

Modelos Aditivos Generalizados (GAM): Uma visão prática

Eduardo Camargo INPE/DPI



O modelo de regressão simples – uma breve recordação

$$Y_i = \beta_0 + \beta_1 X_i + \xi_i$$

em que:

Y: é denotada de variável dependente ou resposta

X: variável independente

 β_0 : intercepto

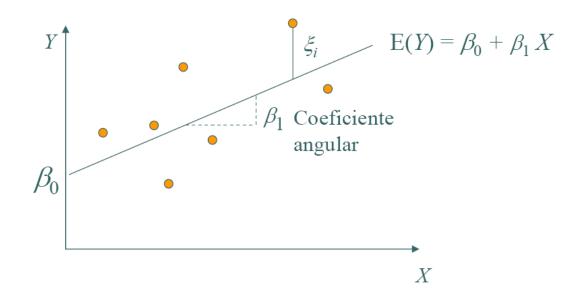
 β_1 : inclinação

 ξ_i : erro aleatório, $\xi_i \sim N(0, \sigma^2)$, $\xi_i = Y_i - \beta_0 + \beta_1 X_i$,



O modelo de regressão simples – uma breve recordação

$$Y_{i} = \beta_{0} + \beta_{1}X_{i} + \xi_{i}$$





- O modelo de regressão simples uma breve recordação $Y_i = \beta_0 + \beta_1 X_i + \xi_i$
- Em geral não se conhece os parâmetros β_0 e β_1 .
- Eles podem ser estimados através de dados obtidos por amostras.
- O método utilizado na estimação de β_0 e β_1 é o método dos mínimos quadrados, o qual considera os desvios dos Y_i de seu valor esperado:

$$\xi_i = Y_i - \beta_0 + \beta_1 X_i$$

• O método dos mínimos quadrados requer que consideremos a soma dos *n* desvios quadrados, denotado por *Q*:

$$Q = \sum_{i=1}^{n} [Y_i - \beta_0 - \beta_1 X_i]^2$$



• O modelo de regressão simples – uma breve recordação $Y_i = \beta_0 + \beta_1 X_i + \xi_i$

• De acordo com o método dos mínimos quadrados, os estimadores de
$$\beta_0$$
 e β_1 são aqueles, denotados por b_0 e b_1 , que tornam mínimo o valor de Q .

• Derivando: $\frac{\partial Q}{\partial \beta_0} = -2\sum_{i=1}^n [Y_i - \beta_0 - \beta_1 X_i] \qquad \frac{\partial Q}{\partial \beta_1} = -2\sum_{i=1}^n [Y_i - \beta_0 - \beta_1 X_i] X_i$

• Igualando as equações a zero obtém-se os valores b_0e b_1 que minimizam Q:

 $b_{1} = \frac{\sum_{i=1}^{n} (X_{i} - \overline{X})(Y_{i} - \overline{Y})}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}}$ $b_{0} = \overline{Y} - b_{1}\overline{X}$ $\hat{Y} = b_{0} + b_{1}X$ $e_{i} = Y_{i} - \hat{Y}_{i} \text{ (resíduo)}$ 29/09/2015



■ Síntese - GAM

Modelo de Regressão Linear Multiplo

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \dots$$

- Modelo Linear Generalizado (MLG). Ingredientes básicos:
 - 1. K valores independentes Y_1 , ..., Y_K , de uma variável resposta que <u>segue</u> <u>uma distribuição da família exponencial</u>, com valor esperado $E(Y_i) = \mu_i$;
 - 2. Uma função de ligação, denotada por $g(\mu_i)$, tal que:

$$g(\mu_i) = \Sigma \mathbf{X} \boldsymbol{\beta}$$

X: vetor das variáveis explicativas.

β: representa o vetor de parâmetros a serem estimados.

$$g(\mu_i) = E(Y_i)$$



Síntese - GAM

Modelo Aditivo Generalizado - GAM

• É uma extensão do *MLG*, em que o termo $\Sigma \mathbf{X} \boldsymbol{\beta}$ é substituído por $\Sigma f(\mathbf{X})$, assim:

$$g(\mu_i) = \Sigma f(\mathbf{X})$$

- f(X) é uma função não paramétrica (i.e. cuja forma não é especificada)
- f(X) é estimada através de curvas de alisamento (ex: splines).
- A curva alisada permite descrever a forma e revelar possíveis não linearidades nas relações estudadas, uma vez que não apresenta a estrutura rígida de uma função paramétrica, como nos MLG's.



■ Síntese – GAM: exemplos em R

Pacote MGCV

require(mgcv)

Modelo sem suavização em *vexp* (variável explicativa) e sem efeito espacial

fit <- gam(vresp ~ vexp, family=binomial, data=gam.data)

Modelo com suavização em vexp e sem efeito espacial

fit <- gam(vresp ~ s(vexp), family=binomial, data=gam.data)

Modelo sem suavização em vexp e com efeito espacial s(x,y)

fit <- gam(vresp \sim vexp + s(x,y), family=binomial, data=gam.data)

Modelo com suavização em *vexp* e com efeito espacial s(x,y)

fit <- gam (vresp \sim s(vexp) + s(x,y), family=binomial, data=gam.data)



SPGAM um modelo alternativo ao GAM

•
$$g(u_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + ... + k(s)$$

- Uma abordagem semiparamétrica.
- k(s) é uma função de Kernel, em que s é um vetor de coordenadas da variável resposta.
- Exemplo:

Distribuição espacial do risco: modelagem da mortalidade infantil em Porto Alegre, Rio Grande do Sul, Brasil (silvia E. Shimakura, Marilia Sá Carvalho, Denise R. G. C. Aerts, Rui Flores)

$$\begin{split} g(u_i) &= \log \left\{ \frac{p(s,x)}{1-p(s,x)} \right\} = \beta_0 + \beta_1 sexo + \beta_2 peso + \\ \beta_3 idade + \beta_4 inst + \beta_5 ges + \\ \beta_6 grav + \beta_7 parto + k(s). \end{split}$$



Projeto



http://www.dpi.inpe.br/eureqa

- Processo FAPESP No. 2006/53922-9
- Objetivo Principal:

Investigar a correlação entre resistência bacteriana e fatores de risco populacionais (em particular uso populacional de antimicrobianos).

Coordenador Geral

Prof. Dr. Antônio Miguel Vieira Monteiro1

Coordenadores Setoriais

Prof. Dr. Antônio Carlos C. Pignatari³

Dr. Caio Márcio Figueiredo Mendes²

Dr. Carlos Roberto Veiga Kiffer⁶

Dr. Eduardo Celso Gerbi Camargo¹

Pesquisadores

Dra. Soraya Andrade³

Dra. Corina da Costa Freitas1

Dra. Virginia Ragoni de Moraes Correia1

Dra, Elisabete Caria Moraes¹

Dra. Jussara de Oliveira Ortiz¹

Dra. Silvia Emiko Shimakura⁷

Dr. Paulo Justiniano Ribeiro Jr.7

Dr. Getúlio Batista⁴

Msc. Dra. Bianca Lucarevschi4

Msc. Paula Celia Mariko Koga²

Msc. Gabriel Pereira1

Amilton Mouro⁵

Tiago dos Santos Agostinho⁴

¹ Instituto Nacional de Pesquisas Espaciais - INPE

² Fleury Medicina e Saúde

³ UNIFESP / Escola Paulista de Medicina (EPM)/ Laboratório Especial em Microbiologia Clínica (LEMC)

⁴ Universidade de Taubaté - UNITAU

⁵ Hospital Israelita Albert Eisten

⁶ Gestão do Conhecimento Científico - GC2

⁷ Laboratório de Estatística e Geoinformação - LEG, Universidade Federal do Paraná

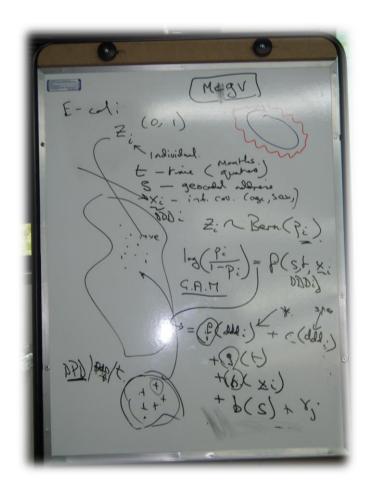


■ 1ª Fase: Construção do Banco de Dados EUREQA

Período de Estudo: 2002 a 2007. Tabelas IMS: contém os endereços dos pontos de vendas de antimicrobianos e do consumo populacional estimado por DDD (Defined Daily Dose) mensal. Tabelas FLEURY e UNIFESP: contém informações de pacientes sobre a resistência **BD** ou sensibilidade das bactérias Haemophilus influenzae, Streptococcus Pneumoniae **EUREQA SP** e Escherichia coli para alguns grupos de antimicrobianos. Dados Geográficos: ruas, unidades territoriais, limites e outros. Dados a serem integrados: informações sócio-econômicas coletadas do censo, dados meteorológicos sobre poluição do ar (H. influenzae, S. Pneumoniae), imagens IKONOS, e outras que possam contribuir para o contexto do EUREQA.



■ 2ª Fase: Discussão e Esboço do Modelo Teórico



Prof. Dr. Trevor Bailey

Prof. Dr. Paulo Justiniano Ribeiro Jr.



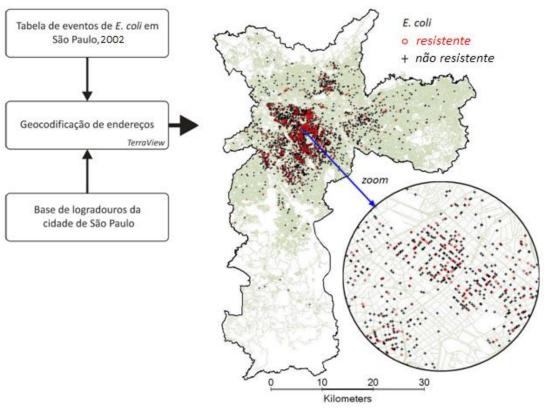
- 3ª Fase: Realização de um Primeiro Estudo
- Investigar se há associação do <u>risco de resistência bacteriana</u>, oriunda de E. coli, com o <u>consumo populacional de ciprofloxaxin</u> expresso pela DDDD.
- > População: mulheres maior que 16 anos.
- Área de estudo: cidade de São Paulo, 2002.
- Covariável desconhecida (DDDD): densidade de uso de antimicrobiano populacional.

$$DDDD = \frac{DDD}{(Pop/1000)*30}$$

"Defined Daily Dose" (OMS) informação conhecida adquirida da IMS

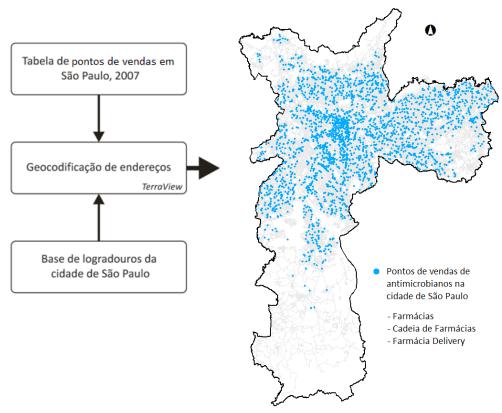


- Preparação dos dados para o GAM
- Espacialização dos eventos, resistentes e não resistentes, decorrente de E. coli



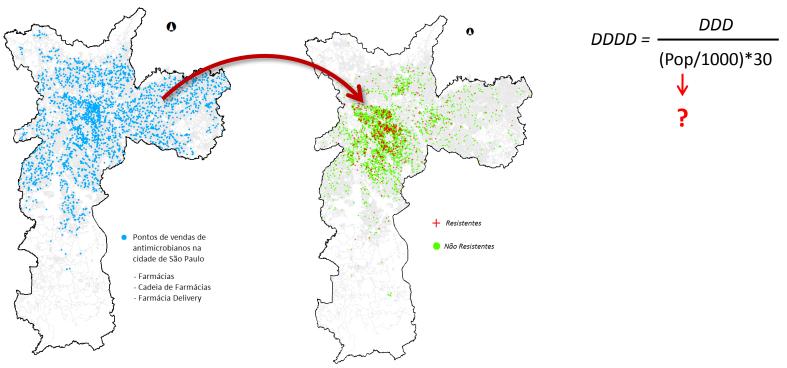


- Preparação dos dados para o GAM
- Espacialização dos pontos de vendas de antimicrobianos na cidade de São Paulo





Preparação dos dados para o GAM
Como calcular DDDD? Como associar valores de DDDD aos eventos resistentes e não resistentes?

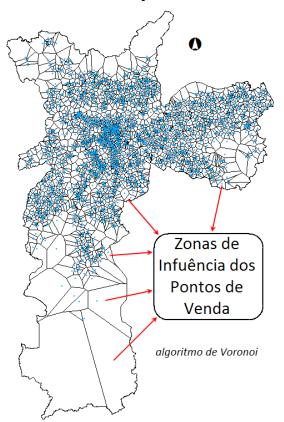




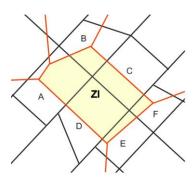
■ Preparação dos dados para o GAM

Solução adotada para o cálculo da DDDD

Parte I: construção de Zonas de Influência (ZI)



Parte II: Interseção dos mapas Zonas de Influência com Unidades Territoriais (Fundação SEADE - SP)



Unidades Territoriais: A, B, C, D, E, F ZI: Zona de Influência

$$\begin{split} Pop_ZI &= Pop_A\left(\frac{\acute{A}rea_ZI_A}{\acute{A}rea_A}\right) + Pop_B\left(\frac{\acute{A}rea_ZI_B}{\acute{A}rea_B}\right) + Pop_C\left(\frac{\acute{A}rea_ZI_C}{\acute{A}rea_C}\right) + \\ &+ Pop_D\left(\frac{\acute{A}rea_ZI_D}{\acute{A}rea_D}\right) + Pop_E\left(\frac{\acute{A}rea_ZI_E}{\acute{A}rea_E}\right) + Pop_F\left(\frac{\acute{A}rea_ZI_F}{\acute{A}rea_F}\right) \end{split}$$

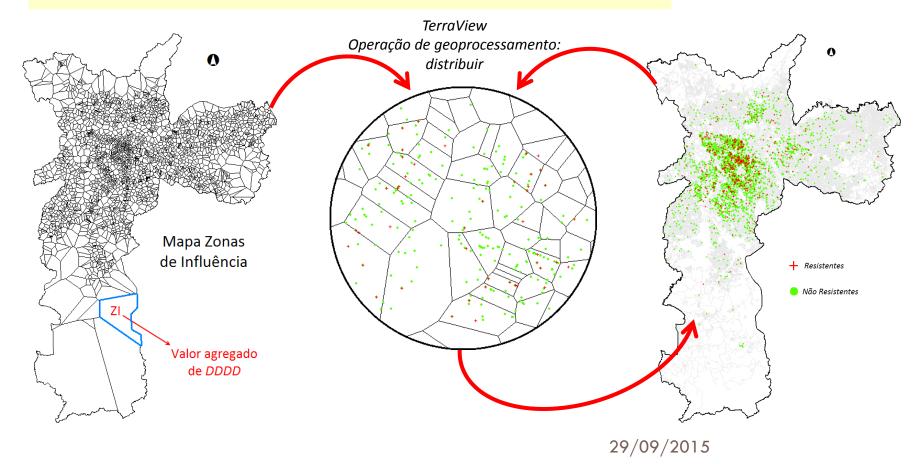
$$DDDD_{-}ZI = \frac{DDD}{\frac{Pop}{1000}*30}$$

29/09/2015



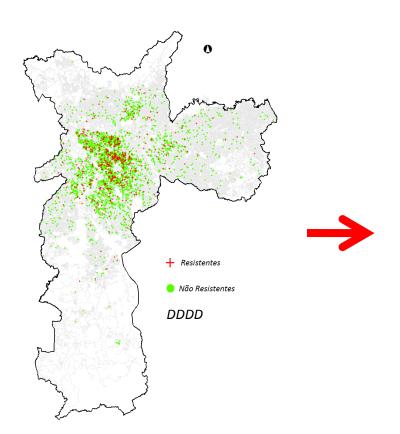
■ Preparação dos dados para o GAM

Associando valores de *DDDD* aos eventos *resistentes* e *não resistentes*





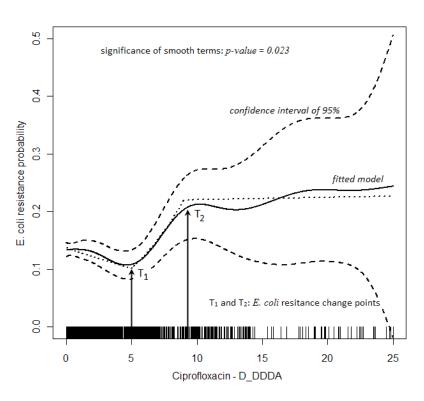
■ Preparação dos dados para o GAM



COORD_X	COORD_Y	Υ	DDDD
330800.361	7395872.440	0	0.342
326544.146	7386268.087	1	0.202
332509.240	7390334.753	1	0.192
329254.299	7389932.687	0	0.416
326670.859	7379104.028	0	0.565
328809.895	7386776.251	0	0.887
329077.650	7386911.517	1	0.887
353109.185	7395543.311	0	0.389
330423.156	7392924.632	0	0.324
332164.392	7388986.068	0	1.212
332358.686	7388984.142	1	1.212
334261.560	7390542.337	0	0.366
326341.197	7393667.842	0	0.567
333549.674	7390623.169	1	0.524
331057.184	7392260.967	0	1.015
331076.514	7392342.553	0	1.015
331292.362	7392105.803	0	1.015
328075.540	7382259.943	1	0.525
325996.190	7387793.954	0	0.480
326128.331	7387648.820	0	0.480
325731.385	7387640.003	0	0.480
325652.046	7387592.818	1	0.480
330605.768	7395941.163	0	0.906
327428.067	7395592.610	1	0.525
333659.078	7392205.810	0	0.197
328254.306	7390217.238	0	0.377
339534.373	7394879.061	1	0.714
^334*33.9*0	7189921,205		0.845



■ Pacote SPGAM (R) com modelo segmentado:

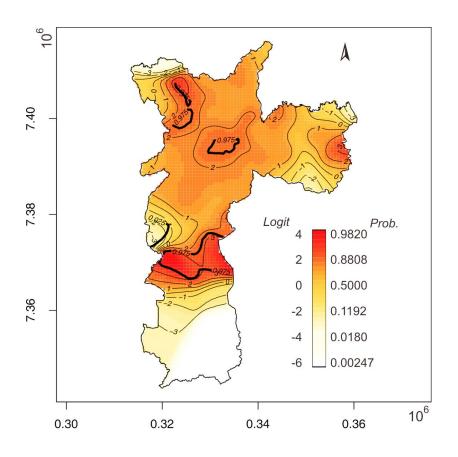


$$log\left\{\frac{p(s,x)}{1-p(s,x)}\right\} = \beta_0 + \beta_1 D_DDDA + \beta_2 I_1 + \beta_3 I_2 + g(s)$$

- β_0 : regression model intercept;
- β_1, β_2 and β_3 : covariates effects;
- D_DDDA: antimicrobial usage mean density;
- $I_1 = 1$ if $T_1 < D_DDDA < T_2$ and $I_1 = 0$ otherwise;
- $I_2 = 1$ if $D_DDDA > T_2$; and $I_2 = 0$ otherwise;
- $T_1 = 5$ and $T_2 = 9$ E. coli resistance change points.

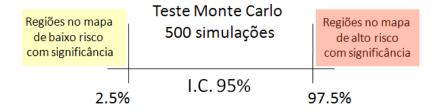


Estimação do mapa de risco



$$log\left\{\frac{p(s,x)}{1-p(s,x)}\right\} = \beta_0 + \beta_1 D_DDDA + \beta_2 I_1 + \beta_3 I_2 + g(s)$$

Os contornos de tolerância e o teste global da hipótese nula de risco não variável espacialmente são obtidos via



p-valor do teste global de risco constante na região

0.001996008