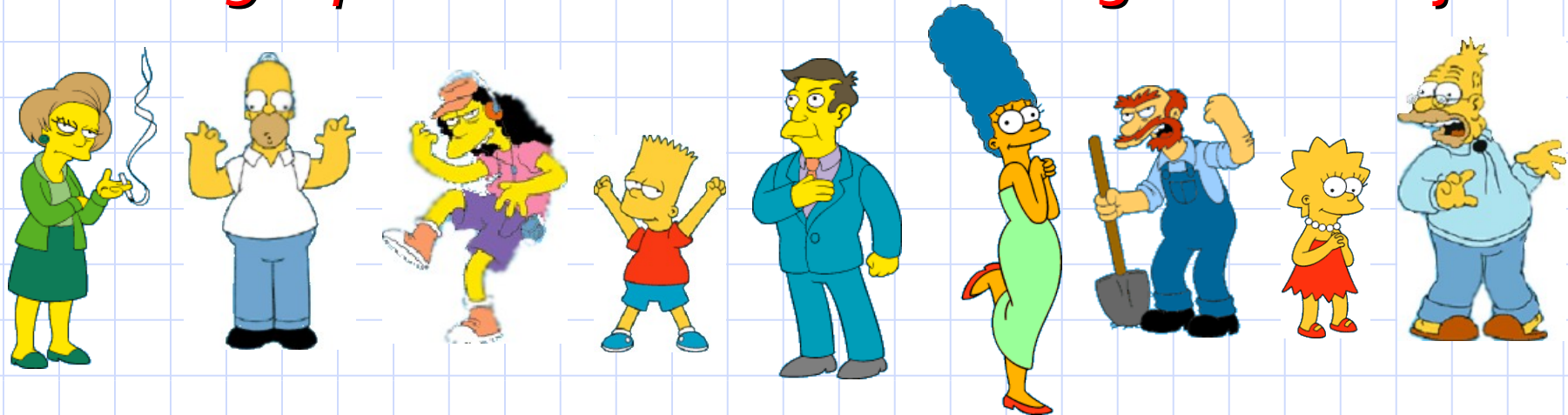


Agrupamento de Dados

Prof. Debora Medeiros

*Slides adaptados de Eduardo Raul Hruschka (USP)

Como agrupar naturalmente os seguintes objetos?



Família

Empregados

Mulheres

Homens

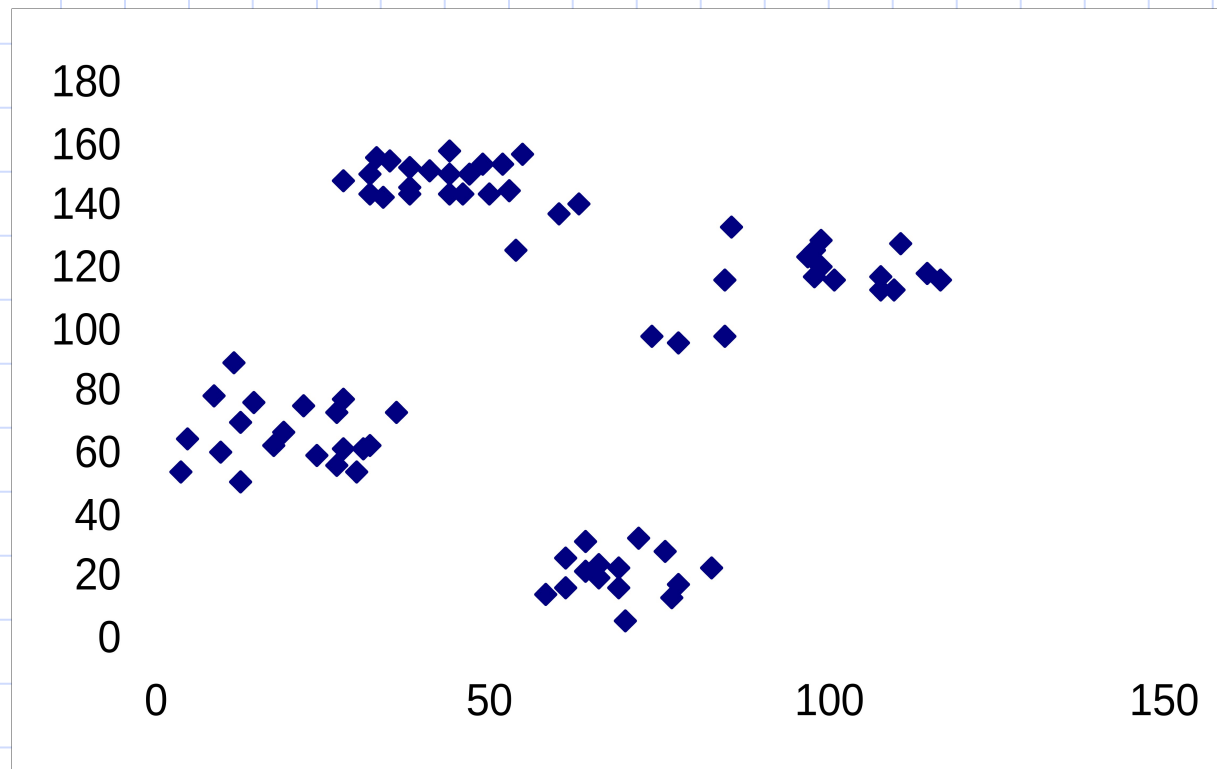
→ **Cluster** é um conceito subjetivo!

O que é um grupo (*cluster*)?

- Definições subjetivas:
 - “Semelhanças entre objetos”.
 - Quais atributos devemos considerar para computar similaridades?



- Numa “abordagem matemática”, critérios numéricos usualmente consideram:
 - Homogeneidade (coesão interna);
 - Heterogeneidade (separação entre grupos);

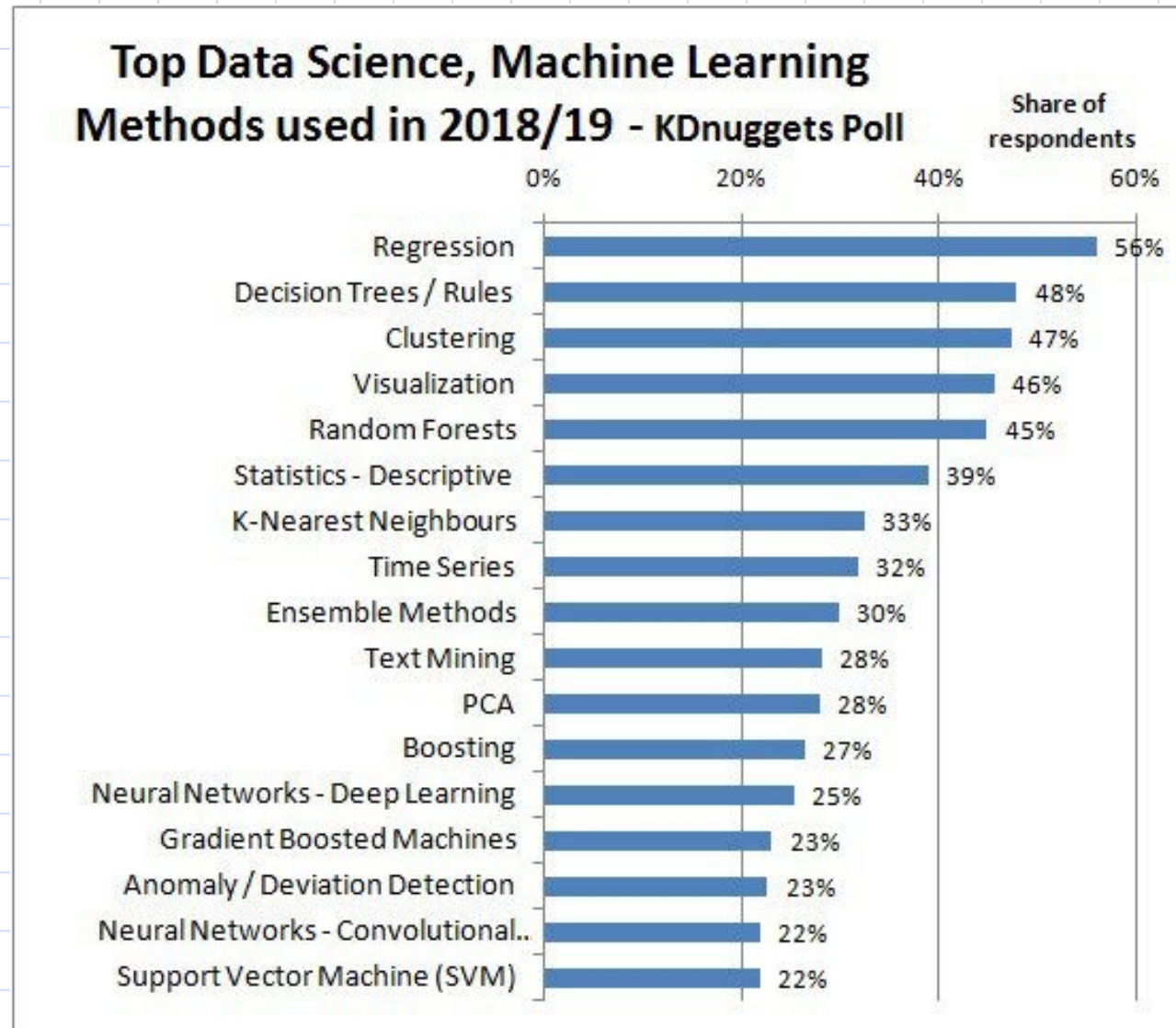


(Ruspini, 1970)

Agrupamento de Dados (ADs) é uma técnica importante para Análise Exploratória de Dados :

- Engenharia;
- Biologia;
- Psicologia;
- Medicina;
- Administração (*Marketing* , Finanças,...);
- Ciência da Computação:
 - Bioinformática;
 - Dados coletados via sensores;
 - Componentes de sistemas inteligentes;
 - Componentes de algoritmos para aprendizado de máquina, ...

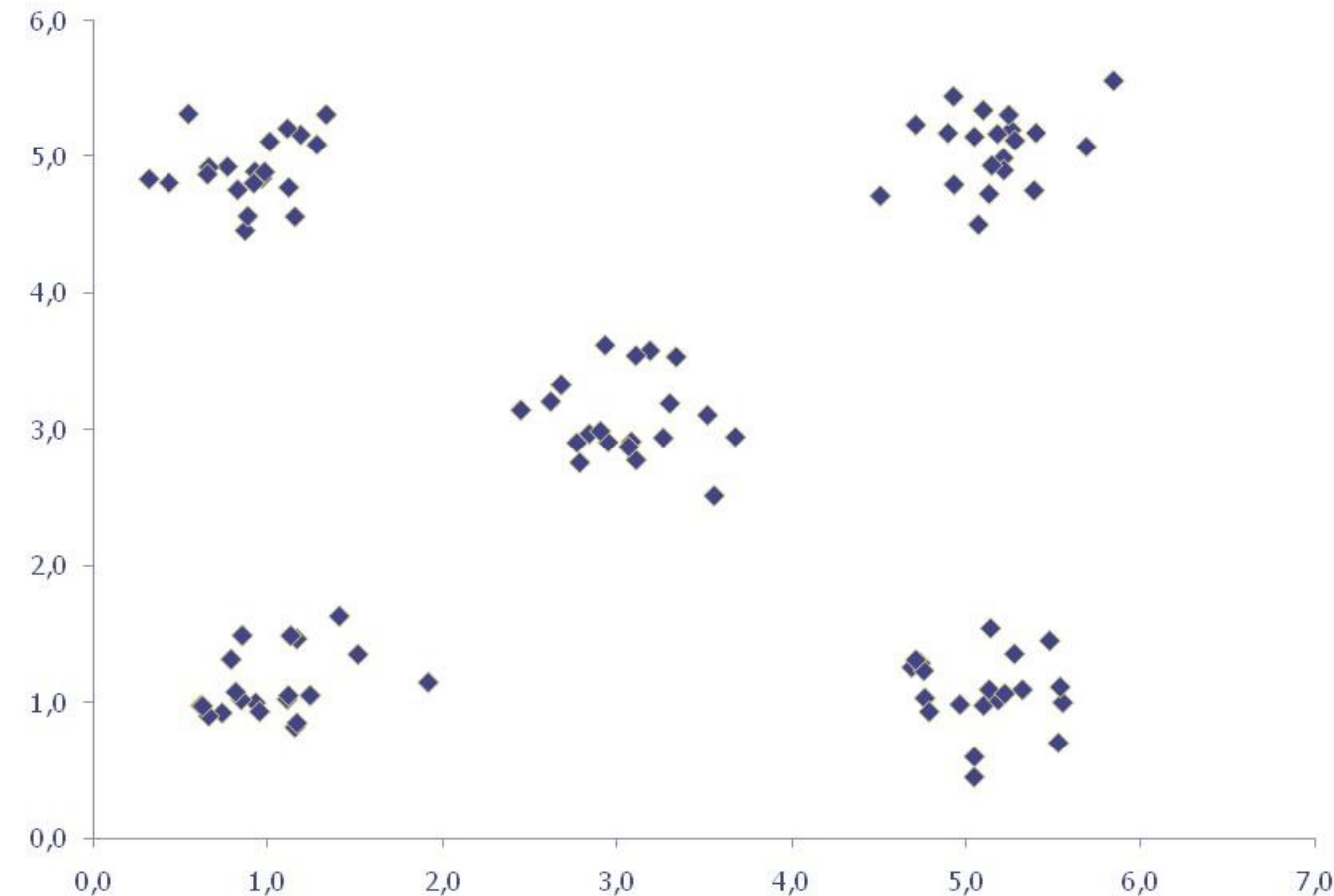
<https://www.kdnuggets.com/2019/04/top-data-science-machine-learning-methods-2018-2019.html>



2. Conceitos Básicos

Algumas Definições (Everitt, 1974):

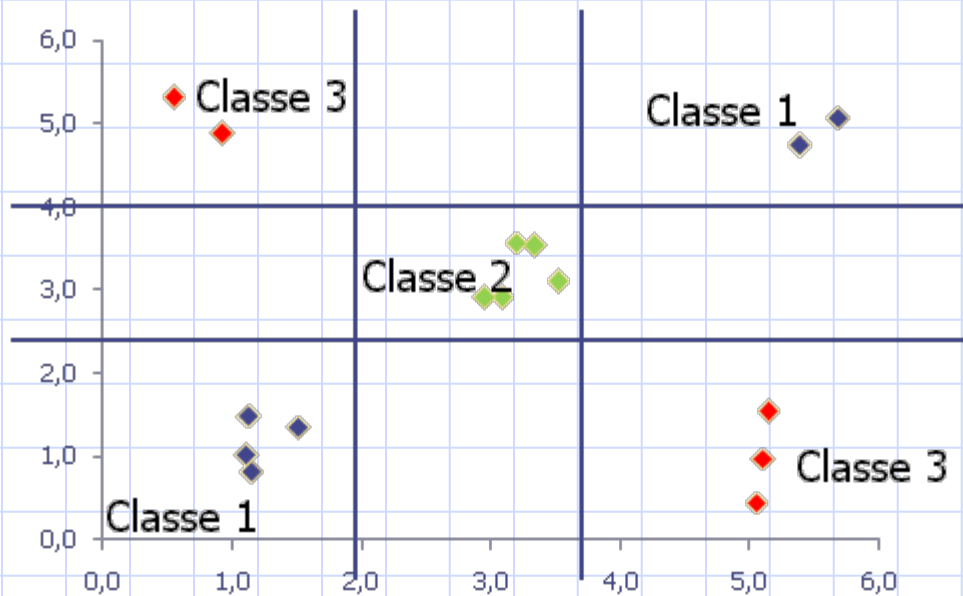
- Um cluster (**grupo**) é um conjunto de **entidades semelhantes**, e entidades pertencentes a **diferentes clusters não são semelhantes**.
- Um grupo é uma aglomeração de pontos no espaço tal que a **distância** entre quaisquer **dois pontos no grupo** é **menor** do que a **distância** entre qualquer ponto no grupo e qualquer **ponto fora** deste.
- Grupos podem ser descritos como **regiões** conectadas de um espaço multidimensional contendo uma **densidade** de pontos relativamente **alta, separada** de outras tais regiões por uma região contendo uma **densidade** relativamente **baixa** de pontos.
- Humanos reconhecem **clusters** no plano quando os vêem, sem saber explicar exatamente ⁸porquê (Jain & Dubes, 1988) ...



Agrupamento:
Indução de
grupos a partir
da base de
dados...

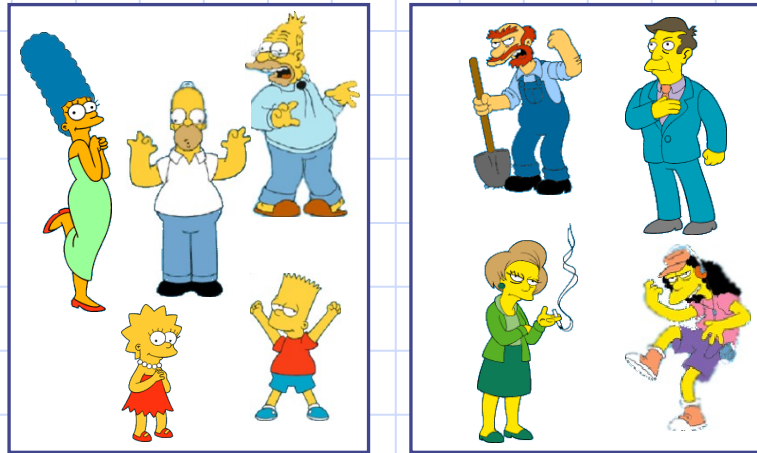
Agrupamento X Classificação?

...

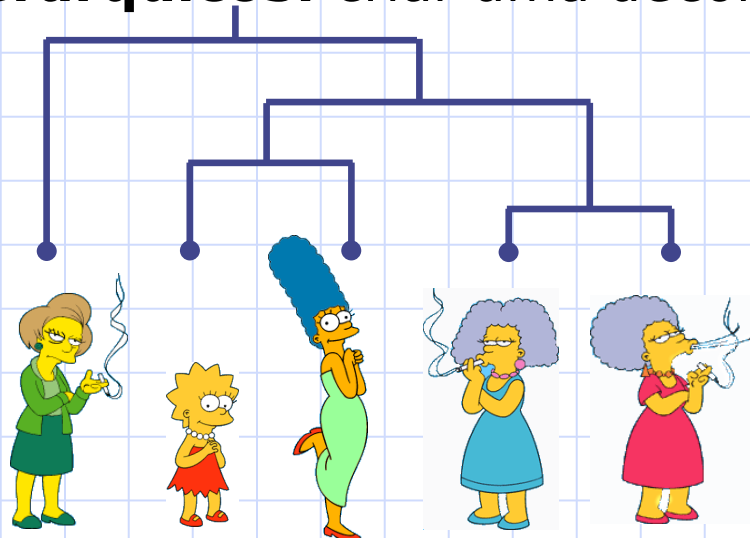


5 grupos

Algoritmos para particionamento: construir partições.

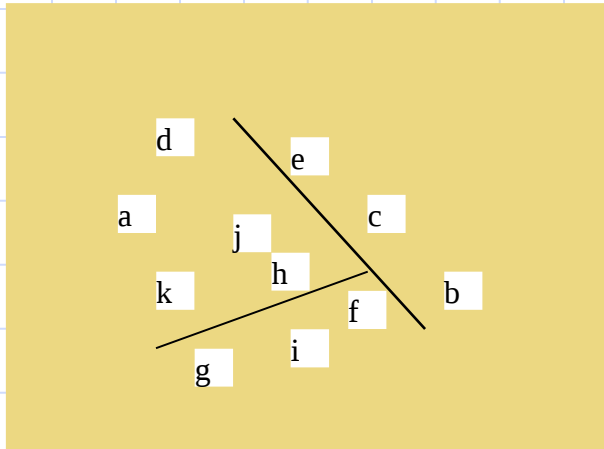


Algoritmos hierárquicos: criar uma decomposição hierárquica.

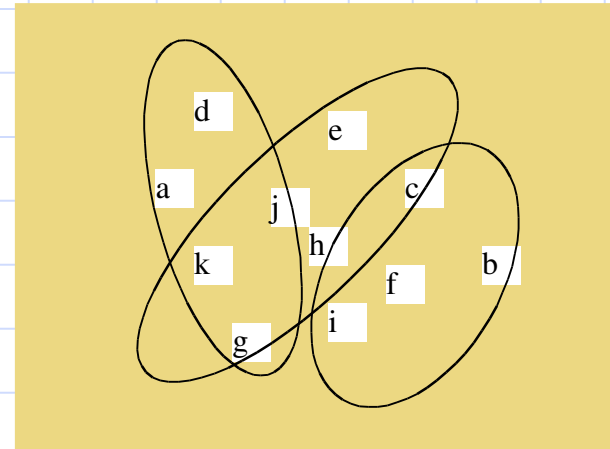


- **Métodos para Particionamento:**

Rígidos



Sobreposição



- Em princípio requerem a definição, a priori, do número de grupos;
- Métodos hierárquicos, por sua vez...

Algoritmos particionais

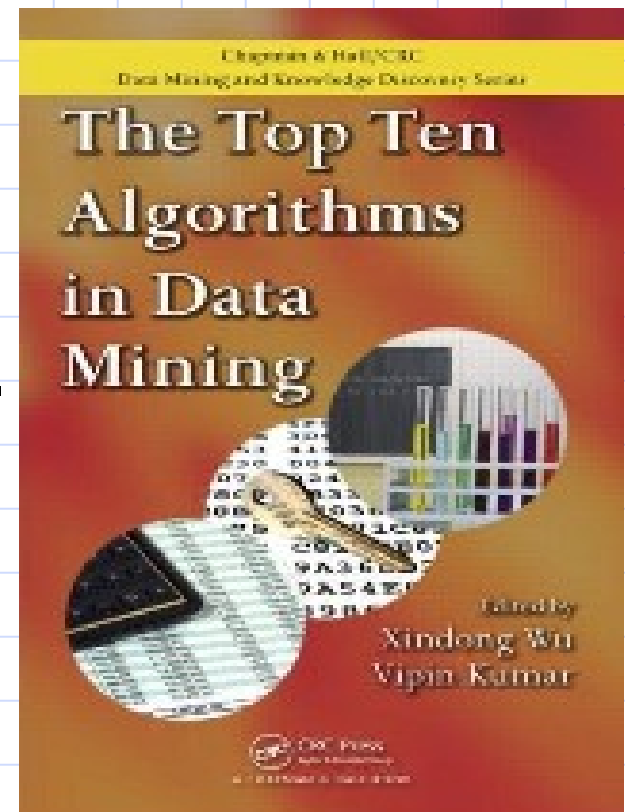
Partições rígidas

Abordaremos um algoritmo amplamente usados na prática:

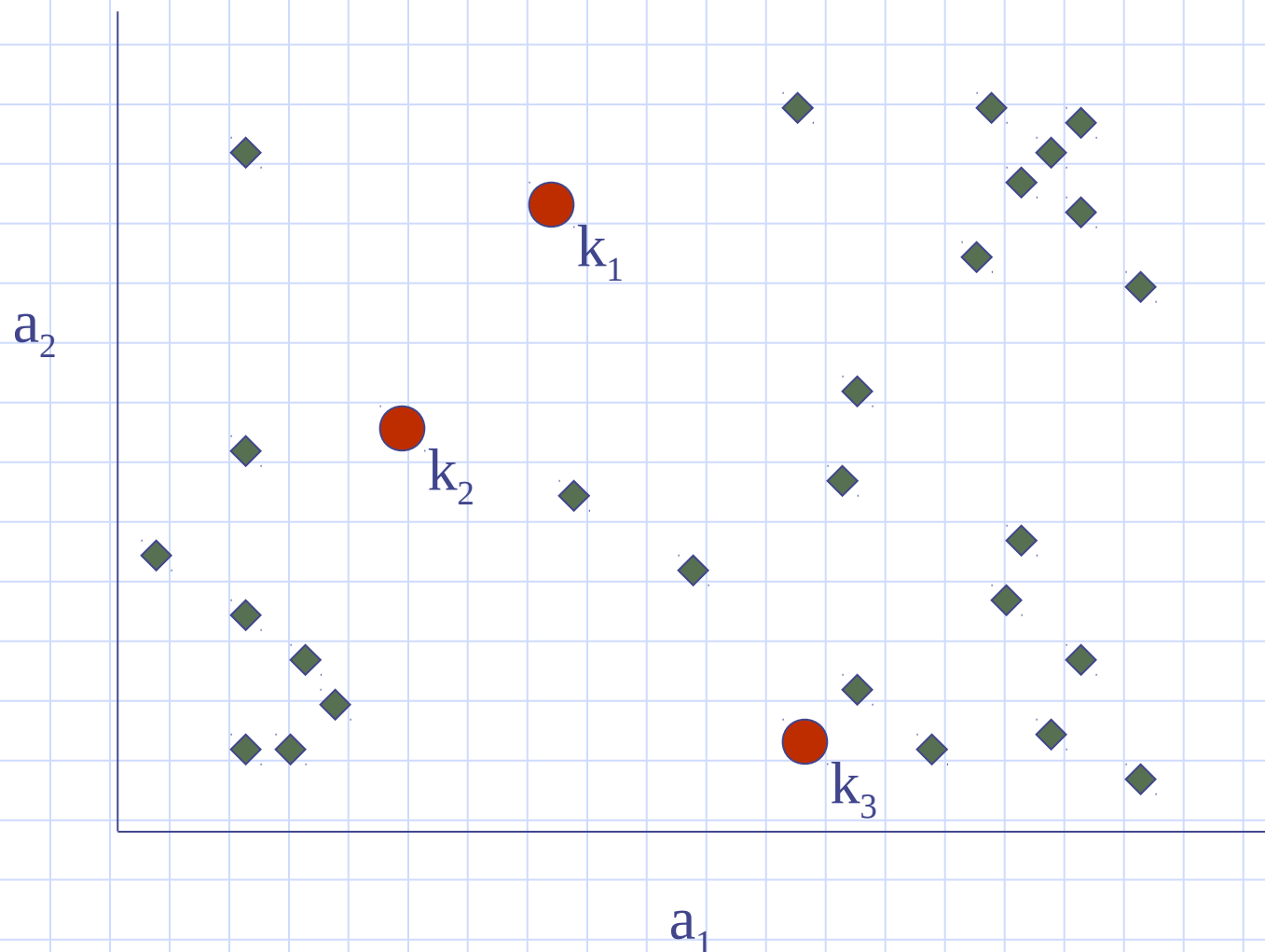
- ***k*-médias (*k-means*);**

Algoritmo *k*-médias (MacQueen, 1967)

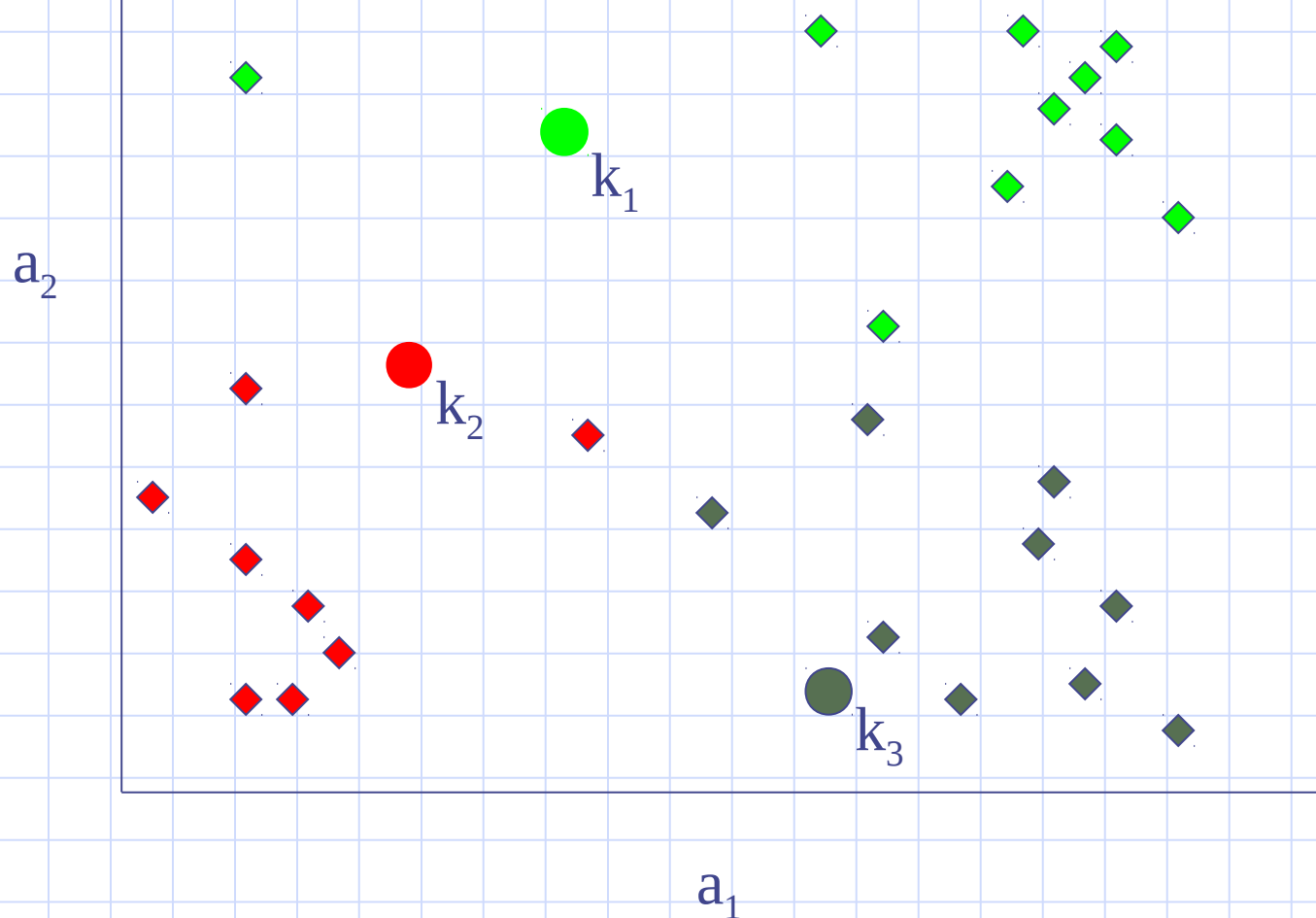
- Amplamente usado na prática:
 - Simplicidade;
 - Interpretabilidade;
 - Eficiência Computacional.



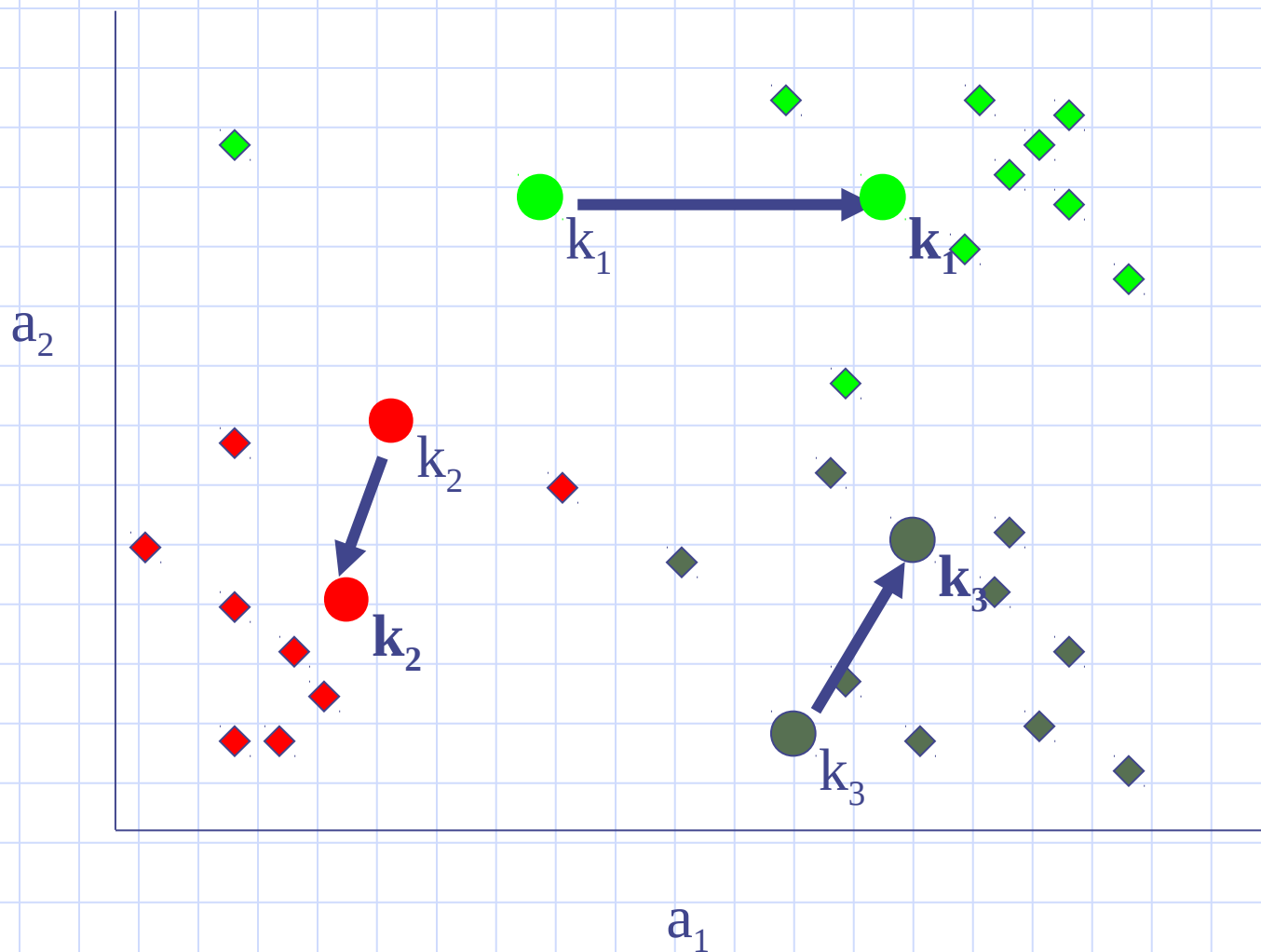
- Assumamos que queremos encontrar três *clusters* ($k = 3$) para uma base de dados bi-dimensional:



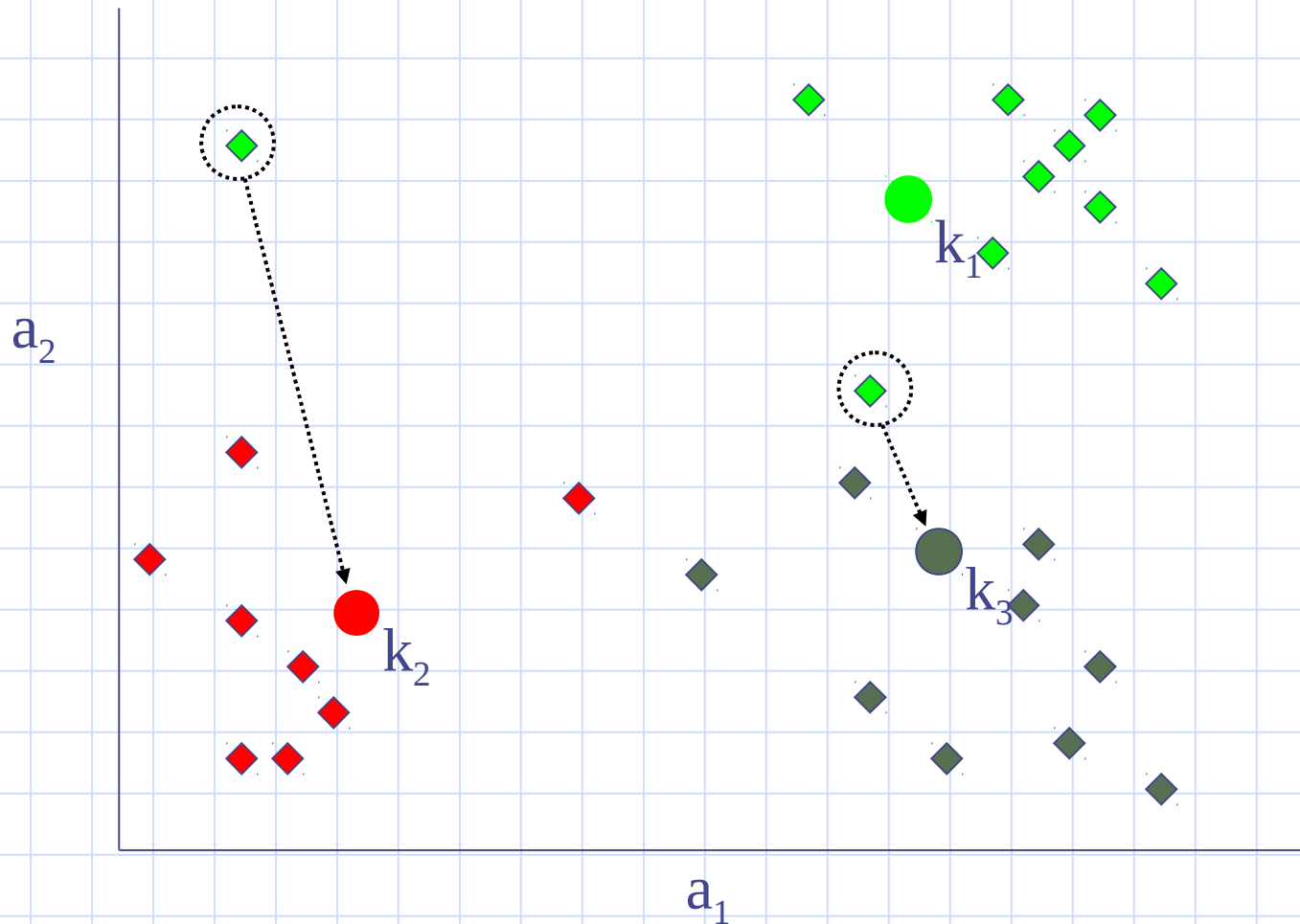
- Calcular dissimilaridades entre objetos e protótipos (k_1, k_2, k_3), encontrando grupos iniciais pela regra do vizinho mais próximo:



- Atualizar os protótipos (centróides) dos grupos:



- Calcular dissimilaridades entre objetos e centróides;
- Atualizar *clusters* (regra do vizinho mais próximo);



- Repetir até convergência/ número de iterações.

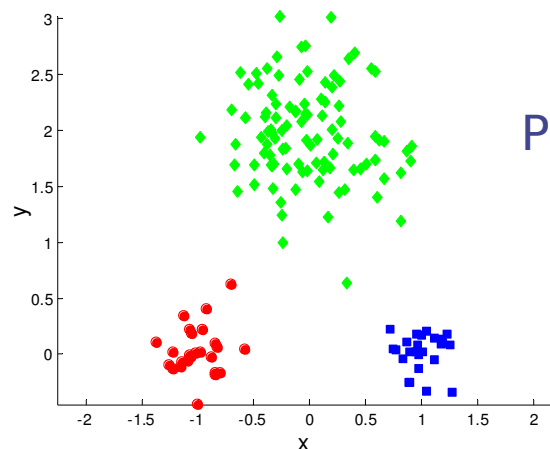
Algoritmo básico:

1. Selecionar k pontos (*centróides* iniciais);
2. Repetir até “convergir”:
 - 2.1 Formar k grupos atribuindo cada ponto ao seu centróide mais próximo;
 - 2.2 Re-computar o centróide de cada grupo;

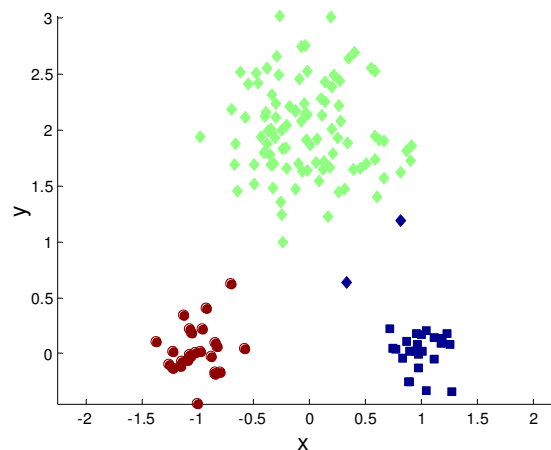
Detalhes sobre o k -médias:

- ◆ Centróides iniciais são frequentemente escolhidos aleatoriamente.
 - *Clusters* obtidos podem variar de uma rodada para outra?
- ◆ *Proximidade* pode ser medida por meio de Distância Euclidiana, co-seno, correlação, etc.
- ◆ k -média converge, geralmente nas primeiras iterações;
- ◆ Vejamos alguns exemplos interessantes...

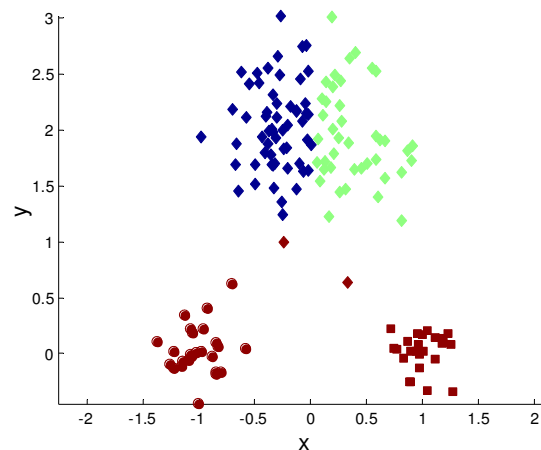
Consideremos duas partições diferentes obtidas para $k = 3$:



Pontos originais

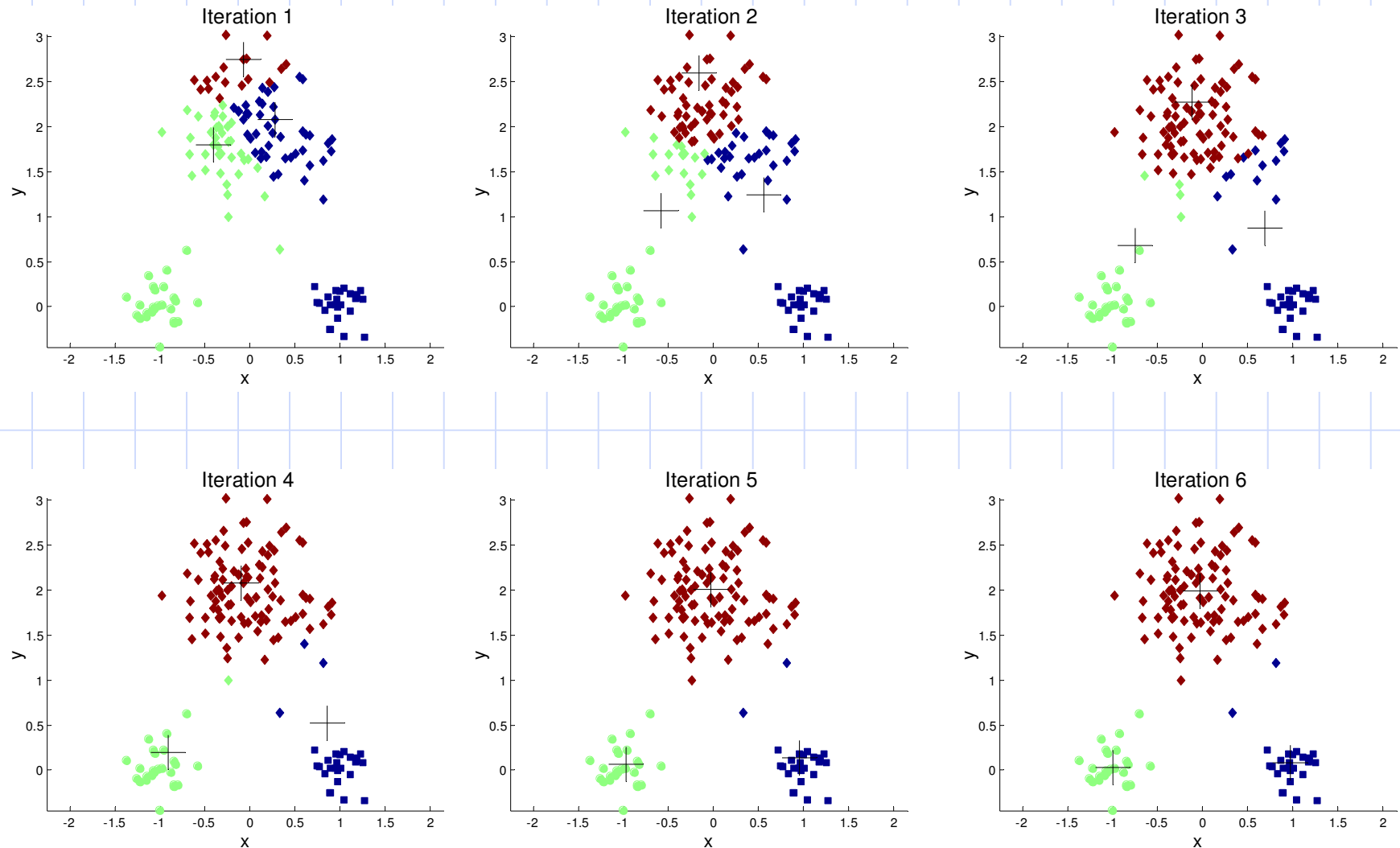


Partição ótima

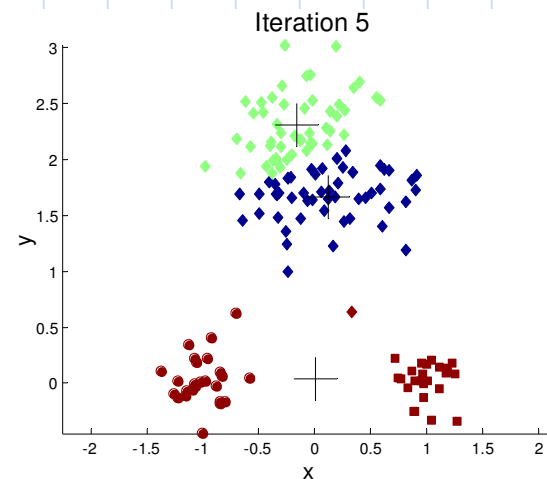
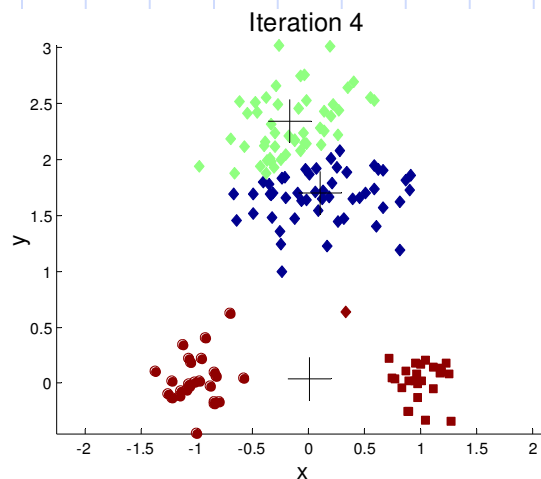
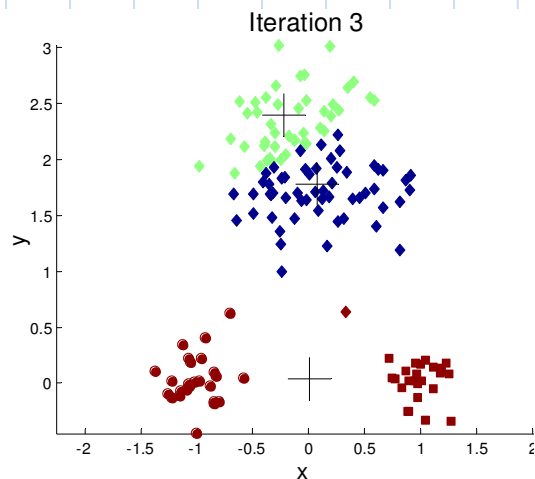
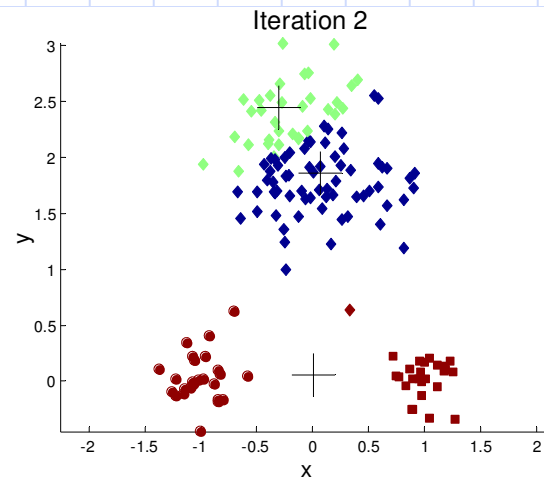
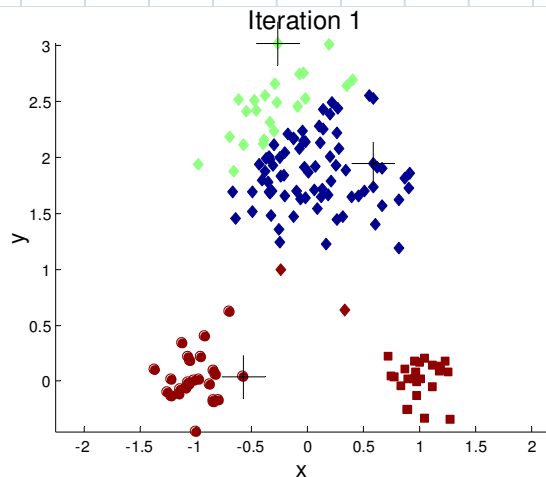


Partição Sub-ótima

Importância da escolha dos centróides iniciais:



Importância da escolha dos centróides iniciais...



Soluções para inicialização?

- ◆ Múltiplas execuções:
 - Ajuda, mas pequena P_{sucesso} ;
- ◆ Amostragem via métodos hierárquicos;
- ◆ Seleção “informada” de centróides distantes entre si;
- ◆ Algoritmos de busca (e.g., evolutivos);

Avaliando os grupos obtidos:

◆ Soma dos erros quadráticos:

$$SEQ = \sum_{i=1}^k \sum_{x \in C_i} dist^2(m_i, x)$$

- Dadas duas partições, escolher aquela que apresenta SEQ menor;
- Aumento de k : tende a diminuir, por si só, SEQ;

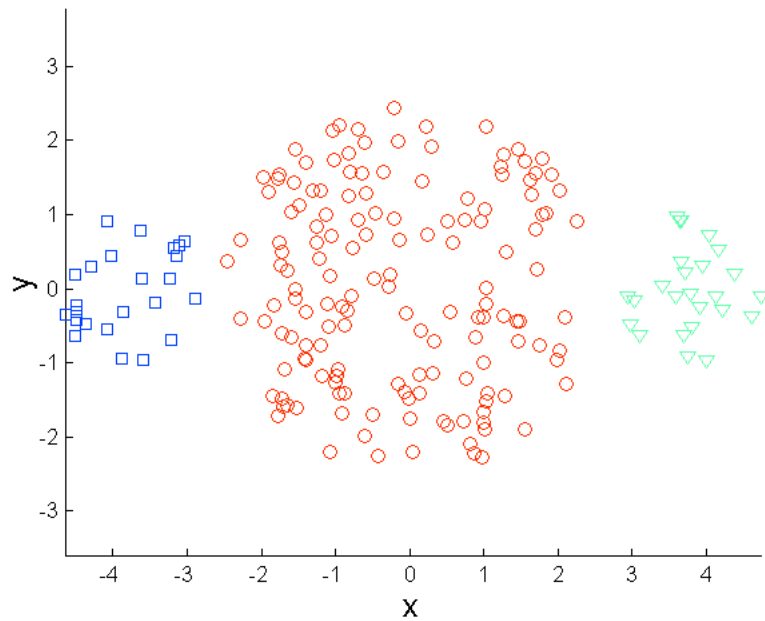
Limitações do k -médias:

◆ Grupos de diferentes:

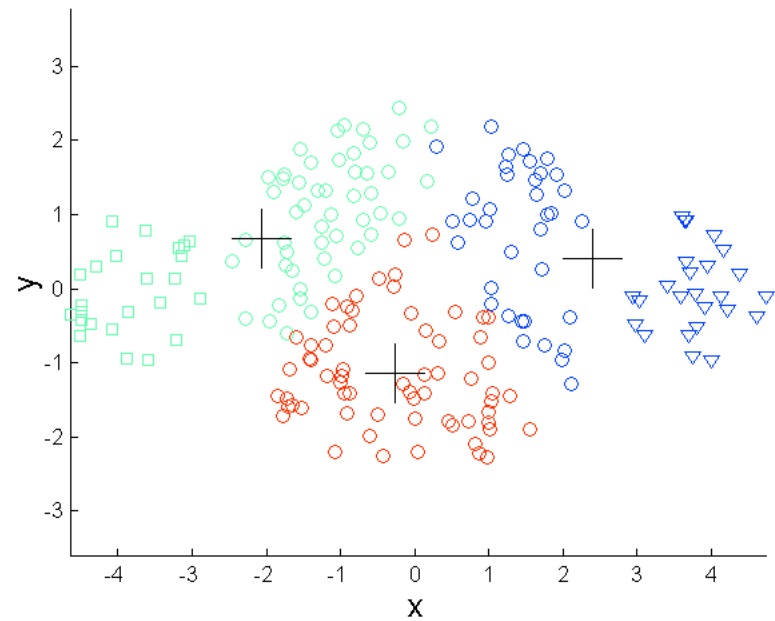
- Tamanhos;
- Densidades;
- Formas não globulares.

◆ *Outliers*.

Limitações do k -médias: grupos de tamanhos diferentes

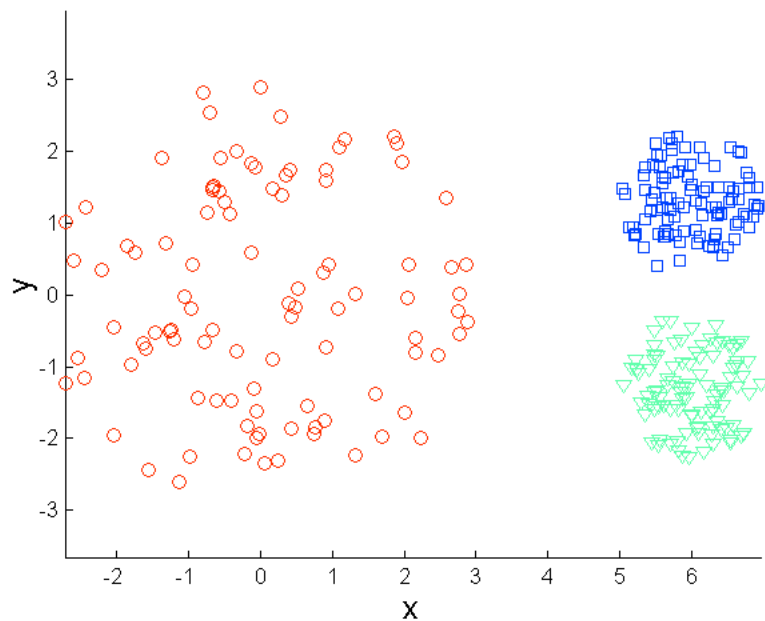


Pontos originais

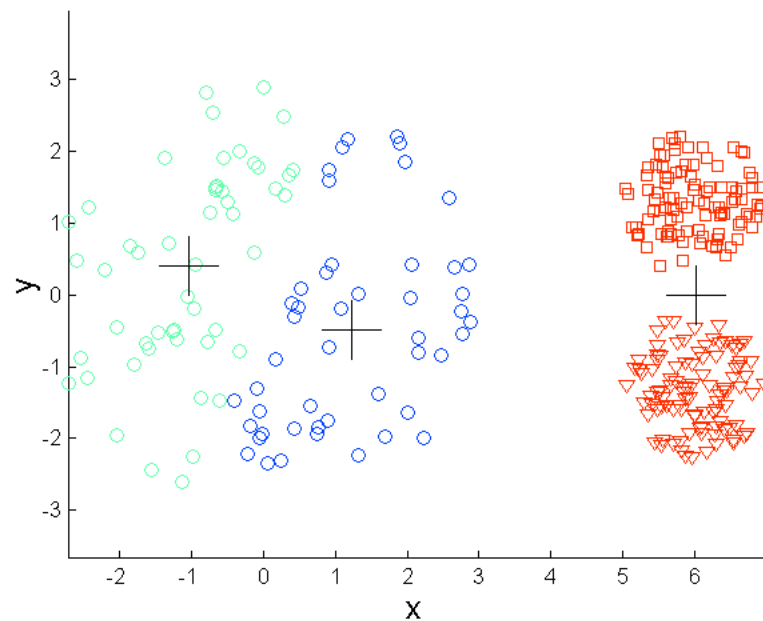


3-médias

Limitações do k -médiás: densidades diferentes

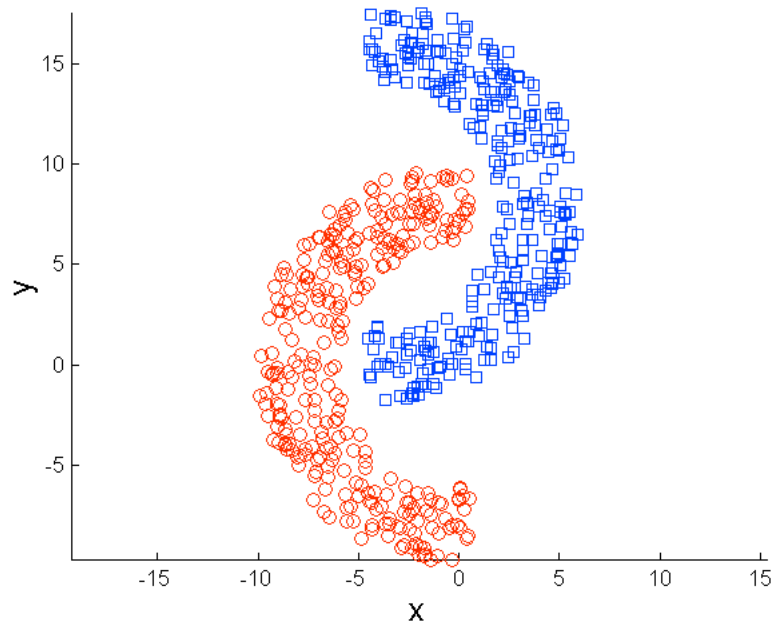


Pontos originais

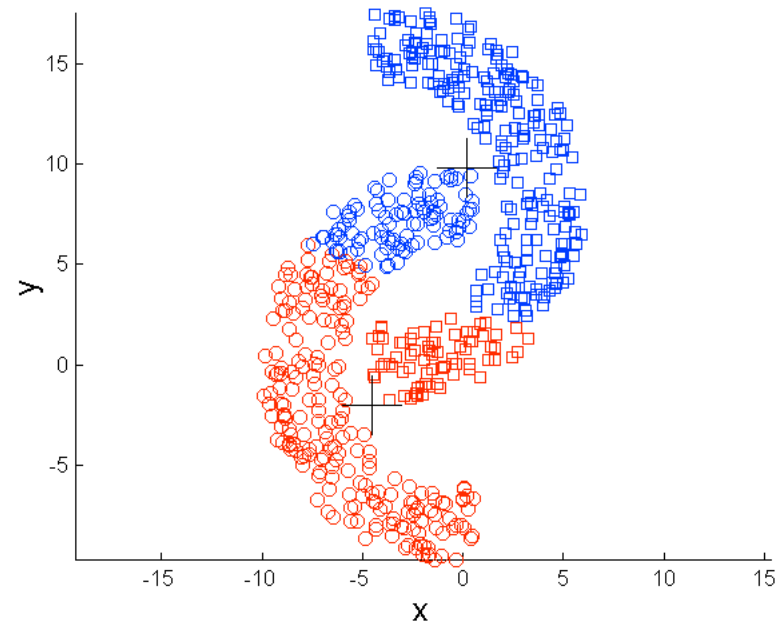


3-médias

Limitações do k -médiãs: formas não globulares

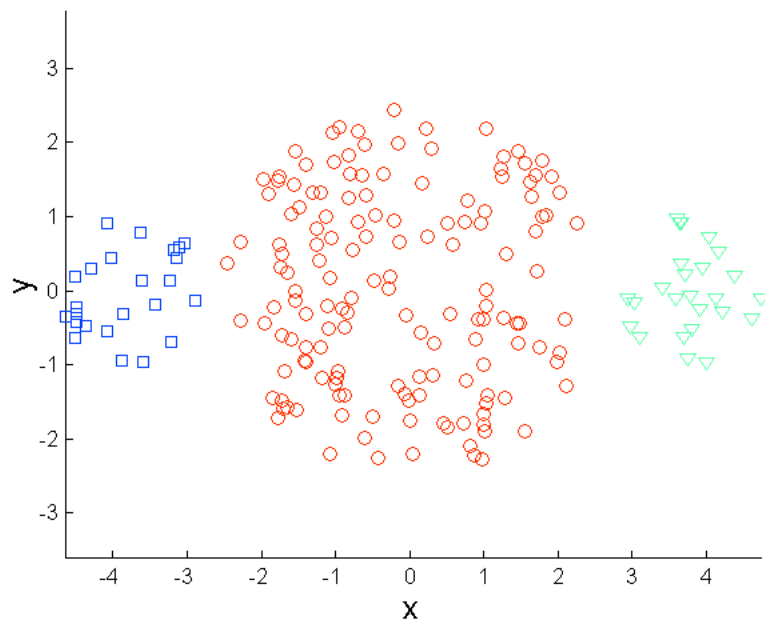


Pontos originais

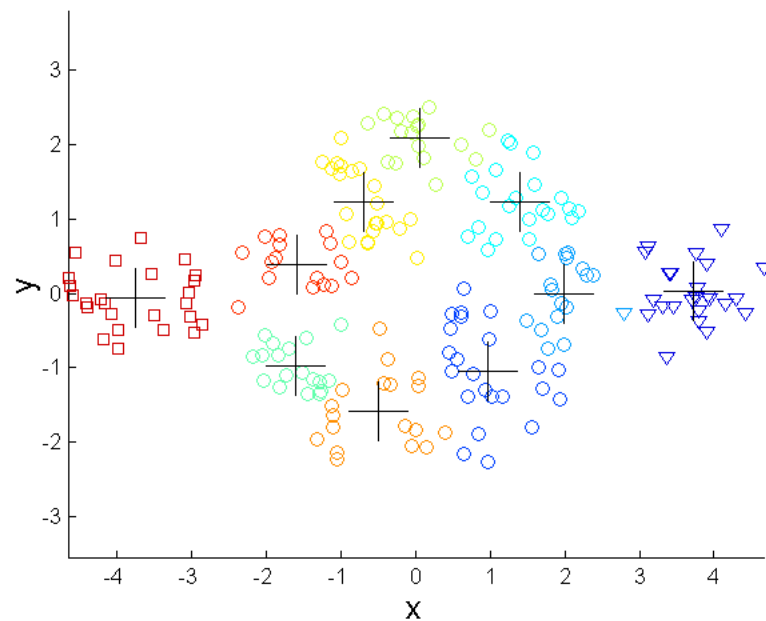


2-médias

Superando algumas limitações do *k-médias*

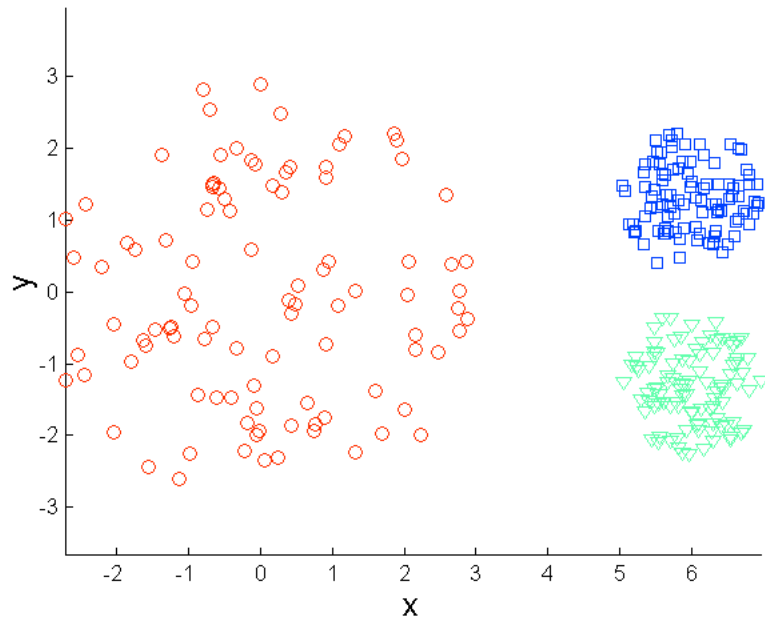


Pontos originais

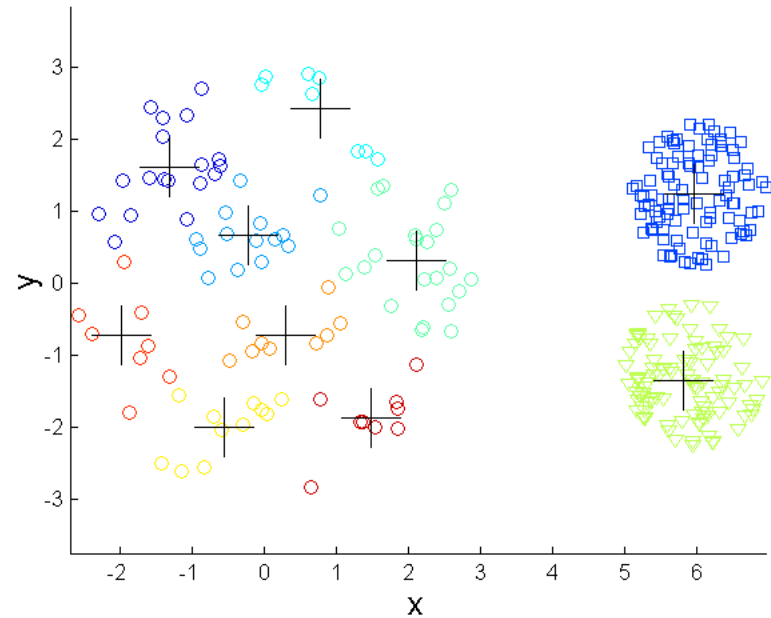


“Mais Grupos”

Superando algumas limitações do *k*-*médias*...

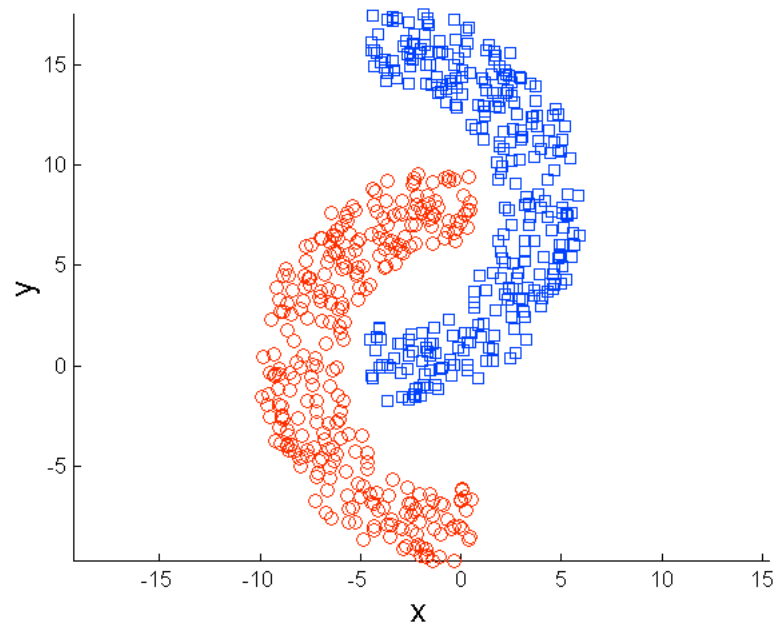


Pontos originais

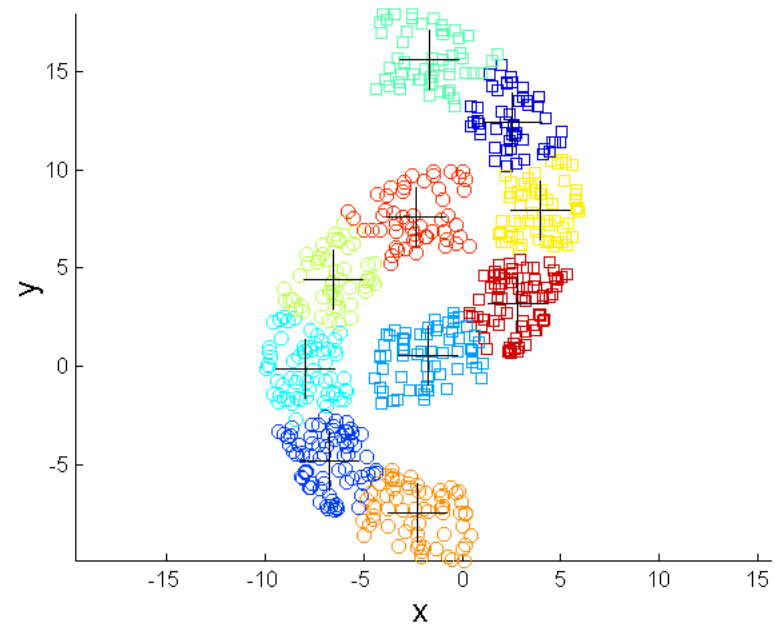


“Mais grupos”

Superando algumas limitações do *k-médias*...



Pontos originais



“Mais grupos”



Validação

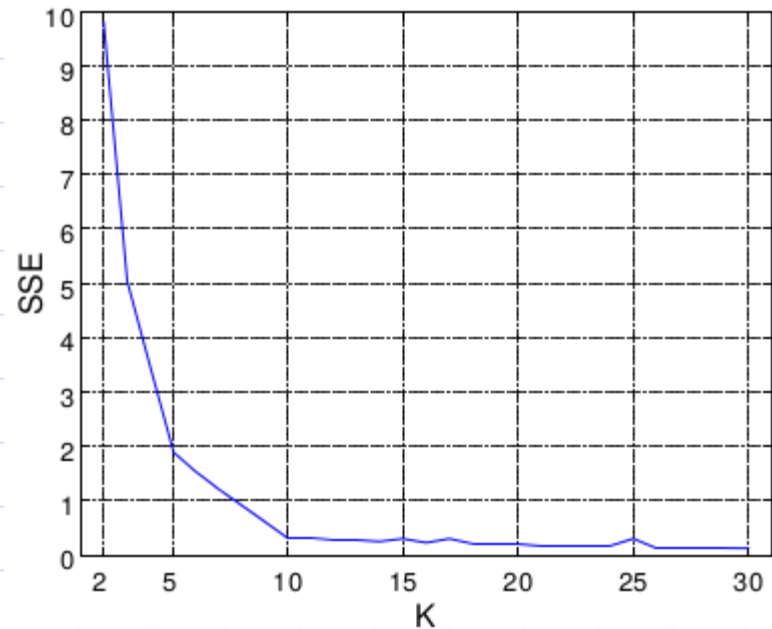
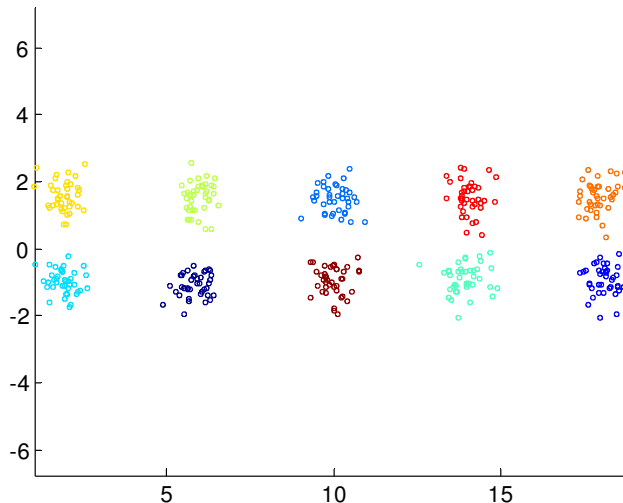
Critérios de validade:

Medidas numéricas para julgar diversos aspectos de validade:

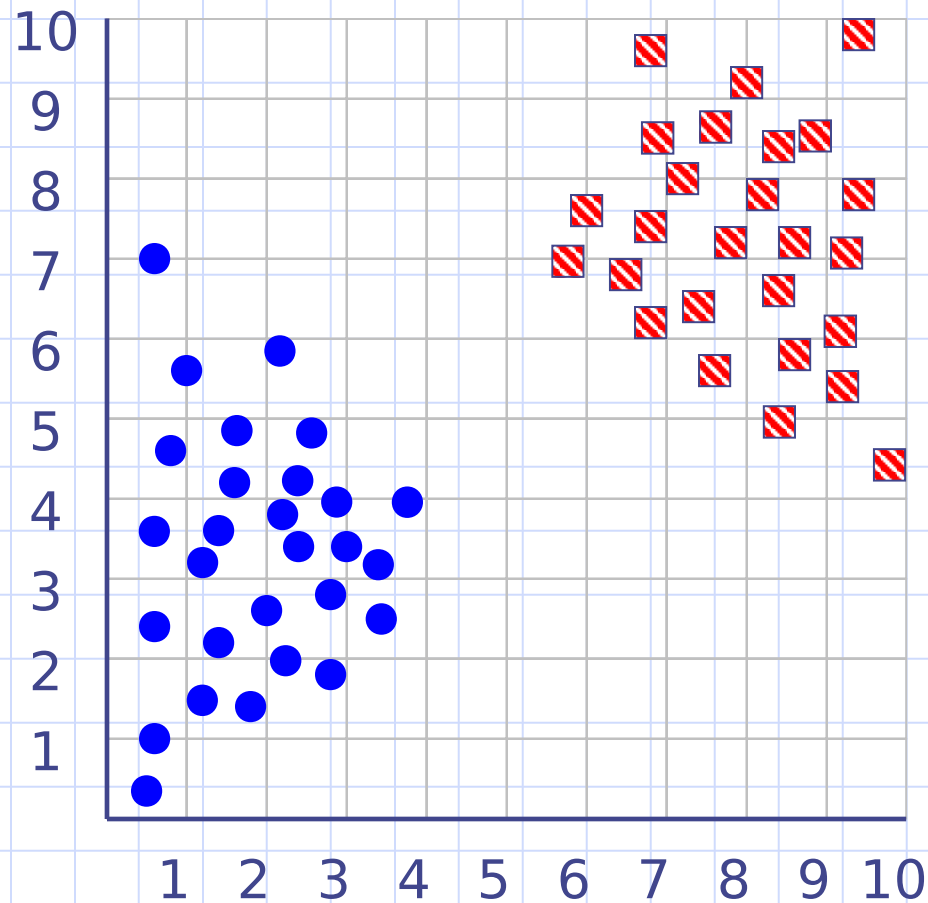
- **Índices Internos:** medem a qualidade de uma estrutura de agrupamento independentemente de informação externa (e.g., EQ);
 - **Índices Relativos:** comparam duas partições (ou grupos).
 - **Índices Externos:** medem o quão bem os rótulos dos grupos representam categorias pré-estabelecidas (e.g., Rand);
- Iniciaremos por um índice interno que já foi “indiretamente” estudado...

Erro Quadrático (EQ):

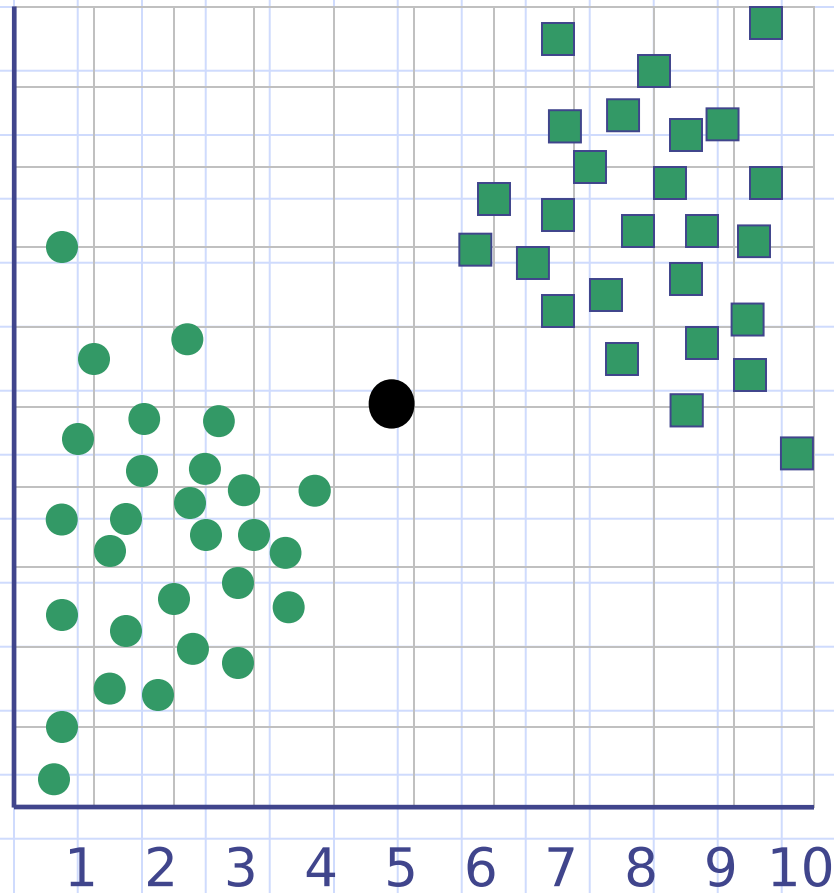
- Usado para medir a qualidade da estrutura de grupos obtida sem se considerar informações externas;
- Comparar partições (mesmo k) e grupos (EQ médio) ;
- Pode ser usado para estimar o número de grupo (com cuidado!):



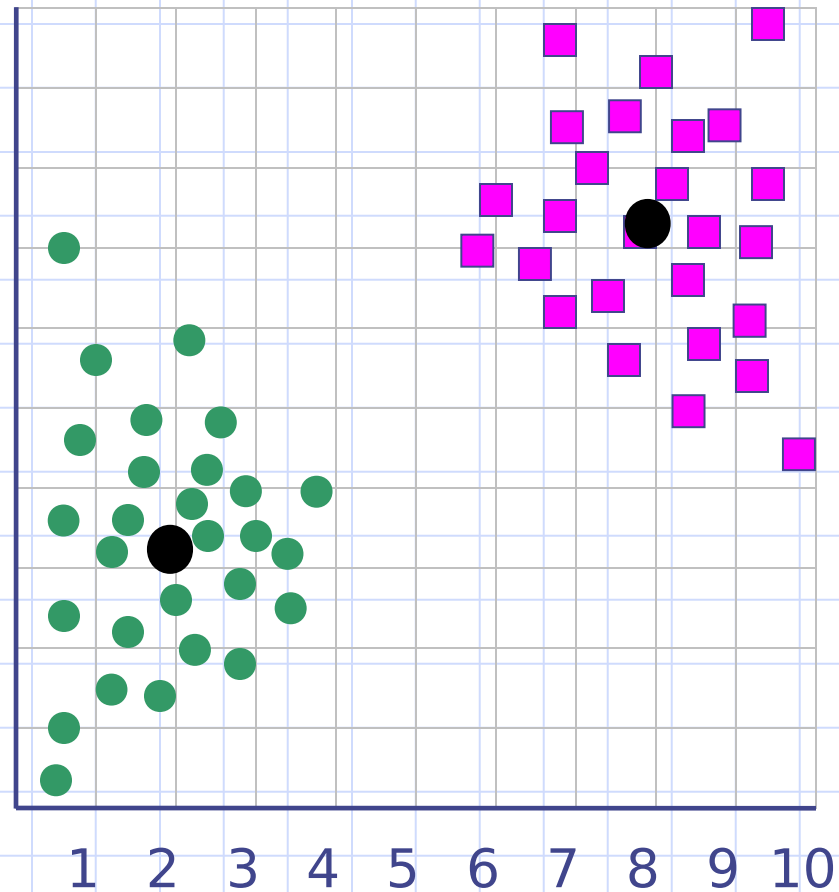
Estimando o número de grupos:



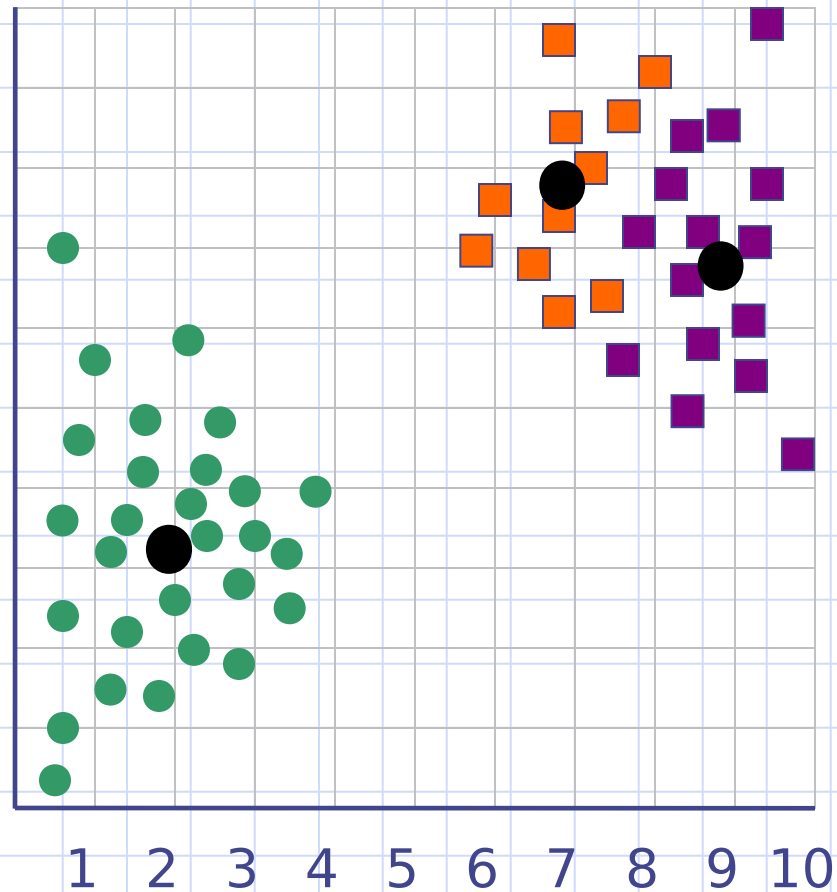
Para $k = 1$, $EQ = 873$.



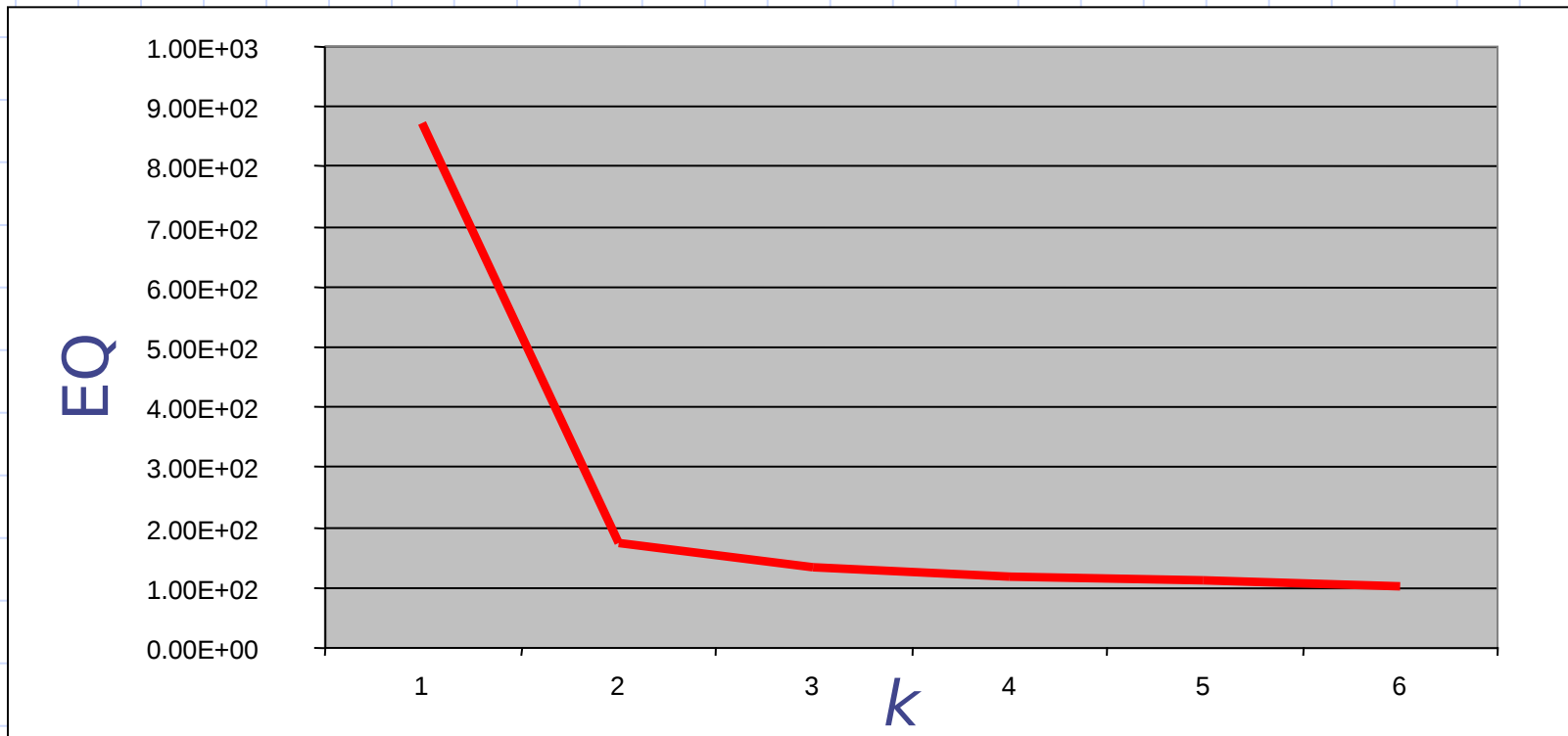
Para $k = 2$, $EQ = 173$.



Para $k = 3$, EQ = 133.



Plotar EQ para $k = 1, \dots, 6$, tentando identificar um joelho:

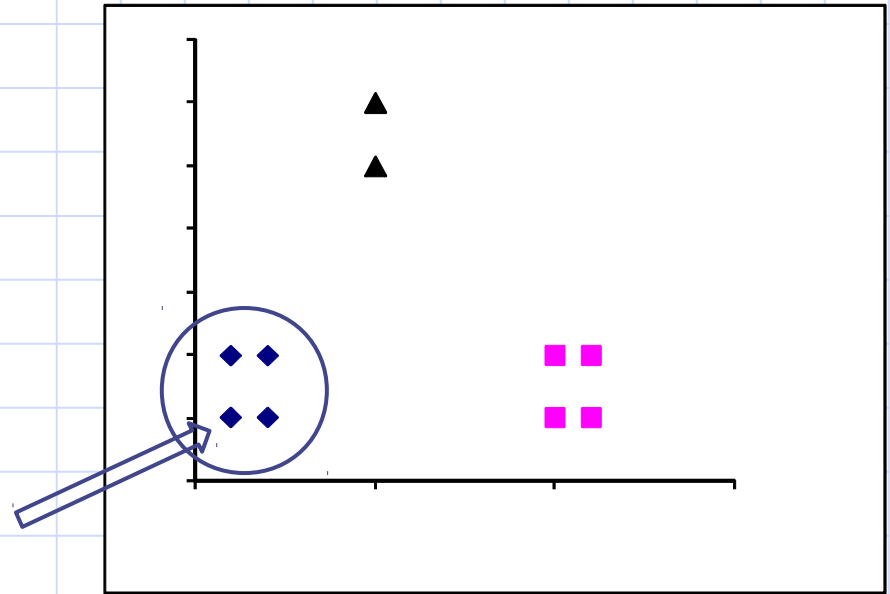


Índices Relativos:

- Normalmente usados para comparar partições diferentes e para estimar o número de grupos de dados;

Silhueta (Kaufman & Rousseeuw, 1990):

i	x	y	$cluster$
1	1	1	C_1
2	1	2	C_1
3	2	1	C_1
4	2	2	C_1
5	10	1	C_2
6	10	2	C_2
7	11	1	C_2
8	11	2	C_2
9	5	5	C_3
10	5	6	C_3



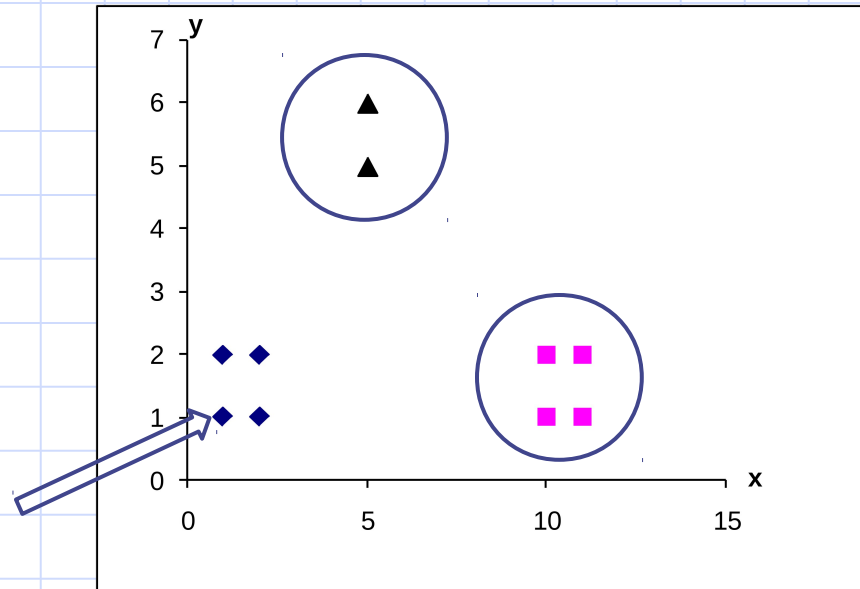
- Quão bem “classificado” está o i -ésimo objeto (e.g., $i=1$)?

1) Dissimilaridade em relação aos objetos do seu grupo:

□ $a(i)$: dissimilaridade média de i aos objetos do seu grupo:

- $(d_{12} + d_{13} + d_{14})/3$.

i	x	y	$cluster$
1	1	1	C_1
2	1	2	C_1
3	2	1	C_1
4	2	2	C_1
5	10	1	C_2
6	10	2	C_2
7	11	1	C_2
8	11	2	C_2
9	5	5	C_3
10	5	6	C_3



2) O que se pode dizer em relação à sua dissimilaridade relativa aos objetos do grupo vizinho?

$$b(i): \min\{ (d_{15} + d_{16} + d_{17} + d_{18})/4 ; (d_{19} + d_{110})/2 \}.$$

Silhueta (Kaufman & Rousseeuw, 1990):

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$$S_M = \frac{\sum_{i=1}^N s(i)}{N}$$

$a(i)$: dissimilaridade média de i aos outros objetos de seu grupo.

$b(i)$: dissimilaridade média de i em relação aos objetos de seu grupo vizinho.

→ Aumentar k tende a diminuir $a(i)$, mas também tende a diminuir $b(i)$.

Existem dezenas de índices reportados na literatura;

Maior parte dos índices é baseada em conceitos de compactação e separabilidade...

Índices Externos:

- Comparar uma partição obtida com uma partição de referência;
- Aplicações:
 - Experimentos controlados;
 - Prática?
- Consideremos:
 - Um par de objetos (x_i, x_j) de $X = \{x_1, x_2, \dots, x_n\}$;
 - Uma partição dos dados $C = \{C_1, C_2, \dots, C_k\}$;
 - Uma partição de referência $P = \{P_1, P_2, \dots, P_p\}$:

Há 4 casos possíveis para atribuir $(\mathbf{x}_i, \mathbf{x}_j)$:

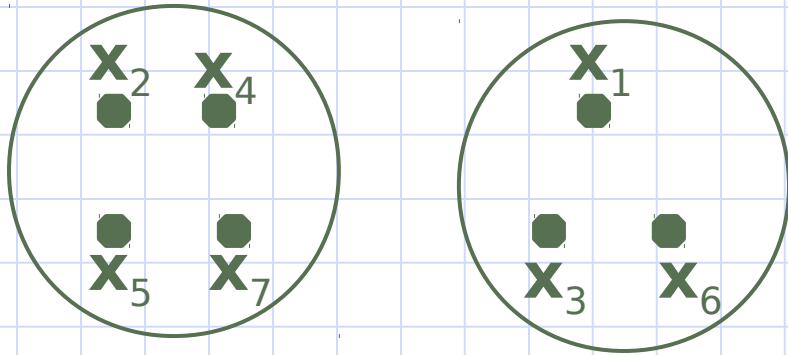
Caso 1: \mathbf{x}_i e \mathbf{x}_j estão no mesmo grupo em **C** e na mesma categoria em **P**; (a)

Caso 2: \mathbf{x}_i e \mathbf{x}_j estão no mesmo grupo em **C** e em categorias diferentes em **P**; (b)

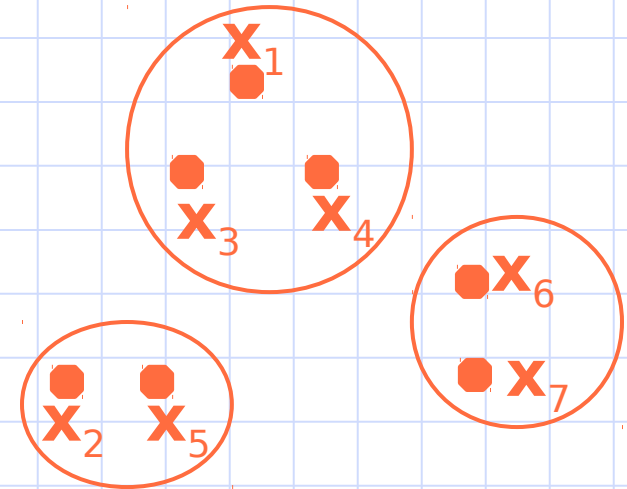
Caso 3: \mathbf{x}_i e \mathbf{x}_j estão em grupos diferentes em **C** e na mesma categoria em **P**; (c)

Caso 4: \mathbf{x}_i e \mathbf{x}_j estão em grupos diferentes em **C** e em categorias diferentes em **P**; (d)

Exemplo (Xu & Wunch, 2009):



Partição de Referência - **P**.



Clusters obtidos - **C**.

$$a = 2 \text{ (13,25)}$$

$$b = 3 \text{ (14,34,67)}$$

$$c = 7 \text{ (16,24,27,36,45,47,57)}$$

$$d = 9 \text{ (12,15,17,23,26,35,37,46,56)}$$

Nº de pares:

$$a+b+c+d = n(n-1)/2 = 21.$$

Alguns índices muito usados:

Jaccard (1908): $J = \frac{a}{a+b+c}$

Principais referências usadas para preparar essa aula:

- Xu, R., Wunsch, D., **Clustering**, IEEE Press, 2009.
 - Capítulo 4.
- Tan, Steinbach & Kumar, **Introduction to Data Mining**, Pearson, 2006.
 - Capítulo 8, pp. 496-515.
- Jain, A. K., Dubes, R. C., **Algorithms for Clustering Data**, Prentice Hall, 1988.
 - Capítulo 3, pp. 89-142.
- Bishop, C. M., **Pattern Recognition and Machine Learning**, 2006.
 - Capítulo 9, pp. 423-439.