Agrupamento por similaridade - Processos clássicos

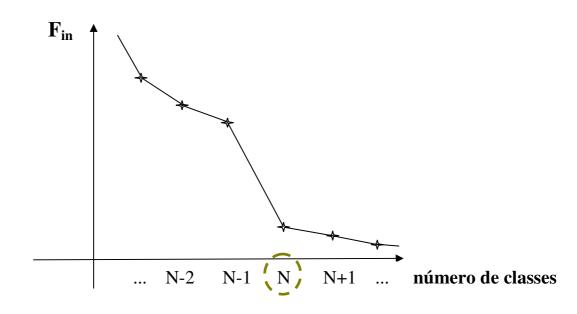
Agrupar, "de forma natural", conjuntos de dados com "similaridade interna"

determinar o número de classes à usar e

alocar os elementos às classes para minimizar F_{in}

1 - Determinação do número de classes a utilizar

Heurística: um agrupamento "natural" (um número N "bom" de classes) é alcançado na vizinhança imediata de uma grande variação em F_{in}.



2 - Alocação dos elementos às classes para minimizar a dispersão total intra classes $F_{\rm in}$

função à minimizar

$$F_{in} = \sum_{\forall C_j} F_j \qquad \text{onde} \qquad \qquad F_j = n_j \sigma_j^2 = \sum_{\forall \vec{x}_i \in C_j} \left\| \vec{x}_i - \vec{m}_j \right\|^2$$

F_j é a dispersão intra classe da classe Cj

2,1 Localização dos padrões das classes

O padrão $\underline{p}_j\,$ da classe C_j deve minimizar a dispersão intra classe da classe C_i

$$F_{j} = \sum_{\forall \vec{x} \in C_{j}} |\vec{x} - \vec{p}_{j}|^{2} = \sum_{\forall \vec{x} \in C_{j}} \sum_{\forall k} (x_{k} - p_{jk})^{2} = \sum_{\forall k} \sum_{\forall \vec{x} \in C_{j}} (x_{k} - p_{jk})^{2}$$

Como visto anteriormente

o padrão que minimiza a dispersão intra classe (ou erro de representação) de uma classe é o seu baricentro .

$$p_{ik} = m_{ik} \implies \vec{p}_i = \vec{m}_i$$

2.2 - Processos de Agrupamento

Como realizar a

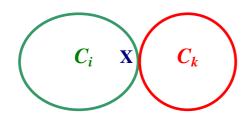
Minimização interativa de F_{in}

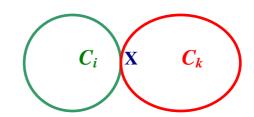
(e simultaneamente a maximização de F_{out})?

$$F_{in} = \sum_{j} F_{j}$$

Rearranjando os elementos entre as classes para reduzir F_{in}

Deslocando um elemento $\hat{\vec{x}}$ da classe C_i para a classe C_k





Após alguma álgebra:



$$\vec{m}_{k} = \frac{1}{n_{k}} \sum_{\forall \vec{x} \in C_{k}} \vec{x} \qquad \Rightarrow \quad \vec{m}'_{k} = \frac{1}{n_{k}} \sum_{\forall \vec{x} \in C_{k}'} \vec{x} = \frac{1}{n_{k} + 1} (\hat{\vec{x}} + n_{k} \vec{m}_{k}) = \vec{m}_{k} + \frac{\hat{\vec{x}} - \vec{m}_{k}}{n_{k} + 1}$$

$$F_{k} = \sum_{\forall \vec{x} \in C_{k}} |\vec{x} - \vec{m}_{k}|^{2} \qquad \Rightarrow \quad F'_{k} = \sum_{\forall \vec{x} \in C'_{k}} |\vec{x} - \vec{m}'_{k}|^{2} = \dots = F_{k} + \frac{n_{k}}{n_{k} + 1} |\hat{\vec{x}} - \vec{m}_{k}|^{2}$$

$$\vec{m}_{i} = \frac{1}{n_{i}} \sum_{\forall \vec{x} \in C_{i}} \vec{x} \qquad \Rightarrow \quad \vec{m}'_{i} = \frac{1}{n_{i}'} \sum_{\forall \vec{x} \in C'_{i}} \vec{x} = \frac{1}{n_{i} - 1} (-\hat{\vec{x}} + n_{i} \vec{m}_{i}) = \vec{m}_{i} - \frac{\hat{\vec{x}} - \vec{m}_{i}}{n_{i} - 1}$$

$$F_{i} = \sum_{\forall \vec{x} \in C_{i}} |\vec{x} - \vec{m}_{i}|^{2} \qquad \Rightarrow \quad F'_{i} = \sum_{\forall \vec{x} \in C'_{i}} |\vec{x} - \vec{m}'_{i}|^{2} = \dots = F_{i} - \frac{n_{i}}{n_{i} - 1} |\hat{\vec{x}} - \vec{m}_{i}|^{2}$$

$$\Delta F_{in} = (F'_k - F_k) + (F'_i - F_i) = \frac{n_k}{n_k + 1} |\hat{\vec{x}} - \vec{m}_k|^2 - \frac{n_i}{n_i - 1} |\hat{\vec{x}} - \vec{m}_i|^2$$

 \mathbf{F}_{in} é reduzido quando deslocamos um elemento $\hat{\vec{x}}$ da classe C_i para a classe C_k se $\Delta F_{in} < 0$, i.e. se

$$\left| \frac{n_i}{n_i - 1} \right| \hat{\vec{x}} - \vec{m}_i \right| > \frac{n_k}{n_k + 1} \left| \hat{\vec{x}} - \vec{m}_k \right|^2$$

como $\frac{n_i}{n_i-1} > 1 > \frac{n_k}{n_k+1}$ **podemos** simplificar a condição, basta que

$$\left| \hat{\vec{x}} - \vec{m}_i \right| > \left| \hat{\vec{x}} - \vec{m}_k \right|^2$$

i.e. que a distância do elemento ao baricentro da nova classe seja menor que ao da antiga.

2.3 Um processo de clusterização interativo simples e eficaz: K-means, C-means, Isodata ou Lloyd

entradas: \vec{x}_{i} i = 1, 2, ..., N

centros de classe: \vec{m}_j j = 1, 2, ..., K

Inicialização: Arbitre K centros \vec{m}_i j = 1, 2, ..., K

sugestão: use as K primeiras entradas \vec{x}_i i=1,2,...,K não muito próximas entre si (à discutir mais tarde)

Loop

1 - Reclassificação das entradas

$$orall \ \vec{x}_i \quad i=1,2,...,N$$
 Se \vec{m}_j é a melhor representação de \vec{x}_i , i.e. se $\left|\vec{x}_i-\vec{m}_j\right|<\left|\vec{x}_i-\vec{m}_k\right| \ \forall\, k\neq j \quad {\bf então}$ $\vec{x}_i\in C_j$

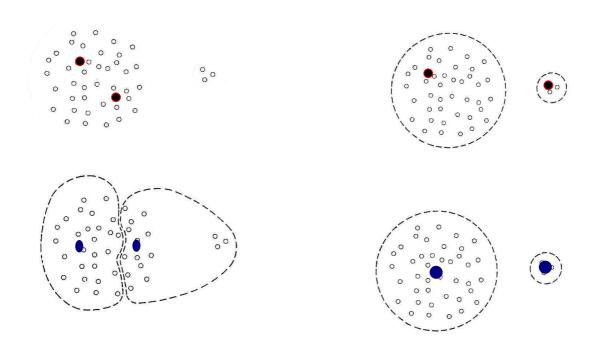
2 - Atualização dos centros de classe

$$\forall \quad j = 1, 2, ..., K$$

$$\vec{m}_j = \frac{1}{n_j} \sum_{\forall \vec{x} \in C_j} \vec{x}$$

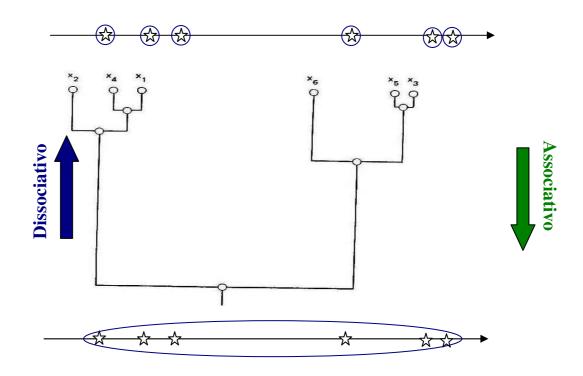
3 - Retorne ao passo 1

Quando nenhum elemento trocar de classe: Fim.



inicializar com entradas não muito próximas

2.4 - Algorítmos Hierárquicos: dendogramas, árvores

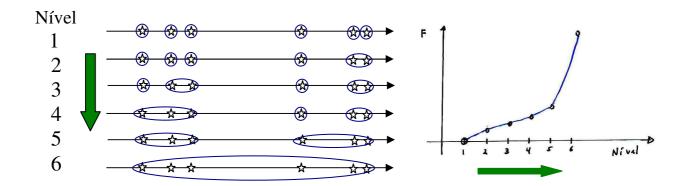


2.4.1 Algorítmo Aglomerativo, Associativo:

n elementos c classes

$$c = n \qquad n-1 \qquad n-2 \quad ... \quad 2 \qquad 1$$

$$F_{in} = F_{n} = 0 \quad < \quad F_{n-1} \quad < \quad F_{n-2} \quad < \quad F_{2} \quad < \quad F_{1}$$

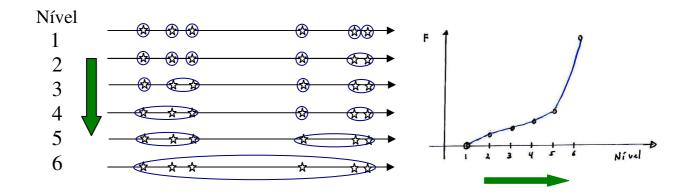


Algoritmo Hierárquico Associativo:

Inicialização: n elementos ⇒ n clusters

Loop:

Junte os dois clusters que aumentam F_{in} o mínimo possível Repetir até chegar ao número de clusters desejado ou F_{in} crescer muito



Obs:

se \underline{x}_i e \underline{x}_j estão em um mesmo cluster, permanecerão juntos.

É dito um algorítmo hierárquico.

Que clusters juntar?

$$C_{j} \qquad \vec{m}_{j} = \frac{1}{n_{j}} \sum_{\forall \vec{x} \in C_{j}} \vec{x} \qquad F_{j} = \sum_{\forall \vec{x} \in C_{j}} \left| \vec{x} \right|^{2} - n_{j} \left| \vec{m}_{j} \right|^{2}$$

$$C_k \qquad \vec{m}_k = \frac{1}{n_k} \sum_{\forall \vec{x} \in C_k} \vec{x} \qquad F_k = \sum_{\forall \vec{x} \in C_k} \left| \vec{x} \right|^2 - n_k \left| \vec{m}_k \right|^2$$

$$C_{jk} = C_j \cup C_k \qquad \vec{m}_{jk} = \frac{n_j \vec{m}_j + n_k \vec{m}_k}{n_j + n_k} \qquad F_{jk} = \sum_{\forall \vec{x} \in C_j, C_k} |\vec{x}|^2 - (n_j + n_k) |\vec{m}_{jk}|^2$$

$$\Delta F = F_{jk} - (F_j + F_k) = \frac{n_j n_k}{n_j + n_k} |\vec{m}_j - \vec{m}_k|^2$$

$$\Delta F = \frac{n_j \ n_k}{n_j + n_k} \left| \vec{m}_j - \vec{m}_k \right|^2 \quad \text{deve ser pequeno}$$

escolher centros próximos:

$$\vec{m}_j \approx \vec{m}_k$$

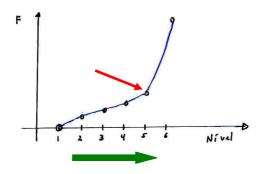
de preferência com populações desiguais:

$$\frac{n_{j}n_{k}}{n_{j}+n_{k}} \approx \begin{cases} n_{j} & se & n_{j} << n_{k} \\ \\ \frac{n_{j}}{2} & se & n_{j} \approx n_{k} \end{cases}$$

Quando parar o processo?

Heurística: um agrupamento "natural" (um número "bom" de classes) é alcançado no algorítmo associativo imediatamente antes de $F_{\rm in}$ apresentar um grande acréscimo em um único passo.

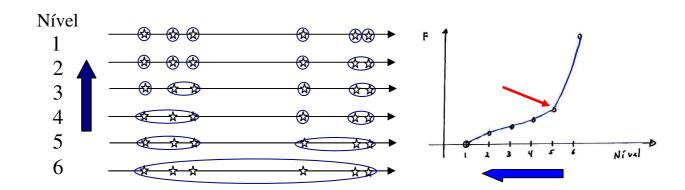
Por este critério o ponto de parada está indicado pela seta vermelha no exemplo anterior.



2.4.2 Algorítmo Divisivo, Dissociativo:

n elementos c classes

$$c = 1 2 3 \dots n-1 n$$
 $F_{in} = F_1 < F_{2} < F_3 < F_{n-1} < F_{n} = 0$



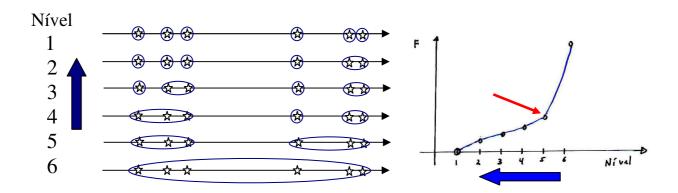
Algoritmo Hierárquico Dissociativo:

Inicialização: n elementos \Longrightarrow 1 cluster

Loop:

Divida o cluster que reduz Fin o máximo possível

Repetir até chegar ao número de clusters desejado ou $F_{\rm in}$ quase não decair



Que cluster dividir?

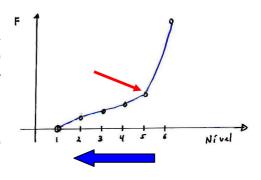
Critério: maior F_j , por exemplo

Método: $\vec{m}_j + \vec{\delta}$ $\vec{m}_i - \vec{\delta}$ **K-means**

Quando parar o processo ?

Heurística: um agrupamento "natural" (um número "bom" de classes) é obtido no algorítimo dissociativo imediatamente após $F_{\rm in}$ apresentar um grande decréscimo em um único passo.

Por este critério o ponto de parada esta indicado pela seta vermelha no exemplo anterior.



2.5 - Algorítmos Mistos

Iniciar com K-means

Associar / Dissociar verificando a variação de $F_{\rm in}$