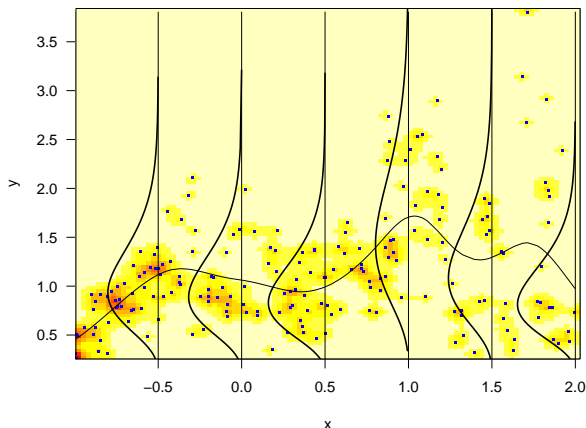


Figura 4: Ilustração de modelo de regressão Gamma com função suave para a média.

Modelo generalizado aditivo

- Modelos aditivos são mais flexíveis do que modelos totalmente paramétricos.
- Os efeitos das covariáveis inseridas ao preditor por meio de funções suaves podem ser interpretados usando gráficos apropriados;
- O pacote `gamlss` oferece diferentes alternativas de funções suaves a serem usadas em modelos aditivos, que serão abordadas posteriormente.
- Na aplicação referente aos preços de aluguel vamos considerar a inclusão de efeitos aditivos (não paramétricos) para a área e o ano de construção do imóvel.

Modelo generalizado aditivo

- O modelo generalizado aditivo a ser ajustado aos dados de preços de aluguel, considerando resposta Gamma, é especificado da seguinte forma:

$$y|\mathbf{x} \stackrel{ind}{\sim} \text{Gamma}(\mu_{\mathbf{x}}, \sigma), \quad y > 0, \mu > 0, \sigma > 0,$$

com

$$\begin{aligned} \log(\mu_{\mathbf{x}}) = & \beta_0 + s_1(Fl) + s_2(A) + \beta_1 \times I(H = 1) \\ & + \beta_2 \times I(loc = 2) + \beta_3 \times I(loc = 3), \end{aligned}$$

em que s_1 e s_2 são suavizadores não paramétricos aplicados às variáveis Fl e A , respectivamente.

Modelo generalizado aditivo

- O modelo ajustado fica dado por:

$$\log(\hat{\mu}_x) = 3.0851 + s_1(Fl) + s_2(A) - 0.3008 \times I(H = 1) \\ + 0.1887 \times I(loc = 2) + 0.2720 \times I(loc = 3),$$

com

$$\log(\hat{\sigma}) = -1.0019,$$

tal que $\hat{\sigma} = 0.33672$.

Modelagem do parâmetro de escala

- Até o momento consideramos apenas modelos em que a média (parâmetro de locação) da distribuição varia conforme os valores das covariáveis;
- Modelos mais gerais permitem incluir covariáveis também na modelagem de outros parâmetros da distribuição (por exemplo, para o parâmetro de escala).
- Para a distribuição Gamma, por exemplo, temos que $Var(y) = \sigma^2 \mu^2$, ou seja, $\sigma = \sqrt{Var(y)}/\mu$ é o coeficiente de variação de y .
- Podemos modelar σ em função de covariáveis, permitindo avaliar se o coeficiente de variação muda conforme os valores de alguma (ou algumas) delas.

Modelagem do parâmetro de escala

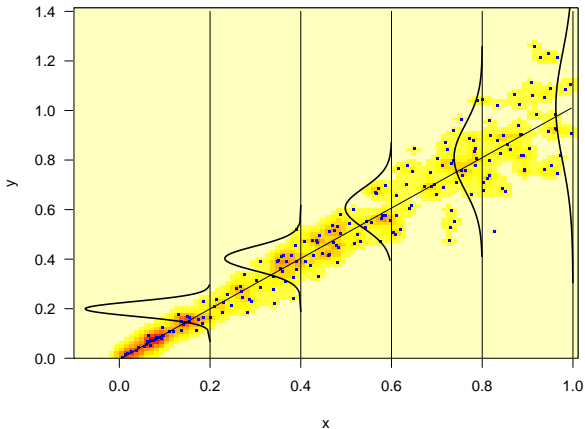


Figura 5: Ilustração - distribuição normal com média e variância dependentes de x .

Modelagem do parâmetro de escala

- Uma formulação geral para o modelo, neste caso, seria da seguinte forma:

$$y|\mathbf{x} \stackrel{ind}{\sim} D(\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}),$$

onde

$$\begin{aligned} g_1(\mu_{\mathbf{x}}) &= \beta_{10} + \beta_{11}x_1 + \dots + \beta_{1j_1}x_{j_1} + \dots + s_{j_1+1}(x_{j_1+1}) + \dots + s_{r_1}(x_{r_1}) \\ g_2(\sigma_{\mathbf{x}}) &= \beta_{20} + \beta_{21}x_1 + \dots + \beta_{2j_2}x_{j_2} + \dots + s_{j_2+1}(x_{j_2+1}) + \dots + s_{r_2}(x_{r_2}), \end{aligned}$$

em que D representa alguma distribuição com dois parâmetros (μ e σ), que são funções de covariáveis inseridas por meio de termos paramétricos ou não paramétricos.

Modelagem do parâmetro de escala

- Neste contexto, vamos retomar a análise dos dados de preços de aluguel com o ajuste de modelos com resposta Gamma e normal inversa, inserindo covariáveis também na modelagem do parâmetro de escala.
- Em ambos os casos vamos considerar função de ligação logarítmica tanto para μ quanto para σ , produzindo os seguintes modelos:

$$\log(\mu_x) = \beta_{10} + s_{11}(Fl) + s_{12}(A) + \beta_{11} \times I(H = 1) \\ + \beta_{12} \times I(loc = 2) + \beta_{13} \times I(loc = 3),$$

e

$$\log(\sigma_x) = \beta_{20} + s_{21}(Fl) + s_{22}(A) + \beta_{21} \times I(H = 1) \\ + \beta_{22} \times I(loc = 2) + \beta_{23} \times I(loc = 3).$$

Modelagem do parâmetro de escala

- O modelo ajustado com resposta Gamma (que produziu menor AIC do que com resposta normal inversa) produziu os seguintes resultados:

$$\log(\hat{\mu}_x) = 2.8844 + s_{11}(FI) + s_{12}(A) - 0.2918 \times I(H = 1) \\ + 0.1938 \times I(loc = 2) + 0.2734 \times I(loc = 3),$$

com

$$\log(\hat{\sigma}_x) = 5.9225 + s_{21}(FI) + s_{22}(A) + 0.0659 \times I(H = 1) \\ - 0.1166 \times I(loc = 2) - 0.1702 \times I(loc = 3).$$

Modelos generalizados aditivos para locação, escala e forma

- Os modelos considerados anteriormente baseiam-se em distribuições com dois parâmetros, permitindo modelar locação (média) e escala (dispersão) da distribuição da variável em função de covariáveis.
- A metodologia GAMLSS contempla distribuições de probabilidade com até quatro parâmetros, permitindo modelar outros parâmetros da distribuição, associados, por exemplo, à assimetria e curtose.
- Desta forma, dispõe-se de uma ampla variedade de modelos, flexíveis e úteis para a análise de dados com distribuições de diferentes formas.

Modelos generalizados aditivos para localização, escala e forma

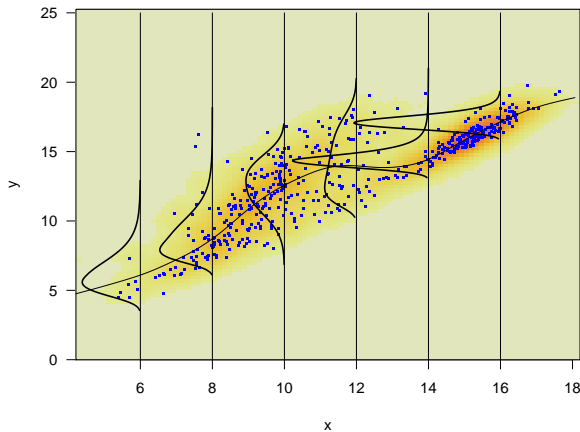


Figura 6: Ilustração - distribuição com localização, dispersão e forma dependentes de

Modelos generalizados aditivos para localização, escala e forma

- Um modelo generalizado aditivo para localização, escala e forma pode ser definido, de forma geral, por:

$$y|\mathbf{x} \stackrel{ind}{\sim} D(\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}, \nu_{\mathbf{x}}, \tau_{\mathbf{x}}),$$

onde

$$\begin{aligned} g_1(\mu_{\mathbf{x}}) &= \beta_{10} + \beta_{11}x_1 + \dots + \beta_{1j_1}x_{j_1} + \dots + s_{j_1+1}(x_{j_1+1}) + \dots + s_{r_1}(x_{r_1}) \\ g_2(\sigma_{\mathbf{x}}) &= \beta_{20} + \beta_{21}x_1 + \dots + \beta_{2j_2}x_{j_2} + \dots + s_{j_2+1}(x_{j_2+1}) + \dots + s_{r_2}(x_{r_2}) \\ g_3(\nu_{\mathbf{x}}) &= \beta_{30} + \beta_{31}x_1 + \dots + \beta_{3j_3}x_{j_3} + \dots + s_{j_3+1}(x_{j_3+1}) + \dots + s_{r_3}(x_{r_3})', \\ g_4(\tau_{\mathbf{x}}) &= \beta_{40} + \beta_{41}x_1 + \dots + \beta_{4j_4}x_{j_4} + \dots + s_{j_4+1}(x_{j_4+1}) + \dots + s_{r_4}(x_{r_4}) \end{aligned}$$

em que $D(\mu, \sigma, \nu, \tau)$ é uma distribuição de quatro parâmetros e ν e τ são parâmetros de forma.

Modelos generalizados aditivos para locação, escala e forma - Distribuições

- O pacote `gamlss` dispõe de aproximadamente 100 distribuições implementadas. Além disso:
 - É possível ao usuário implementar novas distribuições;
 - Versões truncadas ou censuradas podem ser definidas a partir das distribuições originais;
 - Pode-se criar novas distribuições a partir de misturas das distribuições originais;
 - Distribuições discretas podem ser criadas a partir de distribuições originalmente contínuas;
 - Distribuições com suporte nos intervalos $(0, \infty)$ ou $(0, 1)$ podem ser criadas a partir de variáveis com suporte no intervalo $(-\infty, \infty)$.

Modelos generalizados aditivos para locação, escala e forma - Termos aditivos

- Termos aditivos podem ser adicionados ao modelo de diferentes formas, usando, por exemplo:
 - Penalized B-splines (P-splines);
 - Monotone P-splines;
 - Cycle P-splines;
 - Varying coefficient P-splines;
 - Cubic smoothing P-splines;
 - Loess curve fitting;
 - Fractional polynomials;
 - Random effects;
 - Ridge regression;
 - Nonlinear parametric fits.

Modelos generalizados aditivos para locação, escala e forma - Estimação

- Quando o modelo especificado contém apenas termos paramétricos, o método da máxima verossimilhança é aplicado na estimação dos parâmetros;
- Em situações mais gerais (ex: para modelos incluindo suavizadores) o método da máxima verossimilhança penalizada é aplicado;
- O pacote `gamlss` dispõe de dois algoritmos distintos para o ajuste dos modelos: CG (Cole e Green) e RS (Rigby e Stasinopoulos), que podem ainda ser usados de forma combinada.

Modelos generalizados aditivos para locação, escala e forma - Aplicação

- Dando sequência à análise dos dados de preços de aluguel, vamos considerar, como alternativa à distribuição Gamma, a distribuição Box-Cox Cole e Green (BCCGo);
- A distribuição BCCGo tem três parâmetros (μ , σ e τ), sendo τ o parâmetro de forma da distribuição;
- Cada um dos parâmetros pode ou não ser modelado em função de covariáveis. Além disso, diferentes conjuntos de covariáveis podem ser incluídos para cada parâmetro;
- Modelos com diferentes especificações (distribuições, termos aditivos, covariáveis...) podem ter seus ajustes comparados usando o critério de informação de Akaike (AIC) ou Akaike generalizado (GAIC), dentre outros.

Modelos generalizados aditivos para locação, escala e forma - Aplicação

- No modelo especificado para a análise dos dados, todas as covariáveis foram incluídas na modelagem dos três parâmetros (μ , σ e ν).
- Como funções de ligação para cada parâmetro adotou-se o default do pacote `gamlss` para a distribuição BCCGo (log, log e identidade, respectivamente);
- Novamente, funções suaves foram consideradas para as variáveis FI e A .

Modelos generalizados aditivos para locação, escala e forma - Aplicação

- Dessa forma, o seguinte modelo foi proposto:

$$y|\mathbf{x} \stackrel{ind}{\sim} BCCGo(\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}, \nu_{\mathbf{x}}),$$

onde

$$\begin{aligned} \log(\mu_{\mathbf{x}}) = & \beta_{10} + s_{11}(Fl) + s_{12}(A) + \beta_{11} \times I(H = 1) \\ & + \beta_{12} \times I(loc = 2) + \beta_{13} \times I(loc = 3), \end{aligned}$$

$$\begin{aligned} \log(\sigma_{\mathbf{x}}) = & \beta_{20} + s_{21}(Fl) + s_{22}(A) + \beta_{21} \times I(H = 1) \\ & + \beta_{22} \times I(loc = 2) + \beta_{23} \times I(loc = 3), \end{aligned}$$

$$\nu_{\mathbf{x}} = \beta_{30} + s_{31}(Fl) + s_{32}(A) + \beta_{31} \times I(H = 1)$$

Modelos generalizados aditivos para locação, escala e forma - Aplicação

- Como resultado temos o seguinte modelo ajustado:

$$\log(\hat{\mu}_x) = 2.0285 + s_{11}(Fl) + s_{12}(A) - 0.3213 \times I(H = 1) \\ + 0.1853 \times I(loc = 2) + 0.2742 \times I(loc = 3),$$

$$\log(\hat{\sigma}_x) = 6.6534 + s_{21}(Fl) + s_{22}(A) + 0.0819 \times I(H = 1) \\ - 0.0851 \times I(loc = 2) - 0.1410 \times I(loc = 3),$$

$$\hat{\nu}_x = -3.1601 + s_{31}(Fl) + s_{32}(A) - 0.2866 \times I(H = 1) \\ - 0.1711 \times I(loc = 2) - 0.0845 \times I(loc = 3).$$

Modelos generalizados aditivos para locação, escala e forma - Resumo

- Algumas propriedades dos GAMLSS:
 - GAMLSS configura uma metodologia flexível para modelos de regressão;
 - Permite assumir diversas distribuições de probabilidades para a variável resposta;
 - Todos os parâmetros da distribuição podem ser modelados em função de covariáveis;
 - Diferentes tipos de termos aditivos podem ser incluídos nos preditores de cada parâmetro;
 - GAMLSS estende diversas outras metodologias (como LM, GLM e GAM) permitindo modelar dados com super dispersão, excesso de zeros, diferentes níveis de assimetria e curtose. . .