

# Artigo 1: Algoritmos de Agrupamento Particionado Rígido Baseados em Múltiplas Matrizes de Dissimilaridade

O **agrupamento particionado rígido** organiza um conjunto de objetos em clusters, de forma que os objetos dentro do mesmo cluster são mais semelhantes entre si e os objetos de clusters diferentes são mais distintos. A abordagem proposta neste artigo se aplica a cenários onde os dados são descritos por múltiplas matrizes de dissimilaridade. Estas matrizes podem ser geradas a partir de diferentes variáveis ou funções de dissimilaridade, proporcionando representações complementares dos dados, o que ajuda a segmentar os dados de maneira mais robusta.

Este trabalho amplia o conceito de algoritmos de agrupamento relacional rígido dinâmico, como o **SRDCA** (Self-Relational Dynamic Clustering Algorithm), permitindo o uso simultâneo de múltiplas matrizes de dissimilaridade. Além disso, introduz o conceito de **pesos de relevância** para cada matriz de dissimilaridade, que podem ser ajustados local ou globalmente durante o processo de agrupamento. A abordagem proposta é conhecida como **MRDCA-RWL** (Multiple Relational Dynamic Clustering Algorithm with Relevance Weight for Each Dissimilarity Matrix Estimated Locally), e visa melhorar a qualidade do agrupamento em cenários mais complexos, nos quais múltiplas fontes de dados (ou múltiplas matrizes de dissimilaridade) são usadas para descrever as mesmas entidades.

## Objetivo

O objetivo deste artigo é detalhar a implementação do algoritmo **MRDCA-RWL** para agrupamento particionado rígido, levando em consideração múltiplas matrizes de dissimilaridade e pesos de relevância ajustados localmente. Além disso, o artigo examina os **resultados experimentais** do algoritmo, utilizando dados sintéticos e reais para demonstrar a eficácia do método.

## Metodologia

### Representação dos Dados

- **Dados Relacionais:** Os dados são representados por matrizes de dissimilaridade  $D=[d(e_i, e_l)]$   $D = [d(e_i, e_l)]$ , onde  $d(e_i, e_l)$  representa a dissimilaridade entre os objetos  $e_i$  e  $e_l$ . As matrizes de dissimilaridade podem ser derivadas de diferentes variáveis ou funções de dissimilaridade.

### Algoritmos Propostos

O artigo apresenta três variações do algoritmo de agrupamento relacional rígido, todas baseadas em múltiplas matrizes de dissimilaridade:

1. **MRDCA:** Extensão do SRDCA para múltiplas matrizes de dissimilaridade. Assume que todas as matrizes têm a mesma relevância.
2. **MRDCA-RWL:** Estima pesos de relevância localmente, ou seja, por cluster. A relevância de cada matriz de dissimilaridade pode variar entre clusters e é ajustada iterativamente.
3. **MRDCA-RWG:** Estima pesos globalmente, ou seja, aplica os mesmos pesos para todos os clusters.

Os algoritmos têm como objetivo otimizar um critério de adequação JJ, que mede a correspondência entre os clusters e seus protótipos, levando em consideração as múltiplas matrizes de dissimilaridade e seus pesos de relevância.

## Cálculo de $\lambda_{kj}$ (Pesos de Relevância)

O cálculo dos pesos de relevância  $\lambda_{kj}$  é central no algoritmo **MRDCA-RWL**. Os pesos são ajustados iterativamente para minimizar o critério JJ, que é uma medida de adequação do agrupamento. A fórmula para o cálculo de  $\lambda_{kj}$  é dada por:

$$\lambda_{kj} = \frac{\sum_{e_i \in C_k} \sum_{e \in G_k} d_j(e_i, e) \left( \prod_{h=1}^p \sum_{e_i \in C_k} \sum_{e \in G_k} d_h(e_i, e) \right)^{1/p}}{\sum_{e_i \in C_k} \sum_{e \in G_k} d_j(e_i, e) \left( \prod_{h=1}^p \sum_{e_i \in C_k} \sum_{e \in G_k} d_h(e_i, e) \right)^{1/p}}$$

Onde:

- $C_k$  é o  $k$ -ésimo cluster.
- $G_k$  é o protótipo do cluster  $k$ .
- $d_j(e_i, e)$  é a dissimilaridade entre o ponto  $e_i$  e o protótipo  $G_k$  usando a matriz  $j$ -ésima de dissimilaridade.
- $p$  é o número de matrizes de dissimilaridade.

Os pesos  $\lambda_{kj}$  refletem a importância relativa de cada matriz de dissimilaridade para o cluster  $C_k$ . A normalização multiplicativa garante que a soma dos pesos de todas as matrizes seja igual a 1, o que impede que uma matriz domine excessivamente o processo de agrupamento.

## Implementação e Explicação do Código

A seguir, apresentamos uma implementação do algoritmo **MRDCA-RWL** em Python e C. O código em Python usa **NumPy** para lidar com as matrizes de dissimilaridade e calcular os pesos de relevância de forma eficiente.

### Código Python - MRDCA-RWL

```
import numpy as np

# Matrizes de dissimilaridade
D1 = np.array([
    [0, 1, 2, 3],
    [1, 0, 2, 2],
    [2, 2, 0, 1],
    [3, 2, 1, 0]
])
D2 = np.array([
    [0, 2, 1, 4],
    [2, 0, 3, 3],
    [1, 3, 0, 2],
    [4, 3, 2, 0]
])

D_matrices = [D1, D2]

# Inicialização
n_objects = 4
n_clusters = 2
clusters = {1: [0, 1], 2: [2, 3]} # Índices dos objetos em cada cluster
prototypes = {1: 0, 2: 3} # Protótipos iniciais
```

```

p = len(D_matrices) # Número de matrizes de dissimilaridade

def calculate_lambda(cluster, prototype, D_matrices):
    """Cálculo dos pesos lambda"""
    lambdas = []
    total_d = []
    for D in D_matrices:
        sum_d = sum(D[obj, prototype] for obj in cluster)
        total_d.append(sum_d)
    product = np.prod(total_d)
    for sum_d in total_d:
        lambdas.append((product ** (1 / p)) / sum_d)
    return lambdas

def update_clusters(D_matrices, lambdas, prototypes):
    """Atualização dos clusters baseada na dissimilaridade ajustada"""
    n_objects = D_matrices[0].shape[0]
    new_clusters = {1: [], 2: []}
    for obj in range(n_objects):
        min_distance = float('inf')
        best_cluster = -1
        for k, prototype in prototypes.items():
            distance = sum(
                lambdas[k-1][j] * D[obj, prototype]
                for j, D in enumerate(D_matrices)
            )
            if distance < min_distance:
                min_distance = distance
                best_cluster = k
        new_clusters[best_cluster].append(obj)
    return new_clusters

# Algoritmo MRDCA-RWL
def mrdca_rwl(D_matrices, clusters, prototypes, max_iter=100, tol=1e-4):
    for iteration in range(max_iter):
        print(f"\nIteração {iteration + 1}:")

        # Atualizar pesos lambda
        lambdas = {}
        for k, cluster in clusters.items():
            lambdas[k] = calculate_lambda(cluster, prototypes[k], D_matrices)
            print(f"Lambdas para Cluster {k}: {lambdas[k]}")

        # Atualizar clusters
        new_clusters = update_clusters(D_matrices, lambdas, prototypes)

        # Atualizar protótipos
        for k, cluster in new_clusters.items():
            prototypes[k] = cluster[0] # Simplesmente escolhemos o primeiro
            elemento

        # Verificar convergência
        if new_clusters == clusters:
            print("Convergiu!")
            break
        clusters = new_clusters

    return clusters, lambdas

# Rodar o algoritmo
final_clusters, final_lambdas = mrdca_rwl(D_matrices, clusters, prototypes)
print("\nClusters finais:", final_clusters)
print("Lambdas finais:", final_lambdas)

```

### Explicação do Código Python:

1. **Matrizes de Dissimilaridade:** As matrizes  $D1D1$  e  $D2D2$  representam duas diferentes formas de dissimilaridade entre os objetos. Essas matrizes são usadas no cálculo de dissimilaridade entre os objetos e seus protótipos.
2. **Função `calculate_lambda`:** Esta função calcula os pesos de relevância ( $\lambda_{kj}$  `lambda_{kj}`) para cada matriz de dissimilaridade, ajustando-os iterativamente com base na distância entre os objetos e os protótipos.
3. **Função `update_clusters`:** Após calcular os pesos de relevância, esta função atualiza os clusters, atribuindo os objetos ao cluster cujo protótipo minimiza a dissimilaridade ajustada.
4. **Função `mrdfa_rwl`:** A função

principal que executa o algoritmo **MRDCA-RWL**. Ela segue uma abordagem iterativa:

- Atualiza os pesos de relevância ( $\lambda_{kj}$  `lambda_{kj}`).
  - Atualiza os clusters.
  - Atualiza os protótipos.
  - Verifica a convergência do algoritmo.
5. **Execução:** O algoritmo é executado e os **clusters finais** e **pesos  $\lambda_{kj}$  `lambda_{kj}`** são exibidos após a convergência.

### Resultados do Código Python:

Iteração 1:

Lambdas para Cluster 1: [1.4142135623730951, 0.7071067811865476]

Lambdas para Cluster 2: [1.4142135623730951, 0.7071067811865476]

Iteração 2:

Lambdas para Cluster 1: [1.4142135623730951, 0.7071067811865476]

Lambdas para Cluster 2: [1.4142135623730951, 0.7071067811865476]

### Resultados do Código C:

Lambdas para Cluster 1:

1.41421 0.707107

Lambdas para Cluster 2:

1.41421 0.707107