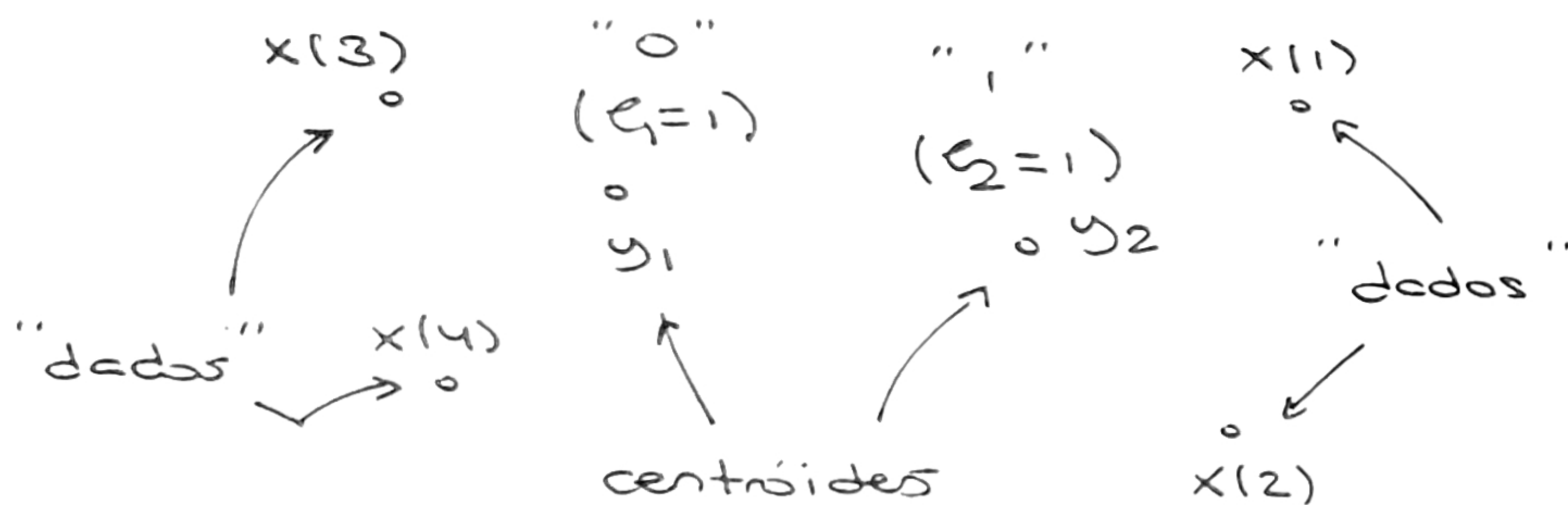


UFRRS — COPPE — PEE — CPE723 — Optimization Natural

Simulated Annealing — Aula 06 — Deterministic Annealing

Voltando à Aula 03, Exemplo 4:



$l_i$  ( $i=1,2$ ) ("length"): comprimento da "palavra binária" transmitida.

$$D = \frac{1}{2} \sum_{n=1}^2 d(x(n), y_{k(n)}) = D(Y)$$

$$R = \frac{1}{2} \sum_{n=1}^2 l_{k(n)} = \frac{1+1+\dots+1}{2} = 1$$

"data rate" (taxa de dados) (1 bit/vetor)

$$\text{Projeto: } Y^* = \underset{Y}{\operatorname{argmin}} D(Y)$$

Minimização de  $D(Y) = \frac{1}{N} \sum d(x(n), y_{k(n)})$  por método gradiente, ou seja, solução via GLA (Generalized Lloyd Algorithm\*) (ou LBG). Duas etapas.

Etapas: (Condição de Parada — "calcular partição, dado  $Y$  fixo")

$$D = \frac{1}{4} ( \|x(1) - y_{\bigcirc}\|^2 + \|x(2) - y_{\bigcirc}\|^2 + \|x(3) - y_{\bigcirc}\|^2 + \|x(4) - y_{\bigcirc}\|^2 )$$

$$\uparrow k(1) = \underset{i}{\operatorname{argmin}} \|x(1) - y_i\|$$

$$\forall n: k(n) = \underset{i}{\operatorname{argmin}} \|x(n) - y_i\|$$

(\*) A. Gersho e R. M. Gray. Vector Quantization and Signal Compression, Ed. Kluwer, 1992.

Y. Linde, A. Buzo e R. M. Gray. An algorithm for vector quantizer design.

IEEE Trans. Commun., vol. 28, pp. 84-95, 1980.

Etapa 2 (Condição do Centróide — "calcular  $\gamma$ , dada uma partição fixa")

$$D = \frac{1}{4} ( \|x(1) - y_2\|^2 + \|x(2) - y_2\|^2 + \|x(3) - y_1\|^2 + \|x(4) - y_1\|^2 )$$

$\frac{dD}{dy_k} = ? \rightarrow$  derivada de função escalar em relação a um vetor (slide seguinte)

$$\frac{d}{dy_2} (\|x(1) - y_2\|^2) = \frac{d}{dy_2} (x(1)^T x(1)) - \frac{d}{dy_2} (2x(1)^T y_2) + \frac{d}{dy_2} (y_2^T y_2) = -2x(1) + 2y_2$$

$$\frac{dD}{dy_2} = \frac{1}{4} (-2x(1) - 2y_2 - 2x(2) - 2y_2) = 0 \quad (\text{otimizando...})$$

$y_2 = \frac{x(1) + x(2)}{2}$

$\swarrow$   $x(1) \in \text{cluster 1}$   
 $\nwarrow$   $x(2) \in \text{cluster 1}$   
 $\nwarrow$   $2$   $\leftarrow$  nº de vetores  $x(n) \in \text{cluster 1}$

Em geral:

$$y_k = \frac{1}{N_k} \sum x(n) \text{ "E } y_k \text{"}$$

4  
Derivada de função escalar em relação a um vetor:

$$\frac{d(y^T y)}{dy} = \begin{bmatrix} \frac{d(y_1^2 + \dots + y_m^2)}{dy_1} \\ \vdots \\ \frac{d(y_1^2 + \dots + y_m^2)}{dy_m} \end{bmatrix} = \begin{bmatrix} 2y_1 \\ \vdots \\ 2y_m \end{bmatrix} = 2y$$

(vetor)  $\nearrow$

$$\frac{d(x^T y)}{dy} = \begin{bmatrix} \frac{d(x_1 y_1 + \dots + x_m y_m)}{dy_1} \\ \vdots \\ \frac{d(x_1 y_1 + \dots + x_m y_m)}{dy_m} \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = x$$

(vetor)  $\nearrow$

2

— Deterministic Annealing (DA) (Rose, 1993 e Rose, 1998)

- Minimização de funções não-convexas (habilidade de evitar mínimos locais)
- Aplicabilidade a diferentes tipos de problemas
- Número mínimo de restrições (máxima entropia)
- Analogia com termodinâmica estatística
- Ver: Rose, 1993 (Seções I e II)

Rose, 1998 (até o final da Seção II.A)

— Deterministic Annealing (DA) (Rose, 1993 e Rose, 1998)

Vamos estudar DA através de um exemplo ("soft clustering")

Matriz  $P_{y|x}$  :

	$x(1)$	$x(2)$	$x(3)$	$x(4)$
$y_1$	0.7	0.8	0.4	0.3
$y_2$	0.3	0.2	0.6	0.7

( $P_{y|x}(y|x)$ )

$$D = \sum_x p_x \sum_y P_{y|x} d_{xy} = D(P_{y|x}) = \frac{1}{4} \left( 0.7 \|x(1) - y_1\|^2 + 0.3 \|x(1) - y_2\|^2 \right. \\ \left. + 0.8 \|x(2) - y_1\|^2 + 0.2 \|x(2) - y_2\|^2 \right. \\ \left. + 0.4 \|x(3) - y_1\|^2 + 0.6 \|x(3) - y_2\|^2 \right. \\ \left. + 0.3 \|x(4) - y_1\|^2 + 0.7 \|x(4) - y_2\|^2 \right)$$

( $D = \sum_x \sum_y p_{xy} d_{xy}$ )



Observação: uma forma alternativa de codificar o estado ("Y") para otimização em problemas de "clustering": ao invés de usar o dicionário Y como estado, podemos usar a partição ("hard" ~~ou~~ "soft") como estado:

$\rightarrow P_{Y|X} =$   
 (hard)  
 $k=8$   
 centróides

$$\begin{matrix}
 & x(1) & x(2) & \dots & x(800) & \leftarrow \text{dados} \\
 y_1 & \begin{bmatrix} 1 & 0 & & 1 \end{bmatrix} \\
 y_2 & \begin{bmatrix} 0 & 1 & & 0 \end{bmatrix} \\
 \vdots & \vdots & \vdots & & \vdots \\
 y_8 & \begin{bmatrix} 0 & 0 & & 0 \end{bmatrix} \\
 & \downarrow & \downarrow & & \downarrow \\
 & \Sigma=1 & \Sigma=1 & & \Sigma=1
 \end{matrix}$$

ou seja, estado = [ 1 2 7 ... 5 8 1 ]

agora o estado fica discreto;  
 e uma perturbação possível  
 seria sortear uma posição do  
 estado e perturbá-la para  
 outro inteiro de 1 a 8. O resto  
 do SA permanece igual. Y passa  
 a ser consequência do estado.

Entropia de uma variável aleatória  $X$ :

$$H(X) = - \sum_{k=1}^K p_k \log p_k \quad (\log_2 \longrightarrow \text{em bits}; \ln \longrightarrow \text{em nats})$$

Variável aleatória uniforme:

$x$	$p(X=x)$	código
0	0.25	00
1	0.25	01
2	0.25	11
3	0.25	10

$$H = -0.25 \log_2 0.25 - 0.25 \log_2 0.25 \\ - 0.25 \log_2 0.25 - 0.25 \log_2 0.25$$

$$H = 2 \text{ bits} \quad (\log_2 0.25 = -2)$$

Outra distribuição de probabilidades:

$x$	$p(X=x)$	código
0	0.5	0
1	0.25	10
2	0.125	110
3	0.125	111

$$H = -0.5 \log_2 0.5 - 0.25 \log_2 0.25 \xrightarrow{(-2)} \\ - 0.125 \log_2 0.125 - 0.125 \log_2 0.125 \xrightarrow{(-3)}$$

$$H = 0.5 + 0.5 + 0.375 + 0.375 = 1.75 \text{ bits}$$



Entropia conjunta de variáveis aleatórias  $X$  e  $Y$ :

$$H(X, Y) = - \sum_x \sum_y p(x, y) \log p(x, y)$$

$$(H(X, Y) = - \sum_x \sum_y p_{xy} \log p_{xy}) \quad (\text{forma "abreviada"})$$

$$H(X, Y) = - \sum_x \sum_y p(y|x) p(x) \underbrace{\log p(y|x) p(x)}_{(\log p(x) + \log p(y|x))}$$

$$H(X, Y) = - \sum_x p(x) \log p(x) \underbrace{\sum_y p(y|x)}_1 - \sum_x p(x) \sum_y p(y|x) \log p(y|x)$$

$$H(X, Y) = H(X) + \boxed{H(Y|X)}$$

$H(Y|X)$

$$\left( H(X, Y) = - \sum_x p_x \log p_x \overbrace{\sum_y p_{y|x}}^1 - \overbrace{\sum_x p_x \sum_y p_{y|x} \log p_{y|x}}^{H(Y|X)} \right)$$

Maximização de  $H(y|x)$ :

(portador com "máxima entropia")

	$x(1)$	$x(2)$	$x(3)$	$x(4)$
$y_1$	0.5	0.5	0.5	0.5
$y_2$	0.5	0.5	0.5	0.5

Maximização de  $H(y|x)$  é minimização de  $D$ : (↖  $p_{y|x}$ )

$$\boxed{J = D - TH}$$

$$J = \sum_x p_x \sum_y p_{y|x} d_{xy} + T \sum_x p_x \sum_y p_{y|x} \log p_{y|x} +$$

parcela para garantir  
que  $\sum_y p_{y|x} = 1$ , caso contrário  
 $p_{y|x} = 0$  leva a  $J \rightarrow -\infty$   
(solução inválida).