

# Uma Estratégia Híbrida Baseada em PSO para a Inicialização de Centróides no Algoritmo K-means

Clênio Eduardo da Silva  
Universidade de Uberaba - UNIUBE  
Campus Via Centro - Uberlândia  
clenio.silva@uniube.br

## RESUMO

O agrupamento de dados é uma técnica fundamental de aprendizado não supervisionado, cuja finalidade é particionar instâncias com base em seu grau de similaridade. O algoritmo K-means destaca-se por sua simplicidade e eficiência computacional; contudo, sua sensibilidade à escolha inicial dos centróides frequentemente compromete a qualidade da solução, levando à convergência prematura em ótimos locais. Neste contexto, este trabalho propõe uma abordagem híbrida que combina o K-means com o algoritmo bio-inspirado de Otimização por Enxame de Partículas (PSO), com o objetivo de aprimorar a seleção dos centróides iniciais e acelerar o processo de convergência. A principal contribuição reside na melhoria da robustez e eficiência do K-means, mantendo sua simplicidade. Experimentos conduzidos com três conjuntos de dados amplamente utilizados demonstram que o algoritmo híbrido proposto alcança soluções de alta qualidade com um número reduzido de iterações, evidenciando sua eficácia frente ao K-means tradicional.

## CCS Concepts

• Hybrid algorithms→Applied to data clusterings.

## Keywords

K-means; Particle Swarm Optimization; hybrid algorithms data clustering

## ABSTRACT

Data clustering is a fundamental unsupervised learning technique aimed at grouping instances based on their degree of similarity. The K-means algorithm is widely recognized for its simplicity and computational efficiency; however, its performance is highly sensitive to the initial selection of centroids, often leading to premature convergence to local optima. To address this limitation, this paper proposes a hybrid clustering approach that integrates Particle Swarm Optimization (PSO) with K-means to enhance the initialization process and accelerate convergence. The main contribution of this study lies in improving the robustness and effectiveness of K-means without compromising its simplicity. Experiments conducted on three well-known benchmark datasets demonstrate that the proposed hybrid algorithm consistently achieves high-quality clustering solutions with fewer iterations compared to the standard K-means, highlighting its practical advantages.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

## 1. INTRODUÇÃO

O aprendizado de máquina (AM) é um ramo da Inteligência Artificial (IA) dedicado ao desenvolvimento de algoritmos capazes de identificar padrões e tomar decisões a partir de dados, mimetizando processos cognitivos humanos [11, 13]. O AM pode ser categorizado principalmente em duas abordagens: supervisionado e não supervisionado [10].

No AM supervisionado, o modelo é treinado com dados rotulados, ou seja, instâncias acompanhadas da resposta correta, permitindo a aprendizagem de uma função de mapeamento entre entrada e saída [19]. Por outro lado, o AM não supervisionado não dispõe de rótulos durante o treinamento; seu objetivo é descobrir estruturas ou padrões intrínsecos, frequentemente por meio da formação de agrupamentos (*clusters*) de dados semelhantes [13]. Essa tarefa é conhecida como agrupamento de dados (*data clustering*) e é amplamente utilizada em diversas áreas, como mineração de dados, análise de imagens e bioinformática [1].

Entre os inúmeros algoritmos de agrupamento, o *k-means*, proposto originalmente em meados do século XX [17], destaca-se por sua simplicidade e eficiência computacional, tornando-se um método padrão para a tarefa. Entretanto, o desempenho do *k-means* é altamente dependente da escolha dos centróides iniciais, o que pode resultar em convergência para mínimos locais e soluções subótimas [4, 2].

Diante desse desafio, heurísticas e metaheurísticas bioinspiradas, como a Otimização por Enxame de Partículas (Particle Swarm Optimization, PSO) [14], têm sido aplicadas para melhorar a seleção dos centróides iniciais, buscando otimizar a qualidade do agrupamento e acelerar a convergência [3, 16]. Neste trabalho, propõe-se um algoritmo híbrido que integra PSO e *k-means*, explorando a capacidade do PSO em explorar o espaço de soluções para determinar inicializações eficazes para o *k-means*. A abordagem visa, assim, mitigar a sensibilidade do *k-means* à inicialização, promovendo convergência mais rápida e soluções de maior qualidade.

Para validar a proposta, foram realizados experimentos com três conjuntos de dados amplamente utilizados na literatura, permitindo comparações diretas com o *k-means* tradicional. Como contribuição principal, o estudo demonstra que o algoritmo híbrido PSO-*k-means* acelera significativamente o processo de convergência e mantém, ou melhora, a qualidade do agrupamento.

A estrutura deste artigo está organizada da seguinte forma: a Seção 2 apresenta os fundamentos teóricos essenciais; a Seção 3 discute trabalhos relacionados; a Seção 4 descreve a metodologia adotada, incluindo modelagem e função de avaliação; as Seções 5 e 6

apresentam os resultados experimentais e as conclusões, respectivamente.

## 2. FUNDAMENTAÇÃO TEÓRICA

Nesta seção é apresentada a teoria dos algoritmos e detalhamento dos componentes necessários para condução da metodologia apresentada.

### 2.1 Agrupamento de dados com k-means

A ideia principal do algoritmo k-means é dividir  $M$  pontos em  $N$  dimensões através de  $k$  clusters para que a soma de quadrados dentro do cluster seja minimizada [12].

O algoritmo requer como entrada uma matriz de  $M$  pontos e  $N$  dimensões, e uma matriz com  $k$  centróides de cluster iniciais em  $N$  dimensões. O número de pontos do cluster  $L$  é denotado por  $NC(L)$ . Onde  $D(X, L)$  é a distância euclidiana entre o ponto  $X$  e o cluster  $L$ , como apresentado na equação 1.

$$D(X_p, L_j) = \sqrt{\sum_{i=1}^d (X_{pi} - L_{ji})^2} \quad (1)$$

O procedimento geral é procurar uma partição  $k$  com a soma de quadrados localmente ideal dentro do cluster, movendo pontos de um cluster para outro [12]. A estratégia do algoritmo é agrupar pontos de dados de maneira que a distância euclidiana entre os pontos pertencentes a cada grupo seja minimizada. Desse modo, o algoritmo k-means tenta encontrar os melhores centróides dos grupos.

O algoritmo é dividido em dois estágios: um estágio inicial e um iterativo. O primeiro estágio envolve a definição dos  $k$  centróides; e o segundo estágio consiste no cálculo dos  $k$  novos centróides [18]. O algoritmo é finalizado quando certo critério de parada é encontrado, podendo ser número de iterações, ou caso não ocorra mais mudanças na posição dos centróides. Desse modo, dado um conjunto  $N$  de amostras, onde o objetivo é classificar os dados em  $k$  clusters, o algoritmo tende a minimizar uma função de erro, tal como o erro médio quadrático, equação 2.

$$E = \sum_{j=1}^k \sum_{i=1}^N \|x_i - c_j\|^2 \quad (2)$$

Em que  $k$  é o número de agrupamentos,  $N$  o número de amostras,  $x$  a entrada de cada amostra e  $c_j$  é o centróide. O Algoritmo 1 resume o treinamento do k-means.

---

#### Algoritmo 1: k-means [18]

---

```

Definição dos  $k$  centros aleatoriamente;
while Não houver mudanças nos centróides do
for Cada  $x_i$  de treinamento do
Atribuir  $x_i$  ao grupo associado com o centróide  $c_j$ 
mais próximo Atualizar os centróides de cada
cluster;
for cada cluster  $c_j$  do
Calcular  $c_j = \sum_{i=1}^k \frac{\|x_i - c_j\|}{N}$ 
end
end
end

```

---

### 2.2 Particle Swarm Optimization (PSO)

O algoritmo de otimização por enxame de partículas (do inglês, *Particle swarm optimization*) é um método de otimização que simula o comportamento de um bando de pássaros procurando alimento [15]. Basicamente um bando de pássaros voando aleatoriamente no espaço de busca, onde cada pássaro é uma solução (partícula).

Considere que um conjunto de partículas voa com uma determinada velocidade e se move para encontrar a melhor posição global em um processo iterativo. A cada iteração do algoritmo, o vetor de velocidade para cada partícula é modificado com base em três parâmetros: o momento da partícula, a melhor posição alcançada pela partícula e a melhor posição de todas as partículas até o estágio atual. Então, com base na velocidade determinada para cada partícula, ocorre uma movimentação para sua próxima posição. Eventualmente, é provável que o enxame, como um todo, se aproxime de um nível ótimo de função de aptidão [15].

Em um espaço  $n$ -dimensional, a posição e velocidade das partículas a cada iteração são definidas pelos vetores  $X_i^{(t)} = (x_{i1}^{(t)}, x_{i2}^{(t)}, \dots, x_{in}^{(t)})$  e  $V_i^{(t)} = (v_{i1}^{(t)}, v_{i2}^{(t)}, \dots, v_{in}^{(t)})$ , respectivamente. Usando a função de aptidão (*fitness*), cada partícula é avaliada em cada estágio do algoritmo. Com isso, a melhor posição da partícula  $i$  a cada iteração do algoritmo é atualizada no vetor  $P_i = (p_{i1}, p_{i2}, \dots, p_{in})$ , denotado por *personal best* (*pbest*). Além disso, a melhor posição global (*gbest*), considerando todas as partículas, também é registrada no vetor  $G = (g_1, g_2, \dots, g_n)$ . Em cada iteração a velocidade e posição das partículas são calculadas de acordo com as equações 3 e 4, respectivamente.

$$V_i^{(t)} = w * V_i^{(t-1)} + c_1 * r_1 * (P_i - X_i^{(t-1)}) + c_2 * r_2 * (G - X_i^{(t-1)}) \quad (3)$$

$$X_i^{(t)} = X_i^{(t-1)} + V_i^{(t)} \quad (4)$$

Onde,  $w$  é chamado peso de inércia, o qual controla o impacto da velocidade da partícula.  $r_1$  e  $r_2$  são duas variáveis aleatórias independentemente distribuídas uniformemente no intervalo  $[0, 1]$ .  $c_1$  e  $c_2$  são constantes positivas chamadas coeficientes de aceleração.

O processo de atualização da posição das partículas é repetido até que um critério de parada seja atingido. Uma condição comum usada para finalizar o algoritmo PSO é um número pré-definido de iterações. O pseudocódigo do PSO é apresentado no Algoritmo 2.

### 2.3 Bases de dados

Na tarefa de agrupamento de dados (*data clustering*) o conjunto de amostras é de grande importância. Tanto no aprendizado supervisionado como no não supervisionado é necessário uma base de dados, a qual o modelo de aprendizado irá aprender correlações com esses dados.

#### 2.3.1 Customers

O dataset *Mall Customer Segmentation Data* [20] possui dados para segmentação de clientes de um shopping, sendo composto por duzentas instâncias contendo informações básicas, como: ID, idade, sexo, renda, pontuação de gastos. A Figura 1 apresenta um *DataFrame* os atributos e dez instâncias em resumo do conjunto de amostras.

**Algoritmo 2:** PSO [15]

```

for cada partícula  $i$  do
  Inicialize a posição  $x_i$  aleatoriamente;
  Inicialize a velocidade  $v_i$  com vazio;
end
while  $N \leq totalN$  do
  for cada partícula  $i$  do
    Calcular valor de aptidão  $f(p_i)$ ;
    if  $f(p_i) \leq p_{best}$  then
       $p_{best} = f(p_i)$ ;
    end
  end
   $g_{best} = p_i$  com melhor  $f(p_i)$ ;
  for cada partícula  $i$  do
    Calcular a velocidade com a Equação 3;
    Atualizar a posição de  $p_i$  com a Equação 4;
  end
end
return  $g_{best}$ ;

```

**Figura 1:** DataFrame com amostra dos dados de customers.

CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
...	...	...	...	...	...
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

### 2.3.2 Íris

O conjunto de dados Íris é talvez a base de dados mais conhecida encontrada na literatura para reconhecimento de padrões. O conjunto contém 3 classes de 150 instâncias, onde cada classe se refere a um tipo de espécie íris e quatro atributos, comprimento e largura das sépalas e pétalas. O *dataset* pode ser obtido pela base de dados UCI Machine Learning Repository [7]. A Figura 2 apresenta os atributos juntamente com as nove primeiras instâncias do *dataset* Íris.

### 2.3.3 Wine

O *dataset* Wine é também uma base bem conhecida na literatura. Os dados do conjunto são resultados de uma análise química de três tipos de vinhos cultivados na mesma região da Itália, onde foram derivadas 13 quantidades de constituintes encontrados nos três tipos analisados. Desse modo, o *dataset* compreende 3 classes, sendo cada uma, um tipo de vinho, e 13 atributos que representam elementos da constituição do tipo de vinho [8]. O *dataset* pode ser obtido pela base de *datasets* UCI Machine Learning Repository [9]. A Figura 3 apresenta os atributos juntamente com as nove primeiras instâncias do *dataset* Wine.

**Figura 2:** DataFrame com amostra do conjunto de dados Íris.

	sepal length	sepal width	petal length	petal width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
5	5.4	3.9	1.7	0.4	Iris-setosa
6	4.6	3.4	1.4	0.3	Iris-setosa
7	5.0	3.4	1.5	0.2	Iris-setosa
8	4.4	2.9	1.4	0.2	Iris-setosa

**Figura 3:** DataFrame com amostra do dataset Wine.

	class	Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	Diluted wines	Proline
0	1	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.640000	1.04	3.92	1065
1	1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.380000	1.05	3.40	1050
2	1	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.680000	1.03	3.17	1185
3	1	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.800000	0.86	3.45	1480
4	1	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.320000	1.04	2.93	735
5	1	14.20	1.76	2.45	15.2	112	3.27	3.39	0.34	1.97	6.750000	1.05	2.85	1450
6	1	14.39	1.87	2.45	14.6	96	2.50	2.52	0.30	1.98	5.250000	1.02	3.58	1290
7	1	14.06	2.15	2.61	17.6	121	2.60	2.51	0.31	1.25	5.050000	1.06	3.58	1295
8	1	14.83	1.64	2.17	14.0	97	2.80	2.98	0.29	1.98	5.200000	1.08	2.85	1045

Bansal e Gupta [3], propuseram uma abordagem híbrida PSO-k-means que melhora a robustez do agrupamento diante de ruídos e outliers, demonstrando melhor desempenho em datasets reais e sintéticos. Liu et al. [16], apresentaram melhorias na função de aptidão do PSO para clustering, resultando em uma maior separação entre os grupos e convergência mais rápida.

Outros trabalhos recentes exploram a combinação do PSO com aprendizado profundo para agrupar dados de alta dimensionalidade, como em Zhang et al. [22], que reportaram avanços significativos em escalabilidade e precisão.

Além disso, algumas pesquisas focam em aprimorar a representação das partículas e a modelagem da função de avaliação para o PSO no contexto do k-means. Por exemplo, Chen et al. [5] propuseram uma notação adaptativa para partículas que melhora a exploração do espaço de centróides, enquanto Wang e Li [21] desenvolveram uma função de aptidão que incorpora medidas de densidade e separação para otimizar a qualidade dos clusters.

A proposta deste trabalho distingue-se das abordagens citadas ao introduzir uma nova modelagem das partículas e uma função de aptidão customizada para a PSO, especificamente projetada para otimizar a inicialização do k-means. Além disso, realizamos uma análise comparativa rigorosa em três *datasets* amplamente usados, ressaltando a eficiência da convergência e a qualidade das soluções, destacando o equilíbrio entre desempenho e custo computacional.

## 3. METODOLOGIA

O algoritmo k-means é amplamente reconhecido como uma das principais técnicas de agrupamento de dados. No entanto, sua principal limitação reside na sensibilidade à escolha inicial dos centróides, que são definidos aleatoriamente. Essa aleatoriedade pode conduzir a soluções subótimas e aumentar o número de iterações até a convergência. Neste trabalho, propomos um algoritmo híbrido, denominado k-means-PSO, que integra o

Particle Swarm Optimization (PSO) ao k-means com o objetivo de aprimorar a inicialização dos centróides e, consequentemente, a qualidade do agrupamento.

Para aplicar o PSO ao problema de agrupamento, é necessário modelar adequadamente as partículas (soluções candidatas) e definir uma função de aptidão (*fitness function*) que avalie a qualidade de cada solução gerada. Nesta seção, detalhamos os principais componentes do PSO, sua adaptação ao contexto de agrupamento, e a forma como ele foi integrado ao algoritmo k-means na abordagem proposta.

### 3.1 Modelagem das Partículas

No contexto do PSO, cada partícula representa uma possível solução para o problema, isto é, um conjunto de  $k$  centróides. A posição da partícula no espaço de busca é definida pelas coordenadas dos centróides, enquanto sua velocidade representa o deslocamento desses centróides a cada iteração. Inicialmente, as posições das partículas são definidas aleatoriamente dentro do espaço de atributos dos dados.

A velocidade é modelada como um vetor com componentes correspondentes aos deslocamentos esperados dos centróides em relação à melhor posição individual encontrada até o momento ( $p_{best}$ ) e à melhor posição global ( $g_{best}$ ). A atualização das posições e velocidades segue as equações tradicionais do PSO (Equações 3 e 4), promovendo o movimento das partículas em direção a regiões do espaço com melhor avaliação de qualidade de agrupamento.

### 3.2 Função de Avaliação

A função de avaliação, ou função de aptidão, tem por objetivo quantificar a qualidade de uma partícula, ou seja, da configuração de centróides por ela representada. A métrica utilizada foi a soma total das distâncias euclidianas entre os dados e seus respectivos centróides. Essa função visa minimizar a dispersão intra-cluster, conforme definido na Equação 5.

$$F = \sum_{i=1}^k \sum_{j=1}^m DE(c_i, p_{ij}) \quad (5)$$

Onde:

- $k$  é o número de clusters;
- $m$  é o número de instâncias atribuídas ao cluster  $i$ ;
- $c_i$  é o centróide do cluster  $i$ ;
- $p_{ij}$  é a  $j$ -ésima instância do cluster  $i$ ;
- $DE$  representa a distância euclidiana definida na equação 1.

Durante o processo de otimização, a função de aptidão é continuamente recalculada após cada atualização da posição das partículas, permitindo a identificação de  $p_{best}$  e  $g_{best}$ .

### 3.3 Integração PSO e K-means

Após um número pré-definido de iterações do PSO, a melhor partícula global ( $g_{best}$ ), que representa o conjunto de centróides mais promissor, é utilizada como ponto de partida para o algoritmo k-means. Dessa forma, o PSO atua como um mecanismo de inicialização inteligente, potencialmente conduzindo o k-means a uma convergência mais rápida e a agrupamentos de melhor qualidade.

O Pseudocódigo do algoritmo k-means-PSO é apresentado no Algoritmo 3.

---

#### Algoritmo 3: k-means-PSO

---

```
Function carrega_dataset(dataset):
|return dataset

Function gera_centroides_aleatorios(k,
dataset):
|return centroides

Function gera_centroides_PSO(k, dataset):
|return gBest

Function executa_kmeans_PSO(k, dataset, pso):
|if pso == True then
|centroides = gera_centroides_PSO(k, dataset)
|else
|centroides = gera_centroides_aleatorios(k, dataset)
|end
|while não houver mudanças nos centróides do
|for cada instância  $x_i$  do
|Atribuir  $x_i$  ao grupo com centróide  $c_j$  mais
|próximo
|end
|for cada cluster  $c_j$  do
|Atualizar  $c_j = \sum_{j=1}^k \frac{\|x_i - c_j\|}{N}$ 
|end
|end
|return clusters
```

---

### 3.4 Implementação

O algoritmo k-means-PSO foi implementado em Python no ambiente Jupyter Notebook. A função principal recebe como parâmetros o número de clusters  $k$ , o conjunto de dados, e um parâmetro booleano  $pso$  que define o método de inicialização. Se  $pso = True$ , os centróides são gerados pelo PSO; caso contrário, são inicializados aleatoriamente. O processo de agrupamento é então executado utilizando o k-means com os centróides definidos pela escolha inicial.

### 3.5 Hiperparâmetros do PSO

A eficácia do PSO depende diretamente da configuração de seus hiperparâmetros, que controlam o comportamento das partículas no espaço de busca. Na implementação deste trabalho, utilizamos os seguintes hiperparâmetros:

- $n_{particulas}$ : número de partículas na população. Valor adotado: 30;
- $n_{iteracoes}$ : número máximo de iterações do PSO. Valor adotado: 100;
- $w$ : fator de inércia, responsável por manter o movimento anterior da partícula. Valor adotado: 0,72;
- $c_1$ : coeficiente de aprendizado cognitivo, que influencia a atração pela melhor posição individual. Valor adotado: 1,49;
- $c_2$ : coeficiente de aprendizado social, que influencia a atração pela melhor posição global. Valor adotado: 1,49.

Esses valores foram selecionados com base em recomendações da literatura [6], visando um equilíbrio entre exploração do espaço de busca e convergência eficiente. Adicionalmente, a inicialização aleatória das partículas é realizada dentro dos limites definidos pelos atributos dos dados de entrada, garantindo diversidade inicial.

### 3.6 Complexidade Computacional

A complexidade computacional do algoritmo híbrido k-means-PSO pode ser analisada a partir de duas fases principais:

1. **Fase PSO:** Para cada uma das  $n_{iterações}$  do PSO, é necessário calcular a função de aptidão para cada uma das  $n_{partículas}$ , o que envolve a atribuição de instâncias aos centróides e o cálculo das distâncias. Isso resulta em uma complexidade de aproximadamente  $O(n_{partículas} * k * m * d)$ , onde  $k$  é o número de clusters,  $m$  o número de instâncias e  $d$  o número de atributos.
2. **Fase K-means:** Após a inicialização via PSO, o k-means é executado até a convergência. Sua complexidade por iteração é  $O(k * m * d)$ , sendo geralmente mais eficiente após uma boa inicialização.

Em comparação com o k-means tradicional, que é sensível à inicialização e pode demandar múltiplas execuções com diferentes sementes para garantir bons resultados, o k-means-PSO busca uma solução de maior qualidade em uma única execução, com custo adicional apenas durante a fase de otimização inicial.

## 4. EXPERIMENTOS E RESULTADOS

Os experimentos foram organizados em três cenários distintos para realizar uma análise comparativa entre os algoritmos k-means-PSO e k-means, utilizando os datasets Customers, Íris e Wine.

Na base Customers, foram considerados os atributos *Annual Income* e *Spending Score*, enquanto nas bases Íris e Wine foram utilizados todos os atributos disponíveis, exceto os atributos de classe. Para cada cenário, o número de clusters  $k$  foi definido de acordo com os dados (e.g., 5 para Customers, 3 para Íris e Wine). Cada configuração foi executada 30 vezes para garantir robustez estatística.

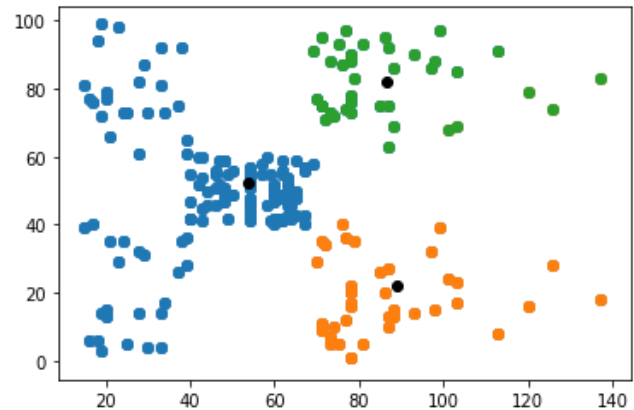
No k-means-PSO, foi realizada uma busca por parâmetros que evitassem a convergência prematura. Os parâmetros testados incluíram  $c_1$  e  $c_2$  no intervalo  $[1.0, 2.5]$ , e  $w$  no intervalo  $[0.5, 0.9]$ . Os melhores resultados foram obtidos com os valores apresentados na Tabela 1.

**Tabela 1: Melhores parâmetros identificados para o K-means-PSO.**

Parâmetros do algoritmo K-means-PSO		
Parâmetro	Descrição	Valor
$T_p$	Número de partículas	100
$N$	Número de gerações	20
$w$	Peso de inércia	0,5
$c_1$	Coefficiente cognitivo	1,0
$c_2$	Coefficiente social	1,0

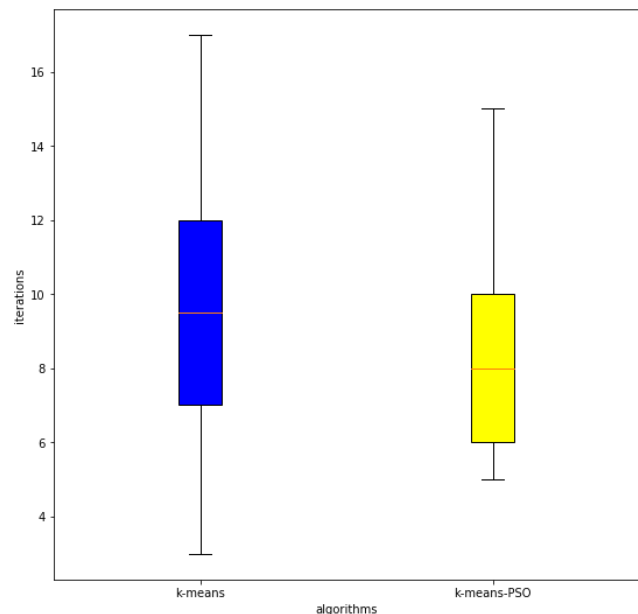
Para avaliar a performance dos algoritmos, foi analisado o número de iterações até a convergência e a função objetivo (*fitness*). A Figura 4 apresenta um gráfico de dispersão para o cenário Customers, ilustrando os agrupamentos formados pelo algoritmo k-means-PSO.

**Figura 4: Gráfico de dispersão com clusters gerados pelo k-means-PSO para a Base de dados: Customers.**



Além disso, foram gerados boxplots comparando os algoritmos quanto ao número de iterações até a convergência, conforme as Figuras 5, 6 e 7.

**Figura 5: Box Plot comparando número de iterações entre k-means e k-means-PSO. (Base de dados: Customers).**

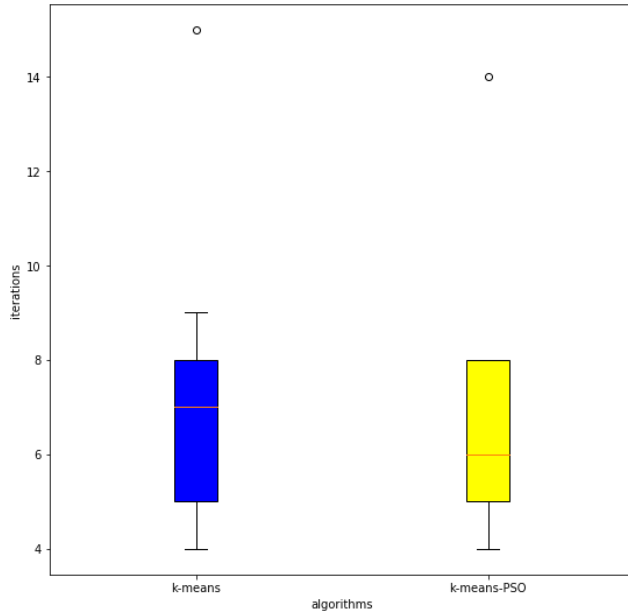


A Tabela 2 apresenta a média do número de iterações até a convergência em cada cenário.

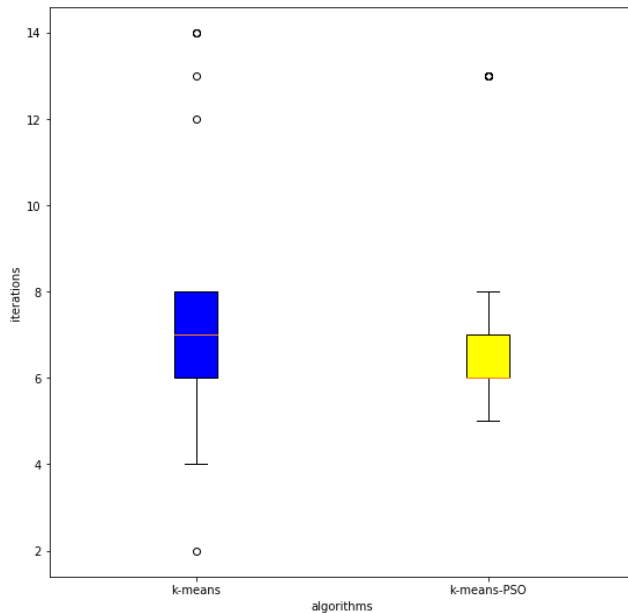
**Tabela 2: Média do número de iterações até a convergência por algoritmo e dataset.**

Dataset	k-means	k-means-PSO
Customers	9,30	8,73
Íris	6,97	6,47
Wine	7,50	7,47

**Figura 6: Box Plot comparando número de iterações entre k-means e k-means-PSO. (Base de dados: Íris).**



**Figura 7: Box plot comparando número de iterações entre k-means e k-means-PSO. (Base de dados Wine).**



#### 4.1 Validação Estatística

Para verificar a significância das diferenças observadas entre os algoritmos, foi realizado um teste *t* pareado (*paired t-test*) com nível de significância  $\alpha = 0,05$ . Os resultados são apresentados na Tabela 3.

Os valores-*p* obtidos indicam que não há diferença estatisticamente significativa entre os algoritmos nos três datasets avaliados. Contudo, a análise gráfica dos boxplots sugere que o k-means-PSO apresenta menor variabilidade nas execuções e tendência a soluções mais consistentes, o que pode ser vantajoso em aplicações onde a estabilidade do resultado é crítica.

**Tabela 3: Resultado do teste *t* pareado entre k-means e k-means-PSO para o número de iterações até a convergência.**

Dataset	Média k-means	Média k-means-PSO	valor-p
Customers	9,30	8,73	0,412
Íris	6,97	6,47	0,214
Wine	7,50	7,47	0,962

## 5. CONCLUSÃO

Este trabalho propôs uma abordagem híbrida para o algoritmo k-means, utilizando o PSO (*Particle Swarm Optimization*) como mecanismo de otimização para a escolha inicial dos centróides. Essa proposta visa mitigar a principal limitação do k-means clássico, cuja inicialização aleatória pode levar a soluções subótimas e maior número de iterações até a convergência.

Para avaliar o desempenho do algoritmo k-means-PSO, foram conduzidos experimentos em três cenários distintos utilizando os datasets Customers, Íris e Wine. Os resultados obtidos mostraram que a abordagem híbrida tende a reduzir, ainda que modestamente, o número médio de iterações até a convergência quando comparada ao k-means clássico. Além disso, os boxplots evidenciaram uma menor variabilidade nas execuções do k-means-PSO, indicando maior estabilidade na geração dos agrupamentos.

A análise estatística, por meio do teste *t* pareado, indicou que as diferenças observadas entre os algoritmos não são estatisticamente significativas ao nível de 5%, embora os resultados visuais e descritivos sugiram vantagens operacionais do uso da meta-heurística PSO, principalmente no que se refere à consistência das soluções.

Por outro lado, foi observado que o PSO pode ocasionalmente ficar preso em mínimos locais, o que compromete a qualidade dos centróides iniciais e limita os ganhos esperados na convergência do k-means. Essa limitação reforça a necessidade de explorar estratégias adicionais para guiar a busca em direção a soluções globais.

Como trabalhos futuros, pretende-se investigar a combinação do k-means com outras meta-heurísticas, como Algoritmos Genéticos, além da incorporação de mecanismos adaptativos no PSO que promovam a evasão de mínimos locais. Também se planeja aplicar a abordagem proposta em conjuntos de dados maiores e mais complexos, provenientes de aplicações reais, de modo a validar sua escalabilidade e eficácia prática. Tais avanços reforçam o potencial das abordagens híbridas na solução de problemas complexos de agrupamento em cenários reais.

## 6. REFERÊNCIAS

- [1] C. C. Aggarwal. Data Mining: The Textbook. Springer, 2014.
- [2] D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, pages 1027–1035, 2007.
- [3] S. Bansal and A. Gupta. A hybrid pso-k-means clustering algorithm for improved performance. Expert Systems with Applications, 165:113962, 2021.

- [4] M. E. Celebi, H. A. Kingravi, and P. A. Vela. A comparative study of efficient initialization methods for the k-means clustering algorithm. *Expert Systems with Applications*, 40(1):200–210, 2013.
- [5] X. Chen, J. Li, and M. Zhao. Adaptive particle representation for pso in data clustering. *Information Sciences*, 15:192–205, 2023.
- [6] M. Clerc and J. Kennedy. The particle swarm: Explosion, stability, and convergence in a multidimensional complex space. *IEEE Transactions on Evolutionary Computation*, 6(1):58–73, 2002.
- [7] R. Fisher. Iris data set. <https://archive.ics.uci.edu/ml/datasets/Iris>, 2019. Online; Acesso em 14 Dezembro 2019.
- [8] M. Forina, S. Lanteri, C. Armanino, et al. Parvus-an extendible package for data exploration, classification and correlation, institute of pharmaceutical and food analysis and technologies, via brigata salerno, 16147 genoa, italy (1988). Av. Loss Av. O set Av. Hit-Rate, 1991
- [9] M. Forina, S. Lanteri, C. Armanino, et al. Wine dataset. <https://archive.ics.uci.edu/ml/datasets/Wine>, 2019. Online; Acesso em 14 Dezembro 2019.
- [10] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.
- [11] A. Géron. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly Media, 2017.
- [12] J. A. Hartigan and M. A. Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100–108, 1979.
- [13] M. I. Jordan and T. M. Mitchell. Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245):255–260, 2015.
- [14] J. Kennedy and R. Eberhart. Particle swarm optimization. *Proceedings of IEEE International Conference on Neural Networks*, 4:1942–1948, 1995.
- [15] J. Kennedy and R. Eberhart. Particle swarm optimization. In *Proceedings of ICNN'95-International Conference on Neural Networks*, volume 4, pages 1942–1948. IEEE, 1995.
- [16] X. Liu, Y. Wang, and J. He. An improved particle swarm optimization algorithm for clustering. *Applied Soft Computing*, 87:105933, 2020.
- [17] J. MacQueen. Some methods for classification and analysis of multivariate observations. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1:281–297, 1967.
- [18] T. M. Mitchell. *Machine learning mcgraw-hill*, 1997. Google Scholar Google Scholar Digital Library Digital Library.
- [19] T. M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [20] Vijay Choudhary. Mall Customer Segmentation Data. <https://www.kaggle.com/vjchoudhary7/customer-segmentation-tutorial-in-python>, 2019. Online; Acesso em 14 Dezembro 2019.
- [21] H. Wang and F. Li. Density and separation based fitness function for pso clustering. *Swarm and Evolutionary Computation*, 61:100801, 2021.
- [22] Y. Zhang, L. Chen, and H. Wang. Deep clustering based on particle swarm optimization for high-dimensional data. *IEEE Transactions on Neural Networks and Learning Systems*, 33(1):75–88, 2022.