Métodos de Agrupamento Para

Dados Simbólicos: Hierarquico Divisivo

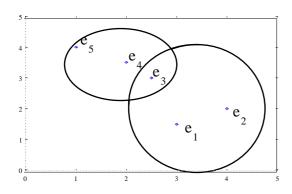




Estruturas classificatórias



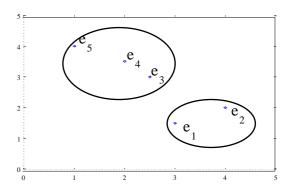
Cobertura



1)
$$\forall \ell = 1, \dots, K \text{ tem - se } P_{\ell} \neq \emptyset$$

$$2)\bigcup_{l=1}^{K}P_{l}=\Omega$$

Partição



3)
$$\forall \ell, m = 1, \dots, K \ e \ l \neq m$$

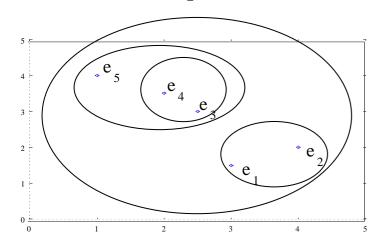
$$ent\tilde{a}o \ P_l \cap P_m = \emptyset$$



Estruturas Classificatórias



Hierarquia

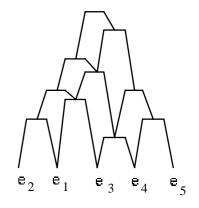


$$1)\Omega \in H$$

- 2) $\forall e \in \Omega \text{ então } \{e\} \in H$
- 3) $\forall h, h' \in H \text{ tem se}$:

$$h \cap h' \neq \emptyset \Rightarrow h \subset h' \text{ ou } h' \subset h$$

Piramide



3) $\forall h, h' \in H \text{ tem - se } h \cap h' = \emptyset \text{ ou } h \cap h' \in H$

4)Existe uma ordem θ tal que

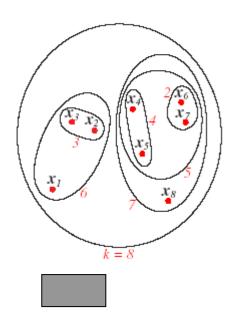
 $\forall h \in H$, $h \in \text{um intervalode } \theta$

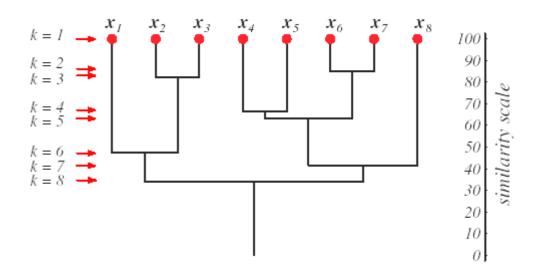
Classificação Hierarquica



Diagrama de Venn sobre os dados bidimensionais

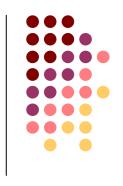
Dendrograma: beaseado em um índice de heterogeneidade





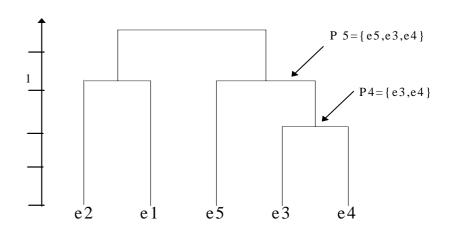


Métodos Hierárquicos



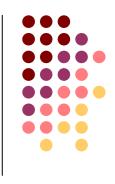
Parte-se de uma tabela de dados e calcula-se uma distância entre os individuos de Ω

Os métodos ascendentes hierárquicos tem por objetivo a construção de uma sequencia de partições encaixadas chamada hierarquia. A representação gráfica dessas hierarquias é realisada por uma arvore hierarquica ou dendrograma.

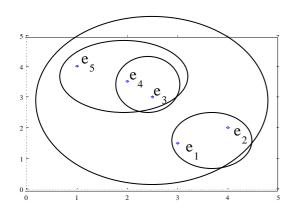




Hierarquia com índice



Hierarquia H



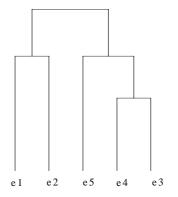
$$1)E \in H$$

2)
$$\forall e \in E \text{ então } \{e\} \in H$$

3) $\forall h, h' \in H \text{ tem - se}$:

$$h \cap h' \neq \emptyset \Rightarrow h \subset h' \text{ ou } h' \subset h$$

Hierarquia com indice (H,f)



$$f: H \to \mathfrak{R}^+$$

- (1) f(h) = 0 se e somente se card(h) = 1
- (2) $\forall h, h' \in H \ h \subset h' \in h \neq h' \Longrightarrow f(h) < f(h')$



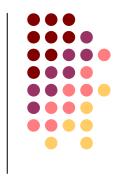
Método Divisivo

Marie Chavent et al (2000)

☐ Um método divisivo consiste em obter grupos de
objetos em diferentes etapas produzindo uma hierarquia.
□ O processo inicia com todos os objetos e divide
sucessivamente cada grupo formando grupos menores
de objetos.
□Uma divisão é definida otimizando um critério de
minimização de variância intra-grupo
☐ Em cada etapa o agrupamento é obtido de forma
monotética.



Notações



 Ω = conjunto de todos os objetos (dados simbólicos)

{Y₁,...,Y_p} – conjunto de variáveis simbólicas (multi-valor ordinal , intervalo ou modal)

 $C = \{C_1, C_2\}$ uma bipartição associada a uma variável Y_i

n – número total de objetos

n_i – número total de objetos da classe i, i=1,2



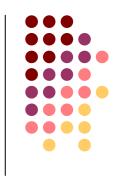




k	peso	altura
1	[30,40]	0.10 (alta) 0.30 (média) 0.60 (baixa)
2	[35,65]	0.20 (alta) 0.50 (média) 0.30 (baixa)
3	[20,30]	0.10 (alta) 0.80 (média) 0.10 (baixa)



Algoritmo



- Uma divisão é obtida por otimizando um critério de variância intra-classe (Inércia).
- O critério é minimizado com respeito a todas as bipartições induzidas por um conjunto de questões binárias de uma variável.
- Seleciona-se a bipartição onde a inércia inter-classe é máxima.



Algoritmo

Passo 1: $P_1 = Ω$; v = 1; $C = P_1$



Passo 2: Faça enquanto v ≤ V-1 (V é o nº de iterações)

<u>Passo 3</u>: Para cada C de P_v escolher $\{C_1, C_2\}$ induzida por questões binárias onde a inércia intra-grupo é mínima

$$W = I(C_1) + I(C_2) = \frac{1}{n \times n_1} \sum_{k,m \in C_1} \sum_{m > k} d_{km}^2 + \frac{1}{n \times n_2} \sum_{k,m \in C_2} \sum_{m > k} d_{km}^2$$

Passo 4: Escolher C de P_v tal que a inércia inter-grupo é máxima

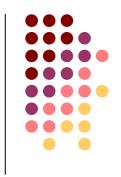
$$\Delta(\mathbf{C}) = \mathbf{I}(\mathbf{C}) - \mathbf{I}(\mathbf{C}_1) - \mathbf{I}(\mathbf{C}_2)$$

<u>Passo 5</u>: $P_{v+1} = P_v \cup \{C_1, C_2\} - \{C\}$

Passo 6: v = v + 1, vá para o Passo 2



Saída do algoritmo



A saída do método é uma hierarquia *H* constituída de V+1 classes

Cada classe $C_i \in H$ é organizada por $\Delta(C_i)$ no dendograma.

$$C_i \subset C_{i'} \Rightarrow \Delta(C_i) \leq \Delta(C_{i'})$$



Extensão do Critério de Variância Intra-grupo



- O critério usado para avaliar a qualidade de uma partição é uma extensão do critério soma de quadrados intra-grupo para o caso de uma matriz de distâncias.
- A extensão da soma de quadrados de erro (inércia) de uma classe

$$I(C_i) = \sum_{k \in C_i} / |x_i - \mu_i|^2$$

para uma matriz de distância D é dada por

$$I(C_i) = \frac{1}{n \times n_i} \sum_{k,m \in C_1} \sum_{m > k} d_{km}^2$$

onde n_i é o número de elementos do grupo i

Uma generalização do critério de uma partição com m classes é dado por

$$W = \sum_{i=1}^{m} I(C_i)$$
 Centro de Informática $V \in \mathbb{R}^{n}$





- O objetivo é encontrar a bipartição C₁= (C₁¹,C₁²) de menor soma de quadrados de erro.
- Uma escolha ótima é realizada avaliando $2^{n_i-1}-1$ bipartições.
- Uma solução é escolher a bipartição entre as bipartições induzidas pelo conjunto de todas as possíveis questões binárias.

Questões Binárias – [$Y_j \le c$]

$$q_c: \Omega \to \{V,F\}$$

A bipartição $\{C_1,C_2\}$ induzida por $q_c = [Y_i \le c]$ é

$$C_1 = \{k \in C / q_c = V\}$$

 $C_2 = \{k \in C / q_c = F\}$



Se Y_j é uma variável intervalo, então $Y_j(k) = [\alpha, \beta]$ e

$$q_c(k) = V \text{ se } m_k \le c$$

$$q_c(k) = F \text{ se } m_k > c$$

$$\alpha + \beta$$

onde $m_k = \frac{\alpha + \beta}{2}$

Existirá n_i -1 questões binárias onde c é o valor entre dois consecutivos valores de m_k

Exemplo : n = 3

$$Y_1 = peso: m_1 = 25 m_2 = 35 e m_3 = 50$$

$$c = 30$$

Para c = 30
$$q_c(1) = V$$
, $q_c(2) = F e q_c(3) = F$

$$C_1 = \{1\} e C_2 = \{2,3\}$$



Se Y_j é uma variável ordinal ou modal então $Y_j(k) = \pi_k$ e



$$q_c(k) = V$$
 se $\sum_{x \le c} \pi_k(x) \ge 1/2$

$$q_c(k) = F$$
 se $\sum_{x \le c} \pi_k(x) < 1/2$

Existirá $|O_j|$ - 1 questões binárias

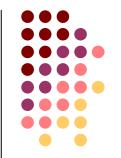
Exemplo: c=baixa para a variável modal da tabela

Para c = baixa
$$q_c(1) = V$$
, $q_c(2) = F e q_c(3) = F$

$$C_1 = \{1\} e C_2 = \{2,3\}$$



Medida de Distância



 ${
m I\hspace{-.1em}I}$ Se ${
m Y}_{
m j}$ é uma variável intervalo, então ${
m Y}_{
m j}$ (k) =[lpha,eta] e

$$d(k,m) = \left[\sum_{j=1}^{p} (\max\{ \left| \alpha_{k}^{j} - \alpha_{m}^{j} \right|, \left| \beta_{k}^{j} - \beta_{m}^{j} \right| \})^{2} \right]^{1/2}$$

Versões normalizadas

$$d(k,m) = \left[\sum_{j=1}^{p} \left[\frac{\max\{\left|\alpha_{k}^{j} - \alpha_{m}^{j}\right|, \left|\beta_{k}^{j} - \beta_{m}^{j}\right|\}}{m_{j}}\right]^{2}\right]^{1/2}$$

a)
$$m_j^2 = \frac{1}{2n^2} \sum_{k=1}^n \sum_{m=1}^n (\max(|\alpha_k^j - \alpha_m^j|, |\beta_k^j - \beta_m^j|)^2$$

b)
$$m_j = |O_j|$$



Exemplo

$$a_1 = [Y_1 \in [30.40]]$$
 $a_2 = [Y_1 \in [35.65]]$

$$d(1,2) = \left[(\max\{ |30 - 35|, |40 - 65| \})^2 \right]^{1/2} = 25$$

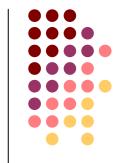
Versão normalizada – $O_1 = [20,70]$

$$d(1,2)=0.5$$



Se Y_i é uma variável modal então $Y_i(k) = \pi_k$

p=número de variáveis



	1	• • •	T	Σ
1	f ₁₁		f _{1T}	p
•	:		:	•
n	$\mathbf{f}_{\mathbf{n}1}$		$\mathbf{f}_{\mathbf{nT}}$	p
Σ	F .1		F .T	np

-				
	1	• • •	T	Σ
1	f ₁₁ /np	•	f _{1T} /np	p _{1.}
:			•	•
n	f _{n1} /np		f _{nT} /np	p _{n.}
Σ	p .1		p .T	1

$$\mathbf{T} = |\mathbf{O}_1| + ... + |\mathbf{O}_p|$$

$$d^{2}(k,m) = \sum_{t=1}^{T} \frac{p..}{p.t} \left[\frac{p_{kt}}{p_{k}} - \frac{p_{mt}}{p_{m}} \right]^{2} \qquad \qquad d^{2}(k,m) = n \sum_{t=1}^{T} \frac{\left[f_{kt} - f_{mt} \right]^{2}}{f_{kt} + f_{mt}}$$

$$d^{2}(k,m) = n \sum_{t=1}^{T} \frac{[f_{kt} - f_{mt}]^{2}}{f_{kt} + f_{mt}}$$

$$p_{kt} = \frac{f_{kt}}{np}$$

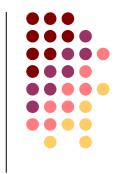
$$p_{k.} = \sum_{i=1}^{T} p_{k}$$

$$p_{.t} = \sum_{k=1}^{n} p_{kt}$$

$$p_{kt} = \frac{f_{kt}}{np}$$
 $p_{k.} = \sum_{t=1}^{T} p_{kt}$ $p_{.t} = \sum_{k=1}^{n} p_{kt}$ $p_{..} = \sum_{k=1}^{n} \sum_{t=1}^{T} p_{kt} = 1$



Exemplo



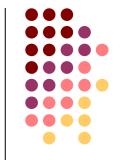
	baixa	média	alta	Σ
1	0.10	0.30	0.60	1
2	0.20	0.50	0.30	1
	0.30	0.80	0.90	2

$$d^{2}(1,2) = 2 \times \left[\frac{(0.10 - 0.20)^{2}}{0.10 + 0.20} + \frac{(0.30 - 0.50)^{2}}{0.30 + 0.50} + \frac{(0.60 - 0.30)^{2}}{0.60 + 0.30} \right] = 0.37$$

$$d(1,2)=0.60$$







k Specific Gravity		Frezzing point	lodine Value	Saponification	
1	[0.930,0.935]	[-27,-18]	[170,204]	[118,196]	
2	[0.930,0.937]	[-5,-4]	[192,208]	[188,197]	
	[0.916,0.918]	[-6,-1]	[99,113]	[189,198]	
4	[0.920,0.926]	[-6,-4]	[104.116]	[187.193]	
	[0.916,0.917]	[-21,-15]	[80,82]	[189,193]	
5 6	[0.914,0.919]	[0,6]	[79,90]	[187,196] [190,199]	
7	[0.860,0.870]	[30,38]	[40,48]	[190,202]	
8	[0.858,0.864]	[22,32]	[53,77]		

1 – linseed oil 5 – camelia oil

2 – perilla oil 6 – olive oil

3 – cottonseed oil 7 – beef tallow

4 – sesam oil 8 – hog fat



Matriz de Dissimilaridade



	1	2	3	4	5	6	7	8
1	-	0.81	1.17	1.08	1.30	1.34	2.62	2.34
2	-	-	1.19	1	0.70	1.25	3.31	2.00
3	-	-	-	0.02	0.09	0.03	0.55	0.98
4	-	-	-	-	0.11	0.06	0.90	1.14
5	-	-	-	-	-	0.11	0.77	0.62
6	-	-	-	-	-	-	0.86	0.79
7	-	-	-	-	-	-	-	0.05
8	-	-	-	-	-	-	-	-

1 – linseed oil

5 - camelia oil

2 – perilla oil

6 – olive oil

3 - cottonseed oil

7 – beef tallow

4 – sesam oil

8 – hog fat



Passo 1



Calcular $W = I(C_1)+I(C_2)$ para todas as questões binárias e escolher a bipartição que minimize W

a) Y = specific gravity

$$m_1 = 0.9325$$
 $m_2 = 0.9335$
 $m_3 = 0.9170$
 $m_4 = 0.9230$
 $m_5 = 0.9165$
 $m_6 = 0.9165$
 $m_7 = 0.8650$
 $m_8 = 0.8610$
 0.86300
 0.89075
 0.91675
 0.92000
 0.92775
 0.93330



Para c = 0.86300

$$q_c(1) = F$$
 C_1 : $Y_1 \le 0.86300$ C_2 : $Y_2 > 0.86300$

$$q_c(2) = F$$
 $q_c(3) = F$
 $q_c(4) = F$
 $q_c(4) = F$
 $C_1 = \{8\} \ e \ C_2 = \{1,2,3,4,5,6,7\}$
 $n_1 = 1 \ n_2 = 7$

$$q_c(4) = F$$
 $n_1 = 1 n_2 = q_c(5) = F$

$$\mathbf{q_c(7) = F}$$

$$\mathbf{q_c(8) = V}$$

 $q_c(6) = F$

$$I\left(C_{1}\right)=0$$

$$I(C_2) = \frac{1}{n \times n_2} \sum_{k,m \in C_2} \sum_{k>m} d^2(k,m) = \frac{1}{8 \times 7} [0.81 + 1.17 + 1.08 + \dots + 0.77 + 0.86] = 0.344$$

$$W = 0.344$$



Para c = 0.89075

$$\mathbf{q_c}(1) = \mathbf{F}$$

$$q_c(2) = F$$

$$q_c(3) = F$$

$$q_c(4) = F$$

$$q_c(5) = F$$

$$q_c(6) = F$$

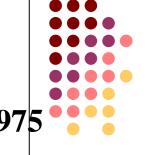
$$q_c(7) = V$$

$$q_c(8) = V$$

$$I(C_1) = \frac{0.05}{8 \times 2} = 0.003$$

$$I(C_2) = \frac{1}{6 \times 8} 10.86 = 0.21$$

$$\mathbf{W} = \mathbf{0.213}$$

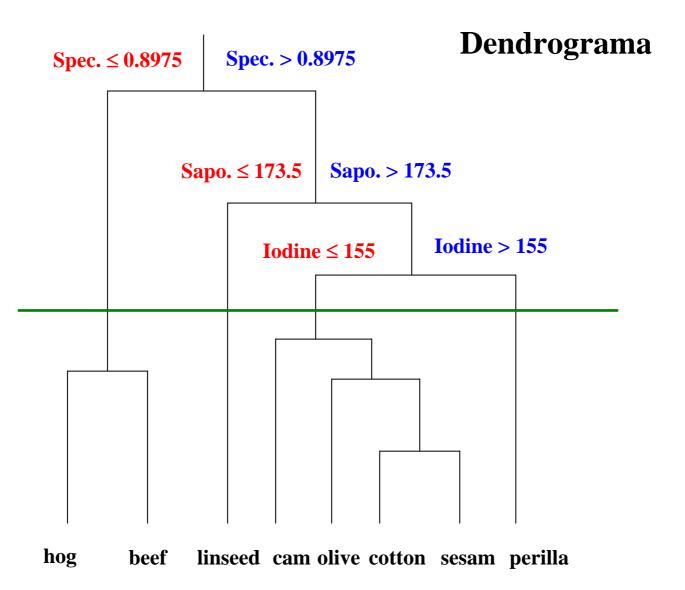


$$C_1$$
: $Y_1 \le 0.8975$ C_2 : $Y_2 > 0.8975$

$$C_1 = \{7,8\} \ \ e \ C_2 = \{1,2,3,4,5,6\}$$

$$n_1 = 2 n_2 = 6$$





Classe3

Classe 2

Classe 1



Classe 4

Saída do algoritmo



 $C_1 = \{ beef tallow, hog fat \}$

 $C_2 = \{linseed\}$

 $C_3 = \{camelia,olive,cottonseed,sesam\}$

 $C_4 = \{perilla\}$

Descrições

 $C_1 = [Spec. \le 0.8975]$

 $C_2 = [Spec. > 0.8975] \land [Sapo. \le 173.5]$

 $C_3 = [Spec. > 0.8975] \land [Sapo. \le 173.5] \land [Iod. \le 155]$

 $C_4 = [Spec. > 0.8975] \land [Sapo. \le 173.5] \land [Iod. > 155]$

