

Aprendizado Não-Supervisionado

Advanced Institute for Artificial Intelligence

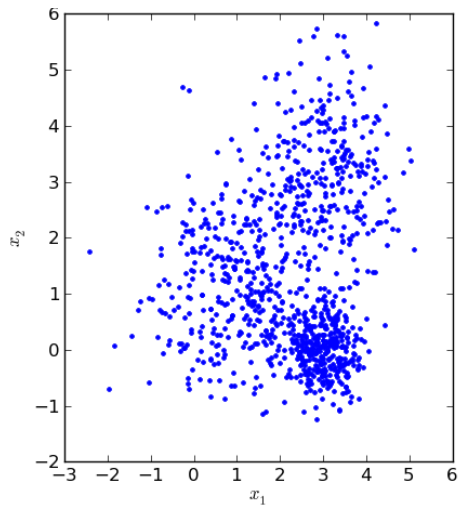
<https://advancedinstitute.ai>

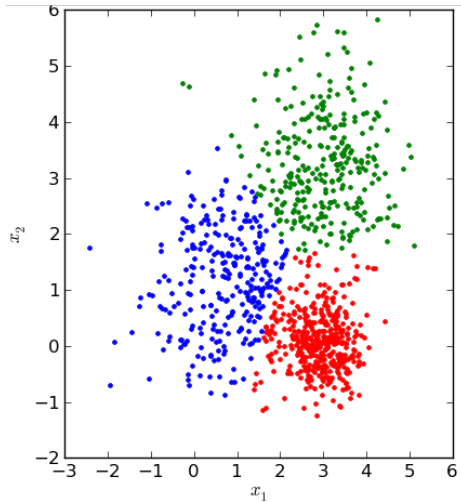
Agrupamento (Clustering): O que é?

- Em alguns domínios, gostaríamos de "classificar" exemplos em grupos semelhantes
- Porém, não possuímos rótulos de classes na base de treinamento

Exemplos

- **Agrupar** consumidores de uma plataforma web
- **Agrupar** documentos de acordo com seu tipo, tópico, etc
- **Identificar** anomalias e fraudes





Conceitos (Everitt, 1974)

- Um grupo (*cluster*) é um conjunto de entidades semelhantes entre si, enquanto entidades pertencentes a grupos diferentes não são semelhantes entre si
- A **distância** (diferença) entre quaisquer 2 instâncias em um mesmo grupo é menor que a distância entre qualquer dupla de instâncias pertencentes a grupos distintos

Mais Formalmente.....

Agrupamento em partições rígidas

- Agrupar objetos, designando instâncias X em grupos $C = \{C_1, \dots, C_k\}$, tal que:
- $C_1 \cup \dots \cup C_k = X$: Todos as instâncias devem ser designadas a um grupo
- $C_i \neq \emptyset$: Todos os grupos possuem pelo menos 1 instância
- $C_i \cap C_j = \emptyset$: Cada instância possui apenas um rótulo
- **Objetivo:** Encontrar mapeamento f :

$$f : X \rightarrow \{1, \dots, k\},$$

atribuindo cada instância a um grupo

- Enumerar e avaliar todas as partições possíveis é inviável para qualquer aplicação não-trivial

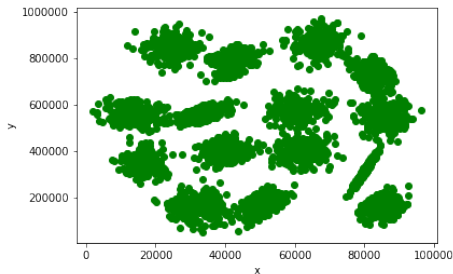
Como definir que instâncias são similares o suficiente para serem colocadas no mesmo grupo?

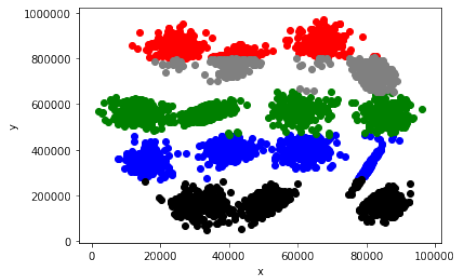
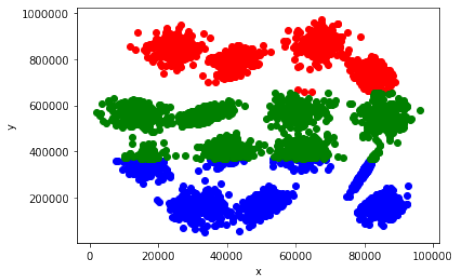
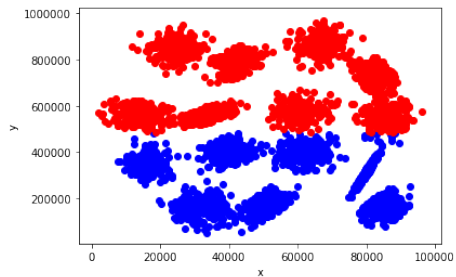
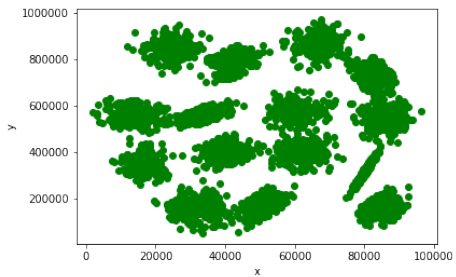
Diferenças podem ser subjetivas



- Assim como na *classificação*, algoritmos de agrupamento se baseiam nos atributos para tentar formar agrupamentos bem definidos

Avaliar qual seria o melhor agrupamento pode ser desafiador





- Assim como para classificação, cada algoritmo de agrupamento segue alguma estratégia e suposições para definir um agrupamento adequado

Algoritmo K-Médias (*K-Means*)

Separa os dados de treinamento em K grupos

Algoritmo amplamente utilizado na prática

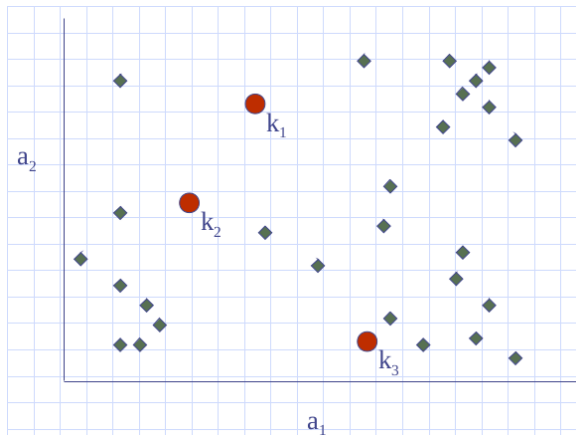
- Simples
- Fácil de se interpretar
- É eficiente computacionalmente

Partindo de um número conhecido de grupos K :

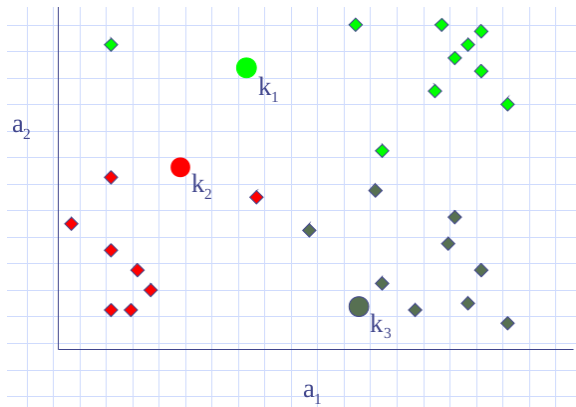
- 1 Selecione K *centróides* (instâncias que representarão cada grupo)
- 2 Atribua cada instância ao cluster que tiver o centróide mais próximo
- 3 Recomputa o centróide de acordo com a nova separação
- 4 Repita os 2 passos acima até a convergência

Exemplo Ilustrativo

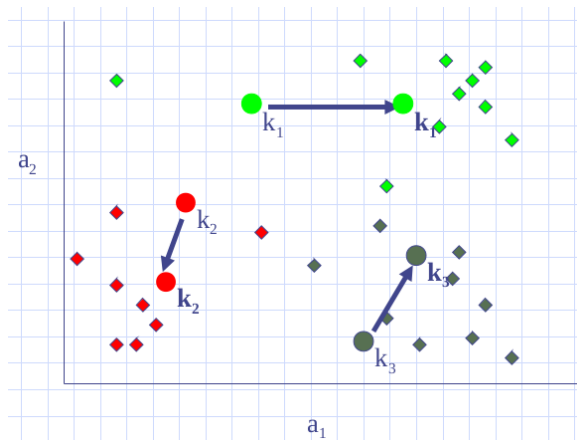
Com $K = 3$, iniciamos definindo 3 centróides



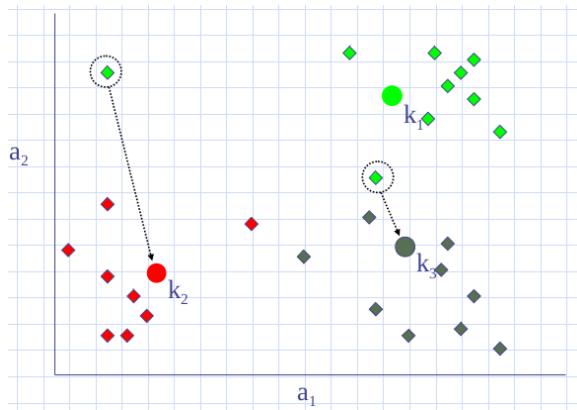
Rotular todos os exemplos de acordo com o centróide mais próximo



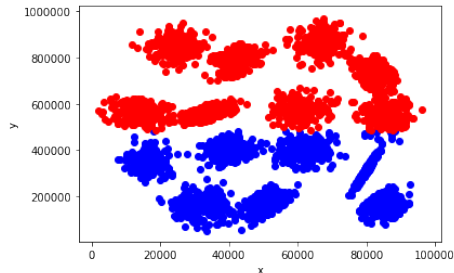
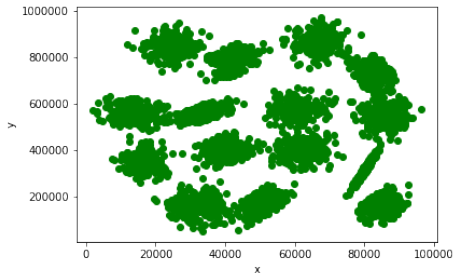
Atualizar todos os centróides



Atualizar grupos, e repetir até convergência



Exemplo com $K = 2$



- ① K
- ② Métrica de distância
- ③ Inicialização dos centróides

Distância Euclidiana

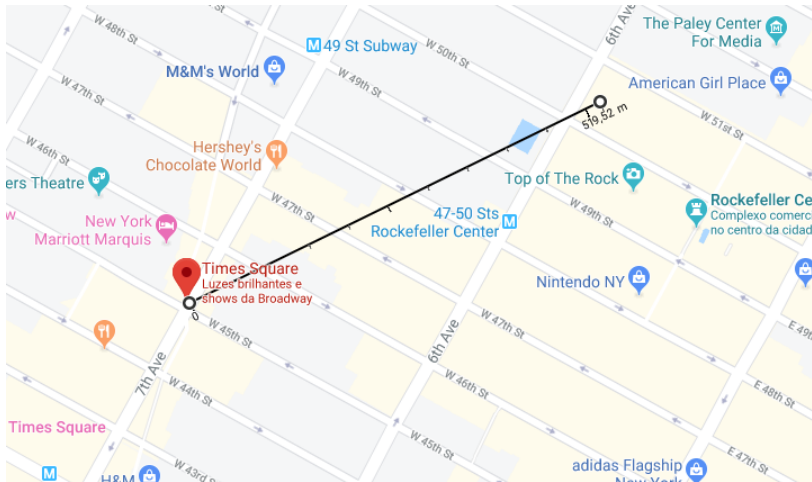
$$d_{euc}(\mathbf{a}, \mathbf{b}) = \sqrt{\sum_{m=1}^M (a_m - b_m)^2} \quad (1)$$

- Distância "comum" entre dois pontos

Exemplo Distância Euclidiana



Exemplo Distância Euclidiana



Distância Manhattan

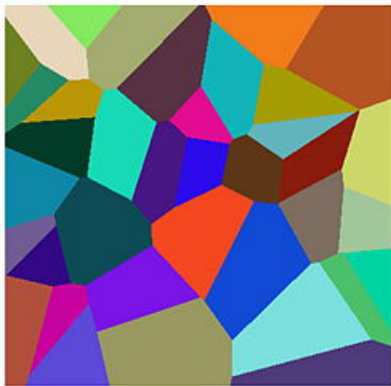
$$d_{man}(\mathbf{a}, \mathbf{b}) = \sum_{m=1}^M |a_m - b_m| \quad (2)$$

- Na analogia com distâncias espaciais, seria como "percorrer as ruas"

Exemplo Distância Manhattan



Comparação entre as distâncias



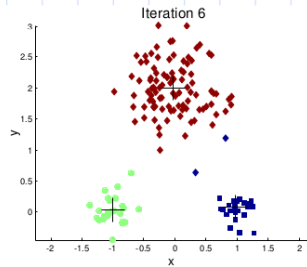
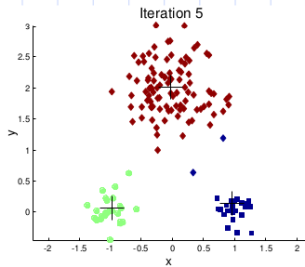
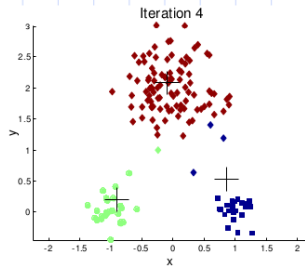
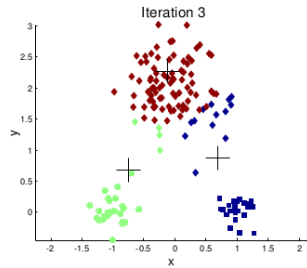
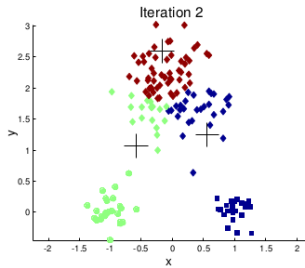
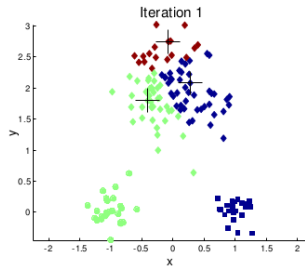
Euclidean

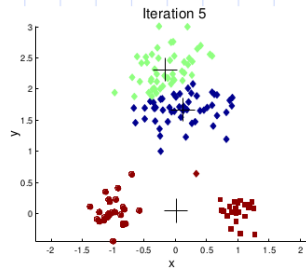
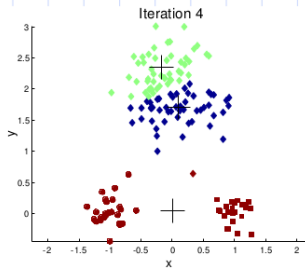
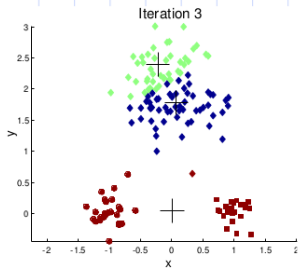
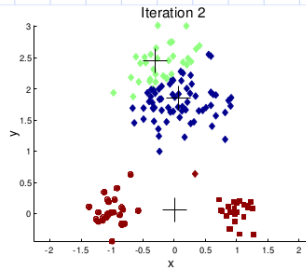
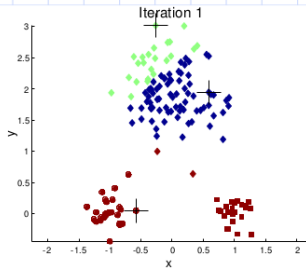


Manhattan

Inicialização dos centróides

- A inicialização dos centróides afeta no resultado final e na velocidade de convergência
- Portanto, a estratégia de inicialização afeta na efetividade do algoritmo



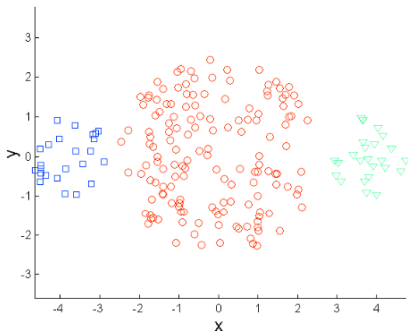


Como inicializar os centróides?

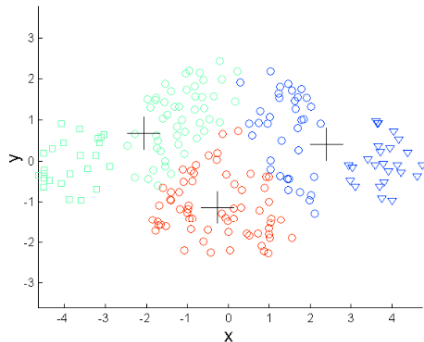
- Executar múltiplas vezes e selecionar os melhores clusters
- Seleção informada de centróides distantes entre si

Limitações do KMeans

Grupos de tamanhos diferentes



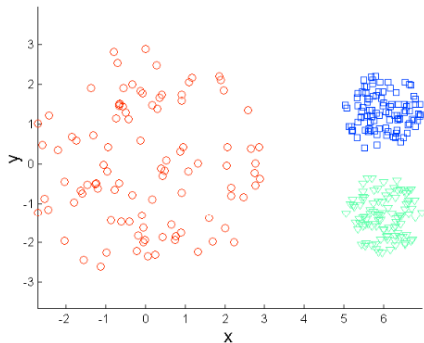
(a) Dados Gerados



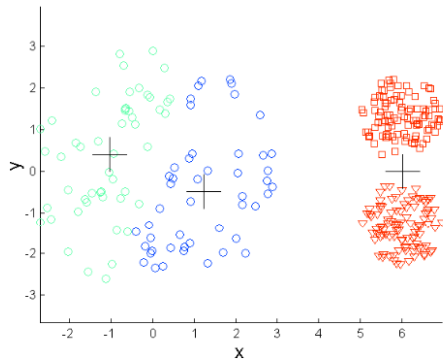
(b) $K = 3$

Limitações do KMeans

Grupos de densidades diferentes



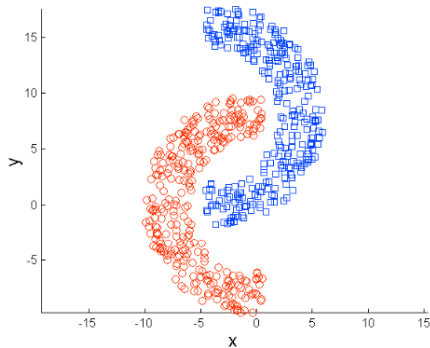
(c) Dados Gerados



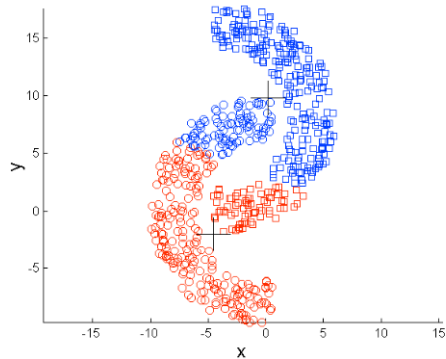
(d) $K = 3$

Limitações do KMeans

Formas não globulares

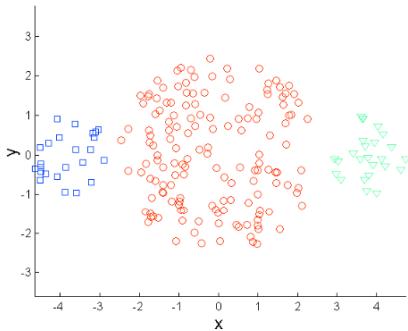


(e) Dados Gerados

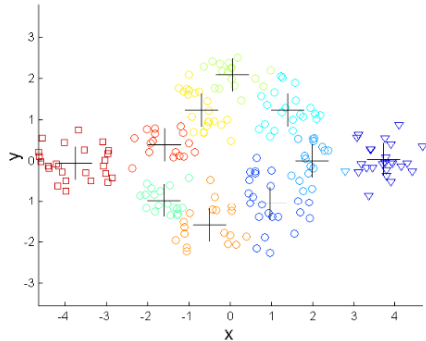


(f) $K = 2$

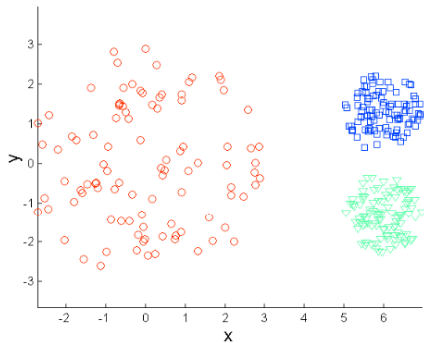
Aliviando esses problemas aumentando o parâmetro K



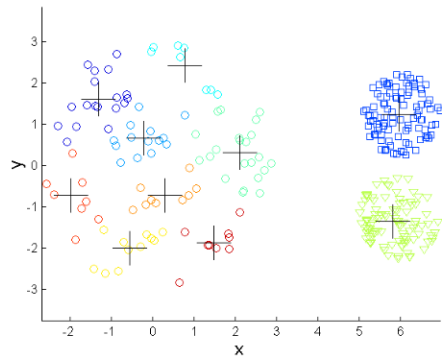
(g) Dados Gerados



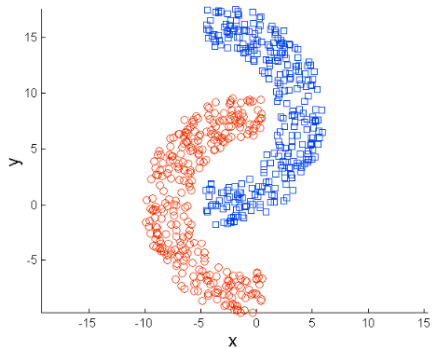
(h) Mais grupos



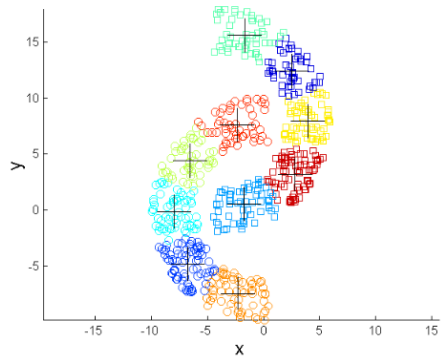
(i) Dados Gerados



(j) Mais grupos



(k) Dados Gerados



(l) Mais grupos

Pré e Pós-processamentos para o KMeans

Pré-processamento

- Normalização
- Eliminação de outliers

Pós-processamento

- Eliminar pequenos clusters
- Dividir grupos com instâncias muito distantes entre si
- Unir grupos com instâncias muito próximas

DBSCAN

- O KMeans é um algoritmo baseado em protótipos
- Algoritmos baseados em densidade definem clusteres separando regiões com alta concentração de instâncias
- O *DBSCAN* é um dos algoritmos mais utilizados dentre os baseados em densidade

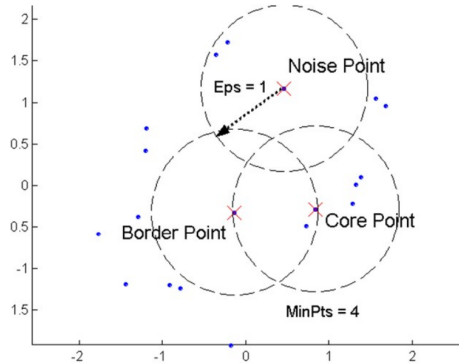
Separa as instâncias da base em 3 tipos

- **Core** - Pontos que estão no centro de uma concentração de instâncias (semelhante à ideia de centróide)
- **Border** - Pontos que fazem parte de um grupo mas não estão no centro dele, formando a "borda".
- **Noise** - Pontos isolados das outras instâncias

- 1 Percorre a base de treinamento e rotula cada exemplo como **core**, **border**, ou **noise**.
- 2 Elimina todos os exemplos rotulados como **noise**
- 3 Insere uma aresta entre cada par de exemplos **core** próximos uns dos outros
- 4 Cada componente conexo resulta em um cluster
- 5 Cada **border** é atribuído ao **core** correspondente

Como exatamente rotular cada instância? DBSCAN possui 2 parâmetros: *eps* e *minSamples*

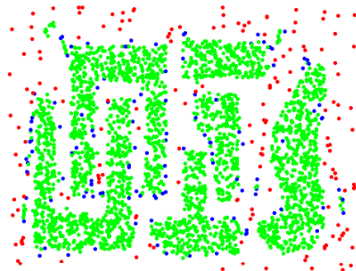
- Todas as instâncias que possuem no mínimo *minSamples* pontos a uma distância máxima de *eps* (vizinhança) são **core**
- Se um ponto possui menos que *minSamples* vizinhos mas está na vizinhança de algum ponto core, ele é **border**
- Todos os pontos não classificados como core ou border são **noise**



Exemplo definição de pontos



Original

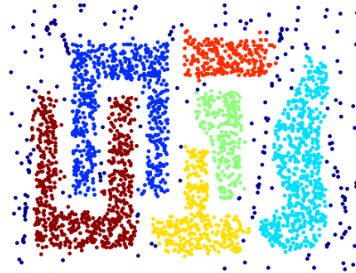


Pontos: **core**, **border**
e **noise**

Exemplo grupos



Original



Clusters

Vantagens do DBSCAN

- Resistente a ruído
- Consegue lidar com grupos de diferentes tamanhos e formas

Desvantagens do DBSCAN

- Não consegue capturar grupos com diferentes densidades