

Projecto de Algoritmos e Modelação Computacional
Agrupamento (*clustering*) para modelos farmacocinéticos

2013/14

MEBiom

Conteúdo

1	Objectivo	5
2	Conceitos básicos	7
2.1	Modelos farmacocinéticos baseados em compartimentos	8
2.2	Algoritmo EM para misturas de Gaussianas	12
2.3	Aplicação a farmacocinética	16
2.3.1	Condições de paragem	21
3	Tipos de dados para a 1ª entrega	23
3.1	Amostra	24

4		
3.2	Grafos de compartimentos	25
3.3	Misturas de Gaussianas	26
4	2ª entrega	27

Capítulo 1

Objectivo

O objetivo do projeto é desenvolver um classificador não supervisionado para curvas que descrevem modelos farmacocinéticos. O classificador é aprendido a partir de dados públicos fornecidos pelo software *monolix* - <http://accp1.org/pharmacometrics/cssolutionmonolix.htm>. O objectivo é **encontrar classes de comportamento** do organismo na absorção de um certo fármaco numa dada população. Após encontrar estas classes, pode-se fazer aprendizagem supervisionada para classificar um dado paciente numa destas classes (a ser estudado nas aulas como aplicação do tipo de dados grafo).

Capítulo 2

Conceitos básicos

2.1 Modelos farmacocinéticos baseados em compartimentos

No projecto em questão será considerado apenas modelos com **um compartimento**. Este modelos servem para estimar a concentração de fármacos no organismo.

No modelo com um compartimento em questão considera-se que a droga é administrada oralmente, e que cada pessoa tem uma capacidade k_a de absorção do fármaco e uma capacidade k_e de eliminar o fármaco. Estas capacidades dão origem às seguintes equações diferenciais para a *consumo de fármaco* $I(t)$ e quantidade deste no organismo $Q(t)$ ao longo do tempo:

$$\begin{aligned} I'(t) &= -k_a I(t) \\ Q'(t) &= -k_e Q(t) + k_a I(t) \end{aligned} \tag{2.1.1}$$

restritas às seguintes condições iniciais

$$I(0) = \text{Dose} \times F \text{ e } Q(0) = 0$$

onde $F \in [0, 1]$ é uma constante que indica a taxa de fármaco que o organismo tem acesso, neste projecto consideramos $F = 0.5$.

As curvas dizem respeito à concentração $C(t)$ do fármaco no sangue (grandeza que pode ser medida na prática), sendo claro que

$$C(t) = \frac{Q(t)}{V}$$

onde V é o volume do organismo. Na prática, os fármacos são administrados assumindo que k_a e k_e são iguais em toda a população, sendo apenas distinguido o V . No entanto, verifica-se que existe grande variabilidade nestas constantes que justifica dar diferentes doses a indivíduos com o mesmo volume (medicina personalizada).

Da solução para a equação (2.1.1) obtemos:

$$C(t) = \frac{(\text{Dose} \times F)}{V} \frac{k_a}{k_a - k_e} (e^{-k_e t} - e^{-k_a t})$$

como na maior parte das vezes V é desconhecido, assume-se que a curva $C(t)$ é da forma

$$C(t) = a_1 e^{-b_1 t} + a_2 e^{-b_2 t}$$

pretendendo-se estimar os parâmetros a_i e $b_i = k_i$, assumindo ainda que $a_1 = -a_2 > 0$.

2.2 Algoritmo EM para misturas de Gaussianas

Como iremos ver, o algoritmo de agrupamento baseia-se no algoritmo de EM para misturas de Gaussianas (no caso em questão são unidimensionais) dado que vamos assumir que o erro na medição da grandeza de concentração tem uma distribuição normal em cada instante do tempo.

No caso de misturas de Gaussianas, os dados são um conjunto de pontos $\{x_i\}_{i=1,\dots,K}$ com $x_i \in \mathbb{R}$ i.i.d. de uma distribuição desconhecida, correspondente a uma mistura de M Gaussianas $p_\theta(x) = \sum_{j=1}^M w_j g_j(x)$ onde cada Gaussiana tem como parâmetros μ_j, σ_j^2 e os pesos w_j estão normalizados tal que $\sum_j w_j = 1$. O conjunto θ inclui todos estes parâmetros. M corresponderá ao número de classes a serem encontradas. Neste projecto consideramos que o M é fixo, no entanto, poderá ser relevante na prática considerar um M variável de acordo com os dados da população.

O objectivo é encontrar os parâmetros que maximizam a verosimilhança dos dados, isto é, encontrar θ tal que $\arg_{\theta} \max p_{\theta}(x_1 \dots x_K)$ onde

$$p_{\theta}(x_1 \dots x_K) = \prod_{k=1}^K p_{\theta}(x_k)$$

o que é equivalente a maximizar o logaritmo da verosimilhança

$$\log(p_{\theta}(x_1 \dots x_K)) = \sum_{k=1}^K \log(p_{\theta}(x_k))$$

A ideia do algoritmo de EM consiste em calcular uma sequência de parâmetros

$$\theta^0 \dots \theta^n$$

partindo de um conjunto de parâmetros inicial θ^0 de tal forma que

$$p_{\theta^m}(x_1 \dots x_K) < p_{\theta^{m+1}}(x_1 \dots x_K)$$

Para este fim, considera-se uma família de variáveis escondida $\{Y_k\}_{k=1,\dots,K}$ que toma valores no conjunto $\{1, \dots, M\}$ tal que

$$P(Y_k = j | X = x_k)$$

é a probabilidade de i -ésima amostra x_i ter sido amostrado de acordo com a j -ésima Gaussiana.

Considera-se de seguida o valor esperado de acordo com a distribuição $p_{\theta^m}(Y|X)$ para $\log(p_{\theta}(x_1 \dots x_K))$

$$\begin{aligned}\log(p_{\theta}(X)) &= E[\log(p_{\theta}(X))] \\ &= E[p_{\theta}(X, Y)] - E[\log(p_{\theta}(Y|X))]\end{aligned}$$

Denota-se por $Q(\theta, \theta^m)$ o valor $E[p_{\theta}(X, Y)]$.

Lema: Seja θ tal que $Q(\theta, \theta^m) > Q(\theta^m, \theta^m)$, então $\log(p_{\theta}(X)) > \log(p_{\theta^m}(X))$.

Após calcular $Q(\theta, \theta^m)$ (E-step) basta encontrar $\arg_{\theta} \max Q(\theta, \theta^m)$ (M-step). Tal é possível ser feito analiticamente para misturas (ver cálculos no quadro).

2.3 Aplicação a farmacocinética

No caso de dados farmacocinéticos temos que a amostra

$$y_i(t) = C(t) + \varepsilon$$

onde $\varepsilon \sim N(0, \sigma^2)$ e logo

$$p(y_1(t_1), \dots, y_1(t_n), \dots, y_K(t_n)) = \sum_{j=1}^M w_j \prod_{i=1}^K \prod_{\ell=1}^n g_j(y_i(t_\ell), t_\ell)$$

onde $g_j(y, t) \sim N(C_j(t), \sigma_j^2)$.

Objectivo: Encontrar $\theta = \{\theta_j\}_{j=1 \dots M}$ e $\theta_j = \{w_j, \sigma_j, a_{1j}, a_{2j}, b_{1j}, b_{2j}\}$ que maximizam a verosimilhança dos dados, restrito a que $a_{1j} = -a_{2j}$ e que $b_{1j} > b_{2j}$ para que $C(t)$ seja positiva.

Recorde que a função f que queremos estimar é

$$f(\theta_j, t) = \sum_{i=1}^2 a_{ij} e^{-b_{ij}t} \quad (2.3.2)$$

que representa a concentração da droga no organismo do componente j no instante t . Utilizando a técnica semelhante ao EM para misturas Gaussianas, temos que:

$$Q(\theta, \theta^{(k)}) = \sum_{j=1}^M \sum_{i=1}^K X_{ij} \log w_j^{(k)} p(y_i | \theta_j^{(k)}) \quad (2.3.3)$$

com

$$p(y_i | \theta_j) = \frac{1}{(2\pi\sigma_j^2)^{\frac{n}{2}}} e^{\frac{-1}{2\sigma_j^2} \sum_{\ell=1}^n (y_i(t_\ell) - f(\theta_j, t_\ell))^2} \quad (2.3.4)$$

e

$$X_{ij} = \frac{w_j p(y_i | \theta_j)}{\sum_{u=1}^M w_u p(y_i | \theta_u)} \quad (2.3.5)$$

Observe que tanto o denominador como o numerados são bastantes baixos, e no Java pode acontecer que ambos tomem o valor 0. Para que tal não aconteça (nos dados oferecidos) devem multiplicar o numerador e o denominador por e^{500} .

Seguindo o algoritmo de EM para Gaussianas, adaptando apenas o valor médio obtemos que a alteração dos w_i em cada iterada do algoritmo que maximiza $Q(\theta, \theta^{(k)})$ é feita de acordo com o seguinte equação:

$$w_j^{(k+1)} = \frac{1}{K} \sum_{i=1}^K X_{ij}^{(k)}. \quad (2.3.6)$$

Para encontrar o valor de σ_j^2 que maximiza $Q(\theta, \theta^{(k)})$ é necessário derivar $Q(\theta, \theta^{(k)})$ a σ^2 e encontrar um zero para o qual a segunda derivada seja negativa. Assim, temos que

$$\frac{\partial Q(\theta, \theta^{(k)})}{\partial \sigma_j^{2(k)}} = \sum_{i=1}^K X_{ij} \left(-\frac{K}{2\sigma_j^{2(k)}} + \frac{1}{2 \left(\sigma_j^{2(k)} \right)^2} \sum_{l=1}^n (y_i(t_l) - f(\theta_j^{(k)}, t_l))^2 \right)$$

Logo $\frac{\partial Q(\theta, \theta^{(k)})}{\partial \sigma_j^{2(k)}} = 0$ sse

$$\sigma_j^{2(k+1)} = \frac{\sum_{i=1}^K \sum_{l=1}^n X_{ij} (y_i(t_l) - f(\theta_j^{(k+1)}, t_l))^2}{\sum_{i=1}^K n X_{ij}} \quad (2.3.7)$$

Note que vai ter de actualizar todos os parâmetros de θ antes de actualizar σ e que $f(\theta_j^{(k+1)}, t_l)$ depende apenas de $a_j^{(k+1)}, b_{1j}^{(k+1)}, b_{2j}^{(k+1)}$. Para simplificar a notação vamos

utilizar y_{il} em vez de $y_i(t_l)$. É relativamente fácil verificar que a segunda derivada é negativa neste ponto.

Para ser possível actualizar o θ é necessário derivar $Q(\theta, \theta^{(k)})$ em ordem a $a_j = a_{1j} = -a_{2j}$ para encontrar o máximo e assim iterar o valor de a_j . Neste caso, $\frac{\partial Q(\theta, \theta^{(k)})}{\partial a_j^{(k)}} = 0$ sse

$$a_j^{(k+1)} = \frac{\sum_{i=1}^K \sum_{l=1}^n X_{ij} y_{il} (e^{-b_{1j}^{(k)} t_l} - e^{-b_{2j}^{(k)} t_l})}{\sum_{i=1}^K \sum_{l=1}^n X_{ij} \left(e^{-b_{1j}^{(k)} t_l} - e^{-b_{2j}^{(k)} t_l} \right)^2} \quad (2.3.8)$$

e mais uma vez, verifica-se que este ponto é um máximo. O caso mais complicado ocorre quando se tenta maximizar b_{1j} e b_{2j} . o que apenas se consegue fazer numericamente. Desta forma:

$$h_1(b_{1j}) = \frac{\partial Q(\theta, \theta^{(k)})}{\partial b_{1j}^{(k)}} = -\frac{a_j^{(k)}}{\sigma_j^{2(k)}} \sum_{i=1}^K \sum_{l=1}^n X_{ij} t_l e^{-b_{1j}^{(k)} t_l} (y_{il} - a_j^{(k)} (e^{-b_{1j}^{(k)} t_l} - e^{-b_{2j}^{(k)} t_l})).$$

Para obter um zero numericamente deve proceder ao método de Newton onde

$$b_{n+1} = b_n - \frac{h_1(b_n)}{h'_1(b_n)} \quad (2.3.9)$$

a segunda derivada de Q em ordem a b_{1j} é $h'_1(b_{1j})$ cuja expressão é:

$$h'_1(b_{1j}) = -\frac{a_j^{(k)}}{\sigma_j^{2(k)}} \sum_{i=1}^K \sum_{l=1}^n X_{ij} t_l^2 e^{-b_{1j}^{(k+1)} t_l} (a_j^{(k)} (2e^{-b_{1j}^{(k+1)} t_l} - e^{-b_{2j}^{(k)} t_l}) - y_{il})$$

De forma análoga, deve ser encontrado o máximo de b_{2j} onde

$$h_2(b_{2j}) = \frac{\partial Q(\theta, \theta^{(k)})}{\partial b_{2j}^{(k)}} = \frac{a_j^{(k)}}{\sigma_j^{2(k)}} \sum_{i=1}^K \sum_{l=1}^n X_{ij} t_l e^{-b_{2j}^{(k)} t_l} (y_{il} - a_j^{(k)} (e^{-b_{1j}^{(k)} t_l} - e^{-b_{2j}^{(k)} t_l}))$$

para a qual a segunda derivada é:

$$h'_2(b_{2j}) = \frac{a_j^{(k)}}{\sigma_j^{2(k)}} \sum_{i=1}^K \sum_{l=1}^n X_{ij} t_l^2 e^{-b_{2j}^{(k+1)} t_l} (a_j^{(k)} (e^{-b_{1j}^{(k)} t_l} - 2e^{-b_{2j}^{(k+1)} t_l}) - y_{il}).$$

2.3.1 Condições de paragem

Falta indicar quais as condições de paragem para o método Newton e para o EM em geral.

A convergência do método de Newton é feita da seguinte forma: Os valores das iteradas de b_1 e b_2 têm de verificar as seguintes condições:

- $0 < b_{1j}^{(k+1)} < b_{2j}^{(k)}$ – condição para $b_{1j}^{(k+1)}$;
- $b_{1j}^{(k)} < b_{2j}^{(k+1)} < 5$ – condição para $b_{2j}^{(k+1)}$ onde o valor 5 é uma heurística.

Se os valores das iteradas verificarem sempre a condição respectiva, então o método deve terminar quando $h_d(b_{dj}) = 0$ ou depois de 10000 (dez mil) iteradas.

Se por outro lado assim que uma iterada de $b_{dj}^{(k+1)}$ ficar fora do intervalo respectivo, neste caso existe um cotovelo (knee choice - apesar de “joelho” ser diferente de “cotovelo”) e é necessário proceder a uma escolha diferente do $b_{dj}^{(k+1)}$ (para os dados fornecido este problema apenas acontece para o b_{2j}). A escolha é feita da seguinte forma

1. Reinicia-se o método de Newton (no máximo com 10000 iteradas)
2. Termina-se o ciclo do método de Newton se o $h'_d(b_{dj}) > \beta = -0.3$ ou os limites b_{dj} sejam violados;
3. Caso os limites sejam violados decrementa-se $\beta = \beta - 0.2$ e voltamos para o passo 1.
4. Caso o β seja menor que -5 então

$$b_{dj}^{(k+1)} = b_{dj}^{(k)}.$$

Quanto à convergência do EM, este termina quando

$$(b_{dj}^{(k+1)} - b_{dj}^{(k)})^2 < 0.000001$$

para $d = 1$ e $d = 2$. Não se faz nenhuma consideração ao a .

Capítulo 3

Tipos de dados para a 1ª entrega

A entregar a 24 de Abril de 2014. As classes a serem utilizados neste projeto são os seguintes:

3.1 Amostra

- `add`: recebe um vector com três campos (índice, tempo e valor) e acrescenta o vector à amostra;
- `length`: retorna o comprimento da amostra;
- `element`: recebe uma posição e retorna o vector da amostra;
- `indice`: retorna uma lista de pares (tempos, valor) para um dado índice;
- `join`: recebe uma amostra e retorna uma nova amostra com as duas concatenadas;

3.2 Grafos de compartimentos

- `grafoo`: método construtor recebe um natural n e retorna o grafo com n nós e sem arestas. Cada nó representa um compartimento.
- `add_edge`: recebe dois nós e adiciona ao grafo uma aresta de um nó para outro com as estimativa inicial para uma mistura Gaussiana entre estes compartimentos.
- `remove_edge`: recebe dois nós e retira ao grafo uma aresta de um nó para outro.
- `up_edge`: recebe dois nós modifica a estimativa da mistura de gaussianas entre estes compartimentos.

3.3 Misturas de Gaussianas

- `mix`: Método construtor que recebe um inteiro (número de misturas), um conjunto de parâmetros θ ;
- `prob`: Recebe uma lista de pontos ao longo do tempo e retorna a probabilidade dessa lista de pontos ser observada pela mistura.
- `theta`: Retorna a lista de parâmetros actual;
- `update`: Método que recebe uma lista de parâmetros θ e actualiza a mesma;

Capítulo 4

2ª entrega

Na segunda entrega deverão ser implementadas a aplicação gráfica pela qual deve ser possível:

- Construir um grafo de compartimentos (número de compartimentos dado pelo utilizado bem como que arestas existem entre eles);
- Associar um ficheiro de dados *.csv e um ficheiro *.theta com a aproximação inicial de θ para uma aresta;
- Fazer a aprendizagem não supervisionada, tal como descrito na Secção 2, dos parâmetros;
- Apresentar quais foram os parâmetros aprendidos para a aresta por intermédio de um ficheiro.

Deverá ser elaborado um relatório em que conste:

- Documentar as opções tomadas no projecto bem como a justificação de alterações à 1a entrega.
- Pequeno manual de utilização.
- Experimentação dos dados fornecidos na página da disciplina onde se deve ilustrar as curvas aprendidas (usando e.g. Mathematica)

Cotação

São avaliadas as opções de velocidade, bem como a documentação do código.

- Tipos de dados (3 val)
 - Grafo (1 val)
 - Amostra (1 val)
 - Mistura de Gaussianas (1 val)
- Algoritmo de Aprendizagem (3 val)
- Input/output de dados/resultados (2 val)
- Aplicação Gráfica (1 val)
- Relatório (1 val)