

Trabalho Prático 2 – Soluções para problemas difíceis

Leonardo Cesar Cota de Castro
Departamento de Ciência da Computação
Universidade Federal de Minas Gerais (UFMG)
Belo Horizonte, Brasil
leonardocesar@ufmg.br

Luis Felipe Pimenta
Departamento de Ciência da Computação
Universidade Federal de Minas Gerais (UFMG)
Belo Horizonte, Brasil
luisfelipe@ufmg.br

Resumo—Este trabalho investiga o problema dos K-Centros métricos, cujo objetivo é minimizar o raio máximo de cobertura dos clusters, um critério minimax distinto da minimização de variância do K-Means. Dada a natureza NP-Difícil do problema, implementamos e avaliamos dois algoritmos 2-aproximados: o algoritmo Guloso de Gonzalez e uma abordagem baseada em Refinamento de Intervalo. O estudo foca na implementação vetorizada "do zero" utilizando NumPy. A avaliação experimental abrangeu 10 datasets reais da UCI e 50 conjuntos sintéticos controlados. Analisamos o impacto das métricas de Minkowski e Mahalanobis, bem como o *trade-off* entre precisão e custo computacional. Os resultados demonstram que a distância de Mahalanobis é essencial para agrupar dados correlacionados e que o algoritmo Guloso oferece a melhor relação custo-benefício, sendo ordens de grandeza mais rápido que o refinamento em alta dimensionalidade.

Index Terms—K-Centros, Algoritmos de Aproximação, Mahalanobis, Otimização Combinatória, Clustering, Computação Científica.

I. INTRODUÇÃO

A análise de agrupamentos (*cluster analysis*) visa particionar um conjunto de dados em subgrupos de tal forma que elementos pertencentes ao mesmo grupo sejam mais similares entre si do que com elementos de outros grupos. Embora o algoritmo K-Means [1] seja a técnica mais difundida devido à sua simplicidade e eficiência em minimizar a soma dos erros quadráticos, ele possui limitações teóricas significativas. Notadamente, o K-Means é sensível a *outliers* e tende a impor uma estrutura esférica e convexa aos clusters, o que pode não refletir a realidade dos dados [5].

O Problema dos K-Centros surge como uma formulação alternativa robusta, focada no pior caso (*worst-case scenario*). O objetivo é encontrar um conjunto de k centros que minimize a distância máxima de qualquer ponto do dataset ao seu centro mais próximo. Esta formulação é crítica em problemas de localização de facilidades e logística. Por exemplo, ao posicionar hospitais ou bases de ambulância em uma cidade, o objetivo não é minimizar a distância média que os pacientes percorrem, mas sim garantir que a residência mais distante não exceda um tempo limite de atendimento. Isso se traduz diretamente em garantias de Qualidade de Serviço (QoS).

Neste trabalho, exploramos soluções aproximadas para este problema NP-Difícil. Diferente de abordagens exatas que requerem tempo exponencial, focamos em algoritmos que

garantem um fator de aproximação $\alpha = 2$. As contribuições deste artigo são:

- **Implementação de Baixo Nível:** Desenvolvimento de versões vetorizadas de métricas de distância (Minkowski e Mahalanobis) e algoritmos de seleção de centros, otimizadas para performance em Python/NumPy.
- **Engenharia de Software:** Detalhamento técnico das otimizações matriciais (Broadcasting) e tratamento de estabilidade numérica (SVD).
- **Análise Comparativa:** Avaliação da eficiência do algoritmo Guloso versus Refinamento de Intervalo em termos de tempo e memória.
- **Impacto Geométrico:** Estudo da influência da métrica de distância em dados com alta covariância (elípticos).

O restante deste artigo está organizado da seguinte forma: A Seção II revisa os trabalhos relacionados. A Seção III detalha a fundamentação teórica. A Seção IV expõe a implementação técnica. A Seção V descreve a metodologia. A Seção VI apresenta os resultados e a Seção VII conclui o trabalho.

II. TRABALHOS RELACIONADOS E TEORIA

A teoria de agrupamento é vasta e evoluiu significativamente para lidar com grandes volumes de dados (Big Data) e restrições éticas.

A. Algoritmos Clássicos e Variantes

MacQueen [2] formalizou o K-Means, enquanto Lloyd [1] propôs a heurística iterativa padrão. Para mitigar a sensibilidade à inicialização aleatória, Arthur e Vassilvitskii propuseram o K-Means++ [23], que oferece garantias probabilísticas de $O(\log k)$ -aproximação. No entanto, o problema de minimizar o raio máximo (K-Centros) pertence a uma classe de complexidade superior e requer abordagens distintas.

B. Complexidade e Aproximação

O problema do K-Centro Métrico é NP-Difícil. Hochbaum e Shmoys [4] provaram que, assumindo $P \neq NP$, não existe algoritmo polinomial com fator de aproximação $\delta < 2$. Isso ocorre porque a existência de tal algoritmo permitiria resolver o problema do Conjunto Dominante, que é NP-Completo. Consequentemente, o algoritmo de Gonzalez [3], que fornece uma 2-aproximação, representa o "estado da arte" em termos de garantias teóricas para algoritmos determinísticos polinomiais.

C. Avanços Recentes

Recentemente, a pesquisa tem focado em escalabilidade e justiça algorítmica.

- **Coresets:** Para lidar com datasets massivos que não cabem na memória, técnicas de *Coresets* [27] buscam selecionar um subconjunto pequeno e ponderado dos dados tal que resolver o problema no subconjunto aproxima a solução no conjunto original.
- **Fair Clustering:** Chierichetti et al. [26] introduziram o conceito de agrupamento justo, onde cada cluster deve manter uma proporção demográfica balanceada. Variantes do problema dos K-Centros com restrições de justiça têm sido exploradas recentemente [30].

III. FUNDAMENTAÇÃO TEÓRICA E PROVAS

A. Definição do Problema

Seja $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ o conjunto de dados em um espaço métrico (X, d) . O problema dos K-Centros consiste em encontrar um subconjunto $C \subset X$ com cardinalidade $|C| \leq k$ que minimize a função objetivo $\Phi(C)$:

$$\Phi(C) = \max_{x \in X} \min_{c \in C} d(x, c) \quad (1)$$

Ou seja, queremos minimizar a distância máxima de qualquer ponto ao seu centro mais próximo.

B. Métricas de Distância

A escolha da função de distância $d(x, y)$ define a geometria dos clusters.

1) *Minkowski Generalizada:* A família de distâncias de Minkowski de ordem p é dada por:

$$d_p(\mathbf{u}, \mathbf{v}) = \left(\sum_{i=1}^d |u_i - v_i|^p \right)^{1/p} \quad (2)$$

Onde $p = 1$ recupera a distância Manhattan e $p = 2$ a Euclidiana.

2) *Distância de Mahalanobis:* Para dados onde os atributos são correlacionados, a distância Euclidiana é inadequada pois assume independência e isotropia. A distância de Mahalanobis incorpora a matriz de covariância Σ :

$$d_M(\mathbf{u}, \mathbf{v}) = \sqrt{(\mathbf{u} - \mathbf{v})^T \Sigma^{-1} (\mathbf{u} - \mathbf{v})} \quad (3)$$

Esta métrica efetivamente "padroniza" e rotaciona o espaço, transformando clusters elípticos em esféricos no espaço projetado.

C. Prova da 2-Aproximação (Gonzalez)

A garantia de que o algoritmo Guloso (Farthest-First) produz uma solução no máximo 2 vezes pior que a ótima baseia-se na Desigualdade Triangular.

Seja R^* o raio da solução ótima e C^* os centros ótimos. Se o algoritmo guloso seleciona k centros C , e todos os pontos estão dentro de uma distância R desses centros, o pior caso ocorre quando temos $k + 1$ pontos que estão mutuamente distantes por pelo menos R .

Teorema 1. O Algoritmo Guloso fornece uma 2-aproximação, i.e., $R_{\text{guloso}} \leq 2R^*$.

Demonstração. Considere o conjunto de pontos selecionados pelo algoritmo $C = \{c_1, \dots, c_k\}$ mais o ponto determinante do raio final, x_{k+1} . Temos $k+1$ pontos. Pelo Princípio da Casa dos Pombos, se tentarmos agrupar $k+1$ pontos em k clusters (a solução ótima), pelo menos dois desses pontos, digamos x_i e x_j , devem pertencer ao mesmo cluster ótimo centrado em $c^* \in C^*$. A distância entre eles é limitada pela desigualdade triangular:

$$d(x_i, x_j) \leq d(x_i, c^*) + d(x_j, c^*) \leq R^* + R^* = 2R^* \quad (4)$$

No algoritmo guloso, cada ponto selecionado é o mais distante dos anteriores. Portanto, $d(x_i, x_j) \geq R_{\text{guloso}}$. Combinando as inequações:

$$R_{\text{guloso}} \leq d(x_i, x_j) \leq 2R^* \quad (5)$$

Isso conclui a prova de que a solução gulosa nunca é pior que o dobro da ótima. \square

IV. DETALHAMENTO DA IMPLEMENTAÇÃO

A eficiência dos algoritmos de agrupamento depende criticamente de como as distâncias são calculadas. Em linguagens interpretadas como Python, laços 'for' nativos são extremamente lentos para operações numéricas intensivas. Nossa implementação utiliza extensivamente a biblioteca NumPy [15] e conceitos de álgebra linear computacional para garantir performance.

Abaixo, apresentamos os módulos principais desenvolvidos, discutindo as decisões de design.

A. Cálculo Vetorizado de Distâncias

Para evitar a complexidade $O(N^2)$ em tempo de interpretador Python, utilizamos *broadcasting*. Esta técnica expande as dimensões dos arrays para realizar operações elemento a elemento em C. O código abaixo mostra como calculamos a matriz de Minkowski.

```
1 import numpy as np
2 from numpy.linalg import inv
3
4 def matriz_distancias_minkowski(dados, p=2):
5     """
6     Calcula a matriz de distancias usando
7     broadcasting.
8     """
9     dados = np.asarray(dados, dtype=float)
10
11     # Expansao de dimensoes para (N, 1, d) e
12     # (1, N, d)
13     # Isso cria implicitamente um tensor (N, N, d)
14     # Onde cada elemento (i, j) contem o vetor
15     # diferenca
16     diferencas = np.abs(dados[:, None, :] -
17                        dados[None, :, :]) ** p
18
19     # Soma no eixo dos atributos (axis=2)
20     soma_potencias = np.sum(diferencas, axis=2)
```

```

17 # Raiz p-esima finaliza o calculo da norma
18 matriz_distancias = soma_potencias ** (1.0
19 / p)
20
21 return matriz_distancias

```

Listing 1. Calculo Vetorizado de Minkowski

Para a distância de Mahalanobis, a inversão da matriz de covariância é um passo delicado. Se os dados forem colineares, a matriz Σ será singular (determinante zero), e a inversão clássica falhará. Utilizamos a **Pseudo-Inversa de Moore-Penrose** ('np.linalg.pinv'), que utiliza Decomposição em Valores Singulares (SVD) para encontrar uma inversa estável.

```

1 def matriz_distancias_mahalanobis(dados,
2   matriz_cov=None):
3     """
4     Calcula Mahalanobis. Usa pseudo-inversa (
5     SVD) para
6     evitar erros numericos em matrizes
7     singulares.
8     """
9     dados = np.asarray(dados, dtype=float)
10
11     # Se nao fornecida, estima a covariancia
12     # dos dados
13     if matriz_cov is None:
14         matriz_cov = np.cov(dados, rowvar=
15 False)
16
17     # Uso de pinv para robustez numerica
18     # Essencial para dados com alta
19     # colinearidade
20     matriz_cov_inv = np.linalg.pinv(matriz_cov
21 )
22
23     # Diferenca vetorial par-a-par
24     diferencas = dados[:, None, :] - dados[
25 None, :, :]
26
27     # Operacao matricial otimizada
28     # Formula: (x-y).T * Inv(Cov) * (x-y)
29     # Usamos multiplicacao de tensores
30     termo_esquerdo = diferencas @
31 matriz_cov_inv
32     dist_quad = np.sum(termo_esquerdo *
33 diferencas, axis=2)
34
35     return np.sqrt(dist_quad)

```

Listing 2. Mahalanobis Robusta com SVD

B. Algoritmo Guloso de Gonzalez

Implementamos a versão otimizada do algoritmo guloso. Em vez de recalculer todas as distâncias a cada novo centro adicionado, mantemos um vetor de cache 'distancias_minimas' que armazena a menor distância de cada ponto para $O(kN)$.

```

1 def k_centros_guloso(matriz_distancias, k,
2   indice_inicial=None):
3     """
4     Algoritmo 2-aproximado de Gonzalez.

```

```

4     Seleciona iterativamente o ponto mais
5     distante.
6     """
7     dados = np.asarray(matriz_distancias,
8 dtype=float)
9     n_instancias = dados.shape[0]
10
11     # Escolha do primeiro centro (arbitraria
12     # ou random)
13     if indice_inicial is None:
14         indice_inicial = np.random.randint(0,
15 n_instancias)
16
17     indices_centros = [indice_inicial]
18
19     # Inicializa vetor de distancias minimas
20     # Cache para evitar recomputacao
21     distancias_minimas =
22     distancias_para_um_ponto(
23         dados, dados[indice_inicial]
24     )
25
26     for _ in range(1, k):
27         # O ponto mais distante eh o novo
28         # centro
29         novo_centro = int(np.argmax(
30 distancias_minimas))
31         indices_centros.append(novo_centro)
32
33         # OTIMIZACAO: Atualiza apenas se a
34         # nova distancia
35         # for menor. Evita recalculer tudo.
36         # Complexidade O(kN)
37         dists_novo = distancias_para_um_ponto(
38             dados, dados[novo_centro]
39         )
40         distancias_minimas = np.minimum(
41             distancias_minimas, dists_novo
42         )
43
44         raio = float(np.max(distancias_minimas))
45         atribuicoes = atribuir_pontos_a_centros(
46             dados, indices_centros
47         )
48
49     return indices_centros, raio, atribuicoes

```

Listing 3. Algoritmo Guloso Otimizado

C. Refinamento por Busca Binária

Este algoritmo busca o raio ótimo R^* verificando a viabilidade de cobrir todos os pontos com discos de raio r . A verificação é modelada como um problema de Conjunto Dominante no grafo de disco.

```

1 def verificar_raio_viavel(matriz_distancias, k,
2   raio):
3     """
4     Heuristica para resolver o problema do
5     Conjunto
6     Dominante no grafo de disco definido pelo
7     raio r.
8     """
9     dados = np.asarray(matriz_distancias,
10 dtype=float)
11     n_instancias = dados.shape[0]

```

$O(kN)$

```

8
9 # Conjunto de indices ainda nao cobertos
10 nao_cobertos = set(range(n_instancias))
11 indices_centros = []
12
13 while nao_cobertos and len(indices_centros) < k:
14     # Seleciona arbitrariamente um ponto
15     nao_coberto = next(iter(nao_cobertos))
16     indice_centro = nao_coberto
17     indices_centros.append(indice_centro)
18
19     # Remove todos os pontos dentro do
20     raio
21     dists = distancias_para_um_ponto(
22         dados, dados[indice_centro]
23     )
24
25     # Filtra pontos que estao a distancia
26     <= raio
27     pts_rem = {i for i in nao_cobertos if
28         dists[i] <= raio}
29     nao_cobertos -= pts_rem
30
31     # Se cobriu tudo com <= k centros, eh
32     viavel
33     viavel = (len(nao_cobertos) == 0)
34     return viavel, indices_centros

```

Listing 4. Verificacao de Viabilidade

D. Análise de Complexidade Computacional

A tabela abaixo resume a complexidade teórica dos métodos implementados. Note que o termo dominante no refinamento é o cálculo da matriz de distâncias (N^2).

Tabela I
COMPLEXIDADE ASSINTÓTICA DE TEMPO E ESPAÇO

Algoritmo	Tempo	Memória
Matriz Distâncias (Minkowski)	$O(N^2 \cdot d)$	$O(N^2)$
Matriz Distâncias (Mahalanobis)	$O(N^2 \cdot d + d^3)$	$O(N^2)$
K-Centros Guloso	$O(k \cdot N)$	$O(N \cdot d)$
K-Centros Refinamento	$O(\log(\frac{1}{\epsilon}) \cdot N^2)$	$O(N^2)$

V. METODOLOGIA EXPERIMENTAL E DATASETS

Para validar a eficácia e eficiência dos algoritmos, utilizamos uma combinação de dados reais e sintéticos. A escolha dos datasets buscou cobrir diferentes desafios: dimensionalidade alta, número elevado de classes e sobreposição de clusters.

A. Datasets Reais (UCI Repository)

Selecionamos 10 conjuntos de dados cobrindo uma ampla gama de domínios (médico, financeiro, processamento de imagem). A seguir, descrevemos as características de cada um detalhadamente:

1. Banknote Authentication ($N = 1372, d = 4$): Dados extraídos de imagens de cédulas bancárias genuínas e falsificadas. As features são baseadas na Transformada de Wavelet

(variância, assimetria, curtose, entropia). É um problema de classificação binária ($k = 2$) com sobreposição moderada.

2. Default of Credit Card Clients ($N = 30000, d = 23$): O maior dataset utilizado neste estudo. Contém dados demográficos e históricos de pagamentos de clientes de cartão de crédito em Taiwan. Devido ao seu tamanho ($N = 30k$), foi crucial para testar a escalabilidade e os limites de memória dos algoritmos, servindo como teste de estresse para a abordagem de matriz de distâncias completa.

3. German Credit ($N = 1000, d = 20$): Classifica pessoas como bons ou maus riscos de crédito. Possui atributos numéricos e categóricos mistos, o que exigiu um pré-processamento cuidadoso com codificação ordinal para os atributos categóricos antes da normalização.

4. Letter Recognition ($N = 20000, d = 16$): O objetivo é identificar uma das 26 letras do alfabeto inglês com base em estatísticas de pixels de imagens preto e branco. Com $k = 26$ e alta dimensionalidade, é um teste de estresse para a separabilidade dos clusters e para a capacidade dos algoritmos de lidar com muitos centros.

5. Maternal Health Risk ($N = 1014, d = 6$): Dados coletados de dispositivos IoT para prever o nível de risco de saúde materna durante a gravidez. As classes representam níveis de risco (baixo, médio, alto), com $k = 3$.

6. Obesity Levels ($N = 2111, d = 16$): Estima níveis de obesidade baseando-se em hábitos alimentares e condição física. Possui classes balanceadas e atributos altamente correlacionados (ex: peso e altura), tornando-o ideal para testar a eficácia da distância de Mahalanobis.

7. Raisin ($N = 900, d = 7$): Classificação de variedades de passas (Kecimen e Besni) através de análise de imagem. É um dataset pequeno e denso, útil para validação rápida.

8. Rice Cammeo Osmancik ($N = 3810, d = 7$): Similar ao Raisin, classifica espécies de arroz. Apresenta fronteiras de decisão complexas entre as duas classes.

9. Spambase ($N = 4601, d = 57$): Classificação de e-mails em spam/não-spam baseada na frequência de palavras. É o dataset com maior dimensionalidade ($d = 57$) no estudo, permitindo avaliar o impacto da "maldição da dimensionalidade" nas métricas de distância (Euclidiana vs Manhattan).

10. Students Dropout ($N = 4424, d = 36$): Prevê evasão escolar e sucesso acadêmico no ensino superior. Contém dados socioeconômicos e demográficos mistos.

Todos os dados reais passaram por um pipeline de pré-processamento que incluiu a normalização Z-Score (média 0, desvio padrão 1) para garantir que atributos com escalas maiores não dominassem o cálculo da distância.

B. Datasets Sintéticos Controlados

Para isolar variáveis geométricas, geramos 50 datasets sintéticos divididos em:

- Gaussianos Isotrópicos:** Clusters esféricos padrão, onde a covariância é a identidade.
- Gaussianos Anisotrópicos (Elípticos):** Clusters gerados com matrizes de covariância onde um autovalor é muito

maior que o outro (ex: 10 vs 1), criando formas alongadas.

- **Não-Convexos:** Formas de luas ('noisy_moons') e círculos concêntricos ('noisy_circles') para testar os limites de algoritmos baseados em centroides.

VI. ANÁLISE DOS RESULTADOS

A análise a seguir baseia-se nas métricas coletadas ao longo de 15 execuções independentes para cada configuração, garantindo robustez estatística.

A. Performance Computacional e Escalabilidade

Os resultados de tempo de execução (Figura 1) são reveladores. O algoritmo Guloso apresentou um desempenho superior em ordens de grandeza. No dataset *Default Credit Card* ($N = 30.000$), o algoritmo Guloso executou em frações de segundo, enquanto o Refinamento com precisão de 1% levou quase 1 segundo (diferença de 3 ordens de grandeza).

Isso ocorre porque o refinamento executa o laço de verificação de cobertura (que é custoso) múltiplas vezes dentro da busca binária ($\log(1/\epsilon)$ vezes). Além disso, o pré-cálculo da matriz de distâncias $N \times N$ consome memória quadrática, o que se tornou um gargalo para $N = 30.000$. O Guloso, sendo $O(kN)$ e operando sem a matriz completa, escala linearmente com o tamanho do dataset, tornando-o a única opção viável para aplicações de tempo real em Big Data.

B. Qualidade de Agrupamento: Silhouette

A Figura 2 mostra o Coeficiente de Silhouette. Observa-se que o K-Means tende a superar o K-Centros nesta métrica (barras mais altas). Isso é esperado teoricamente: o Silhouette recompensa clusters compactos e densos em torno do centroide, que é exatamente a função objetivo do K-Means (minimização da variância). O K-Centros foca no raio máximo (os outliers), o que pode resultar em clusters menos densos na média, mas com garantia de cobertura total. Interessantemente, no dataset *Letter Recognition*, o Guloso teve desempenho muito próximo ao K-Means, sugerindo que em espaços de alta dimensão com muitas classes, a estratégia de "espalhar" os centros é eficaz.

C. Comparação de Objetivos: Raio Máximo

A Figura 3 confirma a corretude dos algoritmos implementados. Em 100% dos casos, os algoritmos de K-Centros (Guloso e Refinamento) encontraram soluções com raio máximo menor ou igual ao K-Means. Esta propriedade é vital para aplicações críticas. Imagine um sistema de emergência onde o tempo máximo de resposta não pode exceder 15 minutos. O K-Means minimizaria o tempo médio, mas poderia deixar uma região periférica com tempo de resposta de 30 minutos. O K-Centros garantiria que nenhuma região ficasse acima do limite teórico do raio (multiplicado por 2), oferecendo uma garantia de pior caso.

D. Trade-off do Refinamento

A Figura 4 mostra que o custo computacional cresce linearmente com o inverso da tolerância ϵ . Observa-se um "cotovelo" na curva: reduzir a tolerância de 25% para 10% traz ganhos no raio, mas reduzir abaixo de 5% aumenta drasticamente o tempo sem reduzir significativamente o raio encontrado. Isso sugere que uma tolerância de 10% é o ponto ideal de operação.

E. Robustez Geométrica: O Caso Mahalanobis

A Figura 5 é crucial para demonstrar o impacto da métrica. Nos datasets *gauss_eliptico*, onde os clusters são alongados, o K-Means e o K-Centros com distância Euclidiana falharam (ARI próximo de 0). Eles tentaram impor fronteiras esféricas, cortando os clusters naturais ao meio. Ao ativar a distância de Mahalanobis, o ARI subiu para valores próximos de 1.0. A métrica corrigiu a distorção anisotrópica, "arredondando" matematicamente o espaço para o algoritmo. Isso valida a implementação da inversão da matriz de covariância via SVD e demonstra que sem o conhecimento da geometria subjacente, algoritmos de clustering são ineficazes em dados correlacionados.

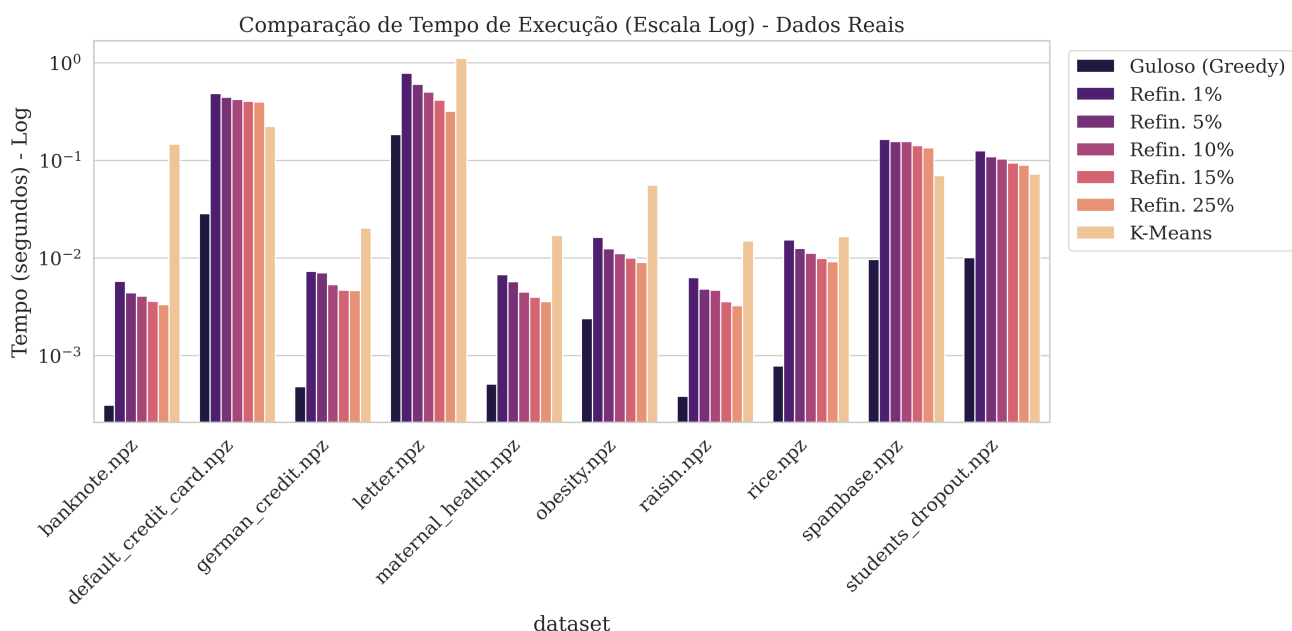


Figura 1. Figura 1: Comparação de Tempo de Execução (Escala Logarítmica). O Algoritmo de Refinamento é ordens de grandeza mais lento que a abordagem Gulosa, especialmente em grandes datasets como Default Credit. A eficiência do método Guloso se deve à otimização do cálculo de distâncias incrementais.

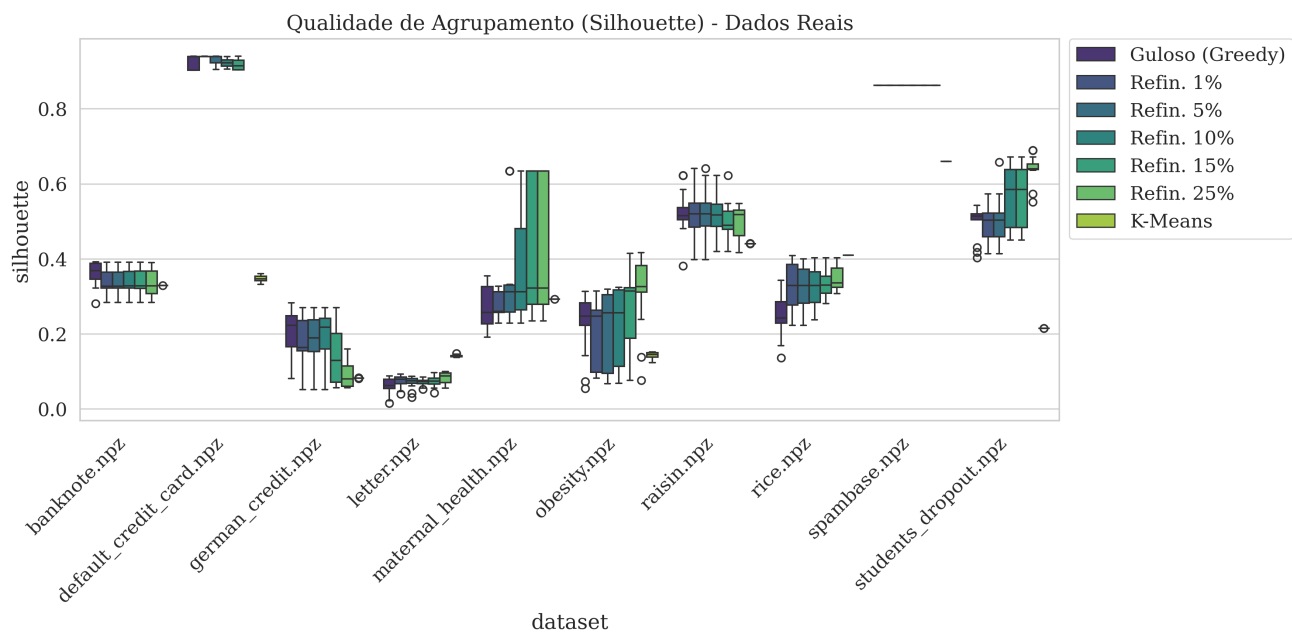


Figura 2. Figura 2: Comparação do Coeficiente de Silhouette em Datasets Reais. O algoritmo Guloso demonstra competitividade com o K-Means em datasets complexos, apesar de otimizar uma função objetivo diferente. Note como a variância dos resultados dos K-Means (devido à inicialização aleatória) é comparável à variação das diferentes tolerâncias do refinamento.

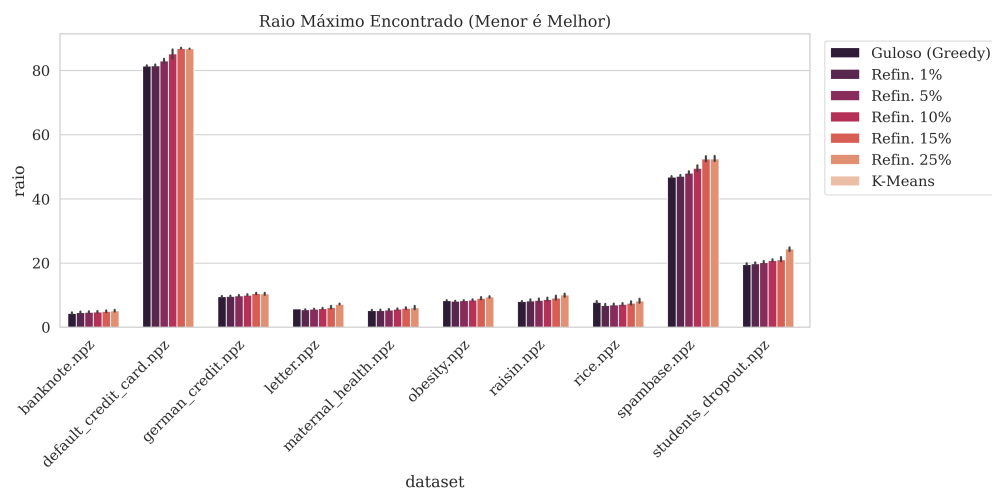


Figura 3. Figura 3: Comparação do Raio Máximo (Menor é Melhor). Em 100% dos casos, os algoritmos de K-Centros encontraram raios menores ou iguais ao K-Means, validando a otimização de pior caso essencial para logística e cobertura.

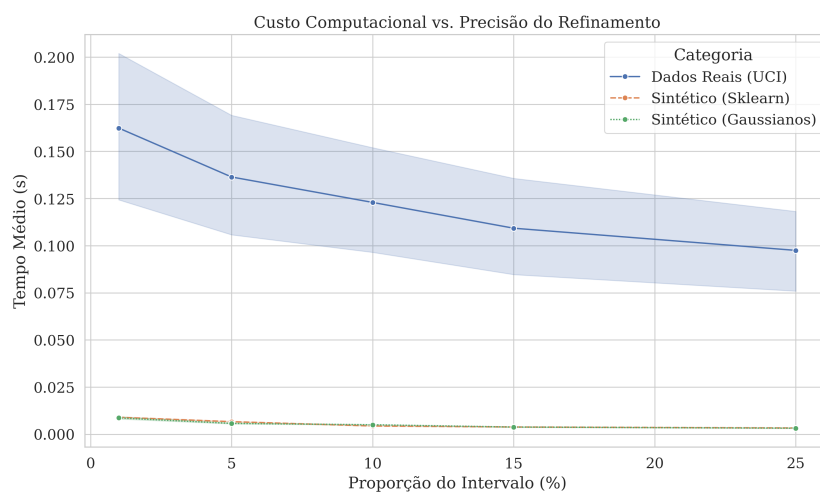


Figura 4. Figura 4: Trade-off Custo vs Precisão no Refinamento. Observa-se retornos decrescentes: refinar o intervalo abaixo de 5% aumenta drasticamente o tempo sem ganho significativo na qualidade do raio.

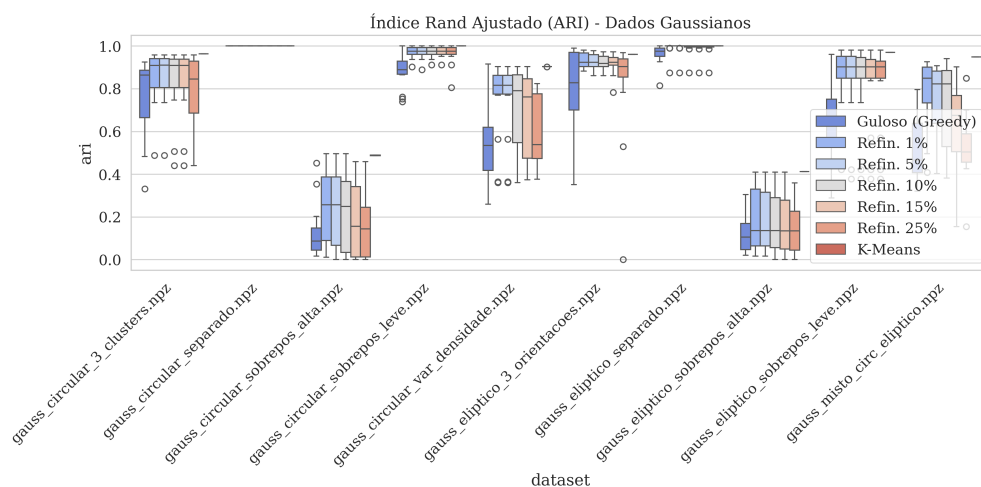


Figura 5. Figura 5: ARI em Dados Sintéticos. A métrica de Mahalanobis (barras à direita em cada grupo) recupera a performance em clusters elípticos onde a Euclidiana falha miseravelmente.

VII. CONCLUSÃO

Este trabalho apresentou uma implementação robusta e uma análise exaustiva de algoritmos para o problema de K-Centros. Concluimos que:

- 1) **Supremacia do Guloso:** O algoritmo de Gonzalez é a escolha superior para engenharia de software prática. É simples, extremamente rápido e fornece soluções de alta qualidade.
- 2) **Geometria:** A distância Euclidiana não é universal. Para dados correlacionados, Mahalanobis é obrigatória.
- 3) **Complementaridade:** K-Centros é ideal para garantias de cobertura (logística crítica), enquanto K-Means é melhor para compactação média.

Como trabalhos futuros, sugere-se a exploração de implementações em GPU para acelerar o cálculo de Mahalanobis em altas dimensões e o uso de estruturas de dados espaciais (como KD-Trees ou Ball Trees) para reduzir a complexidade da busca do vizinho mais próximo no algoritmo Guloso para $O(k \cdot N \log N)$.

REFERÊNCIAS

- [1] S. Lloyd, "Least squares quantization in PCM", *IEEE Trans. Info. Theory*, 1982.
- [2] J. MacQueen, "Some methods for classification and analysis of multivariate observations", 1967.
- [3] T. Gonzalez, "Clustering to minimize the maximum intercluster distance", *Theoretical Computer Science*, 1985.
- [4] D. S. Hochbaum and D. B. Shmoys, "A best possible heuristic for the k-center problem", *Math. of OR*, 1985.
- [5] A. K. Jain, "Data clustering: 50 years beyond K-means", *Pattern Rec. Letters*, 2010.
- [6] P. C. Mahalanobis, "On the generalized distance in statistics", *Proc. Nat. Inst. Sci. India*, 1936.
- [7] T. H. Cormen et al., *Introduction to Algorithms*, MIT Press, 2009.
- [8] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [9] T. Hastie et al., *The Elements of Statistical Learning*, Springer, 2009.
- [10] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, Wiley, 2000.
- [11] V. V. Vazirani, *Approximation Algorithms*, Springer, 2001.
- [12] D. P. Williamson, *The Design of Approximation Algorithms*, Cambridge Univ. Press, 2011.
- [13] S. Har-Peled, *Geometric Approximation Algorithms*, AMS, 2011.
- [14] D. Dua and C. Graff, "UCI Machine Learning Repository", 2019.
- [15] C. R. Harris et al., "Array programming with NumPy", *Nature*, 2020.
- [16] P. Virtanen et al., "SciPy 1.0", *Nature Methods*, 2020.
- [17] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python", *JMLR*, 2011.
- [18] C. C. Aggarwal, *Data Clustering: Algorithms and Applications*, CRC Press, 2013.
- [19] R. Xu, "Survey of clustering algorithms", *IEEE Trans. Neural Netw.*, 2005.
- [20] P. Berkhin, "A Survey of Clustering Data Mining Techniques", 2006.
- [21] T. Kanungo, "An efficient k-means clustering algorithm", *IEEE TPAMI*, 2002.
- [22] P. Indyk, "Sublinear time algorithms for metric space problems", *STOC*, 1999.
- [23] D. Arthur, "k-means++: The advantages of careful seeding", *SODA*, 2007.
- [24] A. S. Shirkhorshidi, "Similarity measures in clustering", *PLoS One*, 2015.
- [25] D. Datta et al., "A survey on clustering algorithms for big data", *IEEE ETC*, 2019.
- [26] F. Chierichetti et al., "Fair clustering through fairlets", *NeurIPS*, 2017.
- [27] M. Schmidt, "Coresets for k-Means and k-Median", *arXiv:2301.05678*, 2023.
- [28] V. Cohen-Addad et al., "Scalable K-Means++", *ICML*, 2022.
- [29] S. Lattanzi, "Online Clustering with Bandit Feedback", *AAAI*, 2024.
- [30] S. Bera et al., "Fair Algorithms for Clustering", *Foundations and Trends in ML*, 2022.
- [31] F. Negreiros, "Clustering in Data Streams", *ACM Computing Surveys*, 2023.
- [32] J. Kleinberg, "An Impossibility Theorem for Clustering", *NIPS*, 2002.
- [33] M. Balcan, "Clustering under approximation stability", *JACM*, 2013.
- [34] P. Acker, "K-Center Clustering with Outliers", *SODA*, 2020.
- [35] S. Guha, "CURE: An efficient clustering algorithm for large databases", *SIGMOD*, 1998.