

Quimiometria e Análise de Dados Funcionais

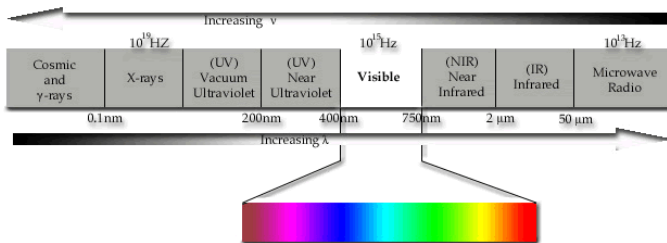
Nancy L. Garcia¹²

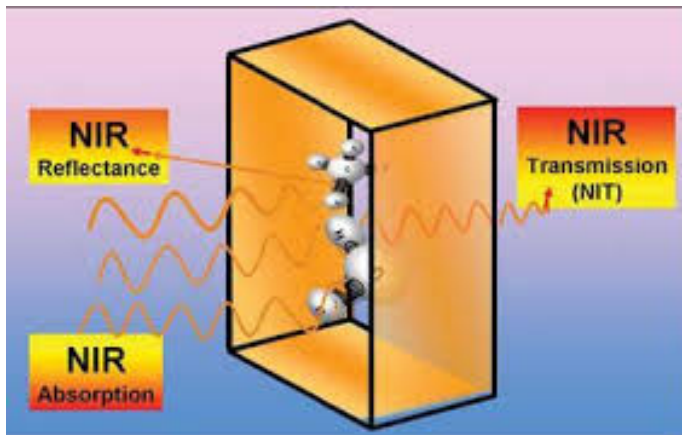
¹IMECC/UNICAMP

²trabalho conjunto com Alex Rodrigo, Alexandra Schmidt, Guilherme Ludwig,
Mariana Rodrigues Motta, Marley Saraiva e Ronaldo Dias

29 de Maio - Dia do Estatístico
Maringá

Comprimento de onda em nanômetros:

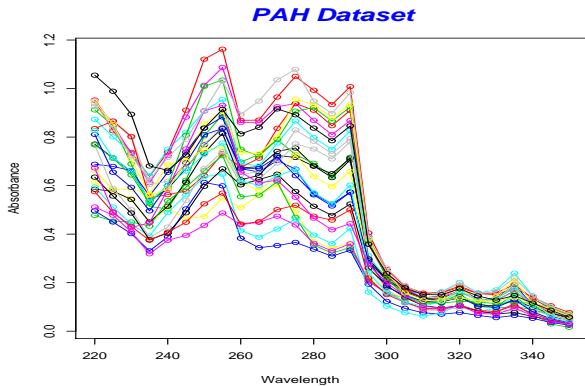




- Quando uma substância é submetida a diferentes comprimentos de onda, os sobretons e combinações na banda do InfraVermelho Próximo (NIR) irá produzir padrões muito complexos que caracterizam os componentes da amostra.
- As diferentes características dos constituintes da amostra se sobrepõem e dão origem a uma curva que é uma soma de várias curvas dependendo da **concentração de cada constituinte**.

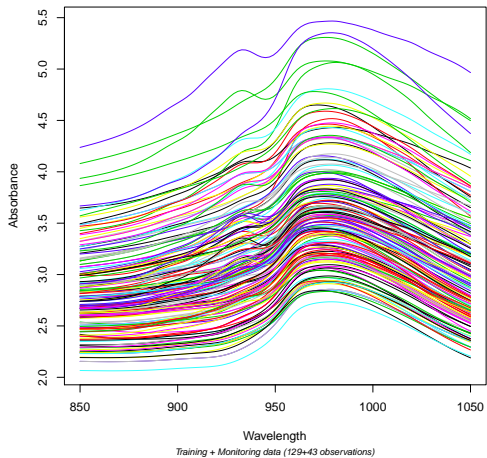
Hidrocarbonetos poliaromáticos – Polyaromatic hydrocarbons(PAH)

- EAS (Electronic Absorption Spectroscopy)
- 25 amostras químicas
- Cada amostra com 10 elementos químicos diferentes (constituintes)
- 27 comprimentos de onda (220 nm – 350 nm)



Tecator

- Cada amostra contém carne moída pura, com diferentes quantidades de umidade, gordura e proteína.
- Conjunto de aprendizado: 129 + 43 amostras
- Conjunto de teste: 43 amostras.
- Cada amostra contém 3 constituintes: umidade, gordura e proteína
- 100 comprimentos de onda (850nm – 1050nm)
- Disponível em domínio público sem qualquer responsabilidade da fonte de dados original
<http://lib.stat.cmu.edu/datasets/tecator>



Problema de calibração

- Existem várias técnicas para medir os constituintes de uma amostra.
- Alguns deles são muito precisos e caros (técnicas analíticas).
- Alguns deles não são tão precisas, mas baratas (por exemplo, NIR espectroscopia).
- Para usar técnicas de baixo custo é necessário **calibrar o instrumento**.
- Por razões práticas, só é possível medir os dados espectrais em um número finito de comprimentos de onda $t_1 < t_2 < \dots < T_T$. Muitas vezes T está por volta de 100–200 ou mais.

Problema de predição

- Depois de calibrar o instrumento
- Tomar novas medidas
- Estimar as concentrações na nova amostra.

- Uma amostra química é constituída de vários constituintes ($\ell = 1, \dots, m$).
- Cada constituinte de interesse é chamado de **analito**.
- y_ℓ : Concentração do analito ℓ
- Amostra é **fechada** se

$$y_1 + \dots + y_m = 1$$

i.e. quando todos os constituintes são analisados

- absorvância do analito ℓ no comprimento de onda t será denotado por

$$\alpha_\ell(t)$$

- Assuma que temos n amostra de composições variadas
- Conjunto de dados Y :
 $y_{i,\ell}$: concentração (medida através de um método de referência - padrão ouro) para analito ℓ na i -ésima amostra
- **Dados espectrais:** Absorbâncias

$$W(t_1), \dots, W(t_T)$$

medidas através de NIR espectroscopia em T comprimentos de onda $t_1 < t_2 < \dots < t_T$.

O modelo: Lei de Beer-Lambert

- Para a i -ésima amostra fechada composta de m constituintes e T comprimentos de onda,

$$W_i(t_j) = \sum_{\ell=1}^K y_{i,\ell} \alpha_{\ell}(t_j) + \epsilon(t_j), \text{ para } j = 1, \dots, T$$

- Geralmente somente um subconjunto de constituintes é analisado. Neste caso, não temos a restrição $y_{i,1} + \dots + y_{i,m} = 1$ mas uma modificação da Lei de Beer-Lambert pode ser usada

$$W_i(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell} a_{\ell}(t_j) + \epsilon(t_j), \text{ para } j = 1, \dots, T$$

onde $\epsilon(t)$ é um processo Gaussiano com função de covariância dada por $\sigma(s, t) = \text{Cov}(\epsilon(s), \epsilon(t))$.

- Modelo de Calibração

$$W_i(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell} \alpha_{\ell}(t_j) + \epsilon_i(t_j), \text{ para } j = 1, \dots, T$$

- Modelo de Predição

$$W_i^*(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell}^* \alpha_{\ell}(t_j) + \epsilon^* i(t_j), \text{ para } j = 1, \dots, T$$

Problema de calibração Dados $y_{i,\ell}$ e $W_i(t_j)$ estimar $\alpha_{\ell}(t)$.

Problema de predição Dado um novo conjunto de dados

$W_i^*(t_j)$, estimar $y_{i,\ell}^*$.

Expansão em bases

Assuma que existe um inteiro positivo K e uma sequência de *knots* ξ tal que

$$\alpha_\ell(t) = \sum_{k=1}^K \beta_{\ell,k} B_k(t),$$

onde $B_k(t)$, $k = 1, \dots, K$ são B-splines cúbicos.

Estrutura de covariância

- Assumimos que a correlação entre os pontos $W_i(t)$ e $W_i(s)$ decai exponencialmente quando $|t - s|$ cresce

$$\sigma^2 \exp(-\phi|t - s|)$$

Aprendizado Supervisionado.

- Dados: $W_i(t)$, $i = 1, \dots, I$ observados em $t = t_1, \dots, t_T$.
- Constantes conhecidas: matriz Y com linhas linearmente independentes contendo as concentrações medidas por um método de referência $y_{i,\ell}$ para $i = 1, \dots, I$ e $\ell = 1, \dots, m$
- Dado o modelo:

$$W_i(t) = \sum_{\ell=1}^m y_{i,\ell} \sum_{k=1}^K \beta_{k,\ell} B_k(t) + \epsilon_i(t)$$

- Estimar: $\beta_{\ell,k}$, $k = 1, \dots, K$ and $\ell = 1, \dots, m$, σ^2 e ϕ .

Basis smoothing:

- K é escolhido de maneira *ad-hoc*
- β pode ser facilmente estimado por mínimos quadrados, com solução explícita

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{W},$$

Smoothing splines:

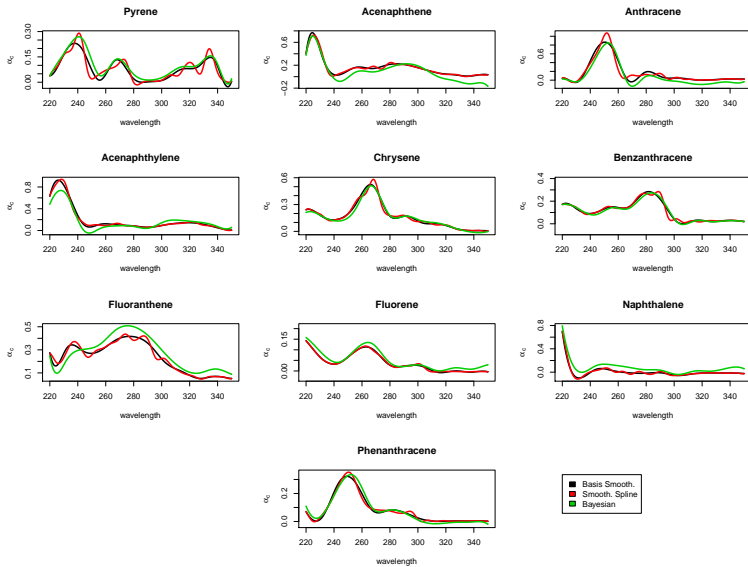
- Cada observação é um *knot*
- Número de coeficientes a serem estimados é maior que o número de observações.
- Suavização: penalizar a norma quadrada da segunda derivada das bases spline no problema de mínimos quadrados:



$$\hat{\beta} = (\mathbf{X}^+ \mathbf{X}^+ + \lambda \mathbf{I}_{m \times m} \otimes \mathbf{R})^{-1} \mathbf{X}^+ \mathbf{W}^+.$$

onde \mathbf{R} é uma matriz com entradas

$\mathbf{R}_{i,j} = \int_{-\infty}^{\infty} D^2 B_i(t) D^2 B_j(t) dt$, e D^2 é o operador diferencial de segunda ordem $\partial^2 / \partial t^2$.



- Abordagem Bayesiana segue bastante próxima de *Basis Smoothing*
- *Smoothing Spline* estimativas são menos suaves e tem mais “bumps” e captura mais variação local dos dados
- Brereton(2003): curvas são amostradas em um conjunto de comprimentos de onda mais esparsos e isto causa uma redução no ruído.
- Portanto, a maior parte da variação local nas curvas são características importantes da variação local dos dados e precisam ser identificadas.

Relembrando:

- Modelo de calibração

$$W_i(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell} \alpha_{\ell}(t_j) + \epsilon_i(t_j), \text{ for } j = 1, \dots, T$$

- Modelo de Predição:

$$W_i^*(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell}^* \alpha_{\ell}(t_j) + \epsilon_i^*(t_j), \text{ for } j = 1, \dots, T$$

Problema de calibração Dados $y_{i,\ell}$ e $W_i(t_j)$ estimar $\alpha_{\ell}(t_j)$.

Problema de predição Dado novo conjunto de dados $W_i^*(t_j)$, estimar $y_{i,\ell}^*$.

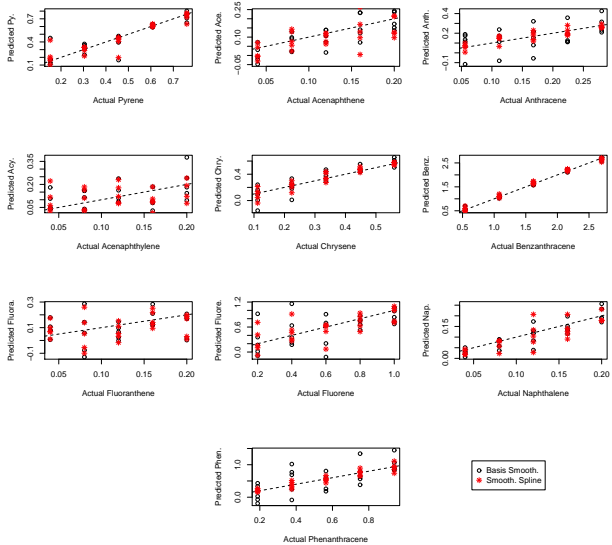
Abordagem simples:

- Use os espectros estimados $\hat{\alpha}_\ell$
- Coloque-os no modelo

$$W_i^*(t_j) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell}^* \alpha_\ell(t_j) + \epsilon_i(t_j)$$

e encontre o conjunto de \mathbf{y}_j^* que minimizam a soma de quadrados dos erros de predição

$$\hat{\mathbf{y}}_j^* = \arg \min_{\mathbf{y}_j^*} \sum_{r=1}^T \sum_{\ell=1}^m \left(W_j^*(t_r) - \hat{\mathbf{y}}_{j,\ell}^* \hat{\alpha}_\ell(t_r) \right)^2.$$



Estimação simultânea

- **Modelo de calibração:**

$$W_i(t) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell} \alpha_{\ell}(t) + \epsilon_i(t), \text{ for } i = 1, \dots, I$$

- **Modelo de predição**

$$W_i^*(t) = \alpha_j + \sum_{\ell=1}^K y_{i,\ell}^* \alpha_{\ell}(t) + \epsilon^*(t), \text{ for } i = 1, \dots, I^*$$

- Considere como dados:

- 1 $W_i(t_j)$ para $i = 1, \dots, I, j = 1, \dots, T$
- 2 $y_{i,\ell}$ para $i = 1, \dots, I$, e $\ell = 1, \dots, C$
- 3 $W_i^*(s_j)$ para $i = 1, \dots, I^*, j = 1, \dots, T^*$

- Considere como parâmetros

- 1 β tais que $\alpha_c(t) = \sum_{k=1}^K \beta_{c,k} B_k(t)$,
- 2 σ^2 e ϕ parâmetros da covariância
- 3 $y_{i,\ell}^*$ para $i = 1, \dots, I^*$, e $\ell = 1, \dots, C$

Erro padrão de predição – Standard Error of Prediction (SEP)

A contribuição do componente ℓ é dada por

$$\text{SEP}_\ell = \left(\frac{1}{J-1} \sum_{j=1}^J (y_{j,\ell}^* - \mathfrak{y}_{j,\ell}^*) \right)^{1/2}$$

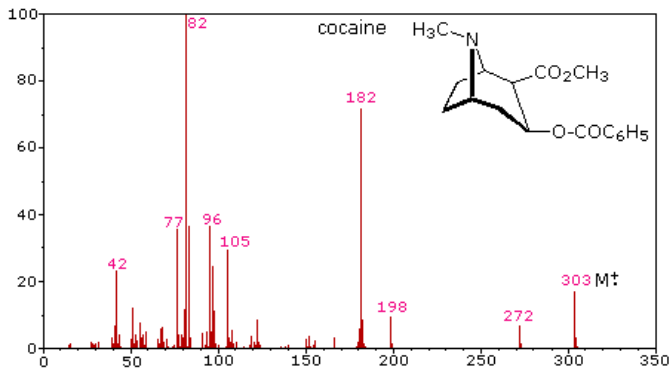
e o SEP total será dado por

$$\text{SEP} = \left(\frac{1}{mJ-1} \sum_{\ell=1}^m \sum_{j=1}^J (y_{j,\ell}^* - \mathfrak{y}_{j,\ell}^*) \right)^{1/2}$$

SEP (mg/l) para PAH.

Constituyente	MLR	PCR	PLS	OLS-K	OLS-SS	ML
Pyrene	0.09	0.09	0.09	0.09	0.09	0.06
Acenaphthene	0.06	0.03	0.03	0.05	0.06	0.03
Anthracene	0.03	0.03	0.03	0.11	0.04	0.01
Acenaphthylene	0.09	0.06	0.06	0.08	0.09	0.03
Chrysene	0.06	0.05	0.04	0.10	0.06	0.02
Benzantracene	0.07	0.07	0.06	0.06	0.07	0.04
Fluoranthene	0.09	0.08	0.08	0.10	0.09	0.04
Fluorene	0.24	0.20	0.15	0.32	0.23	0.07
Naphthalene	0.04	0.03	0.04	0.03	0.04	0.03
Phenanthracene	0.08	0.09	0.08	0.31	0.09	0.06
Total	0.11	0.09	0.08	0.16	0.10	0.05

CSI

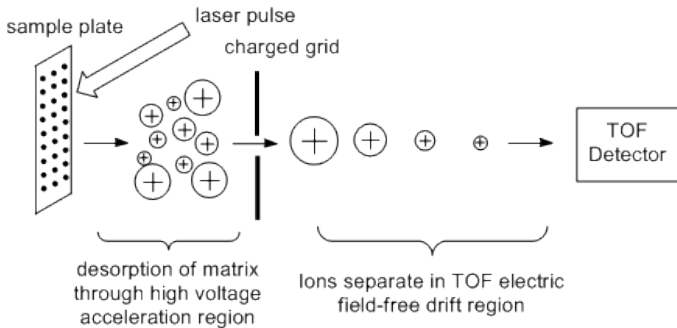


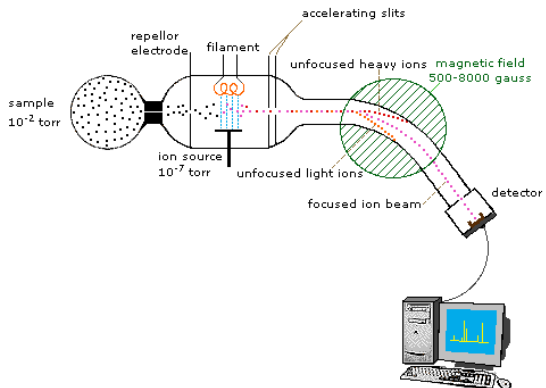
Espectrometria de massa

- A amostra, a qual pode ser sólida, líquida, ou gás, é ionizada, por exemplo, através do bombardeamento com electrons.
- Isso pode fazer com que algumas moléculas da amostra se quebrem em fragmentos carregados eletricamente.
- Estes ions são então separados de acordo com a sua razão massa/carga, tipicamente acelerando-as e submetendo-as a um campo elétrico ou magnético: ions de mesma razão massa/carga vão sofrer a mesma quantidade de deflexão.
- Os ions são detectados por um mecanismo capaz de detectar partículas carregadas, tais como um multiplicador de electrons.

Espectro de massa

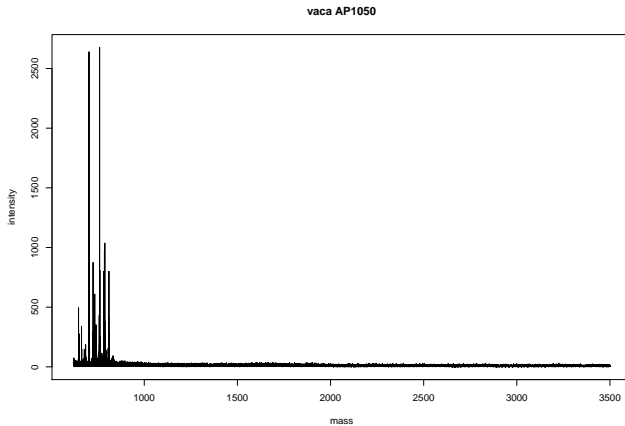
- Os resultados são apresentados como espectros da abundância relativa dos ions detectados como função da razão de massa/carga.
- Os átomos ou moléculas na amostra podem ser identificados por meio da correlação entre massas conhecidas para as massas identificadas ou por meio de um esquema de fragmentação característico.



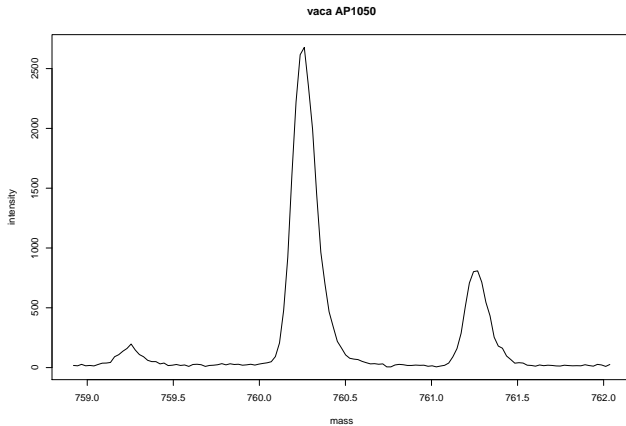


Marcador biológico para progesterona

- 17 espectros de massa para dois grupos de vacas: AP (alta progesterona) e BP (baixa progesterona).
- Para cada vaca temos duas amostras (medidas repetidas)
- Problema muito importante: insiminação artificial
- Prof. Mário Binelli (Faculdade de Medicina Veterinária e Zootecnia da USP de Pirassununga)



Zoom: 759 a 762 m/z



Dados de útero de vaca

- 7 (14) curvas para grupo AP $W_{i,AP}(t), i = 1, \dots, 7$
- 10 (20) curvas for grupo BP $W_{j,BP}(t), j = 1, \dots, 10$

O modelo: Assuma que existam conjuntos C_1 e C_2 tais que

$$W_{i,AP} = f(t) + \sum_{x \in C_1} R_x^{AP} K \left(\frac{t - x}{\sigma_x^{AP}} \right) + \epsilon_{i,AP}(t)$$

e

$$W_{j,BP} = f(t) + \sum_{y \in C_2} R_y^{BP} K \left(\frac{t - y}{\sigma_y^{BP}} \right) + \epsilon_{j,BP}(t)$$

para $t \in [619, 3500]$, onde $\epsilon_{i,AP}(\cdot)$ e $\epsilon_{j,BP}(\cdot)$ são processos Gaussianos independentes, $f(t)$ is é uma função comum (baseline), e $K(t)$ é uma função kernel (eg Gaussiana).

Os objetivos:

- Identificar os conjuntos C_1 e C_2
- Estimar os parâmetros f , R_x^{AP} e σ_x^{AP} , para $x \in C_1$ bem como R_y^{BP} e σ_y^{BP} , para $y \in C_2$.
- Finalmente identificar $C \subset C_1 \cup C_2$ tais que $R_x^{AP} \neq R_x^{BP}$ para $x \in C$

Muito obrigada!